

User's Guide

Provalis Research

1998-2021 All right Reserved - Provalis Research

Table of Contents

Introduction to WordStat 9.0	1
Program Capabilities	2
The Content Analysis & Categorization Process	7
A Quick Tour	9
Performing Content Analysis	
Starting WordStat from QDA Miner or SimStat	12
Creating a Project in WordStat	14
Creating a Project from a List of Documents	
Creating a WordStat Project from Windows Explorer	17
Notes on Importing PDF Files	
Creating a Project from an Existing Data File or Web Service	19
Importing from a Data File	20
Importing from Web Surveys	
SurveyMonkey	24
Qualtrics	
SurveyGizmo	
Voxco	
QuestionPro	
TripleS v.1.2	
Importing from Social Media	
Importing from RSS	
Importing from Twitter	
Importing from Facebook	
Importing from Reddit	
Importing Bibliographic References	
EndNote	
Mendeley	
RID FILE	
Novis I INI	49
Factiva	53
Importing from Email Servers	55
Outlook	
Hotmail	
Gmail	
MBox	
Using the Document Conversion Wizard	63
Starting WordStat Document Explorer from Windows Explorer	65
The User Interface	68
The Data Tab	70
Proiect Properties	70

Appending Documents	74
Appending from a Data File	
Monitoring a File or Folder	
Deleting Cases	
Identifying Duplicate Cases	
Setting the Case Descriptor	84
Filtering Cases	85
Editing Droiget Data	20
Adding New Variables	
Adding New Variables	00
Reordering Variables	
Transformation	
Changing Variable Types	
Recoding Values of a Variable	
Transform Document	
Binning Numerical Variables	
Performing Complex Numeric Transformation	
xBase functions	
Editing Variable Properties	
Variable Statistics	105
Analyze	109
The Text Processing Tab	111
Language	113
Preprocessing	
Notes on Preprocessing	116
Writing a Preprocessing Routine with Python	117
Calling an Executable Program	
Calling a Function in a DL	119
	120
Substitution	10/
Substitution	
Substitution Categorization	
Substitution	
Substitution	124 127 134 138
Substitution Categorization Creating and Maintaining Dictionaries Working with Rules Working with Regular Expressions	124 127 134 138 140
Substitution	124 127 134
Substitution	124 127 134
Substitution	124 127 134 138 140 140 142 144
Substitution	124 127 134 138 140 142 142 144 144
Substitution	124 127 134 138 140 140 142 144 144 145 145
Substitution	124 127 134 138 140 142 144 144 144 145 146 147
Substitution	124 127 134 138 140 142 144 144 144 144 145 146 147 154
Substitution	124 127 134 138 140 142 144 144 144 145 146 147
Substitution	124 127 134 138 140 142 142 144 144 144 145 146 147 154 156
Substitution Categorization Creating and Maintaining Dictionaries Working with Rules Working with Regular Expressions Importing Dictionaries Exporting Dictionaries Printing Dictionaries Using Lexical Tools for Dictionary Building Basic Dictionary-Building Tools Advanced Dictionary-Building Tools Postprocessing The Frequencies Tab Using the Dictionary Panel	124 127 134 138 140 142 142 144 144 145 145 146 147 154 156 159
Substitution Categorization Creating and Maintaining Dictionaries Working with Rules Working with Regular Expressions Importing Dictionaries Exporting Dictionaries Printing Dictionaries Using Lexical Tools for Dictionary Building Basic Dictionary-Building Tools Advanced Dictionary-Building Tools Postprocessing The Frequencies Tab Using the Dictionary Panel Working with the Suggestions Tab	124 127 134 138 140 142 142 144 144 145 145 146 147 154 156 159 161
Substitution Categorization Creating and Maintaining Dictionaries Working with Rules Working with Regular Expressions Importing Dictionaries Exporting Dictionaries Printing Dictionaries Printing Dictionaries Using Lexical Tools for Dictionary Building Basic Dictionary-Building Tools Advanced Dictionary-Building Tools Postprocessing The Frequencies Tab Using the Dictionary Panel Working with the Suggestions Tab Barcharts, Pie Charts, and Word Clouds	124 127 134 138 138 140 142 144 144 144 145 145 146 147 154 156 159 161 162
Substitution Categorization Creating and Maintaining Dictionaries Working with Rules Working with Regular Expressions Importing Dictionaries Exporting Dictionaries Printing Dictionaries Using Lexical Tools for Dictionary Building Basic Dictionary-Building Tools Advanced Dictionary-Building Tools Postprocessing The Frequencies Tab Using the Dictionary Panel Working with the Suggestions Tab Barcharts, Pie Charts, and Word Clouds Keyword Retrieval	124 127 134 138 138 140 142 144 144 145 145 146 147 154 156 159 161 162 165
Substitution Categorization Creating and Maintaining Dictionaries Working with Rules Working with Regular Expressions Importing Dictionaries Exporting Dictionaries Printing Dictionaries Using Lexical Tools for Dictionary Building Basic Dictionary-Building Tools Advanced Dictionary-Building Tools Postprocessing Using the Dictionary Panel Working with the Suggestions Tab Barcharts, Pie Charts, and Word Clouds Keyword Retrieval Running and Creating Post-Processing Scripts	124 127 134 138 138 140 142 144 144 145 145 146 147 154 159 159 161 162 165 170
Substitution Categorization Creating and Maintaining Dictionaries Working with Rules Working with Regular Expressions Importing Dictionaries Exporting Dictionaries Printing Dictionaries Using Lexical Tools for Dictionary Building Basic Dictionary-Building Tools Advanced Dictionary-Building Tools Postprocessing The Frequencies Tab Using the Dictionary Panel Working with the Suggestions Tab Barcharts, Pie Charts, and Word Clouds Keyword Retrieval Running and Creating Post-Processing Scripts	124 127 134 138 138 140 142 144 144 145 145 146 147 154 159 159 161 162 165 170
Substitution Categorization Creating and Maintaining Dictionaries Working with Rules Working with Regular Expressions Importing Dictionaries Exporting Dictionaries Printing Dictionaries Using Lexical Tools for Dictionary Building Basic Dictionary-Building Tools Advanced Dictionary-Building Tools Postprocessing The Frequencies Tab Using the Dictionary Panel Working with the Suggestions Tab Barcharts, Pie Charts, and Word Clouds Keyword Retrieval Running and Creating Post-Processing Scripts The Extraction Tab	124 127 134 138 138 140 142 144 144 144 145 146 147 154 159 159 161 162 165 170
Substitution Categorization Creating and Maintaining Dictionaries Working with Rules Working with Regular Expressions Importing Dictionaries Exporting Dictionaries Printing Dictionaries Using Lexical Tools for Dictionary Building Basic Dictionary-Building Tools Advanced Dictionary-Building Tools Advanced Dictionary-Building Tools Postprocessing The Frequencies Tab Using the Dictionary Panel Working with the Suggestions Tab Barcharts, Pie Charts, and Word Clouds Keyword Retrieval Running and Creating Post-Processing Scripts The Extraction Tab Topics	124 127 134 138 140 142 144 144 144 144 145 146 147 156 159 161 162 170 177 177
Substitution	124 127 134 138 140 142 144 144 144 144 145 146 147 156 159 161 162 170 177 182
Substitution	124 127 134 138 140 142 144 144 144 144 145 146 147 156 159 161 162 170 177 182 187
Substitution Categorization Creating and Maintaining Dictionaries Working with Rules Working with Regular Expressions Importing Dictionaries Exporting Dictionaries Printing Dictionaries Using Lexical Tools for Dictionary Building Basic Dictionary-Building Tools Advanced Dictionary-Building Tools Advanced Dictionary-Building Tools Postprocessing The Frequencies Tab Using the Dictionary Panel Working with the Suggestions Tab Barcharts, Pie Charts, and Word Clouds Keyword Retrieval Running and Creating Post-Processing Scripts The Extraction Tab Topics Phrases Named Entities Misspellings & Unknowns	124 127 134 138 140 142 144 144 144 144 145 146 147 156 159 161 162 170 177 182 187 188
Substitution	124 127 134 138 140 142 144 144 144 144 145 146 147 154 155 156 157 161 162 170 177 182 187 188 195
Substitution Categorization Creating and Maintaining Dictionaries Working with Rules Working with Regular Expressions Importing Dictionaries Exporting Dictionaries Printing Dictionaries Printing Dictionaries Using Lexical Tools for Dictionary Building Basic Dictionary-Building Tools Advanced Dictionary-Building Tools Postprocessing The Frequencies Tab Using the Dictionary Panel Working with the Suggestions Tab Barcharts, Pie Charts, and Word Clouds Keyword Retrieval Running and Creating Post-Processing Scripts The Extraction Tab Topics Phrases Named Entities Misspellings & Unknowns	124 127 134 138 140 142 144 144 144 144 144 145 146 147 154 155 156 159 161 162 170 177 182 183 195 105
Substitution	124 127 134 138 140 142 144 144 144 144 144 145 146 147 154 155 159 161 162 170 177 177 182 183 195 195 100
Substitution	124 127 134 138 140 142 144 144 144 144 145 146 147 154 155 159 161 162 163 170 177 177 182 188 195 195 198 200
Substitution	124 127 134 138 140 142 144 144 144 144 145 146 147 154 155 159 161 162 163 170 177 177 182 188 195 198 200
Substitution	124 127 134 138 140 142 144 144 144 144 145 146 147 154 155 159 161 162 163 170 177 178 182 195 195 198 200 202

Link Analysis	206
Statistics	
Rubble Charts	
Heatmap Plot	
Correspondence Analysis	226
Keyword-In-Context Page	229
The Classification Tab	234
Settings	235
Select Features	236
Learn & Test	
Classification Experiment Dialog Box	
Missellenseus	
wiscenaneous	250
Preparing and Importing Data	250
Preliminary Text Preparation	250
Importing Spreadsheet or Database Files	
Importing Plain Text or Word Processor Files from Simstat	
Creating Comparison Charts	
Edit Case Descriptors	
Filtering Cases	
Geocoding	
Mapping	
Text Editor	269
Publishing Categorization Models	272
Creating and Using Norm Files	272
Brogram Sottings	
Program Settings	
Exporting Frequency Data	276
Performing Multivariate Analysis	278
Performing Analysis on Manually Entered Codes	280
Web Collector	281
Word Frequency Analysis	285
WordStat Software Development Kit (SDK)	288
Report Manager	
Working with the Table of Contents	
Creating a New Item	
Importing Items from Files	
Renaming an Item	291
Deleting an Item	
Moving all liell	
Adding or Editing item Comments	
Editing Documents	292
Editing Tables	

	Editing Charts	293
	Searching and Replacing Text	293
	Exporting items to HTML or Word	294
Re	eferences on Content Analysis	297

Introduction to WordStat 9.0

WordStat is a text analysis module specifically designed to study textual information such as responses to open-ended questions, interviews, titles, journal articles, public speeches, electronic communications, etc. WordStat may be used for automated categorization of text using a dictionary approach or various text mining methods. WordStat can apply existing categorization dictionaries to a new text corpus. It also may be used in the development and validation of new categorization dictionaries. WordStat can now be used as a stand-alone software. When used as in conjunction with QDA Miner, this module can provide assistance for a more systematic application of coding rules, help uncover differences in word usage between subgroups of individuals and assist in the revision of existing coding using KWIC (Keyword-In-Context) tables.

WordStat includes numerous exploratory data analysis and graphical tools that may be used to explore the relationship between the content of documents and information stored in categorical or numeric variables such as the gender or the age of the respondent, satisfaction scores, publication date, etc. Relationships among words or categories as well as document similarity may be identified using hierarchical clustering, topic modeling and multidimensional scaling analysis. Correspondence analysis and heatmap plots may be used to explore relationship between keywords and different groups of individuals.

WordStat can be used as a stand-alone software or it can be run as a module of either of the following base products:

QDA Miner - The text management and qualitative analysis program allows you to create and edit data files, import documents, and perform manual coding and annotation of the documents. Several analysis tools are also available to look at the frequency of manually assigned codes and the relationship between these codes and other categorical or numeric variables. When used with QDA Miner, WordStat can perform content analysis on whole documents or selected segments of the documents tagged with specific user defined codes.

Simstat - This statistical software provides a wide range of statistical procedures for the analysis of *quantitative* data. It offers advanced data file management tools such as the ability to merge data files, aggregate cases, perform complex computation of new variables and transformation of existing ones. When used with Simstat, WordStat can analyze textual information stored in any alphanumeric, plain text and rich text memo variable (or field). It includes various tools to explore the relationship between any numeric variable of a data file and the content of alphanumeric ones.

Stata - Stata is a complete, integrated software package that provides all your data science needs—data manipulation, visualization, statistics, and automated reporting.

The WordStat module may be accessed the first two applications from the **CONTENT ANALYSIS** command in the **ANALYSIS** menu. When installed in Stata, WordStat is accessible from the **WORDAT | CONTENT ANALYSIS** command in the **USER** menu.

A few additional utility programs are also included with WordStat that may be run as standalone applications or be accessed directly through WordStat:

- **Report Manager** This application has been designed to store, edit and organize documents, notes, quotes, tables of results, graphics and images created by QDA Miner and WordStat or imported from other applications.
- Document Conversion Wizard This utility program provides an easy way to import numerous documents and create a project file. It can also be used to split large files into smaller units and to extract various numeric and alphanumeric data from structured documents.
- Document Classifier This utility program is a stand-alone application that may be used to perform content analysis
 and automatic text classification on a text pasted from the clipboard or stored in a file. For more information on this
 utility program, see <u>WordStat Document Explorer</u>.
- Dictionary Builder This tools allows the development of comprehensive categorization dictionary for automatic content analysis. The program may be run as standalone application but also from Text Processing tab of WordStat by pressing the Suggest button. <u>Click here to obtain more information on the Dictionary Builder</u>.

For a detailed listing of WordStat features, see Program's capabilities.

Program Capabilities

Content Analysis Capabilities

- Text analysis of textual documents.
- Full Unicode support allows one to analyze almost any language.
- Performs analyses on alphanumeric variables containing short textual information such as responses to open ended questions, titles, descriptions, etc. as well as on longer documents stored as plain ASCII or as rich text (RTF) document.
- Stemming in 24 human languages (Arabic, Armenian, Basque, Catalan, Czech, Danish, Dutch (Kraaij-Pohlmann), English (Lovins, Porter, Snowball), Finnish, French, German, Hungarian, Irish, Italian, Latin, Norwegian, Portuguese, Romanian, Russian, Slovene, Spanish, Swedish, Tamil and Turkish)
- Automatic lemmatization (available in English, French, Spanish, Swedish, German, Norwegian, Polish and Italian), contact us if you need support of other languages.
- Substitution process for customized lemmatization of words or automatic spell correction of common misspellings.
- Optional exclusion of pronouns, conjunctions, expressions, etc, by the use of existing or user-defined stop word lists (a.k.a. exclusion list).
- Automated categorization of words, word patterns, or phrases using existing or user-defined categorization dictionaries.
- Advanced proximity rules with proximity operators (NEAR, AFTER, BEFORE, etc.) and word distance within sentence, paragraphs or documents may be used for word disambiguation.
- Frequency analysis of words, content categories or concepts.
- Phrase finder allows identification of the most recurring phrases.
- Pattern-based named-entity extraction allows identification of people, product names, places and acronyms.
- Choice between automatic spelling correction and semi-automatic misspelling identification and processing.
- Interactive development and validation of multi-level categorization dictionaries or taxonomies
- Easy drag and drop assignments of words, phrases, or topics to categorization dictionaries.
- Ability to restrict an analysis to specific portions of a text or to exclude comments and annotations.
- Ability to perform a content analysis on a random sample of cases.
- Integrated spell-checking with support for 20 languages.
- Integrated thesaurus and lexical database (English, French, Spanish, German, Norwegian, Polish, Italian) to assist the creation of comprehensive categorization dictionaries.
- Case filtering on any numeric or alphanumeric variable and on keyword occurrence (with AND, OR, and NOT Boolean operators).
- Printing of presentation quality tables.
- Export any output table to SPSS, Stata, Excel, HTML, JSON, XML, ASCII, Tab separated or comma separated value files.

• All graphics may be saved to disk in BMP, JPEG or PNG file format.

Univariate Keyword Frequency Analysis

- Univariate keyword frequency analysis (keyword count and case occurrence).
- Keyword cooccurrence matrix (within documents, paragraphs, sentences, etc.).
- Keyword by case data matrix.

Relationship to Numerical, Categorical Data, or Dates

- Comparison between any textual variable and values of one or two nominal or ordinal variables (such as sex of the respondent, specific subgroups, years of publication, etc.).
- Automatic recoding of date variables into days, weekdays, weeks, months, quarters, years or decades.
- Choice between 12 different association measures to assess the relationship between keyword occurrence and nominal or ordinal variables (Chi-square, Likelihood ratio, Student's F, Tau-a, Tau-b, Tau-c, symmetric Somers' D, asymmetric Somers' Dxy and Dyx, Gamma, Pearson's R, and Spearman's Rho)
- Correspondence analysis allows examination of relationships between words or categories and other nominal or ordinal variables.
- Deviation table for easy identification of characteristic words or content categories for values of categorical or numerical variables.
- Heatmap plot and dual hierarchical clustering may be used to identify functional relationship between keywords and values of categorical or numerical variables.
- Ability to sort keyword matrix in alphabetical order by keyword frequency or keyword occurrence, on the obtained statistics or on its probability.

Keyword cooccurrence and Analysis

- · Integrated clustering and dendrogram display of keyword cooccurrence
- Powerful topic modeling module with unique topic enrichment features for more accurate assessment of topic (including suggestions for topic names. associated phrases, exceptions, and spelling corrections)
- 2-D and 3-D multidimensional scaling on cooccurrence of words or content categories.
- Link analysis using force-based graphs, multidimensional scaling or circular graph displays.
- Proximity plot to easily identify all keywords that cooccurs with one or several target keywords.
- Flexible keyword cooccurrence criteria (within a case, a sentence, a paragraph, a window of n words, a userdefined segment) as well as clustering methods (first- and second-order proximity, choice of similarity measures).
- Easy text retrieval directly from dendrogram or proximity plots.

Advanced Topic Modeling

• Automatic topic extraction with non-negative matrix factorization (NNMF) and factor analysis (additional algorithms such as LDA may also be performed using R or Python scripts).

- Automatic generation of topic names.
- Topics may be merged, deleted and renamed. One may also remove topic items.
- Unique topic enrichment feature for automatic inclusion of associated phrases, and generation of suggested phrases, exceptions and spelling corrections.
- Push topic solutions to the crosstabulation feature of WordStat for computation of comparison statistics and graphics, correspondence analysis, bubble plots, deviation table, etc..
- Perform co-occurrence analysis on topic solutions (hierarchical clustering, multidimensional scaling, proximity plot, etc.)
- Save topic solutions as a categorization model.for further editing or to apply topic model to new datasets.

Analysis of Case or Document Similarity

• Hierarchical clustering, multidimensional scaling and proximity plot may be used to explore the similarity between documents or cases.

Norm Creation and Comparison

- Ability to create norm files based on frequency analysis of words or content categories.
- Comparison of obtained frequencies to previously saved norm files.

Multiple Responses and Comparisons Between Variables

- Can perform a single frequency analysis on information stored in several text variables (documents or short alphanumeric variables)
- Comparison of keyword occurrence between variables.

Automated Text Classification

- Machine learning algorithms (Naive Bayes and K-Nearest Neighbors) for document classification.
- Flexible feature selection for automatic selection of best subsets of attributes.
- Numerous validation methods (leave-but-one, n-fold crossvalidation, split sample).
- Experimentation module allows easy comparison of predictive models and fine-tuning of classification models.
- Classification models may be saved to disk and applied later using either a standalone document classification utility program, a command line program or a programming library (the command line and the programming library are part of a new software developer's kit (SDK) sold separately).

Keyword-In-Context

- Ability to display a Keyword-In-Context (KWIC) table of any included keywords, leftover or user defined word, word pattern or phrase.
- KWIC tables may be sorted in ascending order of case number, context, or on values of categorical or numerical variables.

- Ability to jump from a specific occurrence in the KWIC table to the original text variable in order to view or edit the selected document.
- KWIC tables may be saved as a new data file for further processing.
- Customizable KWIC and report function to display all hits as lists of paragraphs, sentences or user defined segments.
- Perform quick word frequency analysis on context words appearing before, after, or both before and after.

Keyword Retrieval Function

- A powerful keyword retrieval function allows identification of text units (documents, paragraph or sentences) containing one keyword or a combination of keywords with optional filtering of cases.
- Ability to attach QDA Miner codes to retrieved text segments.
- Retrieved segments may be exported to disk in tabular format (Excel or delimited text files) or as text reports (Rich Text Format).
- Perform quick word frequency analysis on retrieve results.

Supports R and Python pre- and post-processing scripts

- Real time text transformation may be performed using Python or R scripts, allowing one to perform various tasks such as cleaning text data, tokenizing Asian languages, performing part-of-speach tagging,
- Post-processing scripts in R and Python may be created to analyze either original sources or document-term matrices, allowing one to extend the analytical capabilities of WordStat with additional machine learning algorithms, predictive models, etc.
- Ability to create dialog boxes for customizing the execution of scripts. allowing the execution of flexible scripts by non-programmers.
- Automatically import text files, tables and images outputs for easy review.

Full Integration with QDA Miner and SimStat

- Document variables are stored in the same file as all other numeric and categorical variables.
- Variable selection and analysis are performed within the main statistical program using a simple 3-step operation:
 - 1. Open the existing data file.
 - 2. Select one or several alphanumeric fields as dependent variables and, optionally, other nominal or ordinal variables to be treated as independent.
 - 3. Execute the **CONTENT ANALYSIS** command from the **ANALYZE** drop-down menu.
- New variables representing frequency or occurrence of words, keywords or concepts can be added to the existing data file or exported to a new data file in order to be submitted to more advanced analysis (such as cluster analysis, correspondence analysis, multiple regression, etc.).
- Data can be imported from and exported to different file format including Excel, SPSS, Stata, comma or tab separated text files, etc.
- Ability to perform numeric and alphanumeric transformation or to apply filters on cases of the data file to restrict the analysis to specific subgroups.

Integration with Stata

- WordStat may be configured to run from the Stata USER menu to analyze text data stored in the currently open dataset.
- The Document Conversion Wizard may be used to create Stata .dta project files out of MS Word, RTF, HTML, TXT and PDF documents.
- When running WordStat as a standalone, it is possible to import Stata files from version 8 to version 17 as well as export any project to corresponding Stata file formats.

The Content Analysis & Categorization Process

The most basic form of content analysis WordStat can perform is a simple frequency analysis of all words contained in one or several text variables of a data file. However, WordStat offers several features that permit the user to accomplish more advanced forms of content analysis that may involve automated categorization, different weighting of words and phrases, inclusion or exclusion of words based on frequency criteria, etc. To fully understand the possibilities offered by the program, you first need to understand the various underlying processes involved in a typical WordStat frequency analysis and how these processes may be combined to achieve various kinds of content analysis tasks.

WordStat's categorization involves up to seven consecutive processes:

1- Text Preprocessing (including stemming, n-grams, etc.)

The preprocessing option allows users to access external text preprocessing routines that are not part of the WordStat program. This option is useful to perform custom transformation on the text to be analyzed. WordStat includes a few sample text processing routines, such as a Porter stemmer which remove common English suffixes and prefixes as well as a a character n-grams routine which decomposes every word into sequences of 3, 4 or 5 characters. Please note that the Porter stemmer routine available in the preprocessing process is for demonstration purpose only and will greatly reduce the processing speed of WordStat. We recommend using instead the integrated stemming routine (see #2 below)

2- Stemming

The stemming process is a natural language processing routine that reduces inflected words to a common stem or root form. The English stemmer, for example, returns "write" for "write", "writes", "writing", and "writings". Stemming routines are available for these languages: Arabic, Armenian, Basque, Catalan, Czech, Danish, Dutch (Kraaij-Pohlmann), English (Lovins, Porter, Snowball), Finnish, French, German, Hungarian, Irish, Italian, Latin, Norwegian, Portuguese, Romanian, Russian, Slovene, Spanish, Swedish, Tamil and Turkish

3- Substitution Process (including lemmatization and automatic spell correction)

The substitution process takes individual words and replace them with another word form or with a sequence of words. Such a process is typically used for lemmatization, a procedure by which all plurals are transformed into singular forms and past-tense verbs are replaced with present-tense versions. It may also be used for derivational stemming in which different nouns, verbs, adjectives and adverbs derived from the same root word are transformed into this single word. Custom substitution process may also be created to perform automatic spelling corrections of common misspelling.

4- Exclusion Process

An exclusion process may be applied to remove words that you do not want to be included in the content analysis. This process requires the specification of an exclusion list. Such a process is used mainly to remove words with little semantic value such as pronouns, conjunctions, etc., but may also be used to remove some words or phrases used too frequently or with little discriminative value.

5- Categorization Process

The categorization process allows you to change specific words, word patterns, or phrases to other words, keywords or content categories and/or to extract a list or specific words or codes. This process requires the specification of a categorization dictionary. This dictionary may be used to remove variant forms of a word in order to treat all of them as a single word. It may also be used as a thesaurus to perform automatic coding of words into categories or concepts. For example, words such as "good", "excellent" or "satisfied" may all be coded as instances of a single category named "positive evaluation", while words like "bad", "unsatisfied" or expressions like "not satisfied" may be categorized as "negative evaluation".

6- Addition of Frequent Words

The fifth process is the application of a frequency criterion that is used to add to the included words and categories words that are used more than a specific number time or that are found in more than a specific number of cases. When a categorization dictionary is used, this option will append to this list of included words or categories, all words that meet the minimum frequency criterion. If no categorization dictionary is used, all words that meet this minimum

requirement and that have not been excluded (see process #3) will be added to the final word/category list. Note that this process can only be used to add new words to the actual list of words and categories found in this categorization dictionary. It cannot be used to remove any of these items (see process #6).

7- Removal or Words or Categories

When this process is applied, all words or categories that do not meet a minimum frequency or case occurrence criterion will be removed from the final list. It can also remove items occurring in too many cases. This process may be combined with the categorization process (#4) to remove infrequent or too common categories. It may also be used in conjunction with the addition criterion (see process #5) to provide a composite criterion of inclusion that involves both a minimum word frequency and a minimum case occurrence.

Since the application of each process is optional, numerous combinations are possible, each combination allowing the researcher to perform different types of content analysis. For example, here are the minimal requirements for different forms of content analysis:

Types or Analysis	Preprocessing & Substitution	Exclusion List	Categorization Dictionary	Add Words	Remove Words	Comments
Simple word frequency analysis (most frequent words)				V*		
Simple frequency analysis of semantically significant words		V		V*		
Word count with lemmatization	V	V		V*		
Word count of specific words			Ń			
Automatic categorization of texts			V			
Frequency analysis on the most frequent categories.			V		V	
Frequency analysis of manually entered codes or keywords			Ń			Codes may optionally be inserted between brackets
Rating of texts on specific attributes			V			Weights may be assigned to different words

* The use of minimum frequency criteria is recommended when you want to perform analysis on only the most frequent words.

A Quick Tour

A Content Analysis on Personal Ads

For this example, we will produce a content analysis on personal ads published in a Montreal cultural newspaper on January 22 and January 29, 1998, and we will examine whether there is a relationship between words used and the gender and age of the person who wrote the ad.

The required data has been stored in a file named SEEKING.PPJ.

Performing Content Analysis

The required data has been stored in a file named SEEKING.PPJ.

Step #1 - Open or create a project.

• Open the data file in either QDA Miner or SimStat, select your variables and run the content analysis module. If you have WordStat as a stand-alone software, create a project in WordStat.

Step #2 - Choose the proper dictionaries

WordStat consists of a single dialog box with seven or eight tabs. The second tab, **Text Processing**, allows you to select, view, and edit the dictionaries used in this specific content analysis. Set the dictionaries to the following values:

Exclusion: ENGLISH

Categorization: SEEKING

and make sure both of them are enabled. (see the checkboxes on the tabs)

Step #3 - Setting the proper options

Still on the **Text Processing** tab **Preprocessing** tab allows you to specify various options such as whether numeric values should be included etc. Move to the **Postprocessing** tab, here you can decide which frequent words should be added etc. Disable all options by removing any check mark in the various check boxes.

Step #4 - Perform a univariate frequency analysis on categories

- Click the third tab (Frequency). The program will perform a categorization of words found in the ads and compute a frequency analysis on the categories.
- By default, the words displayed in the matrix are those specified in the Categorization dictionary. To display words that have been left out, click the **Leftover words** tab.
- To move a word to the categorization dictionary or the exclusion list, right-click the mouse. (Click here for more information on the <u>creation and maintenance of dictionaries</u>).

Step #5 - Examining the relationship between included categories and the gender of the author.

- Press on the sixth tab (Crosstab).
- Click the WITH drop-down list and select GENDER to display a contingency table of category frequency by gender.

The TABULATE option allows you to choose whether the table should be based on the total frequency of included words or on the total number of cases containing those words.

The SORT BY option allows you to sort the table on the word or category name (alphabetical order) or by descending order of keyword frequency. You may also click any column header to sort the grid in ascending or descending order of the values found in this column.

The DISPLAY option allows you to specify the information displayed:

- Count
- Row percent
- Column percent
- Total percent
- Category percent (for case occurrences)
- Percent of total words (for keyword frequency)

Step #6 - Estimating the strength of the relationship

- Use the STATISTIC drop-down list to select an association measure, such as a Chi-square or a Pearson's R statistic.
- To sort the table on the chosen statistic or on its probability, use the SORT BY drop-down list.

Step #7 - Visualizing the relationship between some categories and the gender of the author.

- Use the mouse to highlight cells of the categories you would like to compare.
- Click the 💹 button or right-click the mouse and select the Chart Selected Rows menu item.

Step #8 - Performing correspondence analysis on age groups

- Click the WITH drop-down list box and select AGEGROUP to display keyword counts by age group.
- Click the 🗄 button to access the correspondence analysis dialog box.
- Press on the Mapping tabs to examine a 2-D axis or a 3-D axis solution, or on the STATISTICS tab to browse through the correspondence analysis statistics.
- Click the 划 button to close the dialog box and return to WordStat main window.

Step #9 - Displaying a keyword by keyword matrix or a keyword by case matrix

Click the WITH drop-down list and select <other keywords> to display a keyword by keyword matrix or on <case no> to view a keyword by case matrix.

Step #10 - Viewing a Keyword-In-Context (KWIC) list of specific words or categories

- Click the KEYWORD-IN-CONTEXT tab to access the KWIC table.
- Set the LIST option to Included and select the word or category for which you would like to obtain a KWIC table.
- Click the button to display the KWIC table for this word or category.

- To sort the table by case number, by keyword along with the prior or subsequent words, or by the sex of the respondent, use the SORT BY drop-down list.
- To display KWIC tables of any user-defined word or word pattern, set the LIST option to "User-defined", enter your word pattern (with or without wildcards) in the WORD edit box and click the >> button.

Step #11 - Editing a text from the KWIC list

• To modify the word or keyword or the text surrounding it, select it from the KWIC list, right-click and select the **EDIT** command. (You may also double-click the specific line you wish to edit).

Step #12 - Creating a concordance report

- Make sure the KEYWORD-IN-CONTEXT page is active and that the KWIC table displays the proper information.
- Set the amount of context that should be displayed around each keyword by setting the CONTEXT DELIMITER option.
- Click the button. Note: If this button is inactive, click the button to refresh the content of the KWIC table and then click the button.

Step #13 - Examining relationships between words or content categories using hierarchical clustering and multidimensional scaling

- Go to the Cooccurrence page.
- Click the DENDROGRAM tab to perform a hierarchical cluster analysis on included categories. You may change the number of partitions displayed using the No clusters option.
- Click the MAPPING tab to perform a multidimensional scaling, and display a plot in two or three dimensions. For more information see <u>Hierarchical Clustering and Multidimensional Scaling</u>

Step #14 - Saving or exporting frequency statistics to disk

- Move to the **Frequencies** tab.
- Press the 👒 button.
- Set the options and export to the existing data file or a new one.

Step #15 - Quitting the module and returning to QDA Miner or SimStat

• Click the ≡ button in the upper left-hand corner of WordStat and select the EXIT command or click the X mark in the upper right-hand corner.

Related Topics:

Creating, developing, and maintaining dictionaries

Performing analysis on manually entered codes

Automated Text Classification

Performing multivariate analysis (PCA, factor analysis, correspondence analysis, etc.) on words or categories.

Saving numeric and text results into a data file

Starting WordStat from QDA Miner or SimStat

WordStat can run as a stand-alone software or as a module run from either one of the following two base products:

QDA Miner: The text management and qualitative analysis program allows you to create and edit data files, import documents, and perform manual coding of the documents. Several analysis tools are also available to look at the frequency of manually assigned codes and the relationship between these codes and other categorical or numeric variables. When used with QDA Miner, WordStat can perform content analysis on whole documents or selected segments of the documents tagged with specific user defined codes.

SimStat: This statistical software provides a wide range of statistical procedures for the analysis of quantitative data. It offers advanced data file management tools such as the ability to merge data files, aggregate cases, perform complex computation of new variables and transformation of existing ones. When used with Simstat, WordStat can analyze textual information stored in any alphanumeric, plain text and rich text memo variable (or field). It includes various tools to explore the relationship between any numeric variable of a data file and the content of alphanumeric ones.

To run WordStat from QDA Miner:

• Start QDA Miner. You will be presented with a dialog box like this one:



- Click the **Open an existing project** and select the data file containing the document you want to analyze.
- If you closed or disabled this introductory dialog box, from the main screen select the **OPEN** command from the **PROJECT** menu and select this data file.
- Execute the CONTENT ANALYSIS command from the ANALYZE menu.
- Select in the Variables list box the variable that contains the documents you want to analyze and set the text to analyze to All text.
- To examine the relationship between the content of those documents and any categorical or numerical variable, in the **In relation with** group box, select the **Other variables** radio button and click the drop-down list to select those categorical or numerical variables.

Conter	nt Analysis							÷		×
Text to an	alyze:		_						_	
Variables:	[DOCUMENT]									~
	 All text 									
	O Coded segments:	19	- 0	ľ.						A 15
In relation	with:	only								
 Other 	r variables (nominal, num	neric, date, lo	gical)							
	[GENDER AGEGROUP]	-								~
() Assig	ned codes or category									
	Include category fo	r remaining t	ext							
	Additional variables:	1								-
					-					_
				🗸 Ok		Cance	el .			

• Press the Ok button. You will be taken to WordStat's Text Processing tab, where you can begin your analysis.

To run WordStat from SimStat:

- From within SimStat, select the FILE | DATA | OPEN command sequence and select the Seeking.ppj file
- Execute the STATISTICS | CHOOSE X-Y command
- Set the Variable list box to ALL to view all variable types.
- Move the GENDER and AGEGROUP variables to the Independent list box.
- Move the AD_TEXT variable to the Dependent list box.
- Press the Ok button
- Execute the **STATISTICS | CONTENT ANALYSIS** command. You will be taken to WordStat's **Text Processing** tab, where you can begin your analysis.

Creating a Project in WordStat

WordStat can be run as a stand-alone software. You can now create your projects directly from within WordStat.

Please see the following sections for more information on:

Creating a Project from a List of Documents Importing from a Data File Importing from Web Surveys Importing from Social Media Importing from Email Servers Importing from News Aggregators Importing Bibliographic References Using the Document Conversion Wizard

Creating a Project from a List of Documents

The easiest method to create a new project and start performing analysis in WordStat is by specifying a list of existing documents and importing them into a new project. This method creates a simple project with two-to-four variables: a categorical variable containing the original name of the files from which the data originated, a categorical variable containing the location of the files from which the data originated (optional), a variable containing the date and time that the files were created (optional), and a **DOCUMENT** variable containing imported documents. All text files are stored in different cases, therefore, if 10 files have been imported, the project will have 10 cases with two-to-four variables each. To split long documents into several documents or to extract numerical, categorical, or textual information from documents and store it in additional variables, use the <u>Document Conversion Wizard</u>.

To create a new project from a list of documents:

• Select the **NEW** command. A dialog similar to the one below will appear.



• Click the Create a project from a list of documents button. A dialog box similar to the one below will appear.

File type: Any documents (*.txt;*.rtf;*.doc;	".docx; ".wpd; ".odt; ".htm; ".html; ".pdf;"	*.epub; *.p	pt;*.pptx;*.) ∨ 🖽 Ŧ	
My Data Sources My Provalis Research Projects Dictionaries Exhibits Models Samples Workshop Candidates Election 2008 GOP GOP Debates 2011 GOP Debates 2012 GOP Debates 2012 GOP Debates 2012	Name ABC - December 11, 2011.H FOX - August 11, 2011.pdf FOX - December 15, 2011.pdf FOX - May 5, 2011.pdf FOX - Sept 22, 2011.pdf MSNBC - November 9, 2011.txt MSNBC - Sept 7, 2011.txt	Size 119KB 148KB 132KB 67KB 127KB 102KB 96KB	Type HTML Document Adobe Acrobat Docu Adobe Acrobat Docu Adobe Acrobat Docu Adobe Acrobat Docu Text Document Text Document	Modified 01/19/2013 3:10 PM 01/19/2013 3:24 PM 01/19/2013 3:24 PM 01/19/2013 3:21 PM 01/19/2013 3:22 PM 01/19/2013 3:25 PM
Import Import	Projects Workshop \GOP \GOP Debates 20	Add 11/Bloombe 11/CBS - No 11/CNN - Jo 11/CNN - No 11/CNN - S 11/CNN - S	Remove erg - October 11, 2011.do ovember 12, 2011.docx une 13, 2011.RTF lovember 22, 2011.RTF loctober 13, 2011.RTF ept 12, 2011.RTF	c.

- In the upper right section of the dialog box, WordStat displays all supported document file formats that may be imported, such as MS Word, WordPerfect, RTF, PDF documents, plain text files or HTML etc. To display only files of a specific type, set the **File type** list box to the desired file format.
- Click the file you would like to import. To select multiple files, hold down the Ctrl key while clicking on the other files.
- Click the Add button to add the selected files to the list of files to import, located at the bottom of the dialog box. You may also drag the files from the top right section to this list.
- Click the button to add numerous files contained within a folder. A dialog box will appear giving you the option of including **subfolders** and choosing **file types** to include in the import.
- To remove a file from the list of files to import, select that file name and click the
- Once all files have been selected, click the Create button. A dialog box similar to the one below will appear.

mportation	Document Importation	Лок
Import file location	Remove Images	
☑ Import creation date & time	Remove text formatting	X Cancel
	Fast PDF Importation	

You have the options to **Import file location** and **Import creation date & time**. If you select these options, this information will become variables entitled **LOCATION** and **CREATED**.

At this point your importation options differ depending on the documents you are importing. WordStat will only make available the importation options of the file type(s) that you have chosen. Other options will be grayed out.

Document Importation Options:

button.

The **Remove Images** option has been set by default. When importing formatted documents like Word, HTML or PDF files, images stored in the document may significantly increase the resulting document size and slow down the browsing and text-processing speed of WordStat. To keep the images in the imported documents, disable this option.

Check the **Remove Text Formatting** checkbox to further reduce the size of the imported documents if the text formatting of the existing documents such as the font styles and colors or the paragraph formatting are not relevant. Enable this option to convert all documents into plain-text documents without any formatting or images.

The fast **PDF Importation option** removes both images and text formatting from PDF files and is considerably faster that choosing both **Remove Images** and **Remove text formatting** options at the same time.

- Set the importation options.
- Click the **OK** button.
- Choose a name for your project and save it in the appropriate location. The **Data tab** will appear, containing a table with your imported data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Related Topics:

For more information on importing PDFs, see Notes on Importing PDF Files.

To configure the data set and start analyzing please see The Data Tab.

Creating a WordStat Project from Windows Explorer

The easiest way to create a project from multiple documents or a folder containing documents is to run WordStat through Windows Explorer. This method allows you to create a project without going through the importation process in QDA Miner or WordStat.

To access WordStat from Windows Explorer:

- Select the folder in Windows Explorer containing the subfolders and documents you wish to analyze.
- Right click your mouse, a menu opens, scroll down to WORDSTAT CONTENT ANALYSIS and choose RUN WORDSTAT from the adjacent menu. A dialog box will appear giving you the option of including subfolders and choosing file types to include in the import.

r: C:\Users\Amanda\Documents\My Prov	alis Research Projects\Workshop\GOP	
Include subfolders		
File types		
Text file (*.txt)	Rich Text files (*.rtf)	Acrobat PDF files
MS Word files (*.doc;*.docx)	HTML files	Powerpoint (*.ppt;*.pptx)
Ebooks (*.epub)	OpenOffice (*.odt)	WordPerfect (*.wpd)
XML Paper Specification (*.xps)		

• Once all files types have been selected, click the **OK** button. A Project configuration dialog box opens.

Rather than selecting a folder you can select the individual documents you would like to explore, using the **Shift** or **Ctrl** key to select numerous documents at the same time. Right click your mouse, a menu opens, scroll down to **WORDSTAT CONTENT ANALYSIS** and choose **RUN WORDSTAT** from the adjacent menu.

Project configuration		ć		×
Select the default language:	English	Ŷ		
Number of topics desired:	20		1	ОК

- Select the default language from the drop down menu.
- Enter the **Number of topics desired** for topic extraction by typing the number in the field or clicking on the up or down arrow beside the field until the desired number is reached. A temporary project opens containing three tabs: **Frequencies**, **Phrases** and **Topics**.



- Click the ≡ button and scroll down the MODE. Select EXPERT from the adjacent menu or click the Switch To Export Mode button. You will now have access to seven tabs and WordStat's full capabilities.
- Select the X at the top right of the screen to close the window. A dialog opens.
- You are asked if you would like to save the categorization model. If you would, select Yes, if not, select No.
- You are asked if you you want to save the data into a new project file. Select Yes. A Save As dialog opens.
- Choose a name for your project file and a place to save it. Select Save.

Notes on Importing PDF Files

The PDF file format is designed to display text in the same way on various platforms. Different strategies have been considered in QDA Miner and WordStat to support PDF files. They could simply be imported as is and stored in the QDA Miner project, allowing you to view the documents exactly as they appear in Acrobat. You could then code the PDF documents directly. Such an approach has, however, several inconveniences. First, the document may not be edited, preventing you from removing unnecessary information, such as headers and footers, appearing on each page or unrelated information like advertising. Also, PDF files are often created by outputting each page and each line of text separately, often resulting in the loss of the original document structure. Carriage returns will often be inserted at the end of each line, paragraphs spread over two pages remain physically split up, as do hyphenated words. In multi-column documents, lines from different columns may even be mixed up. All these minor imperfections may prevent you from retrieving full sentences and paragraphs in QDA Miner and may reduce recall results when searching for specific words or phrases separated by hyphens, carriage returns or page limitations. Such issues may also undermine the performance of numerous features of WordStat relying on phrase identification, word cooccurrence analysis or proximity rules.

Therefore, we chose to import PDF documents by converting them into editable documents, storing them in rich-text format, just like any other document file type supported by QDA Miner and WordStat. The conversion engine has been carefully designed to remove hyphens and unnecessary carriage returns, adjust the text flow in multi-column documents, and import tables and images correctly. While headers and footers will still be imported, which will break the flow of text across pages, you may now easily identify those and remove them from the text since the imported document will be fully editable.

Despite all the advanced importation features, some PDF documents may still not be imported properly. Several factors may explain such difficulties. Some PDF files consist of only scanned images of the original document and contain no text at all. These can be easily recognized by the fact that it is not possible to select any text segment or that text searches never return any hits. Other documents may have quite complex layout designs, making their proper importation quite difficult. For those documents, we recommend pre-processing them with full-fledged OCR (Optical Character Recognition) software like Abbyy's <u>FineReader</u> or Nuance's <u>OmniPage</u> or by some PDF to Word conversion utility tools like <u>PDF Transformer</u> by Abbyy. Although the latest version of Acrobat Professional does include some OCR features, its performance was deemed not good enough to preserve the document structure.

We also recommend removing images or scaling down the graphic resolution of images to the lowest setting because images will significantly increase the resulting document size and will slow down the browsing and text-processing speed of QDA Miner and WordStat.

Creating a Project from an Existing Data File or Web Service

The program can read data stored in the following file formats:

- MS Access
- MS Excel
- Comma Separated Values
- Tab Separated Values
- SPSS
- Stata v8 15
- Triple-S 1.2 and 2.0 XML file (an XML standard to exchange survey data)
- RIS files
- MBox and EML email file formats
- MS Outlook accounts and PST data files
- Reference Information System (RIS) data files
- ODBC
- EndNote

It can also import data from various Internet database and social media services such as:

- SurveyMonkey
- Qualtrics
- QuestionPro
- SurveyGizmo
- Voxco
- EndNote
- Zotero
- Mendeley
- Reference Information System (RIS)
- Twitter
- RSS Feeds
- Lexis UNI (from LexisNexis)
- Factiva
- Facebook
- Reddit
- Gmail
- Hotmail
- Outlook

Importing from a Data File

WordStat allows you to directly import data files from spreadsheets and database applications, as well as from plain ASCII data files (comma or tab delimited text).

To import data from a data file:

• Select the New button. A dialog box similar to the one below will appear.



- Click the Import from an existing data files or web service button.
- Select the desired file format from the menu. An import dialog box will appear.
- Select the corresponding file format from the Files of type drop-down list and select the file you want to import.
- Click Open.
- Name your project and save it in the appropriate location.

We will now discuss the specific file formats in greater detail, including formatting requirements (if applicable), supported features and limitations.

Excel Data Files

In an Excel spreadsheet you can enter both numeric and alphanumeric data in the cells of a data grid. WordStat can import spreadsheet files created by MS Excel *.xls and *.xlsx files.

Formatting Spreadsheet Data

The selected range must be formatted so that the columns of the spreadsheet represent variables (or fields), while the rows represent cases. As well, the first row should preferably contain the variable names, while the remaining rows should hold the data, one case per row. WordStat will automatically determine the most appropriate format based on the data that it finds in the worksheet columns. Cells in the first row of the selected range are treated as field names. If no variable name is encountered, WordStat will automatically provide one for each column in the defined range.

To create a new project from an Excel data file:

• Once you have named your project and saved it, a dialog similar to the one below appears.

Impor	t		-	
Sheet:	survey		~	V Import
Range:	• All	ORange		X Cance
				Preview >>

When you select an Excel file format for importation, the program displays a dialog box in which you can specify the spreadsheet tab and the range of cells where the data are located. You must specify a valid range name or provide upper left and lower right cells, separated by two periods (such as A1..H20). If you set the **Range** radio button to **All**, the program attempts to read the whole tab. You can preview the data before you import by selecting the **Preview** button. Your Excel spreadsheet must be closed before you select **Import**. You can only import one tab at a time. If you have more than one tab to import, containing the same column headers, simply append the other tabs.

- In the Sheet drop-down choose the tab you would like to import.
- Select either a Range or All and set you range, if necessary.
- Select Import. An Import Options dialog appears similar to the one below.

			R	OW IDE	NTIFICATION	N - STEP 1 OF 2		
loes yo	ur workshe	et contain:						
P	Variable 1	Names: Start on row: 1	4					
E] Variable [Descriptions: Start on row 2	(÷					
	Data:	Starts on row: 2	\$					
review								
1	ID	WHEDE	GENDER		ETHNICITY	108	SLITAS	-
-	6	Los Appeles CA	Male	25.30	caucacian	antertainment industry	None that I am aware of LOI	
3	8	Scotland	Male	19-24	Scottish	Student	Leading huge PvP groups and	
4	10	Scobard	Male	25-39	black	autoworker	n/a	
5	11	1154	Male	14-18	American	Student	It helps me to think things	
6	12	Moose Jaw, Canada	Male	25-39	Irish	Accountant	Tokes and COH humor cross	
7	13	Birmingham, Alabama-USA	Male	19-24	Caucasian	College Student/Restaurant	I've learned a little more	
8	15	lakenheath, england	Male	25-39	caucasion	retail	typing skills if any	
9	17	Bremerton, Washington USA	Male	40-54	American	Programmer	none, games do not effect	
10	18	Amherst, Massachusetts,	Male	25-39	Caucasian	Molecular biologist	At most, the game has taken	
11	19	Arkansas	Male	25-39	Black	Govt	Humor. Im a real serious	
12	24	Sacramento, CA.	Male	25-39	Caucasion	State Worker	I really don\'t think much I\'ve	
	-	11 7 11 5 1		05.00	1.22		4	

This dialog is a two-step process that helps you properly configure your data for import. On the first page you are asked to identify the location of the variable names and variable descriptions, if present. You are also asked to identify the row on which the data starts.

- If your spreadsheet contains variable names, check the checkbox called **Variable Names** and enter the row number in which they are located.
- If your spreadsheet contains variable descriptions, check the checkbox called **Variable Descriptions** and enter the row number in which they are located.
- In the **Data** field enter the row number on which the data starts.
- Select **Next**. The second page appears.

lect the vi typing in t	ariables that you would I he associated cell. Char	ike to import. You can customize the na nge the variable type by selecting from	ame and description if necessary the dropdown menu.		
bles					
lected	Name	Description	Туре		
	ID	ID	Integer	*	
$\mathbf{\nabla}$	WHERE	WHERE	Short String	*	
\square	GENDER	GENDER	Nominal	~	
\leq	AGEGROUP	AGEGROUP	Nominal	×	
\checkmark	ETHNICITY	ETHNICITY	Document	×	
\square	JOB	JOB	Short String	*	
V	SKILLS	SKILLS	Document	2	

The second page allows you to choose the variables you would like to import. You can change the names and descriptions by typing directly in the cells. To change the data type select from the drop-down list.

- Select the variables you want to import by checking their checkboxes.
- Modify the variable names, descriptions and data types as necessary.
- Select **Import**. The data will be imported and WordStat's **Data** tab will appear containing a table with your imported Excel data. From here you can further tailor your data set if necessary, or start analyzing immediately.

MS Access Database Files

When you select a MS Access data file, the program displays a dialog box in which you can specify the table where the data are located. Once the table has been selected, click the **Import** button to import the data file.

ASCII Data Files

WordStat will read up to 2000 numeric and alphanumeric variables from a plain ASCII file (text file). The file must have the following format:

- Every line must end with a carriage return.
- The first line must include the variable names, separated by spaces, tabs and/or commas.
- Variable names may not be longer than ten characters. Longer strings are truncated at ten characters.
- The remaining lines must include alphanumeric or numeric scores separated by spaces, tabs and/or commas.
- Each line must contain data for one case; variables must be in the same order for all cases.
- All invalid scores and all blanks encountered between commas or tabs are treated as missing values. A single dot can also be used to represent a missing value.

- Comments can be inserted anywhere in the file by putting an asterisk ('*') at the beginning of the line.
- Blank lines can also be inserted anywhere in the file.

Importing from Web Surveys

WordStat allows you to import survey data directly from various web survey platforms, enabling you to analyze open-ended responses, using the text mining and content analysis features of WordStat. Importing open-ended and closed-ended questions permits you to quickly compare results by demographics and find out how they are related to ratings, or any other numeric, categorical or date variables. WordStat supports SurveyMonkey, Qualtrics, SurveyGizmo, Voxco, and QuestionPro survey platforms. It also supports TripleS v1.2 survey files. Importing data from the various survey platforms is generally done in a similar fashion. You will be asked to log in to your account or provide a token before importing survey data. QuestionPro, Qualtrics and Voxco survey platforms require tokens to import data into WordStat.

Creating a New Project Using Survey Data

• Select the New button on the Data tab. This command calls up a dialog box similar to the one below.



- Select the **Import from data files or web services** button. A menu will appear on the right-hand side of the dialog box.
- Scroll down to the SURVEY menu item and choose the survey platform from which you wish to import the data.

For instruction on how to import survey data from the various survey platforms, please see:

Importing from <u>SurveyMonkey</u> Importing from <u>Qualtrics</u> Importing from <u>SurveyGizmo</u> Importing from <u>Voxco</u> Importing from <u>QuestionPro</u> Importing from a TripleS v1.2 file

SurveyMonkey

Connecting to SurveyMonkey

- Select SURVEY | SURVEYMONKEY from the menu. You will be prompted to log in to your SurveyMonkey account.
- Log in to your account. A dialog box requesting authorization for WordStat to access your account will appear.



- Select the I'm not a robot checkbox and perform the requested task.
- Click Authorize.

Importing SurveyMonkey Survey Data

Once a connection has been established an Import a Survey dialog box will appear.

85073769 Préparez-vous pour 9/20/2016 8:36:00 9/21/2016 9:56:00 7 79224068 test 5/6/2016 4:14:00 9/13/2016 6:49:00 6 79165889 Customer Survey 5/5/2016 4:04:00 6/20/2016 2:32:00 116		Name	Date Created	Date Modified	Status	Responses Count
79224068 test 5/6/2016 4:14:00 9/13/2016 6:49:00 6 79165889 Customer Survey 5/5/2016 4:04:00 6/20/2016 2:32:00 116	73769	Préparez-vous pour	9/20/2016 8:36:00	9/21/2016 9:56:00		7
79165889 Customer Survey 5/5/2016 4:04:00 6/20/2016 2:32:00 116	24068 t	test	5/6/2016 4:14:00	9/13/2016 6:49:00		6
	65889	Customer Survey	5/5/2016 4:04:00	6/20/2016 2:32:00		116

The dialog box lists the ID, Name, Date Created, Date Modified, Status and the Response Count of each survey on the platform.

• Choose the survey that you wish to analyze. You have two options. You can import all questions from the survey or you can select specific survey questions that are relevant for your analysis.

To import all survey questions:

- Select the Import All radio button at the bottom of the dialog box.
- Click OK.

• Choose a name for your project and save it in the appropriate location. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported survey data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To import only specific questions:

- Select the Select Questions radio button at the bottom of the dialog box.
- Click OK. A dialog box like the one below will appear.

Variable ID	Selected	Name	Description	Туре
Email		Email	Respondent Email	String
ID		ID	Respondent ID	String
IP_Address		IP_Address	Respondent IP Address	String
Date_Start		Date_Start	Date Started	DateTime
Date_Mod		Date_Mod	Date Modified	DateTime
Duration		Duration	Duration in seconds	Numeric
1017157268	\checkmark	Q1	Seriez-vous des nôtres?	Nominal/Ordinal
1017160556_1059		Q2_1	Parmi les dates proposées laquelle vous convient le plus? - Vendredi 9 décembre 2016	Nominal/Ordinal
1017160556_1059	\checkmark	Q2_2	Parmi les dates proposées laquelle vous convient le plus? - Samedi 10 décembre 2016	Nominal/Ordinal
1017160556_1059	\checkmark	Q2_3	Parmi les dates proposées laquelle vous convient le plus? - Vendredi 16 décembre 2016	Nominal/Ordinal
1017163634	\checkmark	Q3	commentaires	String
1017255475		Q4	Votre nom?	String

- The Variable IDs will become your project variables. Select the checkboxes of the variables you wish to import.
- Click OK.
- Choose a name for your project and save it in the appropriate location. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported survey data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Related Topics:

To configure the data set and start analyzing please see The Data Tab

To append new survey responses to your project please Monitoring Online Resources.

Qualtrics

Connecting to Qualtrics

You must enter a token to link with your Qualtrics survey. To get your Qualtrics token:

- Log in to your Qualtrics account.
- Go to the Account Settings menu item in the User drop-down menu.
- Go to the Qualtrics IDs.
- Click Generate to receive a token.

1. Login to Qualtri 2. Go to Account S	ics Settings in the use	er dropdow	٧n			
My Projects			Project	be Contacts	Library Admin	Help & Feedback ()
Folders				+ Create	Project Q	Contacts
					-	Account Setting
All Projects	Ali Proje	ects				Refresh Account
Shared with Me	0					Logout
	Tach MoDe	sevent!				
			User Dropdown			
3. Go to Qualtrics y Account) Uner Settings @ Dograde Ac	IDS.	Qualitics Ids		Projecta	Contacts Likery	y Admin Nely & Feedback
3. Go to Qualtrics (y Account U Use Setting: (a) Upgede Ao ualtrics IDs aurige	IDs count 🔄 Account Insign	🗘 Qualitatice i das	1944	Projecta	Contacts Library	r Admin Heig & Fens®hack
3. Go to Qualtrics ty Account () Une Setting: (* Logisch Ac- ualtrics IDa Service Nore	IDS.	C Guaittice ids	Une UR-eCOLES/sept04-0	Projecta	Confacts Library	y Admin Tindy & Tendbuck Admin Tindy & Tendbuck Administrative State (State State St
3. Go to Qualtrics (y Account () the Settings (* Logisch Ac ualtrics IDs Aways New	IDS	Qualityice ids User Id Department id	Ver UR «ColgityupiteD averparget	Prosects	Contacts Library	r Admin Tindy & Tendkuck Adm molekuljovenulogannikevozotn B
3. Go to Qualtrics (y Account () the Settings (* logisch Ac ualtrics IDs Nine	IDS	Cualification Unarrig Deganization 18	Vere UR-ciOsylysystem exemplorogid	Projects	Confacts Library	r Admin Heig & Teadlack Am molekulgssine acceptive courts
3. Go to Qualitrics by Account Dues Serings (a) Lagrade Ao ualtrics IDs Note Duese to Carter Anti-		Qualifyice ids User M Organization id	Jaw Uni «Cotylegaster sampinoga	Proot	Confacts Library EXAMPLETORIAL All Documentation Dataset	Admin Heig & Feedback Administry (Control of Control of Control International Control of Control of Control International Control of Control of Control International Control of Control of Control of Control International Control of Control of Control of Control International Control of
3, Go to Qualitrics y Account (a) two terring (a) tegrade Ao waltrice IDa water water water begins des facts that the begins des facts that that design begins des facts that design begins des facts that design begins des facts that design begins des facts that des facts that design begins des facts that des facts that des facts that des begins des begins des facts that des begins des begins des facts that des begins d		Ö Guartske ida User of Dige-Entise id	New UR 4/COstrugateD exemptionged	Proofs	Contarits Library	Adam Initia Eritadeus An An An An An An An An An An An An An
3. Go to Qualitrics V Account V Account U ourse terring U ours		Constitute i dar Lines of Digartetion II	Stee UK-cChully-yeldb skorplangel	Proposts Totato Uriti	Contaction Library EXMIPLETONIAN INFORMATION Design	Admini Hidg & Pawikani An Antini Luguna Roger Verw22111
3. Go to Qualitrics V Account Une terring (a) langed Ac unatrices Ds unatrices Ds unatrices Cost encoded account of the cost o	IDs	Constitute da	law Un «Colydrycetho exemptorget	Francis Total US	Contacts - Livery Local/PLTICeCP2 3/FLDecanation Genes	Adres Initig & Facebook Adres Same Science (Same Science (
3. Go to Qualitrics y Account une terring @ logical A unitries Da merce terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring terring		C Gualitice lets	Var Un «Cotophysetheo averpisorgal	Prosect	Contacts Livery	Adren Inico & Freedowi Adren Edit
3. Go to Qualitrics V Anoount U Anoount U Une Burrys @ Upyth An Upyth Anoount U Une Burrys @ Upyth An Upyth Anoount U Une Burrys @ Upyth An Upyth Anoount Anoount U Une Burrys @ Upyth Anoount	IDS	© Gualities ids User of Digardeton IS	Juni Uni «Cotoposatero sempioropa Qualitrics IDs ur token yet	Total Total	Contracts Library	Admini Help & Financianos Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas Administrativas A

- In WordStat, select **SURVEY | QUALTRICS** from the menu.
- Enter your token in the Connect to Qualtrics dialog box.

Connect to Qualtrics	-	
Token		
	d ok	X Cancel

• Click OK.

Importing Qualtrics Survey Data

Once a connection has been established an Import a Survey dialog box will appear.

		Duce cicuccu	Date Modified	Status	Responses Count
SV_cC3JDz2J	provalisTest		6/17/2016 5:53:40	Active	

The dialog box lists the ID, Name, Date Created, Date Modified, Status and the Response Count of each survey on the platform.

• Choose the survey that you wish to analyze. You have two options. You can import all questions from the survey or you can select specific survey questions that are relevant for your analysis.

To import all survey questions:

- Select the Import All radio button at the bottom of the dialog box.
- Click OK.
- Choose a name for your project and save it in the appropriate location. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported survey data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To import only specific questions:

- Select the Select Questions radio button at the bottom of the dialog box.
- Click **OK**. A dialog box like the one below will appear.

/ariable ID	Selected	Name	Description	Туре	^
D		ID	Respondent ID	String	
P_Address		IP_Address	Respondent IP Address	String	
Date_Start		Date_Start	Date Started	DateTime	
ate_Ended		Date_Ended	Date Ended	DateTime	
Duration		Duration	Duration in seconds	Numeric	
atitude		Latitude	Respondent Geolocation - Latitude	Decimal Number	
ongitude		Longitude	Respondent Geolocation - Longitude	Decimal Number	
22		Q2	Aimez-vous programmer?	Numeric	
24		Q4	Êtes-vous d'accord avec cette affirmation: L'hiver, au Québec, on gèle!	Numeric	
25		Q5	Pensez-vous qu'il serait justifié de décerner le prix du CEO of the month au Canada à votre	Numeric	
Q6		Q6	Décrivez-nous vos obstades quotidiens et expliquez-nous comment vous êtes parvenus à	String	
28		Q8	À titre de statistique seulement, veuillez entrer votre numéro d'assurance sociale. Aucune	String	
Q9_1		Q9_1	Évaluer la vitesse des logiciels suivants:-Nvivo	Numeric	
29 2		Q9 2	Évaluer la vitesse des logiciels suivants:-WordStat	Numeric	

- The Variable IDs will become your project variables. Select the checkboxes of the variables you wish to import.
- Click OK.
- Choose a name for your project and save it in the appropriate location. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported survey data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Related Topics:

To configure the data set and start analyzing please see The Data Tab

To append new survey responses to your project please Monitoring Online Resources.

SurveyGizmo

Connecting to SurveyGizmo

- Select SURVEY | SURVEYGIZMO from the menu. You will be prompted to log in to your SurveyGizmo 3.0 account.
- Log into your account. A dialog box requesting authorization for WordStat to access your account will appear.



• Click Approve.

Importing SurveyGizmo Survey Data

Once a connection has been established Import a Survey dialog box will appear.

(D	Name	Date Created	Date Modified	Status	Responses Count
2725046	Api Integration Test	4/19/2016 2:57:07	4/21/2016 1:42:18	Launched	150

The dialog box lists the ID, Name, Date Created, Date Modified, Status and the Response Count of each survey on the platform.

• Choose the survey that you wish to analyze. You have two options. You can import all questions from the survey or you can select specific survey questions that are relevant for your analysis.

To import all survey questions:

- Select the Import All radio button at the bottom of the dialog box.
- Click OK.
- Choose a name for your project and save it in the appropriate location. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported survey data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To import only specific questions:

- Select the Select Questions radio button at the bottom of the dialog box.
- Click OK. A dialog box like the one below will appear.
| ariable ID | Selected | Name | Description | Туре | 1 |
|------------|----------|------------|-------------------------------------------------------------------------|-----------------|-----|
| D | | ID | Respondent ID | String | 6. |
| P_Address | | IP_Address | Respondent IP Address | String | |
| ate_Start | | Date_Start | Date Started | DateTime | |
| ate_Ended | | Date_Ended | Date Ended | DateTime | |
| atitude | | Latitude | Respondent Geolocation - Latitude | Decimal Number | |
| ongitude | | Longitude | Respondent Geolocation - Longitude | Decimal Number | |
| | | Q5 | Aimes-tu la vie comme moi? | Nominal/Ordinal | |
| | | Q7 | Comment aimes-tu ça?-Poulet | Nominal/Ordinal | |
| | | Q8 | Comment aimes-tu ça?-Carotte | Nominal/Ordinal | |
| | | Q9 | Comment aimes-tu ça?-Olive | Nominal/Ordinal | |
| 0 | | Q10 | Nombre de dents dans ta bouche. | Nominal/Ordinal | |
| 1_10009 | | Q11_10009 | Cochonnerie que tu possèdes-Cellulaire | Nominal/Ordinal | |
| 2 | | Q12 | Quel est le nom de votre animal de compagnie? | String | |
| 3 | | Q13 | Selon vous, quelle est la relation entre le neutrino et l'antineutrino? | Open-ended | . 5 |

• The Variable IDs will become your project variables. Select the checkboxes of the variables you wish to import.

• Click OK.

• Choose a name for your project and save it in the appropriate location. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported survey data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Related Topics:

To configure the data set and start analyzing please see The Data Tab

To append new survey responses to your project please Monitoring Online Resources.

Voxco

Connecting to Voxco

You must enter a token to link to your Voxco survey. To get your Voxco token:

- Log in to your account on Voxco's A4s platform.
- Select User Settings from the User drop-down menu.
- On the User Settings page you will find the API Access Key in the list of Security options. Copy this token.
- In WordStat, select SURVEY | VOXCO from the menu.
- Enter your token into the Connect to Voxco dialog box.

Connect	to Voxco	-		×
Token	Qcawhbfcewuhf;qhfqpNVINRWEG	1)		
		🗸 ок	×	Cancel

• Click OK.

Importing Voxco Survey Data

Once a connection has been established Import a Survey dialog box will appear.

ID	Name	Date Created	Date Modified	Status	Responses Count
2725046	Api Integration Test	4/19/2016 2:57:07	4/21/2016 1:42:18	Launched	150

The dialog box lists the ID, Name, Date Created, Date Modified, Status and the Response Count of each survey on the platform.

• Choose the survey that you wish to analyze. You have two options. You can import all questions from the survey or you can select specific survey questions that are relevant for your analysis.

To import all survey questions:

- Select the Import All radio button at the bottom of the dialog box.
- Click OK.
- Choose a name for your project and save it in the appropriate location. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported survey data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To import only specific questions:

- Select the Select Questions radio button at the bottom of the dialog box.
- Click **OK**. A dialog box like the one below will appear.

ariable ID	Selected	Name	Description	Туре	-
D		ID	Respondent ID	String	6.
P_Address		IP_Address	Respondent IP Address	String	
ate_Start		Date_Start	Date Started	DateTime	
ate_Ended		Date_Ended	Date Ended	DateTime	
atitude		Latitude	Respondent Geolocation - Latitude	Decimal Number	
ongitude		Longitude	Respondent Geolocation - Longitude	Decimal Number	
		Q5	Aimes-tu la vie comme moi?	Nominal/Ordinal	
		Q7	Comment aimes-tu ça?-Poulet	Nominal/Ordinal	
		Q8	Comment aimes-tu ça?-Carotte	Nominal/Ordinal	
		Q9	Comment aimes-tu ça?-Olive	Nominal/Ordinal	
0		Q10	Nombre de dents dans ta bouche.	Nominal/Ordinal	
1_10009		Q11_10009	Cochonnerie que tu possèdes-Cellulaire	Nominal/Ordinal	
2		Q12	Quel est le nom de votre animal de compagnie?	String	
3		Q13	Selon vous, quelle est la relation entre le neutrino et l'antineutrino?	Open-ended	. 5

- The Variable IDs will become your project variables. Select the checkboxes of the variables you wish to import.
- Click OK.
- Choose a name for your project and save it in the appropriate location. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported survey data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To configure the data set and start analyzing please see The Data Tab

To append new survey responses to your project please Monitoring Online Resources.

QuestionPro

Connecting to QuestionPro

You must enter a token to link to your QuestionPro survey. To get your QuestionPro token:

- Log in to your QuestionPro account.
- Select Surveys from the navigation bar at the top of the screen.
- Select the survey you wish to import.
- Click on Developer API in the Integration drop-down menu. The API Access Key is your token.

2 QuestionPro	Policy Communities involution Ct.	
QuestionPro · We	hsite Satisfaction	
Question to . Ne	bare addatactor	
tai multi sunt apoc	The second	
XML API - Real-Time Post Q		
III WITTP / MAR Pysik		
Real-Time Dashboard O		
C Active Ind - pril Publick		
REST//SON API	NETATIVE FOR STATE STATE	
= REST APt Console		
downable API Calls	Http://www.marticlemarg.com/law	
4241 1.005	mp.//www.serifers.monial.edu.pertMantsurvey.arXXV.arvey7a.mov	
Request (SOV		
		8

- In WordStat, select SURVEY | QUESTIONPRO from the menu.
- Copy and paste this key into the Connect to QuestionPro dialog box.



• Click OK.

Importing QuestionPro Survey Data

Once a connection has been established Import a Survey dialog box will appear.

10	Name	Date Created	Date Modified	Status	Responses Count
4478558	ProvalisTest -				
4474485	ProvalisTest		1.2		
4469924	QuestionPro :			1	

The dialog box lists the ID, Name, Date Created, Date Modified, Status and the Response Count of each survey on the platform.

• Choose the survey that you wish to analyze. You have two options. You can import all questions from the survey or you can select specific survey questions that are relevant for your analysis.

To import all survey questions:

- Select the Import All radio button at the bottom of the dialog box.
- Click OK.
- Choose a name for your project and save it in the appropriate location. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported survey data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To import only specific questions:

- Select the Select Questions radio button at the bottom of the dialog box.
- Click **OK**. A dialog box like the one below will appear.

Variable ID	Selected	Name	Description	Туре	^
ID		ID	Respondent ID	String	
IP_Address		IP_Address	Respondent IP Address	String	
Date_Start		Date_Start	Date Started	DateTime	
Duration		Duration	Duration in seconds	Numeric	
Q1		Q1	How often do you conduct surveys?	Numeric	
Q1_234687113		Q1_Other	How often do you conduct surveys? - Other	String	
Q2		Q2	What types of credit cards do you have (Select all that apply)?	Numeric	
Q2_234687120		Q2_Other	What types of credit cards do you have (Select all that apply)? - Other	String	
Q2_234687121		Q2_Other	What types of credit cards do you have (Select all that apply)? - Other	String	
Q3		Q3	How often do you conduct surveys?	Numeric	
Q3_234687122		Q3_Other	How often do you conduct surveys? - Other	String	
Q4		Q4	Comments/Suggestions:	Numeric	
Q4_234687126		Q4_Other	Comments/Suggestions: - Other	String	
Q5		Q5	Name	Numeric	
Q5_234687127		Q5_Other	Name - Other	String	
Q6		Q6	How many months have you been at your current residence?	Numeric	
Q6 234687128		Q6 Other	How many months have you been at your current residence? - Other	String	~

- The Variable IDs will become your project variables. Select the checkboxes of the variables you wish to import.
- Click OK and the selected questions will be imported into your project.
- Choose a name for your project and save it in the appropriate location. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported survey data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Related Topics:

To configure the data set and start analyzing please see The Data Tab

To append new survey responses to your project please Monitoring Online Resources.

TripleS v.1.2

Importing TripleS v1.2 files

• Select SURVEY | TRIPLES V1.2 from the menu. An Import data file dialog box will appear.

			_	
Look in	TripleS		🖌 🕝 🤌 📂 🖽 🗸	
Quick access	Name Speech2.s test_sss.ss	55 5	Date modified 2/12/2016 10:01 AI 9/16/2016 4:26 PM	Type M SSS File I SSS File
Desktop				
Libraries				
This PC	<			
Natural	File name:		~	Open
Network	Files of type:	Triple-S XML files (*.sss)	Ŷ	Cancel
		Open as read-only		

- · Choose the triple-s file that you would like to import.
- Click Open.
- Choose a name for your project and save it in the appropriate location. The data will be imported and WordStat's Data tab will appear, containing a table with your imported survey data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Related Topics:

To configure the data set and start analyzing please see The Data Tab

Importing from Social Media

Social media is a popular networking and marketing medium. Many companies, organizations, public figures etc. have a social media presence. WordStat allows you to import social media data directly from various social media platforms, allowing you to use the text mining and content analysis features of WordStat to analyze and monitor social media data over time. WordStat supports RSS, Twitter, Facebook pages and Reddit. You must have access to a Twitter account to use the Twitter importation feature.

Creating a New Project Using Social Media Data

• Select the New button on the Data tab. This command calls up a dialog box similar to the one below.

Create a project from a list of docume	ents
Timport from datafiles or web services	
Run Document Conversion Wizard	

- Select the **Import from data files or web services** button. A menu will appear on the right-hand side of the dialog box.
- Scroll down to the **SOCIAL MEDIA** menu item and choose the social media platform from which you wish to import the data.

For instruction on how to import survey data from the various survey platforms, please see:

Importing from RSS Importing from Twitter Importing from Facebook pages

Importing from RSS

WordStat allows you to create projects by importing RSS feeds and use the text mining and content analysis features of WordStat to analyze and monitor your RSS sources over time.

Creating a New Project Using a RSS Feed

• Select SOCIAL MEDIA | RSS FEED from the menu. A Define Query dialog box will appear.

History	Variables
Source: RSS V	Valuaties
Max Posts: 0 (0 = unlimited)	

- Select RSS as the Source of your data import.
- Copy the URL of the RSS feed of your choice into the Expression field.

- The Max Posts field will be set to 0 by default, which will give you an unlimited number of results. You can limit the amount of RSS data collected by changing the number in the Max Posts field.
- The **Source**, **Title**, **Link Description** and **PubDate** variables in the **Variables** list box are grayed out as they are imported into the project by default. Check the **Article** checkbox from the **Variables** list box if you would like to import the full text of the article in a document variable.
- The Live Monitoring option instructs WordStat to monitor the RSS feed over time and collect for import RSS posts published after the initial data capture. Enabling this option calls the Web Collector, an external tool, which collects your RSS data. You can access the Web Collector through the system tray and the Windows start menu. When you restart your computer the Web Collector will start automatically and continue collecting data if the query status is active.
- Set the set the **Time Delay** frequency and the **Stop Date**. The Web Collector will collect new RSS posts published according to the selected time interval until the stop date is reached. Data will be collected as long as the Web Collector is running. As the Web Collector is a completely separate application WordStat does not have to be running for it to work.
- Once all the options have been set, click OK.
- Enter the project name and save it in the location of your choosing. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported RSS data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To configure the data set and start analyzing please see The Data Tab.

To append RSS data to your project please Monitoring Online Resources.

For more information on the Web Collector please see the section entitled Web Collector.

Importing from Twitter

WordStat allows you to create a project by importing Twitter data and use the text mining and content analysis features of WordStat to analyze and monitor your Twitter sources over time. You must have access to a Twitter account to use this feature.

Creating a New Project Using Twitter Data

• Select SOCIAL MEDIA | TWITTER from the menu. A Define Query dialog box will appear.

Define/Edit Query	- 0
tings:	Variables
Source: Twitter	Date Created
Expression: Provalis	Number of Retweets
Max Posts: 0 (0 = unlimited)	Geocoordinate Latitude
Geocode Locations	User Name
	User URL User Number of Followers User Number of Friends
Time Delay: 1 Days	User Number of Listed
Stop Date: 11/ 4/2016	User Number of Tweets
	V OK X Can

- Select Twitter as the Source of your data import.
- Enter your Twitter query in the **Expression** field. You can enter a simple Boolean query directly into this field or you can click the *P* button to the right of the field and use the advanced search option to help you create your query. The **Build your expression** dialog box is an advanced search option. It mirrors the advanced search in Twitter and works on the same principles.
- The Max Posts field will be set to 0 by default, which will give you an unlimited number of results. You can limit the amount of Twitter data collected by changing the number in the Max Posts field.
- Check the corresponding boxes to Include Retweets and Geocode locations in the import if you desire. Including Retweets may result duplicate cases. The number of duplicates indicates the virality of the Tweet. Importing Geocode Locations will result in the addition of three variables: Longitude, Latitude and User Location. User Location is the user defined. Including these variables in the import allows you to use the GEOCODING function to obtain the geographic locations of the users and map them using WordStat's mapping feature.
- The **Date Created**, **Tweet ID** and **Tweet Text** variables in the **Variables** list box are grayed out as they are imported into the project by default. Check the variables that you would like to import from the **Variables** list box.
- The Live Monitoring option instructs WordStat to monitor the Twitter search over time and collect for import Tweets published after the initial data capture. Enabling this option calls the Web Collector, an external tool, which collects your Twitter data. You can access the Web Collector through the system tray and the Windows start menu. When you restart your computer the Web Collector will start automatically and continue collecting data if the query status is active.
- Set the **Time Delay** frequency and the **Stop Date**. The Web Collector will collect new Tweets published according to the selected time interval until the stop date is reached. Data will be collected as long as the Web Collector is running. As the Web Collector is a completely separate application WordStat does not have to be running for it to work.
- Click **OK.** If it is the first time you are creating a WordStat project from Twitter data, you will need to authorize the Provalis **Web Collector** to connect to your Twitter account. A dialog box like below will appear.

¥	Sign up for Twitter s
Authorize Provalis WebCollector to use your account?	PROVALIS
Usemame or email	Provalis WebCollector By Provalis Research provalisresearch.com/
Password Remember me - Forgot password?	Imports tweets directly into QDA Miner.
Authorize app Cancel	
This application will be able to:	
Read Tweets from your timeline.	
See who you follow.	
Will not be able to:	
 Follow new people. 	
Follow new people.Update your profile.	
Follow new people.Update your profile.Post Tweets for you.	
 Follow new people. Update your profile. Post Tweets for you. Access your direct messages. 	

- Enter your log in details and click Authorize app.
- Enter the project name and save it in the location of your choosing. The data will be imported and WordStat's Data tab will appear, containing a table with your imported Twitter data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Note: Twitter will automatically go back seven days in its database to retrieve information. The maximum number of Tweets accessible each 15 minutes by a single Twitter account is 18,000. If you have reached your limit, the Web Collector will stop and display a dialog box to inform you that you have reached your limit.

Related Topics:

To configure the data set and start analyzing please see The Data Tab.

To append Twitter data to your project please Monitoring Online Resources.

For more information on the Web Collector please see the section entitled Web Collector.

Importing from Facebook

WordStat allows you to create a project by importing Facebook page data and use the text mining and content analysis features of WordStat to analyze and monitor your Facebook page sources over time. WordStat cannot access data from password protected profiles, only from public pages.

Creating a New Project Using Facebook Page Data

• Select SOCIAL MEDIA | FACEBOOK from the menu. A Define Query dialog box will appear.

Define/Edit Qu	Jery	
ettings:	/Edit Query Source: Facebook ✓ ression: https://www.facebook.com/search/top/?q=provalis%20resea From: 9/ 1/2016 ■▼ Until: 9/30/2016 ■▼ Until: 9/30/2016 ■▼ Stop Date: 11/ 4/2016 ■▼	Variables
Source: Expression:	Facebookhttps://www.facebook.com/search/top/?g=provalis%20resea	 ☑ ID ☑ Username ☑ Post ☑ Comment ☑ Created Time ☑ Like Count
From:	9/ 1/2016	 ✓ Comment Count ✓ Is Liked by Page Owner ✓ Tags ✓ Page TR
Until:	9/30/2016	M Parent 1D
⊡ Live M Time Stop	Ionitoring Delay: 1 Days Date: 11/ 4/2016	

- Select Facebook as the Source of your data import.
- Enter the URL of the Facebook page you would like to analyze.
- Set the From and Until date options by clicking the calendar to the right of the fields and selecting start and end date. WordStat will retrieve posts that were posted during the chosen time span. This time span must not exceed 6 months. The last possible day of the date range is the present day. It is not possible to set this option to capture future posts. Comments associated with the posts will also be imported. Imported comments, however, may have been published after the chosen date range. It is possible to capture future comments by setting the Live Monitoring option.
- Check the variables that you would like to import from the Variables list box.
- The Live Monitoring option instructs WordStat to monitor the Facebook posts published during the chosen time period and collect for import any comments associated with the posts that have been published after the initial data capture. Enabling this option calls the Web Collector, an external tool, which collects your Facebook data. You can access the Web Collector through the system tray and the Windows start menu. When you restart your computer the Web Collector will start automatically and continue collecting data if the query status is active. It is important to note that the Web Collector will not collect any new posts, only new comments associated with the posts from the chosen date range.
- Set the set the **Time Delay** frequency and the **Stop Date**. The Web Collector will collect comments published according to the selected time interval until the stop date is reached. Data will be collected as long as the Web Collector is running. As the Web Collector is a completely separate application WordStat does not have to be running for it to work.
- Once all the options have been set, click **OK**.
- Enter the project name and save it in the location of your choosing. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported Facebook data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To configure the data set and start analyzing please see The Data Tab.

To append Facebook data to your project please Monitoring Online Resources.

For more information on the Web Collector please see the section entitled Web Collector.

Importing from Reddit

WordStat allows you to create projects by importing Reddit data and use the text mining and content analysis features of WordStat to analyze and monitor your Reddit sources over time.

Creating a New Project Using a Reddit

• Select SOCIAL MEDIA | REDDIT from the menu. A Define Query dialog box will appear.

Reddit	- 🗆 ×
Search Subreddits Posts	
Search Expression:	
Search for: Posts URL of a post Subreddits	
Search: Obama Sort: Relevance V From: All time	✓ Max hits: 100
Web Collector	
Retrieve Every: 1 🕞 Days 🗸 Stop Date: 🗹 5/27/2021	
	Search

- Select what you are searching for Posts, URL of a post or Subreddits.
- If you are looking for **Posts** or **Subreddit** type your query into the **Search for** field. If you have the **URL** of a Reddit that you would like to import type that in the **Search for** field.
- Choose the the sorting order of your imported data by selecting from the drop down in the **Sort** field. You can sort by **Relevance**, **Hot**, **Top**, **New** and **Comments**.
- Choose the start date of your import by selecting from the drop down in the From field.
- Set the Max hits field to determine the maximum number of posts for the project.
- The Live Monitoring option instructs WordStat to monitor the Reddit search over time and collect for import Reddit posts published after the initial data capture. Enabling this option calls the Web Collector, an external tool, which collects your Reddit data. You can access the Web Collector through the system tray and the Windows start menu. When you restart your computer the Web Collector will start automatically and continue collecting data if the query status is active.
- Set the time delay frequency in the **Retrieve Every** field and the **Stop Date**. The Web Collector will collect new Reddit posts published according to the selected time interval until the stop date is reached. Data will be collected as

long as the Web Collector is running. As the Web Collector is a completely separate application WordStat does not have to be running for it to work.

 Once all the options have been set, click Search. If it is the first time you are creating a WordStat project from Reddit data, you will need to authorize the Provalis Web Collector to connect to your Reddit account. A dialog box like below will appear.

Log in or sign up to connect your re	ddit account with QDA Miner .	
MUS.		
CREATE A NEW ACCOUNT	LOG IN	
choose a username	username	
password	password	
verify password	remember me reset password	
email	LOG IN	
🗌 remember me		
🗌 email me updates		

- Enter your log in details and click Authorize app.
- Enter the project name and save it in the location of your choosing. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported Reddit data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Related Topics:

To configure the data set and start analyzing please see The Data Tab.

To append Reddit data to your project please Monitoring Online Resources.

For more information on the Web Collector please see the section entitled Web Collector.

Importing Bibliographic References

Reference Management tools are used to store, organize and search bibliographic data, abstracts and complete papers. WordStat allows you to create projects by importing bibliographic data and use the text mining and content analysis features of WordStat to effectively analyze your bibliographic data. WordStat supports EndNote, Mendeley and Zotero. It also supports the RIS (*.ris) file format.

Creating a New Project Using Bibliographic Data

• Select the New button on the Data tab. This command calls up a dialog box similar to the one below.



- Select the Import from data files or web services button.
- Select the **REFERENCE MANAGERS** menu item and choose the desired reference management software from its submenu.

For instruction on how to import bibliographic reference data from various sources, please see:

Importing from <u>EndNote</u> Importing from <u>Mendeley</u> Importing from <u>Zotero</u> Importing a <u>RIS (*.ris)</u> file

EndNote

Creating a New Project Using EndNote

- Select REFERENCE MANAGERS | ENDNOTE from the menu and an Import data file dialog box will appear.
- Select the EndNote file that you wish to import.
- Click Open.
- Enter the project name, save it in the location of your choosing. A Select the variables dialog box will appear.

All References All Groups All Groups Articles_QDA_Wc Articles_pas_sur ProQuest_QDA_V Bons_article_QDA Allref Allref Article_journal_2(Article_journal_2(Article_journal_2(Article_journal_2(Articles_vrai_proc Articles_proquest Temporaries	✓ Id Reference Type Text Styles ✓ Author ✓ Year ✓ Title ✓ Pages Secondary Title Volume Number Number Of Volumes Secondary Author Place Published Publisher Subsidiary Author Edition Keywords Type Of Work ✓ Date ✓ Abstract	*
nportation options		
Import Documents	Tertiary Title	
On disk only	ISBN	
Remove Images	Custom 1	~

- On the right-hand side of the dialog box is the **Groups** list box. A group in EndNote basically separates your library into categories. Select the check-boxes of the groups you want to import from.
- You are given the option to **Import Documents**. Check the **Import Documents** check box if you would like to import documents.
- Documents in EndNote are saved on disk, URLs that access online documents are also saved within the files. You have the option of importing documents that are saved **On disk only** and ignoring the URLs that point to online documents.
- Check the **Remove Images** checkbox if you would like to remove the images contained in the documents you are importing.
- On the right side of the dialog box is a list box of variables you can select for importation. Check the boxes of the variables you wish to import.
- Click **OK**. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported bibliographic data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To configure the data set and start analyzing please see <u>The Data Tab</u>.

To append new bibliographic references to your project please Monitoring Online Resources.

Mendeley

Creating a New Project Using Mendeley

• Select **REFERENCE MANAGERS** | **MENDELEY** from the menu and a dialog box requesting access to your Mendeley account appears.

8	MENDELEY
QDA Miner update data	is requesting the ability to access and from your Mendeley account.
Email	
Password	
	Authorize

- Enter your Email and your Password.
- Click Authorize if you give permission for WordStat to access your account. An Import References dialog box will appear.

Options
Import Documents

- On the left side of the dialog box is a list box of variables you can select for importation. Check the boxes of the variables you wish to import.
- On the right-hand side of the dialog box you are given the option to **Import Documents.** Select the **Import Documents** checkbox if you would like to import documents.

- Click OK.
- Enter the project name and save it in the location of your choosing. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported bibliographic data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To configure the data set and start analyzing please see The Data Tab.

To append new bibliographic references to your project please Monitoring Online Resources.

Zotero

Creating a New Project Using Zotero

- Select **REFERENCE MANAGERS** | **ZOTERO** from the menu and a dialog box asking you to sign in to your Zotero account appears.
- Log in to your Zotero account. A dialog box notifying you that WordStat would like to access your account appears.



• Click Accept Defaults if you give permission for WordStat to access your account. An Import References dialog box will appear.

elect variables to import:	Options
 Title Publication Authors Date Abstract Date Added Date Modified Pages Volume Issue Tags Collection Text Collection Title Journal Abbreviation Language DOI ISSN Short Title URL Accessed Archive Location Library Catalog Call Number Rights Extra 	☐ Import Documents

- On the left side of the dialog box is a list box of variables you can select for importation. Check the boxes of the variables you wish to import.
- On the right-hand side of the dialog box you are given the option to **Import Documents**. Check the **Import Documents** checkbox if you would like to import documents.
- Click OK.
- Enter the project name and save it in the location of your choosing. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported bibliographic data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To configure the data set and start analyzing please see The Data Tab.

To append new bibliographic references to your project please Monitoring Online Resources.

RIS File

Creating a New Project from a .RIS File

Bibliographic data from Reference Management tools and digital libraries is often stored in the RIS (*.ris) file format. RIS files are plain-text files and have by default a TXT file extension. To import such a file, set the **Files of type** list box to Reference Information Management (RIS) and then select the file you wish to import. Another approach is to change the file extension to RIS, and WordStat will automatically recognize this file extension and import the references in this file without the need to specify the proper file type.

• Once a file has been selected for importation, WordStat will display an Import RIS Data File dialog box.

Select variables to import:	Importation options
Authors	Transform publication year into a date
Journal Volume Issue Start page Publication Year Abstract Keywords	Merge start and end pages
City of publication Address Type Web URL Link to document	OK X Cance
Misc VISBN	

The following importation options are available:

Select variables to import: This list box allows you to put check marks beside the variables you want to import. All unchecked items will be ignored.

Transform publication year into a date: The publication dates in RIS files can consist of a full date, with days and months, or only the year, or sometimes the publication month and year. By default, WordStat will import just the publication year and store it in an integer variable. Choosing this option will store all dates into a **Date** variable. If only the year is specified, WordStat will set the day and month to January 1. If the publication date consists of a month and a year only, then WordStat will set the day to the first day of the month.

Merge start and end pages: By default, start-page and end-page numbers are stored in separate variables. Selecting this option will join both numbers with a hyphen character and store in a single string variable.

• Once your options are set, select **OK**. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported bibliographic data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Related Topics:

To configure the data set and start analyzing please see The Data Tab.

To append new bibliographic references to your project please Monitoring Online Resources.

Importing News Transcripts

Factiva and **Nexis UNI** are online database giving their users access to previously published news from newspapers and magazines around the world, TV and radio news transcripts as well as additional information specific to each service such as online blogs, financial and market news, company information, case law, etc. WordStat can import files obtained from those two service, allowing one to analyze media coverage of specific topics, compare it by media outlets, countries, asses changes over time, etc.

Creating a New Project with News Transcripts

• Select the New button on the Data tab. This command calls up a dialog box similar to the one below.

Creat	a project from	a list of docume	ents	
	t from datafiles o	or web services	6.	
RunD	ocument Conver	sion <u>W</u> izard		

- Select the Import from data files or web services button.
- Select the NEWS menu item and choose the desired news aggregation service from its submenu.

For instruction on how to import news transcripts from those two sources, please see:

Importing from <u>Nexis Uni</u> Importing from <u>Factiva</u>

Nexis UNI

Nexis UNI from LexisNexis contains information from over 10,000 news, legal and business sources, including international and domestic newspaper, magazines, broadcast news transcripts, company financial information, industry and market news, law reviews, medical news, etc. It also includes business information on over 80 million public & private international companies and 75 million executives. Non-English language news sources are available in Spanish, French, German, Italian, and Dutch. Campus news from some 400 college/university papers and over 50 wire services are also available.

For the sake of this presentation, we will make references below to news transcripts, yet similar procedures may be followed to import other types of text data from Nexis UNI.

To import news transcripts from Nexis Uni, you first need to access the online service, perform your searches, select the relevant news transcripts and then save those on disk using the download command. A dialog box like this one will appear:

Download



Selected Documents (10)	Basic Options	Formatting Options	Content-specific Options
What do you want to downlo	bad?		
 Full documents (10) Include document attach Full document attachments Results list for 'News': (10) 	ments, where available only (where available)		
□ Include Bibliography			
File type			
O PDF 100 document limit.			
MS Word (docx) 100 document limit.			
 Rich Text Format (rtf) 100 document limit. 			
When downloading multiple	documents		
 Group and save documents Save as individual files 	as a single file (Note: Att	achments will be delivered as	separate documents)
□ Compress files in .ZIP format	What's this?		
Filename			
100-c Files(10)	haracter limit		

Make sure the **Full Documents** option is selected, and choose either the **MS Word (docx)** or the **Rich Text Format (rtf)** file format is selected. It is also recommended to keep documents grouped in a single file rather than individual files. You may need to create several of those files if you need to import more than the maximum number of documents allowed per download.

Once all the news transcripts have been downloaded, you can extract the news transcripts from those files with WordStat by following these steps:

• Select the New button on the Data tab. This command calls up a dialog box similar to the one below.



- Select the Import from data files or web services button.
- Select the NEXIS UNI menu item from the NEWS menu.
- A dialog box similar to the one below will appear allowing you to select the files containing the news transcripts.

🔳 Imp	ort Documents					_	- 🗆	×
F	ile type: LexisNexis Nexis UNI files (*.rtf	;*.docx)			\sim		
Đ	Live_Version_1.2	^	Name	Size	Date modified	Date created	Date a	cessed
Đ	New folder		Files(100) (1).DOCX	594 KB	1/28/2020 8:02 AM	1/28/2020 8:02 AM	5/22/20	21 10:27 /
+) Nexis UNI files]	Files(100) (2).DOCX	599 KB	1/28/2020 8:03 AM	1/28/2020 8:03 AM	5/22/20	21 10:27 A
±	order		Files(100) (3).DOCX	559 KB	1/28/2020 8:04 AM	1/28/2020 8:04 AM	5/22/20	21 10:27 A
±	Provalis		Files(100) (4).DOCX	544 KB	1/28/2020 8:07 AM	1/28/2020 8:07 AM	5/22/20	21 10:27 A
±	Workshop		Files(100) (5).DOCX	523 KB	1/28/2020 8:08 AM	1/28/2020 8:08 AM	5/22/20	21 10:27 A
🛨 🗄	Documents		Files(100) (6).DOCX	522 KB	1/28/2020 8:10 AM	1/28/2020 8:10 AM	5/22/20	21 10:27 A
🗉 🚽	Downloads		Files(100) (7).DOCX	518 KB	1/28/2020 8:11 AM	1/28/2020 8:11 AM	5/22/20	21 10:27 A
	Music		Files(100) (8).DOCX	509 KB	1/28/2020 8:12 AM	1/28/2020 8:12 AM	5/22/20	21 10:27 A
+	Pictures		Files(100) (9).DOCX	520 KB	1/28/2020 8:14 AM	1/28/2020 8:14 AM	5/22/20	21 10:27 A
÷ -	Videos		Files(100) (10).DOCX	524 KB	1/28/2020 8:15 AM	1/28/2020 8:15 AM	5/22/20	21 10:27 A
± =	PR-Files-Backup (\\PR-NAS-01) (B:)		Files(100) (11).DOCX	512 KB	1/28/2020 8:18 AM	1/28/2020 8:18 AM	5/22/20	21 10:27 A
🛨 🟪	System (C:)							
主 - 🚅	DVD RW Drive (D:)	~	<					>
					🖗 Add 🛛 👆 Remov	ve		
	·····							
1								
						v 0	к	Cancel

- Using the upper left panel, navigate to the folder containing the downloaded files.
- In the upper right panel, files in the selected folder that match the file type will be listed. Select the files you want to import and click the Add button to move them to the bottom file list.
- Once the files to import have been moved to the bottom section of the dialog box, click the **OK** button. WordStat will import all files by extracting each news transcript separately. It will also store various information in as many variables such as the title, the body of the article, the publication type, the source, the author, etc. WordStat will also save all indexing tags that have been assigned to the specific news respective categories (subject, geographic, organization, person, industry, etc.). If needed, variables containing superfluous information may then be deleted (see <u>Deleting Existing Variables</u>).

wordstat 9.0.7	- iest.ppj													-	
🔝 Data															
🔁 Open 🔹	🤔 New 🔛	Save 🔚 Sa	ve as 🛛 🙀 Export 🔢 Prop	erties				Analyze]						
A Accend	FILENAME	PUR TYPE	SOURCE	DATEPUB	AUTHOR	TITLE	BODY	LANGUAGE	LOAD DATE	NEWORDS	SUBJECT	GEOGRAPHIC	ORGANIZAT	INDUSTRY	PERSO
10 Happing	Files(100) (9)	Journal	Cameroon Tribune (Vaoundé)	1/27/2010	Brice Mheze	IDOCUM	IDDCUM	FRENCH	2/12/2010	893	IDOCUM	IDOCUMENTI	[document]	[document]	Idocume
lo Delete	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	5/3/2013	Aliance Nuchia	IDOCUM	IDDCUM	FRENCH	5/3/2013	442	IDOCUM	(DOCUMENT)	IDOCLIMENT1	IDOCLIME	Idocume
Dupicates	Files(100) (9)	Journal	Cameroon Tribune (Vaoundé)	11/14/2013	Badiano Ba Nken	IDOCUM.	IDOCUM.	FRENCH	11/15/2013	600	IDOCUM.	(DOCUMENT)	[document]	IDOCUME.	Idocum
(setures	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	11/18/2013	RiaDiba	IDOCUM	IDOCUM	FRENCH	11/19/2013	431	IDOCUM	IDDCLIMENT1	[document]	IDOCUME	Idocum
Descriptor	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	7/28/2010	Josiane Tchakounte	IDOCUM.	IDOCUM.	FRENCH	7/28/2010	404	IDOCUM.	IDOCUMENT1	[document]	IDOCUME	Idocur
Filter	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	9/21/2010	Essama Essomba	IDOCUM.	IDOCUM.	FRENCH	9/21/2010	304	IDOCUM.	(DOCUMENT)	(DOCUMENT)	IDOCUME	Idocur
	Files(100) (9)	Journal	Cameroon Tribune (Vaoundé)	3/7/2014	Alain Tchakounte	IDOCUM	IDOCUM.	FRENCH	3/10/2014	556	IDOCUM.	(DOCUMENT)	[document]	IDOCUME.	Idocur
ow editing	Files(100) (9)	Journal	Cameroon Tobune (Yaoundé)	2/7/2011		IDOCUM	IDDCUM.	FRENCH	2/8/2011	372	IDOCUM.	(DOCUMENT)	[document]	IDDCUME	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	5/21/2012	Alliance Nuchia	IDOCUM.	IDOCUM.	FRENCH	5/22/2012	435	IDOCUM.	(DOCUMENT)	[document]	IDOCUME	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	11/16/2011	Makon Ma Pondi	IDOCUM.	IDOCUM.	FRENCH	11/16/2011	662	IDOCUM.	IDOCUMENTI	[document]	IDOCUME	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	4/6/2010	Alain Tchakounte	IDOCUM.	IDOCUM.	FRENCH	4/6/2010	438	IDOCUM.	IDOCUMENT1	IDOCUMENT1	[document]	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	7/16/2010	Elice Zemine	IDOCUM.	IDOCUM.	FRENCH	7/16/2010	368	IDOCUM.	(DOCUMENT)	[document]	[document]	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	12/18/2014	Monda Bakoa	IDOCUM.	IDOCUM.	FRENCH	12/18/2014	536	IDOCUM.	[document]	[document]	IDOCUME.	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	4/4/2011	Sainclair Mezing	IDOCUM.	IDOCUM.	FRENCH	4/5/2011	336	IDOCUM.	IDOCUMENT1	[document]	IDOCUME	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	3/1/2010	Jeanine Fankam	IDOCUM.	IDOCUM.	FRENCH	3/1/2010	392	IDOCUM.	(DOCUMENT)	[document]	[document]	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	10/27/2010		IDOCUM.	IDOCUM.	FRENCH	10/28/2010	328	IDOCUM.	IDOCUMENTI	[document]	[document]	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	2/3/2010	Simon Pierre Etoundi	IDOCUM.	IDOCUM.	FRENCH	2/12/2010	679	IDOCUM.	IDOCUMENT1	[document]	IDOCUME	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	5/6/2010	Eric Vincent Formo	IDOCUM.	IDOCUM.	FRENCH	5/6/2010	600	IDOCUM.	(DOCUMENT)	[document]	IDOCUME	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	5/25/2010	Alfred Myogo Biveck	IDOCUM.	IDOCUM.	FRENCH	6/4/2010	449	IDOCUM.	IDOCUMENT1	[document]	[document]	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	3/25/2010	Steve Libaro	IDOCUM.	IDOCUM.	FRENCH	3/25/2010	490	IDOCUM.	IDDCUMENT1	[document]	IDOCUME	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	5/21/2012	Walfo Mongo	IDOCUM.	IDOCUM.	FRENCH	5/22/2012	562	IDOCUM.	IDOCUMENT1	[document]	IDOCUME.	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	7/31/2013	Messi Bala	IDOCUM.	IDOCUM.	FRENCH	7/31/2013	302	IDOCUM.	(DOCUMENT)	[document]	IDOCUME	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	5/3/2013	Monda Bakoa	IDOCUM.	IDOCUM.	FRENCH	5/3/2013	588	IDOCUM	(DOCUMENT)	IDOCUMENT1	IDOCUME	Idocu
	Fileo(100) (9)	Journal	Cameroon Tribune (Yaoundé)	2/22/2010	Josiane Tchakounte	IDOCUM.	IDOCUM.	FRENCH	2/22/2010	478	IDOCUM.	(DOCUMENT)	[document]	[document]	Idocu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	10/3/2013	Jean Francis Belbi	(DOCUM.	DOCUM.	FRENCH	10/3/2013	317	DOCUM	(DOCUMENT)	[document]	(DOCUME	[docu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	2/24/2014	Josiane Tchakounte	[DOCUM.	(DOCUM	FRENCH	2/25/2014	624	DOCUM	[DOCUMENT]	[document]	(DOCUME	[docu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	5/10/2010	Steve Libam	(DOCUM.	DOCUM.	FRENCH	5/10/2010	477	DOCUM	[DOCUMENT]	[document]	(DOCUME	[docu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	6/18/2010	Augustin Fogang	[DOCUM.	DOCUM.	FRENCH	6/18/2010	638	DOCUM	[DOCUMENT]	[document]	(DOCUME	[docu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	5/7/2010	Josiane Tchakounte	(DOCUM.	(DOCUM	FRENCH	5/7/2010	255	(DOCUM	(DOCUMENT)	(DOCUMENT)	[document]	[docu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	10/28/2014	Eric Elouga	(DOCUM.	DOCUM.	FRENCH	10/28/2014	438	(DOCUM	[DOCUMENT]	[document]	(DOCUME	[docu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	5/23/2012	Jean Baptiste Ketchateng	(DOCUM.	(DOCUM	FRENCH	5/29/2012	358	(DOCUM	[DOCUMENT]	[document]	(DOCUME	[docur
	Filed(100) (9)	Journal	Cameroon Tribune (Yaoundé)	6/21/2010	Patrice Mbossa	(DOCUM.	DOCUM.	FRENCH	6/22/2010	425	(DOCUM	[DOCUMENT]	[document]	IDOCUME	[docur
	Files(100) (9)	Journal	Cameroon Tobune (Yaoundé)	3/7/2012	Steve Libam	[DOCUM.	DOCUM.	FRENCH	3/7/2012	435	[DOCUM	(DOCUMENT)	[document]	(DOCUME	[docu
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	5/7/2010	Eric Vincent Forso	(DOCUM.	DOCUM.	FRENCH	5/7/2010	462	(DOCUM.	[DOCUMENT]	[DOCUMENT]	(DOCUME	[docur
	Files(100) (9)	Journal	Cameroon Tribune (Yaoundé)	7/16/2014	Angèle Bepede	[DOCUM.	(DOCUM	FRENCH	7/16/2014	651	(DOCUM	[DOCUMENT]	[document]	(document)	[docur
	Eaut1000100	I av ment	Common Table on Manundal	6/20/2011	Calculat Maning	monund	10.00114	COENCU	6 /00 /0011	- 400	monet	IDOCUMENT1	Advances M	monute	Line

Factiva

Factiva is a business information and research tool owned by Dow Jones & Company that aggregates content from both licensed and free sources. Factiva provides access to more than 32,000 sources (such as newspapers, journals, magazines, television and radio transcripts, photos, etc.) from nearly every country worldwide in 28 languages, including more than 600 continuously updated newswires. It can also retrieve press releases and financial information about public and private companies as well as historical and current market indices.

For the sake of this presentation, we will make references below to news transcripts, yet similar procedures may be followed to import other types data retrieved from Factiva.

To import news transcripts from Factiva, you first need to access the online service, perform your searches, select the

relevant news transcripts and then save those on disk in RTF format by click the me button and selecting Article Format.

Once all the news transcripts have been downloaded, you can extract the news transcripts from those files with WordStat by following these steps:

• Select the New button on the Data tab. This command calls up a dialog box similar to the one below.



- Select the Import from data files or web services button.
- Select the FACTIVA menu item from the NEWS menu.
- A dialog box similar to the one below will appear allowing you to select the files containing the news transcripts.

Import D	ocuments						- 🗆 ×
File typ	e: Factiva files (*.rtf)					\sim	
Ξ.	Desktop	^	Name	Size	Item type	Date modified	Date created 🗠
1	± Avi		Factiva-2020060	1.97 MB	Rich Text Format	6/2/2020 7:58 AM	6/2/2020 7:58 AM
	Backup		Factiva-2020060	1.94 MB	Rich Text Format	6/2/2020 7:59 AM	6/2/2020 7:59 AM
	Bannieres		Factiva-2020060	1.69 MB	Rich Text Format	6/2/2020 8:00 AM	6/2/2020 8:00 AM
	Documents		Factiva-2020060	2.42 MB	Rich Text Format	6/2/2020 8:01 AM	6/2/2020 8:01 AM
	E Factiva files		Factiva-2020060	2.01 MB	Rich Text Format	6/2/2020 8:03 AM	6/2/2020 8:03 AM
	E Live_Version_1.2		Factiva-2020060	771 KB	Rich Text Format	6/2/2020 8:04 AM	6/2/2020 8:04 AM
1	E Order		Factiva-2020060	2.29 MB	Rich Text Format	6/2/2020 8:07 AM	6/2/2020 8:07 AM
	± Provalis		Factiva-2020060	1.38 MB	Rich Text Format	6/2/2020 8:08 AM	6/2/2020 8:08 AM
	E Workshop		Factiva-2020060	1.40 MB	Rich Text Format	6/2/2020 10:33 AM	6/2/2020 10:33 AM
+	Documents		Factiva-2020060	784 KB	Rich Text Format	6/2/2020 10:33 AM	6/2/2020 10:33 AM
+	Outlook		Factiva-2020071	1.53 MB	Rich Text Format	7/17/2020 7:47 AM	7/17/2020 7:47 AM
• • • • • • • • •	Part 2	~	Factiva-2021052	12.0 MB	Rich Text Format	5/22/2021 11:03 AM	5/22/2021 11:03 AM ¥
<		>	<				>
					n Add 🔶 Remo	ve	
E: \Normand \De	sktop\Factiva files\Factiva-202	:00602	2-0656.rtf				
						~	OK X Cancel

- Using the upper left panel, navigate to the folder containing the downloaded files.
- In the upper right panel, files in the selected folder that match the file type will be listed. Select the files you want to import and click the Add button to move them to the bottom file list.
- Once the files to import have been moved to the bottom section of the dialog box, click the **OK** button. WordStat will import all files by extracting each news transcript separately. It will also store various information in as many variables such as the title, the body of the article, the publication type, the source, the author, etc. If needed, variables containing superfluous information may then be deleted (see <u>Deleting Existing Variables</u>).

Data								
i Open ▼ Editor Struct	🤗 New 📻 Save 🖥	Save as 🛛 🔒 Export	대. Properties	🍃 Analyze				
Append	FILENAME	SOURCE	DATEPUB	AUTHOR	TITLE	BODY	LANGUAGE	NBWORDS
	Factiva-20200602-0656	The New York Times	6/28/2020	By Javier C. Hernández and Benjamin Mueller	(DOCUM	[DOCUM	Anglais	15
o Delete	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020		DOCUM	DOCUM	Anglais	17
Duplicates	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020		[DOCUM	DOCUM	Anglais	35
-	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020	By David Leonhardt	[DOCUM	DOCUM	Anglais	18
Descriptor	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020	By Giovanni Russonello	[DOCUM	DOCUM	Anglais	13
Filter	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020	By Ross Douthat	[DOCUM	DOCUM	Anglais	10
	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020		[DOCUM	[DOCUM	Anglais	2
ow editing	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020		[DOCUM	[DOCUM	Anglais	13
	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020	By John Branch	(DOCUM	(DOCUM	Anglais	29
	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020	By A.O. Scott	(DOCUM	(DOCUM	Anglais	15
	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020	By Kevin Draper and Julie Creswell	(DOCUM	[DOCUM	Anglais	8
	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020	By Tonya Russell	(DOCUM	[DOCUM	Anglais	10
	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020	By Nick Corasaniti	(DOCUM	[DOCUM	Anglais	15
	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020	By Gail Collins and Bret Stephens	(DOCUM	[DOCUM	Anglais	16
	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020	By Jamelle Bouie	(DOCUM	[DOCUM	Anglais	12
	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020		(DOCUM	[DOCUM	Anglais	5
	Factiva-20200602-0656	NYTimes.com Feed	6/2/2020	By Trish Bendix	(DOCUM	[DOCUM	Anglais	9
	Factiva-20200602-0656	The New York Times	6/2/2020	By Declan Walsh	(DOCUM	[DOCUM	Anglais	10
	Factiva-20200602-0656	The New York Times	6/2/2020		(DOCUM	[DOCUM	Anglais	
	Factiva-20200602-0656	The New York Times	6/2/2020		(DOCUM	[DOCUM	Anglais	8
	Factiva-20200602-0656	The New York Times	6/2/2020		(DOCUM	[DOCUM	Anglais	2
	Factiva-20200602-0656	The New York Times	6/2/2020		[DOCUM	(DOCUM	Anglais	
	Factiva-20200602-0656	The New York Times	6/2/2020	By Evan Hill, Ainara Tiefenthäler, Christiaan Triebert, Drew Jordan,	. [DOCUM	(DOCUM	Anglais	13
	Eactive-20200602-0656	The New York Times	£ /2 /2020	Pu Sanna Mahashuari and Mishaal Corkery	гоосци	гоосции	Anglaio	12

Importing from Email Servers

Email has become a major source of communication and may also be used as a cost-efficient data collection tool. As a result, electronic mailboxes often contain a wealth of useful information which may be the basis of a research project. WordStat allows you to create a project by importing email data and use the text mining and content analysis features of WordStat to analyze your email data. You can create projects from Outlook, Gmail and Hotmail cloud-based email servers. You can also create projects from the Outlook account on your PC and from PST and MBOX files.

Creating a New Project Using Email Data

• Select the New button on the Data tab. This command calls up a dialog box similar to the one below.



- Select the Import from data files or web services button.
- Select the E-MAILS menu item and choose the desired email service from its submenu.

For instruction on how to import email from various sources, please see:

Importing from <u>Outlook</u> Importing from <u>Gmail</u> Importing from <u>Hotmail</u> Importing from an <u>Outlook Account</u> Importing from a <u>PST</u> file Importing from a <u>MBOX</u> file

Outlook

There are three ways to access emails stored in Outlook. You can connect through the Outlook.com email server, connect directly with the outlook application on your computer or import Outlook data contained in PST files. PST files can be imported regardless of whether or not you have access to the account from which they originated. Two examples of when a PST importation could be used are to import backups of old emails or emails that are not yours and for which you don't have access to the email account.

Connecting to Outlook.com

- Select E-MAILS | OUTLOOK | OUTLOOK.COM from the menu to log in to Microsoft Outlook's cloud-based email service. A log in dialog box will appear.
- Log in to your Outlook.com account. A permissions dialog box will appear.

		El ante de	
PROVALIS	et this app a	ccess your i	nfo?
QDA Mine	er needs you to co	onfirm its perm	ission to:
Rea QD, mai	d your mail A Miner will be ab Ibox.	le to read email	in your
You can c	hange these appli	cation permissi	ons at any
time in yo	our account setting	JS.	
-			

• Click Yes if you give permission for WordStat to access your account. An Import Emails dialog box will appear.

💽 Import Emails		\overline{a}		×
Folders	Variables			
Deleted Items (0) Drafts (0) Inbox (2) Junk Email (0) Outbox (0) Sent Items (1)	Image: Section of the	d		
		OK	×	Cancel

- The items in the **Folders** list box on the left side of the dialog box correspond to the folders in your email account. Check the **Folders** you wish to import.
- The items in the **Variables** list box on the right side of the dialog box will become project variables if imported. Check the **Variables** you want to include in your project.
- Click OK.
- Enter the project name and save it in the location of your choosing. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported email data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Importing Data from an Outlook Email Application on Your Computer

• Select E-MAILS | OUTLOOK | ACCOUNT from the menu to import email from the Outlook application that is installed on your computer. An Import from Outlook dialog box will appear.

- The items in the **Folders** list box on the left side of the dialog box correspond to the folders in your email account. Check the **Folders** you wish to import.
- The items in the **Variables** list box on the right side of the dialog box will become project variables if imported. Check the **Variables** you wish to import.
- Click Import.
- Enter the project name and save it in the location of your choosing. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported email data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Importing Data From a PST File

- Select E-MAILS | OUTLOOK | PST FILE from the menu to import a file saved on your PC in the Personal Storage Table (*.pst) format. An Import data file dialog box will appear.
- Select the PST file you want to import.
- Click Open. An Import from PST dialog box will appear.

import riom PS1			-	Ц	>
olders:		Variables			
 Archive Folders □Deleted Items (33) □DiskStation (0) □DiskStation 1 (53) □Trbox (50) □Pierre (17) □Adam (7) □Lida (0) □Dropbox (14) □TMS (11) □Alain (12) □Jonathan (2) □TATUK (14) □Survey (0) □Tech (0) □Facebook (0) □App keys (2) □Cutbox (0) □Sent Items (109) □Calendar (0) □Journal (0) □Notes (0) 	~	 Folder From Name From Email Sent To CC BCC Subject Received Time Sent On Body Importance Sender IP Address 			

- The items in the **Folders** list box on the left side of the dialog box correspond to the folders in your email account. Check the **Folders** you wish to import.
- The items in the **Variables** list box on the right side of the dialog box will become project variables if imported. Check the **Variables** you wish to import.
- Click Import.
- Enter the project name and save it in the location of your choosing. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported email data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To configure the data set and start analyzing please see The Data Tab.

To append new email data to your project please Monitoring Online Resources.

Hotmail

Connecting to Hotmail

- Select E-MAILS | HOTMAIL from the menu. A log in dialog box will appear.
- Log in to your Outlook.com account. A permissions dialog box will appear.

		000
PROVAL	Let this app a	ccess your info?
QDA	Miner needs you to c	onfirm its permission to:
~	Read your mail QDA Miner will be at mailbox.	ble to read email in your
You c time i	an change these appl in your account settin	ication permissions at any gs.
T		

• Click Yes if you give permission for WordStat to access your account and an Import Emails dialog box will appear.

- 0 3
Variables
From CC BCC Subject Body To Received Sent Folder

- The items in the **Folders** list box on the left side of the dialog box correspond to the folders in your email account. Check the **Folders** you wish to import.
- The items in the **Variables** list box on the right side of the dialog box will become project variables if imported. Check the **Variables** you want to include in your project.
- Click OK.

• Enter the project name and save it in the location of your choosing. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported email data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Related Topics:

To configure the data set and start analyzing please see The Data Tab.

To append new email data to your project please Monitoring Online Resources.

Gmail

Connecting to Gmail

- Select E-MAILS | GMAIL from the menu and you will be prompted to log in to your Gmail account.
- Log in to your account. A permissions dialog box will appear similar to the one below.



• Click Allow if you give permission for WordStat to access your Gmail account. An Import Emails dialog box will appear.

Variables
From CC Subject Body To Received Labels IP

- The items in the **Folders** list box on the left side of the dialog box correspond to the folders in your Gmail account. Check the **Folders** you wish to import.
- The items in the **Variables** list box on the right side of the dialog box will become project variables if imported. Check the **Variables** you want to include in your project.
- Click OK.
- Enter the project name and save it in the location of your choosing. The data will be imported and WordStat's **Data** tab will appear, containing a table with your imported email data. From here you can further tailor your data set if necessary, or start analyzing immediately.

To configure the data set and start analyzing please see The Data Tab.

To append new email data to your project please Monitoring Online Resources.

MBox

MailBox (*.mbox) is a file extension for a text file used to store email. It was first implemented for Fifth Edition Unix, and though not formally defined by Internet Engineering Task Force (IETF) and the Internet Society (ISOC), it is the most common format for storing email messages on a hard drive. All messages are stored in a single plain text file. Each message starts with a "From" line and ends with a blank line.

Importing Data From a MBOX File

- Choose the E-MAILS | MBOX from the menu and Import data file dialog box will appear.
- Select the MBOX file that you would like to import.
- Click Open.

• Enter the project name and save it in the location of your choosing. An Import Mbox box will appear.

From Sent To CC		
Subject Date Body		
Sender IP Add	ress	

- The items in the **Select Fields** list box will become project variables when imported. Select the fields you wish to import.
- Click **Import** and WordStat's **Data** tab will appear, containing a table with your imported email data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Related Topics:

- To configure the data set and start analyzing please see The Data Tab.
- To append new email data to your project please Monitoring Online Resources.

Using the Document Conversion Wizard

The Document Conversion Wizard is a utility program used to import one or more documents into a new project file. This tool supports the importation of numerous file formats including plain ASCII, Rich Text Format, MS Word, HTML, Acrobat PDF files, and WordPerfect documents. It may be instructed to split large files into several cases and to extract numeric and alphanumeric data from these files. The Document Conversion Wizard may be run either as a stand-alone application or from within WordStat. When no splitting, transformation or extraction of information from documents are necessary, you can also use the easier to use <u>Creating a Project from a List of Documents</u> importation method.

To run the Document Conversation Wizard from WordStat:

• Select the New button on the Data tab. A dialog box similar to the one below will appear.



• Click the **Run Document Conversion Wizard** button. The program then guides you through the necessary steps to import the documents and extract relevant data.

For more detailed information on the document importing process using this utility program, consult the Document Conversion Wizard help file.

Starting WordStat Document Explorer from Windows Explorer

WordStat Document Explorer can be started directly from Windows Explorer. WordStat Document Explorer, a pared-down version of WordStat, allows you to explore your documents, listing word and phrase frequencies and associated text to quickly reveal themes present in the selected documents. You can also apply a previously saved categorization or classification model when exploring documents.

Words	rases	contents				88	8		Paragraphs O Document
	FREQ	% SHOWN	% WORDS	CASES	% CASES	TF-I ^	DOCUMENT	FREQ	9 Still our disappointments should not obscure
AMERICA	152	7.0%	0.6%	11	84.62	11./	McCain - Foreign Policy	17	expecting too much too soon, but they should
AMERICAN	115	5.3%	0.4%	11	84.62	8.3	Bush - Foreign Policy	10	not intimidate us from making the most of this
PEOPLE	98	4.5%	0.4%	11	84.62	7.1	Gore - Foreign Policy	5	moment to continue building a safer, more bumane world than the one that inaugurated this
CHINA	53	2.4%	0.2%	6	46.15	17.	McCain - Announcement	4	century. I have been pleased to support Clinton
FREEDOM	53	2,4%	0.2%	8	61.54	11.	Bush - Announcement	2	Administration policies that have seized
GREAT	53	2.4%	0.2%	10	76.92	6.0	Buchannan - Announcement	1	and ideals abroad. I strongly supported the
CENTURY	51	2.3%	0.2%	9	69.23	8.1	Forbes - Announcement	1	expansion of NATO, for example, and the two
COUNTRY	50	2.3%	0.2%	10	76.92	5.7	Bradley - Foreign Policy	1	great trade successes NAFTA and the Uruguay Round and the need for stronger action in
WAR	46	2.1%	0.2%	6	46.15	15.			responding to the crisis in Kosovo.
PEACE	44	2.0%	0.2%	6	46.15	14.:			
GOVERNMENT	42	1.9%	0.2%	8	61.54	8.9			Administration policies because they more often
INTERESTS	41	1.9%	0.1%	8	61.54	8.6			than not manifest two closely-related and central
AMERICANS	40	1.8%	0.1%	8	61.54	8.4			flaws: strategic incoherence and self-doubt. The
DIICETA	20	1 00/	0 10/	-	AC 10	10 4			framework that establishes the relationships
CHILDREN	EMOCF	RACY FREE	POLICY RE	AGAN S	TRATEGIC	9 8			second refers to a mystifying uncertainty about how America should act in a world where we are the only superpower.
					ACE ALLES	RPOSE SHIP RPOSE APOSE APONS TH TO TO			supported paying our UN dues, as I have supported sound arms control treaties and international assistance and any policies that effectively address the contemporary threats to our core strategic interests and defining ideals. But if those efforts serve no larger purpose than temporarily placating America's foreign critics - whose number, sadly, never seems to appreciably decline as America's singular contributions to world peace and human

WordStat Document Explorer

WordStat Document Explorer contains three or four tabs. The **Data** tab lists the files you have imported. On the top left of the **Frequencies** tab, WordStat Document Explorer allows you to view the most frequent words and phrases or categories (if a categorization module has been applied) in your selected documents. Below the frequencies list in the left panel is an interactive word cloud. A list of documents containing the selected words, phrases or categories appears in the middle panel. On the right you can toggle between viewing the paragraphs or full documents containing the selected word, phrase or category. The **Document** tab contains the word or category count per document. The **Classification** tab, only visible if you have run a classification model, lists the selected documents classified according to the chosen model.

WordStat Document Explorer with a Categorization Model

WordStat allows you to save and publish categorization models so you can analyze new projects with a previously created model. WordStat Document Explorer will apply the categorization dictionary, exclusion list and substitution list etc. contained in the chosen model, allowing you to see the most frequent categories, the documents containing the words and phrases

and rules in those categories, and the associated text. A categorization model can be chosen when opening WordStat Document Explorer from the Windows Explorer menu or can be applied while WordStat Document Explorer is already open. For more information please see <u>Categorization</u>.

WordStat Document Explorer with a Classification Model (WordStat Document Classifier)

WordStat also allows you to save classification models so that you can classify documents with a previously tested model. You can see the most frequent words and phrases or categories, the documents containing the words and phrases, and the associated text on the **Frequencies** tab. A **Classification** tab has been added that lists the documents classified according to the chosen model. For more information please see <u>Classification</u>. A classification model can be chosen when opening WordStat Document Explorer from the Windows Explorer menu or can be applied while WordStat Document Explorer is already open.

To open WordStat Document Explorer from individual documents in Windows Explorer:

- In Windows Explorer, select the individual documents you would like to explore, using the **Shift** or **Ctrl** key to select numerous documents at the same time.
- Right click your mouse, a menu opens, scroll down to **WORDSTAT CONTENT ANALYSIS** and choose **EXPLORE**, **CATEGORIZE** or **CLASSIFY** from the adjacent menu.
- If you have chosen **CATEGORIZE** or **CLASSIFY**, select from the adjacent menu one of your previously saved categorization or classification models. WordStat Document Explorer opens, pausing briefly on the **Data** tab, moving automatically to the **Frequencies** tab once the data has been processed.

To open WordStat Document Explorer from a folder in Windows Explorer:

- If you with to analyze all documents in a folder, select the folder.
- Right-click your mouse to display the context menu.
- Scroll down to **WORDSTAT CONTENT ANALYSIS** and choose either **EXPLORE**, **CATEGORIZE** or **CLASSIFY** from the adjacent menu.
- If you have chosen **CATEGORIZE** or **CLASSIFY**, select from the adjacent menu one of your previously saved categorization or classification models.
- A dialog box will appear allowing yo to choose which **file types** to import and analyze and giving you the option of including documents in **subfolders** as well.

r: C: Users Amanda Documents My Prov	alis Research Projects \Workshop \GOP	
Include subfolders		
File types		
Text file (*.txt)	Rich Text files (*.rtf)	Acrobat PDF files
MS Word files (*.doc;*.docx)	HTML files	Powerpoint (*.ppt;*.pptx)
Ebooks (*.epub)	OpenOffice (*.odt)	WordPerfect (*.wpd)
XML Paper Specification (*.xps)		
• Once all files types have been selected, click the **OK** button. WordStat Document Explorer opens, pausing briefly on the **Data** tab, moving automatically to the **Frequencies** tab once the data has been processed.

To apply a categorization or classification model while in WordStat Explorer:

- Select the solution. A dialog box will appear containing your previously saved categorization and categorization models.
- Select the model you would like to apply and click the **Open** button.
- Select the button.

The User Interface

The WordStat user interface is built around a multi-tab workspace that provides an integrated environment for project creation, data management, developing, transformation, editing, testing, validating, applying a content analysis dictionary and performing various text mining tasks.

<u>Data</u> - This tab appears only when running WordStat as a standalone application and will not be visible if you run WordStat from QDA Miner, Simstat or Stata. This is where all the data you have imported into your project is displayed. The rows in the data tab represent cases and the columns are variables. You can append and delete cases, as well as edit the data within. The structure tab allows you to view and modify your project variables.

<u>Text Processing</u> - This tab contains various options controlling how the text data will be processed. It also includes options affecting linguistic tools. It allows you to create and modify dictionaries, exclusion and substitution lists, as well as add, remove and edit existing dictionaries.

<u>Frequencies</u> – This tab displays a table of the frequency of keywords or content categories. You may also view a list of leftover words, allowing you to modify the current categorization dictionary, the exclusion list or the substitution list. There is also a suggestions tab displaying leftover words potentially related to the currently selected item and comparison graphs which visualize the relationship between the currently selected item and the chosen variables.

<u>Extraction</u> - This tab allows you to perform topic modeling, find the most common phrases, extract named entities, as well as misspellings and uncommon words, and assign them to the current categorization dictionary, the exclusion list or the substitution list. You may also use the misspelling tab to batch replaced misspellings in the original documents.

<u>Cooccurrences</u> - This tab allows you to explore connections between words, keywords, phrases or content categories using hierarchical clustering, multidimensional scaling, link analysis and proximity plots.

<u>Crosstab</u> – This tab allows you to compare keyword frequencies across values of numerical, categorical or date variables. Along with a table of frequencies, several statistics and graphical techniques may be applied including correspondence analysis, heatmaps, bubble charts, bar charts and line charts.

<u>Keyword-in-Context</u> – This tab allows you to display a concordance table of specific words, word patterns or phrases, or of all items related to a content category. This is very useful to validate a dictionary by allowing you to examine, in context, how words are being used.

<u>Classification</u> - This tab gives access to the automated text classification module that allows you to apply a machinelearning approach to the existing textual database. Options allow you to develop a classification model that can later be used to accurately classify uncategorized documents into predefined classes.

Note: When you are using the stand-alone version of WordStat in the Explore mode only the Data, Frequencies, Phrases, and Topics tabs are visible.

Two additional drop-down menus can be accessed to perform various tasks:

• Clicking the ≡ button in the upper left corner of the main window displays a menu that allows you to leave WordStat or return to the calling application as well as perform various tasks such as:

Editing documents Editing case descriptors Filtering cases Geocoding Accessing the Report Manager Program Setting Mode • The 🐲 button in the upper right corner provides access to this help file, which can also be accessed at any time by pressing the **F1** key. In addition, this menu allows you to check whether you are using the latest version of WordStat and also gives access to specific important information and some useful links to the Provalis Research website.

The Data Tab

The **Data** tab is visible only when you open WordStat directly. It does not appear when you call WordStat from QDA Miner, Simstat, or Stata. This tab allows you to open existing projects, modify those, or create new ones. Once your project had been created or opened, the data will be displayed on this tab. The **Data** tab allows you to append and delete cases, as well as edit the data within. You can also view and modify your project variables. Once you are satisfied with the content and structure of your data you can start your analysis.

The Data Editor Tab

Project data will be displayed in rows and columns. The rows on the **Data Editor** tab represent cases and the columns are project variables. You can append, delete and filter cases on this tab as well as identifying duplicates and changing case descriptors. You can also edit the project data directly on this tab

vvordStat 9.0 - 1	Elect	ion 2008b.pprj					- D	×
E 🛄 Data	- 1	ext Processing	a 📻 Frequencies 🤮	Extraction 🗞	Cooccurrences	Crosstab	III Keyword-In-Context 🛛 <table-cell-columns> Classification</table-cell-columns>	
Pata Editor Struct	Ure	New 🖬 Sa	ave 🔚 Save as 📑 Ex	oport 🏦 Properties			Analyze	
Append		CANDIDATE	FILE	DOCUMENT	PARTY	DELIVERY		^
1. 0.1.1.		Biden	1-Biden20060316	[DOCUMENT]	Democratic	2006		
Vo Delete		Biden	1-Biden20060719	[DOCUMENT]	Democratic	2006		
Duplicates	•	Biden	1-Biden20060907	[DOCUMENT]	Democratic	2006		
N Bernink		Biden	1-Biden20060920	[DOCUMENT]	Democratic	2006		
Descriptor		Biden	1-Biden20061205	[DOCUMENT]	Democratic	2006		1.10
Filter		Biden	1-Biden20070110	[DOCUMENT]	Democratic	Q1-2007		
-		Biden	1-Biden20070111	[DOCUMENT]	Democratic	Q1-2007		
_ Allow editing		Biden	1-Biden20070203	[DOCUMENT]	Democratic	Q1-2007		
		Biden	1-Biden20070205	[DOCUMENT]	Democratic	Q1-2007		
		Biden	1-Biden20070215	[DOCUMENT]	Democratic	Q1-2007		
			* D'1 00030000			~* ~~~~~~		*
	F	ive Year ive years ago, aid: "We will h he belly of a p	on September 10th, 20 nave diverted all that m lane, or are smuggled i	Rethinking / 001, standing at th oney to address th nto a city in the m	America's nis podium, I an ne least likely tl niddle of the nig	Future Se rgued against th hreat while the ght in a vial in a	ecurity his administration's fixation on national missile defense. I real threats come into this country in the hold of a ship, or backpack."	
	Ia	wasn't clairvo dministration	yant. I was making a p has the wrong premise:	oint that was valid s and the wrong p	l then and rem riorities.	ains valid today	: when it comes to America's national security, this	
	L	he President i	s right, as he put it this	week: we're "a n	ation at war." 1	That makes it a	I the more incomprehensible that, five years after 9/11, he has	~

Related Topics:

Appending Documents Appending from a Data File Monitoring a File or Folder Deleting Cases Identifying Duplicate Cases Filtering Cases Editing Project Data

The Structure Tab

The **Structure** tab displays project variables. Each row presents information about a variable. From this tab you can add, delete and transform variables as well as view variable properties and obtain the frequency and cross-frequency distribution of numerical, categorical, date and short-string variables.

WordStat 9.0.7	- Election 2008b	o.pprj			-	×
🔳 🔝 Data	Text Proces	sing 📑 Fr	requencies 🔋 Extraction 🔗 Cooccurren	nces 🛅 Crosstab	I Keyword-In-Context < Classification	0.
🔗 Open 🔻	🕑 New 🔓	Save	Save as 📑 Export 🏦 Properties		Analyze	
Data Editor Struct	ure					
Add	Name	Туре	Description	Missing Values	Values	
EQ Hod	CANDIDATE	Nominal			Biden; Obama; Richardson; Clinton; Edwards; Kucinich; Thompson; McCain;	
Delete	FILE	Nominal			1-Biden20060316; 1-Biden20060719; 1-Biden20060907; 1-Biden20060920;	
te Reorder	DOCUMENT	Document				
Ch. Transform	PARTY	Nominal			Democratic; Republican; Green; Independent	
Properties	DELIVERY	Nominal	Date period when the speech was delivered		2006; Q1-2007; Q2-2007; Q3-2007; Q4-2007; Q1-2008; Q2-2008; Q3-2008;	
THE FE						
243 cases					834 item(s) in 5 categorie(s)	-8

Related Topics:

Adding New Variables Deleting Existing Variables Transforming Variables Recoding Values of a Variable Editing Variable Properties Variable Statistics

The Toolbar

Control

Description



×

Open

This button opens a dialog that allow you to select and open previously created projects. Click the arrow to access a list of recently opened project.

This button opens a dialog that allows you to create a new project. See <u>Creating a Project in</u> <u>WordStat</u>.



This button allows you to save the currently opened project to disk.



This button opens a dialog which allows you to name your file and save the project to disk.



This button opens a dialog that allows you to <u>analyze</u> your data in both Explore mode and Expert mode.

Project Properties

The project properties function allows you to create a description of your project that may be displayed automatically when opening the project file as well as set access privileges to the project to limit the number of users accessing the data and performing specific operations.

To access the project properties function:

• Selecting the ^{11 Properties} button on the tool bar of the data page will display a dialog box similar to the one below.

Project pr	operties	>
Description	Security / User Access	
Description:		
Sample dat republicans	aset of 243 speaches from 10 candidates to the 2008 US Presidential race (6 democrats and 4),	
Show des	scription upon opening	

The Project properties dialog box contains tow tabs: Description and Security/User Access.

To save or edit a project description:

- Select the **Description** tab.
- Enter or edit a project description.
- If you wish to display this description automatically when the project file is opened, select the **Show description upon opening** checkbox.
- Click the **OK** button to save the changes you made. If you choose to quit this dialog box without saving the changes, click the **Cancel** button.

Adjusting security and user access settings

It is possible to limit who can access a specific project as well as limit the type of operations that can be performed by certain individuals by creating user accounts and requiring people to provide a username and password to access the project. This feature is useful for security or confidentiality issues, but it is also useful to prevent the deletion or editing of existing cases or variables.

When setting up a project to support multiple users, one of these users must be able to control access to the project, create and delete user accounts and define passwords. This user is commonly known as the **administrator**. When creating a new project, an administrator account is automatically created. Both the default username and password for this account is **ADMIN**. If you choose to restrict access to some users, it is highly recommended that you change this username and password to prevent unauthorized changes to user access rights.

To change the User Access settings of a project:

• Click the Properties button on the tool bar of the data page to display the project property dialog box and then move to the Security / User Access page. The dialog box should look like this:

	and the second	
Description Security /	User Access	
Users must log		
* Admin	Add	Available features:
Guest		Add or delete cases
	Edit	Add or delete variables & documents
	Delete	Modify existing variables
	Set as admin	Edit or transform documents
		Save modified project file
		Export project, tables or graphics
		Modify categorization model items
		Select a different categorization model
		Run post-processing scripts (Python or R)
		Create or edit scripts

By default, opening a project without using the user log screen gives the user all administrator access rights. Enabling the **Users must log** option will display a **Name and password** dialog box prompting the user to enter a valid username and password to access the project file.

To add a new user account:

• Click the Add button The following dialog box will appear:

Name and password	÷		×
User name:	_		
Enter password:			
Re-enter password:			
J OK	¥ Can	cel	

• Enter the Username.

- In the Enter password edit box, enter this user's password.
- Enter the password a second time in the next edit box to make sure it was entered correctly.
- Click **OK** to save the new user account. Once a user's account has been created, you may specify which specific features this user will be able to access.

To define the user access rights:

- Select the user for which you want to define or edit access rights.
- In the list of available features to the right of the dialog box, select the features you want this user to have access to and clear the features that you do not want them to access.
- To prevent users from modifying existing documents or values stored in variables, clear the check box beside the **Modify variables** and **Edit or transform documents** options. It may also be a good idea to disable the **Add or delete cases** and **Add or delete variables & documents** options.
- Alternatively, you may allow the user to perform any modification they want on the data set but prevent them from saving and exporting the project.
- You may prevent the user from modifying any categorization models accessed from this project, as well force this user to use a specific categorization model by preventing him from selecting a different one.
- Despite the presence of other restrictions such as the inability to export data, create variables, etc, Python and R scripts could circumvent some of those restrictions and allow a user to export sensitive data. Removing the Create or edit scripts will prevent the user from editing existing pre- or post-processing scripts or creating new ones. Please note that if you disable only this option, the user will still be able to execute existing post-processing scripts. Disabling the Run post-processing scripts will prevent this user from running those scripts. It will also prevent the creation or editing of existing post-processing scripts, even if the Create or edit scripts option is enabled. However, the creation and editing of pre-processing scripts will remain possible unless one disables this feature.
- Once the restrictions have been set, Click OK.

To delete a user account:

- In the list of existing accounts to the left of the dialog box, select the account you want to delete.
- Click the **Delete** button.
- To save the changes you made to the accounts and close this dialog box, click the **OK** button. To close this dialog box without saving any changes, click the **Cancel** button.

Appending Documents

To append documents and store them as new cases:

• Select the **Append** button and choose **DOCUMENTS** from the adjacent menu. A dialog box similar to the one below will appear:



- Click on a folder in the folders list on the upper left section of the dialog box to display its contents. If you want to see the contents of a drive, go to the folders list, click **My Computer**, and then double-click on a drive.
- In the upper right section of the dialog box, WordStat displays all supported document file formats that may be imported. To display only documents of a specific type, set the **File Type** list box to the desired file format.
- Click the file you would like to import. To select multiple files, hold down the CTRL key while clicking the other files.
- Click the Add button to add the files to the list of files to import, located at the bottom of this dialog box. You may also drag the files from the top right section to this list.
- Click the button to add numerous files contained within a folder. A dialog box will appear giving you the option of including **subfolders** and choosing **file types** to include in the import.
- To remove a file from the list of files to import, select that file name and click the management button.
- Once all files have been selected, click the **Append** button. If the project contains more than one categorical or document variable, a dialog box similar to this one will appear:

Variable Selection		×
Store file name in:	FILE	V
Store file location in:	LOCATION	Ý
Store document in:	DOCUMENT	~
🗸 ок	X Cancel	

• Select the categorical variable in which the file name will be stored. Set this list box to **<none>** to prevent the program from storing this information in the project.

- If your project contains a LOCATION variable the **Store file location** option will be available. Select the variable in which you want the file location to be stored. This variable is created upon initially importing your documents by selecting the **Import file location** option when creating your project.
- Select the document variable where the imported documents should be stored.
- Click the **OK** button.

Related Topics:

<u>Creating a project from a list of documents</u> <u>Importing an Existing Data File or Web Service</u> <u>Appending from a Data File</u> Monitoring a File or Folder

Appending from a Data File

The **Append from a data file** function allows you to append cases stored in an external data file to the current project. In order for data to be properly imported, both data files need to share variables with identical names and compatible data types (numerous type conversions are supported). If variables do not exist in the external file or if their types do not match or cannot be converted, the value of these variables will be set to missing. WordStat can import cases stored in the following file formats:

- QDA Miner projects (*.ppj)
- QDA Miner 1.0 to 4.0 project file (*.wpj)
- MS Excel spreadsheets (*.xls or *.xlsx)
- MS Access data files (*.mdb)
- Tab delimited files (*.tab)
- Comma separated value files (*.csv)
- Stata 8 to 15 (*.dta)
- SPSS (*.sav)

To append from a data file:

- Select the Append button and choose FROM A DATA FILE from the adjacent menu.
- QDA Miner will first ask you to select the data file containing the cases to be imported. It will then display a dialog box similar to this one:

able			3
OAp	opend all cases		
) Ap	opend new cases only		
	Common variables:		Key variables:
	AUTHOR		DOCNO
		>	
		_	
		1	OK 🗶 Cancel

If the **Append all cases** option is chosen, QDA Miner will import all cases found in this other file and store them as new cases in the current project.

If the **Append new cases only** option is chosen, QDA Miner will require the identification of one or several key variables that will be used to differentiate already existing cases that will be omitted from the append. Only new cases that will be appended to the current project.

To identify a key variable:

- Select the key variable in the Common variables list box
- Click the button to move it to the **Key variables** list box. If a single key variable is identified, then all cases matching existing values in the current project will be ignored while all the other ones containing new values will be imported as new cases. If several key variables are chosen, a case will be considered as existing and be ignored only if it matches values on ALL the key variables.
- Click OK.

Monitoring a File or Folder

QDA Miner allows you to monitor a file, folder or online service for recent additions. You can configure your project to monitor a specific folder for newly added documents and import the documents. You can also monitor changes to data files or online services such as web surveys, email accounts or reference manager tools, allowing you to import any new cases, responses, emails, Tweets etc.

• To access this function select **Append** and choose **MONITORING** from the adjacent menu. A **Monitoring Settings** dialog box similar to the one below will appear.

Monitoring Settings				-		×
O No Monitoring						
Monitor new documents or image	s stored in designated fol	der				
Folder to monitor:	C: \Users \Amanda \Docum	nents\My Prova	lis Research Projects\Wor	kshop\GOP		8
Store in variable:	DOCUMENT	~	Store file name:	LOCATION		~
After importation:	Leave files in place	Y				
Monitor data files for changes:						
Monitor from original online source	2;					
				🗸 Ok	×c	ancel

There are three types of monitoring that can be performed: monitoring a folder for new documents that have been added, monitoring individual data files for changes and monitoring from an online source. The options available to you are dependent on the data source of your project. Options that are not available to you will be grayed out.

· Select the radio button that applies your project.

Monitoring a Folder

Adding relevant documents to an existing project normally requires you to open the project and perform some operations to append the newly found documents in the project. An easier way to import documents is to instruct QDA Miner to monitor a specific folder for any newly stored documents. This allows you to quickly add new documents to existing projects, by simply saving or copying the desired documents to the folder being monitored, without the need to open them or even run QDA Miner. By specifying a folder on a cloud storage service such as Dropbox, Google Drive, or OneDrive, you can even add documents to a QDA Miner project located on a remote computer.

To monitor a folder for new documents or images:

• Click on the 😂 button to the right of the Folder to monitor field. A Document storage folder dialog box appears.

> ≝ □	ocum)	ents	1
2	Ama	anda QDA Miner 5 Training	
	Carr	ntasia Studio	
	My	HelpAndManual Projects	
~	My	Provalis Research Projects	
	1	Dictionaries	
		Models	
		Samples	
		Workshop	
	>	Candidates	
		Election 2008	
	>	GOP	
	1	Paintings	
	- 5	Science Education	
	1	Science Education	`

- Select the folder that you would like to monitor or click the **Make New Folder** button to create a new one, and then click **OK**
- Set the Store In variable list box to the variable in which you would like to store the document.
- Set the **Store file name** list box to the variable in which you would like to store the name of the file.
- Using the After Importation list box, choose what you would like to happen once the new files are appended. You can delete files, move files to another folder or leave files where they are once they have been appended. If you choose to move a file once it is appended another field will appear. Click on the 😂 button and select folder to which you would like the newly appended file to be moved.
- Click **OK**. Upon reopening the project, you will be notified if any new documents have been added to the folder you are monitoring.

Monitoring a Data File

When a project has been created from a data file such as Excel, MS Access, Endnote, etc., QDA Miner stores the name and location of this file as well as its size and the date it was last modified. It can be set to monitor changes to the file and, if needed, import new cases from the files. When such an operation is available, the **Monitoring data file for changes** option becomes enabled and may be chosen.

To monitor changes to individual data files:

The original source file path should appear next to this option. If the source file is not on your PC this option will not be available to you.

• Click the **Check Now** button to see if the file has been modified since you created your project or last appended from the data file. If the file has been modified a dialog box asking you if you want to append the new cases will appear.

The file REFERENCES 3 MONITOR	ING TEST, XLS has be	en modified. Do you want to a	ppend new cases (if any)?
	✓ Yes	S No	

• Select Yes and a Table dialog box will appear.

) Append all cases		
) Append new cases only		
Common variables:	key variables	
AUTHOR DOCNO PUBYEAR		
	и ок и жа	ancel

• Select whether you want to Append all cases or Append new cases only. If you choose to append all cases, this may result in projects containing duplicates, as some of the cases may have been imported previously. If you choose

to **append new cases only**, only new rows that were added to the data file since the last importation will be appended. If you choose to append only new cases you must choose the common variable for the cases to be matched on.

• Click **OK** and the cases will be added to your project.

Monitoring an Online Source

When a project has been created from an online source such as a web survey platform, an online reference management tool or from Email services, QDA Miner stores with the project, the connection information needed to retrieve additional cases. When this information is available, the **Monitoring original online source** option becomes enabled and may be selected.

To monitor from an online source:

- The online source will appear next to this option.
- Click the **Check Now** button to see if new data has been collected by the online source. If so, a dialog box, similar to the one below, asking if you want to append the new cases will appear.

Do you want to in	nport 472 new cases

• Select Yes and the cases will be appended to the project.

Note: If your online source is **Twitter**, **RSS**, **Facebook** or **Reddit** once the data has been imported the **Posts** column in the **Web Collector** for the query in question will be reset to zero. Please see the <u>Web Collector</u> section for more information.

For more information on projects created from online sources please see:

Importing from Web Surveys Importing from RSS Importing from Twitter Importing from email servers Importing bibliographic references

Deleting Cases

To delete the current case:

- On the Data Editor tab, select the case you would like to delete by clicking it.
- Select the **Delete** button on the left of the screen and choose **CURRENT CASE** from the adjacent menu.

To delete multiple cases:

• Select the **Delete** button on the left of the **Data Editor** tab and choose **MULTIPLE CASES** from the adjacent menu. A dialog box appears showing all cases in the project.

Delete i cuse(s)	÷	
ABC - December 11, 2011		
Bloomberg - October 11, 2011		
CBS - November 12, 2011		
CNN - June 13, 2011		
CNN - November 22. 2011		
CNN - October 18, 2011		
CNN - Sept 12, 2011		
F0X - August 11, 2011		
F0X - December 15, 2011		
🔲 FOX - May 5, 2011		
F0X - Sept 22, 2011		
MSNBC - November 9, 2011		
MSNBC - Sept 7, 2011		
ABC - Jan 7, 2012		
CNN - Jan 19, 2012		
CNN - Jan 26, 2012		
🔲 FOX - Jan 16, 2012		
MSNBC - Jan 23, 2012		
MSNBC - Jan 8, 2012		
CNN - Feb 22, 2012		

- Select the checkbox beside the cases you want to delete. By default, the current case is automatically selected.
- Once you have selected all cases you want to delete, click the OK button.

Identifying Duplicate Cases

Duplicate cases may occur for many reasons, including data importation, data entry or data management errors, or may naturally occur in Twitter and other online sources (ads, news lines, press releases etc.). Duplication may affect WordStat's ability to extract relevant features (topics, phrases, etc.). Use the **Duplicates** button to identify, tag, select, filter out or delete duplicate cases.

To identify duplicate cases:

• Select the **Duplicate** button on the left of the Data Editor tab. A dialog box similar to the one below will appear.

tch cases on:	Action
FILE LOCATION CREATED DOCUMENT	 Delete duplicates Review cases before deletion Create an indicator variable Name: DUPLICATE Create an index variable Name: SEQUENCE Filter data: Primary & Duplicate Cases Duplicate Cases Only Unique and Primary Cases

There are two actions that you can perform. You can delete duplicate cases, with or without reviewing them first, or you can create an indicator variable that identifies duplicate, primary or unique cases. Selecting the second option also offers you possibility of creating an index variable for easy identification of primary and duplicated cases, and to filter cases based on the new indicator variable.

To delete duplicate cases:

- Select the items from the **Match cases on** list box on the left side of the dialog box. These items are the variables that are compared to determine if a case is a duplicate.
- Choose the **Delete duplicates** option from the **Action** options.
- If you wish to immediately delete duplicate cases without reviewing them click **OK**. If you wish to **review the cases before deletion** check the corresponding box and click **OK**.
- A List of duplicate cases dialog box, like the one below will appear.

List of I	luplicate cases	- D
🖌 Keep		
ASE NO.		
10126	Eagan, MN (US)	
10386	Eagan, MN (US)	
10085	Eagan, MN (US)	
10069	Eagan, MN (US)	
10131	Eagan, MN (US)	
8114	Eagan, MN (US)	
10219	Eagan, MN (US)	
10152	Eagan, MN (US)	
10199	Eagan, MN (US)	
3212	Los Alamitos, CA (US)	
16482	West Melbourne, FL (US)	
19325	Union City, CA (US)	
10436	Norfolk, VA (US)	
10002	Norfolk, VA (US)	
10263	Norfolk, VA (US)	
10016	Norfolk, VA (US)	
11135	Chennai (India)	
7455	Camden, AR (US)	

- Select cases you wish to retain, if any, and click Keep.
- When you are finished reviewing the list click OK.
- You will be asked to confirm that you want to delete duplicate case. If so, click Yes.

To create an indicator variable that identifies duplicate cases:

- Select the items from the **Match cases on** checkbox list on the left side of the dialog box. These items are the variables that are compared to determine if a case is a duplicate.
- Choose the Create an indicator variable option in the Action options and enter a Name for your new variable. The
 name of this variable should be descriptive so it is easily identifiable. We have chosen DUPLICATE as it clearly
 describes the information in the variable. For each case, this variable will have a value of unique, primary or
 duplicate. This allows you to easily identify duplicates.

When you create an indicator variable you have the options of **Creating an index variable** and **Filtering the data**.

To create an index variable:

- Select the Create an index variable checkbox if you would like to assign a number to each unique case. Assigning
 an index variable can prove useful when using the <u>Case Descriptor</u> function. This will allow you to sort by sequence
 or group items on this variable to easily identify which cases have duplicates.
- Enter a **name** for the index variable in the corresponding field. The name of this variable should be descriptive so it is easily identifiable. We have chosen **SEQUENCE** in the example above as it clearly describes the information in the variable.

To filter cases on the indicator variable:

• Select the **filter data** checkbox to filter the data to upon the creation of an indicator variable. Choose the filtering parameters from the three options:

Primary & Duplicate Cases: Shows which cases have duplicates and how many.

Duplicate Cases Only: Allows you to see only the duplicates, without the primary cases.

Unique and Primary Cases: Filters the duplicates from the dataset.

• Click OK.

For more information on filtering data see Filtering Cases.

For more information on grouping and descriptors see Setting the Case Descriptor and Grouping.

Setting the Case Descriptor

When retrieving results from a project using the keyword retrieval or report features of WordStat, each hit will be associated with a text description of the case it come from. For this reason, it is recommended to adjust these descriptors so that cases are easily identifiable. To build such a descriptor, you can use existing variable values as well as text strings.

To adjust the descriptors :

• Select the **Descriptor** button at the left of the **Data Editor** tab. The following dialog box will appear:

ariables: CASENUM		
FILE LOCATION		
CREATED		
Description St	ripa: /ETLE]	
Description St	ring: {FILE}	

This dialog box allows you to specify a label that will be used to describe each case. The label may be changed by editing the text in the **Description String** edit box. To insert the value stored in a specific variable into the description, simply enter the variable name in uppercase letters and enclose this name between braces. Alternatively, you can insert a variable name at the current cursor location by clicking the corresponding item in the **Variables** list located just above the edit box.

If you enter the following string:

```
{GENDER} subject - {AGE} years old
```

The {GENDER} and {AGE} strings will be replaced with their corresponding value for this specific case. If the current case contains information about a seventeen-year-old male, the above string will be displayed as:

Male subject - 17 years old

It is also possible to insert the following string:

{CASENUM}

This string will display a unique case number, representing the physical order of this case in the project file.

Filtering Cases

The **Filter** button allows you to temporarily select cases according to logical conditions. You can use this command to restrict analysis to a subsample of cases or to temporarily exclude some subjects. The filtering condition may consist of a simple expression, or include up to four expressions joined by logical operators (AND, OR).

To filter cases:

• Select the **Filter** button on the left of the Data Editor tab. A dialog like the one below appears.

		CANDIDATE	~	equals	~	[Bradley:Bush:Buchannan]	~
AND	×	TOPIC	~	equals	~	[Announcement]	~
AND	~		~		~		Y
AND	-		-		-		-

The following table shows the various operators available for each data type

AVAILABLE OPERATORS
Equals Does not equal Is empty Is not empty
Equals Does not equal Is greater than Is lesser than Is greater than or equal to Is lesser than or equal to Is empty Is not empty
Is true Is false
Contains Does not contain

	Is empty Is not empty
DOCUMENT	Is empty Is not empty Is coded Is uncoded
IMAGE	Is empty Is not empty

• Once a filtering expression has been entered, you can apply the filter and leave this dialog box by clicking the **Apply** button. If the filter expression is invalid, a message will appear and exiting from the dialog box will not occur.

To temporarily deactivate the current filter expression, click the **Ignore** button. The filter expression will be kept in memory and may be reactivated by selecting the **Filter button** again and clicking **Apply**.

To store the filtering expression, click the 🔐 button and specify the name under which the filtering options will be saved.

To retrieve a previously saved filter, click the
button and select from the displayed list the name of the filter you would like to retrieve.

To exit from the dialog box and restore the previously active filtering expression, click the **Close** button.

Editing Project Data

To edit project data:

- Check the Read only checkbox on the left side of the screen
- Select the cell you would like to edit.
- Populate the cell with new data.
- Select the Save button at the top of the screen.

Adding New Variables

The Add button on the Structure tab is used to add new variables to the project file.

To add new variables:

• Select the Add button. This command displays a dialog box similar to the one below.

Variable Definition	1				×
Identification:					
Variable name:	PARTY				
Description:	Political party the candid	ate belongs to			
Data type:	Boolean	×			
		Add 🕅 Remove			
VARIABLE NAME	DATA TYPE	DESCRIPTION		1	
CANDIDATE	Nominal/Ordinal	The person who has made known his or her intent	ion to seek, or	Ci.	
DATESPEECH	Date	Date the speech wa given			Ŷ
					+

- Specify the name of the new variables and set the variable types, sizes and descriptions.
- Click the OK button to create these new variables and add them to the end of the current data set.

Deleting Existing Variables

To delete one or more variables:

• Select the Delete button on the Structure tab. The following dialog box will appear.

elete variables		×
Existing variables:	Variables to	delete:
AGE		
SELECTION		
	>	
1	OK X Cancel	1

- Highlight the names of the variables that you want to delete and click the 🔛 button to move them to the Variables to delete list box.
- To delete successive variables, click the first variable, drag the mouse cursor down the list to highlight multiple variables, and then click the button.
- To remove a variable from the list of variables to delete, select the name of this variable in the Variables to delete list box and click the button.
- Click the OK button to delete all selected variables.

Reordering Variables

To reorder variables:

• Select the **Reorder** button on the **Structure** tab. The following dialog box will appear.

FILE LOCATION	🗸 ок
CREATED DOCUMENT	X Cancel
TOPIC	? Help
	+
	+1

- Select the variable you would like to move.
- Click on the up or down arrow buttons until the variable is in the desired position.
- Click the OK button.

Transformation

Certain operations require specific variable types or transformations. The **Transformation** function allows you to alter variables in a number of ways, and either override the current variable or save the transformation as a new variable.

Wordstat provides numerous functions for <u>changing variable types</u> allowing you to change the data type of the selected variable.

The <u>Recode</u> function allows you to apply multiple changes to the values of numeric, categorical or string variables or to create new variables based on a new grouping of the values of an existing variable.

The <u>Transform Document</u> function allows you to transform document variables stored in your project using any one of the preprocessing routines available in WordStat.

The <u>Binning</u> function allows you group numeric or floating point variable variables into different categories of nominal variables that consists of a range of numbers.

The <u>Compute Number</u> command allows transformations using various functions on one or more variables. WordStat offers more than 50 operations and functions including numerical operators, trigonometric transformations (cos, sin, log, etc.), statistical functions (mean, minimum, maximum across variables or cases, etc.), date and random number operations. Conditional transformation can also be performed using an IF-THEN-ELSE logical structure.

Changing Variable Types

Some operations require specific variable types. For example, a change in data type may be required when you need to add decimal values to a numeric variable created as an integer. A date or a numeric variable may have been imported as a string and may need to be transformed into the correct data type for processing.

WordStat offers the ability to change the type of an existing variable or to create a new variable containing the values of the existing variable stored in a different data type. The following transformations are currently supported:

- Float -> Integer
- Float -> String
- Integer -> Float
- Integer -> String
- Integer -> Nominal
- Nominal -> String
- Nominal -> Date
- Nominal -> DateTime
- String -> Nominal
- String -> Document (codable)
- String -> Date
- String -> DateTime
- String -> Integer
- String -> Float
- Document -> String
- Date -> String
- DateTime -> String

To change the type of a specific variable:

- Select the variable you want to transform by clicking on it on the Structure tab.
- Select the Transform button. The allowed transformations will be listed in a submenu.
- Choose the desired data-type transformation. For most transformations, a dialog box similar to the one shown below will appear.

Transform AGE into Nominal/Ordinal	×
Overwrite existing variable	
O Store into a new variable named:	
🖌 OK 🛛 🗶 Cancel	

- To change the type of the selected variable, choose Overwrite existing variable and click the OK button.
- To copy the values of the selected variable into a variable of the new data type, choose **Store into a new variable named**, type the name of the new variable in the edit box and click the **OK** button. If the new variable name already exists, you will be prompted to confirm the overwriting of this variable.

Transforming Strings Into Dates

WordStat can extract date information from a string and store the date into a new variable. It can recognize in a string various date formats such as 01/05/2011, January 5th, 2011, 5 JAN 2011, or 2011-01-05 and will ignore surrounding text. When this transformation is performed, a dialog box similar to the one shown below will appear:

Destination:				
O Overwrite ex	isting variab	le		
Store into a	new variable	e named: Pl	UBDATE]	
Date format: (MDY	ODMY	OYMD	

- To change the selected string variable into a date, choose **Overwrite existing variable**. To store the date values that are extracted from the selected variable into a new date variable, choose **Store into a new variable named** and type the name of the new variable in the edit box and then click the **OK** button. If the new variable name already exists, you will be prompted to confirm the overwriting of this variable.
- The **Date format** option allows you to choose the specific sequence that is used for displaying dates. QDA Miner will recognize only one date sequence format at a time, either month/day/ year (MDY), day/month/year (DMY) or year/month/day (YMD). Select the sequence used in the most common date format. If no date with a specific format is found for a case, the date variable will remain empty.

Transforming Documents Into Strings

WordStat can transform a document variable into a string with a maximum length of 1000 characters. Several situations may justify such a transformation. For example, you cannot filter cases based on the content of document variables, but you may perform such a filter on string variables. Also, case descriptors cannot include text stored in documents, while string variables may be used as part of the case descriptors. When this transformation is requested, a dialog box similar to the one shown below will appear:

Destination:	
Overwrite existing variable	
O Store into a new variable named:	
String length: Set to the longest string (maximum 1000) Maximum length: 10	
VOK X Cancel	

- To change the selected document variable into a string type, choose **Overwrite existing variable**. To store the date values that are extracted from the selected variable into a new date variable, choose **Store into a new variable named** and type the name of the new variable in the edit box. If the new variable name already exists, you will be prompted to confirm the overwriting of this variable.
- WordStat sets the size of the new string variable automatically by selecting the **Set to the longest string** option. WordStat will go through all cases and set the size to the longest string encountered. If a document is longer than 255 characters, the variable size will be set to 255 and the document will be truncated to this length. You can also set the size to a fixed length by setting the **Maximum length** option to the desired length. All documents longer than the set length will be truncated to this length.

Recoding Values of a Variable

The **RECODE** command, available when you select the **Transform** button, provides an easy way to apply multiple changes to the values of numeric, categorical or string variables or to create new variables based on a new grouping of the values of an existing variable.

To recode a variable:

• Select the **Transform** button and choose the **RECODE** command from the adjacent menu. A recoding dialog box will appear with the following regions.

Recode CANDIDATE	+		×
Destination		_	_
Overwrite existing variable			
Store into a new variable named: PARTY	1		
Existing values:			
Bradley Buchannan			
Bush			
Forbes			
<missing></missing>			
2 · · · · · · · · · · · · · · · · · · ·			
Recode as:		v] 10	H
Gore -> Democratic			-
Sole y Senderate			
			_
	-		
	OK	YC	ancel

- At the top of the dialog box the **Destination** section allows you to specify where the computation results will be stored. To transform the values of the selected variable, choose the **Overwrite existing variable** option. To keep the selected variable intact and store the result in another variable, select the **Store in a new variable named** option and type in the name of the new target variable in the edit box. If you enter the name of an existing variable, the program will ask you if you want to replace its values with those produced by the change. If the variable does not exist, the program will ask you to confirm the creation of the new variable.
- The Existing values list box displays the list of all values found in the selected variable. Select one or several items from this list and recode them by either typing a new value in the Recode as edit box or selecting another existing value from its drop-down list. To confirm the recoding, click the Add button. The value transformations to be performed will be listed in the recoding list box. Repeat this operation until all desired transformations have been specified. Untransformed values remaining in the Existing values list box will remain unchanged or will be copied to the destination variable if you selected to store values in another variable. QDA Miner can automatically detect when a nominal variable is needed and perform the variable type transformation and recoding in a single operation. An example of when this may be useful is when you recode ages to age ranges.

The special keyword **<missing>** is used to replace specific values with an empty cell and further treat those values as missing. Cases with missing values are ignored when performing some analysis involving these variables.

To remove a transformation from the list of recoding, select it and press the <Delete> keyboard key.

• Once a valid destination and a recoding list have been entered, perform the recoding by clicking the **OK** button.

Transform Document

Preprocessing routines allow for the live transformation of text prior to the standard text analysis processes provided by WordStat (for more information on this feature see <u>Preprocessing</u>). Live transformation typically results in an increase in

the processing time which may be significant. For example, performing part-of-speech tagging followed by a lemmatization using the NLTK Python library could add several minutes to the analysis of a small text collection. Rather than applying this transformation repeatedly on the original documents, you can apply a transformation only once and either replace the original document with the transformed one, or store the result of the transformation into a new document variable.

The **Transform document** function allows you to transform documents stored in your project using any one of the preprocessing routines available in WordStat. To apply a transformation on a document variable:

- Select the document variable on which to apply the transformation.
- Click the **Transform** button and choose the **TRANSFORM DOCUMENT** command from the adjacent menu. A dialog box similar to this one will appear:.

DocTransformDIg	-	
Destination:		
Overwrite existing variable		
O Store into a new variable named:		
String length:		
oung engen		
Transformation: Porter Stemmer (DLL version)	Edit	New
OK Cancel		

- The **Destination** section of this dialog box allows you to choose if you would like to overwrite the selected variable with the transformed document or to store the transformed text in a new document variable. If you choose this second option, you will need to type the name of the new variable in the edit box.
- The Transformation option allows to you choose an existing preprocessing routine from a drop down list. You can
 also Edit an existing routine or add a New one by selecting the corresponding button. For more information on
 creating a preprocessing routine please see <u>Preprocessing</u>.
- Once your options have been chosen select he **OK** button.

Binning Numerical Variables

The **Binning** function allows you to "bin" or group variables into different categories. This function transforms a numeric or floating point variable to a nominal variable that consists of a range of numbers. For example, if you have a numeric variable with values ranging from 1 to 100, and you want to transform this variable into only 5 different classes, you can use this function to create 5 bins: 1 to 20, 21 to 40, 41 to 60, 61 to 80 and 81 to 100

To perform binning on a numerical variable:

- On the Data tab move to the Structure tab
- Select a numeric or floating point variable
- Select the Transform button.
- Scroll down to **BINNING** in the adjacent menu. A dialog box like the one below opens.

Overwrite existing variable		
O Store into a new variable na	med:	
nning;		
Number of bins: 5 🚔		
Method: () Foual interv	als Start: 0 Inc	rement: 40
O Equal freque	encies	
Boundaries: Todude upo	er bounderv (<-) O Evolude u	ppor boundary (2)
boondariest @ Include upp		pper boundary (<)
	Denerate	
Label	Upper boundary	Frequency
0 - 40	40	22323
>40 - 80	80	1
>80 - 120	120	3
>120 - 160	160	0
>160 - 200	200	1

At the top of the dialog box the **Destination** section allows you to specify where the computation results will be stored. To transform the values of the selected variable, choose the **Overwrite existing variable** option. To keep the selected variable intact and store the result in another variable, select the **Store in a new variable named** option and type in the name of the new target variable in the edit box. If you enter the name of an existing variable, the program will ask you if you want to replace its values with those produced by the change. If the variable does not exist, the program will ask you to confirm the creation of the new variable.

The **Binning** section allows you to choose the **Number of bins** and the binning **Method**. You have two **Method** available:.

- 1. Equal Intervals means the length of the interval for each bin will be the same. In the example above, we have an interval of 20 for each bin. With this option you can set the Start and Increment. These values are automatically generated but you have the choice to change them by typing the chosen value in the field.
- 2. The **Equal Frequencies** option automatically calculates the intervals of each bin and attempt to produce bin with a similar number of cases. The **Boundaries** option determines if the upper boundary is included in the bin or not.
- Decide if you want to **Overwrite the existing variable** or **Store in a new variable**. If you chose the latter option, type the name of your new variable.
- The table displays the currently set bins.
- Choose the number of bins.
- Set your **Method** and your **Boundaries**.

- Select the Generate button.
- Click **OK** to perform the transformation.

Performing Complex Numeric Transformation

The **Compute Number** command allows transformations using various functions on one or more variables. WordStat offers more than 50 operations and functions including numerical operators, trigonometric transformations (cos, sin, log, etc.), statistical functions (mean, minimum, maximum across variables or cases, etc.), date and random number operations. Conditional transformation can also be performed using an IF-THEN-ELSE logical structure. When this command is selected, a dialog box similar to this one is displayed.

Numeric	~		101		2	3	1	ACOS (numeric)	1
ID							1	COS (numeric)	
AGEGROU	P	200		4	5	6	*	CSC (numeric)	
LEADERSH	IP	-	Lun	7	8	9	÷	EXP (numeric)	
SENSEHUN	IOR	1.0	APPL	0			+	LAG (variable)	1.1
YRS_EXPE	RI	2.=	CH-		_			LOG (numeric)	
INPS		-	11101					MAX (varlist) MEAN (varlist)	
								MIN (varlist)	
								MOD10 (numeric)	
								MONTH (date) NORMAL (numeric)	
								RND (numeric) SEC (numeric)	
ansformat	on								
Store in:	FUTURSPEND	8	Condi	tional trai	nsform	ation	C	Integer Floating point	
If:	(NPS > 8) ANI	D YEAR (SL	BSCRIPTION	1)					
Then:	SUM(SPENDIN	IG1SPEN	DING21)+SPI	ENDING2	5				š
	12 33 25 26 20 2		CNIDTNICOA)						

- Store In- This option allows you to specify the target variable where the computation results will be stored. If the target variable already exists, WordStat asks you whether you want to replace its values with those produced by the transformation. If the variable does not exist, it is appended at the right end of the data file. The new numerical variable can be either an **Integer** or a **floating point** numerical variable. If the computed numerical values contains decimal and is stored into a integer data type, the value is rounded to the nearest integer. The TRUNC function (see below) may however be used to store only the integer part of the computed numerical value.
- **Conditional transformation** This option allows you to select whether the same transformation will be performed on all active records or whether different transformations should be performed on specific records according to a logical condition.
- Formula When the Conditional Transformation option is deactivated, this edit box is used to store the transformation expression that will be applied to all active records.
- If When a conditional transformation is requested, this option contains a logical expression up to 250 characters long specifying criteria for cases on which to apply a specific transformation. This expression must be evaluated as true or false. It can be a simple logical expression including the following 3 components: a variable name, a relational operator, and a variable name or a numeric constant. It can also be a complex expression composed of many simple expressions related with logical operators (AND, OR, XOR). Parentheses can be used to specify the order in which expression are evaluated. In the following example, the last two conditions are assessed before the first one.

(GROUP = 1) OR ((GROUP = 2) AND (AGE > 30))

For a list of all available functions and operations available for logical expressions see xBase functions.

Then / Else - If the expression is true the transformation expression in the THEN field is computed, if not the expression in the ELSE field is used.

Expression Operators and Rules

ARITHMETIC OPERATOR

- + Addition
- Subtraction
- * Multiplication
- / Division
- Exponentiation

CONSTANT

PI 3.1415926535897932385

MISSING VALUE FUNCTIONS

SYSMIS	Blank cells
MISSING	Variable missing values

NUMERIC AND TRIGONOMETRIC FUNCTIONS

Syntax FUNCTION (value, variable or expression)

ABS	Absolute
ACOS	ArcCosine
ASIN	ArcSine
ATAN	Arctangent
CSC	Cosecant
COS	Cosine
EXP	Exponential
FACT	Factorial
LN	Natural logarithm
LOG	Base-10 logarithm
MOD10	Modulus
RND	Round
SEC	Secant
SQRT	Square root
SQR	Square
SIN	Sine
TAN	Tangent
TRUNC	Truncate

STATISTICAL FUNCTIONS (ACROSS VARIABLES)

Syntax FUNCTION (Var [Var Var..Var])

MEAN	Mean
COUNT	Count (missing value excluded)
SUM	Sum

SIDEV	Standard	deviation

- VAR Variance
- MIN Minimum
- MAX Maximum

STATISTICAL FUNCTIONS (ACROSS CASES)

Syntax FUNCTION (Variable)

VMEAN	Mean
VCOUNT	Count
VSUM	Sum
VSTDEV	Standard deviation
VVAR	Variance
VMIN	Minimum
VMAX	Maximum
ZSCORE	Normalized score
LAG	Lag

RANDOM NUMBER FUNCTIONS

Syntax FUNCTION (value, variable or expression)

NORMAL	Normal pseudo-random number with mean of 0 and standard deviation of X
UNIFORM	Uniform pseudo-random number between 0 and X

DATE FUNCTIONS

YRMODA	(yy,mm,dd) convert 3 values, variables or expressions into a julian date
YEAR	(julian date) return the year of a julian date
MONTH	(julian date) return the month of a julian date
DAY	(julian date) return the day of a julian date
TODAY	return current julian date

xBase functions

ALIAS()

Returns the Alias name of the current workarea as a string.

ALLTRIM(String)

Trims both leading and trailing spaces from a string. The string may be derived from any valid xBase expression.

```
ALLTRIM(" Provalis ") returns 'Provalis'.
```

AT(SearchString, TargetString)

Determine whether a search string is contained within a target. If found, the function returns the position of the search string within the target string (relative to 1). If not found, the function returns 0 (zero).

```
AT("gh", "defghij") returns 4.
```

CHR(Val)

Converts a decimal value to its ASCII equivalent.

CHR(83) returns 'S'

CTOD(String)

Converts a character string into an xBase date. The string must be formatted according to the Windows date format settings.

DATE()

Returns the system date (today). Use DTOC(DATE()) to retrieve today's date formatted according to the Windows settings.

DAY(Date)

Returns the day portion of an xBase date as an integer.

DELETED()

Returns True if the record is deleted and False if not deleted.

DESCEND(String)

An xBase function that inverts a key value using 2's complement arithmetic. The result of the operation is the arithmetic inverse of the key value. When inverted keys are sorted in ascending sequence, the result is in descending order. A filter expression could be

```
DESCEND(DTOS(billdate)) + CUSTNO
```

DTOC(Date)

Converts an xBase date into a character string formatted according to the Windows settings. For example, if the date format was American and the date field contained March 21, 1995, DTOC (datefield) would return '03/21/1995'.

DTOS(Date)

Converts an xBase date into a string formatted according to standard xBase storage conventions (CCYYMMDD). For example, December 21, 1993 would be returned as '19931221'. Indexes that contain date elements should use the DTOS() function, which naturally collates into oldest date first.

EMPTY(Field)

Reports the empty status of any xBase field. Character and date fields are empty if they consist entirely of spaces. Numeric fields are empty if they evaluate to zero. Logical fields are empty if they evaluate to False.

Memo fields that contain no reference to a memo block in the associated memo file are empty.

IF(Logical, True Result, False Result)

This is the immediate if function. If the Logical expression is true, return the True result, otherwise return the False result. The types of the True Result and the False Result must be the same (i.e., both numeric, or both strings, etc.) The logical expression must of course evaluate as True or False.

IF(DATE() - CTOD("12/31/93") > 0,"This Year", "Last Year")

IIF(Logical, True Result, False Result)

Supported exactly like IF() as noted above.

INDEXKEY()

Returns the current index key as a string. (Same as ORDKEY()).

LEFT(String, Length)

Returns the leftmost characters of the expression for the defined length.

LEFT("xyzabc", 3) returns 'xyz'.

LEN(Expression)

Returns the length of the expression result as an integer.

LOWER(String)

Converts the string expression into lower case.

MONTH(Date)

Returns the month portion of an xBase date as an integer.

ORDER()

Returns the current index order as an integer.

ORDKEY()

Returns the current index key as a string. (Same as INDEXKEY())

PADC(String, Length, Character)

Centers the passed string between a number of the passed character to make the string the specified length.

'[' + PADC("Scott", 9 ,"-") + ']' returns '[--Scott--]'.

PADL(String, Length, Character)

Pads the passed string to the specified length with the specified characters. If the string is longer than the value specified by Length, the string is truncated to this length.

```
'[' + PADL("Scott", 8, "*" ) + ']' returns '[***Scott]'.
'[' + PADL("Loren Scott", 8, " " ) + ']' returns '[Loren Sc]'.
```

PADR(String, Length, Character)

Pads the passed string to the specified length using the specified character. If the string is longer than the value specified by Length, the string is truncated to this length.

'[' + PADR("Scott", 8, " ") + ']' returns '[Scott]'.
'[' + PADR("Loren Scott", 8, " ") + ']' returns '[Loren Sc]'.

RAT(SearchString, TargetString)

Determine whether a search string is contained within a target, starting from the right side of the target string. If found, the function returns the position of the search string within the target string (relative to 1). If not found, the function returns 0 (zero).

RAT("ab", "abzaba") returns 4.

RECCOUNT()

Returns the number of records in the table as a long integer.

RECNO()

Returns the current physical record number as a long integer.

RIGHT(String, Length)

Returns the rightmost characters of the expression for the defined length. RIGHT("xyzabc", 3) returns

'abc'.

SELECT()

Returns the workarea number for the current workarea as a long integer.

SPACE(Length)

Returns a string consisting entirely of spaces for the defined length.

STOD(String)

The inverse of DTOS(). STOD() converts a string formatted according to standard xBase storage conventions (CCYYMMDD) to an xBase Date formatted according to the Windows settings.

STR(Number, Length, Decimals)

Converts a number into a right justified string with decimals digits following the decimal point. The total length of the string is defined by the length parameter. STR(RECNO(), 5, 0) is a common indexing element that ensures creation of unique keys if appended to another field element.

An index key using this expression could be built with NAME + STR(RECNO(), 5, 0)

If the decimals parameter is omitted, the function defaults to zero decimals. If the length parameter is omitted as well, the length of the result is the length of the field.

STRZERO(Number, Length, Decimals)

Converts a number into a, zero-padded right justified string with decimals digits following the decimal point. The total length of the string is defined by the length parameter.

STRZERO(1234, 10, 2) returns '0001234.00'

If the decimals parameter is omitted, the function defaults to zero decimals. If the length parameter is omitted as well, the length of the result is the length of the field.

SUBSTR(String, Start, Length)

Returns a portion of the string expression starting at the defined start location for the defined length.. SUBSTR('xyzabcd', 3, 4) returns 'zabc'.

TIME()

Returns the system time as a string in the form HH:MM:SS.

TRANSFORM(Expression, Picture)

Transform converts strings and numeric values into formatted character strings. The function transforms the result of the first expression in accordance with the second picture string.

The picture string is made up of two parts. The first part is the Function string and it is optional for both strings and numeric values (as long as the second Template string is present).

A character string transformation picture may consist of only a Function string or only a Template or both.

A numeric picture must contain a Template string; the Function string is optional.

A logical value must contain only a Template string with Template characters L or Y.

The Function string consists of a leading @ character followed by one or more formatting characters. If the Function string is present, the @ character must be the first character in the picture string with its formatting characters immediately following and it may not contain spaces.

If a Template string exists as well, it follows the Function string. A single space separates the Function string and the Template string.

Function string characters allowed for numeric values are:

B left justify;
C display CR after positive numbers;
X display DR after negative numbers;
Z blank a zero value;
(encloses negative numbers in parentheses.

Function string characters allowed for strings are:

R inserts unassigned template characters; ! converts all alpha characters to upper case.

The @R Function requires a Template; the ! Function does not.

The Template string describes the format on a character by character basis. The Template string is made up of special characters which have specific results and optional unassigned characters which either replace characters or are inserted in the formatted string depending upon the absence or presence of the @R Function string.

Template assigned characters are as follows:

A,N,X,(,# are place holders and are interchangeable;

L displays logical values as T or F;

Y displays logical values as Y or N;

! converts the corresponding character to upper case;

, (comma) or a space (in Europe) in a numeric template separate the elements of a number;

. (period) or , (comma - in Europe) in a numeric template specify the decimal position;

* fills leading spaces with asterisks in a numeric template;

\$ as the leading character in a numeric template results in a floating dollar sign being placed in front of the formatted number.

Example: Where "phone" is a character field holding a phone number with no formatting characters.

'transform(phone, "@R (###) ###-####")' returns '(909) 699-6776'.

If the formatting characters were actually present in the field, the "@R" function would be omitted

For numeric fields,

'transform(123456.78, "\$9,999,999.99")' returns' \$123,456.78'.

TRIM(String)

Removes trailing spaces from the string expression.

UPPER(String)

Converts the string expression into upper case. Character fields used in index expressions should always be converted to upper case to insure correct collating sequence.

VAL(String)

Converts a string of numeric characters into its equivalent numeric value. The conversion stops at the first nonnumeric character encountered (or the end of the string).

VAL("123ABC") returns a value of 123.

YEAR(Date)

Returns the year portion of an xBase date as an integer.

Editing Variable Properties

WordStat allows you to edit various properties of existing variables in your project. For example, you can attach a short and a long description to a variable, set the number of decimal places for floating-point numerical values, edit values of categorical variables, etc. You can also set an individual variable as "read only" or change its name.

To access the Variable Properties Editor dialog box:

- On the Structure tab, position the cursor on the variable that you want to edit.
- Click the Properties button. This displays the Variable Properties dialog box as shown below.

Variable Properties Editor - CANDIDATE		X
operties Values Missing values		
lame: CANDIDATE Rename Type: Nominal/Ordinal		
Short description:		
Long description:		

Once in the dialog box, you can move through variables by clicking either the is or the button. When viewing or editing the properties of a categorical variable, a second tab appears. You can add new values, as well as edit or delete existing ones on this page (see below).

The Properties Tab

The **Properties** tab of the dialog box offers the following options:

Read only: When selected, this option prevents a variable from being modified. This option is useful to prevent accidental or unauthorized changes to the values of a variable. Setting a variable to "read only" only affects the editing of values in existing records, but still allows users to create new cases and assign values to these new cases.

Type: The type of variable is shown here. If it is a string variable, the number of characters permitted appears beside the variable type in brackets.

Description: This option lets you enter both a short single-line alphanumeric description as well as a detailed description of the variable. The short description is displayed in various locations in the program to remind users of the exact content of this variable.

Rename button: You can use this button to change the name of the current variable. When you click this button, you will be prompted for a new variable name. This new name must not exist in the current data file and should follow the basic rules for valid variable names.

The Values Tab

When viewing the properties of a categorical variable, a second tab is shown. The **Values** tab allows you to add new values, as well as edit or delete existing ones. The **Values** tab allows you to add new values, as well as edit or delete existing ones.
	none> ~	1,		
alue: Gingritch	Edit	Remove	Reorder	
Obama Richardson Cinton Edwards Kucinich Thompson McCain Giuliani Romney Empty Nader McKinney				

To add a new value labels:

- In the Value edit box, enter the string that will be used to describe this new value.
- Click the Add button.

To remove an existing value label:

- In the list box located in the lower half of this page, select the label you want to delete.
- Click the **Remove** button

To edit an existing value label:

- In the list box located in the lower half of the page, select the label you want to edit.
- Click the Edit button.
- Once you have finished editing the label, click the **OK** button to apply the change.

Reordering Value Labels

Ordinal variables assume a clear ordering of values. A good example of such a variable would be responses to a satisfaction questionnaire with values ranging from "Very Unsatisfied" to "Very Satisfied", or an age group variable containing several age ranges going from "18 or less" to "60 or more." To display the values in a proper order in various tables and to be able to apply some of the ordinal statistics available in WordStat, values of these variables should be properly ordered. The **Reorder** button allows you to change the natural order of values of ordinal variables. When this button is clicked, a dialog box like this one will appear:



- To reorder values, simply select the value you would like to move and click the up or down arrow button until the label is in the proper order.
- To confirm the reordering, click **OK**. To return to the original ordering of values, click the **Cancel** button.

Using an Existing Value Labels Definition

In some projects, several variables share the same value labels. For example, a questionnaire may use a common ordinal scale for several questions. Rather than re-entering the same value labels over and over again, WordStat can establish a link between a variable without value labels and an existing one that already contains labels.

- To establish a link, click the down arrow button of the Link to values in list box.
- Then select from the list of variables the one containing the value labels you want to use. These labels will appear in the value list box.
- Once a link has been established, every change made to the value labels list will affect the labels associated with the original variable as well as with all other variables currently linked to this variable.

You may also use this feature to copy value labels from one variable to another. To do this, follow the previous instructions to link the current variable to the one from which you want to copy labels. The labels should appear in the value labels list. Then, remove this link by setting the **Link to values in** option to **none**. You may now edit the newly copied labels without affecting the labels of other variables.

Resizing a String Variable

When you create a string variable you are asked to determine the maximum number of characters the string can contain. If you find it necessary to modify the size of the string, you can do so by using the **Variables Properties Editor**.

- On the Structure tab choose the string variable that you want to edit.
- Click the **Properties** button. This displays the **Variable Properties** dialog box as shown below.

and the second second second second				
roperties Missing values				
Name: AUTHOR Rename	Type: String (98)	Resize		
Short description:				
ong description:				
		I		
		I		
		I		

- The size of the string will be indicated in brackets beside the variable Type.
- Click the Resize button and a dialog box similar to the one below will appear.



- Adjust the number of characters to fit your needs using the **Up** and **Down** arrow buttons.
- Click **OK** when finished.

Variable Statistics

Selecting the **Statistics** button allows you to quickly obtain the frequency and cross-frequency distribution of numerical, categorical, date and short-string variables. The univariate frequency table includes the frequency count for each value of the selected variable as well as the percentage of the count over all cases and over valid cases only. The contingency table describes the distribution of two variables simultaneously by displaying either the frequency or the row, the columns or the total percentages. Several types of charts may be created to illustrate the distribution or cross-frequency distribution of variables. You will find below a list of those charts.

Frequency Distribution

To obtain the frequency distribution of a variable:

• Select the Statistics button on the Structure tab.

Statistics	s - CANDIDA	TE			$\sqrt{\tau}$		×
Frequency	Crosstab						
Variable:	CANDIDATE	~	Search	✓ Total	40	đ	
VALUE	FREQUENCY	TOTAL PERCENT	VALID PER	CENT			
Biden	17	6.9%	6.9%				
Clinton	39	15.9%	15.9%				
Edwards	14	5.7%	5.7%				
Thompson	9	3.7%	3.7%				
Giuliani	15	6.1%	6.1%				
Kucinich	6	2.4%	2.4%				
MCCain	49	20.0%	20.0%				
Obama	68	27.8%	27.8%	6. T			
Richardson	14	5.7%	5.7%				
Romney	14	5.7%	5.7%				
TOTAL	245	100%	100%				

You may obtain a frequency table on any other numerical, categorical, alphanumerical, and date variable by selecting its name in the **Variable** list box and then clicking the **Search** button.

To chart the distribution of a variable:

Clicking the 🔟 button allows you to obtain up to five types of charts to visually display the distribution of specific codes.

Some of the charts are available only for numerical and date variables (histograms and box-&-whiskers plots), while others like the bar charts and pie charts will be available in all situations when the total number of values is less than 100.

- The vertical bar chart is the default chart used to display the frequencies of distinct values of a nominal or ordinal variable. It is especially useful to compare two or more values.
- The horizontal bar chart displays the same information as the vertical bar chart. It is especially useful when the number of values is high and their labels cannot be displayed entirely on the bottom axis.
- The pie chart is useful to display the relative frequency of each value and compare individual values to other values and to the whole. Numerical values displayed in pie charts are always expressed in percentages of either the total frequency or case occurrences.
- The donut chart, similar to the pie chart, is useful to display the relative frequency of each value and compare individual values to other values and to the whole. Donut charts are considered easier to read as you can focus on reading the length of the arcs, rather than comparing the proportions between slices.
- The histogram graphically displays the distribution of a numeric variable. When selected, the program first separates the values into non-overlapping intervals of equal width, and then plots bars that represent the frequencies of each interval.
- The box-&-whiskers plot can be used to examine the distribution of numerical variables. It is especially useful to detect the presence of outliers and asymmetry in the data distribution. The box includes values that fall between the first and the third quartiles (about 50 percent of the values). The line in the middle of the box represents the median value while the whiskers extend to the farthest observations within 1.5 times the interquartile range measured from the nearest quartiles. Values that are situated farther than 1.5 times the interquartile range but within three times this distance are represented by a dot, while values farther than three times the interquartile range from the nearest quartile are represented by the letter X (for extreme).

Joint Distribution of Two Variables

To obtain the joint distribution of two variables:

- Select the Statistics button on the Structure tab, to display the Statistics dialog box (see above).
- Move to the Crosstab tab. The dialog box will be similar to the one below:

Statistics	- CAND	DIDATE									×
Frequency	Crosst	ab									
Tabulate:	CANDID	ATE	~	With: DELIN	VERY	~ (%)				0 8	-
Display:	Row per	rcent	~	☑ Total	With miss	sing values		Search	9		
	2006	Q1-2007	Q2-2007	Q3-2007	Q4-2007	Q1-2008	Q2-2008	Q3-2008	Q4-2008	TOTAL	-
Biden	23.1%	46.2%	23.1%	7.7%						100%	
Clinton	7.7%	2.6%	7.7%	7.7%	23.1%	25.6%	25.6%			100%	
Edwards	14.3%	7.1%	21.4%	21.4%	35.7%					100%	
Thompson			22.2%	44.4%	33.3%					100%	
Giuliani			40.0%	26.7%	26.7%	6.7%				100%	
Kucinich	66.7%	16.7%	16.7%							100%	
MCCain			2.0%	16.3%	12.2%	14.3%	28.6%	20.4%	6.1%	100%	
Obama	2.9%		1.5%	11.8%	11.8%	20.6%	22.1%	22.1%	7.4%	100%	
Richardson			28.6%	42.9%	28.6%					100%	
Romney		35.7%	28.6%	14.3%	14.3%			7.1%		100%	
TOTAL	115%	108%	192%	193%	186%	67%	76%	50%	14%	1001%	

- Select from the **Tabulate** list box the variable you would like to be displayed on the rows of the table.
- Select the variable you would like to be displayed at the top of the table by selecting its name in the With list box.
- In the **Display** list box, select the statistics you would like to be displayed in the table.
- Click the Search button.

To chart the joint distribution of two variables:

Bar charts or line charts are useful for visually comparing the joint distribution of two variables. To produce these types of charts:

- Set the Tabulate, With, and Display options so that the information to be viewed is displayed in the table.
- Click the 🖊 button.

For more information, see bar chart and Pie Chart.

To append a table to the Report Manager:

• Click the 🛄 button. A descriptive title will be provided automatically for the table. To edit this title or to enter a new one, hold down the **Shift** key while clicking this button.

For more information on the Report Manager, see the Report Manager Feature topic.

To exporting retrieved codes to disk:

- Click the 🔚 button. A **Save File** dialog box will appear.
- In the **Save as type** list box select the file format under which you would like to save the table. The following formats are supported: ASCII file (*.TXT); Tab delimited file (*.TAB); Comma delimited file (*.CSV); MS Word (*.DOC); HTML file (*.HTM; *.HTML); XML file (*.XML);Excel spreadsheet file (*.XLS; *.XLSX), and SPSS data file (*.SAV).
- Type a valid filename with the proper file extension.
- Click the **Save** button.
- The table can also be stored in a new project file, where each row becomes a new case and each column is transformed into a variable. This type of project file may be useful for performing a more detailed analysis of coded segments.

To print the table:

Click the button.

Analyze

Once you are satisfied with the content and the structure of the data you have imported into your WordStat project you can begin your analysis. Data can be analyzed in two different modes: **Explore** mode and **Expert** mode. The **Explore** mode lets you see at a glance what is in your data set. It will provide the list of the most frequent words and most frequent phrases, and will extract the most salient topics using topic modeling. Each of those results will allow quick comparisons with up to two numerical, categorical or date variables. For a more comprehensive analysis, run the **Expert** mode. This mode provides, in addition to the basic descriptive tools of the **Explore** mode, additional features such as co-occurrence analysis (clustering, multidimensional scaling, link analysis), crosstabulation (heatmap, correspondence analysis, bubble charts, etc.), machine learning, as well as all the dictionary building features of WordStat (exclusion list, categorization dictionaries, keyword-in-context, etc). It also gives you access to all processing options (stemming, lemmatization, character processing, etc.).

To begin the analysis process :

Select the Analyze button. A Choose Variables dialog bo will appear similar to the one below:

	PARTY
EXPLORE - Extract most f	frequent words, obrases and tonics

The text variables (document and string) are listed on the right side of the dialog. On the left side of the dialog box are listed the categorical, numeric or date variables that you can analyze in relation to the text variables. Two modes of analysis are available, **Explore** and **Expert**.

- Select the text variables by checking the boxes.
- If you want to perform comparisons with numerical, categorical or date variables, select their checkboxes in the In relation with panel. Please note that, by default, cases with a missing value on any of the selected variables are exclude from the text analysis. To include all cases, select the Include cases with missing values option on the Preprocessing tab.
- Select the mode in which you would like to preform your analysis.

To analyze in Explore mode:

• Select Explore. A Project configuration dialog will appear similar to the one below.

Project configuration		~		×
Select the default language:	English	Ŷ		
Number of topics desired:	20		1	ОК

- A default language will be chosen for you. Modify the language by choosing one from the drop-down list.
- The number of topics to be extracted will automatically be calculated for you. You can change this number by clicking on the up or down buttons beside the number field or simply by typing the desired number of topics in the field.
- Select **OK**. The word frequencies will be calculated, and phrases and topics will be extracted. Once the software is done processing you will be taken to the **Topics** tab and you can begin to explore your data.

To analyze in Expert mode:

• Select **Run Expert Mode**. If you are analyzing the project for the first time, a Create Configuration dialog box will appear similar to the one below.

Create Categorization	×
This project does not have an associated categorization process file.	
Would you like to:	
Open an existing categorization process file?	
	Browse:
Create a new categorization process file?	
Default exclusion lis:	
English.exc	V 🔗 Browse
V OK	

Here you have the option of using an existing categorization model file or creating a new one.

To use an exiting categorization model file:

- Choose a categorization file from the drop-down menu or select **Browse** and find the categorization file you wish to use and select it.
- Select OK. WordStat will open and you will be taken to the Text Processing tab to start your analysis.

If you choose to create a new categorization model file:

- Choose an exclusion list from the drop-down or select **Browse** and find the exclusion list you wish to use and select it.
- Select OK. A dialog will open.
- Type a name for your categorization model and a place to save it. Select Save.
- WordStat will open and you will be taken to the **Text Processing** tab to start creating your categorization dictionary, exclusion and substitution lists, and begin analyzing your data.

Related Topics:

For more information on categorization models please see The Text Processing Tab.

The Text Processing Tab

Without further information, WordStat can perform a frequency analysis on each of the words encountered in the chosen document or alphanumeric variables. However, it is also possible to apply various transformations on the words before performing the frequency analysis. The **Text Processing** tab allows you to specify how the textual information will be processed. It contains all the options related to text preprocessing (stemming, lemmatization, etc), text postprocessing (frequency criteria), as well as those involved in the creation and management of exclusion and substitution lists and categorization dictionaries.

🕥 WordStat 9.0 - El	ection 2008 Coded.ppj										×
🚍 🛄 Data 🦂	Text Processing	Frequencies 😗 Extr	action	S Cooccur	rences	Crosstab	Keywo	rd-In-Context	< Clas	ssification	0.
Categorization Mod	el: Brand Personality.wmod	el	~ 🔗 o	Open 🔻	Save	Save as	Publish	New	Delete	1 Options	
🔇 Language 🚳 Pi	reprocessing Substitut	ion Z Exclusion	Categorizati	on 🐮 Po	stprocessi	ng					
Processes:											
Preprocessor:	3 grams & words.ex	p ~	🛃 Edit	🕑 New	n De	lete					
Stemming:	English (porter)	~									
Lemmatization:	English	~	🛃 Edit								
Character recognition	n:										
Case sensitive	Add ch	aracters appearing:	Anywher	e:							
Accept numeric	characters	Embe	dded in word	ds:]					
Text to include or ign	nore;										
Don't process to	ext within braces	Ignore duplicate:	ODocumen	ts Parag	graphs						
Don't process te	ext within brackets	Analyze the first	500 🚔 w	ords.							
Process only te	xt within brackets	Up to the en	d of the para	graph							
Convert Emojis	to text	Ignore URLs (http	, https and f	ф)							
		Ignore speaker de	signations in	transcripts							
Case processing:											
Random sample	e: 1 out of: 5	Weighting variabl	e: None		~						
Include cases v	with missing values										
10/245											

A WordStat categorization model consist of various text processing settings and routines such as lemmatization, exclusion, categorization, etc. They are grouped under six tabs: Language, <u>Preprocessing</u>, <u>Exclusion</u>, <u>Substitution</u>, <u>Categorization</u> and <u>Postprocessing</u>, allowing you control how text will be processed, transformed and reported. All the settings are stored on disk in a single file with a **.WMODEL** file extension. Categorization model files are independent of WordStat projects, allowing you to apply a categorization model to several projects as well as apply, if needed, various categorization models to the same project file.

Initial Setting of the Categorization Model

While models are stored independently of projects files, a project cannot be processed without a categorization model. For this reason, when you attempt to analyze a text analysis project for the first time, you will be asked to either select an existing categorization model or create a new one. A dialog box similar to this one will appear:

Create Ca	ategorization		×
This proje	ect does not have an associated categorization proc	ess file,	
Would yo	u like to:		
	Open an existing categorization process file?		
	golden state killer.wmodel	V 🔗 Brow	se
00	Create a new categorization process file? Default exclusion list:		
	English.exc	🗸 🔁 Brow	se

To use an existing categorization model:

- Select the Open an existing categorization process file radio button.
- Choose a categorization file from the drop-down menu or select **Browse** and find the categorization file you wish to use and select it.
- Select OK. WordStat will open and you will be taken to the Text Processing tab to containing the categorization dictionary, exclusion list and substitution list present in the chosen model. From here you can tailor the lists to suit your current project or start analyzing your data immediately.

To create a new categorization model:

- Select the Create a new categorization process file radio button.
- Choose an exclusion list from the drop-down menu or select **Browse** and find the exclusion list you wish to use and select it. While only one exclusion (stop word) list may be selected, one may later import additional ones.
- Select **OK**. A dialog will open.
- Choose a name for your categorization model and a place to save it. Select Save.
- WordStat will open and you will be taken to the **Text Processing** tab to start creating your categorization dictionary, exclusion and substitution lists, and begin analyzing your data.

While categorization models and project files are independent, WordStat stores in the project file the information about the last categorization model used, so that when this project is reopen, this model will be automatically retrieved.

Using the Categorization Model Toolbar

The toolbar across the top of the Text Processing tab pertains to the selection, creation and publication of categorization models etc. The **Categorization Model** drop-down list box allows you to access previously saved categorization models. Select the arrow on the right and choose the desired model from the drop-down list.

By default, the models available from this list are those stored in the Dictionaries folder, under My Provalis Research

Projects. To open a categorization model stored elsewhere on your computer, use the button to browse through your system and select the proper **.WMODEL** file.

The following table describes the other functions available on this tool bar.

Control: Description:

-	
	Carin
10.0	Save

Press this button to save the currently displayed categorization model to disk.

Save as

Pressing this button allows you to save the categorization model **.WMODEL** under a new name. A Save File dialog box will ask you to type the name and select the location where you would like the model to be saved.

Categorization models stored as **.WMODEL** files can only be used in the WordStat desktop software. To use a categorization model within QDA Miner, from Windows Explorer, in the <u>WordStat</u> <u>Document Explorer</u> utility program, or in any third party application making use of the <u>WordStat</u> <u>Software Developer's kit (SDK)</u>, the categorization model has to be stored on disk in a portable file format. with the **.WCAT** file extension. This file is an encrypted and compressed version of the model file along with additional settings, information and resources relevant for its successful application.

Pressing the Publish button displays a Save File dialog box allowing you to type the name of the file and select the location where it will be saved. However, in order to be used in QDA Miner or from Windows Explorer, the .WCAT file should imperatively be stored in the **Models** folder under your **My Provalis Research Project** folder.

This button allows you to create a new categorization model. You will be asked whether to keep the current exclusion list, the substitutions as well as the categorization dictionary. Select those you would like to use in your new categorization model and then type **OK**. A Save File dialog box will appear, asking you to type the name and select the location where you would like this new model to be saved.



This button allows you to delete the currently selected categorization model.

This button allows you to assign a description to a categorization model. It also allows you lock specific features of your model and prevent either accidental or undesired changes that may affect negatively the accuracy of the categorization process.

Language

The **Language** page allows one to choose which language resources will be available for the analysis, such as the spellchecking dictionaries used for automatic spelling corrections, lemmatization routines, as well as the thesaurus that will be used to provide suggestions. It offers the following options:

Add or remove language resources: WordStat may use various resources (such as spelling dictionaries, thesauri, and lemmatization routines) to process some natural languages more efficiently. By default, WordStat installs resources for processing English documents. Clicking this button allows you to download and install resource files for other languages or remove language resources previously installed.

Active spelling dictionary: WordStat makes use of language dictionaries to spell-check existing textual data and to suggest inflected forms of words found in the user dictionary. This group of options lets you specify which dictionary to use with the current data file.

Active thesaurus: WordStat's **Suggest** feature can use a thesaurus to suggest synonyms of existing words in the text collection or in the user categorization dictionary. This option allows you to select the language of the thesaurus. For more information, see <u>Basic Dictionary-Building Tools</u>

Preprocessing

The following section provides a description of the preprocessing steps involved in the transformation of textual data into keywords or content categories.

WordStat 9.0 - El	lection 2008 Coded.ppj								÷.		×
🗮 🛄 Data 🦂	Text Processing Fr	equencies 🛭 🔁 E	straction	Cooccurre	ences	Crosstab	🏬 Кеуwа	ord-In-Context	: < Clas	ssification	0
Categorization Mod	del: Brand Personality.wmode	U	~ <mark>6</mark> 0	open 🕶 la	Save	Save as	Publish	New	💼 Delete	1 Options	s
🔇 Language 🚳 P	Preprocessing Substitutio	n 🗹 Exclusion	Categorizati	on 🐮 Pos	tprocessi	ng					
Processes:											
Preprocessor:	3 grams & words.exp		🗸 📝 Edit	🔗 New	🏦 De	lete					
Stemming:	English (porter)		~								
Lemmatization:	English		- Z Edit								
Character recognitio	in:										
Case sensitive	Add char	acters appearing:	Anywher	e:	-						
Accept numeric	characters	En	bedded in word	ls:]					
Text to include or ig	nore;										
Don't process t	text within braces	Ignore duplicat	e: O Documen	ts Paragr	aphs						
Don't process t	ext within brackets	Analyze the first	t 500 🚔 w	ords.							
Process only te	ext within brackets	Up to the	end of the para	graph							
Convert Emojis	to text	Ignore URLs (h	ttp, https and f	ф)							
		Ignore speaker	designations in	transcripts							
Case processing:											
Random sample	le: 1 out of: 5	Weighting vari	able: None		~						
Include cases	with missing values		-		-						

The Preprocessing tab contains two tabs: Options and Languages

The **Processes** section of the dialog box offers the following options:

Preprocessor: This option allows for the custom transformation of the text to be analyzed prior to, or in place of the execution of the other three standard processes provided by WordStat: lemmatization, exclusion and categorization. This transformation is accomplished by the execution of specially designed external routines accessible in the form of Python script, an external EXE file or a function in a DLL library. This feature is provided to offer greater flexibility by allowing any user with programming skills or resources to customize the processing of textual information. For more information on this feature see <u>Notes on Preprocessing</u>.

Stemming: This is a natural language processing routine that reduces inflected and derived forms of words to a common root form or word stem. The English stemmer, for example, returns "write" for "write", "writer", "writing" and "writings". Stemming can be especially useful in some exploratory text-mining tasks or when developing automatic document classification models by grouping related words together, and, thus reducing the total number of word forms. However, it may also decrease the precision in the measurement of some topics associated with specific inflected forms. Plural and singular forms of some nouns are often used to refer to different concepts or ideas. The same is true for various tenses of verbs. For example, in a sentiment analysis project, we found that the verb "improve" was often associated with negative comments, while its past-tense form "improved" was generally associated with positive comments. Since stemming is based on a limited number of morphological rules, stemming algorithms are also prone to errors. For example, the English Porter stemmer will group words like "universal", "universe" and "university" into the single word root: "univers". Stemming may also fail to group related words that do not follow typical grammatical rules.

Lemmatization: WordStat provides predefined substitution processes to perform lemmatization on documents. Lemmatization is a process by which various forms of words are reduced to a more limited number of canonical forms. A typical example of lemmatization would be the conversion of plurals to singulars and past tense verbs to present tense verbs. The lemmatization algorithm implemented in WordStat is a dictionary-moderated method, partly inspired by Krovetz's KSTEM suffix substitution algorithm. Since the lemmatization algorithm does not rely on a prior part-of-speech tagging of words, it is much faster than traditional lemmatization routines. It may, however, result in a few invalid word substitutions, but usually, those errors will have no major consequences on the result of an analysis.

You may override specific substitutions by creating custom word substitutions. It is important to remember that lemmatization, like stemming, may decrease the measurement precision of some concepts or topics.

Upon installation, WordStat provides lemmatization for English. It is also possible to install additional language modules, some of which contain lemmatization routines. Seven of those modules currently offer lemmatization: French, German, Italian, Norwegian, Polish, Swedish and Spanish. Additional lemmatization modules may later be available.

The **Character recognition** section of the dialog box offers the following options:

Case sensitive: By default, WordStat internally converts all text to uppercase letters so that processing of words is case insensitive. This may be inappropriate if you want to identify proper nouns or analyze text written in some European languages like German where differences in letter cases may denote different meaning. Enabling this option prevents the internal conversion to uppercase letters and will treat two instances of the same word in different cases (lower or upper case) as two distinct words.

Accept numeric characters: By default, every word consisting of numeric values or of a mix of letters and numbers is excluded from the analysis. This option can be used to include these words.

Add characters appearing: This set of options allows you to specify which characters, besides letters of the alphabet, should be considered as an integral part of a word. For example, the word "ex-wife" can be treated as a single word or as two separate words ("ex" and "wife") if the hyphen is included in the list of valid characters. Two edit boxes may be used to specify additional characters. The **Anywhere** option is used to specify special characters that will be considered as part of a word, no matter where they appear. The **Embedded in words** option should be used to specify characters that should be enclosed within other valid characters and not at the beginning or at the end of a word. For example, adding a period and a comma to the list of characters embedded in words will allow you to retrieve numeric values such as 97.5 or 1,000,000 or domain names like www.google.com as a single unit without the risk of retrieving words immediately followed by commas or periods.

The Text to Include or Ignore section of the dialog box offers the following options:

Don't process text within braces: This option can be used to instruct the program to skip all text found between braces (i.e. { and }). This option is especially useful to insert comments or annotations in the text variable without affecting the analysis of its content. It can also be used to ignore all questions, prompts, and other verbal interventions in an interview transcript made by the interviewer.

Don't process text within brackets: This option can be used to instruct the program to skip all text found between brackets (i.e. [and]). Since WordStat can also be configured to analyze only text found between such brackets (see option below), these two options may be used to toggle between an analysis of keywords entered manually between those brackets and of the surrounding text.

Process only text within brackets: This option can be used to instruct the program to process only the text found between brackets (i.e. [and]). This option is especially useful to perform an analysis on keywords entered manually in the text by one or more coders.

Ignore duplicates: This option can be used to prevent identical **documents** or **paragraphs** to be processed more than once. It may be useful to remove from documents some text segments that appear more than once such as the boilerplate of some forms, ads in magazines, questions in structured interviews, copyright notices, headers and footers, etc.. To identify and either tag or permanently remove identical documents in a project, see <u>Identifying Duplicate Cases</u>.

Analyze the first n words: In some situations, analyzing the full document is not essential for identifying the main topic and one may even get a better representation of those topics by focusing on just a few pages or a few paragraphs. A good example would be a dataset of papers from an academic journal. The first few pages which often include the title, the abstract and part of the introduction and ignore the remaining parts, including the reference section, allowing one to extract more easily and much more quickly the main topic of each paper. Another reason to use such an option may be to standardize the length of documents being analyzed. To enable this option, simply put a check mark and set the number of words you want to include in the analysis. You may also prevent WordStat from stopping in in the middle of a paragraph by enabling the **Up to the End of the Paragraph** option.

Ignore URLs: Selecting this option will prevent the text analysis of URLs starting with either HTTP, HTTPS or FTP by removing those links from the text data prior to the analysis.

Ignore speaker designations in transcripts: In interview transcripts, news transcripts, or debates, the persons speaking are often identified by their initials or by name in uppercase letters followed by a colon. Analyzing the text of those documents without removing those turn-taking indicators often results in a word frequency list containing all those indicators as the most frequent words. Enabling this option will remove those speech turn indicators allowing one to identify more easily what people are talking about rather than who is talking.

The **Case Processing** section of the dialog box offers the following options:

Random sample: When this option is activated, the program will randomly select a fraction of all cases and perform the content analysis on this subsample. The proportion of cases can be specified using the spin button located at the right of the checkbox. This option reduces the processing time for large files and is especially useful during the initial phase of an analysis where dictionaries are constructed, and categorization schema are developed and revised. It also allows you to preview the kind of results that would be obtained on very large data files.

Include cases with missing values: When examining the relationship between textual data and categorical or numerical variables, WordStat will skip any cases with a missing value on any one of these variables. Enabling this option instructs WordStat to include all cases, whether or not values are missing. All missing values are assigned to an additional class labeled as "MISSING." Any analysis involving comparisons between classes of categorical variables (cross-tabulation, correspondence analysis, etc.) will include this additional class.

Weighting variable: This option allows the selection of a variable that will be used to apply weight to the cases. When the program reads a case, the value of the weighting variable for this case is truncated to an integer. This integer value specifies how many times the case will be duplicated. If the value is less than one, the case is excluded from the analysis. This option is especially useful when the textual data to be analyzed has already been reduced to a frequency list, such as when analyzing a list of the most frequent queries on a search engine.

Notes on Preprocessing

The **Preprocessor** option allows users to access external text preprocessing routines that are not part of the WordStat program. This option is useful to perform custom transformation on the text to be analyzed. For example, a routine may be created to remove all foreign accents, to segment a document in a particular way, to perform part-of-speech tagging, word disambiguation, stemming or transforming words into letter n-grams (sequences of letters). These transformations are not applied to the original documents stored in the database but are instead performed live immediately after the textual information has been read into memory and prior to any text processing available in WordStat (lemmatization, exclusion of words, categorization, etc.). The same routines may also be applied permanently on text stored in WordStat using the <u>Transform Document</u> feature.

Such external routines may be written in R, Python, or in any programming language that can create or be called from a stand-alone EXE file or a DLL. The programmer is responsible for matching the formal parameters and result type of the conventions used by WordStat for data interchanges. Technical information on these programming conventions is available on request from Provalis Research.

To create a new preprocessing routine, see:

Writing a Preprocessing Routine in R or Python Calling an Executable program Calling a Function in a DLL

To edit the settings of an external routine:

- Select the preprocessing routine you would like to edit from the drop-down list.
- Click the Letter icon to edit the external routine settings.

To remove an external routine:

• Select the preprocessing routine you would like to remove from the drop-down list.

Click the Delete button.

Writing a Preprocessing Routine with Python

WordStat allows you to use a preprocessing script created in R or Python. This offers you a wide range of text processing options that may not be available in WordStat. There are certain guidelines that must be followed when creating your function:

- The function must be called "preprocessText".
- The function can only have one string argument.
- The return value must be a string.

Other than the restrictions above you are free to compose the functions of your choice.

To add a preprocessing routine:

- Go to the Text Processing tab > Preprocessing tab > Options tab
- To the left of the **Preprocessor** option select the *Preprocessor* button.
- Choose **R or Python function** from the list. A dialog like the one below will appear.

Name	xternal Process	
	lame	

- Enter the name of your process in the Name field.
- Select OK. A Python editor dialog will open.

Python Editor X Sample text: Script: Less than a mile from where I'm standing, near the banks of Test 1 import nltk the Mississippi River, there once stood the Pittsburgh Plate 2 from nltk.stem.wordnet import WordNetLemmatizer Glass Company. It was there for a hundred years. In its from nltk.corpus import wordnet 3 heyday, it employed 4,000 people and turned out thousands of tons of glass a year. It seemed that just about everybody 4 in town worked for what we called PPG. We didn't grow corn 5 or wheat here in Crystal City because for me, this is where 6 def get_wordnet_pos(treebank_tag): the world of possibility and hope all began, a world I want to 7 if treebank tag.startswith('J'): open for all America As a boy, I used to explore the bluffs to the south of town, looking for fossils and arrowheads. When 8 return wordnet.ADJ I was a little older, my grandfather and I sometimes took a 9 elif treebank_tag.startswith('V'): .22 and went down to the river and shot at logs floating by. 10 return wordnet.VERB We watched the great river ebb and flow. We felt its incredible force and marveled at its beauty. had only one 11 elif treebank_tag.startswith('N'): stoplight then, but it had a rich array of ethnic families 12 return wordnet.NOUN Among others, I remember the Auddifreds, the La Prestas, 13 elif treebank_tag.startswith('R'): the Trautweins, the Pouliezoses, the Fortneys, the Ryans, the Shapiros, the Cooks, the Salvos, the Evans - families 14 return wordnet.ADV drawn by the factory that used their special skills. Americans American 15 else: child children. 16 return '' 17 18 def preprocessText(aText): 19 tokens = nltk.word tokenize(aText) 20 21 Results: tags = nltk.pos tag(tokens) 22 Less than a mile from where I 'm stand , near the bank of the 23 lemmatizer = WordNetLemmatizer() Mississippi River, there once stand the Pittsburgh Plate Glass Company. It be there for a hundred year. In its heyday , it employ 24 4,000 people and turn out thousand of ton of glass a year. It seem 25 sResult = '' that just about everybody in town work for what we call PPG. We do 26 n't grow corn or wheat here in Crystal City because for me, this be where the world of possibility and hope all begin , a world I want to 27 for tag in tags: open for all America As a boy , I use to explore the bluff to the south 28 sPos = get_wordnet_pos(tag[1]) of town , look for fossil and arrowhead. When I be a little old , my 29 sTag = tag[0]grandfather and I sometimes take a .22 and go down to the river and shot at log float by. We watch the great river ebb and flow. We felt its 30 incredible force and marvel at its beauty. have only one stoplight then 31 , but it have a rich array of ethnic family. Among others , I remember 32 if (sPos != ''): the Auddiffeds , the La Prestas , the Trautweins , the Pouliezoses , the Fortneys , the Ryans , the Shapiros , the Cooks , the Salvos , the 33 sLemm = lemmatizer.lemmatize(tag[0], sPc Evans - family draw by the factory that use their special skill. 34 else: Americans American child child. 35 sLemm = sTag 20 🗸 ок X Cancel

The example above performs a lemmatization using NLTK WordNet lemmatizer using Python. This function tokenizes the text and applies Part-of-Speech (POS) tags to each token. It then lemmatizes each token in relation to its POS tag.

- Enter your source code in the Script window on the left of the dialog box. Be careful to follow the guidelines listed above.
- Enter sample text in the sample text window.
- Select the Test button
- Ensure the results in the **Results** window are as desired.
- Select OK.

To apply your R or Python function:

- Select the check box beside the Preprocessing option.
- Choose the R or the Python function from the drop-down menu.

Calling an Executable Program

The second method by which an external routine may be integrated within WordStat is through the calling of an executable program (typically a console application with an EXE file extension). With such a method, information is transferred between the two programs by way of temporary files generated on the fly and stored in a default temporary folder. WordStat first creates a text file (WORDSTAT.IN) in the temporary folder containing the text to be processed. It then calls the external routine, with optional parameters (which may or may not include the input and output file name and their locations). Once the external program ends, WordStat retrieves another text file (WORDSTAT.OUT) created by the external program and containing the text to be processed in place of the original one. The two temporary files are then deleted.

To configure WordStat to call an EXE file:

- Click the D button located on the right of the **Preprocessor** list box.
- Type the name you would like to give to this preprocessing routine and click **OK**. This name will be added to the list box of available routines.
- Enter the name of the program file including the full path or click the 🚔 button to display a dialog box that allows browsing through folders and then select the appropriate program file.
- In the **Working Dir** edit box, specify the working directory for the program if necessary. Specifying \$TEMP as the working directory instructs the program to set the working directory to the temporary folder.
- Enter the **parameters** to transfer to the program at start-up. Typically, you will transfer the input and output file names, as well as any command line options needed for the external routine, to perform the required transformation. You can specify multiple parameters and can use any one of these three string constants:

CONSTANT	STANDS FOR
\$TEMP	The system temporary folder.
\$IN	The temporary text file created by WordStat and to be processed by the external routine (the actual file name used is WORDSTAT.IN)
\$OUT	The name of the text file created by the external routine and retrieved by WordStat. (the actual file name used is WORDSTAT.OUT).

To apply your new function:

- Select the checkbox beside the **Preprocessor** option.
- Choose the new function from the dropdown menu.

Calling a Function in a DLL

The third method used to call a preprocessing routine is to execute an external function stored in a dynamic library (DLL). This will manipulate the text stored in the computer memory at a specific memory address. Since no file input and output operations are necessary to transfer information, this method is often much faster than running an EXE file. However, because this external routine has access to the same memory space as WordStat, great care should be taken when running or creating such a routine. To minimize the risks involved in calling external functions, a programming convention has been imposed where the name of the function to be called **must** begin with the two uppercase letters 'WS', thus reducing the risk of calling a function that has not been written specifically for WordStat.

To configure WordStat to call a DLL function:

• Click the D button located on the right of the **Preprocessor** list box.

- Type the desired name for this preprocessing routine and click **OK**. This name will be added to the list box of available routines.
- Enter the name of the DLL file including the full path or click the *button* to display a dialog box that allows browsing through folders and then select a DLL.
- Once a DLL file has been entered, the dialog box will provide a list of functions that are likely to be compatible with WordStat (with names starting with "WS"). Select the function containing the transformation routine needed.
- WordStat must set apart, in advance a "buffer size" that will contain the transformed text. By default, the memory space is equal to the length of the original document. For many text transformation routines, such as stemming or lemmatization, which often result in shorter text, this space should be large enough. However, for other types of text preprocessing (e.g., part-of-speech tagging or transformation of words into n-grams), the size of the transformed text may be twice or three times larger than the original. The **Buffer Size** option allows you to specify how much larger the memory space should be in order to hold the transformed text. A numerical value between 1 and 10 can be used to represent the value by which the original text should be multiplied. For example, if this option is set to 3 and the size of the original text into memory is 10 kilobytes, then 30 kilobytes of memory will be reserved for holding the transformed text. The calling routine should run a test to see whether the reserved space is sufficient and if not should return an error message, allowing WordStat to respond with an error message to the user.
- Once all the options have been set, click the **OK** button.

To apply your new function:

- Select the checkbox beside the **Preprocessor** option.
- Choose the new function from the dropdown menu.

Exclusion

The **Exclusion** tab contains an exclusion dictionary (also known as a stop list). It is used to remove all words that are not to be included in the frequency analysis. It is used mainly to remove words with little semantic value, such as pronouns and conjunctions. It may also be used to remove phrases.

WordStat 9.0	.7 - Election 200	18 Coded.ppj	- 0	×
🔳 🕅 Data	🔫 Text Proc	essing 📻 Frequencies 👏 Extraction 🗞 Cooccurrences 🛅 Crosstab 🏢 Keyword-In-Context 🔫 C	Classification	0.
Categorization	Model: Brand P	ersonality.wmodel 🗸 🧭 Open 🔻 🕁 Save 🕁 Save as 🛼 Publish 🌛 New 👔 Delet	te 🔃 Options	
🔇 Language	Preprocessing	Substitution Z Exclusion Categorization 🐑 Postprocessing		
+ Add		A ABOUT	Starts with: APPEAR*	
- Remove		ACTUALLY AGAIN	APPEAR WILL MATCH:	
Z Edit		AGAINST AIN'T		
🔊 Undo	🔁 Import	ALL ALLOW	APPEAR* WOULD ALSO MATCH:	
Print	Sexport 2	ALLOWS ALMOST		
		ALDING ALREADY ALREADY ALTHOUGH ALTHOUGH ALTHOUGH AMONGG AMONGST AN AND AND AND AND AND AND AND AND AND	APPEARING	
		ARE AREN'T AROUND AS ASIDE AT AWAY AX		
245 cases		552 stop word(s	5)	

The currently opened exclusion dictionary may be deactivated by removing the check mark in the check box in the tab. Wildcard symbols - such as * , ?, # and [] are supported. One may also specify case sensitive entries.

The * character (also call the "asterisk" or "star"): matches zero or more characters.

For example, the following expression:

REPORT*

will exclude all words beginning with REPORT (such as, REPORT, REPORTS, REPORTER),

while this expression:

COLO*R

will match both COLOR and COLOUR.

The question mark character (?): matches exactly one character.

For example:

WOM"N

will match both WOMEN and WOMAN.

The number sign (#) wildcard: stands for a single numerical digit.

For example:

B##

will match words like B12 (the vitamin) as well as B52 (the aircraft), while:

#####

will match any five digit number, typical of US zip codes. Please note that the use of the # wildcard will work as long as you set the Accept Numeric Characters check box on the **Preprocessing** tab.

The square brackets "[" and "]": are used to match a single character out of a list of characters. For example, [AEIOU] will match any one of these vowels, while [A-E] will match any letter between A and E.

The following pattern:

[A-Z][0-9][A-Z]_[0-9][A-Z][0-9]

will match typical Canadian postal codes, like H3B 1W9. Note that the underscore character in the above example represents a space.

An expression or a phrase that includes several words: may also be excluded by joining the various words with underline characters.

For example:

NOT_*

Will exclude all words preceded by the word "not".

The ^ character allows you to type a case sensitive entry:

For example, the following two entries:

^it ∕\t

Will match the pronoun "it" whether or not it appears at the beginning of the sentence. These two entries will still allow WordStat to recognize "IT", which stands for "information technology".

To add words to the exclusion list:

Press on the button and select Wo

button and select Words / Phrases...

- Add the words or phrases into the box on the Add words dialog.
- Select OK.

To remove a word from the exclusion list:

• Select the words you would like to remove.

Click the ______ button.

To edit an entry in the exclusion list:

- Select the words you would like to edit.
- Click the button.
- Modify the entry in the Edit selected words dialog box.

To search for an entry in the exclusion list:

- Click anywhere in the exclusion list.
- Press Ctrl F. A search dialog box will appear.
- Type the desired search string in the **Find What** edit box. To restrict the search to items starting with the search string, enable the **Match Beginning of Item** option and to restrict it to whole words matching the search string, enable the **Match Whole Word Only**.
- Click the **Find** button to search the first item matching the entry. Clicking this button again finds the next occurrence of the search string, starting at the currently selected item.

To undo a change in the exclusion list:

- Select the button. A dialog box appears listing changes.
- Select the operation you would like to undo.
- Select the button in the dialog box.

To import and exclusion list:

- Select the <a>Import button. A dialog box containing files available for import will open. WordStat can import from three file formats: An ANSI exclusion list file with a .exc file extension, a UNICODE stop word list file with a .stop file extension, or an existing categorization file (.wmodel).
- Select the file you would like to import and select Open button.
- You will be asked if you want to append the words to the exclusion list.
- Select Yes.

To export an exclusion list:

- Click the Export button. A Save File dialog box will appear
- Select the proper file format under which you would like to save the file. WordStat can export to an ANSI exclusion list file with a .exc file extension or a UNICODE stop word list file with a .stop file extension.
- Type the name of the file, and click **OK** to create it.

Adding Words to an Exclusion List from other Places in WordStat

From the Frequencies page:

- Select the rows containing the words you would like to add.
- Right click your mouse, a menu will appear.
- Select **TO EXCLUSION LIST** option.

Note: You may also drag and drop a word into the dictionary panel to the right of the frequency table (see <u>Using the</u> <u>Dictionary Panel</u>).

From the Crosstab page:

- Select the rows containing the words you would like to add.
- Right click your mouse, a menu will appear.
- Select the TO EXCLUSION LIST option.

Note: You may also drag and drop a word into the dictionary panel to the right of the frequency table (see <u>Using the</u> <u>Dictionary Panel</u>)

From the text editor:

- Select the word you would like to add.
- Right click your mouse, a menu will appear.
- Select the **TO EXCLUSION LIST** option.

If you choose to add a word to the exclusion list, the word will automatically be stored in this file without any dialog box.

Using the Dictionary Assistance panel

Taking into account the impact of words in the exclusion list on the recognition of other entries in the categorization dictionary may be difficult especially when some items contain a wildcard. A Dictionary Assistance panel displayed on the right can provide some guidance about possible interaction between the currently selected entry and other entries in the exclusion list and the categorization dictionary. For more information on how to use such a feature, see <u>Using the Dictionary Assistance panel</u>.

Substitution

The **Substitution** tab may be used to automatically replace specific words with other word forms. It may also be used to substitute common misspellings. You can also use this process to perform a simple type of categorization where specific words are replaced with keywords.



The currently opened substitution list may be deactivated by removing the check mark in the check box in the tab.

To add a new substitution:

Click the button. The following dialog box appears:

Add Substitution		×
Original:	Replace with:	
	V OK	X Cancel

• Type in the **Original** edit box the word you would like to replace, then type the replacement word in the **Replace with** edit box and click **OK** to create the new substitution rule. This new rule is automatically added to the substitution process.

To remove a substitution:

• Select the rule and click the ______ button.

To edit an entry in the substitution list:

• Select the words you would like to edit.

- Click the button.
- Modify the entry in the dialog box.

To undo a change in the substitution list:

- Select the button. A dialog box appears listing changes.
- Select the operation you would like to undo.
- Select the select the select the dialog box.

To import and add entries to the substitution list:

- Select the Provide the button. A dialog box listing categorization files (.wmodel) available for import will open.
- Select the file from which you would like to import substitution entries and select Open button.
- You will be asked if you want to append the words to the substitution list.
- Select Yes.

To search for an entry in the substitution list:

- Click anywhere in the substitution list
- Press Ctrl-F. A search dialog box will appear.
- Type the desired search string in the **Find What** edit box. To restrict the search to items starting with the search string, enable the **Match Beginning of Item** option and to restrict it to whole words matching the search string, enable the **Match Whole Word Only**.
- Click the **Find** button to search the first item matching the entry. Clicking this button again finds the next occurrence of the search string, starting at the currently selected item.

Adding Words to an Substitution List from other Places in WordStat

From the Frequencies page:

- Select the rows containing the words you would like to add.
 - Right click your mouse, a menu will appear.
 - Select SUBSTITUTE WITH option.

Note: You may also drag and drop a word into the dictionary panel to the right of the frequency table (see <u>Using the</u> <u>Dictionary Panel</u>).

From the Crosstab page:

• Select the rows containing the words you would like to add.

- Right click your mouse, a menu will appear.
- Select the SUBSTITUTE WITH option.

Note: You may also drag and drop a word into the dictionary panel to the right of the frequency table (see <u>Using the</u> <u>Dictionary Panel</u>)

Related Topics:

Misspellings & Unknowns

Categorization

The categorization process allows you to change specific words, word patterns, or phrases to other words, keywords or content categories and/or to extract a list or specific words or codes. This process requires the specification of a categorization dictionary. This dictionary may be used to remove variant forms of a word in order to treat all of them as a single word. It may also be used as a thesaurus to perform automatic coding of words into categories or concepts. For example, words such as "good", "excellent" or "satisfied" may all be coded as instances of a single category named "positive evaluation", while words like "bad", "unsatisfied" or expressions like "not satisfied" may be categorized as "negative evaluation".

A categorization dictionary may also contain rules delineating the conditions under which specific words or phrases should be categorized. The rules may consist of complex expressions involving Boolean (AND, OR, NOT) and proximity operators (NEAR, BEFORE, AFTER). These kinds of rules allow you to eliminate basic ambiguity in words by taking into account the presence of other words that may alter the meaning. A good example would be the presence of a negative word form (such as "rarely" or "never") close to an adjective. Another example would be the differentiation of the various meanings of the word "bank" by identifying other words like "river", "money" and "deposit" surrounding "bank". For more information on rules, see section <u>Working with Rules</u>.



The categorization dictionary is structured as a hierarchical tree where words, word patterns, phrases, and rules are grouped in a folder that represents a category name. Categories and individual items may also be included in a higher order category, allowing you to create multi-level dictionaries like the following one:

🚞 COUNTRY

🚞 NORTH-AMERICA

- CANADA (1)
 - UNITED_STATES (1)
 - USA(1)
- MEXICO (1)
- SOUTH-AMERICA
 - BRAZIL (1)
 - CHILI (1)

In the above example, words like CANADA, USA or MEXICO may be coded as either NORTH_AMERICA or COUNTRY, depending on whether the categorization is performed up to the first or second level of the dictionary (see <u>Level of Analysis</u> below).

Wildcards

Wildcards such as *, ?, # are supported.

For example, the following item under the support category:

🚞 SUPPORT

• SUPPORT*

will change SUPPORT, SUPPORTS, SUPPORTING, SUPPORTIVE, SUPPORTER, etc. into a single word SUPPORT,

while the following word pattern:

🚞 SUPPORT

• *SUPPORT*

will also substitute all words with the substring "SUPPORT" in it, such as UNSUPPORTEDLY, UNSUPPORTED, etc.

An expression that includes several words may also be substituted by joining the various words with underline characters. For example, you may change the expression "going out" with the category "NIGHTLIFE" by specifying the following item:

🚞 NIGHTLIFE

• GOING_OUT

You may also use wildcards in expressions such as:

🚞 NIGHTLIFE

• GO*_OUT

to substitute several forms of an expression at once.

Integer weights can be assigned to specific items so that a specific word or word pattern may count for more than one instance of the category. For example, in order to compute an aggressiveness score on specific texts, you may

choose to assign a weight of 5 points to word patterns such as KILL* or MURDER* but only a single point to word patterns like INSULT*.

Case Sensitive Entries

While by default, WordStat text analysis is case insensitive, sometimes you need to perform a case sensitive match. For example when you want to differentiate the pronoun "us" from the abbreviation of United States "US", or identify references to "Information technology" using the "IT" abbreviation. To create a case sensitive entry, simply type the ^ character as the beginning of the entry.

For example:

^Bill

will match the given name Bill, but will not match this same word without the uppercase letter B. It also won't mach this word if in all uppercase ("BILL") or bill as in "dollar bill".

Regular Expressions

Beside the above wildcards and case sensitive entries, WordStat categorization dictionaries may contain entries consisting of regular expressions, allowing you to detect specific patterns such as zip codes, phone numbers, or email addresses. For more information on how to enter, edit and test regular expression, see <u>Working with Regular</u> <u>Expressions</u>.

Categorization Settings

Level: The level option allows you to specify up to which level the categorization process should be performed. Two ways of setting levels is available. Setting this option to **Up to level** allows you to specify a fixed level up to which the coding should be performed. For example, in the following dictionary:

🚞 COUNTRY

if a level of 1 is specified, all words that are stored at a higher level than the root level will be coded as the parent category at this first level. For example, words like CANADA and MEXICO will be coded as COUNTRY along with other country names like BRAZIL. Setting the level of analysis to a numeric value of 2 will results in the coding of those two words as NORTH-AMERICA, while BRAZIL will be coded as SOUTH-AMERICA. Items stored at the same or at a lower level will remain unchanged. In the above example, increasing the level of analysis to 3 will return the frequency of each country name.

Setting the **Level** option to **As shown** instructs WordStat to match the level of categorization performed to the level of details currently displayed in the tree view of the categorization dictionary. This option allows one to set different levels of categorization by expanding broad categories that should be broken down and by collapsing categories for which finer details are not needed. For example, if we modify the above tree by collapsing the NORTH-AMERICA category, WordStat will display it the following way:

🚞 COUNTRY



- BRAZIL (1)
- CHILI (1)

The program will report frequencies of individual countries like BRAZIL or CHILI but will categorize every instance of CANADA, UNITED-STATES, USA and MEXICO as NORTH-AMERICA.

Please note that it is possible to prevent a category from being broken down into subcategories or items, even if the level of analysis is set to a higher setting or if it is set to AS SHOWN and the items contained in this category are visible. Such a feature is useful when the content of a category consists of different ways of referring to the exact same thing (for example UNITED_STATES, UNITED_STATES_OF_AMERICA, US and USA) or consists of various misspellings.

To make a category unbreakable, select the category in the dictionary tree, click the button, and put a check mark in the Unbreakable box. The folder icon normally used to represent categories will be transformed into a folder icon with a key inside. You may also select the category, right click, and then select UNBREAKABLE | YES from the pop-up menu. To unlock the folder, follow the previously described steps for editing the category and remove the check mark in the **Unbreakable** box.

Categories only: When the **Level** option is set to a value higher than one, this option instructs WordStat to limit the level increase to the coding of the last category at or below the specified level. This option is especially useful when working with unbalanced hierarchical categorization systems where individual words are stored at different levels. For example, in the following dictionary:

😂 SENSATION

Setting the level of analysis to 2 without enabling this option would code words like AROMA or BREATH as ODOR, but would include in the final results individual words like TREMOR or AFRAID. Enabling the **Categories only** option ensures that individual words won't be included but will be coded as their parent category.

Use full path as category name: When the Level option is set to a value higher than one, this option instructs WordStat to substitute the full path of an item as the category name. The slash (/) character is used to separate the various levels. For example, in the above example, enabling this option and setting the level of analysis to 2 will code the word AROMA as SENSATION/ODOR. Increasing the level of analysis up to 3 will return SENSATION/ODOR/AROMA.

Allow overlaps: By default, categories are mutually exclusive such that a word can only be entered in a single category. Enabling this option allows you to create overlapping categories where words can be classified simultaneously into two or more categories. However, please take note that current multivariate techniques available in WordStat such as clustering, correspondence analysis and multidimensional scaling as well as other multivariate statistical procedures make the assumption that categories are statistically independent. Using overlapping categories creates data that clearly violates this assumption and may yield dubious results.

Show warnings: Some items in an exclusion list or categorization dictionary may remain undetected in documents because of their incompatibility with some analysis options. This occurs, for example, when an item is found both in the categorization dictionary and the exclusion list, or when this item includes non-alphabetic characters that have

not been specified as valid. The following table displays the various types of problems that may be identified by WordStat:

ТҮРЕ	DESCRIPTION
Item includes invalid characters	WordStat identifies individual words using alphabetic characters and other special characters specified by the user in the <u>Valid Characters</u> option. So, to make sure any item containing non-alphabetic characters is properly recognized, this special character must be added to the list of valid characters.
Item includes numeric characters	An item in the categorization dictionary or the exclusion list that includes numeric characters cannot be recognized since the <u>Accept</u> <u>Numeric Characters</u> option is currently disabled.
Item also in the exclusion list	An item found in a categorization dictionary cannot be recognized if it matches an item found in the exclusion list.
Phrase starts with an excluded word	In order to be recognized, a phrase cannot start with a word found in the exclusion list. Therefore, this excluded word should preferably be removed from the exclusion list in order for the phrase to be recognized.

Enabling the **Show Warnings** option instructs WordStat to identify potential compatibility problems affecting items in a dictionary, and it displays a list of those problems in a special dialog box. This dialog box is displayed prior to the application of dictionaries for a content analysis.

Matching Priority in a Dictionary

Dictionaries or lexicons consisting of words or word patterns only often fail to achieve accurate measurements of the topics or concepts they are supposed to assess. This results from their inability to disambiguate words with multiple meanings. WordStat offer a way to disambiguate words and achieve higher accuracy by using an established priority when matching items in the dictionary. The evaluation order is based on a single principle that consists of giving priority to more specific items over less specific ones. Since phrases are more specific than words, and the longer the phrase, the more specific it is, and since whole words are more specific than word patterns (words with wildcards), the evaluation priority has been set this way:

- 1. Long phrases over short phrases
- 2. Phrases over words
- 3. Case sensitive words over non case sensitive ones
- 4. Words over word patterns
- 5. Longer word patterns over shorter ones

Such a matching priority list, not only allows one to measure topics more accurately, but it also prevents double-counting items when more than one candidate item may be matched by a specific text segment, such as when a phrase in the text correspond to a phrase in the dictionary as well as a word or a word patterns also in the dictionary.

For example if your dictionary contains the words and phrases:

🚞 ACCESSIBILITY

- HEALTH_CARE_COST (1)
 - COST_OF_HEALTH_CARE (1)

🛅 TOPICS

- HEALTH_CARE (1)
- CHILD_CARE (1)
- CARE
- HEALTH

and your document contains the sentence "If I'm elected, I promise to reduce the health care cost for middle-class families", WordStat will only match the item HEALTH_CARE_COST, and won't match entries such as HEALTH_CARE, CARE, or HEALTH, since those are less specific than the longest phrase. This will allow you to disambiguate words such as "bank", and remove false positives by excluding or categorizing elsewhere phrases such as "river bank" or "Gaza bank", or making sure than the dictionary entry JOBS won't be triggered when someone mention "Steve Jobs". If you use the word pattern TAX* to identify references to "tax", "taxes" or "taxation", it become possible to restrict hits to only those words by adding entries like TAXI, TAXONOMY or TAXIDERMY to the exclusion list or to another category. Finally, since case sensitive entries are more specific than words, it become possible to differentiate IT (information technology), US (country), or WHO (World Health Organization), from the pronouns "it" or "us", or "who".

Using the Dictionary Assistance Panel

Taking into account the above matching priority and envisioning possible interactions between items in a large dictionary may be challenging. It may be difficult to identify for a specific dictionary item whether it will be superseded by another item or, on the contrary, prevent some other items from being matched.

Another challenge when building categorization dictionaries is to assess the impact or using a * wildcard at the end of a word. Using such a wildcard may generate false positives, consisting of unrelated words matching the item. One may also be unaware of the full potential benefit of using such a wildcard resulting in some false negatives, or relevant items that could have been matched if a wildcard had been used. Also, if one decided to use such a strategy to match multiple words, how far can one go in the truncation of a word to get more relevant words without also matching irrelevant ones. For example, if you have a dictionary item COMPASSIONATE, changing this entry to COMPASSIONAT*, will bring three additional items: COMPASSIONATED, COMPASIONATELY, and COMPASSIONATING, but truncating four more letters to COMPASSI* will also bring COMPASSION and COMPASSING. However, truncating a single additional character will also match COMPASSES, two false negatives.

The Dictionary Assistance panel located on the right of the Categorization and the Exclusion sub-pages provides a way to identify potential interaction between the currently selected item and other dictionary entries. For single word entries it will also allow one to identify the potential impact of using a * wildcard at the end of the word or at different levels of truncation by displaying new entries that would be matched if such a word pattern were used. The following example provides an example of the type of information one may obtain when selecting a single word entry. For this example, we selected the item ENVIRONMENTAL in the Forest Value Dictionary.

=	Starts with: ENVIRONMENTAL*
EN	VIRONMENTAL WILL MATCH:
E] ENVIRONMENTAL
BU	T NOT IF IT ALSO MATCHES:
	ENVIRONMENTAL_CONCERN in category LIFE SUPPORT
	ENVIRONMENTAL_COST in category LIFE SUPPORT
	ENVIRONMENTAL_DEGRADATION in category LIFE SUPPORT
	ENVIRONMENTAL_VALUE in category LIFE SUPPORT
	ENVIRONMENTAL_QUALITY in category LIFE SUPPORT
	ENVIRONMENTAL_IMPACT in category LIFE SUPPORT
Π	WILL SUPERSEDE:
	ENVIRON* in category LIFE SUPPORT
EN	VIRONMENTAL* WOULD ALSO MATCH:
E	ENVIRONMENTALISM
E	ENVIRONMENTALIST
[ENVIRONMENTALISTS
[ENVIRONMENTALLY
E] ENVIRONMENTS
BU	T NOT IF IT ALSO MATCHES:
	ENVIRONMENTALLY_SUSTAINABLE in category LIFE SUPPORT
	ENVIRONMENTALLY_SENSITIVE in category LIFE SUPPORT

The panel allows us to see that ENVIRONMENTAL in the LIFE SUPPORT category would not be matched in the situation where it is followed by any one of the following words: CONCERN, COST, DEGRATION, VALUE, QUALITY or IMPACT. It will also prevent the matching of the item ENVIRON* since it is more specific that this word pattern. Now, the last two groups of items on this panel allow us to see that by adding a * wildcard at the end of ENVIRONMENTAL, we would also match five additional words. However, phrases such as "ENVIRONMENTALLY SUSTAINABLE" and "ENVIRONMENTALLY SENSITIVE" would have precedence over this word pattern or any of the suggested words (phrases being more specific than words and word patterns). Moving the **Start With** cursor located on the top of the panel tool bar will allow one to interactively identify the impact of various level of truncation on the matching items.

Several operations are available from such a list.

To add those words to a content category

- Put a check mark beside the words you want to add to your categorization dictionary.
- Click the \equiv button or right-click to display a contextual menu.
- Choose Add Selected to Categorization
- Select the category in which you want those words to appear and click OK.

To treat those words as exceptions by adding them to the exclusion list

- · Put a check mark beside the words you want to exclude
- Click the \equiv button or right-click to display a contextual menu.
- Choose Add Selected to Exclusion List

To replace the currently selected dictionary item with the word pattern

- Click the ≡ button or right-click to display a contextual menu.
- Select the Replace menu item and then select the displayed word pattern.

To replace the currently selected dictionary item with other words

- Put a check mark beside the words you want to replace the current entry with. You may select the original entry if you wish.
- Click the \equiv button or right-click to display a contextual menu.
- Select the Replace menu item and then choose With Selected Words.

When a phrase is selected, the panel will either be empty, indicating the absence of interaction with other dictionary entries, or it will contain a list of entries which the selected phrase may interact. In the example below, the phrase "ENVIRONMENTAL DEGRADATION" will have precedence over the words "ENVIRONMENTAL" and "DEGRADATION" as well as over the word pattern ENVIRON*, no matter to which content category they belong.

\equiv	Starts with:	
ENV	IRONMENTAL_DEGRADATION WILL SUPERSEDE:	_
0	EGRADATION in category LIFE SUPPORT	
E	NVIRON* in category LIFE SUPPORT	
E	INVIRONMENTAL in category LIFE SUPPORT	

Related Topics:

For more information on how to open, activate or deactivate a dictionary or how to add, edit or remove an entry in a dictionary, see <u>Creating and maintaining dictionaries.</u>

Or jump directly to one of the following topics:

- Opening an existing dictionary Creating or copying a dictionary Adding a word Adding a category Removing words or categories Editing a word or a category Moving words or categories Using lexical tools for dictionary-building
- Working with rules

Creating and Maintaining Dictionaries

Opening a Dictionary

To open an existing dictionary:

• Select the dictionary from the **Categorization Model** drop-down list. If the dictionary is not listed, click the button to display a dialog box that will let you browse through folders and select a dictionary.

To create or copy a dictionary:

• Click the D button located on the right of the exclusion list or of the categorization dictionary. A dialog box enables specifying the name and location of the new dictionary file. If a dictionary file is already active, it will ask whether existing entries should be copied to the new dictionary file. If you answer **Yes**, all entries in the previously opened dictionary will be retrieved and stored in the new one. Answering **No** will result in an empty dictionary.

To add words to a dictionary:

• Press on the **button** and select **Words / Phrases...** The program will display a dialog box similar to the following:



To add new items to an existing content category:

- Type the words or phrases you would like to add in the edit box, one item per line. Spaces are automatically replaced by an underscore character.
- Select the proper category from the right panel.
- Click the Add button.

To add new items to a new content category:

- Type the words or phrases you would like to add in the edit box, one item per line. Spaces are automatically replaced by an underscore character.
- Set the check box on the right panel to **New category**.
- Type the name of the new category and select the existing category under which this new category should be created. To create a main category, select the **<ROOT>** item.
- Click the Add button.

Wildcards such as * and ? are supported.

Weights can also be assigned to specific categorization, so that a specific word, word pattern, or expression may count for more than one instance of the concept. The default value for this option is 1. To use a lower or higher value, edit the Weight option either by entering a new numeric value in the edit box or by clicking the spin buttons to increase or decrease this value. Valid weights can be any floating point value higher than zero.

If you want to add a word to a non-existing category, you first need to create a category (see below) and then follow the above steps to add the word to this new category.

To add categories to the categorization dictionary:

- Press on the button and select the Categories command.
- Select the **Main Category** or the **Sub-Category** radio button depending on whether you want this new category to appear at the main level or whether you want it to be created under an existing category. If you choose to create a sub-category, you then need to select from the **Location** outline the category under which you would like to store it.
- Type the category names you would like to add in the edit box, one category per line, and click the Add button located on the right of this edit box.

To remove a keyword or a content category from a dictionary:

- Select the dictionary from which you would like to remove words or categories.
- Select the words or categories you would like to delete and click the button. If a non-empty category is selected, you will be asked to confirm its deletion. If you answer **Yes**, all words and subcategories belonging to this category will be erased.

To edit an entry in a dictionary:

• Select the item you want to modify and click the **select** button.

To search for an entry in a dictionary:

- Right-click anywhere in the categorization dictionary.
- Select the **FIND** command from the pop-up menu. A search dialog box will appear.
- Type the desired search string in the **Find What** edit box. To restrict the search to items starting with the search string, enable the **Match Beginning of Item** option and to restrict it to whole words matching the search string, enable the **Match Whole Word Only**.
- Click the **Find** button to search the first item matching the entry. Clicking this button again finds the next occurrence of the search string, starting at the currently selected item.

Moving Words or Categories

The easiest way to change the structure of a categorization dictionary is by using drag-and-drop operations. Using the mouse, you can move a word to a different category, or move an existing category or sub-category to another location on the main level or below an existing category. There are two ways to perform such a drag-and-drop operation:

For single-move operations:

- Press and hold the Alt key.
- Click the item in the dictionary you would like to move and hold down the mouse.
- Drag the item to the desired location and release the mouse button.

Another way to achieve multiple drag-and-drop operations is by enabling the **Drag & Drop Editing** option located to the left of the dictionary.

For multiple-move operations:

- Check the **Drag & Drop Editing** check box.
- Click the item you want to move, and hold down the mouse button.
- Drag the item to its new location, and release the mouse button.

By default, the dragged item is stored in the category at the cursor's position. To move a word or a category to the main level or to the same level as the category at the cursor's position, simply hold the Alt key while dropping the dragged item.

To move a word to a distant category:

- Select the item you would like to move.
- Click the right button of the mouse to display the contextual menu.
- Select the **MOVE TO** command.
- Choose the category where the selected item should be moved and click OK.

Adding Words to an Exclusion List from other Places in WordStat

From the Frequencies page:

- Select the rows containing the words you would like to add.
- Right click your mouse, a menu will appear.
- Select TO CATEGORIZATION DICTIONARY option.

Note: You may also drag and drop a word into the dictionary panel to the right of the frequency table (see <u>Using the</u> <u>Dictionary Panel</u>

From the Crosstab page:

- Select the rows containing the words you would like to add.
- Right click your mouse, a menu will appear.
- Select TO CATEGORIZATION DICTIONARY option.

Note: You may also drag and drop a word into the dictionary panel to the right of the frequency table (see <u>Using the</u> <u>Dictionary Panel</u>)

From the text editor:

- Select the word you would like to add.
- Right click your mouse, a menu will appear.
- Select TO CATEGORIZATION DICTIONARY option.

Related Topics:

Using Lexical tools for dictionary-building Merging categorization dictionaries Printing the categorization dictionary Using the Dictionary Panel

Working with Rules

The WordStat Rules editor may be used to define complex coding rules allowing you to specify under which conditions a particular item or category of items should be coded. Such a feature may be useful to differentiate between numerous meanings of a single word (disambiguation). For example, you may limit the coding of the word "bank" to situations where the word refers to the financial institution. This can be done by restricting the coding of "bank" to documents containing vocabulary related to monetary or financial transactions ("cash," "money," "mortgage," "investment," etc.) or by excluding alternate meanings such as when "bank" appears in close proximity to words like "river" or "canoe." Rules may also be used to measure various forms of a phrase. For example, the idiom "TURN OFF" may be expressed in many different ways ("turn it off," "turned off," "turned this off," "turned his radio off"). While figuring out all the possible forms of such an idiom may be very difficult, if not impossible, a single coding rule to look for the word pattern "TURN*" followed by "OFF" within the same sentence could very well cover most of these situations. Rules can also take into account the presence of words that may alter the power of an adjective, such as negations or qualifiers like "rarely," "numerous," "few," etc. Rules may even be used to identify sequences of events or complex actions.

In WordStat, a rule can refer to individual words, word patterns, or phrases. It may also refer to several items belonging to a content category of the current dictionary. A reference to a category is always preceded by the number sign character ('#'). For example, in the rule:

SATISFIED NEAR #PROFESSOR

the first item, SATISFIED, refers to a single word while #PROFESSOR will match any item found in the PROFESSOR content category.

Just like words or phrases, rules may be stored anywhere in a categorization dictionary. A rule consists of a target item and from up to 20 conditions, each condition consisting of another item linked to the first item using a Boolean (AND, NOT) or a proximity operator (NEAR, BEFORE, AFTER) or their negative forms (NOT NEAR, NOT BEFORE, NOT AFTER). The context in which these conditions will be tested also needs to be specified, allowing you to either consider the content of the entire document or restrict the test to a single paragraph or a single sentence. When a proximity operation is used, you also have to specify the maximum distance (in number of words) that must separate the two items in order for this proximity condition to be tested as true or false.

To create a rule:

Click the button and

button and select the RULES menu item. A dialog box similar to the one below will appear:
WOMENS_RIGHTS arget word, phrase or category: EQUALITY Operator: Word, phrase, or category: Within: Distance: Add AND WOM2N Same paragraph 5 0 3 NOT AFTER PHEGATIONS Same sentence 4 0 3 NEAR Solution	Name:								
arget word, phrase or category: EQUALITY Operator: Word, phrase, or category: Within: Distance: Add AND WOM?N Same paragraph 5 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	NOM	ENS_RIGHTS							
EQUALITY Operator: Word, phrase, or category: Add AND WOM?N Same paragraph 5 INDT AFTER #NEGATIONS Same sentence 4 Image: Store in: CULTURE CULTURE CULTURE CULTURE CULTURE Store in: PRO ENVIRONMENT PRO ENVIRONMENT PRO ENVIRONMENT PRO ENVIRONMENT SROUPS ETHNIC WOMEN	arget	word, phrase or catego	ory:						
Operator: Word, phrase, or category: Within: Distance: Add AND WOM?N Same paragraph 5 3 NOT AFTER #NEGATIONS Same sentence 4 1 NEAR #NEAR Same document 50 3 Match All O Match Any Weight: 1 1 Store in: CULTURE CULTURE-HIGH CULTURE-POPULAR SPORT ENVIRONMENT ON ENVIRONMENT PRO ENVIRONMENT PRO ENVIRONMENT ETHNIC WOMEN	EQUAL	.TTY		~					
Add AND VWM?N Same paragraph 5 \$ 3 NOT AFTER PREATIONS Same sentence 4 \$ 3 NEAR VS Same document 50 \$ 3 Match All OMatch Any Weight: 1 \$ Store in: So CULTURE CULTURE-HIGH CULTURE-HIGH CULTURE-POPULAR SPORT ENVIRONMENT CON ENVIRONMENT SO CON ENVIRONMENT SO		Operator:	Word, phrase, or category:		Within:		Distar	nce:	
NOT AFTER	Add	AND 🗸	WOM?N	~	Same paragraph	~	5		x
Match All Match Any Weight: 1 Store in: <a href="https://www.www.www.www.www.www.www.www.www.w</td> <td></td> <td>NOT AFTER ~</td> <td>#NEGATIONS</td> <td>~</td> <td>Same sentence</td> <td>~</td> <td>4</td> <td></td> <td>x</td>		NOT AFTER ~	#NEGATIONS	~	Same sentence	~	4		x
Match All Match Any Weight: 1 Store in: ROOT> CULTURE-HIGH CULTURE-POPULAR SPORT ENVIRONMENT CON ENVIRONMENT GROUPS ETHNIC WOMEN		NEAR V		v	Same document	v	50		x
Store in: Store in: CULTURE CULTURE-HIGH CULTURE-POPULAR SPORT SPORT CON ENVIRONMENT CON ENVIRONMENT PRO ENVIRONMENT SROUPS ETHNIC WOMEN) Mati	ch All O Match Any			9				
CULTURE-POPULAR SPORT SPORT CON ENVIRONMENT PRO ENVIRONMENT GROUPS ETHNIC WOMEN	Mati Weij	ch All O Match Any			9				
SPORT SPORT SOURCE SPORT CON ENVIRONMENT PRO ENVIRONMENT SOUPS ETHNIC WOMEN	Mati Weij Store	dh All ○ Match Any ght: 1 e in: 중 <root> 중 CULTUF</root>	RE TURE-HIGH		9				
CON ENVIRONMENT CON ENVIRONMENT PRO ENVIRONMENT GROUPS ETHNIC WOMEN	Mati Weij Store	ch All O Match Any ght: 1 e in: S <root> CULTUF CULT CULT CULT</root>	RE TURE-HIGH TURE-POPULAR		9				ľ
CON ENVIRONMENT	Mati Weij Store	ch All O Match Any ght: 1 e in: SPOE CULTUR CULT SPOE SPOE	RE TURE-HIGH TURE-POPULAR RT		9				
	Mati Weij Store	ch All O Match Any ght: 1 e in: S <root> CULTUF CULT CULT SPOF ENVIRO</root>	RE TURE-HIGH TURE-POPULAR RT NMENT		9				-
	Mati Weij Store	ch All O Match Any ght: 1 • • • • • • • • • • • • • • • • • •	RE TURE-HIGH TURE-POPULAR RT INMENT ENVIRONMENT		8				
WOMEN	Mati Weig Store	ch All O Match Any ght: 1 + e in: - < ROOT> CULTUF CULT CULT SPOF ENVIRO CON PRO	RE TURE-HIGH TURE-POPULAR RT NMENT ENVIRONMENT ENVIRONMENT		8				
	Mati Weij Store	ch All O Match Any ght: 1 * e in: <a <br=""> CULTUF	RE TURE-HIGH TURE-POPULAR RT INMENT ENVIRONMENT ENVIRONMENT S		9				
	Matu Wei	ch All O Match Any ght: 1 * e in: * <root> CULTUF CULT CULT SPOF ENVIRO CON PRO GROUP ETHY</root>	RE TURE-HIGH TURE-POPULAR RT NMENT ENVIRONMENT ENVIRONMENT S NIC		9				

The minimum requirements for a rule to be valid are:

- A unique name,
- At least one statement consisting of a target item, a Boolean or proximity operator and a second item,
- The text unit within which this rule should be tested, and optionally the maximum distance in words.
- The weight given to items meeting these conditions,
- The content category where the rule will be stored be stored.

The ampersand or at sign ("@") is used as a prefix to denote the presence of a rule and will be added automatically to the rule name. In the above dialog box, the rule @WOMENS_RIGHTS will be stored under the content category WOMEN and will be considered as true if the word EQUALITY occurs in the same paragraph WOMEN or WOMAN and if there is no item in the #NEGATIONS category in the same sentence up to four words before the target word (i.e. EQUALITY).

To enter a specific word, word pattern or phase, simply type the desired item. Spaces between words are automatically converted to underscore characters. To enter a content category, type the number or pound sign ('#') immediately followed by the name of the category. An existing category may also be selected from a drop-down list by clicking the down arrow located to the right of the edit box and clicking the appropriate category name, listed in alphabetical order.

The following operators may be used in a rule:

RULE CONDITION IS TRUE IF...

item1 AND item2	both items occur in the same document, paragraph or sentence.
item1 NOT item2	the first item occurs in the document, paragraph or sentence but not the second one.
item1 NEAR item2	both items occur in the same document, paragraph or sentence, and are no more than n words apart.
item1 BEFORE item2	both items occur in the same document, paragraph or sentence, and the second item appears after the first one within the next n words.
item1 AFTER item2	both items occur in the same document, paragraph or sentence and the first item appears after the second one within the next n words.
item1 NOT NEAR item2	the first item occurs in a document, paragraph or sentence, and is not found within n words of the second item.
item1 NOT BEFORE item2	the first item occurs in a document, paragraph or sentence, and is not followed within n words by the second item.
item1 NOT AFTER item2	the first item occurs in a document, paragraph or sentence, and does not occur within n words after the second item.

To add an additional condition, click the **Add** button located on the left of the first condition. To remove a condition, click the **X** button located on its right. When more than one condition is set, you will be asked to specify whether you want to match all criteria or match any one of them.

Once the rule has been properly defined, click the button located in the lower right corner of the dialog box to append the rule definition to the selected content category and to clear the form. Once you have finished entering rules, click the close button to quit this dialog box and return to the WordStat main screen.

Please note, in order to prevent any recursive or cross-reference problems in rules, content categories can only refer to words, word patterns or phrases stored in categories and will ignore the presence of other rules. For example, if a category named #NEGATION contains 10 word patterns and three rules, any reference to this category in a rule will take into account those 10 words and will ignore instances where any one of the three rules have been found to be true.

Working with Regular Expressions

A regular expression (RegEx) is a formal language made up of a sequence of characters or a text string for describing a search pattern that can be used to extract information from text. Using RegEx is is much more powerful that using wildcards because of the number or types of patterns that can be extracted. Regular expressions can be used in WordStat categorization dictionaries.

RegEx can be used to find the following patterns (this list is not exhaustive):

- floating point number
- password
- username
- email address
- IP Address
- credit card numbers in documents for a security audit
- · duplicated words
- HTML tag

- numeric ranges
- same word with different spelling
- two words near each other
- URL
- valid dates

RegEx can range form simplostic to complicated. You have probably used *.doc to find all word files in a file manager. The regular expression (regex) equivalent is [^]+\.doc. A more elaborate example would be a pattern for matching an email address. The following example is not complete but will convey the idea ([A-Za-z0-9-\._]+)@([A-Za-z0-9-\._]+)\.([A-Za-z] $\{2,6\}$).

A free and useful tool for writing RegEx is called can be found at <u>https://regex101.com</u>. It is a PCRE-based regular expression debugger with real time explanation, error detection and highlighting.

To add RegEx expressions:

• Click the button and select the **RegEx** menu item. A dialog box similar to the one below will appear:



In the example above we are looking for email addresses

- Enter your regular expression in the Regular Expression field.
- Enter sample text in the sample text window.
- Select the Test button

- Ensure the results in the **Results** window are as desired.
- Choose the location in the dictionary panel on the left where you would like the extracted text to be stored.
- Select OK.

Importing Dictionaries

The **Import** feature allows you to append categories and items contained in one categorization dictionary into another dictionary. WordStat categorization dictionaries may be imported from WordStat .CAT and .WMODEL files, Excel, CSV or tab-delimited files, as well as from XML files. The exportation process supports multilevel dictionaries as well as weighting.

To import from WordStat .cat and .wmodel files:

- From the dictionary page of WordStat, open the dictionary into which you would like to import new categories.
- Click the <u>Marginal Import</u> button. An **Open** dialog box is displayed.
- Locate the dictionary containing the items you would like to import and click **Open.** A dialog box similar to the one below is displayed:

Select categories	-		×
ut a check mark beside the	e categor	ies to imp	ort:
PRIMARY			-0
ORALITY			
TOUCH			
TASTE			
DODOR			
GENERA	L SENS	ATION	
HARD			
SOFT			
DEFENSIVE	SYMBO	E.	
PASSIVIT	Υ		*
_	1	-	
✓ Ok	Xc	ancel	

Select the categories you would like to import by clicking in the box beside the desired category(ies) and clicking OK. To select all items, right click anywhere on the list of categories and select Check All. Choosing Uncheck All removes all check marks previously entered.

If an imported category already exists in the currently active dictionary, WordStat ignores duplicate items and only imports new items not already found in the original category. New categories are appended to the existing structure along with all their items.

Other File Types

To import a dictionary stored in Excel, CSV, Tab-delimited and XML files the data file has to be formatted such that the names of the categories are listed in one-to-four columns and the items such as words and phrases are listed in an

additional column. The data file may also contain an additional column containing weights to be used for each item, this weight being represented by a positive integer or floating-point numerical value. The first row must contain a header that will help you identify the content of each column. A typical dictionary table consisting of items stored in a hierarchical dictionary with two levels (main categories and subcategories) along with individual weights may look like this:

MAIN CATEGORY	SUBCATEGORY	ITEMS	WEIGHT
EMOTION	ANXIETY	NERVOUS	0.6
EMOTION	ANXIETY	EMBARRASS*	2.1
EMOTION	ANXIETY	DISCOMFORT*	2.7
EMOTION	SADNESS	ABANDON*	1.0
EMOTION	SADNESS	ALONE	1.3
BEHAVIOR	AGGRESSION	ABUS*	1.2
BEHAVIOR	AGGRESSION	HUMILIAT*	0.8
BEHAVIOR	AGGRESSION	BEAT*	2.5
BEHAVIOR	ASSERTION	CLAIM*	1.0
BEHAVIOR	ASSERTION	REQUEST*	3.2

Another version that omits the repetition of category names may look like this:

MAIN CATEGORY	SUBCATEGORY	ITEMS	WEIGHT
EMOTION	ANXIETY	NERVOUS	0.6
		EMBARRASS*	2.1
		DISCOMFORT *	2.7
	SADNESS	ABANDON*	1.0
		ALONE	1.3
BEHAVIOR	AGGRESSION	ABUS*	1.2
		HUMILIAT*	0.8
		BEAT*	2.5
	ASSERTION	CLAIM*	1.0
		REQUEST *	3.2

To import from Excel, CSV, Tab-delimited and XML files:

- Click the Import button.
- An **Open** dialog box will appear, allowing you to select the file containing the dictionary you want to import.
- Once selected, click **Open**. A dialog box similar to this one will appear:

Import Dictionary	×
Main Categories	Words & Phrases
1- CATEGORY V	3- ITEM 🗸 🗸
Subcategories Level 2 (optional)	Weight (optional)
2- SUBCATEGORY 2 V	<none></none>
Subcategories Level 3 (optional)	
<none> ~</none>	
Subcategories Level 4 (optional)	
<none> ~</none>	Sparse entry of categories
	V OK X Cancel

The drop-down lists (to the left in this dialog box) allow you to identify the columns containing the category and subcategory names. Up to four levels of category can be specified this way. It is, however, possible to import dictionaries with higher levels of category by separating the deeper levels (fifth or more) with a backslash character and storing these at the fourth level as part of this subcategory name. For example, a fourth level subcategory that would be stored as: SENTIMENT/NEGATIVE/UNFAIR would create two additional levels underneath SENTIMENT, with NEGATIVE as a fifth level subcategory of SENTIMENT and UNFAIR as the sixth level (after NEGATIVE).

If the names of the categories are not fully specified, as in our first example, but instead resemble the second table where only the first occurrence of the category name has been entered, then you must select the **Sparse Entry of Category** check box to make sure the items will be imported and stored in their proper category.

- Once all the options have been set, click the **OK** button.
- A **Save File** dialog box will ask you to enter a dictionary file name. WordStat will point to the default folder where other dictionaries are stored. You may browse to a different folder to store the new dictionary at another location. Once saved, the newly created dictionary will become the default one.

Exporting Dictionaries

WordStat categorization dictionaries may be exported to WordStat .CAT, Excel, CSV or tab-delimited files, as well as from XML files. The exportation process supports multi-level dictionaries as well as weighting.

To export a classification dictionary:

- Click the sector button. A Save File dialog box will appear.
- Select the proper file format under which you would like to save the file, type the name of the file.
- Click OK to create it.

Printing Dictionaries

To create a printed version of the categorization dictionary:

• Click the print button. A print dialog box like the one below will appear, allowing you to set various options of what and how the dictionary should be printed.

Printer:	HP Officejet P	ro 8610 [A1D6D7]	¥	Properties
ontent:	All	V Suppre	ess weights	
Title:	Customer Satis	faction Dictionary		
ont size:	8	Number of colu	mns: 2 🛊	
Start	new page on roo	ot category		
Footer	-			
Content	File name	Page number	Date	

Printer: Select the desired printer. To adjust the printer settings, such as the page size and orientation or printer resolution, select the **PROPERTIES** button.

Content: This option allows you to specify what should be printed. Selecting **All** prints the entire dictionary along with all its items. Selecting **As Shown** will print only currently visible items. When the **Select Items** option on the **Categorization** tab is enabled, a third option **Selected Items** becomes available allowing you to restrict the printing to previously selected categories and items.

Suppress weights: When this option is disabled, weights of each item are printed within parentheses. Enabling this option prevents those weights from being printed.

Title: Use this option to specify a particular line of text that will appear at the top of each page.

Font size: This option may be used to adjust the font size used to print dictionary items.

Number of columns: Dictionaries may be printed with up to seven columns per page, allowing you to print large dictionaries on fewer pages. Please note that when increasing the number of columns per page, it may be necessary to decrease the font size to prevent the overlapping of items in adjacent columns.

Start new page on root category: Root categories are the dictionary categories still visible when the dictionary is fully collapsed. Selecting this option instructs the program to start the printing of all items starting from this root category at the top of a new page.

Footer: Enable this option to print a footer at the bottom of each page. A footer can consist of up to three items: The **Filename** (printed on the left margin of the footer), the **Page Number** (located at the bottom center of the page), and the **Date** (printed on the right margin of the footer).

Using Lexical Tools for Dictionary Building

Creating a comprehensive categorization dictionary is quite often a difficult, time-consuming and subjective task. WordStat can assist you in finding words that may be related to existing words in your categories by the use of several lexical tools:

- A spelling dictionary is used to propose inflected forms of existing words already in your dictionary. Several dictionaries are currently available for different human languages such as English, French, Italian, Dutch, etc.
- Two English thesauri are also used to propose synonyms of words already in your dictionary.
- A WordNet based lexical database is used to find synonyms, antonyms as well as hypernyms, hyponyms, coordinate terms, holonyms, meronyms, etc. This database contains over 150,000 root words (including many proper nouns) and offers over 120,000 synonym sets. The availability of word sense definitions allows for manual as well as automatic filtering of proper word senses.

These three tools are available through the auto suggest panel on the frequency list, as well a through two dictionarybuilding commands.

- The <u>Basic</u> command uses the selected spelling dictionaries and any available thesauri to identify related synonyms and inflected forms.
- The <u>Advanced</u> command gives you access to a more powerful dictionary-building tool that uses a WordNet based lexical database to find, not only synonyms, but all related words such as hypernyms, hyponyms, holonyms, meronyms, coordinate terms as well as the selected spell-checking dictionaries to find inflected forms of those words.

Basic Dictionary-Building Tools

To access the basic dictionary-building tool:

- Select the Categorization tab.
- Press on the Suggest button.
- Select the **Basic** command. WordStat will immediately start looking for synonyms and inflected forms of all words in your categorization dictionary and will report them in a dialog box like this one:

ynonyms Inflected forms			Show exi	sting wor	rds or
	Score	Frequency	*	Ì	Ac
> \APPEARANCE			=	-	Clas
ARTS					Clos
COMMUNICATION					
CI LEDUCATION					
C \FAMILY					
FAMILY	5	21			
ANCESTRAL	3				
BABY	3				
CHILD	3	12			
DAM	3				
GRANDMOTHER	3	1			
✓ INFANT	3				
☐ LINE	3	5			
LINEAGE	3				
PAPA	3				
PARENT	3				
PATRIARCHAL	3				
☐ STOCK	3				
AGE	2	11			
ANCESTOR	2				
ANCESTRY	2	1			
F BANTI ING	2		-		

This dialog box displays on the first page a list of synonyms that were found to be related to existing words in the various categories. Synonyms for a specific category are sorted so that those that were related to several existing words in this category are located at the top of the list while synonyms related to only a single word are located at the bottom. The numeric value under the **Score** column indicates the number of existing dictionary entries to which it was related, while the value under the **Frequency** column indicates how often this word has been found in the current text collection.

The second page lists all words whose spelling begins with the same letters as existing words and that were not already included in the actual dictionary. For example, if the word "understanding" is found in the dictionary, the program will suggest words like "understandings", "understandingly", "understands", "understand", "understandable", and "understandably". The frequency score indicates how often this word has been found in the current text collection.

To display only words existing in the current text collection, select the **Show existing words only** option, in the upper right-hand corner of the dialog box. Please note that if this dialog is accessed prior to any WordStat text analysis, this option will be grayed out. Running a simple frequency analysis on the current text collection will collect the frequency information needed to allow this option to be used.

To add suggested words to the dictionary, place a check mark beside the words you would like to add and click the **Add** button.

Click the **Close** button to return to WordStat.

Advanced Dictionary-Building Tools

The advanced dictionary-building tool can be accessed either as a stand-alone application of from within WordStat.

To run the stand-alone version:

• Point to the Programs folder in the Windows' Start menu, then select Provalis Research and then click Dictionary Builder.

To access the advanced dictionary-building tool from within WordStat:

- Select the Categorization tab.
- Press on the Suggest button.
- Select the Advanced command. A dialog box like this one will appear:

WordStat Dictionary Builder - v1.3	- 🗆 ×
Dictionary Definitions Words Inflected Forms	
Search for: Synonyms Antonyms Antonyms Similar terms Hypernyms (is a type of) Hyponyms (types of) Coordinate terms (same hypernym) Holonyms (is a part of) Meronyms (parts of) Members (same holonym) Attributes Others Inflected forms Match partial word Suggested words Existing leftover words only Existing leftover words only The provide only Wence words only <th></th>	

The first tab is used to set various dictionary and search options. The second and third pages are used to find words and idioms semantically related to existing entries in the dictionary, while the last page is used to find derived forms of those entries.

Dictionary Tab

The first tab of the dictionary builder program allows you to select or change the WordStat dictionary, specify the words and categories you want to work with, along with the type of relationship to look for. It also allows you to specify how the program will search for inflected forms of existing words in your dictionary.

To select a dictionary:

- Click the 🖻 button. A standard Open dialog box will appear.
- Select the WordStat dictionary file you want to work with.

Selecting words and/or categories: By default, the dictionary-building program will search for related words and idioms for all existing words and categories in your WordStat dictionary. To restrict the search to specific categories or words within a category, simply deselect the words and categories you want to exclude by removing the check mark beside them. Clicking a category check box to change its state also changes the check box state of all words and subcategories within this category.

Specifying the type of relationship to look for: The **Search for** group box allows you to specify what type of relationship the program will look for. For example, you may choose to search only for synonyms and similar terms or decide to also search for hypernyms, hyponyms, coordinate terms, etc.

Setting how Inflected forms will be retrieved: The **Match Partial Word** option affects how inflected forms are found. When this option is deactivated, the program only retrieves words that start with the whole word. For example, if the dictionary includes the word INTELLIGENT, the program will suggest words like INTELLIGENTLY and INTELLIGENTSIA. If the Match Partial Word option is activated, the program will also suggest words like INTELLIGENCE, INTELLIGENCES, INTELLIGIBLE, and INTELLIGIBLY.

Showing existing words only: By default, suggested words and phrases are presented whether or not they were found in the current text collection. Selecting the Existing Leftover Words Only option restricts the list of suggestions to those present in the text collection and not yet captured by the dictionary. This option will be grayed out if no text processing has been done yet in WordStat. Running a simple frequency analysis in WordStat prior to using running this program will collect the frequency information needed to allow this option to be available.

Definitions Tab

Using a comprehensive lexical database such as WordNet to find related words and phrases has one major drawback. Searching for numerous types of relationship for even small WordStat dictionaries can yield a huge number of suggested words. For example, when searching for suggested words for a dictionary containing 129 words grouped under 13 categories, more than 12,000 new words and phrases were obtained, many of them unrelated to the existing categories. Browsing through such a huge number of suggestions to find the most relevant ones can be an overwhelming task. The **Definitions** tab was created to reduce this burden by providing an intermediary step where the user can select, for each of the words, the word senses that are the most relevant to the containing category. The program offers both manual and automatic selections of word senses and also allows you to combine both methods.

Select: All	Compute relevance				8
Vord	Definition	Words	Relevance *	Word	list
ATHLETIC	relating to or befitting athletics or athletes vigorously active	0 3	0.0 0.0	BEAUTIFUL BEWITCHING CAPTIVATING CHARISMATIC CUNNING CUTE	
	pleasing to the eye or mind especially through beauty or charm	34	1.0	DINK DINKY	
ATTRACTIVE ATTRACTIVE	having power to arouse interest having the properties of a magnet; the ability to draw or pull	1	0.0 0.0	ENGAGING ENTHRALLING ENTRANCING	
ATTRACTIVELY ATTRACTIVELY	in a beautiful manner in a becoming manner	2 2	0.0 0.0	FASCINATING FETCHING GLOSSY HYPNOTIC	
ATTRACTIVENESS ATTRACTIVENESS	the quality of arousing interest; being attractive or something that attracts sexual allure	15 18	0.0 1.0	INVITING IRRESISTIBLE MAGNETIC MESMERIC	
BEAUTIFUL BEAUTIFUL BEAUTY	delighting the senses or exciting intellectual or emotional admiration (of weather) highly enjoyable	31 1	3.0 0.0	MESMERISING MESMERIZING PERSONABLE	
BEAUTY BEAUTY BEAUTY	the qualities that give pleasure to the senses a very attractive or seductive looking woman an outstanding example of its kind	23 14 8	1.0 0.0 0.0	PHOTOGENIC PIQUANT PLEASING PREPOSSESSING	
BODY BODY	the entire structure of an organism (an animal, plant, or human being) a group of persons associated by some common tie or occupation and regar	51 46	1.0	SEDUCTIVE SHOWY SPELLBINDING	
BODY	a natural object consisting of a dead animal or person	12	0.0	UNATTRACTIVE	

Automatic selection of word senses: WordStat dictionary builder uses a basic disambiguation algorithm to try to identify, among all word senses, those that are the most likely to be related to the containing category. This algorithm involves the computation for each word sense of a relevance score. The higher this score, the more likely the word sense will be related to the category, while a score equal to zero suggests that this word sense is unrelated to the category. Once those relevance scores have been computed, the program can use one of three different rules to select proper word senses.

Best: This rule instructs the program to select for each word, the sense that obtained the highest relevance score. When selecting the highest score, a 20% tolerance is used so that, sometimes, more than one word sense will be selected. This selection rule is the most conservative one and ensures that relevant word senses are the most likely to be selected. However, we have also found that this selection method may lack some sensitivity and may fail to select other relevant word senses (false negatives).

Relevance > 0: This rule instructs the program to select all word senses that have been found to be related, even slightly, to the category. This selection rule is very liberal in that it is the most likely to select most relevant word senses at the cost of a lack of specificity (too many false positives).

Relevance > 0.1: This rule is slightly more conservative than the previous one, in that it also rejects all word senses that have obtained a score of 0.1. Besides a score of zero, 0.1 is the lowest score that may be obtained. Experience has shown that, very often, word senses with such a low score are unrelated to the category. Removing those word senses thus results in an increase in specificity along with only a marginal decrease in sensitivity.

The application of any of these three rules is performed by selecting the proper rule from the **Select** drop-down list. This list box may also be used to select or unselect all definitions.

Manual selection of word senses: Manual selection of word senses can be carried out either alone or after an automatic selection has been made by the program.

Manual selection is performed simply by browsing through the list of all definitions and selecting those that are related to the current category while making sure unrelated definitions are unselected. The decision to include or exclude a specific word sense may rely on the displayed definition, on the relevance score, and also on the examination of all words that have been found to be related to this specific word sense. Those suggested words are automatically displayed in the right panel of the Definition page when the word definition is highlighted.

Selected word senses may be saved on disk by clicking the 🛱 button, and later retrieved by clicking the 🛣 button.

Once the word senses have been chosen, activating the Words page will start the search, extract all words and phrases related to the selected word senses, and will display them by categories and by the type of relationship (synonyms, antonyms, etc.).

Words Tab

The Words page displays a list of suggested words and idioms that were found to be related to existing words in the various categories and allows you to select suggestions and add them to the existing dictionary. The **All words** tab includes a list of all words and idioms that were suggested, irrespective of their relationship with the existing entries. The remaining pages allow you to examine those same words by the nature of their relationship with existing entries.

Compute specificity Add words (8)	View	definition		۵ 🛁
words Synonyms Antonyms Is a t	ype of Types Coordin	nates Parts	Similar	1
	Releva	Specificity		
EDUCATION TRAIN	1046 Words	1 00		
	5.00	1.00		
	4.08	1.00	COMMUNICATION 1 00	
	4.08	1.00	COMMUNICATION 1.00	
	4.00	0.80	WORK 1 00	
	7.00	1.00	WORK 1.00	
CREATIVE THINKEP	2.00	1.00		
	3.00	1.00		
LICHBROW	3.07	1.00		
MENTOR	3.06	1.00		
EXPONENT	3.06	1.00		
RATIONAL	3.01	1.00		
GRADE	3.00	1.00		
FDUCATION	3.00	0.75	WORK 1.02	
DRIL	3.00	1.00	Tronk The	
ASSEMBLAGE	3.00	1.00		
WONDERER	2.06	1.00		
WISE MAN	2.06	1.00		
	2.06	1.00		

Relevance ranking and sorting: For each suggestion, a Word relevance score is computed that takes into account the number of times an item has been suggested as well as the relevance score obtained by the word senses from which it was derived. These suggestions are presented in descending order of relevance so that the suggestions that are the most likely related to the containing category are located at the top of the list while suggestions that are less likely to be relevant are found at the bottom of the list.

Specificity index: Very often, a word is suggested in more than one category. This is especially true when the dictionary includes categories that are semantically close to each other. One good example of such a categorization system is the Lasswell dictionary that tries to differentiate ten different forms of power relations (power gain, power loss, cooperation, authoritative, conflict, doctrine, etc.). When making a decision on whether a word should be added to a given category, it is important to consider whether this word is specific to this category or whether it has also been suggested in other categories. The **Compute Specificity** button allows you to obtain a specificity index as well as a list of all the other categories in which this item appears. This specificity index is computed by making the sum of all relevance scores obtained by this word in the various categories and computing the proportion of this total score that is related to the current category. A specificity of 1.0 indicates that this item has only been suggested for this category. When the item has been found to be related to more than one category, a list of all other categories in which it also appears is displayed in the **Other Categories** column along with the relevance score obtained in each of those categories. You can use this information to decide to which category this word should be added.

To add words or idioms to categories

- Place check marks beside the item you would like to add.
- Click the **Add** button.

Inflected Forms Tab

The Inflected Form page lists all words whose spelling begins with the same letters as existing words and that were not already included in the actual dictionary. For example, if the word "understanding" is found in the dictionary, the program will suggest words like "understandings", "understandingly". If the **Match Partial Word** option is enabled (see <u>Categorization</u>), this same word will also yield words like "understands", "understands", "understandable", and "understandably". The **From** column displays the original word from which the inflected form has been derived.

 WordStat Dictionary Builder - v1.3 Dictionary Definitions Words Infle 	cted Forms	- 🗆 X
Add words Auto select		A 4
Word	From	
COMMUNICATION		
	COLLECE	
	DEGREE	
	INTELLECT	
INTELLECTUALIZE	INTELLECT	
T INTELLECTUALIZED	INTELLECT	
INTELLECTUALITY	INTELLECT	
INTELLECTUALISM	INTELLECT	
INTELLECTUALIST	INTELLECT	
INTELLECTS	INTELLECT	
INTELLECTION	INTELLECT	
INTELLECTUALNESS	INTELLECTUAL	
INTELLECTUALS	INTELLECTUAL	
INTELLECTUALLY	INTELLECTUAL	
INTELLECTUALIZING	INTELLECTUAL	
INTELLECTUALIZE	INTELLECTUAL	
INTELLECTUALIZED	INTELLECTUAL	
0 Forms	0.15	

To add suggested words to the dictionary:

- Place a check mark beside the words you would like to add.
- Click the Add button.

The **Auto Select** button allows you to automatically select from all the suggested forms with specific suffixes such as all suggested forms ending with 's' or 'ed'. When searching inflected forms of English words, it is also possible to use WordNet to automatically select words that share the same meaning as the original word from which it was derived. As an example, the program will automatically select words like BEHAVIORS and BEHAVIOURS as valid forms derived from BEHAVIOR since all three forms will yield the same WordNet definitions. You can also set this feature to accept any new word form for which there is at least one WordNet definition containing the original word. For example, when enabling this option the word COMPETING would be automatically selected as a valid inflected or derived form of COMPETITION since one of WordNet definitions associated with COMPETING (i.e. "Being in competition") contains the original word from which it was derived.

• Click the 🛁 button to return to quit the dictionary builder program and return to WordStat.

Postprocessing

The **Postprocessing** tab offers different options that control how the textual information should be processed once all other processes (such as stemming, substitution, exclusion, categorization, etc.) have been applied.

WordStat 9.0.7 - Election 2008 Coded.ppj	é	×
😑 🏢 Data 💡 Text Processing 🍵 Frequencies 🚯 Extraction 🗞 Cooccurrences 🛅 Crosstab 🚛 Keyword-In-Context < Class	ification	
Categorization Model: Laver & Garry.wmodel 🗸 🔗 Delete	14 Options	
🔇 Language 🚳 Preprocessing 👘 Substitution 🗹 Exclusion 🗹 Categorization 🐯 Postprocessing		
Post processing options: Add words with: Frequency case occurrence higher or equal to: Case occurrence less than: Case occurren		
Automatically perform spell correction Add spelling corrections to the substitution list (then disable automatic spelling correction) Ckeep active automatic spelling correction Current language(s): <u>American, British</u>		
245 cases		_

Post-Processing Options:

Add Words: When the categorization dictionary is disabled, all words that are not found in the exclusion list will be included in the final keyword frequency analysis. This option allows you to restrict the number of words included to the most frequent ones by setting a minimum **Frequency** or **Case Occurrence** criterion for inclusion. This option may also be used while the categorization dictionary is active to add to this list, other words that are used at a high frequency. However, this option can only be used to add new words to the list of words and categories found in this categorization dictionary and cannot be used to remove any items. To remove items in this dictionary based on a frequency or case occurrence criterion see the **Remove Words** option below.

Remove Words: This option allows you to restrict the number of included words or categories to the most frequent ones by setting a minimum **Frequency** or **Case Occurrence** criterion for inclusion. This criterion is applied both to items in the categorization dictionary and words that meet the criterion specified with the **Add Words** option.

Examples:

- If no categorization dictionary is used and you want to include any word that appears at least 10 times, but in no less than 5 different cases, you need to activate the Add Words option and set its criterion to a minimum Frequency of 10. You then have to set the Remove Words criterion to a minimum Case Occurrence of 5. Only words that meet both criteria will be included.
- When a categorization dictionary is used to lemmatize words, but you only want to obtain frequency information on those words that appear a specific number of times, you have to activate the dictionary and set the minimum frequency criterion of both the Add Words and Remove Words options to the required frequency.

• When a categorization dictionary is used to categorize words, but you only want to analyze the most frequent categories, you have to activate the dictionary and set the Remove Words option to the required frequency. In this situation, the Add Words option should be deactivated.

Leave Categories Equal to Zero: By default, WordStat removes from the frequency table any keyword or category in the categorization dictionary that was not found in the analyzed text. Enabling this option instructs the program to leave those items with a zero frequency in the table. This option is especially useful when comparing obtained frequencies to normative data or to other samples. This option should also be enabled when creating norm files (see Creating and Using Norm Files).

Remove Items Occurring in More than *n* **Percent of Cases:** This option allows you to remove keywords or categories appearing in more than a specified percentage of cases. This criterion is applied both to items in the categorization dictionary and to words that meet the criterion specified in the **Add Words** option. This option is especially useful to remove words that are too common to have any informative or discriminative value.

Keep a Maximum of *n* **Items:** This option allows you to restrict the number of included words or categories to a maximum number of items, based either on their **total frequency**, number of **case occurrences**, or on the computed **TFxIDF** index. This selection occurs only after all the previous frequency options have been assessed and only if the total number of remaining items is higher than the specified maximum. If the cutting point falls on a frequency or a case occurrence shared by many items, those with the highest **TFxIDF** values will be selected.

Automatic Spelling Correction:

Enabling the automatic spelling correction feature will instruct WordStat to automatically identify misspellings of words or dictionary entries that have met prior selection criteria to be included in the frequency analysis. It will automatically identify misspelling of known words in the currently active spell-checking dictionaries. It will also combine frequency information along with phonetic and string similarities to correct unknown words such as technical terms or proper names. When this automatic spelling correction is enabled, one can choose between two types of outcomes:

Add Spelling Corrections to the Substitution List - The first option will add all spelling replacements to the substitution list, allowing one to review all replacements made during the last text analysis. Once finished, the automatic spelling correction is disabled, preventing the software from applying spelling corrections on the same set of words. Please note that if the number of words to be analyzed is increased, if additional entries are added to the current categorization model or if one changes the categorization model, it is recommended to reenable this automatic spelling correction feature in order to identify and correct additional misspellings.

Keep Active Automatic Spelling Correction - The second option will perform automatic spelling correction every time WordStat will need to process the original text documents. Please note that this option will not store spelling corrections in the substitution list. Consequently, it will not be possible to review, override, or modify corrections that have been made. Such an option is especially useful when frequent changes are made to the text processing options.

The **currently language(s)** option lists the active spell-checking dictionaries. To change the active dictionaries, click on dictionary names displayed on the right to display a dialog box that will allow you to enable or disable dictionaries. One may also move back to the <u>Language</u> page to select different dictionaries.

The Frequencies Tab

The **Frequency** tab is used to display a frequency table of words or content categories. This tab can be used to perform a univariate frequency analysis on words or categories and to modify any of the active dictionaries or word lists.

WordStat 9.0.7 - Election 2008 C	Coded.ppj									
E Data Text Process	sing 🛒 Frequencies	5 Extraction		nces 🛅	Crosstab	Keyword-	In-Context	dassificati	ion	
= 🛛 E 🕸 📽 🗖 谢 :	0 0 8 9 5									
\$ 9 «	Included Leftover v	vords							» U Comparisons 💡 Suggestions	
S EXCLUSION LIST		FREQUENCY	% SHOWN %	PROCESS	ED % TOTA	NO. CASES	% CASES	TF . IDF	La E CANDIDATE	a 18
SUBSTITUTION	NEUTRAL	3968	23.80%	1.63%	0.67%	235	97.51%	43.4		
+ NEW CATEGORY	CONSERVATIVE	3422	20.52%	1.40%	0.58%	238	98.76%	18.6	Biden Biden	-
ROOT>	CULTURE	2250	13.49%	0.92%	0.38%	231	95.85%	41.4	Clinton-	
CULTURE	LAW-CONSERVATIVE	2246	13.47%	0.92%	0.38%	228	94.61%	54.1	Edwards	
CULTURE-HIGH	RADICAL	1878	11.26%	0.77%	0.32%	228	94,61%	45.2	Ciulian	
CULTURE-POPULAR	PRO ENVIRONMENT	822	4,93%	0.34%	0,14%	156	64,73%	155.3	Kucinich I	_
SPORT SPORT	LIBERAL	702	4.21%	0.29%	0.12%	159	65.98%	126.8	MCCain MCCain	
ENVIRONMENT	WOHEN	577	3.46%	0.24%	0.10%	156	56 4356	143.4	Obama	
CON ENVIRONMENT	CON ENVIRONMENT	292	1.75%	0.12%	0.05%	05	30.42%	118.1	Richardson	
PRO ENVIRONMENT	ETHNIC	214	1 29%	0.00%	0.04%	87	26 1094	04.7	Romney	
GROUPS	RIIRAL	158	0.05%	0.05%	0.03%	81	23 6186	74.8	0 5 10 5	15
ETHNIC	HDRAN	26	0.45%	0.00%	0.0370	52	31 6004	50.6		
WOMEN	CHI TUPE HICH	70	0.40%	0.0376	0.01%	22	21.3070	30.0	🚽 🚩 🖄 Rate per 10,000 words 🤟 BY DELIVERY	- E
INSTITUTIONS	CULTURE-HIGH	38	0.23%	0.02%	0.01%	23	9.54%	38.8		
CONSERVATIVE	CULTURE-POPULAR	28	0,1/%	0.01%	0.00%	24	9.96%	28.1		
NEUTRAL RADICAL LAW_AND_ORDER LAW-CONSERVATIVE	LAW-LIBERAL	2	0.01%	0.00%	0.00%	2	0.83%	4.2		\sum
LAW-LIBERAL RURAL URBAN VALUES CONSERVATIVE LIBERAL	L8								2000 Cr.2001 Cr.2001 Cr.2001 Cr.2000 Cr.2000 Cr.	008 04-2008
AUTHORITY CONTINU" DISRUPT* INSPECT* JARISOLICION* LEGITIVATE MARAGC* MARAGCALM RUL* STRIKE* WHITEHALL										CAL

By default, the table shows the included keywords or categories in descending order of frequency.

The table includes the following statistics:

FREQUENCY	Number of occurrences of the keyword					
% SHOWN	Percentage based on the total number of keywords displayed in the table					
% PROCESSED	Percentage based on the total number of words encountered during the analysis					
% TOTAL	Percentage based on the total number of words that have not been excluded					
NO OF CASES	Number of cases where this keyword appears					
% CASES	Percentage of cases where this keyword appears					
TF∗IDF	Term frequency weighted by inverse document frequency. Such a weighting is based on the assumption that the more often a term occurs in a document, the more it is representative of its content yet, the more documents in which the term occurs, the less discriminating it is.					

Tabs at the top of the table allow you to access (1) a frequency table of all **Included** content categories or keywords or (2) a frequency table of **Leftover words** consisting of individual words that have not been categorized or included in the analysis.

To the left of the main grid, a panel lists the content categories of the current dictionary following additional entries representing the substitution and exclusion processes. This panel may be used to quickly assign items in the table to any one of these locations using the drag-and-drop operation. For more information on how to use this panel, see Using the Dictionary Panel.

The ≡ button can be used to move one or several words to the exclusion list or to add or remove a word from the categorization dictionary. The permitted moves depend on the items currently displayed. This button may also be used to display a Keyword-In-Context table of the selected word.

It is also possible to quickly access the pop-up menu invoked by this button by pressing the right button of the mouse anywhere on the grid.

To the right of the main grid is a panel containing two tabs: **Comparisons** and **Suggestions**.

The Comparisons Tab

The **Comparisons** tab allows you to look at the distribution of the selected words or categories among values of up to two structured variables. You may display this distribution using either a vertical bar chart, a horizontal bar chart or a line chart, by clicking on the corresponding button.

Four statistics may also be represented on the charts:

Case Occurrence	Number of cases in this subgroup containing at least one of these words.
Category Percent	Percentage of cases in this subgroup containing at least one of these words.
Word Frequency	Total number of these words in this subgroup.
Rate per 10,000 Words	Rate of words in this subgroup per 10,000 words.

The bottom chart contains the distribution of keywords or categories that can be represented as a word cloud, a vertical or horizontal bar chart , a pie chart and a donut chart. You can display the distribution with the same statistics seen on the frequency table: Frequency, % Shown, % Processed etc.

Right-clicking anywhere in the chart areas displays a pop-up menu that allows you to edit the chart, save it to disk or in the **Report Manager**, or copy it to the clipboard. Clicking a specific bar or a data point of a line chart also allows one to retrieve text segments associated with the selected class and containing words of the selected topic.

The Suggestions Tab

The **Suggestions** tab enables you to view an automatic suggestion panel displaying leftover words potentially related to the currently selected item. For more information on the auto-suggest panel, see <u>Working with the Auto-Suggest Panel</u>.

The Frequencies Toolbar

The button allows you to produce bar charts or pie charts to visually display the distribution of specific keywords or categories.

To produce charts:

- Set the **Sort By** option to the order in which you wish the values be shown graphically.
- Select the rows you would like to plot (multiple but separate rows can be selected by clicking while holding down the CTRL key)
- Click the 👪 button.

For further information see Bar Charts, Pie Charts, and Word Clouds

The kappaties button allows you to create normative frequency data from the current file, to store them on disk and to compare currently displayed frequencies with previously saved norms. See <u>Creating and Using Norm Files</u> for more information on this topic.

The solution allows you to access a keyword retrieval feature to retrieve all documents, paragraphs or sentences containing a specific keyword or a combination of keywords. See <u>Keyword Retrieval</u> for more information on this topic.

The solution allows you to automatically attach QDA Miner codings to all paragraphs or sentenced associated with currently displayed content categories or keywords. When clicked, a dialog box asks whether the coding should be applied to whole paragraphs or to individual sentences. If some WordStat keywords have no corresponding QDA Miner codes, new codes will be created under a special codebook category before the autocoding process begins.

The button is used to draw color guidelines on alternate rows in order to facilitate the reading of large tables. When clicking this button color guidelines are shown. Clicking this button again removes the color guidelines.

The **O** button displays various statistics on the text categorization process, such as the total number of words processed, the number and proportion of words that have been excluded and of those that have been categorized. The dialog box also displays document statistics - such as the average length in words of sentences, paragraphs and documents - as well as several coverage statistics, including the percentage of cases, paragraphs and sentences containing at least one keyword and the proportion of words that have been categorized or included. The coverage statistics are especially useful when you apply a content analysis dictionary developed to describe a specific data set on new data sets. A significant decrease in coverage may indicate the need to update a dictionary in order to better reflect changes over time or specific differences in this new data set.

The C button is used to reapply the content analysis process on the current data set. This button is disabled by default, and it becomes enabled when changes are made to any one of the currently active text analysis processes (such as the categorization dictionary, the exclusion list, or the substitution process). Clicking this button will instruct WordStat to reprocess the text collection and update the current table.

The *button* is used to access the mapping function, which allows you to analyze the spatial distribution of various terms.

The 🗱 button is used to export the selected data to Tableau, a software used for interactive data visualization.

The button may be used to append frequency information to the current data file or save to disk a matrix of word or keyword frequency by cases. For more information on one of these topics see <u>Exporting Frequency Data</u>.

The button is used to apply post-processing scripts on WordStat data files, For more information on how run an existing script or create your own script, see <u>Running and Creating Post-Processing Scripts</u>.

To append a copy of the table in the Report Manager:

• Click the 🛄 button. A descriptive title will be provided automatically for the table. To edit this title or to enter a new one, hold down the **Shift** keyboard key while clicking this button.

(for more information on the Report Manager, see the Report Management Feature topic).

To export the frequency table to disk:

- Click the 🔚 button. A Save File dialog box will appear.
- In the **Save as Type** list box select the file format in which to save the table. The following formats are supported: ASCII file (*.TXT), Tab delimited file (*.TAB), Comma delimited file (*.CSV), HTML file (*.HTM;*.HTML), and Excel spreadsheet file (*.XLS).

- Type a valid file name with the proper file extension.
- Click the **Save** button.

To copy the entire table to the clipboard:

- Right-click anywhere in the frequency table.
- Select the COPY TO CLIPBOARD | TABLE command from the pop-up menu.

To copy selected rows to the clipboard:

- Select the rows you would like to copy (multiple but separate rows can be selected by clicking while holding down the **Ctrl** key).
- Right-click anywhere in the frequency table.
- Select the COPY TO CLIPBOARD | SELECTED ROWS command from the pop-up menu.

To search for a specific item:

- Right-click anywhere in the frequency table.
- Select the **FIND** command from the pop-up menu. A search dialog box will appear.
- Type the search string in the **Find What** edit box. To restrict the search to items starting with the search string, enable the **Match Beginning of Item** option. To restrict the search to whole words matching the search string, enable the **Match Whole Word Only**.
- Click the **Find** button to search the first item matching the typed string. Clicking this button again finds the next occurrence of the search string, starting at the currently selected item.

Using the Dictionary Panel

The dictionary panel provides an easy way to assign words or phrases to the current categorization dictionary and to the exclusion or substitution lists. It may also be used to remove or edit existing items. This panel is located to the left of the **Frequencies**, **Crosstab** and **Phrase Finder** pages and looks like the one shown below.

O WordStat 9.0.7 - Election 2008 C	oded.ppj									- 0	×
E Data Text Process	ang 🛒 Prequencies	C Extraction	n S. Cooccu	imences 💼	Crosstab	Keyword-In	n-Context	< Classifica	abon		0
= 💵 🗟 🤹 🗰 🛈 🕯	0 0 8 5 5									1	66
\$ \$	Included Leftover w	ords							» 🏭 Comparisons 🦞 Suggestions		
EXCLUSION LIST		FREQUENCY	% SHOWN	% PROCESSE	D % TOTAL	NO. CASES	% CASES	TF . IDF	Suggestions: () Less () More		-=
SUBSTITUTION	NEUTRAL	3968	23.80%	1.63%	0.67%	235	97.51%	43,4	maiolinais		_
+ NEW CATEGORY	CONSERVATIVE	3422	20.52%	1.40%	0.58%	238	98.76%	18.6	SYNONYMS		
<pre><root></root></pre>	CULTURE	2250	13.49%	0.92%	0.38%	231	95.85%	41.4	[] BROAD (34)		- 1
CULTURE	LAW-CONSERVATIVE	2246	13.47%	0.92%	0.38%	228	94.61%	54.1	BROADER (27)		
CULTURE-HIGH	RADICAL	1878	11.26%	0.77%	0.32%	228	94.61%	45.2	BROADEST (1)		
CULTURE-POPULAR	PRO ENVIRONMENT	822	4.93%	0,34%	0.14%	156	64.73%	155.3	DAUGHTER (25)		
SPORT	LIBERAL	702	4.21%	0.29%	0.12%	159	65.98%	126.8	DAUGHTERS (39)		
ENVIRONMENT	WOMEN	577	3,46%	0.24%	0.1096	136	56,43%	143.4			
CON ENVIRONMENT	CON ENVIRONMENT	292	1.75%	0.12%	0.05%	95	39,42%	118.1			
PRO ENVIRONMENT	ETHNIC	214	1.28%	0.09%	0.04%	87	36,10%	94.7			
GROUPS	RURAL	158	0.95%	0.06%	0.03%	81	33,61%	74.8			
ETHNIC	URBAN	76	0.46%	0.03%	0.01%	52	21 58%	50.6			
WOMEN	CHI TURE-HIGH	20	0.72%	0.02%	0.0196	22	0 54%	20.0			
INSTITUTIONS	CIII TURE-POPULAR	20	0.17%	0.02%	0.01%	23	0.06%	20.1			
CONSERVATIVE	I AW-I TRERAL	20	0.0104	0.00%	0.00%	24	0.0206	4.7	(16) MISSED (16)		
C NEUTRAL	Contractions.	2	0.01.70	0.00%	0.0074	2	0.0378	7.2	MISSING (12)		
RADICAL	-										
LAW_AND_ORDER									ANTONYMS		
LAW-CONSERVATIVE									BOYS (10)		
LAW-LIBERAL									(10015 (10)		
RURAL									MANNED (3)		
URBAN											
VALUES									MEN (237)		
CONSERVATIVE									RELATED WORDS		
LIBERAL									(140) T (74)		
CONTINU*									AMAZON (2)		
DISRUPT*									BARIES (8)		
INSPECT=									E BEAUTY (1)		
LEGITIMATE											
MANAG*											
RUL*									CI ASSES (14)		
STRIKE"									DEVE (6)		
WHITEHALL											

The main section of this panel is a tree representing the structure of the current content analysis along with the substitution and exclusion processes (if active). When an item on this list is selected, its content is listed in a resizable window below the tree structure. At the top of the panel, a button allows you to undo the last change made to any one of the lists.

To move items from the frequency list to a list or to an existing content category:

- In the table to the right of this panel, click the word or phrase you wish to assign to a list or content category. To select a group of adjacent entries, move the mouse cursor over the first item in the list, click and hold the mouse button, drag the mouse to the last entry to highlight the block of rows you want to assign, and then release the mouse button. To select disjointed items, hold down the **Ctrl** key while clicking each one of the items.
- Once the words or phrases have been highlighted, drag and drop them on the category of your choice. To add an item to the root folder of a categorization dictionary, simply drop it into the **<ROOT>** folder.

To move items to a new category:

- Select the words or phrases you would like to assign to this new category.
- Drag and drop the items on the + NEW CATEGORY item. An Add Word/Category dialog box will appear, allowing you to type the name of this new category.

To rename or delete an existing content category:

- In the tree representation of the dictionary, right-click the category you would like to rename or delete.
- Select the appropriate command.

To rename, delete or move an item in a list or in a content category:

- In the window below the tree display, right-click the item you would like to modify.
- Select the appropriate command.

To cancel a modification:

• If you want to undo the last assignment, deletion or modification made using this panel, click the model button.

Note: If you leave the mouse cursor over this button a hint window will appear showing you which modification will be canceled by clicking this button.

Working with the Suggestions Tab

One of the biggest challenges of quantitative content analysis lies in the fact that a single idea may be expressed in many different ways, like using synonyms, paraphrases, or idioms. When you need to identify all instances of such an idea, a critical task becomes identifying these numerous forms. WordStat provides several tools to support this task such as the Suggest feature on the Text Processing tab that can be used to retrieve a list of all known synonyms, related words and inflected forms of the items already in a dictionary from another source such as a thesaurus or a lexical database. The Suggestions tab on the Frequencies tab also lists suggestions. It differs, however, from the Suggest feature in several ways. First, it only displays suggested words that were found in the current text collection and that are not already in the categorization dictionary (leftover words). It may also be used not only on content categories but also on any words extracted by WordStat and displayed in the Included Words and Leftover Words lists. It may thus be used from the very beginning of the dictionary construction process to guickly identify potential groupings of words and assign the relevant ones to existing or new content categories. To display this panel, select the **Suggestions** tab to the right of the frequency table. To display suggestions, simply select the appropriate row in the frequency grid. When the selected item is a word, the panel will display synonyms, antonyms, related items and words with a similar beginning (potentially related items, inflected and misspelled forms). When the selected row consists of a content category, then the panel displays the same information but for all words currently in this content category. Selecting more than one row will also result in a compound view of suggestions of all selected items.

At the top of the panel, radio buttons allow you to choose the extent of the suggestions. By default, the panel returns the most likely synonyms, while related words consist of hypernyms, hyponyms, holonyms and meronyms. Choosing the **More** option retrieves other potential synonyms, as well as coordinate terms and possible attributes of the selected item.

WordStat 9.0.8 - Election 2008b.pp	rj							- 🗆 X
E Data Text Processing	Frequencies	🔧 Extraction	S Coocc	urrences [Crosstab	Keyword-Ir	-Context	< Classification
= 💵 🗟 🖉 🗖 🛈 🕽	of 22 🖪 🛸							ui 🖬 🖨
\$ 9	Included Lefton	ver words						» 📙 Comparisons 💡 Suggestions
EXCLUSION LIST		FREQUENCY	% SHOWN	% PROCESS	ED % TOTAL	NO. CASES	% CASES	Suggestions: Less
SUBSTITUTION	STAND	369	0.18%	0.15%	0.06%	140	57.61%	
+ NEW CATEGORY	MIDDLE	364	0.18%	0.15%	0.06%	136	55.97%	SYNONYMS
C <root></root>	POLICY	362	0.18%	0.15%	0.06%	138	56.79%	COMMITTEES (8)
HUMAN RIGHTS	TROOPS	362	0.18%	0.15%	0.06%	96	39.51%	
EMPLOYEE	FIGHT	350	0.18%	0.14%	0.06%	132	54.32%	LEGISLATURES (3)
SOCIAL AND COMMUNITY	CHALLENGES	347	0.17%	0.14%	0.06%	138	56.79%	
ENVIRONMENT	BILLION	342	0.17%	0.14%	0.06%	111	45.68%	
G FAMILY	CONGRESS	335	0.17%	0.14%	0.06%	122	50.21%	
NOT RELEVANT	THINGS	335	0.17%	0.14%	0.06%	117	48.15%	
	IRAN	334	0.17%	0.13%	0.06%	65	26.75%	
	ACROSS	331	0.17%	0.13%	0.06%	145	59.67%	: ROOTED (22)
	POLITICAL	330	0.17%	0.13%	0.06%	136	55.97%	ROOTING (6)
	REAL	330	0.17%	0.13%	0.06%	138	56.79%	ROOTS (8)
	TAXES	319	0.16%	0.13%	0.05%	85	34.98%	SENATE (162)
	CAMPAIGN	317	0.16%	0.13%	0.05%	128	52.67%	SAME START
	JOB	317	0.16%	0.13%	0.05%	132	54.32%	CONGREGATIONALISTS (1)
	HISTORY	314	0.16%	0.13%	0.05%	143	58.85%	CONGRESSES (2)
	COSTS	312	0.16%	0.13%	0.05%	96	39.51%	CONGRESSIONAL (26)
	PUBLIC	311	0.16%	0.13%	0.05%	121	49.79%	CONGRESSIONALLY (1)
	HARD	309	0.15%	0.12%	0.05%	143	58.85%	CONGRESSMAN (17)
ABORIGINAL_PEOPLES	LOOK	309	0.15%	0.12%	0.05%	135	55.56%	CONGRESSMEN (3)
ABUSE	FOREIGN	308	0.15%	0.12%	0.05%	121	49.79%	CONGRESSWOMAN (3)
ACCEPT	PERCENT	299	0.15%	0.12%	0.05%	78	32.10%	CONGRESSWOMEN (1)
ACCOMMODATING	OPPORTUNITY	294	0.15%	0.12%	0.05%	136	55.97%	
ACCOMMODATION ACCOUNTABILITY	<						,	•
243 cases					Shown: 14,8	73 Types: 15	,838 Token	is: 592,770 Time: 1.9s

To assigning suggestions to a dictionary:

There are two methods to assign suggested words to the categorization dictionary or to the exclusion list. The first method will perform these operations on items in this panel only, while the second method will also include selected items in the main frequency list.

- To perform one of the above-mentioned operations on the suggested items only, select the check boxes of one or several suggestions, right-click your mouse, and choose the appropriate command.
- To include with the suggestions, currently selected words in the frequency table, click the ≡ button on the toolbar and choose the appropriate command. You may also drag and drop words from the frequency table to the appropriate location in the **Dictionary** panel on the left. When dropping items from the frequency table, all selected suggestions will also be moved to the selected location.

Barcharts, Pie Charts, and Word Clouds

WordStat allows you to produce bar charts or pie charts to visually display the distribution of specific keywords or categories.

To produce charts:

- Move to the Frequencies page.
- Set the Sort By option to the desired graphic order of the values.
- Select the rows you would like to plot (multiple but separate rows can be selected by clicking while holding down the Ctrl key)

Click the 🕌 button.

Types of Charts

Three types of charts may be used to depict the distribution of keywords or content categories:

- The word cloud is an image composed of keywords or content categories, sometimes phrases in which the size of each items indicates its frequency.
- The vertical bar chart is the default chart used to display absolute or relative frequencies of keywords or content categories.
- The horizontal bar chart displays the same information as the vertical one but is especially useful when the number of keywords is high and their labels cannot be displayed entirely on the bottom axis.
- The pie chart is useful to display the relative frequency of each keyword and compare individual values to other values and to the whole. Numerical values displayed in pie charts are always expressed in percentages of either the total frequency or case occurrences.
- The donut chart, similar to the pie chart, is useful to display the relative frequency of each value and compare individual values to other values and to the whole. Donut charts are considered easier to read as you can focus on reading the length of the arcs, rather than comparing the proportions between slices.

The **Plot** option allows you to select the values that will be used as the scale for the length of bars in bar charts or as the percentage base for pie charts.

For bar charts the options are:

% SHOWNPercentage based on the total number of keywords displayed in the table% PROCESSEDPercentage based on the total number of words encountered during the analysis% TOTALPercentage based on the total number of words that have not been excludedNO OF CASESNumber of cases where this keyword appears% CASESPercentage of cases where this keyword appears	FREQUENCY	Number of occurrences of the keyword
% PROCESSEDPercentage based on the total number of words encountered during the analysis% TOTALPercentage based on the total number of words that have not been excludedNO OF CASESNumber of cases where this keyword appears% CASESPercentage of cases where this keyword appears	% SHOWN	Percentage based on the total number of keywords displayed in the table
% TOTALPercentage based on the total number of words that have not been excludedNO OF CASESNumber of cases where this keyword appears% CASESPercentage of cases where this keyword appears	% PROCESSED	Percentage based on the total number of words encountered during the analysis
NO OF CASESNumber of cases where this keyword appears% CASESPercentage of cases where this keyword appears	% TOTAL	Percentage based on the total number of words that have not been excluded
% CASES Percentage of cases where this keyword appears	NO OF CASES	Number of cases where this keyword appears
	% CASES	Percentage of cases where this keyword appears

For pie charts, two options are available to specify how percentages will be computed:

FREQUENCY	Percentage based on the total frequency of keywords
NO OF CASES	Percentage based on the total number of case occurrences

The **View Others** option displays an additional bar or slice representing all items in the frequency table that have not been selected.

The Toolbar

The following table provides a short description of available buttons and controls on this page:

Control: Description:

1	Press this button to append a copy of the graphic in the Report Manager. A descriptive title will be
	provided automatically. To edit this title or to enter a new one, hold down the SHIFT keyboard key while
	clicking this button (for more information on the Report Manager, see the Report Management Feature
	topic).

- Press this button to retrieve a chart previously saved on disk.
- Press this button to save a chart on disk. Charts may be saved in BMP, JPG or PNG graphic file format or may be saved in a proprietary format (.WSX file extension) that may later be edited and customized using the Chart Editor.
- Pressing this button prints a copy of the displayed chart.
- Click this button to turn on/off the 3-D perspective for the current chart.
- This button allows the editing of various features of the chart such as the left and bottom axis, the chart and axis titles, the location of the legend, etc. (see below for available options)
- This button is used to create a copy of the chart to the clipboard. When this button is clicked, a pop-up menu appears allowing you to select whether the chart should be copied as a bitmap or as a metafile.
- Pressing this button closes the chart dialog box and returns to WordStat's main screen.

Customizing the Chart

Clicking the *local content* bottom displays a dialog box that allows you to customize the appearance of bar charts and pie charts. The options available in this dialog box represent only a small portion of all settings available.

Left or Bottom Axis

Minimum / Maximum: WordStat automatically adjusts the vertical axis scale to fit the range of values plotted against it. To manually set these values, type the desired minimum and maximum.

Increment: Increasing or decreasing this value affects the distance between numbers as well as tick marks. Horizontal grid lines are also affected by the modification of this value.

Grid: This option turns horizontal grid lines on and off. Grid lines extend from each tick mark on an axis to the opposite side of the graph. To increase or decrease the number of grid lines or the distance between these lines, change the increment value of the axis. A list box also allows a choice among five different line styles to draw these grid lines.

Titles

Proper titles and axis labels are of utmost importance when describing the information displayed in a chart. By default, WordStat uses variable names and labels as well as other predefined settings to provide such descriptions.

The title page allows you to modify the **top** title, as well as the labels on the **left**, **bottom** and **right** axis. To edit the title, select the proper radio button. Enter several lines of text for each title by pressing the **<Enter>** key at the end of a line before entering the next line.

The Font button on the right side of the edit box allows changing the font size or style of the related title.

3D View

Orthogonal: Turning this option off disables the free elevation and rotation of the 3-D chart.

Zoom: This option zooms the whole chart. Expressed as a percentage, increasing the value positively will bring the chart towards the viewer, increasing the overall chart size as the Zoom value increases.

3-D Percent: The 3-D Percent property indicates the size ratio between chart dimensions and chart depth by specifying a percent number from 1-to-100.

Perspective: Use this property with Orthogonal unchecked to modify the 3-D perspective of the Chart. Larger values add more depth perspective.

Bar shadow: Enabling this option adds a shadow to the sides of 3-D bars. Turning it off will color the sides of the bar the same as the front.

Bar width: This option determines the percent of total bar width used. Setting this value to 100 makes joined bars.

Bar depth: Use this property to limit the depth that each bar series uses. By default, bars will take up the part proportional to the number of bar series in the chart so that the back of a bar will join the front of the bar immediately behind it. To insert a gap between a series of bars, decrease this value.

Pie depth: Use this property to change the thickness of the pie chart.

To further customize the chart, modify data points or value labels, click the ⁴⁴ button located on the right side of the dialog box.

Keyword Retrieval

The **Keyword Retrieval** feature can retrieve any document, paragraph or sentence containing a specific keyword, a combination of keywords or no keyword at all. Text units corresponding to the search criteria are returned in a table on the **Results** tab. This table may then be printed or saved to disk. It may also be used to create tabular or text reports as well as to attach QDA Miner codes to the retrieved text segments.

To start the Keyword Retrieval feature:

• Go to the Frequencies page and click the 🐒 button. This dialog box can also be accessed from other parts of the program to retrieve text units associated with a cluster or with a specific association. It may also be accessed by selecting an item in the frequency page, right-clicking and selecting **Keyword Retrieval** from the pop-up menu. When calling this function, a dialog box similar to this one appears, allowing you to specify the desired search criteria:

C Keywo	ord Retrieval			- 🗆 ×
Criterion	Results			
Retrieve:	Paragraphs 🗸 🗸			
Keyword	filtering			
	If: Include at least 🗸	1 🖨 of these keywords: [M	ILITARY, POWER, SECURITY]	~
				1.00
More	2			
Variable (filtering			
Variable 1	Variable	Onersten	Citation	
		operator:	[02-2009_04-2009]	
-	DELIVERI	equais	[Q3-2000,Q+2000]	
ANU	<none></none>			×
Report				
Add v		tvenvl		107
Aug Vi	CANDIDATE, DEL	IVERT		· · ·
Source	e variable 🗌 Paragraph	& sentence numbers 🛛 Wor	d count	
				Search

Retrieve: This option determines the text unit on which the search will be performed as well as what will be retrieved. You can select three different text units. The **Documents** search unit allows WordStat to apply the search expression on each document associated with a specific case and, if a specific document meets the search condition, its location will be displayed. The **Paragraphs** search unit allows WordStat to display any paragraph meeting the search condition. The **Sentences** search unit will instruct WordStat to return sentences meeting the search condition.

Keyword filtering: This group of options allows you to select the keywords on which the retrieval will be based. Setting the first list box to **No Keyword** will retrieve all text units for which no keyword has been found. This option is especially useful to identify topics or themes that have not been covered by the current categorization dictionary or new keywords that should be added to enhance the coverage of existing categories in the dictionary. Setting the filtering to **Include at least** allows you to retrieve text units containing a minimum number of keywords from a selected list. To select the keywords, click the arrow button to show all available keywords and then click the desired items. The minimum number of keywords a specific unit must contain in order to be retrieved is specified using a small edit box with spin buttons on the right. Setting this numerical value to the total number of keywords selected will force the program to retrieve only those units containing all those keywords. Setting this number to a lower value will retrieve all text units containing at least this number of keywords from the selected list. In the example shown above, any unit containing keywords from at least one of the three categories AMENDMENT, CIVIL, or DIPLOMACY will be retrieved.

To enter a second filtering condition, click the More button. You can choose to link the two filtering conditions using either one of the three Boolean expressions: AND, OR or NOT. Choosing AND will retrieve all text units fulfilling both criteria; selecting OR will result in a retrieval of text units meeting either the first or the second condition, or both, while choosing the NOT Boolean operator will retrieve text units meeting the first condition but not the second one.

To limit the filtering conditions to a single statement, click the **Less** button.

Variable filtering: The second group of options allows you to restrict the retrieved text units to specific cases selected according to some logical condition. This filtering condition may consist of a simple expression, or may include up to two expressions joined by a logical operator (i.e., AND, OR). In the above screen shot, only text units from cases where the variable CANDIDATE is equal to PARTY will be retrieved.

The following table shows the various operators available for each data type:

DATA TYPE	AVAILABLE OPERATORS
NOMINAL / ORDINAL	Equals Does not equal Is empty Is not empty
NUMERIC and DATE	Equals Does not equal Is greater than Is lesser than Is greater than or equal to Is lesser than or equal to Is empty Is not empty
BOOLEAN	Is true Is false
STRING	Contains Does not contain Is empty Is not empty

Add variables: This drop-down checklist box may be optionally used to add the values stored in one or more variables to the table of retrieved segments for the specific case from which a text segment originated.

• Once all the search options have been set properly, simply click the search button to retrieve the selected text units.

Working with the Retrieved Text Units

The retrieved text units are displayed in a table found on the **Results** tab.

CODE:	🖂 🖉 🕂 🏢 🗹 Multilines grid						9991 Q2	2 - 10	
Case #	Text	Matching	Case	CANDIDATE	DELIVERY	Variable	Paragraph	Word Cour	it ^
63	Nost importantly, this plan will ensure that we control the energy we use with resources and technology that are available today. The steps I just spoke about are not far-off, pie-in-the-sky solutions, they are now. Today, there are waiting lists for fuel-efficient cars. There's an old steel mill in Pennsylvania that has become the home of a new wind turbine factory. Two seen a small business in Nevada powered entirely by solar power. Across the planet, countries like Germany and the United Kingdom have already implemented dean energy polices that are reducing their carbon emissions right now, and leaders like Tony Blair and Angela Merkel have done a great job of raising the visibility of climate change within the G8. Now it's our turn to lead – to show that this future is possible for America.	POWER	1-Obama20080711	Obama	Q3-2008	DOCUMENT	23	143	
64	Qaeda, the Taliban, and all of the terrorists responsible for 9/11, while supporting real security in Afghanistan.	SECURITY	1-Obama20080715	Obama	Q3-2008	DOCUMENT	2	16	
64	Our men and women in uniform have accomplished every mission we have given them. What's missing in our debate about Iraq - what has been missing since before the war began - is a discussion of the strategic consequences of Iraq and its dominance of our foreign policy. This war distracts us from every threat that we face and so many opportunities we could seize. This war diminishes our security, our standing in the world, our military, our economy, and the resources that we need to confront the challenges of the 21st century. By any measure, our single-minded and open-ended focus on Iraq is not a sound strategy for keeping America safe.	MILITARY; SECURITY	1-Obama20080715	Obama	Q3-2008	DOCUMENT	11	113	
64	I am running for President of the United States to lead this country in a new direction - to seize this moment's promise. Instead of being distracted from the most pressing threats that we face, I want to overcome them. Instead of pushing the entire burden of our foreign policy on to the brave men and women of our military, I want to use all elements of American power to keep us safe, and prosperous, and free. Instead of alienating ourselves from the world, I want America - nore apain - to lead.	POWER; MILITARY	1-Obama20080715	Obama	Q3-2008	DOCUMENT	12	91	~

This table contains the case number and the variable from which the segment originates, as well as the value of all additional variables selected by the user (see the **Add variables** option above). When searching for paragraphs or sentences, the table also displays the text associated with the retrieved unit and its location (its paragraph and sentence number). By using arrow keys or by clicking a row the associated text is displayed in, a separate window at the bottom of the screen with all keywords in bold. Selecting specific words or phrases in this text window by right-clicking displays a pop-up menu to assign them to a content category, the exclusion list or to obtain a keyword-in-context list.

To sort the table of retrieved units in ascending order in any column, simply click the column header. Clicking the same column header a second time sorts the rows in descending order. Tables may also be printed, stored as a text report, or exported to disk in various file formats such as Excel, ASCII, or HTML.

To remove a search hit from the hit list:

• Select its row and then click the in button.

To assign a QDA Miner code to a specific search hit:

- In the table of search hits, select the row corresponding to the text segment you want to code.
- Use the **CODE** drop-down list located above this table to select the code to assign.
- Click the A button to assign the selected code to the highlighted text segment.

To assign a QDA Miner code to all search hits:

- Use the **CODE** drop-down list located above this table to select the code to assign.
- Click the *button* to assign the selected code to all text segments matching the search expression.

NOTE: To automatically attach QDA Miner tags to all paragraphs or sentences associated with all currently displayed content categories or keywords, you may use the autocoding feature by clicking the substantiation on the **Frequencies** tab.

To create a new QDA Miner code:

 Click the + button. A dialog box will appear allowing you to create a new QDA Miner code (for more information see Adding a QDA Miner codes).

To obtain a word cloud and perform a word frequency analysis on retrieved text segments:

• Click the we button. The analysis will be performed on all text segments retrieved and displayed in the table. (See <u>Word Frequency Analysis</u> for more information on this feature).

To save the table to the Report Manager:

• Select the 🚺 button. A copy of the table will be appended in the **Report Manager**.

To create a report of coded segments:

• Click the button. The sort order of the current table is used to determine the display order in the report. This report is displayed in a text-editing dialog box and may be modified, stored on disk (in RTF, HTML or plain text format), printed, or cut-and-pasted into another application. Graphics and tables may also be inserted anywhere in this report.

To export the table to disk:

- Click the 🖬 button. A Save File dialog box will appear.
- In the **Save As Type** list box, select the file format under which to save the table. The following formats are supported: ASCII file (*.TXT), Tab delimited file (*.TAB), Comma delimited file (*.CSV), HTML file (*.HTM; *.HTML) and Excel spreadsheet file (*.XLS).
- Type a valid file name with the proper file extension.
- Click the SAVE button.

To print the table:

Click the button.

To close the keyword retrieval dialog box:

Click the 💐 button.

Running and Creating Post-Processing Scripts

The post-processing scripting feature allows one to extend the capabilities of WordStat by the writing of additional computation or data transformation scripts. Those scripts may be written in Python or R using custom codes or existing libraries for NLP, machine learning, statistical processing, visualization, and more. Such a feature allows data scientists to focus on the crucial processing step of interest, leaving to WordStat the various tasks associated with the data preparation and preprocessing such as document importation, stemming, lemmatization, spelling correction, removal of stop words, and so on. It also allows one to combine WordStat's powerful categorization features (with word patterns, phrases, disambiguation rules, etc.) with new analytics techniques not available in WordStat.

The additional possibility to design dialog boxes allows data analysts to create simple graphic users' interfaces that may be used for customizing script execution and facilitate their use by other WordStat users with no programming skills.

Running an existing Python or R post-processing script:

On the Frequencies page, applying an existing post-processing script can be as simple as these steps:

Click the
 button to open the main post-processing script window.

The **Script** dropdown menu lists all existing Python and R scripts. Select the desired script. Beneath the **Script** dropdown menu is the **Description** section. Any text describing the script will be displayed here. Resizing the post-processing script window may be required to display all the description text,

Script:	R Sample R - Readability Statistics	~	Ne Run	
cription:	Classification Binary Logistic Regression with sklearn Classification Binary Neural Network with sklearn Classification Binary Random Forest with sklearn			
	Classification Multilabel SVM with sklearn			
	R Document Clustering - K-Mean Detection POS Tagger Sample Python - Descriptive Statistics			
	Sample Python - Readability Statistics	ł		-
	R Sample R - Readability Statistics	ľ		
	R Sample R - Word Statistics			
	Sentiment Analysis			
	Topic Modeling - LDA with Stikit Learn Topic Modeling - Structural Topic Model			

• Click the button. For some scripts, no other steps are required, and the Report Manager will display the output of your script. Running other scripts will result in the appearance of a dialog box like the one below, allowing you to set some analysis options.

Options	- 🗆 X
Integer:	2
Floating Point:	0.0002
Section	
String:	
	Boolean
List of options:	choose_one 🗸 🗸
	a second second
	Protocol and a second sec

• While the script is running, a Python console or R console will appear, displaying what is logged or printed as output (stdout) or any error messages.

Output

- If a script is executed successfully, WordStat will import any output stored in any supported file format and display them in the Report Manager, which will open automatically. WordStat will import text output stored in any file with a .TXT file extension, assuming an UTF-8 text encoding, or documents stored in .DOCX, .RTF or HTML files. It will also import table output files with either .CSV, .TAB, .TSV, or .XLS file extension, and any graphic produced with a .JPG, .JPEG, .PNG, .GIF, or BMP file extension. Once imported, those files are deleted.
- If a table created by the script has a column named **RECNO** containing record numbers, a dialog box will appear, giving the possibility to append data contained in this table to the current project data.

The file RESULTS-EXAMPLE.TAB of	ontains record specificic information
Do you want to append 1 new var	riables and replace 0 existing ones?
✓ Yes	No No

To append data in this table, click **Yes**. If the table's column names correspond to existing variables, data into those variables will be overwritten, while new column names will result in new variables being appended. Clicking **No** will append the table to the report manager.

Writing a new Python or R post-processing script

- To create a new post-processing script, open the main post-processing script window by clicking the **Section** button on the **Frequencies** page.
- Then click the button next to the **Run** button and select **New script** from the dropdown menu.
- Select the desired programming language for this script (Python or R)
- A Script Options window will appear. Type the Name of the new script.
- Optionally, add a **Description**. This description will appear when the script is selected from the script selection dropdown box.

100 million (1			
Name:	Example Script		_
Description:	Describe the script here		
Document x Tern	n matrix (input.tab): Frequencies Case occurrences Rate per 10,000	11 ∑	XIDF

- The options below allow you to select the types of input files that will be generated and processed by the script. Up to three separate data files can be generated from a choice of seven options.
 - The input.tab data file contains numerical values resulting from the quantification of text data by WordStat
 where each row represents a document, and columns consist of frequency information for every item
 displayed in the Frequencies tab (i.e., either words or content categories). When no categorization
 process is applied, such a data file corresponds to a Document x Term matrix. A choice of one of four
 statistics may be stored in this data file: the term's frequencies, the case occurrences (i.e., either 0 or 1),
 The rate of this term per 10,000 words, or its TF-IDF score. If more than one of these metrics is checked,
 the user will be asked to select from these options upon execution.
 - The **sources.tab** is the original text data stored in a single file, one document per line. To store a document in a single line, carriage return (ASCII #13) and line feed characters (ASCII #10 have to be replaced with other characters (ASCII 30 and 31).
 - The **tokens.tab** will hold the result of the project's text processing steps including preprocessing, word replacements, stop word removal and categorization. By default, each document is stored on a single line. One may also segment this file by paragraphs or by sentences.
- Once the script options have been set, click **OK** to open the script editor window.

- show	Modal					-	×
File Edit	Run						
		A T					
Variable	Туре	Prompt	Description	Options	Returned value		
				-annot-			_
1	_			mm			
1	-			-01001			
1 1							
1							
1				-99900			

- The Variables section at the top allows you to define variables that can be used in the script to customize its execution. Upon execution of a script, if variables have been defined, a dialog box will be presented to set those options.
- To add a new variable to the script, click the 🔜 button. The variable definition window will appear.

Variable D	efinition					-		×
Type:	Floating point	t 🗸 🗸						
Name:								
Prompt:								
Description:								_
Range:	Minimum:	0,0000 🜲	Maximum:	1.0000 🗘	Default:	0.000	00 ≑	
Range:	Minimum:	0,0000 🖨	Maximum:	1.0000 🖨	Default:	0.000	00 🜲]

- From the **Type** dropdown menu, select what will be the type of the variable. You may select among six types of variables: an integer, a floating point, a string, a Boolean, a list of options or a project variable. An additional **Section** type may also be used to group various options into distinct sections. When such an option is selected, you will be asked to enter a string that will be displayed in bold character as the title of the section.
- In the **Name** edit box, type the name of this variable. This is the name that should be inserted into the script source code to customize its execution.
- The **prompt** edit box is the text that will be displayed on the left of the data editing control. The **Description** edit box may also be used to provide additional information about this option. If a description or instructions are provided, a hint window will display this information upon hovering over the variable data entry control.
- Depending on the variable type, different specifications can be added:
 - For a **Floating Point** or an **Integer** variable, the range section allows you to set a **Minimum**, **Maximum** and/or **Default** value. First check the box for the desired range specification, then increase or decrease the value via the up/down arrows, or by typing in the numerical value.
 - A Boolean variable can be enabled by default, by checking the default box under the **Description** textbox.
 - For a List of options, the strings that will appear as options to be selected must be specified in the **Options** textbox. Each option should be added on a separate line without quotation marks.
 - If a script includes a **Project variable**, choose the type of data that will be expected via the **Data type** dropdown list. Choosing a data type will restrict the list presented to the selected data type. Selecting **Any type** will present a list of all variables on the project. Checking the **Optional** checkbox allows the Project variable to be left blank.
- To remove an existing variable, first select the variable to be deleted and then click the **s** button.

- Select a variable and click the so edit variable button to open the variable definition window and make any changes to the existing specifications.
- Click the 4 down arrow to change the order of the variables and bring the selected variable down in the list of variables, or the 1 up arrow to move the variable up.
- The enview button shows the **Options** dialog box as it would display when the script is run, without having to run the script. The dialog box can also be previewed by selecting **Test Dialog Box** under the **Run** menu.

In the example below, the dialog box consists of six variables grouped under two sections: Classifier Parameters and Evaluation.

		Т				
Variab	le	Туре	Prompt	Description	Options	Returned value
		Section	Classifier Parameters			1.00
depend	ent_variable	Project variable	Dependent variable		String or Norminal/Ordinal	
classifie	r_type_parameter	List of options	Classifier type		SVC;NuSVC;LinearSVC	
dassifie	r_kernel_parameter	List of options	Classifier kernel type		rbf;poly;sigmoid;linear	
		Section	Evaluation			
shuffle	_samples_parameter	Boolean	Shuffle each class's samples		True	
nb_split	s_parameter	Integer	Number of folds		Min=2 Max=10 Def=5	
test_siz	e_parameter	Floating point	Proportion to include in the test split		Min=0.1000 Max=0.5000 Def=0.7000	
1 2 3 4 5 6 7 8 9 10	import time import numpy import panda import pickl from sklearn from sklearn from sklearn	<pre>v as np s as pd e import svm .model_sele .model_sele .model_sele</pre>	n, metrics ection import GridSearch import make_pipeline ection import train_test	CV, Strat.	ifiedKFold, learning_curve I	

Running the script will cause the following dialog box to appear:

7		×
CANDIDATE	*	
SVC	~	
rbf	~	
Shuffle each o	dass's sample	s
5		
0.7		
	CANDIDATE SVC rbf Shuffle each o 5 0.7	− □ CANDIDATE ✓ SVC ✓ rbf ✓ Shuffle each class's samples 5 € 0.7
The bottom section of the script editor window should contain valid Python or R code, and the formatting will reflect
the syntax of the script's language. Proper commands for importing needed libraries should be specified at the
beginning of the script. Names of the specified variables may be used in the script to customize its execution. The
Python script below illustrates how one may use the defined variables in a script (highlighted in yellow) and their
associated values set by the user using the dialog box.



- Clicking the ▶ button will run the script. Any changes made will prompt a dialog box asking if changes should be saved before running. Running a script from within the script editor opens a Python or R console. In contrast to running a script from the main post-processing scripts window, the console will also display the script's code in the upper half of the window. The code will include all the variable assignments. Once the script terminates, the console window will have to be closed for the output to be processed. The script can also be run by selecting **Run Script** under the **Run** menu.
- A new script can also be created by selecting **New** under the **File** menu of the script editor. There is an option to choose between Python or R.
- Save a new script with the labor by selecting **Save** from the **File** menu. Selecting **Save** as will open the **Script Options** window allowing you to enter a new name for the script.

Editing a script

• To edit an existing script, first select it from the **Script** dropdown menu, then click the button next to the **Run** button and select **Edit Script**.

- Variables can be added, deleted, or edited (See details above)
- Script Options of an existing script can also be edited by selecting **Settings** under the **Edit** menu.
- Changes can be made to the code of the script in the Code section of the script editor. Any edits can be typed into this section or using the Edit menu. Under the Edit menu, you will find common edits such as Undo, Cut, Copy, Paste, Select all.
- Click Find to search for any particular sequence in the code.

Find		×
Find what:		End Vest
Match whole word only	Direction	Cancel
Match case	O <u>U</u> p ● <u>D</u> own	

• Select Replace in the Edit menu to substitute a sequence in the code.

Replace	2
Find what:	Find Neid
Replace with:	Replace
Match whole word only	Replace All
Match case	Cancel

Importing/Exporting a script

Post-processing scripts created in WordStat (whether Python or R) are saved as '.wscr' files in a distant folder not easily accessible. The import and export features have been designed to easily copy an existing script to this folder or create a copy outside of this folder, allowing one to share scripts with other users.

- To import a script, click the button next to the run button and select Import. Select an .wscr script from the file explorer and click Open. A script can also be imported by selecting the **Import** command from the **File** menu in the script editor.
- To export a script, first select it, click the 🖾 button next to the RUN button and select Edit. Then choose the Export command from the file menu of the script editor. Choose its destination in the file explorer and click Save

Console

Any libraries that are imported in a post-processing script will automatically be installed as the script is run. In some cases, a package requires extra steps for installation, which Wordstat cannot perform under-the-hood. A Python or R console can be opened to perform these installations or any other tasks.

- From the Script Editor window, select Console from the Run menu,
- Depending on the language of the script, a Python or R console will open

The Extraction Tab

The Extraction tab groups tools that are used to extract useful features from the text collection.

The <u>Topic</u> modeling tool will automatically extract the most important topics from a text collection using factor analysis. Results may be saved as a content analysis dictionary or may be further examined using cooccurrence analysis or crosstabulation.

The <u>Phrase</u> extraction feature will identify idioms and common phrases and will allow you to add them to a content analysis dictionary as well as perform cooccurrence analysis and comparison analysis of the phrases.

The <u>Named Entities</u> extraction feature can identify proper nouns, names of people, locations or organizations as well as acronyms. You can select relevant items and move them to the categorization dictionary.

The <u>Misspellings & Unknowns</u> extraction feature provides a tool that identifies misspellings and some technical terms by comparing the list of word forms encountered in the entire text collection against a list of common words. Extracted words may be added to the current categorization dictionary or to a substitution process. They may also be replaced in the original documents with properly spelled words.

Topics

The **Topics** extraction feature, the first tab on the **Extraction** tab, attempts to uncover the hidden thematic structure of a text collection by applying a combination of natural language processing and statistical analysis. The main statistical procedure used for topic extraction in WordStat is a factor analysis. Technically speaking, the extraction is achieved by computing a word by document frequency matrix, or alternatively by segmenting documents into smaller chunks and computing a word by segment frequency matrix. Once this matrix is obtained, Wordstat performs either a non-negative matrix factorization (or NNMF) or a factor analysis with Varimax rotation, in order to extract a small number of factors. All words with a loading higher than a specific criterion are then retrieved as part of the extracted topic. While in hierarchical cluster analysis, a word may only appear in one cluster, topic modeling using factor analysis may result in a word being associated with more than one factor, a characteristic that more realistically represents the polysemic nature of some words as well as the multiple contexts of word usage.

To ensure the stability of the factoring solution, low frequency items should preferably be excluded. It is thus strongly recommended to remove any word occurring less than 10 times on smaller data sets, ideally less than 30-to-50 times on larger ones. Stemming, lemmatization or the creation of a categorization dictionary may also be used to group words or phrases, including less frequent ones, prior to the topic extraction.

To extract topics:

Select the ¹¹/₁ button. An Extraction Setting dialog opens similar to the one below.

extraction Settings	×
Topics Phrases Named Entities Misspellings	
Method: NNMF O Factor Analysis	
Segmentation: by Paragraph	
Loading: 0.22	
Seed: Random O Fixed: 1958625732	
Topic Enrichment	
Default confidence level: High 🗸 🗸	
Check for misspelling	
Use topic weights (when saving or analyzing topics)	
	Close

The **1**¹ button is present on all of the tabs on the **Extraction** tab. Each tab of this dialog contains extractions options for the currently select tab.

Method: WordStat implements two methods for topic extraction. The **NNMF** method uses non-negative matrix factorization to extract topics from a word x word correlation matrix computed by WordStat, while **Factor Analysis** perform a similar extraction using a principal component analysis with VARIMAX rotation. The NNMF method is faster and can handle larger matrices than factor analysis. It is however probabilistic in nature, yielding different, yet likely similar, solutions every time you run it, while factor analysis will always produce the exact same results.

Segmentation: This option allows you to specify whether the data to be used for topic modeling will be based on the cooccurrence of words in the same **document**, or whether it will be based on cooccurrence within **paragraphs**, within **sentences**, or within a **windows of words**. The choice of segmentation should ideally reflect how topics are being distributed in a typical document and across documents, as well as the objective of the analysis. When the text collection consists of long documents containing multiple topics (such as long political speeches) and you need to identify all topics in order to compare their relative frequencies, then performing a segmentation by paragraph or by sentence may be more sensitive than computing cooccurrences by documents. Alternatively, if you attempt to differentiate documents by identifying domains or disciplines, or to identify the dominant issue of documents, then performing the analysis at the document level may be more appropriate. When analyzing responses to open-ended questions, which may include several topics listed in a single paragraph, segmenting by sentence may also result in a more precise extraction of the various topics they contain. Finally, if long documents have been imported without any paragraph delimiters, setting a windows of words may be needed to extract meaningful topics.

Loading: This option allows you to set a minimum topic loading a word should reach in order to be retained in the factor solution. By default, this value is set to 0.3. Increasing the cutoff value will reduce the number of words, keeping only the more representative ones, while reducing it may include words that are somewhat less characteristic of the extracted topic.

Seed: This option, which is available only when NNMF method if selected, allows you to set a specific seed value in the pseudo random number generator of WordStat, allowing you to reproduce the exact same results on every run, overriding the probabilistic nature of NNMF. By default, it is set to **Random**, so that repeated analysis on the same data set with the same settings will yield different results. You can achieve replicability of the results by specifying a

Fixed seed value in the WordStat's random number generator. Clicking the *button* generates a random number that will become the new fixed seed value.

Pruning: Factor analysis on large correlation matrices can be time consuming. The factor analysis routine in WordStat is limited to a maximum of 2500 words. To speed up the extraction of topics using factor analysis or to analyze number of words beyond this limit, WordStat implements a pruning technique by which words that are less likely to be included in a factor solution will be removed from the matrix prior to the analysis. Disabling this function prevents the removal of those words and attempts a topic extraction on all words, up to the upper limit of 2500 words.

Topic Enrichment: Topic modeling solutions usually consists of a series of words listed in descending order of topic specificity. The combined presence of those words in a text is used to identify whether a topic is present or not. This "bag-of-word" approach may result in imprecise topic identification and measurement since words often have multiple meanings and are often used in very different contexts. The topic enrichment feature allows you to go beyond this "bag-of-word" approach by attempting to identify some phrases highly associated with the extracted topic, as well as other phrases that may represent exceptions and help disambiguate the various meanings of words. WordStat topic enrichment feature also attempts to reduce false negative that may result from the presence of misspellings by identifying misspelled forms of words in the original topic solution or that are part of any of the suggested phrases. Enabling the **Topic Enrichment** feature gives access to the following options:

Default confidence level: This option allows you to set a confidence level that will be used when adding phrases to topics. Selecting a lower setting (such as **Low** or **Moderate**) tends to bring more phrases than if this option is set to a higher setting. Phrases that are not automatically added to the topics will still be accessible in the **Suggestions** panel on the left from which they can be manually selected and added to the topic.

Check for misspelling: When this option is enabled, WordStat will identify potential misspellings of words in the topic or that are part of any added or suggested phrases. This feature helps reduce the number of false negatives resulting from spelling errors in the original data set.

Use topic weights: Words in extracted topics are presented in descending order of relevance of specificity to the topic. The first words have higher weights than the last ones. When using factor analysis for topic extraction,

the weight of each word in a topic corresponds to its factor loading, while in NNMF it corresponds to a coefficient obtained from the product of the W and H matrices. By default, when running a cooccurrence analysis or performing a crosstab on the extracted topics, or when saving topics to a categorization dictionary, each word will be saved along with its weight, while phrases will get a weight equal to one. Disabling this feature will assign a weight of one to all entries.

- Set the extraction settings and click on the Close button.
- Set the number of topics you would like to extract by typing the number in the **No. topics** field or by using the up or down arrows.
- Once the options have been set, click the button to perform the analysis. Please note that extracting topics on more than a few thousand words can take several minutes. Once extracted, the **Topics** tab should looks like this:

Da	ata Text Processing	Frequencies 🚯 Extraction 🕥 Cooccurrences 🏫 Gro	osstab 🏢 Key	word-In-C	ontext	Classification	1		0
No. t	topics: 50		ve					1	4 44
NO	TOPIC	KEYWORDS	COHERENC	FREQ	CASES	% CASES	🔪 🛀 Comparisons 😳 Suggestions		
38	RENEWABLE ENERGY	SOLAR; WIND; REINEWABLE; SOURCES; COAL; CLEAN; EMERGY; POWER; ELECTRICITY; TECHNOLOGIES; INVEST; CLEAN ENERGY; REINEWABLE ENERGY; SOURCES OF ENERGY; REINEWABLE SOURCES; NUCLEAR POWER; EMERGY NDEPENDENCE; SOLAR POWER; CLEAN COAL; ENERGY PUTURE ENERGY EFFICIENCY; ENERGY EFFICIENT; ENERGY POLICY; ALTERNATIVE ENERGY; ENERGY EFFICIENT; ENERGY POLICY; ALTERNATIVE ENERGY; ENERGY SOURCES; STRATEGIC ENERGY FUND; RENEWABLE SOURCES OF ENERGY;	0.912	840	136	56.43%	Biden Clinton Edvards	~	-
14	DOCTORS	DOCTORS; PATIENTS; NURSES; QUALITY; PATIENT; MEDICAL; CARE; HEALTH CARE; HEALTH CARE SYSTEM; QUALITY CARE; MEDICAL RECORDS; UNIVERSAL HEALTH CARE: MEDICAL CARE; MEDICAL RECORDS; QUALITY HEALTH; QUALITY OF CARE; MEDICAL RECORDS; HEALTH; QUALITY OF CARE; MEDICAL RECORDS; HEALTH CARE PLAN; MEDICAL RECORDS; HEALTH CARE PLAN;	0.870	604	109	45.23%	Giuliani Kucinich MCCain Obama-		
21	HEALTH CARE	HEALTH; CARE: INSURAINCE, COVERAGE; COSTS, UNIVERSAL; ARFORDABLE: COVER: PLAN; QUALITY; COST; AMERICANS; COMPANIES; HEALTH CARE; HEALTH INSURANCE: INSURANCE COMPANIES; HEALTH CARE; SYSTEM; UNIVERSAL HEALTH CARE; HEALTH CARE COSTS;	0.832	2519	204	84.65%	Romney- 0 20 40 60 80 100	120	140
43	TAX CUTS INCOME	INCOME; TAXES; TAX; RAISE; CUT; LOW; CUTS; CAPITAL; RELIEF; CREDIT; TAX CUTS; TAX CREDIT; TAX BREAKS; TAX CUT; TAX CODE; INCOME TAX; BUSH TAX CUTS; CAPITAL GAINS; RAISE TAXES; TAX RATES; TAX RELIEF; TAX RATE; INCOME TAXES;	0.779	864	133	55.19%	240 220 200 180 160 140		/
26	NUCLEAR WEAPONS IRAN	TRAN; NUCLEAR; WEAPONS; TRANIAN; DIPLOMACY; ISRAEL; NORTH; RUSSIA; INTERNATIONAL; NUCLEAR WEAPONS; NORTH KOREA; NUCLEAR AMBITIONS; NUCLEAR POWER; NUCLEAR PROGRAM; INTERNATIONAL COMMUNITY; RANIAN REGIME;	0.776	707	96	39.83%		4	
19	AL QAEDA	QAEDA; AL; PAKISTAN; AFGHANISTAN; TERRORISTS; AL QAEDA; FIGHT AGAINST AL QAEDA; AFGHANISTAN AND PAKISTAN; FINISH THE FIGHT;	0.772	470	75	31.12%	2000 01-2001 02-2001 03-2001 04-2001 01-2008 02-2008	2-2008 of	A-2008

The table to the left contains the following information:

TOPIC	WordStat uses an algorithm to automatically provide a label for the extracted topic.
KEYWORDS	Lists all words meeting the factor loading cutoff criteria in descending order of factor loading, along with associated phrases.
COHERENCE	The coherence is a weighted average of the correlations of words associated with the topic.
EIGENVALUE	When performing topic modeling using factor analysis, this column contains the eigenvalue of each factor.
FREQ	Displays the total frequency of all items listed in the keywords column.
CASES	Shows the number of cases containing at least one of the items listed in the keywords column.

Displays the percentage of cases with at least one of the items listed in the keywords column.

Running multiple topic models

Choosing the number of topics to extract using topic modeling techniques remains a question for which there is, to our knowledge, no definitive answer. We may even raise doubts about the fact that such an optimal number exists. In fact, one may even suggest that information obtained using different settings may well serve different purposes or reveal different aspects of a reality. In such a context of uncertainty, researchers often want to compare various solutions. The batch processing feature allows one to compute multiple topic models by systematically varying the number of topics to extract, and to perform several runs using the same settings in order to assess the stability of the results.

To run multiple topics, click the I button, a dialog box similar to this one will appear:

Topic model gen	erator	×
Method:	NNMF (Segmentation: by paragraph, Loading: 0.30)	
Extract from:	5 🔹 to: 50 🔹 topics by increment of: 5	÷
Replications:	3 For NNMF with random seeds	
	VOK X Cancel	

The **Method** option displays the current topic extraction technique and its associated settings. To change any of those values, click this description to display the Topic Extraction Settings dialog box.

The next line of options allows one to set the smallest and largest number of topics to extract as well as the increment to be used. In the example above, the settings will produce topics of 10 different sizes from 5 topics up to 50 topics, incrementing the number of topics by five up until it reaches 50 topics.

Topic extraction using Non-negative matrix factorization (NNMF) is probabilistic in nature and will never give the exact same results unless a fixed random seed is provided. For this reason, one may be interested in assessing the stability of some topic solutions by performing several runs with the same settings and the same number of topics. The **Replications** option allows one to specify how many replications will be performed for each setting. This option is available only when the NNMFmethod is selected and when no fixed random seed has been set. In the above example, since the user requested topic solutions with 10 different numbers of topics, the total number of topic solutions computed will be 30.

Once the desired options have been set, click the **OK** button to start the computation of all topic models. All topic model solutions will be aggregated in the report manager allowing one to compare solutions obtained in multiple runs using different settings. A summary table with various coherence statistics will be displayed. Those statistics are currently experimental measures of coherence and will later be either dropped or documented. For NNMF, a seed value is also printed, allowing one to replicate a specific topic solution by entering its associated seed value.

Topic Modeling Buttons

- Select to edit the topic name and to remove or exclude words and phrases from the topic.
- Allows you to delete the topic on the selected row.
- Click to merge a topic into another one. You first needs to select the row containing the first topic you would like to merge, and then click this button. A dialog box will appear with a list of all other topics. Select the second topic and click **OK**.

To retrieve segments associated with a topic, select it and click this button. All text segments OC) containing at least two keywords of the selected topic will be retrieved and presented in a table format. You may however change both the type of segments retrieved (paragraphs, sentences or full documents) or the minimum number of topic words needed for retrieval. When using Factor Analysis as an extraction method, this button open a dialog box containing a summary of the statistical calculations used in the topic extraction including a scree plot, percentage communalities, percentage of trace and factor pattern. Accesses the mapping function, which allows you to analyze the spatial distribution of various topics. Stores the extracted topics currently displayed into a new categorization dictionary where folders at Save the first level correspond to different topics, and where each of those folders contains the associated words. A dialog box allows you to save Click to export the selected data to Tableau, a software used for interactive data visualization. 240 83 Allows you to perform cooccurrence analysis of all the extracted topics including clustering and multidimensional scaling, and to create proximity plots as well as link charts. For more information on the various features available, see the <u>Cooccurrence Tab</u> topic. Allows you to perform full crosstabulation analysis of all the displayed topics with structured data, to apply statistical analysis, and to create various charts such as correspondence plots, heatmaps, bubble charts, and bar charts. For more information on the various features available for crosstabulation analysis, see the Crosstab Tab topic. Press this button to append a copy of the topic table in the Report Manager. A descriptive title will be U1 provided automatically. To edit this title or to enter a new one, hold down the SHIFT keyboard key while clicking this button (for more information on the Report Manager, see the Report Management Feature topic). Allows you to store the topic table to disk in various formats, including Excel, tab and comma delimited files, plain text, HTML, XML, SPSS or Stata files. Allows you to print a copy of the displayed chart

The Comparisons Tab

To the right side of the table is a panel containing two tabs. The **Comparisons** tab allows you to look at the distribution of the selected topic among values of up to two structured variables. You may display this distribution using either a vertical bar chart, a horizontal bar chart or a line chart, by clicking on the corresponding button. Four statistics may also be represented on the charts:

- Case Occurrence number of cases in this subgroup containing at least one of these words.
- Category Percent percentage of cases in this subgroup containing at least one of these words.
- Word Frequency total number of these words in this subgroup.
- Rate per 10,000 Words rate of words in this subgroup per 10,000 words.

Right-clicking anywhere in the chart areas displays a popup menu that allows you to edit the chart, save it to disk or in the Report Manager, or to copy it to the clipboard. Clicking a specific bar or a data point of a line chart also allows one to retrieve text segments associated with the selected class and containing words of the selected topic.

The Suggestions Tab

The **Suggestions** tab displays in its top panel suggested phrases, potential exceptions and spelling corrections that can be added to currently selected topic. The bottom panel contains segments related to the selected suggestion, allowing you to

judge its relevance by looking at the context in which it appears. Suggestions are highlighted in yellow, while keywords from the topic are highlighted in green.

To add any of those to the selected topic, check the boxes beside the entries you want to add, right-click the mouse and choose the appropriate destination on the adjacent menu.

Items listed as **Exceptions** result from an attempt to perform automatic disambiguation of words in the topic by identifying phrases containing topic words that are statistically unrelated to the other words. Automatic disambiguation is a very difficult task and often results in incorrect classification of relevant phrases as exceptions. For this reason, right-clicking an phrase classified as a exception gives the possibility to move it either to the list of exceptions or to the list of related phrases. Since WordStat gives precedence to phrases over single words, adding a phrase representing an exception to a topic will ensure that the word, when part of this phrase, won't be considered as an indication of the topic. When a topic containing exceptions is saved, phrases representing exceptions will be stored in the topic itself but with a weight equal to zero.

Phrases

To accurately represent the meaning of a document, it is sometimes not good enough to rely on words alone but you should look at idioms and phrases. While obtaining a comprehensive list of words is easy, finding common phrases in a specific text corpus is often much more difficult. The phrase finder feature, on the second tab (**Phrases**) of the **Extraction** tab, provides such a tool. It will scan an entire text corpus and identify the most frequent phrases and idioms and allow you to easily add them to the currently active categorization dictionary. In order to reduce redundancy in such a table, short phrases that are part of larger ones are automatically removed from the list, provided that their frequency is lower than or equal to the frequency of a longer version.

To extract phrases:

• Select the ¹¹/₁ button. An Extraction Setting dialog opens similar to the one below.

	×
Topics Phrases Named Entities Misspellings	
Min words: 2 🗘 Max words: 5 🗘	
Remove phrases ending with excluded words	
Remove phrases in categorization dictionary	
Maximum number of phrases: 1000	

The **1**¹⁰ button is present on all of the tabs on the **Extraction** tab. Each tab of this dialog contains extraction options for the currently select tab. Before scanning for phrases, you have to set various options that will be used to determine the extent of the scanning process. The first two options that need to be set are the minimum and maximum number of words a phrase can have (**Min words** and **Max words**). These two values determine both the processing time, the memory requirement as well as the number of resulting phrases. The larger the range between these minimum and maximum values, the longer it will take to collect all possible sequences of words. You can also **remove phrases ending with excluded words** or **phrases in categorization dictionary** by selecting the checkboxes beside the options.

• Set the extraction settings and click on the **Close** button.

The **Min. Frequency** or **Min. Cases** options at the top of the **Phrases** tab allow you to eliminate from the list phrases that appear only a few times by setting a minimum frequency criterion. When set to **Min. Frequency**, the criteria specifies the minimum number of times a phrase must appear regardless of whether it comes from a single document or from multiple documents. Setting it to **Min. Cases** allows you to require these occurrences to appear in a minimum specified number of cases.

• When these options have been set properly, click the button to perform the search. Once extracted, the **Phrases** tab should looks like this:

WordStat 9.0.7 - Election 2008 Coded.ppj							140		×
Topics IP Processing Frequence Topics IP Phrases Ry Named Entitles	ies 🐞 Extraction 🔣 Cooccurrences S Misspellings & Unknowns	Crosstab	III Key	word-In-Conte	xt 🔺	Classification			*
🚍 🕅 Min frequency, S 🔅 🍾 Search	8) E 9) 7 🛄 🛃 8) 🖪	1 42						6	ы 🖨
× 9[] *		FREQUENCY N	O. CASES	% CASES L	ENGTH	TF . IDF	> Comparison (5) Overlans		
~ 4	HEALTH CARE	708	138	56.33%	2	176.5	- Companson Es orenaps		
S EXCLUSION LIST	UNITED STATES	476	165	67.76%	2	80.5	📊 📂 🖄 Rate per 10,000 words 🗸 BY CANDIDAT	TE	~ =
SUBSTITUTION	AMERICAN PEOPLE	434	138	56.33%	2	108.2			
NEW CATEGORY	MIDDLE CLASS	-227	82	39,47%	2	107.9	Biden L	-	_
<pre><root></root></pre>	AL QAEDA	189	58	23.67%	2	118.3	Edwards -	-	_
CULTURE	MEN AND WOMEN	160	82	33 47%	3	76.1	Thompson -	-	
CULTURE-HIGH	SENATOR MCCAIN	158	37	15 10%	2	129.7	Giuliani- Ice		1.1
CULTURE-POPULAR	GEORGE BUSH	156	50	24 08%	2	96.5	HCCain - Changes and	-	
SPORT	WAR IN IRAO	152	75	20.61%	2	79.7	Obama -	-	-
ENVIRONMENT	ST CENTURY	127	71	30.07%	2	73.7	Romney		_
CON ENVIRONMENT	YEARSAGO	137	02	40.00%	2	54.5		+	
CROUDS	PRESIDENT BUSH	120	64	26 1296	0	74.6	0 1 2 3 4 5	0	
	NATIONAL SECURITY	107	67	20.1270		24.0			- 55
WOMEN	SENATOR ORAMA	120	20	10 0 100	2	100.4	LE DA Kate per 10,000 words V BY DELIVERY		× -
THETTUTIONS	WHITE HOUSE	120	50	05 2414	4	109,4	10		
CONSERVATIVE	WALL STREET	115	02	20.31%	2	08.0	8	-	-/
NEUTRAL		112	4/	19.18%	2	60.3	6	-	1
RADICAL	TAY CLITS	108	00	22.80%	2	09.2	4	-	-
LAW AND ORDER	IOHN MCCAIN	105	51	20.82%	2	12.2	2 0 0	-	-
LAW-CONSERVATIVE	MIDDU E EART	104	30	14.69%	2	80.0	9-4	1	
LAW-LIBERAL	DUDLE EAST	104	58	23.67%	2	05.1	000 001 001 001 001 008 008	000	008
C RURAL	BOSHADMINISTRATION	88	48	19.59%	2	70.1	The share and share shar	2	A.P
URBAN	PRESIDENT OF THE UNITED STATES	99	64	26.12%	5	57.7	a a a a a a a		-
VALUES	FEDERAL GOVERNMENT	98	66	26.94%	2	55.8			-
- CONSERVATIVE	NUCLEAR WEAPONS	97	37	15.10%	2	79.6	Yest ME E 🗿 🍘 Mot: Frequency 🔗		-
CIBERAL	FOREIGN POLICY	95	46	18.78%	2	69.0			
AUTHORITY	GLOBAL WARMING	95	44	17:96%	2	70.8	and states of the state of the		
CONTINU®	RUNNING FOR PRESIDENT	88	62	25.31%	3	52.6	UNITED STATES		
DISRUPT*	INSURANCE COMPANIES	86	33	13.47%	2	74.9	HEALTH CARE		
JURISDICTION*	TAX CREDIT	86	45	18.37%	2	63.3	TIERETTI VAILE	-	
LEGITIMATE	AMERICAN DREAM	84	46	18.78%	2	81.0	K CARDA AMERICAN PEOP	LE	
MORATORIUM	YOUNG PEOPLE	84	48	19.59%	2	59.5			
RUL"	CLIMATE CHANGE	83	45	18.37%	2	61.1			

By default, found phrases are presented in descending order of frequency. You can sort on any column of the table by clicking its header once to sort the rows in ascending order and a second time to sort the rows on the same column in descending order.

To add a phrase or idiom to the currently selected categorization dictionary or to the exclusion list, simply drag it to the proper location in the **Dictionary panel** located to the left of the screen (see <u>Working with the Dictionary Panel</u>). You may also select a phrase, right click your mouse and select the desired location.

You should keep in mind that phrases in a dictionary are always treated as a single unit and always have precedence over single words or word patterns. This means if a word is included in one content category ("A") and some phrases containing this word have already been added to another category ("B"), the word will only be categorized into "A" if it is not already part of a phrase found in "B". This feature is essential to perform disambiguation of words, since it allows you to remove some false positives associated with a word by identifying phrases associated with those false positives. For example, if a content category measuring references to money contains the word "BILL," then adding phrases like "BILL OF RIGHTS" or "BILL CLINTON" to another category will prevent those instances of "BILL" being categorized as "money." Note: When two phrases start with the same words, longer phrases have precedence over shorter ones, since they are likely more specific. However, when two phrases partially overlap, the first one encountered in the text will be categorized, preventing the second one from being recognized.

The Comparisons Tab

To the right side of the table is a panel containing two tabs. The **Comparisons** tab allows you to look at the distribution of the selected phrase among values of up to two structured variables. You may display this distribution using either a vertical bar chart, a horizontal bar chart or a line chart, by clicking on the corresponding button. Four statistics may also be represented on the charts:

- Case Occurrence number of cases in this subgroup containing at least one of these words.
- Category Percent percentage of cases in this subgroup containing at least one of these words.
- Word Frequency total number of these words in this subgroup.
- Rate per 10,000 Words rate of words in this subg roup per 10,000 words.

The bottom chart contains the distribution of phrases that can be represented as a word cloud, a vertical or horizontal bar chart, a pie chart and a donut chart. You can display the distribution with the same statistics seen on the frequency table: Frequency, No. of cases, % of cases.

Right-clicking anywhere in the chart displays a popup menu that allows you to edit the chart, save it to disk or in the Report Manager, or copy it to the clipboard. Clicking a specific bar or a data point of a line chart also allows you to retrieve text segments associated with the selected class and containing words of the selected topic.

The Overlaps Tab

While WordStat tries to reduce redundancy in the list of phrases by automatically removing short phrases that are part of longer ones, the resulting list may still contain items that are not independent of each other such as phrases that sometimes overlap. In order to allow users to take into account potential overlaps when selecting phrases, WordStat provides a display option that allows you to see when a selected phrase includes a shorter one, is part of a longer one, or sometimes overlaps other phrases. Such information is especially useful when you need to identify idioms that are more specific, often found in longer phrases or more generic ones usually composed of shorter phrases.

Selecting a phrase in the table automatically shows all other items that overlap this selected item on the **Overlaps** panel. Each phrase is accompanied by a ratio indicating the total number of times this other phrase occurs and how many times it overlaps with the selected item. For example, if you select the phrase I'M LOOKING FOR in the table showing it occurs 26 times in a document collection, you may notice that it overlaps with another phrase, LOOKING FOR SOMEONE, with a ratio of 11 out of 12. This suggests that LOOKING FOR SOMEONE occurs 12 times, but on 11 occasions, both phrases overlap (I'M LOOKING FOR SOMEONE). This ratio also indicates that on one other occasion, this second phrase occurs without overlapping the first one. It is also useful to compare the total number of overlaps with the total frequency of the target phrase. In the above example, we can conclude that the phrase I'M LOOKING FOR - occurring 26 times - is followed by SOMEONE on 11 occasions. Thus, on 15 other occasions, it is followed by something else.

Assigning overlapping phrases to a dictionary or obtaining a KWIC table:

There are two methods of assigning phrases listed in the overlap panel to the categorization dictionary or the exclusion list or of producing a keyword-in-context table. The first method performs these operations on items in this panel only, while the second method includes selected items in the main phrase list.

- To perform one of the above-mentioned operations on the overlapping items only, select the checkboxes of one or several overlapping phrases, right-click your mouse and choose the appropriate command.
- To include phrases in the main table, select one or several phrases in the main table, then the overlapping phrases listed in the overlap panel, and right click the mouse and choose the appropriate command. You may also drag and drop phrases from the main table to the appropriate location in the Dictionary panel to the left. All selected overlapping phrases will be added.

Filtering the Table

Extracting phrases from a large collection of documents can result in a very large table containing thousands of phrases.

• Clicking the \Im button brings a dialog box offering filtering options that allow you to view only phrases containing either a key word or phrases that are characteristic of a specific class. Filtering conditions are specified in a dialog box similar to this one:

Filter phrases	>
Phrases containing:	HEALTH
	Whole word only
Scoring high for:	Clinton 😔

Enabling the **Phrase containing** option and entering a string in the edit box allows you to display only phrases containing the specified string. If a comparison has been performed between classes of a categorical variable, you may also view phrases that are characteristic of a class by enabling the **Scoring high for** option and selecting the value associated with this class. In the above example, both filtering options were used, restricting the phrases displayed in the table to those containing the string HUMOUR and found to be characteristic of the 30-39 age group.

• To apply the filtering condition, click the **Apply** button. When a filtering condition is active, the **S** button is down. To remove filtering conditions and display all extracted phrases, click this button again.

Phrase cooccurrence Analysis and Crosstabulation

WordStat offers the possibility to perform cooccurrence analysis (clustering, multidimensional scaling, proximity plot) and crosstabulation on defined content categories as well as on words meeting specific frequency criteria. To apply these operations to phrases, you could add them to a user-defined dictionary, enable this dictionary and then move to the **Cooccurrences** or **Crosstab** tabs in order to access the appropriate command. The **Phrases** tab allows you to perform these operations on extracted phrases without the need to assign them to a dictionary. Special dialog boxes also allow you to add frequent words, define additional selection criteria or analysis options, and save the resulting list of items into a new content-analysis dictionary.

To perform a cooccurrence analysis on phrases:

• Click the ³/₅ button. A dialog box similar to this one will appear:

Options	×
Processing options	
Add words with:) frequency) case occurrence higher or equal to: 30	
Remove phrases with:) frequency Case occurrence less than:	
Remove words or phrases occuring in more than: 70 戻 percent of cases	
Keep a maximum of 500	
Saving options	
Save into a new dictionary: All phrases	
 All items meeting the above criteria 	
OK. X Cancel	

The options are almost identical to those found on the <u>Postprocessing</u> tab, allowing you to add to the extracted phrases, single words occurring more than a specific number of times or in a specific number of cases, or to remove items too frequent or not frequent enough. You may also restrict the analysis to a specific number of items.

The **Saving into a new dictionary** options may be used to store phrases and words into a new dictionary file. Enabling this option and selecting **All Phrases** will give you the opportunity to store all phrases extracted using the **Phrases** tab, whether or not they meet the specified frequency criteria. Selecting the other option will store only phrases meeting these criteria as well as all words that were also included in the analysis.

The **Saving into a new dictionary** options may be used to store phrases and words into a new dictionary file. Enabling this option and selecting **All Phrases** will allow you the opportunity to store all phrases extracted using the **Phrases** tab, whether or not they meet the specified frequency criteria. Selecting the other option will store only phrases meeting these criteria as well as all words that were included in the analysis.

For further information on the various tools available for annualizing cooccurrences, see <u>Hierarchical Clustering and</u> <u>Multidimensional Scaling</u>.

To perform a crosstabulation analysis:

• Click the 🛄 button. A dialog box similar to this one will appear:

Crosstabulation of Phrases	×
Processing options	
Add words with: I frequency Case occurrence higher or equal to: 30	
Remove phrases with: () frequency () case occurrence less than:	
Remove words or phrases occuring in more than: 70 🚔 percent of cases	
Keep a maximum of 500 🔿 items based on TF*IDF	
Saving options	
Save into a new dictionary: All phrases	
○ All items meeting the above criteria	
OK Cancel	

The **Processing** and **Saving** options are identical to the ones available when performing cooccurrence analysis on phrases (see above).

For more information on crosstabulation, see Crosstab Tab.

Other Table Operations

To copy the entire table to the clipboard:

- Right-click anywhere in the frequency table.
- Select the COPY TO CLIPBOARD | TABLE command from the pop-up menu.

To copy selected rows to the clipboard:

- Select the rows you would like to copy (multiple but separate rows can be selected by clicking while holding down the CTRL key).
- Right-click anywhere in the frequency table.
- Select the COPY TO CLIPBOARD | SELECTED ROWS command from the pop-up menu.

To search for a specific item:

- Right-click anywhere in the table.
- Select the **FIND** command from the pop-up menu. A search dialog box will appear.
- Type the desired search string in the **Find What** edit box. To restrict the search to items starting with the search string, enable the **Match Beginning of Item** option and to restrict it to whole words matching the search string, enable the **Match Whole Word Only** option.

Click the **Find** button to search the first item matching the typed string. Clicking this button again finds the next occurrence of the search string, starting at the currently selected item.

Named Entities

The **Named Entities** extraction tool uses pattern-based algorithms to identify people, locations, organizations and acronyms. This feature relies on the existence of mixed-case words and all upper case words in text, and will ignore sentences or paragraphs in full upper case. This approach to named-entity recognition has the benefit of working in many languages but may not be appropriate for some languages, such as German, where upper case letters are used in common nouns. It may also miss some named entities if some of their significant words are not capitalized. For example, "Ministry of education" will not be recognized, while "Ministry of Education" will. The phrases extraction tool may, however, be able to retrieve these named entities if they occur more than once.

To extract named entities:

Clicking the Settings button displays a dialog box similar to the one below:

Extraction	Settings			×
Topics	Phrases	Named Entities	Misspellings	
Rem	ove items in o	categorization diction	ary	
Rem	ove known co	ommon words		
Min tota	frequency:	3		
				Close

Remove items in categorization dictionary: This option removes, from the final list, named entities that are already in the active categorization dictionary. By default, this option is enabled.

Remove known common words: Enabling this option removes extracted single word entities that are also common words. This option is especially useful to remove capitalized words in titles or words in social-media data that have been put in full upper case, typically to stress their importance. It may, however, miss the identification of named entities that are also common nouns such as "Apple", "Windows" or some acronyms.

Minimum total frequency: This option allows you to eliminate, from the list, named entities that appear only a few times.

- When these options have been set properly, select the
- Click the Search button to perform the search.

WordStat 9.0.7 - Election 2008 Coded.ppj					- o ×
E Data Text Processing	Requencies 🐉 Extraction 🖄 Cooccurrences 🛅	Crosstab 🏭 Keywi	ord-In-Context	Classification	0-
Topics Phrases Named Ent	ibes Thisspellings & Unknowns				
🗏 🔢 Settings 🍾 Search 🕅					6 6 6
5.91	« ENTITY	TOTAL	UNIQUE		*
	United States	289	103		
EXCLUSION LIST	U.S.	167	167		
SUBSTITUTION	Obarra	162	59		
T NEW CATEGORY	George Bush	129	124		
2 <roo 1=""></roo>	Senator McCain	118	118		
CULTURE	Wall Street	107	103		
CULTURE-HIGH	President Bush	106	106		
COLTURE-POPULAR	White House	105	98		
ENVIRONMENT	Middle East	104	104		
CON ENVIRONMENT	John McCan	77	69		
PRO ENVIRONMENT	Senator Obama	75	75		
GROUPS	Cold War	70	70		
ETHNIC	New York	70	22		
WOMEN	Al Oaeda	68	66		
INSTITUTIONS	New Mexico	60	57		
CONSERVATIVE	President of the United States	59	45		
D NEUTRAL	Latin America	51	51		
C RADICAL	Ronald Reagan	51	40		
LAW_AND_ORDER	Social Security	50	40		
LAW-CONSERVATIVE	United States of America	50	21		
LAW-LIBERAL	Bush Administration	40	31		
RURAL	World War II	45	45		
URBAN	Democratic Party	45	43		
VALUES VALUES	North Korea	43	41		
- CONSERVATIVE	Nove Hamachina	72	42		
LIBERAL	Canatar Cintan	41	30		
AUTHORITY	American Dream	40	40		
CONTINU"	South Carolina	37	37		
INSPECT*	Didge	37	35		
JURISDICTION*	Carlet Union	37	8		
MANAG"	Soviet Union	36	36		
MORATORIUM	New Orleans	36	36		
ctube:	w new York City	33	31		~

Retrieved items are listed in descending order of frequency. The **Total** column indicates the total frequency of this named entity, while the **Unique** column reports how many do not overlap with others.

Adding Named Entities to a Dictionary or Obtaining a KWIC Table

To add a named entity to the currently selected categorization dictionary or to the exclusion dictionary simply drag it to the proper location in the dictionary panel left of the screen (see <u>Working With the Dictionary Panel</u>). To put them into a new category, drop them into the **New Category** item near the top of the dictionary panel. You may also right click your mouse and select the desired location.

To produce a keyword-in-context table of a specific named entity, first select it and then right click your mouse and select **Keyword-in-Context** list.

Misspellings & Unknowns

The **Misspellings & Unknowns** feature of WordStat offers a tool that identifies common misspellings by comparing the list of word forms encountered in the entire text collection against a list of common words. This feature may also retrieve single words that represent or are part of technical terms, company and product names, as well as abbreviations. Thus, it partly

overlaps the <u>Named Entity</u> extraction tool that will often provide a more complete list of the entities by adding multi-word terms representing the entities as well as common words with irregular capitalization.

When a categorization dictionary is active, this feature will also attempt to match unknown words to existing items in this dictionary and will present the extracted words in two separate lists: the **Categories** list will contain all unknown words that are potentially related to existing items in your dictionary, presented in a tree view organized by categories while all other unknown words will be listed on the **Others** tab.

Once extracted, identified misspellings may be added to the categorization dictionary, to the exclusion list, or a substitution process. You may also replace these words in the original documents with the proper spelling.

To extract misspelled and unknown words:

Clicking the Settings button displays a dialog box similar to this one:

X
Close

Several options are available to control how unknown words are extracted and how they are matched to potential replacements.

Dictionary: By default, the extraction is performed in reference to common English words (British and American English). To identify unknown words in documents written in another language or to exclude technical terms from a specific domain, click the language name(s) listed beside the **Dictionary** option or modify the settings on the <u>Languages</u> section located on the main **Text Processing** tab of WordStat.

Remove words already in the categorization dictionary: Some misspellings or uncommon words may have already been added to an exclusion list or a categorization dictionary. This option automatically removes these items from the lists of retrieved words.

Find match for words in the categorization dictionary: WordStat will attempt to match unknown words with existing ones in the language spelling dictionary and identify the most likely replacements. Enabling this option will also attempt to match uncommon words in the categorization dictionary, such as technical terms, proper names, or even other misspelled words.

Confidence of suggestions: This option allows you to adjust the confidence level used by WordStat when identifying potential replacements. Setting it to **Highest** will result in less suggestions by applying a stricter matching criterion, while setting it to **Moderate** will provide more suggested replacements at the risk of increasing the number of irrelevant suggestions.

Minimum frequency for uncategorized words: This option allows you to eliminate, from the **Others** list, items that appear only a few times, thus allowing you to focus on the most frequent items. Please note that this setting has no effect on the **Categories** list, since matches will be reported even for words appearing only once.

Identify misspelling of words in categorized phrases: For a phrase in a categorization dictionary to be recognized, all its words must be properly spelled. It is thus important to pay close attention to those misspellings that

may result in a failure to recognize such a phrase. Setting this option will add an additional column in the **Others** list to indicate whether it matches a word that is part of one of the phrases in the categorization dictionary.

Batch replacements settings: WordStat remembers text replacements made in prior sessions and can reapply them in batches. Clicking this button brings a dialog box that allows you to edit or export this list of text replacements as well as to import a list created by another user or on another computer. (see <u>Managing the Batch Replacement</u> <u>List</u> for more information).

• Once your settings have been established, click the Search button. A dialog similar to this one will appear.

WordStat 9.0.7 - Aerospace.ppj				-		×
E Data Text Processing Frequencies	Extraction 🗞 Cooccurrences 🛅 Cro	sstab 🏦 Keyword-In-Context < Classifica	ation			
👬 Topics 🛄 Phrases 🖳 Named Entities 🍄 Miss	pellings & Unknowns					
11 Settings 📏 Search					10 k	
*	Categories Others		-	Actions to be performed:	-	
A EXCLUSION LIST	CATEGORIES & WORDS	FREQ	-	SUBSTITUTE: groth -> growth		
12 SUBSTITUTION	VACATION TIME		-	SUBSTITUTE: growth -> growth		
+ NEW CATEGORY	• OVERTIME			SUBSTITUTE: grwoth -> growth		
	EVERYTIMNE	1		REPLACE: limited -> limited		
	C ? OVERTME	1				
	T ? OVETIME	3				
PERFORMANCE REVIEW	- VACATION					
CAREER GROW TH	VACATION	1				
LARGE COMPANY	T ? VACTION	1				
GOOD BENEFITS	· WEEKS	e				
FLEXIBLE WORK SCHEDULE	T ? WEEKSS					
COMPETITIVE SALARY	PERFORMANCE REVIEW					
OPPORTUNITIES	- DEDEODMANCE					
GOVERNMENT CONTRACTS		2				
GREAT PLACE TO LEARN		2				
VEAR AFTER		-				
FRIENDLY AND HELPFUL & WORK ENVIRONMENT	I PERFORMACE	3	Add to Substitution 🔫			
SOFTWARE ENGINEER	I PERFORMANACE	1		1		
SLOW MOVING & CHANGE	I PERFORMENCE	1	Categorize 👄			
HR DEPARTMENT & PROGRAM MANAGER	PERFORMINCE	1	A DESCRIPTION OF A DESC			
LIFE BALANCE	I PERFROMANCE	1	Replace in text			
INTERESTING PROJECTS	PRFORMANCE • PROMOTIONS	1	Add to Spell Dictionary 🔫			
GREAT PRODUCTS	PRMOTOTIONS	1	Demove Artist			
Carmour service a strater. V	2 PROMOTIOINS	1	nemers nests)			
******	REVIEW		III Keyword-In-Context			
FLEXIBLE	REVEIW	1				
PLEXIBLE_PROOKS	 SYSTEM 					
FLEXIBLE_WORK	SSTEM	1				
FLEXIBLE_WORK_HOURS	T ? SYSTME	1				
FLEXIBLE_WORK_SCHEDULE	CAREER GROWTH	2				
HOURS	· ADVANCEMENT					
SCHEDULE	ADAVANCEMENT	1		1		
WORK_HOURS	ADVANCEMETN	i		and mint	in actions	
HODITAL HOLDS			*	Eeu	and accord	

WordStat will retrieve all unknown words and display them in two lists located in the middle panel of the dialog box: the **Categories** list and the **Others** list. You can move from one list to the other by selecting the corresponding tab.

The Categories List

If a content analysis dictionary is available and active, WordStat will attempt to match unknown words with existing items in this dictionary and present them in a tree list view, with the unknown forms stored below the original dictionary form with which it was matched. The suggestions are identified with a "?" symbol and are followed by a numerical value representing their frequency in the current text collection. Existing items and suggested forms are also grouped under their containing content category. Check boxes on the left are used to select multiple items in order to perform some operations on them. To select a single item, click its check box. Click again to remove this check mark. Selecting a content category, subcategory or an existing dictionary item, automatically selects all items below it, If there is no active categorization dictionary, the **Categories** tab will be disabled.

The Others List

The second list contains all words that could not be matched to an existing item of the content analysis dictionary, either because no dictionary was currently active or because they were not similar enough to be matched confidently. On the right of the word, WordStat shows the most likely word replacement that has been found, This column will remain empty if no replacement was similar enough. A third column contains a numerical value representing the frequency of the unknown

form in the current text collection. Finally, a fourth column will contain the word "Yes" if the word may be matched to a word that is part of a phrase in the content dictionary. Such an indication may be useful to identify situations where a phrase may go undetected because one of its words has been misspelled.

Available Operations

Up to four types of transformation or processing are allowed on the words: 1) You can replace all instances of a selected word in the original document by another word or phrase; 2) You can assign the words to existing content categories; 3) you can instruct WordStat to automatically replace this word with another one by adding it to a live substitution process; 4) you can add this word to a custom list of valid words causing the program to ignore these words the next time there is a search for vocabulary words. You may also obtain a keyword-in-context list associated with a specific word in order to decide how that word should be treated.

Except for the Keyword-in-context list, none of the other four operations are performed immediately. Instead they are added to an action list allowing you to review, modify or cancel previously defined actions prior to the application of all the specified changes.

To add a word to the substitution process:

- Select the word you want to add.
- Click the Add to Substitution → button.
- The selected word followed by the word it will be substituted for will be added to the **Actions to be performed** panel on the right.
- Once all desired actions are present in the Actions to be performed list, select the button. The substitution will now be listed on the Substitution tab on the Text Processing tab.

To replace words in the original documents:

- Select the word to be replaced.
- Click the

button. A dialog box similar to this one will appear.

Suggestions:	Replace
humane	huvane
human huddle humanize hulas	With:
Huron	

Replace in text 🔿

- Type the new replacement word or phrase or choose from the **Suggestions** list box on the left side of the dialog box.
- Click the **OK** button to confirm this replacement and add this operation to the list of **Actions to be performed**.

192 WordStat User's Guide

• Once all desired actions are present in the **Actions to be performed list**, select the button.

To add a word to a category of your content analysis dictionary:

- Select the word you want to categorize.
- Click the Categorize button. A dialog box will appear with a list of all categories and subcategories in the current content analysis dictionary.
- Select the category under which this item should be stored and then click OK.
- Once all desired actions are present in the Actions to be performed list, select the button.

You may add words from this list to the current categorization dictionary or to the exclusion list, or assign it to the substitution process in order to have it replaced automatically by another word. To perform any one of these assignments, simply select and drag the item into the proper location in the dictionary panel to the left of the table (see Using the Dictionary Panel).

To add a word to the custom list of words to ignore:

- Select the word you would like to add to the custom dictionary.
- Click the Add to Spell Dictionary
 button.
- The selected word to be added will appear the Actions to be performed panel on the right.
- Once all desired actions are present in the Actions to be performed list, select the button.

To remove operations previously defined:

- Select the operations that you want to remove.
- Click the Remove Action button. All words associated with the removed actions are moved back to the list of unknown words and positioned at the bottom of the list.

Managing the Batch Replacement List

WordStat remembers text replacements made in prior sessions and can reapply them in batches. You can edit entries in this list of replacements, remove items, merge lists by importing lists created by other users, or you can export this list to a file.

To manage the Batch Replace list:

- Click the **Settings** button. The Extraction Settings dialog will appear.
- Select the Batch Replacements Settings button. A dialog box similar to this one will appear:





Perform actions

	MISSPELLED FORM	REPLACED WITH	
+ Add	americian	american	
- Delete	cacasian	caucasian	
Delett	cacasion	caucasian	
Edit	caucaisian	caucasian	
Import	caucasin	caisattion	
	caucassion	caucattion	
Export	caucation	caucasian	
To Substitution	causcasian	caucasian	
	cocasian	caucasian	
	culttural	cultural	
	definately	definitely	
	detai	detail	
	docuament	document	
	dont	don't	
	everytime	every time	
Save	informaiton	information	
Class	knoweldge	knowledge	
-u close	preservationâ	preservation	

To manually add a replacement:

• Click the Add button. A dialog box like this one will appear:

Substitution			×
Substitute:	With:		
	🗸 ок	×	Cancel

- Type the word to be replaced in the **Substitute** edit box.
- Type the word that will replace it in the With edit box.
- Click the OK button to add the replacement to the existing list.

To edit a replacement:

- Select the row containing the item you want to edit.
- Click the Edit button. A dialog box similar to the one shown above will appear.
- Edit any one of the entries in the Substitute or the With edit box.
- Click the **OK** button to confirm the change.

To delete replacements:

• Select the rows containing the items you want to delete. You can hold down the Ctrl key to select disjointed rows, or the Shift key to select a range of successive rows.

• Click the **Delete** button.

To import a replacement list:

- Click the Import button. An Open File dialog box will appear.
- Select the file containing the replacements you want to import, and click **Open**. All replacements in the imported file will be appended to the current replacement list. Duplicate items won't be imported, while conflicting items will be shown, allowing you to either keep the existing replacement or overwrite it with the one in the imported file.

To export a replacement list:

- Click the Export button. A Save File dialog box will appear.
- Type a valid file name and select the location where the file should be stored.
- Click the **Save** button.

To move automatic replacements to the substitution process:

- Select the words you would like to add to the substitution process.
- Select the **To Substitution** button. The substitution will now be listed on the **Substitution** tab on the **Text Processing** tab.

The Cooccurences Tab

The **Cooccurrences** tab allows you to perform hierarchical cluster analysis and multidimensional scaling on all keywords. Results are displayed in the form of dendrograms, concept maps and proximity plots, based on item cooccurrence. It also allows you to compute the similarity of cases based on keyword use. For further information see <u>Hierarchical Clustering and</u> <u>Multidimensional Scaling</u>.

The **Cooccurrences** tab consist of six interior tabs.

The <u>Options tab</u> allows you to specify whether the clustering should be performed on keywords or on cases and to set various analysis and display options.

The <u>Dendrogram tab</u> uses average-linkage hierarchical clustering method to create clusters from a similarity matrix.

The <u>Mapping tab</u> provides a graphic representation of the proximity values computed on all included keywords using multidimensional scaling.

The Link Analysis tab allows you to visualize the connections between keywords or dictionary items using a network graph.

The <u>Proximity Plot tab</u> is the most accurate way to graphically represent the distance between objects by displaying the measured distance from one or several target objects to all other objects.

The <u>Statistics tab</u> displays the cooccurrence and similarity matrices used for building the dendrogram, MDS plots, as well as the proximity plot.

The <u>Cooccurrence matrix tab</u> provides an interactive matrix that allows one to focus on specific co-occurrences, retrieve associated text segments, and see how they are associated with other variables.

Hierarchical Clustering and Multidimensional Scaling Options

WordStat allows you to further develop categorization by providing various graphic tools to assist in the identification of related words or categories. The tools are obtained by the application of hierarchical cluster analysis and multidimensional scaling on all included words or categories and are displayed in the form of dendrograms and concept maps. Cases or documents may also be clustered based on their content similarity using the same statistical and graphic tools.

The first tab, the **Options** tab, is used to specify whether the clustering should be performed on keywords or on cases and to set various analysis and display options.

WordStat 9.0.7 - Election 2008 Coded.ppj		- 0	×
🚍 🏢 Data 🚽 Text Processing 🍵 Frequencies 🚯 Extraction 🗞 Cooccurrences 💼	Crosstab 🏢 Keyword-In-Context < Classification		0
🔃 Options 😫 Dendrogram 😢 Mapping 💥 Link Analysis 🕎 Proximity Plot Σ Statistic	s p ^{re} Cooccurence Matrix		
Clustering			
Keywords/Categories Cases/Documents			
Occurrence: Same paragraph			
Index: Jaccard's coefficient (occurrence)			
Type: Word co-occurrence - First order			
Remove single word dusters			
Multidimensional scaling options			
Real time animation			
Tolerance: 0.000001 🖨 Maximum Iterations: 500 😴			
Initial configuration:			
Classical scaling			
O Randomized location			
Seed:			
History:			
245 cases	Shown: 15 Types: 15,813 Tokens: 590,701 Time: 1.8s (4 cases excluded)		

Clustering Cases/Documents

When the clustering is set to be performed on cases or documents, the distance matrix used for clustering and multidimensional scaling consists of cosine coefficients computed on the relative frequency of the various keywords. The more similar two documents are in terms of the distribution of keywords, the higher the coefficient. The case label that is used to identify the various cases can be set by choosing the <u>Edit Case Descriptors</u> command from the WordStat main menu.

Clustering Keywords

When clustering keywords or content categories, several options are available to define cooccurrence and choose which similarity index will be computed from the observed cooccurrences.

Cooccurrence: This option allows you to specify how a cooccurrence will be defined. By default, a cooccurrence is said to happen every time two words or two categories appear in the same case (**by case** option). You may also restrict the definition of cooccurrence to entries that appear in the **same paragraph** or the **same sentence**, or to words or categories that are located in the same **user defined section** (delimited by a *I* character). Finally, you may restrict even further the definition of cooccurrences by limiting the cooccurrence to a small **window of words** of specified length. Such a small window is especially useful when doing an analysis directly on words (rather than categories) since it allows identifying idioms or phrases that may need to be added to the categorization dictionary. cooccurrence on larger text segments such as cases or paragraphs may be more appropriate to identify the cooccurrence of themes in individual subjects.

Index: The Index option allows the selection of the similarity measure used in clustering and in multidimensional scaling. Four measures are available. The first three measures are based on the mere occurrences of specific words or categories in a case and do not take into account their frequency. In all these indexes, joint absences are excluded from consideration.

Jaccard's Coefficient: This coefficient is computed from a fourfold table as a/(a+b+c) where a represents cases where both items occur, and *b* and *c* represent cases where one item is found but not the other. In this coefficient equal weight is given to matches and non matches.

Sorensen's Coefficient: This coefficient (also known as the Dice coefficient) is similar to Jaccard's but matches are weighted double. Its computing formula is 2a/(2a+b+c) where *a* represents cases where both items occur, and *b* and *c* represent cases where one item is present but the other one is absent.

Phi Coefficient: This is a measure of association for two binary variables. It is similar to the Pearson correlation coefficient in its interpretation.

Cosine Theta: This measures the cosine of the angle between two vectors of values. It ranges from -1 to +1. This coefficient takes into account not only the presence of a word or category in a case, but also how often it appears in this case.

Inclusion Index; This index measures the conditional probability that a document that contains an item X will also contain an item Y. It will take the maximum value of 1 when one of these items always appears when the second one appears, even if the reverse is not necessarily true. The Inclusion Index is optimal for analyzing fields that are organized hierarchically.

Association Strength: This measures the cooccurrence of items taking into account the possibility that two items will sometimes cooccur by chance.

Clustering type: Two broad types of keyword clustering are available. The first method is based on **keyword cooccurrences (First Order Clustering)** and will group together words appearing near each other or in the same document (depending on the selected cooccurrence window). The second clustering method is based on **cooccurrence profiles (Second Order Clustering)** and will consider that two keywords are close to each other, not necessarily because they cooccur but because they both occur in similar environments. One of the benefits of this clustering method is its ability to group words that are synonyms or alternate forms of the same word. For example, while TUMOR and TUMOUR will seldom or never occur together in the same document, second order clustering may find them to be pretty close because they both cooccur with words like BRAIN or CANCER. Second order clustering will also group words that are related semantically such as MILK, JUICE, and WINE because of their propensity to be associated with similar verbs like DRINK or POUR or nouns like GLASS (for more information, see Grefenstette, 1994).

Remove single word clusters: One way to extract potentially interesting knowledge from dendrograms is to focus on the aggregation of items at an early stage of the clustering process. However, when clustering hundreds or thousands of items, the identification of those items requires the user to scroll through a very long dendrogram which includes many clusters of isolated items. Enabling this option simplifies the use of the dendrogram by hiding all single item clusters and allowing you to concentrate only on the strongest associations. Setting this option also removes isolated items from multi-dimensional scaling plots, greatly enhancing their value when analyzing a large number of items. Please note, however, that when this option is enabled, changing the number of clusters while viewing a 2-D or 3-D MDS plot will cause the program to recompute the distance and location of remaining items.

The options below apply whether clustering cases or keywords. They will affect either the computation or the display of multidimensional scaling charts.

Real time animation: When this option is enabled, the multidimensional plots are updated after every iteration allowing the user to monitor the progress made during the analysis at the cost of higher computing time.

Tolerance: This option specifies the tolerance factor that is used to determine when the algorithm has converged to a solution. Reducing the tolerance value may produce a slightly more accurate result but will increase the number of iterations and the running time.

Maximum iterations: This option allows you to specify the maximum number of iterations that are to be performed during the fitting procedure. If the solution does not converge to the limit specified by the TOLERANCE option before the maximum number of iterations is reached, the process is stopped and the results are displayed.

Initial configuration: This option allows you to specify whether the multidimensional scaling will be applied on a random configuration of points or on the result of a classical scaling.

Selecting the **Classical Scaling** option instructs WordStat to perform a classical scaling first on the similarity matrix, and then use the derived configuration as initial values for the ordinal multidimensional scaling analysis.

Selecting the **Randomized Location** option instructs WordStat to perform the multidimensional scaling analysis on a random configuration of points. By default, WordStat initializes the random routine before each analysis with a new

random value. The seed value used for the creation of this initial configuration is stored along with the final stress value in the history list box, located at the bottom of the dialog box. The **Seed** option may be used to specify a starting number that will be used to initialize the randomization process and produce a fixed random sequence. To recall a specific seed value used previously, double-click the proper line in the history list box.

For more information on the available options, click the links below:

Dendrogram <u>2-D and 3-D Concept Maps</u> <u>Proximity Plot</u> Table and Statistics

Dendrogram

WordStat uses an average-linkage hierarchical clustering method to create clusters from a similarity matrix. The result is presented in the form of a dendrogram (see below), also known as a tree graph. In such a graph, the vertical axis is made up of the items and the horizontal axis represents the clusters formed at each step of the clustering procedure. Words or categories that tend to appear together are combined at an early stage while those that are independent from one another or those that don't appear together tend to be combined at the end of the agglomeration process.



AGGLOMERATION ORDER: JACCARD'S COEFFICIENT

- **No clusters** This option allows the setting of the number of clusters that the clustering solution should have. Different colors are used both in the dendrogram and in the 2-D and 3-D maps to indicate membership of specific items to different clusters. However, if the option to remove single item clusters is enabled, an increase in the number of clusters may, in fact, result in a decrease in the number of clusters displayed and in the overall height of the dendrogram since all single item clusters will be hidden.
 - **Display** This option lets you choose whether the vertical lines of the dendrogram represent the agglomeration schedule or the similarity indices.
 - When clustering keywords or content categories, clicking this button displays bars beside each dendrogram item to represent their relative frequencies.
 - Use this button to increase the dendrogram font size and focus on a smaller portion of the tree.
 - Q Use this button to reduce the dendrogram font size and view a larger portion of the tree.

- This button allows you to perform full crosstabulation analysis with structured data, apply statistical analysis, and create various charts such as correspondence plots, heatmaps, bubble charts and bar charts. A dialog box allows you to restrict the analysis to specific clusters containing a minimum number of items, and cluster names are automatically generated using characteristic words and phrases of each cluster. For more information on the various features available for crosstabulation analysis, see the <u>Crosstab</u> topic.
- This button stores the cluster solution currently displayed into a new categorization dictionary where folders at the first level correspond to different clusters, and where each of those folders contains the associated words or expressions. A dialog box allows you to save only clusters containing a minimum number of items. Cluster names are automatically generated using characteristic words and phrases from each cluster. You may then edit these cluster names and their content from the <u>Dictionary tab</u>.
- Press this button to append a copy of the graphic in the Report Manager. A descriptive title will be provided automatically. To edit this title or to enter a new one, hold down the **Shift** key while clicking this button (for more information on the Report Manager, see the <u>Report Management Feature</u> topic).
- This button allows you to store the displayed dendrogram into a graphic file. WordStat supports three different file formats: .BMP (Windows bitmap files), .PNG (Portable Network Graphic compress files) and .JPG (JPeg compressed files).
- To retrieve text segments or documents associated with a specific cluster, click anywhere on a cluster to select it (the selected cluster is displayed using thicker black lines), and then click this button to retrieve the associated documents. When performing first order clustering on keywords, this operation retrieves all text segments containing at least two keywords of the selected cluster. When performing second order clustering of keywords, all text segments containing a single one of those keywords will be retrieved.
- The slide ruler provides another way of quickly changing the number of clusters included in the clustering solution. Moving the slider to the left increases the minimum distance required to form a cluster and thus produces a dendrogram with more clusters. Moving the slider to the right aggregates smaller clusters into bigger ones. However, if the option to remove single-item clusters is enabled, an increase in the number of clusters may, in fact, result in a decrease in the number of clusters displayed and in the overall height of the dendrogram.

Using the Right Panel

On the right side of this table, a panel allows you to look at the distribution of the selected topic among values of up to two structured variables as well as a network graph representing the relationship between all items in the selected cluster.

For the first two panels, you may choose to display the distribution of this cluster (using either a vertical bar chart, a horizontal bar chart, or a line chart) by clicking the corresponding button. Four statistics may also be represented on the charts:

- Case Occurrence number of cases in this subgroup containing at least one of these words.
- Category Percent percentage of cases in this subgroup containing at least one of these words.
- Word Frequency total number of these words in this subgroup.
- Rate per 10,000 Words rate of words in this subgroup per 10,000 words.

Right-clicking anywhere in the chart area displays a pop-up menu that allows you to edit the chart, save it to disk or in the **Report Manager**, or copy it to the clipboard. Clicking a specific bar or a data point of a line chart allows you to retrieve text segments associated with the selected class and containing words of the selected cluster.

In the lower-right panel, the relationships between words within the selected cluster are represented using a network-type graph. The layout of the nodes can be based on a multidimensional scaling analysis, a force based graph, or may be spread out in a circle. A track bar allows you to hide the weaker connections, while other buttons may be used to zoom in and out.

Clicking the X button allows you to jump to the full-sized Link Analysis tab.

2D and 3D Concept Maps

The concept maps are graphic representations of the proximity values computed on all included keywords using multidimensional scaling (MDS). In these maps, a point represents an item (keyword or content category) and the distances between pairs of items indicate how likely those items tend to appear together. In other words, items that appear close together on the plot usually tend to occur together, while words or categories that are independent from one other or that don't appear together are located on the chart far from each other. Colors are used to represent membership of specific items to different partitions created using hierarchical clustering. An option also allows you to vary the size of each data point in order to take into account the observed frequency of each item. The resulting maps are useful to detect meaningful underlying dimensions that may explain observed similarities between items.

Please note that since multidimensional scaling attempts to represent the various points into a two- or three-dimensional space, some distortion may result, especially when this analysis is performed on a large number of items. As a consequence, some items that tend to appear together or are parts of the same cluster may still be plotted far from each other. Also, performing a multidimensional scaling on a large number of items usually produces a cluttered map that is hard to interpret. For these reasons, interpretation of concept maps may be feasible only when applied to a relatively limited number of items.



2-D and 3-D Map controls:

No. clusters: This option allows the setting of the number of clusters that the clustering solution should have. Different colors are used both in the dendrogram and in the 2-D and 3-D maps to indicate membership of specific items to different clusters. The slide ruler located on the top toolbar of the dialog box may also be used to quickly change the number of clusters. Please note that when the **Remove single word clusters** option is enabled, changing the number of clusters either way often causes the program to recompute MDS maps to take into account the different number of items displayed.

The actual orientation of axes in the final solution is arbitrary. The map may be rotated in any way you want provided the distances between items remain the same. The rotating knob can be used to adjust the final orientation of axes in the plane or space in order to obtain an orientation than can be most easily interpreted.

- Clicking this button enables you to zoom in on a plot. To zoom an area of the plot, hold the left mouse button and drag the mouse down/right. A rectangle indicates the selected area. Release the left mouse button to zoom.
- Clicking this button restores the original viewing area of the plot.
- Clicking this button performs another multidimensional scaling on a new random configuration of points. This button is visible only when the initial configuration is set to Random Location.

This button is used to create a copy of the chart to the clipboard. When this button is clicked, a pop-up menu appears, allowing you to select whether the chart should be copied as a bitmap or as a metafile.

This button allows editing of various features of multidimensional scaling plots such as the appearance of value labels and data points, the chart and axis titles, the location of the legend, etc. (see <u>Multidimensional Scaling Plot Options</u>)

Pressing down this button creates a constrained multidimensional scaling. This mapping algorithm allows you to preserve the clustering structure in multidimensional scaling plots, making the interpretation of 2-D and 3-D MDS maps a lot easier and more consistent with the clustering solutions. Enabling this option allows you to use the MDS module to create maps of concepts similar to those suggested by Trochim, in its Concept Mapping procedure.

Pressing down this button displays lines to represent relationships between data points of the multidimensional scaling plot. When the button is down, a cursor will appear in a tool panel below the plot, allowing you to select the minimum association strengths to be displayed.

Clicking this button creates a bubble plot where the areas of data points are proportional to the relative frequency of those items. This type of display is especially useful when you need to take into account a third variable, in this specific case the frequency of items, when interpreting the distance between data points.

Press this button to append a copy of the graphic in the Report Manager. A descriptive title will be provided automatically. To edit this title or to enter a new one, hold down the **Shift** key while clicking this button (for more information on the Report Manager, see the <u>Report Management Feature</u> topic).

- This button allows storing the displayed multidimensional scaling plot into a graphic file. WordStat supports four different file formats: .BMP (Windows bitmap files), .PNG (Portable Network Graphic compress files) and .JPG (JPeg compressed files) as well as .WSX a proprietary file format (WordStat Chart file). Charts stored in the latter format may be opened, further edited and customized using the Chart Editor external utility program.
- Clicking this button allows you to print a copy of the displayed chart.

3-D Map buttons:

This button can be used to show or hide left, bottom and back walls. Clicking this button exchanges the data of the X, Y and Z axis. Stol) Clicking this button draws anchor lines from the floor to the data point to better locate data points in all 3 dimensions. Clicking this button allows changing the viewing angle of the chart. To rotate the chart, make sure ÷ this button is selected, click any area of the chart, hold the mouse button and drag the mouse to apply the desired rotation. Ø Locating a data point on the depth dimension of a 3-D plot can be very difficult, especially when the plot remains static. You often have to rotate this plot constantly on the various axes to get an accurate idea of where the data point is located on this third axis. Clicking this button forces WordStat to rotate the plot automatically. To disable the automatic rotation, click a second time.

Multidimensional Scaling Plot Options Dialog (COPY)

The various options in this dialog box are used to customize the appearance of multidimensional scaling plots. These options represent only a small portion of all settings available.

To further customize the chart, edit data points, value labels, etc., click the ¹⁴/₄ button located at the right side of the dialog box.

Data Points

View data points: By default, multidimensional scaling plots display both the data points and their associated labels. This option toggles on and off the display of the data points.

Style: This option defines the shape used to display the data points. Nine different pointer shapes may be used to represent data points (e.g. square, triangles, dots, cross, etc.).

Shadow: This option toggles on and off the shade on the sides of data points. This option only affects square, triangle and down triangle data points when viewed in three-dimension.

Size: This option specifies the width and height of data points in pixels.

Transparency: This option sets the transparency level of the data points in bubble plots.

Data Labels

Transparent labels: This option allows you to specify whether data labels are to be displayed on an opaque background or whether the background should be invisible.

Background: This option allows you to select the background color of data labels when not transparent

Border: This option enables or disables the rectangle surrounding the label background.

Vertical gap: The vertical gap property determines the vertical distance in pixels between the top of the data point and the bottom of its label.

Font: Click this button to adjust the font properties used to display data labels (i.e., style, size, and color).

Walls and Frame Tab

Wall visible: Use this option to toggle the display of left, back and bottom walls to simulate a 3-D effect.

Transparent: The Transparent property controls whether the wall background will be opaque or transparent.

Dark 3-D: This option shades the sides of the walls.

Title

Show title: By default, multidimensional scaling plots have no title. To display a title, enable the **Show Title** option and enter the desired title in the edit box to the right of this check box. Enter several lines of text by pressing the **<Enter>** key at the end of a line before entering the next line.

The **Font** button at the bottom of this tab is used to change the font size or style of this title.

Proximity Plot

Cluster analysis and multidimensional scaling are both data reduction techniques and may not accurately represent the true proximity of keywords or cases to each other. In a dendrogram, while keywords that cooccur or cases that are similar tend to appear near each other, you cannot really look at the sequence of keywords as a linear representation of those distances. You have to remember that a dendrogram only specifies the temporal order of the branching sequence. Consequently, any cluster can be rotated around each internal branch on the tree without, in any way, affecting the meaning of the dendrogram. The best analogy here is to think of a Calder mobile. Different photos of such a mobile will yield different distances between hanging objects. While multidimensional scaling is a more accurate representation of the distance between objects, the fact that it attempts to represent the various points into a two- or three-dimensional space may result in distortion. As a consequence, some items that tend to appear together or be very similar may still be plotted far from each other.

Proximity Plot Tab

The **Proximity Plot** is the most accurate way to graphically represent the distance between objects by displaying the measured distance from one or several target objects to all other objects. It is not a data reduction technique but a visualization tool to help you extract information from the huge amount of data stored in the distance matrix at the origin of the dendrogram and the multidimensional scaling plots. In this plot, all measured distances are represented by bars, the closer an object is to the selected one, the longer the bar will be.



To select a keyword or a case that will be used as the point of reference, you can choose from the **Target Items:** drop-down checklist located at the top of the tab. You can also freely browse through different keywords or cases by double-clicking its bar on the **Proximity Plot**. The cooccurrence or similarity to more than one target item may be displayed in a single chart allowing for easy comparisons. When several target items are selected, the proximity plot may consist of bars clustered side-by-side (clustered bars) or stacked, representing either the total amount (stacked bars) or the relative distribution of scores (100-percent stacked bars). When two target items are selected, it is also possible to display the bars on both sides of a central axis like the sample chart above (mirrored bars).

By default, the chart displays the proximity of the target items to a maximum of 30 related items. Clicking the we button located in the upper left-hand corner of the chart, displays a dialog box that allows you to either manually choose items to be plotted or to automatically select a specific number of items.

When looking at keyword cooccurrences, selecting a bar enables the button. Clicking this button retrieves every document or text segment containing both keywords, allowing you to further explore the factors that may explain this cooccurrence. When examining the similarity of documents rather than keywords, clicking this button retrieves both documents and displays them side-by-side in a dialog box.

Right-clicking any existing bar displays a menu that allows you to remove the selected item, move it to the list of target items either by adding it to the existing bars or replacing one of them. You may also retrieve documents or text segments using this popup menu.

Table Tab

The **Table** tab allows you to examine in more detail the numerical values behind the computation of these plots. When the distance measure is based on cooccurrences, the table provides detailed information, such as the number of times a given keyword cooccurs with another one (**COOCCURS**) and the number of times it appears in the absence of this selected keyword (**DO NOT**). Such a table also includes the number of times the selected keyword appears in the absence of the

given keyword (**IS ABSENT**). In the example below (computed using the paragraph as the frequency criteria), we can see on the highlighted line that the category **WORK LIFE BALANCE** cooccurs 115 times with **BENEFITS**, but this word is encountered in 230 paragraphs without the word **BENEFITS**, while **BENEFITS** is found in 914 paragraphs in the absence of **WORK LIFE BALANCE**. The Jaccard coefficient of 0.091 indicates that of all paragraphs that contain either one of these words, only 10.9 percent contains both words. Note, however, that not all proximity measures can be interpreted this easily. To facilitate the interpretation of this table, the status bar provides a textual interpretation of some of the statistics.

WordStat	9.0.7 - Election 200	8 Coded.ppj							1	×
E Dat	a 🚽 Text Proc	essing F	Frequenc	ies 😗 Ext	raction	S Cooccu	nces 🛅 Crosstab 🚛 Keyword-In-Context < Classification			2
11 Options	B Dendrogram	п 🚼 Мар	ping 🔆	Link Analysis	T Pr	oximity Plot	Σ Statistics E ^{rr} Cooccurence Matrix			
Target I	tems: [MILITARY;T	ERRORISTS]						đ		
Proximity Plot	Table									
TARGET	KEYWORD	CO-OCCURS	DO NOT	IS ABSENT	Jaccard	1				1
MILITARY	IRAQ	111	621	284	0.109					1
MILITARY	WAR	103	697	292	0.094					
TERRORISTS	OAEDA	26	129	130	0.091					
TERRORISTS	AL	26	147	130	0.086					
TERRORISTS	NUCLEAR	31	218	125	0.083	IIII				
MILITARY	TROOPS	52	246	343	0.081					
TERRORISTS	WEAPONS	21	138	135	0.071					
MILITARY	AEGHANISTAN	36	131	359	0.068					
TERRORISTS	AEGHANISTAN	18	149	138	0.059		N			
MILITARY	PEACE	32	147	363	0.059		12			
MILITARY	SECURITY	50	401	245	0.055					
MILITARY	OAEDA	30	106	345	0.056					
MILITARY	QAEDA	29	120	200	0.050					
MILLIART	AL	29	144	300	0.054					
TERRORISTS	WAR	4/	/53	109	0.052					
MILITARY	WOMEN	31	288	364	0.045					
TERRORISTS	IRAQ	36	696	120	0.042					
TERRORISTS	SECURITY	26	515	130	0.039					
MILITARY	ADMINISTRATION	28	315	367	0.039					
MILITARY	ECONOMIC	27	367	368	0.035					
TERRORISTS	TROOPS	15	283	141	0.034					
MILITARY	VETERANS	17	102	378	0.034					
MILITARY	NUCLEAR	20	229	375	0.032					
MILITARY	IRAN	17	234	378	0.027	1				
MILITARY	BUSH	22	439	373	0.026	1				
MILITARY	BILLION	17	276	378	0.025	1				
TERRORISTS	ISRAEL	6	90	150	0.024	1				
MILITARY	ECONOMY	23	572	372	0.024	1				
TERRORISTS	PEACE	8	171	148	0.024	1				
TERRORISTS	TECHNOLOGY	8	183	148	0.024	1				
TERRORISTS	IRAN	9	242	147	0.023	1				
TERRORISTS	LAW	9	228	147	0.023	1				
MILITARY	NATO	10	46	385	0.023	1				
MILITARY	WEAPONS	12	147	383	0.022	1				
TERRORISTS	FREEDOM	8	208	148	0.022	1				
TERRORISTS	INFORMATION	6	135	150	0.021	1				
	CLORAL		202	201	0.000	1				2

The following list provides a brief description of the buttons found on these two tabs

Proximity plot controls:

- By default items in the proximity plot are sorted in descending order of proximity scores. Clicking this button sorts items in alphabetical order.
- This button is used to create a copy of the chart to the clipboard. When this button is clicked, a pop-up menu appears, allowing you to select whether the chart should be copied as a bitmap or as a metafile.
- This button allows editing of various features of the proximity plot, such as the appearance of value labels and bars, the chart and axis titles, and the location of the legend.

Proximity plot and proximity table controls:

Press this button to append a copy of the chart or the proximity table in the Report Manager. A descriptive title will be provided automatically. To edit this title or to enter a new one, hold down the SHIFT keyboard key while clicking this button (for more information on the Report Manager, see the <u>Report Management Feature</u> topic).

- When the proximity plot is displayed, this button allows storing the chart on disk in one of the supported graphic file formats. When the proximity table is shown, the table can be saved to disk in an Excel, text delimited, XML, HTML or SPSS file.
- Clicking this button allows you to print a copy of the displayed chart or table.

Link Analysis

The **Link Analysis** tab allows you to visualize the connections between keywords or dictionary items using a network graph. It offers a high level of interactivity, allowing you to explore relationships as well as detect underlying patterns and structures of cooccurrences using three layout types: a multidimensional scaling, a force-based graph, and a circular layout.

This tab is initially associated with the clustering features and the **Dendrogram** tab such that selecting a specific cluster in the dendrogram will result in a network view of its elements - where each element is represented as a node, while their relationship will be represented as a line connecting these nodes (also called an "edge"), the thickness of this line representing the strength of this relationship. Nodes and edges can be edited, deleted or moved, and new nodes and edges may also be added, either manually or using statistical criteria.



Layout Selection

The following three buttons allow you to select one of three layout types.

Clicking this button assigns the location of nodes in a multidimensional space such that nodes that cooccur more often are plotted close together, while those cooccurring less often are plotted far from each other.

- Clicking this button draws a force-directed graph in which links are more or less of equal length and there are as few crossed links as possible.
- Clicking this button draws nodes in a circle. Nodes that cooccur often are plotted close together.

Nodes and Links Selection

- Clicking this button allows you to select nodes and links, either manually or based on statistical criteria. See the Selecting Nodes and Links section below for more information on this feature
- This track bar can be used to hide links. By default, the cursor is positioned to the right. Moving it to the left gradually removes the weakest links, allowing you to more clearly identify the strongest ones.
- By default, a graph includes links reaching a specific statistical threshold. Clicking this button adds additional links by gradually reducing this statistical threshold.
- V Clicking this button removes all links that have been hidden using the track bar control (see above), deletes all nodes that are no longer connected, and refreshes the display to take into account the lower number of elements in the graph.

Navigation Mode

The following three buttons allow you to select the default mouse behavior.

Click this button to select specific nodes and links. To select multiple items, hold the Ctrl key while clicking additional items or select the rectangular region of the graph containing the items to select.

To select a node along with all nodes connected to it (neighbors), hold the Alt key while clicking the node. You may also select a node, right-click and choose GET NEIGHBORS.

Once items are selected, they can then be moved, deleted or edited. You may also search for associated text segments either by right-clicking to call up a contextual menu or by clicking the 🕄 button on the toolbar.

- Clicking this button lets you control which part of the image is visible in the image window.
- Clicking this button allows you to zoom into a specific region. Once activated, click the upper-left corner of the rectangular region you want to zoom in on, drag the mouse down the lower-right corner of this region and then release the mouse button.

Zooming In and Out

- Clicking this button zooms in or out so you can see the entire map.
- Clicking this button increases the zoom level of the map.
- Clicking this button decreases the zoom level of the map.



This list box allows you to set the zoom level to predefined levels. You can also type in the desired zoom level.

Editing Buttons

This button allows you to display or hide the numerical value associated with links.

- This button allows you to change the size of the nodes, links and their associated text. When clicked, a dialog box with four track bars appears. Moving the cursor to the right increases the size of its associated elements while moving it to the left decreases its size.
- Clicking this button deletes the selected nodes and links.
- Z Click this button to display a dialog box to change the style, color and size of the selected nodes and links.
- To retrieve segments associated with a node, select the node and click this button. When more than one node or when a link is selected, all text segments containing at least two keywords will be retrieved and presented in a table format. You may, however, change the type of segments retrieved (paragraphs, sentences or full documents) or the minimum number of topic items needed for retrieval.
- This button is used to send a copy of the graph to the clipboard. When this button is clicked, a pop-up menu appears, allowing you to select whether the graph should be copied as a bitmap or as a metafile.
- Click this button to append a copy of the graph into the Report Manager. A descriptive title will be provided automatically. To edit this title or to enter a new one, hold down the Shift key while clicking this button. (For more information on the Report Manager, see the <u>Report Management Feature</u> topic.)
- This button allows you to store the displayed network graph in a graphic file. WordStat supports three different file formats: .BMP (Windows bitmap files), .PNG (Portable Network Graphic compressed files) and .JPG (JPeg compressed files).
- Clicking this button allows you to print a copy of the displayed graph.

Selecting Nodes and Links

There are four methods available to select nodes to be plotted:

1. Selecting a Cluster on the Dendrogram tab

The link analysis graph is initially connected to the **Dendrogram** tab, so that selecting one of its clusters results in its items being represented as nodes on the link graph. The similarity criterion is automatically set to the value connecting the last item of this cluster to the other elements. To map connections between items of another cluster, simply move back to the **Dendrogram** tab and click anywhere on the new cluster. A smaller version of the link map should appear in the lower-right panel of the **Dendrogram** tab. You may then move back to the **Link Analysis** tab, or click the **S** button in the lower-right panel.

2. Double-Clicking a Node

When the graph is in selection mode, double-clicking a node will set the node as the target item and will display all nodes associated with it according to the last criterion settings (primary links). All currently displayed nodes not associated with this new target item will be removed.

If the Shift key is held down while double-clicking a node, the graph will also include nodes associated by a second degree to the target node (secondary links) as long as they also match the same criterion.

3. Selecting Nodes and Right-Clicking

When one or several nodes are selected, right-clicking brings a popup menu that allows you to either **ADD** to the current graph, primary links and secondary links or **REPLACE** the currently displayed nodes with those associated directly (one level) or indirectly (second level) with the selected nodes.

4. Using the Nodes and Link Filtering Dialog Box

• Clicking the Y button on the toolbar displays a dialog box that allows you to manually select items to be graphed. Such a dialog box looks like this:

Selection dialog		-	-		×
ource List:		Destination List:			
		Short Target Items Onl	y		
CLEAN A CLIMATE CLINTON COLLEGE COMPANY CONGRESS COST CREATE CREDIT CREATE CREDIT CRISIS CUT V	♦ ₹	AFGHANISTAN IRAN IRAQ JOB ISRAEL			
Selection Level: O Target Only O Primary Links					
Secondary Links	Direct	Links Only			
Node Selection: Similarity: 0.108 🖨 — ()-	-	Number of Nodes: 26	*		
Link Selection:					
Similarity: 0.108	-		_	0	Apply
		Number of Links: 27	÷	-	Close

The following groups of options are available for selecting nodes and links to be plotted:

- **Source List** This list box contains all items that are available to be plotted. To include them in a graph, simply move them to the **Destination List** using the arrow button. A **Search** edit box at the top of this list allows you to quickly find specific items.
- **Destination** This list box contains items that are currently plotted. When the **Show Target Items Only** option is clicked, the dialog box will offer the possibility of selecting additional nodes and links based on similarity or proximity criteria. When not checked, the graph will display the connection of only the nodes currently listed. Arrow buttons may be used to remove or add additional items.
- Selection When the destination list is set to include only target items, you may choose whether additional nodes should be plotted and how those will be selected. When set to **Target Only**, no additional items will be plotted,. Clicking **Primary Links** will allow the selection of additional nodes directly linked to target items based on a similarity criterion. When the **Secondary Links** option is chosen, additional items may be included if they are indirectly connected to the target nodes through the primary links, as long as they reach the same similarity criteria. By default, all links reaching such a criteria will be plotted, including links among the primary and secondary linked nodes. Selecting **Direct Links Only** will omit most of those links keeping only the ones that connect them to the target items.
- Node Selection This set of options allows you to adjust the similarity criterion to be used for selecting additional items. The similarity criteria displayed is the one currently set on the <u>Options</u> tab of the cooccurrence section. You can type a specific value (typically between 0 and 1) in the edit box, increase or decrease the value using either the **UP** or **DOWN** arrow buttons, or by sliding the cursor on the track bar on the right side of the edit box. The **Number of Nodes** option is automatically adjusted to indicate how many items meet the current similarity criterion. Editing its value will reciprocally adjust the similarity criterion.
- Link This set of options allows you to adjust the similarity criterion to be used for selecting links to be plotted. The similarity criterion for displaying links should always be equal to or less than the one used for selecting nodes. Reducing this criterion will plot additional links. You can type a specific value in

the edit box, and increase or decrease this value using either the **UP** or **DOWN** arrow buttons, or the slide track bar on the right of the edit box. The **Number of Links** option is automatically adjusted to indicate how many links will be plotted. Editing its value will reciprocally adjust the similarity criterion.

• Once the options have been set, click *CAPPLY* to update the graph and display the selected nodes and links.

Statistics

The **Statistics** tab displays the cooccurrence and similarity matrices used for building the dendrogram, MDS plots, as well as the proximity plot and table.

Data	Text Pro	cessing	Freq	uencies	8 Extr	action	S Cooce	currence	s 🗖 C	rosstab		eyword-Ir	n-Context	< d	lassification	1		1
			-	Serie	la Auralianta	Ŧ			Chattan	.E c								
Options S De	narogra	im 💽 i	Mapping	ter Lin	K Analysis	I P	roximity Pic	ot 2	Statistics	BC	ooccure	nce Matro	s					
Full Matrix																8	1	- 6
occurrences Similari	ty Agg	omeration																
	ACCESS	ACROSS	ACT	ADMINISTRATION	AFFORD	AFGHANISTAN	AFTER	AL	AMERICA	AMERICANS	BASED	BENEFITS	BIG	BILL	NOITION	BLACK		
ACCESS																		
ACROSS	6																	
ACT	2	5																
MINISTRATION	9	12	19															
AFFORD	4	8	5	6														
AFGHANISTAN	0	6	10	9	1													
AFTER	4	21	13	24	15	16												
AL	0	7	11	6	1	46	18											
AMERICA	22	107	45	52	36	22	72	27										
AMERICANS	8	41	25	22	39	8	47	18	158									
BASED	4	3	4	6	1	3	7	6	25	18								
BENEFITS	6	7	1	6	5	2	7	3	31	23	6							
BIG	3	10	4	5	8	1	8	1	39	28	2	6						
BILL	7	9	11	9	10	2	24	3	25	27	4	8	9					
BILLION	11	13	5	12	13	8	8	2	55	30	3	4	6	11				
BLACK	0	8	0	1	3	2	6	2	15	14	2	0	1	0	0			
BUSH	6	15	19	119	13	11	54	13	64	51	6	4	12	10	14	3		
BUSINESS	6	7	7	14	7	0	14	0	50	22	4	4	17	11	13	1	6	
BUSINESSES	2	10	7	2	9	0	5	0	46	30	6	10	8	4	13	2	2	
																		>

The first two tabs display by default, only the lower triangular part of the matrix of cooccurrence and similarity. Selecting the **Full Matrix** option will display data on both sides of the diagonal.

To export a matrix to social network analysis software tools:

- Select the matrix you would like to import by activating either the cooccurrences or the similarity tab.
- Click the Management
 button.
- In the **Save As Type** list box, select the file format under which to save the table. The following formats are supported: UCINET file (*.DL), Pajek Network file (*.NET), NetDraw file (*.VNA) and NetMiner file (*.SNF).
- Type a valid file name with the proper file extension.
- Click the **SAVE** button.

To append a copy of the table in the Report Manager:

• Click the 🛍 button. A descriptive title will be provided automatically for the table. To edit this title or to enter a new one, hold down the **Shift** key while clicking this button.

For more information on the Report Manager, see the Report Manager topic.

To export the table to disk:

- Click the 🔚 button. A Save File dialog box will appear.
- In the **Save As Type** list box, select the file format under which to save the table. The following formats are supported: ASCII file (*.TXT), Tab delimited file (*.TAB), Comma delimited file (*.CSV), MS Word (*.DOC), HTML file (*.HTM; *.HTML), XML files (*.XML) and Excel spreadsheet file (*.XLS).
- Type a valid file name with the proper file extension.
- Click the **SAVE** button.

To print the table:

Click the button.

Cooccurrence Matrix

The co-occurrence matrix feature allows one to focus on specific co-occurrences. The main results consist of a table displaying a choice from various co-occurrence statistics. Such a matrix is also highly interactive allowing one to transform specific rows into new columns or vice versa using simple drag-and-drop operations. A charting panel on the left also allows one to assess the distribution of a specific co-occurrence across other variables. One may also obtain a quick view of all text segments associated with a specific co-occurrence. This feature may also be called from the frequency list by selecting target items (words or content categories) that should be displayed as columns, right-clicking, and selecting Co-Occurrence Matrix.

Calling the cooccurrence matrix from the Frequency table

- Select the rows containing the words or content category you would like to use as references. You can select multiple disjunct rows by keeping the CTRL keyboard key down while clicking on the desired rows.
- Right-click to display the popup menu, or click the ≡ button on the left of the Frequencies page toolbar and select the **co-occur with...** command. A dialog box will appear with the selected items as column items and all other words or content categories of the frequency table as row items.

Accessing the co-occurrence matrix feature from the cooccurrences page.

The cooccurrence matrix is also accessible as a separate tab from the main **Cooccurrences** page. When accessing this feature from this location, the unit of text on which cooccurrences are computed is automatically adjusted to the general setting used for clustering, multidimensional scaling and other related features. Moving directly to the **Cooccurrence Matrix** page will result in all clustered items to be displayed as rows without any target item. Selecting a cluster on the **Dendrogram**

page prior to moving to this page, will assign all words of the selected cluster as column items allowing you to see how all items of the selected clusters are related to the remaining items.

WordStat 9.	0.7 - Election 2008	Coded.ppj	nguancias 🌒 Extra	notion O Consciu	ransas 🔲 Grassit	ab III Konword	In Contaxt	- Class	fication		-		×
+1 Options		Mappin	a Ve Link Analysis	T Dravinity Diat		Conscurance Mat	-in-context	Cidos	incauon				
19 Options					2 Statistics	cooccurence mat	ла						
Count: S	entence 🗸	Display: Die	ce coefficient	~ &) 🛄 🔃	_							(a r	
	ECONOMY	EDUCATION	TECHNOLOGY		^		WOMEN	with EDUCA	TION				
SCHOOL	0.000	0.081	0.010		Comparisor	ns 🔊 Segments							
HEALTH	0.042	0.072	0.040			Case Occurence	V RV	CANDIDATE	~				=
SCIENCE	0.015	0.065	0.115					CANDIDATE					_
CARE	0.032	0.062	0.024		Biden-		1				-	-	
QUALITY	0.009	0.060	0.015		Obama-		-						-
JOBS	0.114	0.039	0.019		Richardson-							-	
GLOBAL	0.140	0.033	0.022		Romney-							-	_
ENERGY	0.072	0.043	0.050		Clinton -								
SYSTEM	0.027	0.055	0.025		Edwards-							1	_
TAX	0.041	0.021	0.005		Thompson-	and the second							_
FAMILY	0.006	0.022	0.003		Giuliani -							1	
WOMEN	0.006	0.034	0.003		MCCain								
RESEARCH	0.010	0.021	0.077		McCall		-					-	
LAW	0.000	0.017	0.008			0 0.5	1	1.5 2	2 2	.5	3	3.5	4
BILLION	0.024	0.028	0.038		1 to 🖶 Inc.	10			-				-
WORKERS	0.055	0.016	0.012			Case Occurence	~ BY	DELIVERY	~				=
PLAN	0.037	0.018	0.013		6			1					
VETERANS	0.007	0.016	0.000					*					
MCCAIN	0.046	0.013	0.004		4			/					
COSTS	0.017	0.020	0.008		2-		/		-	-		1	-
SECURITY	0.049	0.015	0.017		0	1 1	×	i	1	1		-	
SENATOR	0.032	0.008	0.006		2008	2008 2008	2008	2007	2001	2007	200	n .	2000
GENERATION	0.013	0.007	0.026		a* a3	02	a	der (2	2	ar		
			0.010		v 1								

Adjusting rows and column items

To transform a row item as a column item:

- Move the mouse over the item name (first column), click on it, and keep the mouse down
- Move the mouse cursor up to the first row at the location where you want this item to appear and release the mouse button.

or

- Click on any cell on the row of the item you want to move to a column. You may also select a range of cells.
- Right-click to display the contextual menu, select **Move** and then choose the **Selected Rows as Columns** command.

To transform a column item as a row item:

- Move the mouse over the item name (first row), click on it, and keep the mouse down
- Move the mouse cursor down to the first column at the location where you want this item to appear and release the mouse button.

or

- Click on any cell on the column of the item(s) you want to move to rows..
- Right-click to display the contextual menu, select **Move** and then choose the **Selected Columns as Rows** command.

To adjust the order of column or row items:

- Move the mouse over the item name you want to move, click on it, and keep the mouse down and drag it to its new location. Green arrows will indicate the location where the item will be inserted.
- One over the desired location, release the mouse button.

To sort rows:

- Double-click on the column header containing the values on which you want rows to be sorted. When clicking on the first column header, items will be sorted alphabetically, while clicking on any other column will sort the rows in ascending order of the numerical value in this column.
- Double-clicking again on the same column header will sort the rows in descending order.

Adjusting display options

The **Count** listbox allows you to specify the size of the text segment that will be used for computing co-occurrences. You can obtain cooccurrence statistics computed on the sentence, the paragraph, the full document, and when there is more than one document selected, on the cases. When only one document variable is analyzed, computed statistics for documents and cases will be identical.

The **Display** list box allow you to create a table using any one of the following statistics:

- Frequency
- % of column item
- % of row item
- Jaccard coefficient
- Dice coefficient
- Adjusted Phi
- Inclusion index

Charting cooccurrences of selected cells:

The **Comparisons** tab on the right panel allows one to compare how co-occurrences of two items are distributed among the selected independent variables using a bar chart or a line chart. Please note that while the charting is linked to the table selection, only one pair of items is being plotted at a time even if multiple cells are selected. The elements of this pair appear at the top of the right panel. One also must remember that what is being charted is not how often those two words are being used, but the absolute or relative frequency by which those two words are being used together. The **Case Occurrence** statistics simply plots how many cases such a cooccurrence is present, irrespective of its frequency while the **Category Percent** computes the percentage of cases for this specific value of the independent variable containing this cooccurrence. All other statistics are computed to represent either the absolute frequency or the relative frequency of one word relative to the other word, or, for the **Jaccard** or **Dice Coefficient**, relative to the presence of either one of them.

One may also represent graphically the co-occurrence statistics of several items by selecting the range of cells to display and by clicking the **iii** button. Alternatively, one may select cells, right-click and select the **Chart Selected Cells** command. For more information on the charting options of this dialog box, see <u>Bar Chart and Line Chart Options</u>.

Creating a bubble chart

Bubble charts are graphic representations of statistical tables where the displayed statistics are represented by circles of different diameters. One can create a bubble chart for representing the entire grid or only selected cells. To create a bubble charts for selected cells only:

- Select the cells you want to plot.
- Right click and choose Create Bubble Chart or click the 🔛 button, and then chose Selected Cells..
- Select Entire grid... instead to create a bubble chart to represent the values contained in the entire table.

For more information, see **Bubble Charts**.

Retrieving associated text segments

The **Segments** tab on the right panel allows one to quickly retrieve all text segments where both items of the selected cooccurrence appear. Please note that while this feature is linked to the table selection, only one pair of items is being searched for at the time even if multiple cells are selected. The elements of this pair appear above the right panel. The top panel displays the extracted text segment while the bottom part displays the entire document from which this segment come from.

One may further explore the context of such a cooccurrence by clicking the substant to produce a word frequency table and a word cloud from all the retrieved text segments. (See <u>Word Frequency Analysis</u> for more information on this feature).

WordStat 9	9.0.7 - Election 2008	Coded.ppj				X
= 🛄 Data	a Text Proces	ssing 📻 Fre	equencies 😗 Extra	ction 🗞 Cooccurrence	s Cr	sstab 🏢 Keyword-In-Context < Classification
11 Options	B Dendrogram	💦 Mapping	💥 Link Analysis	\mathbf{E} Proximity Plot Σ	Statistics	E ^m Cooccurence Matrix
E Count:	Sentence 🗸	Display: Dic	e coefficient	85 🔟 🔛		M H A
	ECONOMY	EDUCATION	TECHNOLOGY	^		WOMEN with EDUCATION
SCHOOL	0.000	0.081	0.010		Compa	risons 🔊 Segments
HEALTH	0.042	0.072	0.040		in compa	
SCIENCE	0.015	0.065	0.115		Word	12 hits
CARE	0.032	0.062	0.024		Case	Text
QUALITY	0.009	0.060	0.015		1-Richardso	WOMEN like my EDUCATION Secretary Dr.
JOBS	0.114	0.039	0.019	F	1-Richardso	Unleashing the economic power of WOMEN through EDUCATION can be the silver bullet that
GLOBAL	0.140	0.033	0.022		n20071018	makes every problem easier to fight.
ENERGY	0.072	0.043	0.050		1-McCain20	Wade and undermining the constitutional protections that decision provided, but by preventing
SYSTEM	0.027	0.055	0.025		0/1011	care and services, empowering WOMEN to make decisions.
TAX	0.041	0.021	0.005		1-McCain20	And I want to invest in the idea-leaders of tomorrow by improving math and science EDUCATION,
FAMILY	0.006	0.022	0.003		080325	expanding fellowships for graduate study, and encouraging more WOMEN and minorities to choose careers in science and engineering.
WOMEN	0.006	0.034	0.003		1-McCain 20	whether it's afferschool funding or healthcare and FOLICATION for WOMFN here in America or
RESEARCH	0.010	0.021	0.077	t i	080331	for WOMEN in Afghanistan— we are a force to be reckoned with.
LAW	0.000	0.017	0.008		1-McCain20	Now, some might say that their work is finished in America since WOMEN no longer face legal
BILLION	0.024	0.028	0.038	-	080331	obstacles to EDUCATION or employment or the right to vote.
WORKERS	0.055	0.016	0.012		1-McCain20 080403	Q when WOMEN continue to be raped as a casualty of conflict, trafficked for commercial advantage, denied EDUCATION and health care and family planning, not given access to credit.
PLAN	0.037	0.018	0.013	-		denied their rights as citizens, that not only affects them and their countries, that compromises
VETERANS	0.007	0.016	0.000		school pri	ncipal, a superintendent, and a reporter before joining the women's suffrage
MCCAIN	0.046	0.013	0.004		movemen	t and realizing her talents as a gifted organizer and dynamic orator.
COSTS	0.017	0.020	0.008		helped to	ro stints as President of the National American Woman Suffrage Association and she found the International Woman Suffrage Alliance to reach out to WOMEN across the
SECURITY	0.049	0.015	0.017		world. In t	ne end, it was Catt who devised the "Winning Plan" for the suffrage movement a plan
SENATOR	0.032	0.008	0.006		to campai	gn simultaneously for suffrage at both the state and national levels. The rest, as they tory
GENERATION	0.013	0.007	0.026		And it real	ly is quite a history, isn't it? There has never been a better time to be a woman in
reanaute_				×	America I	i's almost hard to explain to young WOMEN today how much things have changed **
245 cases					Shown	: 60 Types: 15,813 Tokens: 590,701 Time: 1.8s (4 cases excluded)

The Crosstab Tab

The **Crosstab** tab is used to display a contingency table of words or categories. This contingency table is computed only on items that have been included. If a categorization dictionary has been specified, this grid will display only the words or keywords in this list. If no dictionary has been specified, the grid will display all words that have not been explicitly excluded. Along with absolute and relative frequency of keyword occurrence or keyword frequency, several statistics may be displayed to assess the relationship between independent variables and word usage or to assess the reliability of coding made by several human coders or a single coder at different times.

🛄 Data 🛛 🗧 Tex	Process	ing F	Frequencies	Extract	tion 🗞	Cooccurrer	nces 💼	Crosstab	Keyword-	In-Context	< Classific	ation					
abulate: Total Frequency	· ~	Displa	ay: Rate per	10,000 ~	Statis	tic: None		~ =	85 🔟 🗄	ti						đ	
With: CANDIDATE	~	Sort	2y: Keyword	~	1				0 11 11	6							
	Biden	Clinton	Edwards	Thompson	Giuliani	Kucinich	MCCain	Obama	Richardson	Romney	^	» Filter: All		~			
ADMINISTRATION	17.29	10.09	8.84	3.29	1.50	22.54	3.78	3.90	11.79	0.81		L. E. ha	1				-
AFGHANISTAN	10.20	1.14	2.95	1.64	3.76	2.25	4.31	4.31	4.82	0.54			Rate per 10,000 wo	ords ~	BY CANDIDAT	E	~
AL	5.76	0.57	2.49	1.23	2.63	2.25	4.31	4.08	13.67	1.88		Biden-			-		
BILLION	9.31	8.54	2.27	0.00	2.63	5.63	2.63	6.59	10.45	4.56		Clinton			_		
BUSH	7.54	9.52	18.35	2.47	2.25	4.51	1.47	10.78	17.42	5.37		Edwards-		-			-
CARBON	0.00	2.28	0.23	0.00	0.00	4.51	3.36	1.52	4.56	0.00		Thompson-					-
CARE	7.98	35.96	10.65	4.11	2.63	25.92	24.07	19.00	13.67	5.90		Giuliani -					
CHANGE	13.75		13.82	4.52	3.38	3.38	8.62			12.08		MCCain -	-				_
COMPANIES	0.44	12.86	5.44	0.00	0.38	2.25	6.31	7.75	4.82	1.88		Obama-					-
COSTS	0.44	9.93	1.59	1.64	2.25	0.00	6.52	4.72	4.29	3.22		Richardson-			_		-
CRISIS	0.89	7.57	1.81	0.00	3.38	0.00	4.10	6.53	2.14	2.95		Romney		-	-		-
ECONOMIC	1.77	6.18	5.21	7.40	2.63	4.51	11.67	5.83	10.19	15.30		0	5	10	15		20
ECONOMY	0.89	14.16	10.42	8.22	4.88	11.27	11.56	13.29	9.11	14.22		1.0.00					-
EDUCATION	4.43	4.07	2.27	2.47	1.88	9.01	3.68	7.58	15.82	5.64			Rate per 10,000 wo	ords 🗠	BY DELIVERY		Y
ENERGY	3.99	22.13	4.53	2.47	5.64	6.76	12.51	13.81	25.47	9.66			1 1	1	1 1	1	-
FAMILY	1.77	11.47	2.95	1.23	3.76	4.51	5.57	7.46	4.02	16.10		15		-		-	-
FREEDOM	2.66	0.73	2.95	7.81	7.89	6.76	8.93	3.50	1.88	13.42		10	1				
FUEL	0.00	1.95	0.68	0.00	0.00	1.13	2.63	4.08	5.90	2.15		5-0-0					
GENERATION	2.22	3.99	7.02	8.22	9.02	0.00	2.10	5.89	2.14	8.05		0-4		1	1	- 1	
GLOBAL	0.44	6.83	10.88	1.23	1.88	6.76	7.57	4.95	15.82	2.15		-00 ⁶ .00 ¹	100, 100,	001	000 000	008	S
HEALTH	5.76	30.67	10.65	2.47	7.14	23.66	17.98	16.49	18.50	4.29		012	02 03° 04	01	at a	*	QA'L
INFORMATION	3.99	9.28	0.45	0.41	8.27	1.13	2.42	0.82	1.34	0.27						-	-
INSURANCE	2.66	11.23	1.59	1.64	4.13	4.51	9.78	6.47	3.48	1.61		the late F	Plot:	Frequenc	v v		
IRAN	5.32	1.30	9.97	2.06	6.01	9.01	5.26	3.61	21.18	11.27		Cloud L		requerk	-1		
IRAQ	74.05	7.81	21.53	6.99	8.64	13.52	17.76	18.71	23.59	4.83							
ISRAEL	0.89	0.08	0.91	0.00	3.01	0.00	3.89	4.55	1.88	4.29		3EN	ATOR COMPANIES	BUSH	HELRANCE SYST		-
JOBS	3 10	16.44	12 69	123	3.01 N	7 89	12 72	15.85	10.45	5.64			ENERGY TA	OL.	WAR M	210	8
LAW	9.31	5.61	3.17	11.92	3.01	12.39	5.47	3.09	2.41	6.44		ECO	NOMY JOBS CA	RE	CHANGE a	AL	
MCCAIN	0.00	2 20	0.68	0.00	0.00	0.00	1.37	13.99	0.00	1.07		MILE	TARY IRAQ PL	AN HE	ALTH ADA	NISTR	ATION
MEDICAL	0.44	7.57	0.00	0.82	1 13	1 13	3.05	1.92	4 29	0.27		PERCENT W	ASHINISTON SECUR	YTIS	WORKER!	-	
INCO IONE	0.00	2.02	0.00	0.00	0.20	1 13	3.47	1.57	0.00	10.73	~	na state	104 101 10 00 00	Station EDUC	NTON		

Tabulate: This option allows you to choose whether the values in the table should be based on the total frequency of keywords or the number of cases containing the keywords.

With: This drop-down list allows choices on how the keyword count should be broken down. The following options are available:

- <other keywords>: display a square table showing the number of cooccurrence of keywords in the same case.
- <case number>: display the keyword occurrence or frequency for each individual case.
- Any independent variable: If numeric, categorical or date variables were selected as independent variables, their names will appear in this list box. Selecting any of those variable names will display a contingency table allowing for the assessment of the relationship between this variable and the keywords or content categories. When a date variable is selected, a dialog box like the one below will appear allowing you to automatically recode those dates into various time periods or date units such as week days or months.

Transform TWEET_DATE	×
Transform dates into:	
OHours	
O Date & Time	
Days	
O Week days (7)	
O Weeks (starting day)	
O Months (12)	
O Months & Years	
O Quarters & Years	
() Years	
ODecades	
ОК	

- **<variables>:** When content analysis is performed on several alphanumeric variables, this option will allow for comparison of the occurrence of keywords according to variable name.
- Select variables...: By default, the variables listed in the With list box are the independent variables that have been selected when calling WordStat. Choosing this option displays a dialog box that allows you to select other numeric, categorical or date variables.
- **Combine variables...:** This option allows you to compare the frequency of keywords or content categories among the combined values of two variables. When this item is selected, a dialog box appears allowing you to choose the two variables whose values will be combined.

Sort by: This option presents the opportunity to sort the table by word or category names (alphabetical order) or by descending order of frequency or case occurrence. When a statistic is displayed (see option **Statistic**), the table can also be sorted based on the value of this statistic or on its statistical probability. It is also possible to sort on the values of any specific column by clicking this column heading. Clicking several times on the same column heading toggles between ascending and descending orders.

Display: This list box allows you to specify the information displayed in the table. The following options are available:

- Count
- Row percent
- Column percent
- Total percent

When the Tabulate option is set to Case Occurrence, two additional statistics are also available:

- Percent of cases (percentage of all cases or individuals)
- Category percent (percentage of cases or individuals in this subgroup)

Statistics: When keyword frequency or occurrence is broken down by an independent variable (see **With** option), a drop-down list will appear. This list box allows you to choose among 12 association measures to assess the relationship between this independent variable and the utilization of each word or category.

Nominal level statistics

- Chi-square
- Likelihood ratio
- Student's F

Ordinal or internal level statistics

- Tau-a
- Tau-b
- Tau-c
- Somers' D (symmetric)
- Somers' Dxy (asymmetric)
- Somers' Dyx (asymmetric)
- Gamma
- Spearman's Rho
- Pearson's R

Probability: The probability option allows you to select whether the probability value should be computed using a 1-tailed or 2-tailed test. Probabilities of Chi-square, Likelihood ratio, and Student's F are always computed using a 2-tailed test.

To display column labels vertically, click the $\frac{3}{2}$ button located to the right of the **Tabulate With** list box.

The C button is used to reapply the content analysis process to the current data set. This button is disabled by default and becomes enabled when changes are made to any one of the currently active text-analysis processes, such as the categorization dictionary, the exclusion list or the substitution process. Clicking this button will instruct WordStat to reprocess the text collection and update the current table.

The Comparisons Panel

The **Comparisons** panel on the right allows you to look at the distribution of the selected words or content categories among values of up to two structured variables. It may be used to get a quick graphic representation of the distribution of the variable currently used in the crosstabulation table but may also be adjusted to display the distribution on other variables. You may display such distributions using either a vertical bar chart, a horizontal bar chart or a line chart, by clicking on the corresponding button. Four statistics may also be represented on those charts:

- Case Occurrence Number of cases in this subgroup containing at least one of these words or category items.
- Category Percent Percentage of cases in this subgroup containing at least one of these words or category items.
- **Word Frequency** Total number of selected items in the subgroup.

Rate per 10,000 Words Rate of words in this subgroup per 10,000 words.

The bottom chart contains the distribution of selected items in the crosstab page using either a word cloud, a vertical of horizontal bar chart, a pie chart, or a donut chart. Right-clicking anywhere in the chart areas displays a pop-up menu that allows you to edit the chart, save it to disk or in the **Report Manager**, or to copy it to the clipboard. Clicking a specific bar or a data point of a line chart also allows one to retrieve text segments associated with the selected class and containing words of the selected topic.

The comparison panel may also be used to filter data on one or several values of the variables currently used for comparison in the crosstab page.

To filter cases and compare distribution on other variables:

Adjust the Filter list box at the top of the comparison panel to the value you want to filter on. The two graphs below
will be automatically adjusted to display the distribution of the selected items on the chosen variables for only cases
corresponding to the selected value. The word cloud at the bottom will represent the distribution of all items for cases
corresponding to the selected value. Changing from one value to another may allow one to identify meaningful
differences in distributions associated with those values.

- It is also possible to examine the distribution on a combination of values by selecting **Multiple..** and then putting check marks beside all values that you want to combine and filter on.
- To disable the filter on values, simply set the Filter list box to All.

Creating Bar Charts or Line Charts

The Crosstab tab also allows you to produce bar charts or line charts to visually compare the distribution of specific words or categories among values of an independent variable such as subgroups of individuals (male vs. female) or time periods.

To produce such charts:

- Set the **Tabulate** and **Display** options so that the information to visualize is displayed in the table.
- Using the mouse, select the rows you would like to display. Multiple but separate rows can be selected by clicking while holding down the Ctrl key.
- Click the *is button or right-click the mouse and select the Chart Selected Rows command.*

Creating Bubble Charts

Bubble charts are graphic representations of contingency tables where relative frequencies are represented by circles of different diameters.

To create a bubble chart:

- Set the **Tabulate**, **Count** and **Display** options so that the information to view is displayed in the table.
- Click the 👪 button.

For more information, see **Bubble Charts**.

Creating Heat Maps with Clustering of Rows and Columns

A heatmap plot is a graphic representation of crosstab tables where cell frequencies are represented by different color brightness or tones. When combined with clustering of rows and/or columns, this exploratory tool allows you to identify functional relationships between specific keywords and subgroups defined by values of the independent variable.

To create a heatmap:

- Set the With option to an independent variable.
- Set the Tabulate option to either Case Occurrence or Keyword Frequency.
- Click the button to access the heatmap dialog box.

For more information on heatmaps, see <u>Heatmap plot</u>.

Performing Correspondence Analysis

Correspondence analysis is an exploratory technique that provides a graphic overview of relationships in large crosstabulation tables of frequency.

To perform a correspondence analysis:

- Set the With option to an independent variable.
- Set the Tabulate option to either Case Occurrence or Keyword Frequency.
- Click the the button to access the correspondence analysis dialog box.

For more information on correspondence analysis, see Correspondence Analysis.

Other Tasks

To export the table to disk:

- Click on the 🖬 button. A Save File dialog box will appear.
- In the **Save as Type** list box, select the file format in which to save the table. The following formats are supported: ASCII file (*.TXT), Tab delimited file (*.TAB), Comma delimited file (*.CSV), HTML file (*.HTM;*.HTML), and Excel spreadsheet file (*.XLS). Type a valid file name with the proper file extension.
- Click the Save button.

To append the table to the Report Manager:

- Click the 11 button. A descriptive title will be provided automatically for the table.
- To edit this title or to enter a new one, hold down the **Shift** key while clicking this button.

(for more information, see the <u>Report Manager</u> topic).

To copy the entire table to the clipboard:

- Right-click anywhere in the frequency table.
- Select the COPY TO CLIPBOARD | TABLE command from the pop-up menu.

To copy selected rows to the clipboard:

- Select the rows you would like to copy (multiple but separate rows can be selected by clicking while holding down the **Ctrl** key).
- Right-click anywhere in the frequency table.
- Select the COPY TO CLIPBOARD | SELECTED ROWS command from the pop-up menu.

To search for a specific item:

- Right-click anywhere in the table.
- Select the FIND command from the pop-up menu. A search dialog box will appear.

- Type the desired search string in the **Find What** edit box. To restrict the search to items starting with the search string, enable the **Match Beginning of Item** option and to restrict it to whole words matching the search string, enable the **Match Whole Word Only** option.
- Click the **Find** button to search the first item matching the typed string. Clicking this button again finds the next occurrence of the search string, starting at the currently selected item.

Bubble Charts

Bubble charts are graphic representations of contingency tables where relative frequencies are represented by circles of different diameters. This type of graphic chart allows you to quickly identify high-and low-frequency cells and is especially useful for presentation purposes. Many features of the chart can be customized to highlight specific findings. Rows and columns can be moved freely or deleted, and you can adjust the color of each cell as well as the fonts used in the chart.

The bubble chart represents graphically the underlying table and displays the corresponding measure used (code occurrences or frequencies, word counts or percentages) and will also reflect the currently selected display settings (such as count, percentage of rows or columns, etc.).

To create a bubble chart:

- Move to the **Crosstab** tab.
- Set the Tabulate option to either Case Occurrence or Total Frequency.
- Set the **With** option to the desired independent variable.
- Set the **Display** option to specify how this information will be displayed.
- Click the 👪 button. A dialog box similar to this one will appear:



To adjust the size of the bubbles:

• In the Bubbles group box at the top of the window, adjust the Size option to Small, Medium or Large.

To adjust the color of the bubbles:

- Select the cell or group of cells you would like to alter.
- Right click and select the SET COLOR command. A dialog box will appear letting you choose a specific color value.

To move a row or a column:

- Click anywhere on the cell header containing the row or column label and hold the mouse down.
- Drag the mouse cursor over the desired new location and release the mouse button.

To delete selected rows or columns:

- Select a cell range covering the rows or columns you would like to delete.
- Right-click to display a popup menu.
- Select the REMOVE command and then choose SELECTED ROWS or SELECTED COLUMNS depending on your desired action.

To adjust the font size and style of the title and labels:

• Click the **Change** button beside the **Title** or **Labels** options. A Font setting dialog box will appear, letting you change the font, the font size, style, and color.

To edit the title:

• Click anywhere in the title region and type in the new title. The height of the title region is automatically adjusted to the number of lines in the title.

The following table provides a short description of additional buttons:

Control Description

- Pressing this button transpose the grid so that rows become columns while columns are transformed into rows. Such a feature may be useful to flip rectangular charts and choose between a portrait or a landscape display.
- Press this button to append a copy of the graphic in the Report Manager. A descriptive title will be provided automatically. To edit this title or to enter a new one, hold down the **Shift** key while clicking this button (for more information on the Report Manager, see the <u>Report Management Feature</u> topic).
- Click this button to save a chart on disk. Charts may be saved in BMP, JPG or PNG graphic file format or may be saved in a proprietary format (.WSX file extension) that can later be edited and customized using the Chart Editor.
- Clicking this button prints a copy of the displayed chart.
- This button creates a copy of the chart to the clipboard. When this button is clicked, a shortcut menu appears allowing you to select whether the chart should be copied as a bitmap or as a metafile.

Heatmap Plot

Heatmap plots are graphic representations of crosstab tables where relative frequencies are represented by different color brightness or tones and on which a clustering is applied to reorder rows and/or columns. This type of plot is commonly used in biomedical research to identify gene expressions. When used for text mining, such an exploratory data analysis tool facilitates the identification of functional relationships between related keywords and a group of values of an independent variable by allowing the perception of cell clumps of relatively high or low frequencies or of outlier values.

The heatmap plot used in WordStat allows you to graphically examine the relationship between keywords (rows) and values of an independent variable (columns). While the clustering available through the WordStat Frequency tab (see <u>Hierarchical</u> <u>clustering and multidimensional scaling</u>) is performed on individual cases and documents, the clustering analyses used in the heatmap plot are performed directly on the crosstabulation tables. As a consequence, the similarity index computed for two keywords and used for clustering does not represent their cooccurrences within cases but measures the similarity of their distribution among the various groups of the independent variable. Likewise, two subgroups defined by values on the independent variable will be considered near to each other if the distributions of keywords in those two groups are similar.

The heatmap plot can be performed either on keyword frequency or occurrences within cases (present or absent). For both types of analysis, the observed frequency is transformed into a percentage by dividing the frequency by either the total number of words in a specific subgroup (keyword frequency) or by the total number of cases in this subgroup (case occurrences).

To create a heatmap:

- Select the Crosstab tab.
- Set the Tabulate option to either Total Frequency or Case Occurrence.
- Set the With option to the desired independent variable.
- Click the 🗱 button.

Heatmap Tab

The main section on the right of this tab as shown in the screen shot below, displays the heatmap grid representing the relative frequencies of each cell (row and column intersection) using different brightness or color tones. Optional dendrograms are displayed at the top and to the left margin of this grid. The size of these dendrograms may be adjusted by moving the mouse cursor over the bottom edge (upper dendrogram) or the right edge (dendrogram on the left) of the dendrogram and dragging its limit to the desired size.



The font size used to display the row and column values may also be adjusted by clicking the \bigcirc or \bigcirc buttons located on the top toolbar. The size of cells and the distance between dendrogram leaves are automatically adjusted to the new font size.

To identify which specific cases or text documents are associated with a cell or group of cells, simply select the rectangular area of the list you would like to obtain, click the solution and select **documents**, **paragraphs** or **sentences**. WordStat locates the documents or text segments associated with these cells and displays them in the Keyword Retrieval dialog box.

Rows option group:

Clustering of rows: Enabling this option reorders the rows according to the result of a cluster analysis and displays a dendrogram at the left of the heatmap plot. Keywords that are distributed across the various subgroups in a similar way will tend to be grouped under the same cluster.

Sort by: When the clustering of keywords is disabled, rows may be sorted in alphabetical order, in descending order of frequency or case occurrence or displayed in their original categorization dictionary order.

No clusters: This option allows setting how many clusters the clustering solution should have. When two clusters or more are selected, horizontal red lines are drawn in the heatmap to delineate those clusters.

Columns options group:

Clustering of columns: Enabling this option reorders the columns according to the result of a cluster analysis and displays a dendrogram at the top of the heatmap. Columns with similar distribution of keywords will tend to be grouped under the same cluster. When the clustering option is disabled, the columns are presented in their ascending order of values. Disabling the clustering of columns is especially useful to preserve the ordinal nature of the values. For example, when looking for a relationship between some words or keywords and publication years, then disabling the clustering of columns will automatically sort the columns in chronological order allowing the identification of temporal trends.

No clusters: This option allows setting how many clusters the clustering solution should contain. When two clusters or more are selected, vertical red lines are drawn in the heatmap to delineate those clusters of keywords.

Heatmaps options group:

Colors: The colors list box allows the selection of various color schemes to represent differences in percentages. Monochrome schemes will express differences using levels of brightness of a single color where black always represents the lower limit of the selected range. Multicolor spectrums may be used to represent greater nuances in the selected range of percentages.

Minimum / Maximum: The minimum and maximum slide bars allow the setting of the range of values that will be used to display variations of color tones or brightness. By default, the minimum value is set to zero while the maximum value is set to the highest observed percentage. To increase the contrast at a specific location of this range, minimum and maximum limits may be adjusted. All cells with values lower than the minimum will be represented with the color located on the left end of the selected color spectrum while cells with percentages higher than the maximum limit will be displayed with the color locate on the right end of this spectrum. Reducing the range of values increases the contrast in a specific region of percentage values.

Keyword and Group Dendrogram Tabs

The second and third tabs of the heatmap dialog box allow a more detailed examination of the clustering of words or categories (Keyword Dendrogram) and of all values on the independent variable (Group Dendrogram). WordStat uses an average-linkage hierarchical clustering method to create clusters from a similarity matrix. The result is presented in the form of a dendrogram (see below), also known as a tree graph. In such a graph, the vertical axis is made up of the items and the horizontal axis represents the clusters formed at each step of the clustering procedure. Items that tend to be distributed similarly against the other variable appear together are combined at an early stage while those that have dissimilar distributions tend to be combined at the end of the agglomeration process.

Heatmap on keyv	vord frequency					07		×
Heatmap Keywor	d Dendrogram	Group Dend	ogram	Clustering Statist	ics			
Nb clusters: 30 🕃	Display: Agglo	meration Order	~ Q.	3. 10		1		-
AFGHANISTAN							- 23	^
IRAQ		-		-				
PLAN	-	-		011110000111				
TROOPS				_				
REFORM	-							
SECURITY						A		
AL	1		-					
QAEDA					7			
EDUCATION								
TPAN		_						
SCHOOLS								
FUEL						[0,1,2,-1] + [0,1] + [0,1] + [0,1]		
MEXICO								
CARS								
ISRAEL								
TAX							1.13	
ECONOMIC	-							
TECHNOLOGY							1 11	
ECONOMY								
FUTURE	-							
FAITH						1	£ 13	
FAMILY								
CHANGE							1 2 1	1
DREAM	-							
GENERATION								
WASHING TON								
FREEDOM	-			1				

No clusters: This option sets how many clusters the clustering solution should have. Different colors are used in the dendrogram to indicate membership of specific items in different clusters.

Display: This option sets whether the vertical lines of the dendrogram represent the agglomeration schedule or the similarity indices.

At times, it is desirable to see some areas of a dendrogram. The 🔍 and 🔍 buttons may be used to zoom in or out of the dendrogram.

Clustering Statistics Tab

The fourth tab of the heatmap dialog box allows the close examination of the matrices of similarities among keywords or among groups as well as the statistics associated with the agglomeration processes of both the keywords and the groups.

Heatmap on ke	eyword f	requency										-	
Heatmap Keyv	vord Den	drogram Gr	oup D	endrogra	m	Clustering Sta	tistics						
												85 🖬	
eyword x Keyword	Group x	Group Cluster	ing										
	AFFORD	AFGHANISTAN	AL	BILLION	BUSH	BUSINESSES	CARBON	CARE	CARS	CENTURY	CHANGE	CLASS	COLLEGE
ADMINISTRATION	0.53	0.76	0.67	0.84	0.71	0.43	0.75	0.77	0.32	0.55	0.64	0.73	0.61
AFFORD		0.72	0.73	0.82	0.73	0.78	0.59	0.70	0.80	0.76	0.83	0.83	0.92
AFGHANISTAN	0.72		0.80	0.82	0.70	0.54	0.52	0.60	0.45	0.64	0.80	0.79	0.66
AL	0.73	0.80		0.84	0.79	0.67	0.73	0.56	0.54	0.66	0.69	0.68	0.72
BILLION	0.82	0.82	0.84		0.83	0.73	0.74	0.81	0.68	0.76	0.83	0.91	0.87
BUSH	0.73	0.70	0.79	0.83		0.70	0.61	0.69	0.62	0.90	0.85	0.88	0.89
BUSINESSES	0.78	0.54	0.67	0.73	0.70		0.69	0.81	0.89	0.81	0.80	0.76	0.76
CARBON	0.59	0.52	0.73	0.74	0.61	0.69	-	0.84	0.51	0.53	0.51	0.56	0.57
CARE	0.70	0.60	0.56	0.81	0.69	0.81	0.84		0.66	0.72	0.73	0.82	0.73
CARS	0.80	0.45	0.54	0.68	0.62	0.89	0.51	0.66		0.79	0.82	0.69	0.75
CENTURY	0.76	0.64	0.66	0.76	0.90	0.81	0.53	0.72	0.79		0.91	0.86	0.84
CHANGE	0.83	0.80	0.69	0.83	0.85	0.80	0.51	0.73	0.82	0.91		0.90	0.85
CLASS	0.83	0.79	0.68	0.91	0.88	0.76	0.56	0.82	0.69	0.86	0.90		0.93
COLLEGE	0.92	0.66	0.72	0.87	0.89	0.76	0.57	0.73	0.75	0.84	0.85	0.93	
COMPANIES	0.77	0.50	0.52	0.75	0.75	0.88	0.68	0.92	0.75	0.81	0.77	0.87	0.85
CONGRESS	0.51	0.73	0.66	0.75	0.69	0.66	0.78	0.84	0.52	0.67	0.74	0.69	0.53
COSTS	0.75	0.51	0.56	0.74	0.65	0.92	0.65	0.87	0.78	0.76	0.73	0.82	0.76
CREATE	0.58	0.63	0.69	0.82	0.77	0.62	0.89	0.87	0.51	0.71	0.68	0.72	0.65
CREDIT	0.88	0.64	0.73	0.88	0.80	0.93	0.72	0.86	0.91	0.85	0.88	0.87	0.89
CRISIS	0.82	0.55	0.51	0.74	0.66	0.89	0.55	0.83	0.85	0.85	0.81	0.85	0.81
DEMOCRACY	0.25	0.55	0.42	0.54	0.57	0.32	0.58	0.63	0.25	0.59	0.54	0.49	0.31
DEMOCRATS	0.54	0.69	0.57	0.70	0.65	0.66	0.52	0.69	0.65	0.81	0.80	0.68	0.52
DREAM	0.78	0.75	0.55	0.78	0.74	0.64	0.37	0.65	0.71	0.88	0.89	0.86	0.78
ECONOMIC	0.57	0.57	0.69	0.71	0.69	0.86	0.65	0.72	0.78	0.76	0.78	0.66	0.59
ECONOMY	0.67	0.58	0.61	0.77	0.78	0.86	0.72	0.88	0.81	0.87	0.85	0.79	0.72

Correspondence Analysis

Correspondence analysis is a descriptive and exploratory technique designed to analyze relationships among entries in large frequency crosstabulation tables. Its objective is to represent the relationship among all entries in the table using a lowdimensional Euclidean space such that the locations of the row and column points are consistent with their associations in the table. The correspondence analysis procedure implemented in WordStat allows you to graphically examine the relationship between keywords or content categories and subgroups of an independent variable. The results are presented using a 2 or 3-dimensional map. Correspondence analysis statistics are also provided to assess the quality of the solution. WordStat currently restricts the extraction to the first three axes. This ensures that the results remain easily interpretable. To further differentiate among some values of the independent variable, we recommend applying a filter to restrict the analysis to a subset of values.

The first two tabs of the dialog box, provides graphical displays of correspondence maps. When the number of words or categories or the number of subgroups is less than four, only the 2-D Map tab can be accessed. When the comparison is restricted to two groups of individuals, only one axis can be extracted. In this situation, the axis is plotted diagonally in the two-dimensional space. We have done this mainly for readability reasons: plotting all the keywords on a single horizontal axis would have produced a cluttered list of keywords that would have made the graph useless. (The horizontal axis represents, by convention, the first extracted dimension).

The 2-D correspondence plot offers the ability to remove a keyword or a class of the independent variable and to automatically recompute the analysis on the remaining items. It also allows you to obtain a KWIC table or perform a keyword retrieval for a specific keyword. To perform any of these operations, simply click the desired item to display a popup menu and choose the proper command. If more than one item is close by the cursor, the menu will offer the possibility to choose on which keyword or class the operation should be performed.

You will find below a description of the various controls in the dialog box, followed by some guidelines for interpreting correspondence plots.

2-D Map control

Plot - When more than two axes have been extracted, this control allows the selection of all the possible axis combinations that can be graphed on the two axes of the plot.

2-D and 3-D Map Controls

- **Keywords** This checkbox displays or hides the row points (i.e., words or category names).
- **Groups** This checkbox displays or hides the column points (i.e., subgroup labels)
 - Clicking this button enables you to zoom in a plot. To zoom an area of the plot, hold the left mouse button and drag the mouse down/right. You'll see a rectangle around the selected area. Release the left mouse button to zoom.
 - Clicking this button restores the original viewing area of the plot.
 - Items distributed in a very similar way among subgroups of the categorical variables may be plotted on top of each other, making them hard to differentiate. Clicking down this button adds some random noise to the location of individual words or keywords, allowing you to clearly identify those that overlap. To remove the random noises, click the button again to raise it up.

In 2D correspondence plots, there is no easy way to identify where items are located on the third dimension, while in 3D it remains difficult to position an item on the depth axis without applying a rotating motion to the chart or inserting anchor lines. When this button is pressed gradients of colors are applied on fonts in order to differentiate items on the front end of the depth axis from those in the back end and those in between. Up to four colors may be used to create those gradients. To adjust the color of those gradient, use the chart customization button (see below).

Moving this slider allows the adjustment of the scaling used to plot groups or classes of the categorical variable on the correspondence plot. While the positions of keywords and groups relative to the origin of the axes are significant, it is important to remember that the scaling used for keywords and groups are independent of each other. This slider may be used as a reminder of this fact. It may also be used to increase the readability of some plots, by moving group names so they do not overlap keyword names.

- Press this button to append a copy of the correspondence plot in the Report Manager. A descriptive title will be provided automatically. To edit this title or to enter a new one, hold down the SHIFT keyboard key while clicking this button (for more information on the Report Manager, see the <u>Report Management Feature</u> topic).
- This button is used to create a copy of the chart to the clipboard. When this button is clicked, a pop-up menu appears, allowing you to choose whether the chart should be copied as a bitmap or as a metafile.
- Clicking this button brings up a dialog box to customize the appearance of the correspondence plots (click here to obtain more information on the various settings that may be changed).
- Clicking this button prints a copy of the displayed chart.

3-D Plot controls

 \square

st.l

- This button can be used to show or hide left, bottom and back walls.
- Clicking this button draws anchor lines from the floor to the data point to better locate data points in all 3 dimensions.



Locating a data point on the depth dimension of a 3-D plot can be very difficult, especially when the plot remains static. You often have to rotate this plot constantly on the various axes to get an accurate idea of where the data point is located on this third axis. Clicking this button forces WordStat to rotate the plot automatically. To disable the automatic rotation, click a second time.

Interpreting Correspondence Analysis Results

Interpretation of correspondence analysis maps can be somewhat tricky and should be made with great care, especially when examining the relationship between row points and column points. Here are some basic rules that should help you interpret such maps:

Relationship among words or categories (row points):

- The more similar the distribution of a keyword or a content category among subgroups is to the total distribution of all words within subgroups, the closer it will be to the origin. Words or categories that are plotted far from this point of origin have singular distributions.
- If two words or computed categories have similar distributions (or profiles) among subgroups of the independent variable (columns), their points in the correspondence analysis plot will be close together. For example, if the words consist of artist names and the studied subgroups represent different age groups, then if the form of the distribution of two different artists among those age groups is similar, then they will tend to appear near each other. Words with different profiles will be plotted far from each other. Please note, however that two points may appear close to each other on a two-or-three axes solution, but may, in fact, be far apart when taking into account an additional dimension.

Relationship among subgroups (column points)

- The more singular a profile of words/categories for a subgroup is, compared to the distribution of those words/categories for the entire sample, the farther this subgroup will be from the point of origin.
- If two subgroups of individuals have similar profiles of word usage or content categories, they will be plotted near each other. Subgroups with different profiles will be plotted far from each others. Again, it is important to remember that two points may appear close to each other on a two-or-three axes solution, but may, in fact, be far apart when taking into account an additional dimension.

Relationship between words/categories and subgroups (row and column points):

- Great caution should be taken when interpreting the distances between two sets of points (row and column points). The fact that the name of a subgroup is near a specific keyword or content category should not necessarily be interpreted as an indication that they are closely related.
- While the distance between words or content categories and subgroups has no interpretable meaning, the angle between such a keyword point and a subgroup point from the origin is meaningful:

An acute angle indicates that the two characteristics are correlated.

An obtuse angle, near 180 degrees, indicates that the two characteristics are negatively correlated.

In the example below, Metallica and Korn could be viewed as characteristic of males between 15 and 17 years old. However, Rob Zombie is much more specific to those listeners since it is much farther from the origin. Marilyn Manson and Black Sabbath would come next, followed by Korn and then Metallica,



• Words or categories closely associated with two subgroups will be plotted at an angle from the origin that will lie between those two groups. In the above example, the Beastie Boys seem to be characteristics of both female age groups.

For a more comprehensive description of this method, its computation and applications see <u>Greenacre (1984)</u>. For an application of correspondence analysis to the analysis of textual data see <u>Lebart, Salem, and Berry (1998)</u>.

Keyword-In-Context Page

The **Keyword-In-Context** (KWIC) technique allows you to display in a table the occurrences of either a specific keyword, or of all items related to a category, with the textual environment in which they occur. The text is aligned so that all keywords appear aligned in the middle of the table.

This technique is useful to assess the consistency (or lack of consistency) of meanings associated with a word, word pattern or category. In the example below, we can see that the word pattern KILL*, which may have been assigned to a category like "aggressiveness", refers to words with different meanings, some of them quite distant from the concept of "aggressiveness":

I have decided to **kill** few hours before... He said that he would **kill** if I call the police. Too much garlic **kills** taste of the meat. The Black Death was a disease that **killed** millions of people. When displaying rules, only the keywords or key phrases associated with the first item of the rules are displayed. For example, in a rule like:

#SATISFACTION near #TEACHER and not after #NEGATION

the KWIC list will contain only items in the SATISFACTION category meeting the conditions specified by this rule.

Once an inconsistency has been detected, it becomes possible to reduce it by making changes to the textual data or to the dictionaries. For example, the researcher may change all occurrences of the word KILL in the original text for either KILL1 or KILL2 in order to differentiate the different meanings and then add only one of these modified words (say KILL1) to the categorization dictionary. The word KILLY may also be added to the dictionary of excluded words. The categorization of phases may also be used to distinguish various meanings of a word. For example, the use of KIND to refer to the adjective ("considerate and helpful nature") may be reliably differentiate from the use of KIND as a noun ("category of things") or as an adverb by categorizing the phrase "KIND OF" as instances of this word used as a noun or as an adverb and by categorizing the remaining instances of KIND as the adjective. Disambiguation may also be performed by identifying words in close proximity that are associated with specific meanings, as well as by creating categorization rules (see <u>Working with Rules</u>).

The KWIC technique is also useful to highlight syntactical or semantic differences in word usage between individuals or subgroups of individuals. For example, candidates from two different political parties may use the word "rights" in their discourses at the same relative frequency, but we may find that these two groups use this word with quite different meanings. We may also find that the meaning of a word like "moral" evolves with the age of a child.

The **Keyword-In-Context** tab has been designed to facilitate the various tasks involved in content analysis. The tab looks like this:

🛛 🛄 Data 🛛 🔶 Text Proce	essing 📑 F	requencies	🚯 Extraction 🗞 Cooccurrences 🔚 Crosstab 💷	Keyword-In-C	ontext <pre> Classification</pre>			
List: Induded	~	So	rt by: Keyword & After 🗸 🍾 🖺 🔕 🚟				6	
Keyword: ADMINISTRATION	~	Context deli	miter: Paragraph 🗸					6
	FREQ -	CASENO		KEYWORD		CANDIDATE	DELIVERY	-
SHOW ALL ITEMS		207	ke right away. My administration will create a Veterans'	Care	Access Card to be used by veterans with illness or injury ir	MCCain	Q3-2008	
care	1175 🔳	206	make right away. My administration will create a Veterans'	Care	Access Card to be used by veterans with illness or injury ir	MCCain	Q3-2008	
E care for	115	207	est, most straightforward terms that the Veterans Health	Care	Access Card will expand existing benefits. I don't expect	MCCain	Q3-2008	
	109	207	antly or exclusively affect women. And here the Veterans	Care	Access Card will prove especially valuable, affording wom	MCCain	Q3-2008	
care of	68	206	behind in the services it provides. And here the Veterans	Care	Access Card will prove especially valuable, affording women	MCCain	Q3-2008	
 care system 	63	101	man, woman and child in America spends \$412 on health	care	administration, nearly six times as much as other countries	Clinton	Q2-2007	1
± care costs	46	185	f children from online predators, in helping to make health	care	affordable and accessible to the least fortunate among us	MCCain	Q1-2008	
± care in	36	41	; ill. She needs us to finally come together to make health	care	affordable and available for every American.	Obama	Q1-2008	
t care plan	35	49	dn't. But they believe it's finally time that we make health	care	affordable and available for every single American; that we	Obama	Q2-2008	
Care they	33	33	ca. Thank you. I'll be a President who finally makes health	care	affordable and available to every single American the same	Obama	Q1-2008	
	20	46	stiree making less than \$50,000 per year. To make health	care	affordable for all Americans, we'll cut costs and provide co	Obama	Q1-2008	
	23	54	a wealthy we just can't afford. Rather than making health	care	affordable for every American, like I've proposed, he's offe	Obama	Q2-2008	
	18	57	ver costs for ordinary Americans. We need to make health	care	affordable for every single American, and that's what I'll d	Obama	02-2008	
+ care about	16	57	an average family health care plan, and won't make health	care	affordable for the hardworking Americans who need help	Obama	02-2008	
care reform	16	42	en invested in job training and child care: in making health	care	affordable or putting college within reach.	Obama	01-2008	
E care we	15	105	amatic savings which I will use to continue to make health	care	affordable. I outlined my cost-saving measures in a speech	Clinton	03-2007	
E care the	15	101	study found this model could reduce the cost of diabetes	care	alone by 3 percent, saving us \$4 billion dollars.	Clinton	02-2007	
care crisis	13	81	v've been in a decade, at a time when the cost of health	care	and college have never been higher. It's getting harder	Obama	04-2008	
care i	12	73	sue that would make a difference in your lives - on health	care	and education and the economy - Senator McCain has t	Obama	03-2008	
care coverage	12	109	mored to work throughout my career on issues like foster	care	and adoption, family leave, equal pay and preschool for ou	Clinton	04-2007	-
	12	176	ity, my Florida neighbors looked after my family with great	care	and affection. I think I know the people of Florida pretty	MCCain	04-2007	
🗉 care it	11	94	a aid from loans to grants. We will focus on primary health	care	and affordable vaccines.	Richardson	04-2007	
	11	188	urse care managers to make sure they receive the proper	care	and avoid unnecessary treatments and emergency room s	MCCain	02-2008	
	11 🗸	44	Senate Veterans Affairs Committee – working to improve	care	and benefits for wounded warriors and their families, and	Obama	Q1-2008	- v
erade need to have a chance scade need to have a chance ome. And their children should le need to have an approach usek to overtum the Patriot Ac overmment out of people's be ave true marriage equality – (c his country needs a healing ha ave prenatal care , post-natal (to have a pat have basic rig where a presid tas being unc drooms. We w theers,) ma nd, one that u are, child care	h to citizensh hts. (Cheers, lent will take constitutional. vill make sure rriage equalit will look at th e (cheers,	The next excesses proper normal opportunity, in nove "MUM", jp. They should not be told, after they've made their contr) And we have to make sure that children have a right to a healing hand on this country and restore this country's lo (Cheers, .) We will work to cancel the Military Commissions that people have a right to the private sphere again. (Cheer y, e issue of abortion and understand that we can reconcile it) - universal health care , a living wage; help create the clim	bution to our education and ss of civil liberti Act. It's not v ers, continue.) by saying we o ate that will fo	The nice pair of a tax's and pair of a second occur pair economy, "No, we don't want you anymore; go home," bec health care as well. es. In my first week in office, I will have our Justice Departm vorthy of America. (Cheers, continue.) We will stop govern We will make sure our brothers and sisters who are gay or l can make abortions less necessary, but with a universal healt ster people to make choices that will be confirming of life.	ause America is 1 ent go to federa nent spying. We esbian, transgen n care program,	al court and will get dered, bisexu we can finally	
					11 - 1 - Cit			

The upper-right part of the screen provides a list of all instances of keywords associated with a dictionary category or of a specific word or phrase, along with its surrounding text. The panel on the left shows a tree view of items and their context in descending order of frequency, which may be used to browse through and filter through the KWIC list on the right. The text panel at the bottom of the screen displays the full document from which the selected keyword comes from and highlights the selected keyword. The text panel can be used to examine the full context of a keyword, but may also be used to add words and phrases to the current categorization dictionary or to the exclusion list. To assign a word or a phrase to a list or content category, position the text cursor on the word you want to assign, or select one or several words with the mouse or and right-click to display a contextual menu. Select the **To Categorization Dictionary** or the **To Exclusion List** menu item.

List: This option allows for specifying whether the words for display in the KWIC table either should be selected from the list of included keywords or from the list of all remaining words that have not been explicitly excluded. The option **User Specified** allows you to enter a word or word pattern at the keyboard and search for all instances of this expression.

Keyword: This option allows for choosing among all words belonging to the list of **Included** keyword or **Leftover** words (see above). When the **List** option is set to **User Specified**, this option becomes an edit box where you can type a word or word pattern. (Wildcards such as ***** and **?** are supported as well as phases).

Sort by: This option allows for sorting the keyword-in-context table in either ascending order on any of the following options:

- Case number: The KWIC table is sorted in ascending order of case position.
- Keyword & Before: The KWIC table is sorted on the keyword as well as the words appearing immediately before it.
- Keyword & After: The KWIC table is sorted on the keyword and the words appearing immediately after it.
- Keyword & Variable: When several text variables have been selected, the KWIC table includes a column indicating in which variable the keyword was found. When this option is selected, the KWIC table is sorted so that all words associated with a category or matching a word pattern are displayed in alphabetical order. Lines with identical words are sorted on the variable name from which they come. This display is useful to examine whether specific words are used with the same meaning in different variables.
- Variable & Keyword: When several textual variables have been selected, the KWIC table includes a column indicating in which variable the keyword was found. This option displays a KWIC table where all lines are sorted on the variable name from which they originate. Lines extracted from a single variable are sorted by keywords. This display is useful to establish whether different variables contain different information. For a more detailed analysis of difference in usage of specific words, use the Keyword & Variable sort order.
- Keyword & VARNAME: This option displays a KWIC table where lines are sorted by words. Lines with identical words are sorted on the value of the selected independent variable. This display is useful to highlight differences between subgroups in the meanings associated with a specific word.
- VARNAME & Keyword: This option displays a KWIC table where lines are sorted by the values of the selected independent variable. Lines with identical values on this variable are sorted by keywords. This display is useful to establish whether subgroups differ on the use of words associated with a category. For a more detailed comparison of usage of specific words, use the Keyword & VARNAME sort order.

Context delimiter: This option allows selecting the amount of context displayed in the KWIC table as well as in the concordance report. In the KWIC table, context strings, either before or after the keyword, are limited to 255 characters.

- None: This option instructs WordStat to display as much context as possible, up to a limit of 255 characters.
- **Paragraph:** When this option is selected, the program will limit the context displayed to the paragraph in which a specific keyword appears.
- Sentence: When this option is selected, the program will limit the context displayed to the sentence in which a specific keyword appears. A sentence must end with a period followed by a space or a hard return, or by an exclamation or a question mark.

• User defined: When this option is selected, the program will retrieve text found before and after the keyword until a slash character is encountered.

Once the settings have been set, click the button to start searching all instances of the selected keyword.

Using the Tree List

The upper left panel displays a tree list view where each keyword is presented in descending order of frequency along with either its context or some associated variable. The content by which keywords are sorted depends on the sort setting, so that when the **Sort By** option is set to **Keyword & Before** or **Keyword and After**, this tree list will be sorted and broken down by up to five words of its surrounding text (similar to the screen shot above). When the **Sort By** option is set to the keyword and a variable, the tree list will be sorted on keywords first, and each keyword will be broken down by the values of the selected variable.

Selecting any item in this tree list will filter the keyword-in-context list on the right to display only the selected item. Moving to SHOW ALL ITEMS will remove any filtering conditions and will show all the hits.

By default, the tree list displays all items. You can filter this list and display items meeting a minimum frequency criterion. To adjust this criterion, right-click anywhere in this panel and select **SET MINIMUM**. A submenu allows you to increase the minimum frequency criterion, to reset it to 1, or to adjust to the current value the default minimum frequency criterion to be used for all other KWIC analysis.

Any KWIC table may be saved to disk in Excel, plain ASCII, text delimited, or HTML format by clicking the 🖬 button. To export the content of the table to a new Simstat data file, press the **Export** button located at the top of the Frequencies tab.

To print the KWIC table, click the 🚔 button.

Clicking the *button* produces a concordance report on the keywords currently displayed in the KWIC table. The sort order and context delimiter of the current KWIC table are used to determine the display order and the amount of context displayed in this concordance report. This report is displayed in a text editor dialog box (see below) and may be modified, stored on disk in RTF, HTML or plain text format, printed, or cut and pasted to another application. Graphics may also be pasted anywhere in this report.

© Report			- 🗆 ×
File Edit Search			
「小 Arial ∨ 9 ∨ 🛍	i ≦ B Z U A _S A ^S S € ≣ Ξ		♡-] 函語 號 ♥
	7 8 9 10 11 12	13 14 15 16 17 17 17 17 17 17 17	18 19 20 21 21 22 23
[Case #154 CANDIDATE = Kucinich DEL	VERY = Q2-2007]		
(Cheers, continue.) We will make sure or marriage equality. This country needs a h- abortions less necessary, but with a unive care, a living wage; help create the clima imperative of our Constitution.	ur brothers and sisters who are gay or lesbi aling hand, one that will look at the issue of rsal health care program, we can finally have ite that will foster people to make choices th	ian, transgendered, bisexual have abortion and understand that we c e prenatal care, post-natal care, ch hat will be confirming of life. We r	true marriage equality (cheers,) an reconcile it by saying we can make ild care (cheers,) universal health ueed a president who understands the
[Case #77 CANDIDATE = Obama DELIVE	ERY = Q4-2008]		
That's why I believe that every single Ame industry profiteering, and that should never this.	rican has the right to affordable, accessible he be purchased with tax increases on middle	ealth care - a right that should nev class families, because that is the	ver be subject to Washington politics or last thing we need in an economy like
[Case #105 CANDIDATE = Clinton DELIV	ERY = Q3-2007]		
I believe that is America's choice, to do sor	nething about health care - America's choice f	to tackle problems of cost, quality,	and coverage.
[Case #104 CANDIDATE = Clinton DELIV	ERY = Q3-2007]		
It hurts doctors, who aren't rewarded for pro	oviding the best care - and are often punished	d for it, financially at least.	
[Case #148 CANDIDATE = Edwards DEL	VERY = Q4-2007]		
Achieving universal health care will help ex costs, and better care to the rest of us.	veryone - by covering the 47 million American	is who lack health care - and bring	ing more choices, more security, lower
[Case #101 CANDIDATE = Clinton DELIV	ERY = Q2-2007]		
Also, when you insure everyone, it will may as well as cutting administrative costs. Our	imize the impact of the prevention programs I present system is outdated, ineffective, and ur	I have recommended with earlier nsustainable.	care as opposed to emergency care
[Case #77 CANDIDATE = Obama DELIVE	ERY = Q4-2008]		
In other words, the question isn't how we ca	an afford to focus on health care - but how we	can afford not to.	
[Case #202 CANDIDATE = MCCain DEL]	VERY = Q3-2008]		

A word cloud along with a word frequency analysis may be useful to identify context words associated with the current

target item displayed in the KWIC table. To perform such an analysis, click the set button and choose whether you want to analyze context words appearing **before only**, those appearing **after only** or any words appearing **before and after** the target item. The analysis will be performed on text currently displayed in the corresponding grid column(s) and will thus take into account the context delimiter option and any filtering option currently being applied using the tree list. (See <u>Word</u> <u>Frequency Analysis</u> for more information on this feature).

The Classification Tab

Performing Automated Text Classification

The automated text classification module allows you to apply a machine-learning approach to the existing textual database in order to develop a classification model that can later be used to accurately classify uncategorized documents into predefined classes.

What is Automatic Text Categorization?

Automated text classification is a supervised machine-learning task by which new documents are classified into one or several predefined category labels based on an inductive learning process performed on a set of previously classified documents. This machine-learning approach of classification has been known to achieve comparable if not superior accuracy than classification performed by human coders, yet at a very low cost in manpower. It has been used to automatically classify documents into proper categories or to find relevant keywords describing the content and nature of a document. It has also been used to automatically file or re-route documents or messages to their appropriate destinations, to classify newspaper articles into proper sections or conference papers into relevant sessions, to filter emails or documents (like spam filtering), or to route a specific request in an organization to the appropriate department. Automated text classification may also be used to identify the author of a document of unknown or disputed authorship. For a good overview of automated text classification, see <u>Sebastiani (1999)</u>.

The automated text classification module in WordStat allows you to apply either Naive Bayes or K-Nearest Neighbors learning algorithms on an existing textual database in order to develop a classification model (or classifier). The program also provides features to test the accuracy of the classification and to optimize the various parameters. Once optimized, the obtained classification model may be used immediately to classify classification documents or may be saved on disk to be applied later outside WordStat using the <u>WordStat Document Classifier</u> utility program. The classifier may also be incorporated into a desktop or web application or within a document management system using the <u>WordStat Software Developer's Kit</u>.

The development and application of a text classifier often involve the following steps:

- 1. **Removal** of function words, words that appear in only a few documents and words that appear too often.
- 2. **Dimension reduction**, through lemmatization, stemming, categorization, word clustering or other dimension-reduction techniques.
- 3. Feature selection, which consists of a selection of terms based on their capability to discriminate between categories of documents.
- 4. Training the classifier on the train set.
- 5. Testing the accuracy of the classification on a test set.
- 6. Applying the classifier to new documents.

While the basic content analysis features of WordStat may be used to deal with the first two steps, the Automated Text Classification dialog box allows you to accomplish tasks related to the last four steps. This dialog box consists of five tabs:

- The <u>Settings</u> tab is used to select the variable to predict and the validation method to use to test the accuracy of the classifier.
- The <u>Select Features</u> tab allows you to apply various feature selection methods to select a subset of terms to be used by the classifier.
- The Learn & Test tab is the location where machine-learning algorithms are set and tested. This tab also allows the storing of classification models to disk.

- The <u>History</u> tab keeps track of every learning test performed during a session allowing you to choose the best setting and algorithm for a specific classification task. It also gives access to a batch experiment dialog box that may be used to define numerous tests and perform them all at once.
- The <u>Apply</u> tab is used to apply a classifier to an external document, a list of documents or to the current data file.

Accessing the Automated Text Classification feature

To develop a classifier for a specific categorical variable, you need to select a categorical variable containing the values you want to predict. It can be done in WordStat and from QDA Miner or SimStat while calling WordStat. In WordStat, when you select the **Analyze** button on the **Data** tab, on the left side of the dialog box are listed the categorical, numeric or date variables that you can analyze in relation to the text variables. Choose the variables and then move to the **Classification** tab. In QDA Miner, you have to choose this categorical variable in the **In Relation With** section. In SimStat, this variable should be assigned to the **Independent** list box, and the text variables on which the prediction should be based should be assigned to the **Dependent** list box. Once in WordStat, set the various text-processing options (such as the lemmatization, the exclusion and categorization lists, and all the other analysis options needed) to obtain the desired list of keywords or content categories. Then move to the **Classification** tab.

Settings

The Settings tab allows you to select which variable contains the values to predict and choose a validation method.

Selecting the Variable to Predict

S WordStat 9.0.7 - Election 2008 Coded.ppj					-		×
🚍 🛅 Data 👋 Text Processing 📑 Frequencies	🐴 Extraction	🗞 Cooccurrences	Crosstab	Keyword-In-Contex	t < Classification		
Settings Select Features Learn & Test History	Apply						
Classification options:							
Variable to predict: CANDIDATE							
Validation method: Leave one out	RUN						
		Shown: 3	00 Types: 15,81	3 Tokens: 590,701 Tir	ne: 1.9s (4 cases excl	uded)	

To select the variable containing the values to predict, set the first list box to the name of this variable. If its name is not listed, choose the **<Select Variables>** item to display a list of all available variables, select the variable to predict and click **OK**. Then set the drop-down list to this newly added variable.

Selecting the Validation Method

The evaluation of a classifier consists of measuring its effectiveness at classifying documents that have already been classified. Those documents should, however, not be part of the training set used to develop the classification model, since it would likely overestimate the real performance of the classifier. Yet, training a classifier on only a portion of the available training set may result in a less than optimal classifier. Cross-validation methods have been proposed as a compromise solution that allows you to develop a classification model on all the available documents in the training set yet provides a somewhat more realistic estimate of the classifier performance. WordStat offers three broad types of validation methods:

Leave-one-out - This cross-validation method consists of 1) removing a document from the training set, 2) developing a classification model on the remaining documents, 3) applying this model to predict the membership of this single document and 4) comparing the decision made by the classifier to the actual class to which this document belongs. This procedure is then repeated for each document in the training set and the different decisions are combined to estimate the performance of the classifier. While this method logically involves the computation of a large number of models and may seem to be time consuming, in practice the classification model is computed only once but adjusted analytically to remove the contribution of the test document prior to its classification. This cross-validation method will often overestimate the performance of a classifier if the training set includes duplicate documents or if included documents are not totally independent from one another.

n-folds - This method consists of splitting the training set into smaller partitions and testing each partition on the classification performance obtained by a model developed on the remaining ones. For example, when using a five-fold cross-validation method, the training set is divided randomly into five subsets, each containing approximately 20% of the documents. For each subset, the program tests the accuracy obtained by a classification model developed on the remaining 80% of the original training set. The performances obtained on all five classifiers are then used to estimate the performance of the classifier computed on the full training set. WordStat provides a choice between five-fold, 10-fold and 20-fold cross-validation.

External file - A more conventional method for assessing the performance of a classifier is to test the accuracy of the classifier on an entirely different set of documents that have also been classified but are totally independent of the training set on which the categorization model is based. To perform such a test, WordStat requires the test set to be stored in a different data file. When this option is selected, an Open File dialog box is displayed allowing you to identify the file containing the external set. WordStat then displays a dialog box like the one below allowing you to choose the text variable containing the documents to be used for classification and the numerical variable containing the class to which this document belongs. Once set, click **OK** to return to the classification tab.

Test	database					-		×
Date file:	C: Users Amanda Documer	nts\My Provalis I	Research Proj	jects\San	nples\Ca	andidates -	coded.pp	ni 🔗
Text:	SPEECH	 Class: 	CANDIDATE		Ŷ			
			1	ОК	×	Cancel	?	Help

• Once the variable and the validation method have been set, click the button to continue. WordStat will compute all statistics needed, and will then automatically move to the <u>Select Features</u> tab.

Select Features

The **Select Features** tab allows you to view the strength of the relationship between words, keywords or content categories and classes of the selected categorical variable. It also allows you to select a subset of items based on the statistics.

Data Text Pi	ocessing	Frequent	cies 🧌 Extra	action 8	5 Cooccurrer	ices 🛅 Crossta	b Keyword-In-Context	Classification		
ettings Select Features	Learn & T	est Hist	tory Apply							
Select Filter tabl	e: call dasses >		Compute st	atistics on: (Occurrence	v			1	
election: 1000 of 1000 featur	es				Jeconence				<u></u> ₹↓ [
Name	Global Chi ²	Р	Max Chi2	Р	Biserial	Predict				
MASSACHUSETTS	120.33	0.0000	115.95	0.0000	23.4396	Romney				- 1
MEXICO	88.34	0.0000	82.95	0.0000	19.2704	Richardson				
MAYOR	75.15	0.0000	63.84	0.0000	17.9483	Giuliani				
RONALD	82.22	0.0000	62.17	0.0000	17.1629	Romney				
ILLINOIS	58.38	0.0000	57.46	0.0000	12.5487	Obama				
MARRIAGE	59,48	0.0000	48.92	0.0000	16.2763	Romney				
HEALTHCARE	60.53	0.0000	48.92	0.0000	16.2763	Romney				
CHICAGO	50.97	0.0000	48.39	0.0000	11.5166	Obama				
SECTARIAN	48.85	0.0000	45.89	0.0000	17.1747	Biden				
PAGE	49.48	0.0000	45.77	0.0000	11.8721	Obama				
STREETS	49.73	0.0000	45.73	0.0000	11.1956	Obama				
VIOLENT	56.37	0.0000	45.48	0.0000	15.1500	Romney				
MCCAIN	53.72	0.0000	43.40	0.0000	9.9615	Obama				
INVISIBLE	42.97	0.0000	41.50	0.0000	13.2828	Clinton				
REAGAN	61.60	0.0000	41.24	0.0000	13.9787	Romney				
BIDEN	48.99	0.0000	40.88	0.0000	17.1070	Biden				
INTEND	49.13	0.0000	40.68	0.0000	10.9154	MCCain				
RESPONSIBILITIES	42.74	0.0000	39.30	0.0000	11.0006	MCCain				
			Pe	ercentag	e of case	s with REAGA	AN		_	
Romney									78.6	%
Thompson						44.4%				
Giuliani					33.3%					
MCCain*			20.4%							
Obama	11	.8%								
Richardson	7.1%									
Clinton 2.6%										
Biden" 0.0%										
Edwards 0.0%										

The strength of the relationship between an item and the classes of the categorical variable can be computed either on the occurrence (present or absent), on the frequency of items in each class, or on the percentage of words. To change the base statistic used for assessing differences among classes, set the **Compute statistics on** list to the proper option.

The discriminative strength of each item is assessed using three statistics and is presented in a table containing the following information:

Name	The word, keyword or content category.
Global Chi ²	The overall chi-square value computed on all classes of the categorical variable.
Р	The probability of the above chi-square value.
Max Chi ²	The chi-square value computed on the class with the highest case occurrence or frequency against all the other classes.
Р	The probability of the Max Chi ² value.
Biserial	The biserial correlation computed between the class of the categorical variables with the highest case occurrences and the remaining classes. This coefficient assumes that the presence or absence of a class is determined by a trait normally distributed. Contrary to the standard correlation coefficient, this measure of association may yield a value lower than -1.0 or higher than +1.0.

Predict Indicates the class in which the item most frequently occurs. When the highest case occurrence **or frequency** appears for more than one class, the column includes the labels of all those classes.

Clicking any column header sorts the table in ascending order of the data in that column. Clicking the same column header a second time sorts its content in descending order. The check boxes in the first column show, by default, that all items are to be included in the classification model. To manually remove an item, simply click in the box to remove the check mark.

The **Filter table** option allows you to display only the terms characteristic of a specific class. It may also be set to **<selected>** to display only items that have been selected for inclusion in the categorization model. To display all items set this option to **<all classes>**.

The lower portion of the tab displays a bar chart with the percentage of cases in each class of the categorical variable containing the selected item. Classes are presented in descending order of case occurrence. This graph is synchronized with the above table so that changing the selected item in this table results in the display of the corresponding distribution chart.

To display the chart on the right side of the table, click the 🔟 button. To bring the bar chart back to the bottom of the table, click the 🗮 button.

By default, bars in the chart are displayed in descending order of frequency. Moving from one feature to another will cause the order of the values to be rearranged according to the observed frequencies. To display bars associated with specific values at a fixed location, click the $2 \downarrow$ button.

Performing Automatic Feature Selection

It has been shown that reducing the number of terms in classification models can reduce not only the cost of recognition by reducing the number of features that need to be collected, but also, sometimes, can improve the classification accuracy and can reduce the risk of overfitting. Several methods have been proposed to select a subset that yields the best performance for a specific classification system. WordStat allows you to apply a filter on available items in order to select a specified number based on the computed association with the classes to predict.

• To access the feature selection dialog box, click the <u>Select.</u> button. The following dialog box will appear:



To include all items, set the selection criterion to **All** and click **OK**. To remove all the check marks beside the items, set the selection criterion to **None** and click **OK**. When the **Select** criterion is chosen, specify how many items should be selected and which statistic will be used to select them. Three statistics are available for performing such a selection: The global chi-square value, the max chi-square and the bi-serial correlation (see above for a description of these statistics). By default, the selection is performed by extracting items with the highest values on the selected statistic. Enabling the **Optimization** option instructs the program to apply a special algorithm that will take into account this statistic as well as current contrasts between classes. This option has been found, in most situations, to improve the performance of the classifier for an equal number of selected features.

To close this dialog box and apply the selection criteria, click **OK**. You may also leave this dialog box without affecting the current selection by clicking the **Cancel** button.

Once features have been selected, move to the Learn & Test tab to select a learning algorithm and test it.

Learn & Test

The Learn & Test tab is where machine-learning algorithms are chosen and tested and from which the classifier may be exported to disk.

Data	Text Prod	essing 🛒	Frequencies	👏 Extrac	tion 🗞 Coo	ccurrences	Crosstab	Keyword-In-Context	dassification	(?
Settings Se	lect Features	Learn & Tes	t History	Apply						
Learning:									ut	A
Method: Naive	e Bayes	~	<u>U</u> se:	Keyword free	quency	~				-
🗹 Uni	iform Prior	Feat	ture weighting:	Max-chi2		~	Run 🙀 P	ublish		
onfusion matrix	Confusion List	Review Error	s							
Correct = Incorrect =	= 135 = 13	Average precis Average re	sion = 0.914 ecall = 0.881	4 Non 7	ninal Accuracy = F1 =	0.9122 0.8977	Ordinal Acc	curacy = 0.9574		
Frequency Row Pct Col Pct Tot Pct	Biden	Clinton	Obama	ed Edwards	Richardson	TOTAL	PRECISION RECALL			
Biden	10 76.92 100.00 6.76	0 0.00 0.00 0.00	3 23.08 4.23 2.03	0 0.00 0.00 0.00	0 0.00 0.00 0.00	13 8.78	1.0000 0.7692			
Clinton	0 0.00 0.00 0.00	35 89.74 92.11 23.65	1 2.56 1.41 0.68	1 2.56 7.69 0.68	2 5.13 12.50 1.35	39 26.35	0.9211 0.8974			
Obama	0 0.00 0.00 0.00	2 2.94 5.26 1.35	65 95.59 91.55 43.92	0 0.00 0.00 0.00	1 1.47 6.25 0.68	68 45.95	0.9155 0.9559			
Edwards	0 0.00 0.00 0.00	0 0.00 0.00 0.00	2 14.29 2.82 1.35	12 85.71 92.31 8.11	0 0.00 0.00 0.00	14 9.46	0.9231 0.8571			
Richardson	0 0.00 0.00 0.00	1 7.14 2.63 0.68	0 0.00 0.00 0.00	0 0.00 0.00 0.00	13 92.86 81.25 8.78	14 9.46	0.8125 0.9286			
	10	38	71	13	16	148	0.9144			

The panel at the top of the tab allows you to select which machine-learning algorithm to use, set various analysis options and choose how the classification model will be tested.

Methods: Two machine-learning algorithms are available for text classification:

The **Naive Bayes** algorithm classifies text by estimating the probability of a class, given the presence or absence of specific words or keywords in the document to be classified. It first computes the probability of each term to occur in documents of specific classes in the training set. It then combines the probabilities associated with words found in the document to classify to estimate the probability that this document belongs to different classes. Finally, it assigns the document to the class with the highest probability. A multinomial Naive Bayes model has been chosen to handle both binomial and multinomial classification tasks as well as binary and numerical weights of items.

The **k-Nearest neighbor** classification method compares a document to be classified to all documents in the training set, retrieves the k most similar documents, and then assigns the new document to the most common classes in this retrieved set. This method is usually known to provide accurate classification when the training set is large enough, yet it can be very time-consuming because of the need to compare and rank the entire training set for similarity with the test document. It also usually requires a larger storage space since it must keep frequency information for all documents in the training set rather than just a few classification rules or mathematical formulas like many other machine-learning methods. However, WordStat uses a very efficient K-NN algorithm that drastically improves the computing speed and reduces

the disk space and memory requirement. When this method is chosen, a **NO** edit box appears below the Method list box, allowing you to set the number of similar documents on which this classification will be based. Values higher than 20 or 30 are typically used in text classification tasks.

Use: This option is used to select the item statistic to be used in training and classification. Choosing Case Occurrence results in the use of binary weights, indicating whether or not a word or keyword occurs in the document. Selecting Keyword Frequency allows you to use additional information related to how often this item occurs in each document. Percentage of Words and Percentage of Keywords provide two methods to normalize the obtained frequency to take into account the document length. Such normalization is performed by dividing the frequency either by the total number of words found in the document or the total number of keywords that has been extracted by WordStat.

Feature weighting: Feature weighting has been presented as an alternative to feature selection or as a way to further improve classification accuracy from selected item sets. This method consists of giving more weight to items that are rather good at differentiating documents from distinct classes and negligible weight to those that are distributed evenly among classes. The most frequently used weight in information retrieval is the TF*IDF measure where the frequency of an item is adjusted to take into account the number of documents containing this item. However, such a weighting can be considered to be only a crude approximation of the capacity of the item to differentiate documents from distinct classes. More accurate performance of the classifier can be expected from using a weight based on a more direct indicator of this discriminative capability such as the **Global Chi-Square** or the **Max Chi²** described previously.

• Once the analysis options have been set, click the button to perform the training and test the performance of the obtained classifier.

Results

A common way of assessing the accuracy of a classifier is by comparing the accuracy of predicted class membership against actual membership. Such information is provided by the **Confusion Matrix** where each predicted class is plotted against the actual class. Accurate predictions are plotted in the diagonal going from the top left to the bottom right of the table. Values in this diagonal are printed in bold characters for easy identification. Values in cells below or above this diagonal represent classification errors. Besides the actual number of documents in each cell, the table shows the row, column and total percentages. Row percentages represent the number of documents in a class that have been classified in a specific way, while column percentages express the percentage of a specific prediction actually belonging to a known class. This table may be used to identify which classes are the easiest or hardest to predict, as well as which classification errors are the most common. To facilitate comparisons across the classes of the categorical variable, two related statistics are printed on the right of the table: **Precision** is the probability that documents identified as belonging to a class are correctly classified and **Recall** is the probability of documents in a class to be correctly identified.

Several statistics are provided to assess the global performance of the classifier. The **Nominal Accuracy** measure is the proportion of documents correctly classified. It is considered a micro-average statistic since it gives equal weight to documents regardless of how they are distributed among classes of the categorical variable. The **Average Precision** and **Average Recall** measures are macro-average statistics obtained by computing the mean precision and recall obtained for every class. The **Ordinal Accuracy** measure weights disagreements so that errors in prediction will be considered higher when the predicted value is far from the original value, while predictions that are closer to the original value will be counted as partial disagreements.

The **Confusion List** tab presents information already found in the confusion matrix but in the form of a single list that allows you to identify more easily the most common errors. The table may be sorted on the actual class of the documents, the predicted classification, the number of times such a classification error occurred, or the proportion of documents that have been misclassified this specific way. By default, the table is sorted in descending order of frequency. To sort the table on values in another column, simply click this column header. Clicking the same column header a second time sorts its content in descending order.

The **Review Errors** window displays a list of all documents that have been misclassified, allowing you to examine for each document the classification error made by the classifier as well as the computed values associated with every class of the categorical variable. A text window in the bottom of the list also allows you to review the text on which the classification has been made and potentially identify some of the reasons why the document had been misclassified.

Exporting a Classifier to Disk

Once developed and optimized, a document classification model may be saved to disk and later be used outside the WordStat main program to categorize new documents. The saved categorization model includes all word exclusion, extraction and categorization settings needed to accurately retrieve items used in the classification model as well as either the classification rules (Naive Bayes) or the keyword indexing of documents in the training set (K-nearest neighbors) needed to perform document classification.

To save the classifier on disk:

- Click the working button. A standard Save File dialog box will appear.
- Enter the file name under which you would like to save the classifier and then click the **Save** button. WordStat will automatically provide a .wclas file extension.

Document classifiers may be retrieved and applied to new documents using the <u>WordStat Document Classifier</u> utility program. A special <u>Software Developer's Kit</u> (SDK) is also available upon request from Provalis Research allowing any programmer to integrate WordStat categorization and classification technologies into one's own database or document management system.

History

There is, unfortunately, no rule or rationale to help you chose a classifier and its options and parameters, and to indicate how many or which items should be selected for inclusion in the development of a classification model. For this reason, the selection of classifiers and their optimization must rely on experimentation. In other words, in order to obtain the best classifier, you needs to perform numerous tests involving systematic variations of the various settings and compare the results. The **History** tab offers a way to keep track of previously performed trials, the options used, as well as the obtained performances. From this tab, you can also access an **Experiment** dialog box that provides a convenient way to define and run a large number of classification experiments on the same data set.

E Da	ata 🦷 Text Pro	cessing 📑 Fre	equencies 🔋 E	Extraction	Cooccurrences	Crosstab	E Keyword	d-In-Cor	ntext	Classification		0
Settings	Select Features	Learn & Test	History App	ly								
Contraction (Contraction)	eriment 前 Cl	ass: <all dasses=""></all>	~							1	1 6	0
able Gra	iph	· · · ·				a. t. 1.		n 1				
Method	Parameters	No. of	Selection	validation	Nominal Accuracy	Ordinal Accuracy	Precisio	Recal	FI			2
Naive	Occur x max	2592	Chi-O opt	Leave one	0.8986	0.9568	0.91/9	0.840	0.877			
Naive	% Keywords x chi	2592	Chi2-O opt	Leave one	0.8919	0.9583	0.8850	0.818	0.850			
Naive	Occur x chi	2592	Chi2-O opt	Leave one	0.8649	0.9377	0.8999	0.771	0.830			
Naive	% Reywords x ldr	2592	Chi2-O opt	Leave one	0.7452	0.8827	0.0049	0.024	0.055			
Naive	Occur	2592	Chi2 O opt	Leave one	0.6824	0.8364	0.4762	0.404	0.439			
Naive	Occur x Iui	2392	Chi2 O opt	Leave one	0.0757	0.0560	0.4705	0.940	0.977			
Naive	Pé Kayworda y chi	2300	Chi2-O opt	Leave one	0.0900	0.9568	0.91/9	0.040	0.077			
Naive	Occurs v chi	2300	Chi2-O opt	Leave one	0.8784	0.9324	0.0004	0.012	0.000			
Naive	% Keywords x idf	2300	Chi2-O opt	Leave one	0.8043	0.9377	0.3333	0.771	0.850			
Naive	Occur	2300	Chi2-O opt	Leave one	0.6245	0.9412	0.4900	0.421	0.440	1.1		
Naive	Occur x idf	2300	Chi2-O opt	Leave one	0.6824	0.8417	0.4776	0.416	0.445			
Naive	Occur x max	2000	Chi2-O opt	Leave one	0.9054	0.9589	0.9220	0.954	0.997			
Naive	% Keywords y chi	2000	Chi2-O opt	Leave one	0.8784	0.9524	0.8564	0.812	0.833			
Naive	Occur x chi	2000	Chi2-O opt	Leave one	0.8784	0.9489	0.9112	0.800	0.852			
Naive	% Keywords x idf	2000	Chi2-O opt	Leave one	0.8716	0.9416	0.8221	0.844	0.833			
Naive	Occur	2000	Chi2-O opt	Leave one	0.7027	0.8490	0,4855	0.440	0.461			
Naive	Occur x idf	2000	Chi ² -O opt	Leave one	0.6892	0.8447	0,4768	0.432	0.453			
Naive	Occur x max	1700	Chi ² -O opt	Leave one	0.9122	0.9626	0.9288	0.869	0.898			
Naive	Occur x chi	1700	Chi2-O opt	Leave one	0.8986	0.9586	0.9255	0.842	0.882			
Naive	% Keywords x idf	1700	Chi2-O opt	Leave one	0.8784	0.9455	0.8306	0.870	0.850			

Data from prior trials are presented either in the form of a table (**Table** tab) or as a line chart (**Graph** tab). Clicking any column header of the table sorts it in ascending order of the data in this column. Clicking the same column header a second time sorts its content in descending order. By default, the displayed statistics are computed for all classes of the categorical variable. To restrict the display of either the table or the line chart to statistics related to a single class, set the **Class** list box to the desired class. Setting this option to **<all classes>** brings back the micro- and macro-average statistics computed on all classes.

Data from specific trials may be deleted from the table by selecting their rows and pressing the method.

The **Graph** tab allows you to compare the performance of various settings and the relationship between those settings using a line chart like the one shown below.



By default, the chart displays the relationship between accuracy and the number of features in the classification model. You may, however, examine the relationship between other parameters (such as precision versus recall) by changing the information plotted either on the horizontal or the vertical axis of the chart.

The following table provides a short description of buttons available:

Control Description

- Press this button to save either the table or the chart on disk. The table may be saved to disk in Excel, plain ASCII, text delimited, or HTML Charts may be saved in BMP, JPG or PNG graphic file format or may be stored on disk in a proprietary format (.WSX file extension) that may later be edited and customized using the Chart Editor.
- Pressing this button prints a copy of the displayed table or chart.
- This button allows the editing of various features of the chart such as the left and bottom axis, the chart and axis titles, the location of the legend, etc.
- This button is used to create a copy of the chart to the clipboard. When this button is clicked, a popup menu appears allowing you to select whether the chart should be copied as a bitmap or as a metafile.

The History tab also gives access to the Experiment feature where you can quickly perform a series of classification experiments.

• To access this feature click the ^{Seg Experiment} button. For more information on this feature, see <u>Classification</u> <u>Experiment Dialog Box</u>

Classification Experiment Dialog Box

The Classification Experiment feature allows you to quickly perform a series of classification experiments in order to choose among classification methods, analysis settings and selected sets of items (or features).

To access this dialog box:

• Click the ^{Seperiment} button. A dialog box similar to the one below will appear:

	lent					- 0	×
Feature sel	ection						-
Selec	t: 50 100 150 200 300						
Statistic	: Max chi-square 🗸 🗸	Based on: Pct of word	ds 🗸 🗌 Optimizat	ion			6
Learning: Method:	Naive Bayes 🗸 🗸	Use: Percentag	e of words 🛛 🗸	Testing: Method: Le	ave one out	~	
	Feature	weighthing: x IDF	Ŷ				
			E×E ♦ 🛍 🖡	Run			
Method	Parameters	No. of Features	문x문 🛉 🏦 🛛	Run Testing method	Executed		
Method Naive Bayes	Parameters Case occurrence x IDF	No. of Features	Selection method Max Chi ² -O (opt)	Run Testing method Leave one out	Executed Yes		,
Method Naive Bayes Naive Bayes	Parameters Case occurrence x IDF Keyword frequency	No. of Features 50 100 150 200 300 50 100 150 200 300	Selection method Max Chi ² -O (opt) Max Chi ² -O (opt)	Run Testing method Leave one out Leave one out	Executed Yes Yes		,
Method Naive Bayes Naive Bayes Naive Bayes	Parameters Case occurrence x IDF Keyword frequency Keyword frequency x IDF	No. of Features 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300	Selection method Max Chi ² -O (opt) Max Chi ² -O (opt) Max Chi ² -O (opt)	Run Testing method Leave one out Leave one out Leave one out	Executed Yes Yes Yes		,
Method Naive Bayes Naive Bayes Naive Bayes Naive Bayes	Parameters Case occurrence x IDF Keyword frequency Keyword frequency x IDF Percentage of words	No. of Features 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300	Selection method Max Chi ² -O (opt) Max Chi ² -O (opt) Max Chi ² -O (opt) Max Chi ² -O (opt) Max Chi ² -O (opt)	Run Testing method Leave one out Leave one out Leave one out Leave one out	Executed Yes Yes Yes Yes		^
Method Naive Bayes Naive Bayes Naive Bayes Naive Bayes Naive Bayes	Parameters Case occurrence x IDF Keyword frequency Keyword frequency x IDF Percentage of words Percentage of words x IDF	No. of Features 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300	Selection method Max Chi ² -O (opt) Max Chi ² -O (opt)	Run Testing method Leave one out	Executed Yes Yes Yes Yes Yes		^
Method Naive Bayes Naive Bayes Naive Bayes Naive Bayes Naive Bayes Naive Bayes	Parameters Case occurrence x IDF Keyword frequency Keyword frequency x IDF Percentage of words Percentage of words x IDF Case occurrence	No. of Features 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300	Selection method Max Chi ² -O (opt) Max Chi ² -O (opt)	Run Testing method Leave one out	Executed Yes Yes Yes Yes Yes Yes		^
Method Naive Bayes Naive Bayes Naive Bayes Naive Bayes Naive Bayes Naive Bayes Naive Bayes	Parameters Case occurrence x IDF Keyword frequency Keyword frequency x IDF Percentage of words Percentage of words x IDF Case occurrence Case occurrence x IDF	No. of Features 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300	Selection method Max Chi ² -O (opt) Max Chi ² -O Max Chi ² -O	Run Testing method Leave one out	Executed Yes Yes Yes Yes Yes Yes Yes		^
Method Naive Bayes Naive Bayes Naive Bayes Naive Bayes Naive Bayes Naive Bayes Naive Bayes Naive Bayes	Parameters Case occurrence x IDF Keyword frequency Keyword frequency x IDF Percentage of words Percentage of words x IDF Case occurrence Case occurrence x IDF Keyword frequency	No. of Features 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300 50 100 150 200 300	Selection method Max Chi ² -O (opt) Max Chi ² -O Max Chi ² -O Max Chi ² -O Max Chi ² -O	Run Testing method Leave one out	Executed Yes Yes Yes Yes Yes Yes Yes		^

The basic principle of this dialog box is to create a list of classification experiments involving different settings and then to instruct WordStat to perform all those experiments one after the other. Experiments first need to be defined in the upper part of the dialog box and then moved to the table of experiments located at the bottom of the dialog box. After several experiments have been defined and added to the list, you can execute all of them at once.

The first steps involve setting the experiment options. The **Feature Selection** option located at the top of the dialog box provides automatic feature selection options similar to those available from the first tab of the document classification dialog box with only one exception: while the original dialog box allows you to select only one feature set size at a time, the current dialog box allows you to set numerous feature set sizes at once. For example, by entering the following string in the **Select** edit box: 50 100 150 200 300

Five classification experiments will be performed using the same analysis settings but on features set sizes of 50, 100, 150, 200 and 300, picked out using the chosen selection method. For more information on the **Statistics**, **Base On** and **Optimization** options, see <u>Performing Automatic Feature Selection</u>.

The **Learning** and **Testing** groups of options also provide similar settings as those available on the <u>Learn & Test tab</u>. Please refer to this section for information on the available algorithms and options.

To add a classification experiment to the list:

• Set the various classification experiment settings.

- Click the # button to move the defined experiment to the list.
- Clicking the state button displays a dialog box that allows you to quickly define numerous variations of the current classifier and to add all those at once to the list of experiments to be performed.

To remove an experiment from the list:

- In the list of previously added experiments, select the one to delete.
- Click the 🔳 button.

To run experiments in the list:

• Click the button. The program performs all the experiments in the list that had not been executed before. Once an experiment is completed, the **Executed** column for this item is set to **Yes** preventing the program from executing the same experiments twice. Results of every experiment are automatically appended to the History tab.

To exit:

• Click the 划 button to close the dialog box and return to the History tab of the classification dialog box.

Apply

The **Apply** tab allows you to use the most recently tested classifier to categorize either a single document, a list of files or documents stored in the current data file or another Simstat/QDA Miner data file. To perform similar tasks using previously saved classification models, use the <u>WordStat Document Classifier</u> utility program.

	.oaea.ppj												<u>Ц</u>	X
😑 👖 Data 🛛 Text Processir	ng 📑 Freq	luencies	6 Đ	traction	8 C	ooccurre	nces 🛅	Crosstal	Key	word-Ir	n-Contex	t		< F
Settings Select Features Lea	arn & Test	History	Apply	r										
Text to classify: List of documents	v	3 A	dd files	Meth	od: Naive	e Bayes								
:: Hetection (1-Biden 20060920.TXT :: Hetection (1-Clinton 20060613.TXT :: Hetection (1-Clinton 20060613.TXT :: Hetection (1-Clinton 20070823.TXT :: Hetection (1-Clinton 20070803.TXT :: Hetection (1-Clinton 20070803.TXT :: Hetection (1-Clinton 2007019.TXT :: Hetection (1-Giuliani 20070719.TXT :: Hetection (1-Giuliani 200719.TXT :: Hetection (1-Giuliani 200719.TXT :: Hetection (1-Giuliani 200719.TXT :: Hetection (1-Giuliani 200719.TXT)														
: VElection \1-Kucinich20070319.TXT : VElection \1-McCain20070718.TXT : VElection \1-McCain20070820.TXT : VElection \1-McCain20070820.TXT														
E: \Election\1-Kucinich20070319.TXT : \Election\1-McCain20070718.TXT : \Election\1-McCain20070820.TXT : \Election\1-Obama20080103.TXT					Classi	fy						1	6	
E:\Election\1-4\u0012classes	Classification	Biden	Romney	Clinton	Classi Obama	fy Edwards	Richardson	Kucinich	Thompson	Giuliani	MCCain	¢	6	1.6
E: VElection (1-4/ucinich20070319, TXT E: VElection (1-4/uCinich20070718, TXT E: VElection (1-4/uCinic20070820, TXT E: VElection (1-0bama 20080103, TXT Document E: VElection (1-Biden 20060316, TXT	Classification Obama	Biden 0.0948	Romney 0.0822	Clinton 0.1719	Classi Obama 0.1929	fy Edwards 0.1834	Richardson 0.0660	Kucinich 0.0056	Thompson 0.0216	Giuliani 0.0146	MCCain 0.1670	1	G	1.8
:: Yelection\1-4\udink120070319.TXT :: Yelection\1-4\udink120070319.TXT :: Yelection\1-4\udink20070718.TXT :: Yelection\1-4\udink20080103.TXT :: Yelection\1-0bama20080103.TXT 20cument :: Yelection\1-8iden20060316.TXT :: Yelection\1-8iden20060920.TXT	Classification Obama Biden	Biden 0.0948 0.9458	Romney 0.0822 0.0008	Clinton 0.1719 0.0011	Classi Obama 0.1929 0.0045	fy Edwards 0.1834 0.0118	Richardson 0.0660 0.0220	Kucinich 0.0056 0.0000	Thompson 0.0216 0.0002	Giuliani 0.0146 0.0005	MCCain 0.1670 0.0134	đ	6	
:: VElection (1-4ucinich20070319, TXT :: VElection (1-4ucinich20070718, TXT :: VElection (1-4ucinich2070820, TXT :: VElection (1-4ucinich2070820, TXT :: VElection (1-4ucinich200603103, TXT :: VElection (1-4ucinich20060316, TXT :: VElection (1-4ucinich20060920, TXT :: VELection (1-4ucinich20060920	Classification Obama Biden Clinton	Biden 0.0948 0.9458 0.0149	Romney 0.0822 0.0008 0.0024	Clinton 0.1719 0.0011 0.7668	Classi Obama 0.1929 0.0045 0.0904	fy Edwards 0.1834 0.0118 0.0168	Richardson 0.0660 0.0220 0.0662	Kucinich 0.0056 0.0000 0.0000	Thompson 0.0216 0.0002 0.0006	Giuliani 0.0146 0.0005 0.0007	MCCain 0.1670 0.0134 0.0413	1	te	
E: \Election\1-4ucinich20070319.TXT E: \Election\1-McCain20070718.TXT E: \Election\1-McCain20070820.TXT E: \Election\1-Obama20080103.TXT Document E: \Election\1-Biden20060316.TXT E: \Election\1-Biden20060920.TXT E: \Election\1-Clinton20060613.TXT E: \Election\1-Clinton20060524.TXT	Classification Obama Biden Clinton Clinton	Biden 0.0948 0.9458 0.0149 0.0000	Romney 0.0822 0.0008 0.0024 0.0000	Clinton 0.1719 0.0011 0.7668 0.9931	Classi Obama 0.1929 0.0045 0.0904 0.0000	fy Edwards 0.1834 0.0118 0.0168 0.0000	Richardson 0.0660 0.0220 0.0662 0.0046	Kucinich 0.0056 0.0000 0.0000 0.0000	Thompson 0.0216 0.0002 0.0006 0.0000	Giuliani 0.0146 0.0005 0.0007 0.0000	MCCain 0.1670 0.0134 0.0413 0.0022	1		
Election (1-4ucincich20070319, TXT Election (1-4CCain20070820, TXT Election (1-4CCain20070820, TXT Election (1-4) Cain20070820, TXT Election (1-4) Election	Classification Obama Biden Clinton Clinton Edwards	Biden 0.0948 0.9458 0.0149 0.0000 0.0000	Romney 0.0822 0.0008 0.0024 0.0000 0.0000	Clinton 0.1719 0.0011 0.7668 0.9931 0.0033	Classi Obama 0.1929 0.0045 0.0904 0.0000 0.0662	fy Edwards 0.1834 0.0118 0.0168 0.0000 0.9268	Richardson 0.0660 0.0220 0.0662 0.0046 0.0004	Kucinich 0.0056 0.0000 0.0000 0.0000 0.0000	Thompson 0.0216 0.0002 0.0006 0.0000 0.0000	Giuliani 0.0146 0.0005 0.0007 0.0000 0.0000	MCCain 0.1670 0.0134 0.0413 0.0022 0.0033	1	Ŀ	
E: \Election\1-4ucinic30070319.TXT E: \Election\1-4CCain20070718.TXT E: \Election\1-4CCain20070718.TXT E: \Election\1-4Ccain20070820.TXT E: \Election\1-0bama20080103.TXT Document E: \Election\1-8iden20060316.TXT E: \Election\1-8iden20060920.TXT E: \Election\1-8iden20060920.TXT E: \Election\1-6iden20060920.TXT E: \Election\1-6iden20070524.TXT E: \Election\1-6iden20070823.TXT	Classification Obama Biden Clinton Clinton Edwards Edwards	Biden 0.0948 0.9458 0.0149 0.0000 0.0000 0.0000	Romney 0.0822 0.0008 0.0024 0.0000 0.0000 0.0000	Clinton 0.1719 0.0011 0.7668 0.9931 0.0033 0.0169	Classi Obama 0.1929 0.0045 0.0904 0.0000 0.0662 0.1529	fy Edwards 0.1834 0.0118 0.00168 0.0000 0.9268 0.8064	Richardson 0.0660 0.0220 0.0662 0.0046 0.0004 0.0015	Kucinich 0.0056 0.0000 0.0000 0.0000 0.0000 0.0000	Thompson 0.0216 0.0002 0.0006 0.0000 0.0000 0.0000	Giuliani 0.0146 0.0005 0.0007 0.0000 0.0000 0.0000	MCCain 0.1670 0.0134 0.0413 0.0022 0.0033 0.0222	1		
E: VElection (1-4ucinich20070319, TXT E: VElection (1-4CCain20070820, TXT E: VElection (1-4CCain20070820, TXT E: VElection (1-40bama 20080103, TXT E: VElection (1-40bama 20080103, TXT E: VElection (1-40bama 20060316, TXT E: VElection (1-40bama 20060920, TXT E: VElection (1-40bama 20060920, TXT E: VElection (1-40bama 20070823, TXT E: VELection (1-40bama 20070820, TXT E: VELection (1-40bama 20070823, TXT E: VELection (1-40bama 20070820,	Classification Obama Biden Clinton Clinton Edwards Edwards Thompson	Biden 0.0948 0.9458 0.0149 0.0000 0.0000 0.0000 0.0000	Romney 0.0822 0.0008 0.0024 0.0000 0.0000 0.0000 0.0000	Clinton 0.1719 0.0011 0.7668 0.9931 0.0033 0.0169 0.0020	Classi Obama 0.1929 0.0045 0.0904 0.0000 0.0662 0.1529 0.0018	fy Edwards 0.1834 0.0118 0.0168 0.0000 0.9268 0.8064 0.0009	Richardson 0.0660 0.0220 0.0662 0.0046 0.0004 0.0015 0.0007	Kucinich 0.0056 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000	Thompson 0.0216 0.0002 0.0006 0.0000 0.0000 0.0000 0.0000 0.09762	Giuliani 0.0146 0.0005 0.0007 0.0000 0.0000 0.0000 0.0000	MCCain 0.1670 0.0134 0.0022 0.0033 0.0222 0.0176	1	le:	
Election\1-4ucinich20070339.TXT Election\1-4ucinic0070718.TXT Election\1-4ucini20070820.TXT Election\1-4ucini20070820.TXT Election\1-4ucini20070820.TXT Election\1-4ucini20070816.TXT Election\1-4ucini200708216.TXT Election\1-4ucini2007080103.TXT Election\1-4ucini20060920.TXT Election\1-6ucini200708013.TXT Election\1-6ucini20070824.TXT Election\1-Clinton20070823.TXT Election\1-6uvards20070823.TXT Election\1-6uvards20070823.TXT Election\1-6uvards20070823.TXT Election\1-6uvards20070823.TXT Election\1-filediani20070823.TXT	Classification Obama Biden Clinton Clinton Edwards Edwards Thompson Giuliani	Biden 0.0948 0.9458 0.0149 0.0000 0.0000 0.0000 0.0001 0.0155	Romney 0.0822 0.0008 0.0024 0.0000 0.0000 0.0000 0.00007 0.0007	Clinton 0.1719 0.0011 0.7668 0.9931 0.0033 0.0169 0.0020 0.0566	Classi Obama 0.1929 0.0045 0.0904 0.0000 0.0662 0.1529 0.0018 0.0224	fy Edwards 0.1834 0.0118 0.0168 0.0000 0.9268 0.8064 0.0009 0.0055	Richardson 0.0660 0.0220 0.0662 0.0046 0.0004 0.0015 0.0007 0.0164	Kucinich 0.0056 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000	Thompson 0.0216 0.0002 0.0006 0.0000 0.0000 0.0000 0.0000 0.09762 0.0475	Giuliani 0.0146 0.0005 0.0007 0.0000 0.0000 0.0000 0.0002 0.7312	MCCain 0.1670 0.0134 0.0013 0.0022 0.0033 0.0222 0.0176 0.0542	1	le.	
Election\1-4ucinich20070339.TXT Election\1-4ucinic0070718.TXT Election\1-4ucini20070218.TXT Election\1-4ucini20070218.TXT Election\1-4ucini20070218.TXT Election\1-4ucini20070218.TXT Election\1-4ucini20070216.TXT Election\1-8ucini20060316.TXT Election\1-8ucini20060920.TXT Election\1-8ucini20060920.TXT Election\1-1clinton20060931.TXT Election\1-1clinton20070823.TXT Election\1-1clinton20070824.TXT Election\1-1clinton20070824.TXT Election\1-1clinton20070823.TXT Election\1-1clinton20070824.TXT Election\1-1clinton20070824.TXT Election\1-1clinton20070824.TXT Election\1-1clinton20070824.TXT Election\1-1clinton20070824.TXT Election\1-1clinton20070824.TXT Election\1-1clinton20070824.TXT Election\1-1clinton20070820.TXT Election\1-1clinton20070820.TXT Election\1-1clinton20070820.TXT Election\1-1clinton20070820.TXT	Classification Obama Biden Clinton Clinton Edwards Edwards Thompson Giuliani Giuliani	Biden 0.0948 0.9458 0.0149 0.0000 0.0000 0.0000 0.0001 0.0155 0.1061	Romney 0.0822 0.0008 0.0024 0.0000 0.0000 0.0000 0.0007 0.0505 0.1108	Clinton 0.1719 0.0011 0.7668 0.9931 0.0033 0.0169 0.0020 0.0556 0.0559	Classi Obama 0.1929 0.0045 0.0904 0.0000 0.0662 0.1529 0.0018 0.0224 0.0643	fy Edwards 0.1834 0.0118 0.0000 0.9268 0.8064 0.0009 0.0055 0.0892	Richardson 0.0660 0.0220 0.0662 0.0046 0.0005 0.0007 0.0154 0.0031	Kucinich 0.0056 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0003 0.0027	Thompson 0.0216 0.0002 0.0000 0.0000 0.0000 0.0000 0.09762 0.0475 0.0843	Giuliani 0.0146 0.0005 0.0007 0.0000 0.0000 0.0000 0.0002 0.7312 0.2586	MCCain 0.1670 0.0134 0.0022 0.0033 0.0222 0.0176 0.0542 0.1251	C ¹	Ŀ	
E: \Election\1+4ucinich20070319.TXT E: \Election\1+4Ccini2070319.TXT E: \Election\1+4Ccin20070820.TXT E: \Election\1-4McCan20070820.TXT E: \Election\1-0bama20080103.TXT Document E: \Election\1-8iden20060316.TXT E: \Election\1-8iden20060920.TXT E: \Election\1-6iden20060920.TXT E: \Election\1-6iden20070823.TXT E: \Election\1-6iden20070823.TXT E: \Election\1-6iden20070823.TXT E: \Election\1-6iden20070820.TXT E: \Election\1-6iden20070820.TXT E: \Election\1-6iden20070820.TXT E: \Election\1-6iden20070820.TXT E: \Election\1-6iden20070820.TXT E: \Election\1-6iden20070719.TXT E: \Election\1-6iden20070719.TXT	Classification Obama Biden Clinton Clinton Edwards Edwards Thompson Giuliani Giuliani	Biden 0.0948 0.9458 0.0149 0.0000 0.0000 0.0000 0.0001 0.0155 0.1061 0.0005	Romney 0.0822 0.0008 0.0024 0.0000 0.0000 0.0000 0.00007 0.0505 0.1108 0.0001	Clinton 0.1719 0.0011 0.7668 0.9931 0.0033 0.0169 0.0020 0.0566 0.0559 0.0000	Classi Obama 0.1929 0.0045 0.0004 0.0000 0.0662 0.1529 0.0018 0.0224 0.0643 0.0000	fy Edwards 0.1834 0.0118 0.0000 0.9268 0.8064 0.0009 0.0055 0.0892 0.0015	Richardson 0.0660 0.0220 0.0662 0.0046 0.0045 0.0007 0.0154 0.0031 0.0831	Kucinich 0.0056 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0003 0.0227 0.0000	Thompson 0.0216 0.0002 0.0000 0.0000 0.0000 0.0000 0.09762 0.0475 0.0843 0.0004	Giuliani 0.0146 0.0005 0.0007 0.0000 0.0000 0.0002 0.7312 0.2586 0.9973	MCCain 0.1670 0.0134 0.0022 0.0033 0.0222 0.0176 0.0542 0.1251 0.0001	Už.	le	

The document classification feature supports numerous file formats such as plain ASCII text files as well as HTML, Rich Text, MS Word, WordPerfect, Acrobat PDF files. Detailed results of classifications are displayed in a table at the bottom of the dialog box and may be either saved to disk or printed. When applied to the current database or another database, the automatic classification feature may be useful to categorize unclassified documents or to review existing classifications based on the results of the new classifier.

To classify a single document:

- Set the Text To Classify list box to Single Document.
- Click the **Open File** button to locate and import the file containing the text to be classified. You may also type directly in the text-editing window or paste a text previously copied to the clipboard (by moving to the text-editing window and pressing **Ctrl-V**).
- Click the Classify button to apply the current classifier to the displayed text.

To classify a list of documents:

- Set the Text To Classify list box to List of Documents.
- Click the Edit List button to display a dialog box like the one below that allows you to browse through your computer and select certain documents. You may add documents located in different folders by successively adding documents located in a specific folder and then moving to a new location where other documents are located. Click OK to confirm the changes to the file list.
| Đ | Fla | c | ^ | Name | Size | Item type | Date modified | Date created | Date accessed |
|----------------------------------------------|------------------------------|------------------------------------------------------------------------------------------------------------------|-------------------|----------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------|-----------------------|---------------------|-------------------|--------------------|
| Đ | IP | | | ABC - December | 119 KB | Chrome HTML Documen | t 1/19/2013 1:10 PM | 2/17/2021 8:22 AM | 5/27/2021 10:30 A |
| Ð | Liv | e_Version_1.2 | | Bloomberg - Octo. | 158 KB | Microsoft Word 97 - 2 | . 1/19/2013 1:14 PM | 2/17/2021 8:22 AM | 5/29/2021 2:06 PM |
| Ð | Or | der | | CBS - November | 63.7 KB | Microsoft Word Docu | 1/19/2013 1:23 PM | 2/17/2021 8:22 AM | 5/29/2021 2:06 PM |
| Đ | Pro | ovalis | | CNN - June 13, 2. | 215 KB | Rich Text Format | 1/19/2013 1:16 PM | 2/17/2021 8:22 AM | 5/29/2021 12:47 P |
| E | Wo | orkshop | | CNN - November . | 202 KB | Rich Text Format | 1/19/2013 1:16 PM | 2/17/2021 8:22 AM | 5/29/2021 12:47 P |
| • | | Candidates | | CNN - October 1 | 215 KB | Rich Text Format | 1/19/2013 1:19 PM | 2/17/2021 8:22 AM | 5/29/2021 12:47 PI |
| + | | Election 2008 | | CNN - Sept 12, 2. | 181 KB | Rich Text Format | 1/19/2013 1:19 PM | 2/17/2021 8:22 AM | 5/29/2021 12:47 P |
| | | GOP | | FOX - August 11,. | 148 KB | Adobe Acrobat Docu | 1/19/2013 1:24 PM | 2/17/2021 8:22 AM | 5/29/2021 12:48 Pf |
| | Ŧ | GOP Debates 2011 | | FOX - December | 132 KB | Adobe Acrobat Docu | 1/19/2013 1:24 PM | 2/17/2021 8:22 AM | 5/29/2021 12:48 PI |
| | Ŧ | GOP Debates 2012 - Part | | FOX - May 5, 201. | 67.8 KB | Adobe Acrobat Docu | 1/19/2013 1:21 PM | 2/17/2021 8:22 AM | 5/29/2021 12:48 Pl |
| | Œ | GOP Debates 2012 - Part | | TOX - Sept 22, 2 | 127 KB | Adobe Acrobat Docu | 1/19/2013 1:22 PM | 2/17/2021 8:22 AM | 5/29/2021 12:48 PI |
| (+) | | Paintings | | MSNBC - Novemb. | 102 KB | Text Document | 1/19/2013 1:24 PM | 2/17/2021 8:22 AM | 5/29/2021 12:48 PI |
| | ocum | oads | ÷ | | | rexcolument | 192013 123 FM | 2/17/2021 0.22 AM | 5/25/2021 12:40 FI |
| rmand\C | Desk | top\Workshop\GOP\GOP Deba | tes | 2012 - Part 1\ABC - Jan | 7, 2012.HTML | P Add | Remove | | |
| rmand \C
rmand \C
rmand \C
rmand \C | Desk
Desk
Desk
Desk | top/Workshop/GOP/GOP Deba
top/Workshop/GOP/GOP Deba
top/Workshop/GOP/GOP Deba
top/Workshop/GOP/GOP Deba | tes
tes
tes | 2012 - Part 1(CNN - Jan
2012 - Part 1\CNN - Jan
2012 - Part 1\FOX - Jan
2012 - Part 1\MSNBC - J | 19, 2012.RTF
26, 2012.doc
16, 2012.doc
16, 2012.pdf
Jan 23, 2012.txt | | | | |

• Click the Classify button to apply the current classifier to all documents in the list.

To classify documents in the current data file:

- Set the Text To Classify list box to Current Data File.
- Click the Classify button to apply the current classifier to all documents contained in the selected text variables. Please note that these are the same text variables used for developing the current classifier. However, some documents may have been ignored during the classification phase, either because they had been filtered out or because they had not been classified before and contained missing values in the categorical variable. These documents, as well as all other previously classified documents, will be categorized by the current classifier and the result of this classification will be displayed in a result table.
- To store the predicted class or the computed score obtained for every class, click the button. The following dialog box will appear:

ave classification data	3
Save predicted dass	🗸 ок
Variable name: PREDICTED	Cancel
Save scores	_
Variable prefix:	

- To save the predicted class put a check mark beside Save Predicted Class and enter a variable name.
- To save the scores associated with each class and upon which the classification has been made, put a check mark beside Save Scores and enter a variable prefix (up to 7 characters). Variable names are created by adding successive numeric values to this prefix. For example, if the edit box at the right of the Variable Prefix option is set to "CLASS_", the variable names will be CLASS_1, CLASS_2, CLASS_3, etc.
- If any one of the specified variables does not exist, WordStat will create new ones and store the numerical values associated with either the predicted class or the class scores. A confirmation dialog box will ask for confirmation of the creation of these new variables as well as the overwriting of any existing ones.

To classify documents in another data file:

- Set the Text to Classify list box to External Data File.
- Click **Open File** to locate the Simstat/QDA Miner data file containing the documents to be classified. A dialog box similar to the following one will appear:

elect document or text variables	
ou would like to analyze:	
NEW	

- Select one or several text or document variables that will be used for classification purposes and click **OK**. The content of the data file is displayed in a table, while the text to be classified is displayed on its right. You can resize this text window by dragging its left border.
- Click the Classify button to apply the current classifier to all documents contained in the selected text variables.
- To store the predicted class or the computed score obtained for every class, click the button. A dialog box similar to the following will appear:

Save classification data	>
Save predicted dass	🗸 ок
Variable name: PREDICTED	X Cancel
Save scores	
Variable prefix:	

- To save the predicted class put a check mark beside Save Predicted Class and enter a variable name.
- To save the scores associated with each class upon which the classification has been made, put a check mark beside Save scores and enter a variable prefix (up to 7 characters). Variable names are created by adding successive numeric values to this prefix. For example, if the edit box at the right of the Variable Prefix option is set to "CLASS", the variable names will be CLASS1, CLASS2, CLASS3, etc.
- If any one of the specified variables does not exist, WordStat will create new ones and store the numerical values associated with either the predicted class or the class scores. A confirmation dialog box will ask you to confirm the creation of these new variables, as well as to overwrite any existing variables.

To export the table to disk:

• Click the 🖬 button. A Save File dialog box will appear.

- In the **Save as type** list box select the file format in which to save the table. The following formats are supported: ASCII file (*.TXT), Tab delimited file (*.TAB), Comma delimited file (*.CSV), HTML file (*.HTM;*.HTML), and Excel spreadsheet file (*.XLS).
- Type a valid file name with the proper file extension.
- Click the **Save** button.

Miscellaneous

Preparing and Importing Data

This section provides general information on how to prepare textual data for analysis in WordStat and specific instructions on how to import data into QDA Miner and SimStat.

Preliminary Text Preparation

While interview transcripts, responses to open-ended questions, or any other kind of textual information may be typed directly within SimStat or QDA Miner, there are many situations where electronic versions already exist either in the form of text files or as data files accessible only through specific applications such as word processor, spreadsheet or database programs. This information must be transferred into a QDA Miner project or Simstat data file for further processing. Projects can also be created directly in WordStat from documents and spreadsheets etc. Prior to using WordStat for content analysis, some modification or adjustments may need to be made.

Uppercase and lowercase letters

By default, WordStat is case-insensitive and therefore accepts files in either upper-or-lowercase.

Check spelling of documents

The automatic content analysis feature of WordStat involves numerous operations of word recognition and generally requires each word to be spelled correctly. Any misspelled word may be left uncoded and leads to imprecise or invalid conclusions. Two strategies may be used to deal with misspellings:

- One may run documents through a spell-checker to make sure all words are spelled correctly. WordStat provides spell-checking for more than 20 languages. The spell-checking may be performed in QDA Miner or through <u>Text</u> <u>Editor</u> feature of WordStat. An even more efficient approach is to use the <u>Misspellings and Unknown Words</u> feature of WordStat to quickly retrieve all potentially misspelled words and to replace them all at the same time.
- An alternative approach would be to build a content analysis process that would take into account the misspelling
 of words. To achieve this, one may use the <u>Substitution</u> feature to automatically replace misspelled words with
 their correct forms or add the most commonly misspelled keywords to the content-analysis dictionary.

Remove hyphenation

While WordStat can be configured to accept compound words with dashes, it cannot differentiate between dashes and hyphens. As a consequence, a hyphenated word will often be treated as two separate words. It is recommended to revise the text to ensure no hyphenation is present.

Add or remove square brackets ([]) and braces ({ })

Square brackets and braces have special meanings for WordStat. For example, braces are often used to remove a section of the text that you don't want to process while square brackets may be used to restrict the analysis to specific portions of text. If these symbols are used in a text for other purposes, they should be replaced with other symbols.

If there are specific parts of your text that you do not what to process, such as explanation notes, interviewer questions and probes, comments, etc.), enclose them in braces (ex. comment). Also, if you want to perform a content analysis on only a small portion of the entire text, such as on manually entered codes, enclose this portion of text in square brackets. QDA Miner's coding feature may also be used to restrict the analysis to some sections or exclude specific text segments from the content analysis process. Once the text segments have been manually tagged in QDA Miner, one could then specify, when

calling WordStat, to ignore sections tagged with specific codes or to only analyze segments associated with one or several codes.

Importing Spreadsheet or Database Files

Spreadsheet Data Files

Most spreadsheet programs allow for entry of both numeric and alphanumeric data into cells of a data grid. QDA Miner, SimStat and WordStat can import spreadsheet files produced by EXCEL (*.xls; *.xlsx).

To import an Excel spreadsheet from QDA Miner:

- Choose the NEW command from the PROJECT menu.
- Click the Import from an Data File button.
- Select the file format using the List File of Type drop down list.
- Select the file you want to import and click the OK button.

To import an Excel spreadsheet from Simstat:

- Choose the DATA | IMPORT command from the FILE menu.
- Select the file format using the List File of Type drop down list.
- Select the file you want to import and click the OK button.

The program displays a dialog box where you can specify the spreadsheet tab and the range of cells where the data are located. You must specify a valid range name or provide upper left and lower right cells, separated by two periods (such as A1..H20). If you set the Range Name list box to ALL, the program attempts to read the whole tab.

To create project from a Excel spreadsheet in WordStat:

• Select the New button. A dialog box similar to the one below will appear.



- Click the Import from an existing data files or web service button.
- Select the desired file format from the menu. An import dialog box will appear.
- Select the corresponding file format from the Files of type drop-down list and select the file you want to import.
- Click Open.

- Name your project and save it in the appropriate location.
- Once you have named your project and saved it, a dialog similar to the one below appears.

Import	¢.		-	
Sheet:	survey		*	V Import
Range:	• All	ORange		X Cancel
				Preview >>

When you select an Excel file format for importation, the program displays a dialog box in which you can specify the spreadsheet tab and the range of cells where the data are located. You must specify a valid range name or provide upper left and lower right cells, separated by two periods (such as A1..H20). If you set the **Range** radio button to **All**, the program attempts to read the whole tab. You can preview the data before you import by selecting the **Preview** button. Your Excel spreadsheet must be closed before you select **Import**. You can only import one tab at a time. If you have more than one tab to import, containing the same column headers, simply append the other tabs.

- In the Sheet drop-down choose the tab you would like to import.
- Select either a Range or All and set you range, if necessary.
- Select Import. An Import Options dialog appears similar to the one below.

			P	OW IDE	NTIFICATION	- STEP 1 OF 2		
_	_				The loan of	- STEP TOT E		
oes yo	ur workshe	et contain:						
	Variable N	Names: Start on row: 1	*					
E	1) (ariable [Descriptione: Clock on your	-					
-	1 valiable r							
	Data:	Starts on row: 2	-					
review	-							
1	ID	WHERE	GENDER	AGEGROUP	ETHNICITY	JOB	SKILLS	
2	6	Los Angeles, CA	Male	25-39	caucasian	entertainment industry	None that I am aware of. LOL	- 1
3	8	Scotland	Male	19-24	Scottish	Student	Leading huge PvP groups and	
4	10		Male	25-39	black	auto worker	n/a	
5	11	USA	Male	14-18	American	Student	It helps me to think things	
6	12	Moose Jaw, Canada	Male	25-39	Irish	Accountant	Jokes and COH humor cross	
7	13	Birmingham, Alabama-USA	Male	19-24	Caucasian	College Student/ Restaurant	I\'ve learned a little more	
8	15	lakenheath, england	Male.	25-39	caucasion	retail	typing skills if any	
9	17	Bremerton, Washington USA	Male	40-54	American	Programmer	none, games do not effect	
10	18	Amherst, Massachusetts,	Male	25-39	Caucasian	Molecular biologist	At most, the game has taken	
11	19	Arkansas	Male	25-39	Black	Govt	Humor. Im a real serious	
12	24	Sacramento, CA.	Male	25-39	Caucasion	State Worker	I really don\'t think much I\'ve	
	Inc	11 11 11 10 1	59.2	DE DO	1.9		at	

This dialog is a two-step process that helps to properly configure your data for import. On the first page you are asked to identify the location of the variable names and variable descriptions, if present. You are also asked to identify the row in which the data starts.

• If your spreadsheet contains variable names, check the checkbox called **Variable Names** and enter the row number in which they are located.

- If your spreadsheet contains variable descriptions, check the checkbox called Variable Descriptions and enter the row number in which they are located.
- In the Data field enter the row number in which the data starts.
- Select Next. The second page appears.

lect the va	ariables that you would li	ike to import. You can customize the n	ame and description if necessary		
typing in t	he associated cell. Char	nge the variable type by selecting from	the dropdown menu.		
ables					
elected	Name	Description	Туре		
	ID	ID	Integer	~	
\square	WHERE	WHERE	Short String	*	
	GENDER	GENDER	Nominal	~	
\square	AGEGROUP	AGEGROUP	Nominal	~	
\checkmark	ETHNICITY	ETHNICITY	Document	*	
\square	JOB	JOB	Short String	*	
	SKILLS	SKILLS	Document	2	

The second page allows you to choose the variables you would like to import. You can change the names and descriptions by typing directly in the cells. To change the data type select from the drop-down list.

- · Select the variables you want to import by checking their checkboxes.
- Modify the variable names, descriptions and data types as necessary.
- Select **Import**. The data will be imported and WordStat's **Data** tab will appear containing a table with your imported Excel data. From here you can further tailor your data set if necessary, or start analyzing immediately.

Formatting Spreadsheet Data

The selected range must be formatted such that the columns of the spreadsheet represent variables (or fields) while the rows represent cases. Also, the first row should preferably contain the variable names while the remaining rows hold the data, one case per row. QDA Miner and Simstat will automatically determine the most appropriate format based on the data it finds in the worksheet columns. Cells in the first row of the selected range are treated as variable names. If no variable name is encountered, QDA Miner as well as SimStat will automatically provide one for each column in the defined range.

When reading the data for analysis, all blank cells and all cells that do not correspond to the variable type (e.g., alphanumeric entries under a numeric variable, or a numeric value under a string variable) are treated as missing values.

Database Files

MS Access, dBase and Paradox files

QDA Miner and Simstat can directly import MS Access, dBase and Paradox data files. For the last two file formats, the length of alphanumeric variables should not exceed 256 characters and memo variables are not supported. If you do have this kind of data, you may use the exporting capabilities of your database program to create a date file more compatible with Simstat (such as a Visual FoxPro data file or a tab delimited text file).

Other database files

QDA Miner and SimStat offer ways to import from various database formats by connecting directly to the database system using an ODBC connection. For database systems for which no ODBC driver exists, it is still often possible to import data by exporting the data to a common file format that can be read by QDA Miner or SimStat. The recommended file formats are, in descending order of preference, Excel, Tab-delimited text files, or CSV files.

Importing memo variables

Memo variables that have not been successfully imported may be transferred to the data file either by using cut and paste operations or by retrieving text files from disk. For more information on this topic, see <u>Importing plain text and word</u> processor files.

Importing Plain Text or Word Processor Files from Simstat

QDA Miner provides an easy way to import documents stored in various formats including MS Word, WordPerfect, Rich Text, HTML, PDF and plain text files. When using SimStat as the base module, such a task is not as obvious. One way to transfer data from a word processor document into Simstat is to open simultaneously both applications and use cut and paste operations to transfer data through the clipboard. However, this may not be the most efficient way, especially when one needs to import a large amount of information. The following section presents four additional methods to transfer text information into memo variables:

- Using the Document Conversion Wizard program
- Retrieving a text file into a memo variable
- Importing comma or tab delimited text file
- Importing page delimited memo files

While the first method can read textual data stored in word processor documents, the last three methods require the data to be stored on disk in plain ASCII files without any formatting or typesetting code. Most word processors offer an option to save a document as a plain text file. If you don't know how to create such a text file, please refer to your word processor manual.

Using the Document Conversion Wizard program

WordStat includes a conversion utility program that can assist you in the importation of text files stored in either word processor documents such as MS Word, MS Write, WordPerfect, RTF or Acrobat PDF files, but also of text stored in ASCII (plain text), HTML or even Excel spreadsheet files. To run this program:

• Point to the Programs folder in the Windows' Start menu, then select Provalis Research and then click Document Conversion Wizard.

This utility program will guide you through the process of importing one or numerous text files.

Retrieving a Text File into a Memo Variable

This method should be used to retrieve a single unit of text into a memo variable for a specific case. If textual data for several cases need to be retrieved, they should be stored in different text files. To retrieve the text file from SimStat:

- Open the data file where the information should be stored.
- Position the cursor on the cell in which you would like to store the text. A memo editor should appear at the bottom of the data sheet.
- Click inside the memo editor or press F2.
- Click the Import Text Into Memo button 1, select the text file you wish to retrieve and click OK.

Importing Comma or Tab Delimited Text Files

If you wish to retrieve a text file containing several numeric and alphanumeric variables, you may have to transform this file into a comma or tab delimited text file. There are, however, several limitations to this transfer method. If commas are used as delimiters, then all existing commas within text variables should ideally be removed. If a tab delimited format is chosen, all tab characters already present in a text variable should be removed. Another important limitation is that all the information of a single case must be stored in a single line. For this reason, hard returns in long texts should be removed so that the entire text is stored on a single line. (There is no limitation on the total number of columns per line, so it is possible to store very long texts on a single line).

QDA Miner as well as Simstat can read up to 2000 numeric and alphanumeric variables from a plain ASCII file (text file). The file must have the following format:

- Every line must end with a carriage-return.
- The first line must include the variable names, separated by tabs or commas.
- Variable names may have a length of not more than 10 characters. Longer strings are truncated to 10 characters.
- The remaining lines must include the numeric or alphanumeric values, separated by tabs or commas.
- Each line must contain data for one case and variables must be in the same order for all cases.
- All invalid data and all blanks encountered between commas or tabs are treated as missing values. A single dot can also be use to represent a missing numeric value.
- Comments can be inserted anywhere in the file by putting a * at the beginning of the line.
- Blank lines can also be inserted anywhere in the file.
- Comma delimited text files require a .CSV extension while tab delimited files require a .TAB extension.

Importing Page Delimited Memo Files

SimStat provides a simple method to import numerous records of text by the use of page delimited memo files This file format consists of a plain text file which contains the textual data of numerous individuals for a single memo variable. The text for each record must be separated by page break characters (ASCII 12) or by the ^P string. The file name extension of this text file should be .MMO. To import such a file:

- Choose the DATA | IMPORT command from the FILE menu.
- Set the file format to Page Delimited Memo using the List File of Type drop down list.
- Select the file you want to import and click the OK button.

The resulting file consists of a SimStat data file with two variables: RECNO, a numeric variable containing a sequential number going from 1 up to the total number of cases encountered in the input file, and TEXT, a memo variable containing the textual data for this case.

Note: Importation of numerous text variables may be achieved by performing successive importations of page delimited memo files and then using the APPEND VARIABLES command to merge the resulting files into a single one. In order to achieve this, great care should be taken to give unique names to the various TEXT variables and to assure the case sequence of the various text files is identical.

Creating Comparison Charts

The chart window allows you to graphically examine the relationship between specific items and values of an independent variable. The bar chart should preferably be used to display the distribution of various categories within subgroups as defined by a nominal independent variable, while the line chart may be preferred for the examination of the relationship between these categories and an ordinal or quantitative variable. Eight types of chart can be obtained from this dialog box.

- The vertical bar chart can be used to display the distribution of various categories within subgroups as defined by a nominal independent variable. Bar charts are easy to read since one can easily compare the heights of the bars.
- The **Horizontal bar chart** has the same advantage as the vertical bar chart. The horizontal orientation provides a way to accommodate longer category labels without the need to change their text orientation.
- The **stacked bar chart** allows one to display relative or absolute frequency of items by stacking them for each other class, It allows one to quickly show the relationship of parts to the whole or to emphasize a sum of several codes
- The **100% stacked bar chart** is designed to show the relative percentage of multiple data series in stacked bars, where the total of each stacked bar always equals 100%. Like a pie chart, a 100% stacked bar chart shows a part-to-whole relationship.
- The **line chart** may be preferred for the examination of the relationships between these categories and an ordinal or quantitative variable. It is especially useful when looking for trends over time.
- The stacked area chart is similar to a line chart with the areas below the lines filled with colors. The stacked version of the area chart allows one to display the contribution of various elements to a total over time.
- The **100% stacked area** chart is similar to the stacked area chart with the difference that each value of the horizontal axis is adjusted so that the cumulative area always represents 100%. It is thus appropriate to display the evolution of the contribution of various elements to a total over time.
- The **Radar Chart** represents all categories of a variable on a different axes. The origin of each axis of the radar chart is located at the center of the chart so that the higher frequency extends the plotted area farther toward the edge of the chart. The resulting geometric shapes illustrate the overall distribution of all categories.



A quick way to retrieve documents, paragraphs or sentences associated with a specific bar or pie slice is by rightclicking and selecting the keyword retrieval command.

To search for selected text segments:

- Right-click on the bar corresponding to the item and the specific value for which you would like to retrieved associated text segments.
- Select the **search...** command. All retrieve segments will be displayed in a separate dialog box. For more information on the various features available from this dialog box, see <u>Keyword Retrieval</u>.

To remove specific series:

- Right-click on any bar corresponding either to the code or to the specific value you would like to remove.
- The popup menu will display two **REMOVE** commands, one of them for removing a specific item, while the other one for removing all bars associated with a specific value. Select the one associated with what you would like to delete.

The Toolbar

The following table provides a short description of the available buttons and controls:

Control Description

Press this button to vertically display the labels on the bottom axis.

H	Clicking this button allows you to manually reorder series and bars within a series. It may be used to place bars in ascending or descending orders of frequency, move high frequency series in a back of a 3D chart, so as to make the other lower frequency series more visible. When used with radar charts, it may allow one to produce easier to read geometric shapes.
	Press this button to change the default color palette used for various series of the chart.
<u>7</u> 0	This button allows the editing of various features of the chart such as the left and bottom axes, the chart and axes titles, the location of the legend, etc. (see <u>Chart Options Dialog Box</u>)
3 D	Click this button to turn on/off the 3D perspective for the current chart.
*}	Clicking this button causes the values represented on the bottom axis to be switched with those represented by different lines or bars (legend).
	This button creates a copy of the chart to the clipboard. When this button is clicked, a shortcut menu appears allowing you to select whether the chart should be copied as a bitmap or as a metafile.
V2	Press this button to append a copy of the graphic in the Report Manager. A descriptive title will be provided automatically. To edit this title or to enter a new one, hold down the SHIFT keyboard key while clicking this button.
	Click this button to save a chart on disk. Charts may be saved in BMP, JPG or PNG graphic file formats.
	Clicking this button prints a copy of the displayed chart.

Barchart and Line Chart Options

The various options on this dialog box allow you to customize the appearance of bar charts and line charts. The options available in this dialog box represent only a small portion of all of the settings available.

To further customize the chart, modify data points, value labels, or series order, click the ¹⁴/₄ button located on the right side of the dialog box.

Left Axis

Minimum / Maximum: WordStat automatically adjusts the vertical axis scale to fit the range of values plotted against it. To manually set these values, type the desired minimum and maximum.

Increment: Increasing or decreasing this value affects the distance between numbers as well as tick marks. Horizontal grid lines are also affected by modification of this value.

Horizontal Grid: This option turns horizontal grid lines on and off. Grid lines extend from each tick mark on an axis to the opposite side of the graph. To increase or decrease the number of grid lines or the distance between those lines, change the Increment value of the axis. A list box also allows a choice among five different line styles to draw the grid lines.

Legend

Location: This option positions the legend. Legends may be placed at Top, Left, Right and Bottom side of the chart.

From top: When the legend is displayed on the left or the right side of the chart, this option specifies the legend's top position in percent of total chart height.

From left: When the legend is displayed on the top or the bottom of the chart, this option specifies the legend's top position in percent of total chart width.

Titles

Proper titles and axis labels are of utmost importance when describing the information displayed in a chart. By default, WordStat uses variable names and labels as well as other predefined settings to provide such descriptions.

The title page allows you to modify the **top** title, as well as the labels on the **left**, **bottom** and **right** axis. To edit the title, select the proper radio button. Enter several lines of text for each title by pressing the **<Enter>** key at the end of a line before entering the next line.

The Font button on the right side of the edit box allows changing the font size or style of the related title.

3-D View

Orthogonal: Turning this option off disables the free elevation and rotation of the 3-D chart.

Zoom: This option zooms the whole chart. Expressed as a percentage, increasing the value positively will bring the chart towards the viewer, increasing the overall chart size as the Zoom value increases.

3-D Percent: The 3-D Percent property indicates the size ratio between chart dimensions and chart depth by specifying a percent number from 1-to-100.

Perspective: Use this property with Orthogonal unchecked to modify the 3-D perspective of the Chart. Larger values add more depth perspective.

Bar shadow: Enabling this option add a shadow to the sides of 3-D bars. Turning it off will colors the sides of the bar the same as the front.

Bar width: This option determines the percent of total bar width used. Setting this value to 100 makes joined bars.

Bar depth: Use this property to limit the depth that each bar series use. By default, bars will take up the part proportional to the number of bar series in the chart so that the back of a bar will join the front of the bar immediately behind it. To insert a gap between series of bars, decrease this value.

Edit Case Descriptors

The **Edit Case Descriptor** command allows you to define a case descriptor that will be used to identify each case based on the values it contains in one or several variables. Such a string is used to identify cases when retrieving segments as well as when analyzing case similarities using clustering, multidimensional scaling or proximity plots.

To edit case descriptors:

• Click the ≡ button in the upper left corner of the main window, scroll down to EDIT CASE DESCRIPTORS. A dialog box similar to this one will appear:

riables: ASENUM			
OPIC			
		-	
200000000000000000000000000000000000000	CANDIDATE	1 montal	

This dialog box allows you to specify a label that will be used to describe each case. The label may be changed by editing the text in the **Description String** edit box. To insert the value stored in a specific variable into the description, simply enter the variable name in uppercase letters and enclose this name between braces. Alternatively, you can insert a variable name at the current cursor location by clicking the corresponding item in the **Variables** list located just above the edit box.

If you enter the following string:

{GENDER} subject - {AGE} years old

The {GENDER} and {AGE} strings will be replaced with their corresponding value for this specific case. If the current case contains information about a seventeen-year-old male, the above string will be displayed as:

Male subject - 17 years old

It is also possible to insert the following string:

{CASENUM}

This string will display a unique case number, representing the physical order of this case in the project file.

Filtering Cases

DATA TYPE

The **Filter Cases** command temporarily selects cases according to some logical conditions. You can use this command to restrict analyses to a subsample of cases or to temporarily exclude some group of subjects. Two types of filtering are available in WordStat: a basic filtering dialog box that provides some guidance for easily building filters, and a powerful xBase filtering tool that allows you to create more complex filtering expressions that contain more advanced mathematical functions, as well as Boolean and relational operators.

To filter cases:

• Click the ≡ button in the upper left corner of the main window, scroll down to **FILTER CASES**. A dialog box similar to this one will appear:

	Variable:		Operator:	_	Criteria:	
	FOLLOWERS	~	is lower than	~	500	~
AND	/	*		~		~
AND				-		
AND				-		14

AVAILABLE OPERATORS

This dialog box consists of two major sections. The upper section of the filtering dialog box allows you to specify up to four filtering conditions joined by logical operators (such as AND, OR). Each condition consists of a variable, an operator and, if needed, some numerical, categorical or string values. The following table presents the various operators available for each data type.

NOMINAL / ORDINAL	Equals Does not equal Is empty Is not empty
NUMERIC and DATE	Equals Does not equal Is greater than Is lesser than Is greater than or equal to Is lesser than or equal to Is empty Is not empty
BOOLEAN	Is true Is false
STRING	Contains Does not contain Is empty Is not empty
DOCUMENT	Is empty Is not empty Is coded Is uncoded
IMAGE	Is empty Is not empty

The lower section allows you to filter cases based on the presence of words or phrases associated with specific content categories in the document being processed by WordStat. This section is disabled by default and becomes active as soon as a content analysis has been performed. The AND, OR and NOT Boolean operators are used to determine how those two categories should be combined. For example:

APPEARANCE AND ART Selects only cases that contain words in both categories.

APPEARANCE OR ART Selects cases that contain words in either one of these two categories.

APPEARANCE **NOT** ART Selects all cases that contains a word in category APPEARANCE but that do not contain any word found in the ART category.

- Once a filtering expression has been entered, you can apply the filter and leave this dialog box by clicking the **Apply** button. If the filter expression is invalid, an error message will appear and exiting from the dialog box will not occur.
- To temporarily deactivate the current filter expression, click the Ignore button. The filter expression will be kept in memory and may be reactivated by selecting the FILTER CASES command again and clicking Apply.
- To exit from the dialog box and restore the previously active filtering expression, click the Close button.

Geocoding

Geocoding is the process by which geographical information such as place names, postal codes or IP addresses are transformed into corresponding geographic coordinates. WordStat can create geographic maps by relating coded segments with existing geographic coordinates such as latitude and longitude, or projected coordinates. While it cannot create maps directly from a city name or a postal code, it offers a geocoding service allowing you to transform such kinds of information into latitude and longitude.

WordStat can geocode four types of data:

- Cities, States/Provinces & Countries
- Country names
- Postal code or zip codes
- IP addresses

This transformation process is performed using a geocoding server. It requires an Internet connection. The result is then stored in the project as numerical variables for the corresponding latitude and longitude. When an IP address is geocoded, you can also obtain a city name, a region and a country name and store the information into string variables. Once the variables are created, mapping can be done without the need of an Internet connection.

To access the geocoding feature:

• Select the **GEOCODING** command from the ≡ menu. You may also access this feature by going through the mapping feature and then through the geographic coordinates setup dialog box. The geocoding dialog box looks like the one below:

Source variable(s);			
Cities:	CITY	v	Default value:
States/Provinces:	2	v	
Countries:		~	Canada
Output variables:			
Latitude: LATIT	UDE	Lon	gitude: LONGITUDE
Do not overv	rite existing v	alues (if exis	sts)

- The From list box allows you to select the type of geographic information that can be geocoded. Once set, the dialog box automatically adjusts the **Source variable(s)** option box to display the information that needs to be specified. The above example displays the options for geocoding city names. Set the list boxes to the string variables containing the required information. You may also set default values for some fields. These values will be used if no variables are selected, or if a specific case has a missing value for the selected variable. For example, if you want to geocode US zip codes, you may simply set the postal code variable to the string variable containing the zip code and type 'USA' or 'US' in the country default value text box.
- The Output variables group box allows you to specify the name of the numerical variables in which the Latitude and Longitude information will be stored. If the variables do not exist, the software will ask to confirm their creation. By default, if the variables already exist, their values will be replaced by new values. If you set the Do not overwrite existing values check box, only cases containing empty values will be geocoded. Such a feature may be useful when a secondary type of geographic information is used to complement a first geocoding. For example, if you have, for each case in a dataset, a postal code and an IP address, you may start by geocoding all cases with postal codes

and then complement the cases with no postal codes or those for which no geographic coordinates could be found by geocoding their IP addresses.

Once all the required options have been set up, click the OK button to perform geocoding. WordStat will use the current Internet connection to send all the information to be geocoded to a server that will return to WordStat the corresponding latitudes and longitudes. A report dialog will display a summary of the geocoding process, indicating how many items and cases have been successfully geocoded and how many items could not be geocoded. Clicking the View button will bring a list of all items that could not be geocoded. This table can then be saved on disk or printed, allowing you to easily identify cases that may need revision.

Mapping

The mapping feature of WordSat allows you to analyze the spatial distribution of terms, phrases, content categories or topics. Mapping codes requires relating the codes to geographic location information stored in the project variables in the form either of longitude and latitude variables or of geographic projected coordinates (X and Y). In order to obtain such geographic information, WordStat allows you to compute latitude and longitude data by transforming existing variables such as postal codes, city names, country names as well as IP addresses into their corresponding geographic coordinates. Please see <u>Geocoding</u> for more information.

The mapping of codes function is available anywhere you see the *context* icon. You can find this icon on the Frequencies tab, Extraction tab (topics and phrases) and the Crosstab tab.

To map specific text features:

- Select rows of all items you want to map. For the Topic extraction tool, only one topic may be selected, yet several topics may be selected later.
- Click the distance button. A dialog box similar to this one will appear:

	Law and the second					1
Topics to map:	[Globalization,	Local Economics, Ethic, Power, I	Race/Ethnic relation,Fre	edom,Environ	ment,Nov	e v
Geographic Coordinates:	Latitude & long	gitude	🗸 💮 Setup			
Base map:	OpenStreetMa	p Tiles (online)	~			
Shapefile:	USA-states		~			
Additional variables:	٥					~
lacemark Content						
Case descriptor		Text segment	Aggregate text	segments by	case	
Document Variable		Code name				
Values of additional var	iables	Coder name	Coding date			
Markers Setup						
	1					

• The **Topics to Map** checklist box displays all selected topics. You may add additional topics by clicking the down arrow key at the right end of the list box. You will be presented with a list of all available topics. Select all topics you want to map.

• The **Geographic Coordinates** allows you to select which set of variables contains the geographic location information to be used for mapping items. Two types of coordinates can be used: 1) latitude & longitude or 2) projected coordinates. Clicking the **Setup** button will bring a dialog box similar to this one:

Geographic information					×
Latitudes:	LATITUDE	~ Lon	gitudes:	LONGITUDE	· ~
Projected X Coordinates:	<none></none>	✓ Y Coor	dinates:	<none></none>	~
Projection:	Ĺ				
	Ge	eocode Variables		🗸 ок	X Cancel

You can associate existing numerical variables to geographical properties such as latitudes, longitudes or projected coordinates. In the latter case, you also need to choose the appropriate geographic projection. The dialog box also offers the possibility to geocode existing variables, computing latitude and longitude variables from other types of geographic information (see <u>Geocoding</u>).

- The **Base Map** list box allows you to choose the background map on which data points will be plotted. It may consist of a web map service (WMS) or a geolocalized raster file (or world files). Web map services (identified with the 'online' suffix) require an Internet connection in order to retrieve relevant map files, while raster files may be used without a connection. Setting this option to **Other...** will allow you to select an image file located on your current PC. This image file may consist of a satellite photo, a street map, a topographical map, or any other representation of a geographic region as long as it is accompanied by a World file that contains the appropriate information to convert the image coordinates to real-world geographic coordinates. When a valid image file is selected, the user will be offered a choice to move this file to the default mapping resource folder or to leave the file where it is.
- The Shapefile list box displays the currently available shapefiles. Shapefiles are geospatial vector data format for GIS systems. They typically contain series of polygons, lines and points representing geographical objects, such as countries, states, counties, lakes, etc. Shapefiles are used in WordStat to create thematic maps in which areas are shaded or color-coded in proportion of specific measurements. WordStat provides a few basic shapefiles for representing either all countries or specific continents, as well as some detailed shapefiles for some cities or countries. Thousands of additional shapefile maps can be downloaded for free from various websites. Setting this option to Other... will allow you to select a shapefile located on your PC. When a valid shapefile is selected, the user will be offered a choice to move this file (along with all associated files) to the default mapping resource folder or to leave the file at its current location.
- The Additional variables checklist box allows you to transfer to the mapping module additional information stored in numerical, categorical, string or date variables. Values of these variables may then be used for analysis, displayed or exported.
- The Placemark Content group of options allows you to select how data points will be created as well as what additional information will be used to describe each data point. The following information can be transferred: The Case Descriptor which is a user-defined string used to describe the case from which a specific text segment comes from, the Document Variable from which this text segment has been extracted, as well as the Values of Additional Variables (if any additional variable has been selected). You can also transfer the associated Text Segment. Selecting the Aggregate text segments by case option will merge all the text into a single text segment and produce a single data point per topic. Additionally, you can transfer the Code Name, the Coder Name and the Coding Date.
- The Markers Setup section gives access to two dialog boxes that may be used to configure the appearance of markers (or placemarks). Selecting Assign Manually calls a dialog box that allows you to associate specific topics to specific symbols. The associations are stored in the project so that they can be reused later. Selecting Generate Automatically allows you to configure the sequence in which symbols are being generated. For more information on manual or automatic marker assignments see <u>Customizing Markers</u>.
- Once all the settings have been adjusted, click the OK button to call the mapping module.

Customizing Markers

By default, WordStat automatically assigns various symbols to different topics, cycling through a list of symbols and colors to reduce the likelihood that two topics will have the same marker. You can customize the process by which WordStat assigns symbols and can adjust their visual properties such as their size and outlines. You can also customize symbols for specific topics, so that the same symbol will be used whenever this topic is mapped.

To customize the way symbols are automatically assigned:

• Click the Automatic Generation button. A dialog similar to this one will appear.

Automatic Generation		×
Generate:	○ Symbols first ○	Colors first
Color palette: Victorian	~	
Marker size: 16		
Outline: Black	~	
Overwrite existing symbols		
		VOK X Cancel

Generate: This option specifies how WordStat will cycle through symbols and colors. By default, the software changes both the symbol and the color for every consecutive topic. Selecting **Symbols first** will cycle through the seven symbols but keep the color constant, and will only change the colors when all symbols have been used. Selecting **Color First** will keep the symbols constant until all colors of the select color palette have been used.

Color palette: This option allows you to choose a set of colors to be used in the automatic marker generation process.

Marker size: This option lets you adjust the size of the marker.

Outline: When this option is checked, markers will be surrounded by a single line. The color can be set using the color list box.

Overwrite existing symbols: When custom symbols have been assigned to specific topics, their symbols and colors will remain unchanged even if the automatic marker generation process is used for assigning symbols and colors to other topics. Enabling this option will cause the custom markers to be overwritten with new symbols and colors.

- Set your automatic generation options.
- Click OK.

By default, markers are automatically generated. However, custom markers can be assigned to specific topics, content categories, or words and saved within the project file so that the same marker will be used every time the topics are plotted. The customization feature also allows you to assign a bitmap or a letter to a text element.

To assign specific markers to text features:

• Click the Assign Manually button. A dialog similar to the one below one will appear.

Categories	Symbol	Symbol:
Globalization	•	Cirde ~
Local Economics		Fill:
Ethic		Red
Power	•	
Race/Ethnic relation	•	Size:
Freedom		
Environment	*	Outline:
Novelty		📕 Gray 🗸
Representativeness		1
Family		
Protectionism		
Tradition	×	Conorato
Patriotism	×	Generate
		Л ок
		- OIC

• To assign a custom symbol to a word, a content category or a topic, you first need to select it from the list on the left of the dialog box and then adjust any of the following options:

Symbol: This list box displays the various symbols that could be used as markers for this item. **Bitmap...** will bring an open file dialog box that will allow you to choose a bitmap image file (*.bmp; *.bmp or *.bmp). When an image is selected, the following three options are disabled since they only apply for symbol markers.

Fill: This option allows you to choose a color to be used for the marker.

Size: This option lets you adjust the size of the marker.

Outline: When this option is checked, the current marker will be surrounded by a single line. The color of this line can be set using the color list box below.

• To automatically assign initial markers to all items listed, click the **Generate** button.

Text Editor

WordStat's integrated text editor allows browsing and editing alphanumeric variables and documents submitted to content analysis, as well as spell checking of text found in a specific case or in the entire data file. When viewing plain text documents, the editing window consists of a single editing view where text can be edited and keywords are highlighted. When viewing rich text documents, the editing and keyword highlighting features are accessed through two separate views. The keyword highlighting feature allows you to identify all words or phrases that have been coded as well as those belonging to specific coding categories. The text editor can also be used for dictionary maintenance tasks by allowing the addition of words or expressions to active dictionaries. You can also jump directly from a selected word to a keyword-in-context table of all instances of this word.



It is also possible from this dialog box to examine and edit all numeric and alpha numeric values stored in other variables of the data file. To view and edit the values for the current case, the **Show all variables** check box should be selected. When enabled, the screen will be split vertically. On the left side, a panel with a list of all variables with their values for this case will be shown. To edit any of the values, press the **F2** key or double-click the value to edit.

The following table provides a short description of buttons and controls of the text editor dialog box:

Control Description

- This button allows the importation of text from various file formats including plain text file, RTF, MS Word, WordPerfect, MS Write or HTML. If the variable containing the document supports only plain text documents then all formatting options and unsupported features, such as bullets, graphics or headers are removed.
- Export the current document to disk. Plain text document may only be saved as plain ANSI document while RTF documents may be saved in plain text or RTF format.
- Print the current document.
- \gtrsim Cut the selected text to the clipboard.
- Copy the selected text to the clipboard.
- Paste text from the clipboard at the current cursor position.
- Reverse the last action made to the text.
- A Search for a specified word or phrase.
- Search for and replace a specified word or phrase.
- Spell-check documents for all cases.
- Spell check only the current document.
- Pressing this button allows you to add the selected word to an active dictionary. It may also be used to produce a KWIC table of the currently selected word or expression.
- Variable: This drop-down list box allows you to select the alphanumeric or memo variable displayed in the edit box.
- **Highlight:** WordStat's text editor displays all words that have been coded using bold characters while words belonging to the active category are shown in blue. To change the active category and highlight all words that belong to a selected category, simply choose the proper category from this list box. To select all categories using different colors, set this option to <all categories>.
 - Clicking this button accesses the color coding dialog box that lets you assign to each category in the dictionary specific font and background colors (see below).
 - Move to the first case of the data file.
 - Move to the last case of the data file.
 - Move to the previous case
 - Move to the next case.

Assigning Color Codes to Categories

The color code dialog box allows you to assign specific font and background colors to each category of the current categorization dictionary.

Font:	Background:
ATHLETIC	^
ATTRACTIVE	10
ATTRACTIVELY	
ATTRACTIVENES	S
BAR	
BAR-HOPPING	
BARS	
BASEBALL	
BEAUTIFUL	
BEAUTY	
BIKING	
BODY	
BOWLING	
ROXING ≪	>
Clear 🗶 (Cancel 🖌 🗸 OK

To access the color code dialog box:

- Set the Highlight drop-down list to <all categories>.
- Click the 🖪 button.

To change the colors of a category:

- Select the category in the categories list box.
- Use the **Font** and **Background** color selectors to choose from a list of predefined colors or click the **Other** button to define a custom color.

Publishing Categorization Models

A typical text categorization process may involve any one of the following steps:

- User-defined text preprocessing.
- Automatic lemmatization.
- Exclusion of words and phrases.
- Substitution of words and phrases.
- Categorization of words, word patterns, phrases and coding rules into content categories using a categorization dictionary.
- Specific settings, such as the inclusion of special characters or numerical digits, may also need to be set in order to collect relevant information.

Publishing a model file allows you to access the full text processing features of a WordStat elsewhere including:

- QDA Miner, through the **KEYWORD RETRIEVAL** feature.
- WordStat Document Explorer utility program,
- Or any other application making use of the WordStat Software Developer's Kit.

To publish a categorization model:

- Set the various analysis options on the **Text Processing** tab necessary to reproduce the required categorization process.
- Click the button located at the top of the tab. A dialog box will appear asking you for a file name.
- Enter the file name of the model you want to create and click Save.

Published categorization model files are saved with a .WCAT file extension in the My Provalis Research Projects\Models folder .

NOTE: While the information in the exclusion list and the categorization dictionary is all stored in the categorization file, running a categorization model from outside WordStat may still require the availability of some resource files such as language dictionaries or preprocessing libraries (EXE or DLL). This should not cause an inconvenience when applying the models on the same computer as the one used to create the model, since information about the original locations of the resource files is always stored within the model file. However, when attempting to apply these categorization models on another computer, the calling application may have some difficulty locating the needed resource files. The files should be stored either under paths identical to those on the original computer, in the application folder or under specific subfolders. When an attempt is made to apply a saved categorization or classification model for which some resource files are missing, an error message will be displayed providing the list of all missing files, their original location and alternate locations where they might be.

For information on how to apply a saved categorization model, please refer to the sections on <u>WordStat Document Explorer</u>, or to the <u>WordStat Software Developer's kit</u>.

Creating and Using Norm Files

A useful element in interpreting the results of a content analysis is the possibility of comparing the obtained results to some normative data and identifying how similar or dissimilar the observed frequencies are compared to those norms. For example, you may wish to compare the vocabulary of an adult victim of a brain injury to vocabularies of normal adults or the mission statement of a business to a collection of mission statements of Fortune 500 companies. You may also establish the reading level of a school manual by comparing its vocabulary to collections of books read by children of various ages. Normative data are typically computed on a large sample of documents and represent either general norms with data from a wide variety of sources or are computed on a more specific text corpus related to the channel, the domain area or the specific situation being studied. For example, you could compare the speeches of candidates of a presidential election to a large collection of English text from different sources (newspapers, novels, technical documents, etc.) or to a more specific corpus of spoken English or to a collection of political speeches. Comparison of word frequencies to norms established on a general corpus may be especially useful to identify the specific terminology of a set of documents. On the other hand, using a more specific collection may allow one to identify subtler differences or nuances.

WordStat allows you to create normative data on a collection of documents based on either the content of a categorization dictionary or on the frequency of individual words. The norms may be stored on disk and later be compared to the results of a content analysis performed on other documents. When comparing results to norms established using a categorization dictionary, it is highly recommended to use the same dictionaries and the same analysis options as the ones used to create the norms. Using different settings may result in invalid comparisons. To prevent such a situation, WordStat will detect any difference in settings and issue a warning message, pinpointing all differences in the settings. On the other hand, comparing the results to a norm file based on a comprehensive list of word frequencies may provide a more flexible solution than using a norm established using a categorization dictionary since a single word frequency norm file may be used to compare results obtained using various categorization systems. In such a situation, WordStat automatically computes from the words in the norm file the expected frequencies for each content category. However, it is important to remember that if the categorization dictionary contains phrases or rules, the expected frequency will likely be underestimated and may be invalid.

When a comparison to a norm file is performed, WordStat appends four columns to the right of the frequency table and computes each item's expected frequency, the deviation from the observed frequency, the Z value (standardized deviation) and its two-tailed probability.

To create a norm file based on content categories:

- Open WordStat.
- Select the categorization dictionary on which the norms should be computed and set the various analysis options (exclusion list, lemmatization, etc.).
- Move to the **Frequencies** tab to force the computation of frequencies on the content categories.
- Click the key button and select the SAVE AS A NORM FILE command. A file-saving dialog box will be displayed.
- Enter the name of the file under which you would like to store the norms (by default, the .wnorm file extension is added to the file name), then click **Save** to create the file.

To create a norm file based on word frequencies:

- Open WordStat.
- If a categorization dictionary is active, disable it and set the various analysis options (exclusion list, lemmatization, etc.). It is recommended to set the minimum frequency or record occurrence to "1" and to disable the option to remove words under a specific frequency or occurrence in order to obtain a detailed frequency list of all words in the normative sample.
- Move to the Frequencies tab to force the computation of word frequencies.
- Click the kit button and select the SAVE AS A WORD FREQUENCIES command. A file-saving dialog box will be displayed.

• Enter the name of the file under which you would like to store the norms (by default, the .wfreq file extension is added to the file name), then click **Save** to create the file.

To compare the obtained frequencies with existing norms:

- Set the required dictionary and options and move to the **Frequencies** tab to instruct WordStat to compute the frequencies of words or of content categories.
- Click the button and select COMPARE TO NORM FILE or COMPARE TO WORD FREQUENCIES, depending on whether you would like to compare the obtained frequencies to norms established on content categories or on words.
- Select the norm file to which you would like the comparison to be made and click OK.

To remove comparison statistics:

• Click the key button and select the **REMOVE NORM STATISTICS** command.

Program Settings

The **Program Setting** functions cover a variety of options that are program-wide rather than project-specific. These options include display options, dictionary logs and internet access etc.

To access program settings:

• Select the **PROGRAM SETTING** command from the ≡ menu in the upper left side of the Win. A dialog box will appear similar to the one below.

Program Settings	×
Display Options	
Show paragraph marks as ¶ (KWIC list)	
Percent decimal places: 3	
Color scheme: Office 2007 Luna V Flat tables (without grid lines)	
Misc.	
Treatment of items not found in norm files:	
Leave cells empty	
O Set expected frequency to lowest frequency	
Clear crosstab cells equal to zero	
└ Log all dictionary changes	
Disable all Internet access (except for activation)	
Display startup screen	
Check for updates (once per month)	
Temporary data folder: Windows default Custom path	
	3

Display Options: These options include the possibility of displaying **paragraph marks** by selecting the checkbox beside this option. You can choose the number of **decimal places** you would like when displaying percentages but typing the number in the field or by selecting the up or down arrow until the desired number is reached. The **color scheme** can be

chosen from the drop-down list and you can choose whether or not to have flat tables by selecting the checkbox beside this option.

Treatment of items not found in norm files: When comparing frequencies of words or content categories in a project to some normative frequency data, some words may be unique to the corpus being analyzed. In such situation, a comparison of the observed frequency to the expected one, based on a reference corpus is not possible. WordStat offers two ways to deal with such a situation: 1) It can leave the cells associated with the expected frequency, deviation, z value and probability empty, or 2) it may set the observed frequency to the lower possible frequency. For example, if the normative frequencies were computed on a reference corpus containing 950,000 words, then the ratio used to assess the expected frequency will be set to 1/950,000.

Clear crosstab cells equal to zero: In a crosstabulation table, empty cells are represented as either 0 or as a percentage equal to 0.0%. Enabling this option leave those cells empty.

Log all dictionary changes: When this option is enabled, all changes to the exclusion list, the substitution list or the categorization dictionary are stored in a log file, allowing one to review changes made previously. This feature is especially useful when more than one person work on a single dictionary file and one needs to review changes made previously.

Disable internet access: WordStat is a desktop text analysis tool. Your text data never have to leave your computer in order to be analyzed. The only situation where data may be sent through the internet is if you use the geocoding features of WordStat to transform text fields such as city names, country names or postal code into latitude and longitude. WordStat may however access the internet to retrieve information about potential updates to the software, to get additional resources, or to access a web mapping service used to generate maps in the GISViewer modele. It also use internet for the initial activation of the software and further validation of the software license. Selecting this feature disable all access to internet, except for the activation and validation of the software licenses.

Display startup screen:

Check for updates:

Temporary data folder:

Mode

Data can be analyzed in two different modes: **Explore** mode and **Expert** mode. The **Explore** mode lets you see at a glance what is in your data set. You will have access to extracted topics and phrases as well as a frequency list. If you would like a more comprehensive analysis, run the **Expert** mode. In **Expert** mode, in addition to all the features available in the **Explore** mode, you can also create dictionaries and crosstabs, and perform cooccurrence analysis etc. All of WordStat's capabilities are available in **Expert** mode.

To switch between Expert and Explorer mode:

- Select the **MODE** command from the \equiv menu.
- Choose either **Expert** or **Explorer** form the adjacent menu.

Exporting Frequency Data

Various case statistics may be appended to the existing data file or exported to disk in different file formats including SPSS, Stata, Excel, HTML, XML, and tab or comma delimited text files, allowing this data to be further analyzed. The resulting data consist of a matrix where each row represents a case and where the statistics on content categories will be stored in columns along with a few additional variables.

To append or export to disk content category statistics (data matrix):

• Click the solution located at the top of the **Frequencies** tab. A dialog box similar to this one will appear:

Export to data file	n x
Saving options Destination: O Append to current file O Create a new data file	🗸 ок
Data to save: Case occurrence ~	X Cancel
Variable names: Keyword OPrefix: WORD Save empty categories	? Help
Variable type (occurrence): Multiple dichotomous variables Multiple polynomial variables Multiple string variables Single compound string Indude item frequency 	
Other variables	
Add variables: [GENDER,AGEGROUP]	
Save total number of words	Preview *

Destination: This option allows you to choose whether the new variables should be appended to the current data file or in to a new file. If this last option is selected, a dialog box will appear allowing you to specify the name and location of the new file. When data are saved to a new data file, additional variables are created to store the case number and the numerical values of each independent variable. When data is saved in to a new file, clicking **OK** will call a dialog box that lets you specify the name and location of the new data file along with its type. WordStat can export to several types of files, including SPSS and Stata data files, Excel spreadsheets, tab and comma delimited files, as well as HTML and XML files.

Data to save: This option allows you to choose four different kinds of data that may be saved:

- Keyword frequency
- Case Occurrence (i.e., a dummy variable with 0 when absent or 1 when present)
- Percentage of words (i.e., the frequency of the keyword divided by the total number of words in the case)
- TF*IDF (i.e., the keyword frequency weighted by inverse document frequency).

Variable names: This option sets what method should be used by WordStat to create new variable names. When set to KEYWORD, the program will attempt to use each keyword as the name of a new variable. Illegal characters are

automatically removed and long names are truncated to the first 10 characters. Duplicated variable names are distinguished by the substitution of numerical digits at the end of the name. When this option is set to PREFIX, variable names are created by adding successive numeric values to a user-defined prefix. For example, if the edit box at the right of the prefix option is set to "WORD_", the variable names will be WORD_1, WORD_2, WORD_3, etc. The order of creation of the variables corresponds to the sort order used in the FREQUENCY tab.

Variable type: By default, WordStat saves keyword statistics in as many variables as there are keywords or content categories listed on the frequency tab. For example, if the frequency table contains 100 items, then 100 variables will be necessary to store the statistics associated with each item. When you choose to store the occurrence of codes, WordStat offers you the possibility of storing the observed occurrences in a limited number of polynomial (or multinomial) variables. For example, if the maximum number of different content categories per case is no more than 10, then you may instruct WordStat to create 10 numeric variables and store, in each of those, a numeric value representing one of the content categories. If less than 10 categories are found in a specific case, then the remaining variables are left empty. To store values in a limited set of nominal variables, choose the Multiple Polynomial Variables option and enter the Maximum Number of Variables that should be used for storing the values representing the content categories. To store the name of those categories found in a single case is higher than the specified number of variables, then a warning message will appear to let you know that some information has been lost and to indicate the maximum number of content categories encountered in the project. To export occurrences as zeros and ones in as many variables as there are codes, select the Multiple Dichotomous Variables option.

Add variables: This drop-down checklist box may be used to add the values stored in one or more variable to the exported data file along with the statistics.

Save total number of words: This option appends a numeric variable named TOTWORDS that contains the total number of words processed in each case.

Clicking the **Preview** button displays a grid allowing one to see what the data file will look like.

- Once your options have been chosen select OK.
- If you are creating a data file, choose and name, type and place to save your file and then select **Save**.

Performing Multivariate Analysis

One benefit of the integration of a content analysis module within an existing statistical program is the ability to easily perform on numerical results of content analysis, various statistical analyses such as frequency, crosstabulation, multiple regressions, reliability analysis as well as various multivariate analysis techniques (cluster analysis, factor analysis, etc.). The table below illustrates some types of analysis that may be performed with SimStat and, if needed, the required module to perform these analyses.

TYPE OF ANALYSIS	REQUIRED SIMSTAT MODULE
Multiple regression analysis	None
n-way ANOVA/ANCOVA	None
Reliability analysis	None
Factor Analysis (with or without varimax rotation)	None
Principal Component Analysis	MVSP
Advanced cluster analysis	MVSP

First step - Saving numerical results into a data file

In order to perform a statistical analysis on categories or words, you first need to create new numeric variables that will contain, for each case in the data file, the occurrence or frequency of specific words or categories.

To create variables:

- Set the various options of the Text Processing tab.
- Perform the frequency analysis by clicking the **Frequencies** tab.
- Click the button to access the **Export Data** dialog box.
- Set the Destination option to Append To Current File.
- Set the **Date to save** list box to **Keyword frequency**.
- Set the Variable names option to Keyword to instruct WordStat to use the names of the categories (or included words) as the names for the new variables.
- Click the **OK** button to proceed to the saving of the data and return to WordStat.
- Click the OK button of the WordStat dialog to return to SimStat.

For more information on how to store numeric or textual results into the current data file or into a new data file see <u>Exporting</u> <u>Frequency Data</u>.

Performing a cluster analysis of keywords:

- Select the **CHOOSE X-Y** command from the **STATISTICS** menu and assign all the newly created variables to the list of independent or dependent variables. (The distinction between dependent and independent variables is not relevant for this kind of analysis. However, all variables assigned to a single category will be processed together)
- Choose the OTHER | CLUSTER ANALYSIS command from the statistics menu to display the option dialog box.
- Set the various analysis options to your preferences. (Please take note that in order to perform a cluster analysis on words rather than on cases, the Transpose Data option should be left deactivated).
- Click the **OK** button to perform the statistical analysis.

Performing a factor analysis of words or categories:

- Select the **CHOOSE X-Y** command from the **STATISTICS** menu and assign all the newly created variables to the list of independent or dependent variables. (The distinction between dependent and independent variables is not relevant for this kind of analysis. However, all variables assigned to a single category will be processed together.)
- Choose the OTHER | FACTOR ANALYSIS command from the statistics menu.
- Set the various options to your preferences.
- Click the **OK** button to perform the statistical analysis.

Performing Analysis on Manually Entered Codes

The QDA Miner software is a computer-assisted qualitative coding tool specifically designed to manually code documents. While WordStat was designed mainly to perform automatic content analysis of textual data, you may also use it to manually assign codes to text. When used for such a purpose, codes need to be manually typed into the text itself. WordStat may then be used to extract the codes and perform various types of analyses (frequency, cross-tabulation, cooccurrences). At least two approaches of coding may be used to achieve this:

Unique Keywords

Unique keywords or code names may be inserted anywhere in the text. These keywords should preferably not be an existing dictionary word. For example, it may consist of an abbreviation of one or several words, includes special symbols (such as #, & ^, _ etc.) or numeric digits. The retrieval of these codes can then be achieved by adding all these keywords to the categorization dictionary. If special symbols have been used, they should also be specified as valid characters on the <u>Preprocessing</u> tab.

Codes in square brackets

WordStat provides an option to process only the text found between square brackets (i.e. [and]). Codes corresponding to categories may be typed directly into the text and placed within square brackets. By disabling all dictionaries and setting this option, the program simply ignores all text outside the brackets and performs a frequency count of all words found inside square brackets. This method has several advantages over the previous method. First, there is no need to enter all existing codes in the categorization dictionary, since they will be processed automatically. Misspelled keywords will always be extracted and may be easily identified and changed. Also, the processing time and memory required can be much lower than the other approach since WordStat won't process text found outside the brackets.

Besides the possibility to extract manually entered codes and perform frequency and comparison on those codes, there are several other ways in which the program may be used to assist the work of human coders. Here are just a few examples of possible uses:

- During the exploratory phase of the analysis, keyword frequency and crosstabulation may be used to identify differences between subgroups in word usages, differences that may have remained unidentified.
- The **KWIC** list may be used to locate and visualize text associated with specific codes either to validate the coding or identify associated themes. The **KWIC** list may also be used to perform a systematic search of words frequently found along with a specific code. The examination of all cases containing those words may then help identify instances where the code should have been used.
- A dendrogram of keyword cooccurrences or a crosstabulation of manually entered codes against all words in a text may allow identification of specific words associated with codes and permit the development of dictionaries. Such dictionaries may then be used either to ensure a more systematic application of manually entered codes or to develop and validate a dictionary that will later be used for automatic content analysis.

Web Collector

The **Web Collector** is an external tool that can be accessed through QDA Miner WordStat, through the system tray or the Windows start menu. It monitors RSS feeds, Reddit and Twitter posts and Facebook comments. If you have chosen to live monitor your social media search, the Web Collector will collect newly published posts or comments, which can then be appended to your project.

Setting the Web Collector

When you create an RSS, Twitter or Facebook project you have the option of monitoring the project over time. This means data will be collected by the Web Collector according to a user defined time interval until a stop date is reached or until the monitoring is disabled or deleted. Data will continue to be collected as long as the Web Collector is running. As the Web Collector is a completely separate application, QDA Miner or WordStat do not have to be running for it to work. Please see the corresponding section on how to import from <u>RSS</u>, <u>Twitter</u> or <u>Facebook</u> for more information on creating live monitoring queries.

Define Query		- E] >
ettings:		Variables	
Source: Expression: Max Posts: Web Collector	Twitter v obama 500 (0 = unlimited) Geocode Locations Indude Retweets	Date Created Tweet To Tweet To Device Source Number of Retweets Wimber of Favorites Geocoordinate Longitude Geocoordinate Latitude User ID User Soreen Name User Location User Description User Number of Fieldowers User Number of Fieldowers	5
Sto	p Date: ₩11/ 4/2016	User Time Zone User Number of Tweets User Profile Image URL	

If you have checked **Live Monitoring** query the Web Collector will run in the background. If you restart Windows, the Web Collector will start automatically and continue collecting data as long as the query remains active.

Accessing the Web Collector

• You can access the Web Collector by clicking on the system tray. Alternatively, click on the Web Collector menu item in your Windows start menu and the Web Collector will appear.

Retweets	Status	Posts	Delay	Last Search	New Posts
No	Activo		4	10/28/2016 4:30 PM	4
	No			No Active 4 1m	No Active 4 1 m 10/28/2016 4:30 PM

The Web Collector displays a list of all of the RSS, Twitter and Facebook projects that you are currently monitoring. It displays the following information:

Project name: The project associated with your query.

Application: The social media platform that is being monitored.

Expression: Either the URL of the RSS feed or Facebook page or your Twitter search query.

Retweets: For Twitter searches, this column will read **Yes** or **No** depending on whether you have chosen to include Retweets when setting up your project.

Status: Lets you know if the data is being collected. Will either read Active, Inactive or Searching.

Posts: The number of posts collected. This will be reset to zero if the Tweets, Facebook or RSS posts are appended to a project.

Delay: The time interval you set for your capture.

Last search: The date and time of the last data capture.

New posts: The number of new posts captured since you last viewed the data.

Twitter resources: The percentage of Twitter resources you have available is displayed in the bottom left of the screen.

Facebook Resources: The percentage of Facebook resources you have available is displayed in the bottom left of the screen.

Web Collector Functions

You can **Edit**, **Add**, **Delete**, **Pause** and **Resume** queries in the Web Collector by selecting the corresponding icon from the tool bar or by selecting the corresponding command from the **QUERIES** menu. You can also **view the data** associated with each query in the same manner.

To add a query:

- Select the query you want to modify.
- Click the Geboutton.
- A Define/ Edit Query dialog box will appear.
- Add your query.
- Click OK.

Please see the corresponding section on how to import from <u>RSS</u>, <u>Twitter</u> or <u>Facebook</u> for more information on creating live monitoring queries.

To edit a query:

- Select the query you want to modify.
- Click the *S* button
- A Define/ Edit Query dialog box will appear.
- Add your query.
- Click OK.
To remove a query:

- Select the query you want to modify.
- Click the 🖾 button. A warning dialog will appear.

Warning		2
Do you want to delot	e the selected se	arch?
Do you want to delet	e une selectico se	car cart:
Do you want to delet	e une selectico se	arch

• Click Yes if you want to continue with the deletion. If you do not want to delete the query click No.

To pause a query:

- Select the query you want to modify.
- Click the
 button
- The query status will change from Active to Inactive and the Web Collector will cease collecting data.
- To pause all queries at the same time select **PAUSE ALL QUERIES** from the **QUERIES** menu.

To resume a query:

- Select the query you want to modify.
- Click the Dutton.
- The query status will change from Inactive to Active and the Web Collector will resume collecting data.
- To resume all queries at the same time select **RESUME ALL QUERIES** from the **QUERIES** menu.

To view data:

- Select the query for which you would like to view the data.
- Click the 🙆 button and a Data View dialog will open.

~				
Date Created	Tweet ID	Tweet Text	Number of Retweets	Number of Favori
10/7/2016 3:29 PM	784415384250224 640	@RMarcelc impeach obamai don't care if he has less than 30 days left or whateverCONGRESS , WHAT'S YOUR PROBLEM? SCARED SHITLESS?	0	0
10/7/2016 3:29 PM	784415400113168 389	The 6 key questions for the Obama administration about the Yahoo-NSA cooperation. By @neemaguliani @ACLU https://t.co/eDX3s0BWlh	0	0
10/7/2016 3:29 PM	784415402281533 441	Seven Years Ago This Week: Obama's Nobel Peace Prize https://t.co/6k5uaQA7QJ	0	0
10/7/2016 3:29 PM	784415402692489 220	@realityinACTION @RepMcCaul @SenRonJohnson @PRyan @CNM_ Michael Welcome to Clinton/Obama/Soros Urban Area Terrorist Relocation Program	0	0
10/7/2016 3:29 PM	784415404265381 888	@grumpynaomi @purposehoIy Well actually this isn't widely known but Obama was actually being Satan at the time.	0	1
10/7/2016 3:29 PM	784415405792264 192	https://t.co/Ufnz5lRaCpبالتحقيق يطالب : كيري جون الحرب لجرائم بالنسبة روسيا مع	0	0
10/7/2016 3:29 PM	784415407679627 264	#WOWTHIS IS WONDERFUL#AWESOME29 PRESIDENT OBAMA WITH HIS GRAND MOTHER IN KENYAWATCH THIS https://t.co/eiuapY0dJZ	0	0
10/7/2016 3:29 PM	784415408732241 920	Obama alcanza un nuevo récord de aprobación en su presidencia https://t.co/tyuHJIhE5W https://t.co/PX7cW5graq	0	0
10/7/2016 3:29 PM	784415415942287 360	Newly disclosed emails prove Obama officials were in close contact with Hillary Clinton re potential fallout from her emails @KellyannePolls	0	0
10/7/2016 3:29 PM	784415420308656 129	#US needs a LEADER that Knows 2 Win Enemies Not Elections ONLY!(e.g.Obama) @DTAP4NSecurity @DomusUSA https://t.co/nriAFNNETG	0	0
10/7/2016 3:29 PM	784415422850314	Just remember, under Obama and the democrats, the US lost its AAA for the first time ever.	0	0

The table contains the variables you selected to include when you initially created your project. Posts captured since you last viewed the project will be in bold. Posts that have already been viewed will **not** be in bold.

- Select the Q button to enable the search function at the bottom of the screen.
- Type your search term into the field and select Find next or Find previous button.
- Check the Match case checkbox if you want the text case to be exactly as you have written it in the search box.

To close the Web Collector:

• Click the X at the top right of the dialog box.

Word Frequency Analysis

WordStat can perform basic word frequency analysis allowing one to quickly get a glance of what are the most common terms in large amounts of text. A word frequency analysis may be useful to identify themes as well as associated terms that may later be used to retrieve segments related to those themes. Such a feature may also be useful to identify co-occurring words, allowing one to expand the search terms used for retrieving relevant text segments. The results are presented both visually, using a word cloud graphical display, and in a table format, for more precise frequency information.

In WordStat, the word cloud represents the 200 most frequent words, where the size of each word is proportional to its relative frequency of use in the text.



The **Frequencies** tab allows one to display the results of all the words in a table format. By default, words are displayed in descending order of frequency. A second column also displays this frequency using a rate per 10,000 words. Such a transformation facilitates comparisons across different text collections by eliminating the effect of different text sizes. To sort the table in alphabetical order, click on the first column header. To sort the table by frequency, click the second column header once for sorting the table in ascending order of frequency and click a second time to sort the list in descending order of frequency.

0			(×
Word Cloud Free	quencies				
×				u	
Word	Frequency	Rate per 10K			^
PRESIDENT	106	42.2			
WAR	73	29.1			
GOVERNMENT	72	28.7			
AMERICA	71	28.3			
PEOPLE	69	27.5			
YEARS	56	22.3			
IRAQ	55	21.9			
CARE	49	19.5			
AMERICAN	46	18.3			
MAKE	46	18.3			
SECURITY	46	18.3			
WORLD	43	17.1			
NATIONAL	39	15.5			
HEALTH	36	14.3			
CONGRESS	.35	13.9			
MILITARY	35	13.9			
TIME	35	13.9			
END	32	12.8			
PLAN	31	12.4			
TRADE	31	12.4			
ECONOMIC	30	12.0			
FUTURE	30	12.0			

In WordStat, quick word frequency analysis may be performed on the text segments obtained through several text search tools, such as the keyword retrieval, the keyword-in-context page, as well as from a few additional locations.

Working With the Word Cloud

Items may be removed temporarily from the word cloud. Its appearance may also be customized.

To remove a word temporarily and redraw a new word cloud without it

• Click on the word you want to remove, and then click the × button. You may also right-click on the word and select **Remove & Redraw** from the popup menu.

Working With the Word Frequency table

The same operation to remove words may be achieved from the word frequency table (see section above for specific instructions). However, it adds the ability to perform this operation on more than one word at a time.

Customizing the Appearance of the Word Cloud

The word cloud shape, its color palette, as well as the method used to position words may be adjusted. There are two methods to adjust the colors used to display words.

Clicking the **iii** button will allow you to select a color palette to use from a predefine list. With this method, colors are assigned sequentially and results in colors uniformly distributed across the words irrespective of their frequency. One can also select a single color to be used for all words. Clicking the **iii** button brings a dialog box with an additional possibility to assign a gradient of colors that will be applied based on each word frequency. To enable this option, on the **Text Color** page of this dialog box, select the **Gradient Scheme** radio button and select the colors to be used for high frequency (**From**) and for those with low frequency (**To**). An optional third color may be used to represent **Middle** frequency. The following table presents additional options available.

Font This option allows you to specify the font to be used to display words.

Shadow Enabling this option displays a small drop shadow behind each word, creating a subtle 3D effect giving the impression that the words are raised above the background.

Form This option allows you to choose to display words in inside a **rectangular** or an **elliptical** shape.

Starting position When drawing a word cloud, WordStat starts positioning high frequency words and then attempts to place all the other words by decreasing order of frequency. When the starting position is set to **from center** it will position the most frequent one at the center of the cloud, and progressively position the following ones around it so that lowest frequency words will most likely be plotted on the edge of the cloud, far from the center. Setting this option to **Random** will position words in a more random fashion such that high frequency words may be more evenly distributed within the cloud.

WordStat Software Development Kit (SDK)

WordStat Text Classification Process.

Text mining, as performed by WordStat, involves some form of quantification of text data. Such quantification is achieved by applying natural language processing techniques (stemming, lemmatization, removal of stop words, etc.), statistical selection criteria, as well as grouping of words and phrases into concepts using either lexicons or custom content dictionaries. All these procedures can be combined to extract numbers representing the presence or frequency of important keywords, or key concepts. We call this the **categorization process**. WordStat also supports another form of quantification: **automatic document classification**, by which documents are categorized in one of several mutually exclusive classes using some form of machine learning.

Why a Software Development Kit?

The categorization and classification processes are performed by WordStat, which offers a graphical user interface that allows a user to create, validate and refine those processes, apply these to various text collections, perform comparisons, explore, relate and create graphical and tabular reports. While categorization and classification models can be saved to disk and reapplied on a different set of documents, a human operator is still required to perform those analysis, limiting the ability to fully automate the text analysis and reporting operations.

The WordStat software development kit (SDK) provides a solution, allowing models developed with the WordStat desktop tool to be used in other applications written in other computer languages such as C++, Delphi, C#, VB.Net and so on.

An example of such integration would be the application of a categorization model on a company data collection system of customer feedback in order to automatically measure references to specific topics and to classify those feedback as either positive, negative or neutral.

Applying the SDK

All the analysis and text transformation settings set in WordStat are stored on disk in the model files (stemming, lemmatization, categorization rules, selection criteria, etc.). This greatly simplifies the integration of such text processing in other applications by reducing the application of those text analysis process to four easy steps:

- 1. Load the categorization or classification model file
- 2. Retrieve the text to categorize or classify
- 3. Apply the model to the text
- 4. Retrieve relevant information (frequencies, probabilities, predicted classes, etc.)

A model only needs to be loaded once (step #1), while steps #2 to #4 may be repeated as often as needed.

There are currently no reporting or graphing functions available in the DLL, so it is the task of the programmer to further process the obtained information. Typically, numerical values would be either stored in a database or cumulated to create reports, dashboards, etc..

Technical Details

The SDK consists of a Windows DLL available in both 32 bits and 64 bits versions. The DLL is multi-thread safe, allowing text quantification of multiple documents concurrently. It also supports the simultaneous application of multiple categorization and classification models, allowing one to perform several quantifications of the same documents.

The SDK comes with a sample project with source files illustrating how integration can be achieved. This sample project is currently available in Delphi, C# and VB.NET. Please contact us if you need assistance on how to use the SDK with other computer languages.

Interested?

If you are interested in obtaining more information about the SDK, in getting a trial version (with documentation and sample applications) or to purchase it, please contact <u>developpers@provalisresearch.com</u>.

Report Manager

The Report Manager is a separate application that has been designed to store, edit and organize documents, notes, quotes, tables of results, graphics and images created by QDA Miner or imported from other applications. Items can be added to the Report Manager directly from QDA Miner without needing to run the Report Manager.

The 🔯 button, found in many locations in QDA Miner, may be used to copy entire documents, tables and charts to the Report Manager. Selected text segments or image areas may also be appended by clicking the 두 button.

To access the Report Manager from QDA Miner, run the **REPORT MANAGER** command from the **PROJECT** menu.

The program presents its information as an outline, allowing a hierarchical organization of miscellaneous pieces of information that is ideal for project management, organizing ideas, structuring information, or designing and writing a research report.

The workspace emulates the appearance of Windows Explorer or of a standard Help file with the Table of Contents (TOC) on the left and the Editor on the right.



Table of Contents panel

Report Manager files are made up of items or topics that are like chapters in a book. Each item can be thought of as a separate word processor, a table or a graphic file editor or viewer, all of which are stored together in the QDA Miner project file. This panel provides powerful functions for organizing items and structuring the information in a hierarchical manner.

Item Editor

The largest panel on the right of the program window is the Item Editor, which is like a built-in word processor. This is where the item selected in the Table of Contents can be edited. Clicking a Table of Contents item displays its contents for editing.

Toolbar

The Toolbar provides quick access to the most frequently-used functions. Just position the mouse over a tool button and wait for the display of a brief text describing its function.

Comment panel

The Comment panel below the Item Editor allows the insertion and editing of comments related to the selected topic. When new items are added to the Report Manager from QDA Miner, a default comment is often already present, providing useful information about the origin of this item.

Working with the Table of Contents

Creating a New Item

To create a new item:

- Select the Table of Contents entry that will be the "parent" or "sibling" of the new item.
- Select the **NEW** command from the **ITEMS** menu or click the + button. A dialog box like this one will appear:

New Item	×
Title: {untitled}	
Location: Under current item	O After current item
Content: Document Folder	🗙 Cancel 🛛 🖌 OK

- Enter the title for the new item.
- If the new item should be a "child" of the selected item, click **Under Current Item;** if the item should be positioned after the current item, then set it to **After Current Item**.
- Select whether the new item will be a **Document** or a **Folder**. Folders are empty items that are used as containers for other items.
- Click OK. The new item will become the current one.

Importing Items from Files

To import items from files:

- Select the item under which the imported items will be stored.
- Select the **IMPORT FILES** command from the **ITEMS** menu or click the toolbar button. An Open File dialog box will appear.
- Select the type of data you would like to import by selecting the appropriate Files of Type list box option. The Report Manager can import the following data types:
 - DOCUMENTS Plain text (.TXT), MS Word (.DOC), WordPerfect (.WPD), Rich Text (.RTF) or HTML files (.HTM or .HTML)

- GRAPHICS Windows Bitmap (.BMP), Windows Metafile (.WMF), JPEG files (.JPG or .JPEG) and Portable Network Graphic files (.PNG)
- CHARTS QDA Miner or WordStat Charts (.WSX)
- DELIMITED DATA Tab delimited (.TAB) or Comma Separated Value (.CSV) data files.
- Select one or several files to be imported and click the **Open** button.

Renaming an Item

To rename an item:

- Select the item to be renamed.
- Select the **RENAME** command from the **ITEMS** menu or click the rel toolbar button.
- In the Item Title Dialog, change the title.
- Click OK.

Deleting an Item

To delete an item:

- Select the item to delete in the Table on Contents.
- Select the **DELETE** command from the **FILE** menu or click the toolbar button.
- You will be asked to confirm that you really want to delete the item. If you're sure, then click Yes.

NOTE: Be aware that you cannot undo this if you make a mistake.

Moving an Item

As more items are created and the Report Manager hierarchy grows, it is inevitable that you will want to move items around, either to place one item under another, or to promote one to a higher level.

To move items using drag-and-drop:

The easiest way to move items in the Table of Contents is by using drag-and-drop operations. Using the mouse, you can move an item to a different location or move a group of items stored under a "parent" item by dragging this "parent" item to its new location.

- Select the item to move by clicking and holding down the left mouse button. (Keep the mouse pressed until the dragand-drop operation is completed.)
- Drag the item to its new location and, only then, release the mouse button.
- The dragged item will now become a "child" of the destination item.
- To move the item to the same level as the item under the cursor, simply hold the ALT key while dropping the dragged item.

To move items using the toolbar:

You can also use menu commands and toolbar buttons to move items. To promote an item is to move it to a higher level in the hierarchy, making the item a "sibling" to its former "parent". To demote an item is to move it to a lower level and make it a "child" of its previous "sibling". After selecting the item you want to move, use one of the following four commands:

- To promote the selected item, click the 🗺 button, or select the **PROMOTE** command from the **ITEMS** menu.
- To demote the selected item, click the button, or select the **DEMOTE** command from the **ITEMS** menu.
- To move the selected item up relative to its siblings, click the 🔟 button or select the MOVE UP command from the ITEMS menu.
- To move an item down relative to its siblings, click the **ITEMS** button or select the **MOVE DOWN** command from the **ITEMS** menu.

Adding or Editing Item Comments

The Comment panel below the Item Editor allows one to insert a new comment or edit an existing one related to a selected topic.

To type a new comment or to edit one, simply click in the yellow region of this panel and start to type. The comment is automatically saved as soon as you move to another item or leave the Comment panel. While there is no menu item or toolbar icon associated with this feature, standard clipboard operations are supported for copy and pasting text. A popup menu with standard editing features can also be obtained by clicking the right mouse button.

To search for text in comments, see the information on the <u>Global Search</u> command.

Editing Documents

The Report Manager offers many editing features to create and edit both simple text documents and documents with complex formatting, as well as tables and graphics. When a document item is selected in the Table of Contents, a **DOCUMENT** menu appears, displaying all available formatting and editing options. A similar menu can also be obtained by right-clicking anywhere in this document. The toolbar portion directly over the editing area also displays buttons to access the most often-used editing and formatting functions.

Individual documents may also be printed or exported to disk in various file formats such as plain text, Rich Text or HTML format. An **IMPORT** command is also available to read a document file stored in plain text, Rich Text, MS Word, WordPerfect, HTML and a few additional formats. Executing such a command will replace the existing content with the content of the imported file.

Editing Tables

The Report Manager offers many editing features to customize the appearance of a table, to change the text alignment or font setting, to set the cell background color, or to delete entire rows or columns. When a table item is selected in the Table of Contents, a **TABLE** menu become visible displaying all available formatting and editing options. A similar menu can also be obtained by right-clicking anywhere in this table. The toolbar portion above the editing area also displays buttons to access the most often-used table editing and formatting functions.

Individual tables may also be printed or exported to disk in various file formats such as ASCII (*.TXT), tab delimited file (*.TAB), comma delimited (*.CSV), MS Word (*.DOC), HTML (*.HTM; *.HTML), XML (*.XML) and Excel spreadsheet file (*.XLS).

Editing Charts

Many charts saved in the Report Manager may be edited using many of the same options as those available in QDA Miner, such as the multidimensional scaling plot obtained through the **CODE CO-OCCURRENCE** analysis command, the correspondence analysis plots, and the bar charts and line charts created by the **CODING BY VARIABLES** command, as well as the bar charts and pie charts produced by the **CODING FREQUENCY** command. To obtain information on the display options available for those charts, see their corresponding page in this manual. Other charts - such as dendrograms and heatmaps - are stored as image files, so cannot be modified. However, just like other charts, they may be exported to disk in various file formats such BMP, JPG or PNG graphic files.

Searching and Replacing Text

Two broad types of text search are available in the Report Manager. A local, item-based search-and-replace feature allows one to perform text searches and replacements on individual documents or tables, and a global search engine for searching text patterns in several or all documents, tables and comments in the Report Manager.

To perform a search in a single document or a table:

- Select the document or table you would like to search, by selecting its entry in the Table of Contents.
- Position the editing cursor in the document or select the cell in the table where you want the search to begin.
- Select the FIND command from the SEARCH menu, or click the button.
- Enter a search expression, set the desired search options and then click Find Next.
- To find additional instances of the same text, continue to click Find Next.

To replace text in a single document or table:

- Select the document or table in which you would like to perform the text replacement by selecting its entry in the Table of Contents.
- Position the editing cursor in the document or select the cell in the table where you want the search to begin.
- Select the **REPLACE** command from the **SEARCH** menu, or click the
- In Find What, type the characters or words you want to find. In **Replace With**, type the text you want to replace it with. Set the desired search options and then click **Find Next**. Click **Replace** to change the selected text. To replace all instances of the text, click **Replace All**.

To perform a global search:

Select the GLOBAL FIND method from the SEARCH menu. A dialog box similar to this one will appear:

bal search			
Text to find			
Options	Tre	e Direction:	Feitfler
Case sensitive	C	OUp	
Whole word on	y O Down		Lancei
Scope:			
Documents	Tables	Comments	

The **Text to Find** edit box allows you to specify the text you want to find. The **Case Sensitive** and **Whole Word Only** options function in the same way as in a standard word processor.

The search starts at the current topic item. Select **Forward** to continue searching items below the current one, or **Backward** to move up and search items above the current document or table in reverse order. To search all items, select the top item in the Table of Contents before using the Global Search dialog box.

The **Scope** option box is used to specify what is to be searched. You can restrict the search to **Documents**, **Tables**, or **Comments** attached to items, or any combination of these three.

Once the search options have been set, click the **Find** button to start the search as well as to continue searching for additional instances of this text.

Exporting items to HTML or Word

Individual documents, tables, graphics and images may be exported to disk in numerous formats. Such an exportation can be achieved by clicking the **EXPORT** command from the associated menu.

The Report Manager also offers the possibility of exporting the entire content or selected items into a single HTML or MS Word document. The exportation is achieved by selecting the proper command from the **FILE | EXPORT** menu. For example, to export items to HTML, select the HTML command. A dialog box similar to this one will appear.



By default, all items are marked for export. To prevent some items from being included in the exported file, simply remove the check marks beside them. Clicking a "parent" item affects all "children" items in the same way. To unselect all items, uncheck the project item located at the very top of the tree.

Once the selection process is completed, click the **OK** button. A Save File dialog box will be displayed, allowing you to enter a file name and select the location where the file should be saved. After the file is created, you will be asked if you want to view this file. Clicking **Yes** will open a web browser if the exported file is an HTML document or either MS Word or Wordpad if the exported file is a Word document.

References on Content Analysis

Selected references on content analysis

ALEXA, M. (1997). Computer-assisted text analysis methodology in the social sciences. Mannheim, Germany: ZUMA.

EVANS, W. (1996). Computer-supported content analysis: Trends, tools, and techniques. Social Science Computer Review, 14 (3), 269-279.

KRIPPENDORFF, K. (2019). Content analysis: An Introduction to its methodology (4th ed.). Los Angeles, California: Sage Publications.

LEBART, L., SALEM, A. & BERRY, L. (2011). Exploring Textual Data. Dordrecht, Netherlands: Springer.

NEUENDORF, K.A. (2016). The Content Analysis Guidebook. Beverly Hills, California: Sage Publications.

WEBER, R. P. (2008). *Basic Content Analysis.* Second edition. Quantitative Applications in the Social Sciences, vol 49. Newbury Park, California: Sage Publications.

WEBER, R. P. (1983). Measurement models for content analysis. Quality and Quantity, 17 (2), 127-149.

WEBER, R. P. (1984): Computer-aided content analysis: A short primer. Qualitative Sociology, 7 (1-2), 126-147.

Selected references on text mining

INGERSOLL, G.S., MORTON, T.S. & FARRIS, A.L. (2013). Taming Text: How to Find, Organise, and Nanipulate it. Manning Publications.

IGNATOW, G. & MIHALCEA, R. (2016). Text Mining: A Guidebook for the Social Sciences. Sage Publications.

IGNATOW, G. & MIHALCEA, R. (2017). An Introduction to Text Mining: Research Design, Data Collection, and Analysis. Sage Publications.

WIEDEMANN, G. (2016). Text Mining for Qualitative Data Analysis in the Social Sciences. Springer Fachmedien Wiesbaden.

Multivariate analysis

GREENACRE, M. (1989). Theory and Applications of Correspondence Analysis. London, England: Academic Press.

Other

GREFENSTETTE, G. (1994). Corpus-Derived First, Second and Third-Order Word Affinities. In W. Martin, W. Meijs, M. Moerland, E. ten Pas, P. van Sterkenburg, and P. Vossen, editors, Proceedings of EURALEX'94, Amsterdam, The Netherlands.

SEBASTIANI, F. (2002). Machine Learning in Automated Text Categorization. ACM Computing Surveys, 34(1):1-47.