

Getting Started with SAS[®] Text Miner 12.1



The correct bibliographic citation for this manual is as follows: SAS Institute Inc 2012. *Getting Started with SAS® Text Miner 12.1*. Cary, NC: SAS Institute Inc.

Getting Started with SAS® Text Miner 12.1

Copyright © 2012, SAS Institute Inc., Cary, NC, USA

All rights reserved. Produced in the United States of America.

For a hardcopy book: No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, or otherwise, without the prior written permission of the publisher, SAS Institute Inc.

For a Web download or e-book: Your use of this publication shall be governed by the terms established by the vendor at the time you acquire this publication.

U.S. Government Restricted Rights Notice: Use, duplication, or disclosure of this software and related documentation by the U.S. government is subject to the Agreement with SAS Institute and the restrictions set forth in FAR 52.227–19 Commercial Computer Software-Restricted Rights (June 1987).

SAS Institute Inc., SAS Campus Drive, Cary, North Carolina 27513.

1st electronic book, August 2012

SAS® Publishing provides a complete selection of books and electronic products to help customers use SAS software to its fullest potential. For more information about our e-books, e-learning products, CDs, and hard-copy books, visit the SAS Publishing Web site at

support.sas.com/publishing or call 1-800-727-3228.

SAS® and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Other brand and product names are registered trademarks or trademarks of their respective companies.

Contents

<i>Overview</i>	<i>v</i>
<i>Recommended Reading</i>	<i>vii</i>
Chapter 1 • Introduction to Text Mining and SAS Text Miner 12.1	1
What Is Text Mining?	1
What Is SAS Text Miner 12.1?	2
The Text Mining Process	3
Accessibility Features of SAS Text Miner 12.1	4
Chapter 2 • Learning by Example: Using SAS Text Miner 12.1	5
About the Scenario in This Book	5
Prerequisites for This Scenario	7
How to Get Help for SAS Text Miner 12.1	7
Chapter 3 • Setting Up Your Project	9
About the Tasks That You Will Perform	9
Create a Project	9
Create a Library	10
Explore and Modify the Data	11
Create a Data Source	12
Create a Diagram	13
Chapter 4 • Analyzing the SYMPTOM_TEXT Variable	15
About the Tasks That You Will Perform	15
Identify Input Data	15
Partition Input Data	16
Parse Data	17
Filter Data	18
Cluster Data	19
View Results	20
Examine Data Segments	24
Chapter 5 • Cleaning Up Text	31
About the Tasks That You Will Perform	31
Use a Synonym Data Set	32
Create a New Synonym Data Set	34
Use Merged Synonym Data Sets	37
Chapter 6 • Create Topics and Rules	41
About the Tasks That You Will Perform	41
Create Topics	41
Create Rules	43
Chapter 7 • Create Models and Compare Them	45
About the Tasks That You Will Perform	45
Create Models	45
Compare the Models	47
Chapter 8 • The Text Import Node	49

About the Text Import Node	49
Using the Text Import Node	49
Chapter 9 • The Text Parsing Node	53
About the Text Parsing Node	53
Using the Text Parsing Node	53
Chapter 10 • The Text Filter Node	59
About the Text Filter Node	59
Using the Text Filter Node	59
Chapter 11 • The Text Topic Node	65
About the Text Topic Node	65
Using the Text Topic Node	65
Chapter 12 • The Text Cluster Node	71
About the Text Cluster Node	71
Using the Text Cluster Node	71
Chapter 13 • The Text Rule Builder Node	77
About the Text Rule Builder Node	77
Using the Text Rule Builder Node	77
Chapter 14 • Tips for Text Mining	85
Processing a Large Collection of Documents	85
Dealing with Long Documents	85
Processing Documents from an Unsupported Language or Encoding	86
Chapter 15 • Next Steps: A Quick Look at Additional Features	87
The %TEXTSYN Macro	87
The %TMFILTER Macro	87
Glossary	89
Index	93

Overview

Audience

This book is intended for users who are new to SAS Text Miner. The first seven chapters illustrate how you can use SAS Text Miner nodes in the context of a hypothetical text mining analysis. After completing these chapters, you should be able to create projects, process flow diagrams, understand how you can set properties for SAS Text Miner nodes, run them, and explore the results.

Chapters 8 to 13 provide additional examples about the following SAS Text Miner nodes.

- The Text Import Node
- The Text Parsing Node
- The Text Filter Node
- The Text Cluster Node
- The Text Topic Node
- The Text Rule Builder Node

The final two chapters provide some text mining tips and a quick look at additional features.

See the SAS Text Miner help for more information about the text mining process or any of the SAS Text Miner nodes.

Recommended Reading

- *Many of the concepts and topics that are discussed in additional product documentation for SAS Text Miner 12.1 (<http://support.sas.com/documentation/onlinedoc/txtminer>) and SAS Enterprise Miner 12.1 (<http://support.sas.com/documentation/onlinedoc/miner>) might also help you use SAS Text Miner 12.1.*

For a complete list of SAS publications, go to support.sas.com/bookstore. If you have questions about which titles you need, please contact a SAS Publishing Sales Representative:

SAS Publishing Sales
SAS Campus Drive
Cary, NC 27513-2414
Phone: 1-800-727-3228
Fax: 1-919-677-8166
E-mail: sasbook@sas.com
Web address: support.sas.com/bookstore

Chapter 1

Introduction to Text Mining and SAS Text Miner 12.1

What Is Text Mining?	1
What Is SAS Text Miner 12.1?	2
The Text Mining Process	3
Accessibility Features of SAS Text Miner 12.1	4

What Is Text Mining?

Text mining uncovers the underlying themes or concepts that are contained in large document collections. Text mining applications have two phases: exploring the textual data for its content and then using discovered information to improve the existing processes. Both are important and can be referred to as descriptive mining and predictive mining.

Descriptive mining involves discovering the themes and concepts that exist in a textual collection. For example, many companies collect customers' comments from sources that include the Web, e-mail, and contact centers. Mining the textual comments includes providing detailed information about the terms, phrases, and other entities in the textual collection; clustering the documents into meaningful groups; and reporting the concepts that are discovered in the clusters. Results from descriptive mining enable you to better understand the textual collection.

Predictive mining involves classifying the documents into categories and using the information that is implicit in the text for decision making. For example, you might want to identify the customers who ask standard questions so that they receive an automated answer. In addition, you might want to predict whether a customer is likely to buy again, or even if you should spend more effort to keep the customer.

Predictive modeling involves examining past data to predict results. Consider that you have a customer data set that contains information about past buying behaviors, along with customer comments. You could build a predictive model that can be used to score new customers—that is, to analyze new customers based on the data from past customers. For example, if you are a researcher for a pharmaceutical company, you know that hand-coding adverse reactions from doctors' reports in a clinical study is a laborious, error-prone job. Instead, you could create a model by using all your historical textual data, noting which doctors' reports correspond to which adverse reactions. When the model is constructed, processing the textual data can be done automatically by scoring new records that come in. You would just have to examine the "hard-to-classify" examples, and let the computer handle the rest.

Both of these aspects of text mining share some of the same requirements. Namely, textual documents that human beings can easily understand must first be represented in a form that can be mined by the software. The raw documents need processing before the patterns and relationships that they contain can be discovered. Although the human mind comprehends chapters, paragraphs, and sentences, computers require structured (quantitative or qualitative) data. As a result, an unstructured document must be converted into a structured form before it can be mined.

What Is SAS Text Miner 12.1?

SAS Text Miner is a plug-in for the SAS Enterprise Miner environment. SAS Enterprise Miner provides a rich set of data mining tools that facilitate the prediction aspect of text mining. The integration of SAS Text Miner within SAS Enterprise Miner combines textual data with traditional data mining variables. Text mining nodes can be embedded into a SAS Enterprise Miner process flow diagram. SAS Text Miner supports various sources of textual data: local text files, text as observations in SAS data sets or external databases, and files on the Web.

SAS Text Miner 12.1 includes the following nodes that you can use in your text mining analysis:

- **Text Import** node
- **Text Parsing** node
- **Text Filter** node
- **Text Topic** node
- **Text Cluster** node
- **Text Rule Builder** node

For more information about the SAS Text Miner nodes, see the corresponding chapter in this book, or the SAS Text Miner Help.

Together, the Text Miner nodes encompass the parsing and exploration aspects of text mining and the preparation of data for predictive mining and further exploration when you use other SAS Enterprise Miner nodes. You can analyze structured text information, and combine the structured output of the Text Miner nodes with other structured data as desired. The Text Miner nodes are highly customizable and enable you to choose among a variety of options. For example, the **Text Parsing** node enables you to parse documents for detailed information about the terms, phrases, and other entities in the collection. The **Text Cluster** node enables you to cluster documents into meaningful groups and to report concepts that you discover in the clusters. Sorting, searching, filtering (subsetting), and finding similar terms or documents all enhance the exploration process.

SAS Text Miner also enables you to use a SAS macro that is called %TMFILTER. This macro accomplishes a text preprocessing step and enables SAS data sets to be created from documents that reside in your file system or on Web pages. These documents can exist in a number of proprietary formats.

SAS Text Miner is a flexible tool that can solve a variety of problems. Here are some examples of tasks that can be accomplished using SAS Text Miner:

- filtering e-mail
- grouping documents by topic into predefined categories

- routing news items
- clustering analysis of research papers in a database
- clustering analysis of survey data
- clustering analysis of customer complaints and comments
- predicting stock market prices from business news announcements
- predicting customer satisfaction from customer comments
- predicting costs, based on call center logs

The Text Mining Process

Whether you intend to use textual data for descriptive purposes, predictive purposes, or both, the same processing steps take place, as shown in the following table:

Action	Result	Tool
File preprocessing	Creates a single SAS data set from your document collection. The SAS data set is used as input for the Text Parsing node, and might contain the actual text or paths to the actual text.	Text Import node %TMFILTER macro — a SAS macro for extracting text from documents and creating a predefined SAS data set with a text variable
Text parsing	Decomposes textual data and generates a quantitative representation suitable for data mining purposes.	Text Parsing node
Transformation (dimension reduction)	Transforms the quantitative representation into a compact and informative format.	Text Filter node
Document analysis	Performs classification, prediction, or concept linking of the document collection. Creates clusters, topics, or rules from the data.	Text Cluster node Text Topic node Text Rule Builder node SAS Enterprise Miner predictive modeling nodes

Note: The **Text Miner** node is not available from the **Text Mining** tab in SAS Text Miner 12.1. The Text Miner node has now been replaced by the functionality in other SAS Text Miner nodes. You can import diagrams from a previous release of SAS Text Miner that had a **Text Miner** node in the process flow diagram. However, new **Text Miner** nodes can no longer be created, and property values cannot be changed in imported **Text Miner** nodes. For more information, see the Converting SAS Text Miner Diagrams from a Previous Version topic in the SAS Text Miner Help.

Finally, the rules for clustering or predictions can be used to score a new collection of documents at any time.

You might not need to include all of these steps in your analysis, and it might be necessary to try a different combination of options before you are satisfied with the results.

Accessibility Features of SAS Text Miner 12.1

SAS Text Miner includes accessibility and compatibility features that improve usability of the product for users with disabilities. These features are related to accessibility standards for electronic information technology adopted by the U.S. Government under Section 508 of the U.S. Rehabilitation Act of 1973, as amended. SAS Text Miner supports Section 508 standards except as noted in the following table.

Section 508 Accessibility Criterion	Support Status	Explanation
When software is designed to run on a system that has a keyboard, product functions shall be executable from a keyboard where the function itself or the result of performing a function can be discerned textually.	Supported with exceptions.	<p>The software supports keyboard equivalents for all user actions with the exceptions noted below:</p> <p>The keyboard equivalent for exposing the system menu is not the Windows standard Alt + spacebar. The system menu can be exposed using the following shortcut keys: (1) Primary window — Shift + F10 + spacebar, or (2) Secondary window — Shift + F10 + down key.</p> <p>The Explore action in the data source pop-up menu cannot be invoked directly from the keyboard. There is an alternative way to invoke the data source explorer using the View —> Explorer menu.</p>
Color coding shall not be used as the only means of conveying information, indicating an action, prompting a response, or distinguishing a visual element.	Supported with exception.	Node run or failure indication relies on color. There is also a corresponding pop-up message in a dialog box that indicates node success or failure.

If you have questions or concerns about the accessibility of SAS products, send e-mail to accessibility@sas.com.

Chapter 2

Learning by Example: Using SAS Text Miner 12.1

About the Scenario in This Book	5
Prerequisites for This Scenario	7
How to Get Help for SAS Text Miner 12.1	7

About the Scenario in This Book

The first seven chapters describe an extended example that is intended to familiarize you with SAS Text Miner. Each topic builds on the previous topic, so you must work through these chapters in sequence. Several key components of the SAS Text Miner process flow diagram are covered. In this step-by-step example, you learn to do basic tasks in SAS Text Miner, such as creating a project and building a process flow diagram. In your diagram, you perform tasks such as accessing data, preparing the data, building multiple predictive models using text variables, and comparing the models. The extended example in this book is designed to be used in conjunction with SAS Text Miner software. The remaining chapters focus on each of the SAS Text Miner nodes, and provide additional information that you might find useful for your text mining analysis.

The Vaccine Adverse Event Reporting System (VAERS) data is publicly available from the U.S. Department of Health and Human Services (HHS). Anyone can download this data in comma-separated value (CSV) format from <http://vaers.hhs.gov>. There are separate CSV files for every year since the U.S. started collecting the data in 1990. This data is collected from myriad sources, but most reports come from vaccine manufacturers and health care providers. Providers are required to report any contraindicated events for a vaccine or any very serious complications. In the context of a vaccine, a contraindication event would be a condition or a factor that increases the risk of using the vaccine.

See the following in the **Getting Started Examples** zip file:

- ReportableEventsTable.pdf for a complete list of reportable events for each vaccine
- VAERS README file for a data dictionary and list of abbreviations used

Note: See “[Prerequisites for This Scenario](#)” on [page 7](#) for information about where to download the **Getting Started Examples** zip file.

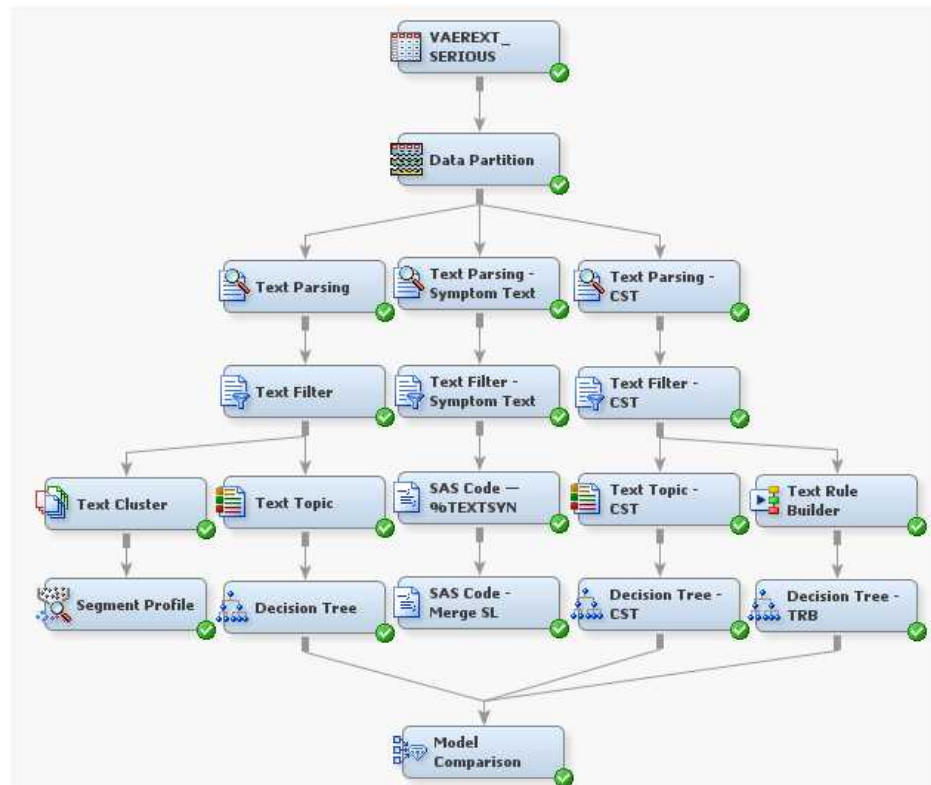
The following figure shows the first 8 columns in the first 10 rows in a table of VAERS data. Included is a unique identifier, the state of residence, and the recipient's age. Additional columns (not in the following figure) include an unstructured text string

SYMPTOM_TEXT that contains the reported problem, specific symptoms, and a symptom counter.

	VAERS_ID	RECVDATE	STATE	AGE_YRS	CAGE_YR	CAGE_MO	SEX	RPT_DATE
1	179605.0	Jan 2, 2002	FL	64.0	64.0	.	F	Dec 26, 2001
2	179606.0	Jan 2, 2002		29.0	29.0	.	F	Dec 26, 2001
3	179612.0	Jan 2, 2002	NJ	40.0	0.0	0.3	F	Dec 23, 2001
4	179613.0	Jan 2, 2002	NY	40.0	1.0	0.6	F	Feb 28, 1998
5	179614.0	Jan 2, 2002	TX	4.0	4.0	.	M	Dec 28, 2001
6	179615.0	Jan 2, 2002	WI	38.0	38.0	.	F	Dec 21, 2001
7	179616.0	Jan 2, 2002	KY	69.0	69.0	.	F	Dec 26, 2001
8	179617.0	Jan 2, 2002	FL	77.0	77.0	.	M	Dec 21, 2001
9	179618.0	Jan 2, 2002		7.0	7.0	.	M	Nov 23, 2001
10	179619.0	Jan 2, 2002		50.0	.	.	F	Dec 20, 2001

As you go through this example, imagine you are a researcher trying to discover what information is contained within this data set. You also want to know how you can use it to better understand the adverse reactions that children and adults are experiencing from their vaccination shots. These adverse reactions might be caused by one or more of the vaccinations that they are given, or they might be induced by an improper procedure from the administering lab (for example, a non-sanitized needle). Some of them will be totally unrelated. For example, perhaps someone happened to get a cold just after receiving a flu vaccine and reported it. You might want to investigate serious reactions that required a hospital stay or caused a lifetime disability or death.

When you are finished with this example, your process flow diagram should resemble the one shown here:



Prerequisites for This Scenario

Before you can perform the tasks in this book, administrators at your site must have installed and configured all necessary components of SAS Text Miner 12.1. You must also perform the following:

1. Download the Getting Started Examples zip file under the SAS Text Miner 12.1 heading from the following URL:

`http://support.sas.com/documentation/onlinedoc/txtminer`

2. Unzip this file into any folder in your file system.
3. Create a folder called **vaersdata** on your C:\ drive.
4. Copy the following files into **C:\vaersdata**:
 - vaerext.sas7bdat
 - vaer_abbrev.sas7bdat
 - engdict.sas7bdat

Note: The preceding list of files might or might not be capitalized depending on which environment you are viewing them in.

How to Get Help for SAS Text Miner 12.1

Select **Help** ⇒ **Contents** from the main SAS Enterprise Miner menu bar to get help for SAS Text Miner.

Chapter 3

Setting Up Your Project

About the Tasks That You Will Perform	9
Create a Project	9
Create a Library	10
Explore and Modify the Data	11
Create a Data Source	12
Create a Diagram	13

About the Tasks That You Will Perform

To set up your project, perform the following main tasks:

1. Create a new project where you will store all your work.
 2. Create a library to store your data sources.
 3. Explore and modify the VAERS data.
 4. Define the SAS Enterprise Miner data source VAEREXT_SERIOUS based on the VAERS data.
 5. Create a new diagram in your project that you will use to interact with nodes.
-

Create a Project

To create a project:

1. Open SAS Enterprise Miner.
2. Click **New Project** in the SAS Enterprise Miner window.
The Select SAS Server page appears.
3. Click **Next**.
The Specify Project Name and Server Directory page appears.
4. Enter a name for the project, such as *Vaccine Adverse Events*, in the **Project Name** field.

5. Enter the path to the location on the server where you want to store data for your project in the **SAS Server Directory** field. Alternatively, browse to a folder to use for your project, or accept the default directory path that appears.

Note: The project path depends on whether you are running SAS Enterprise Miner as a complete client on your local machine or as a client and server application. If you are running SAS Enterprise Miner as a complete client, your local machine acts as its own server. Your SAS Enterprise Miner projects are stored on your local machine in a location that you specify, such as **C:\EM_Projects**. If you are running SAS Enterprise Miner as a client/server application, all projects are stored on the SAS Enterprise Miner server. If you see a default path in the SAS Server Directory box, you can accept the default project path, or you can specify your own project path.

6. Click **Next**.

The Register the Project page appears.

7. Click **Next**.

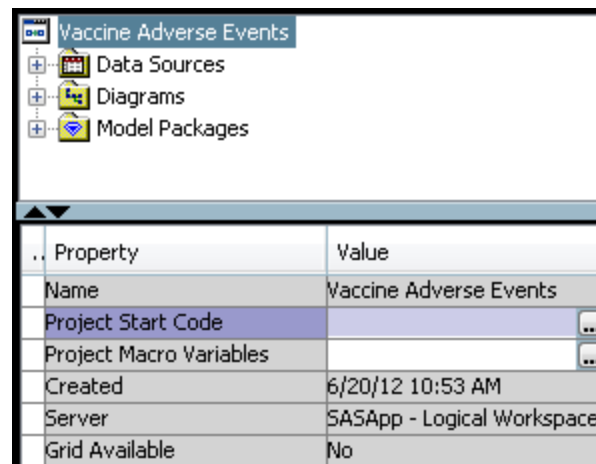
The New Project Information page appears.

8. Click **Finish** to create your project.

Create a Library

To create a library:

1. Select the project name **Vaccine Adverse Events** to display the project Properties Panel.



2. Click  for the **Project Start Code** property.

The Project Start Code dialog box appears.

3. Enter the following code on the **Code** tab:

```
libname mylib "c:\vaersdata";
```

Note: The location depends on where you have stored the data for this tutorial on your system.

4. Click **Run Now**.

- Click **OK** to close the Project Start Code dialog box.

Note: An alternate way to create a library is to use the library wizard. To use the library wizard, select **File** ⇒ **New** ⇒ **Library** from the main menu.

Explore and Modify the Data

After you create a library, you can view the data before creating a data source to use in SAS Enterprise Miner. For example, consider the situation where you want to create another variable to summarize the values of multiple variables.

Perform the following steps to view available SAS data files and create a new SAS data file that you will use as a data source.

- Select **View** ⇒ **Explorer** from the main menu.

The Explorer window appears.

- Select **Mylib** in the SAS Libraries tree.

The contents of the Mylib library appear, and include the three files **Engdict**, **Vaerext**, and **Vaer_abbrev**.

- Double-click **Vaerext**.


The contents of the **Vaerext** file appears in a new window.

- Scroll to the right to view the available variable names that appear as column headings.

Notice the following variables:

- **DISABLE** — a binary variable that has a value of 'Y' if there was a disability
- **DIED** — a binary variable that has a value of 'Y' if there was a death
- **ER_VISIT** — a binary variable that has a value of 'Y' if there was an emergency room visit
- **HOSPITAL** — a binary variable that has a value of 'Y' if there was a hospitalization

Consider that we want to create a new data set, **vaerext_serious**, that includes a binary variable **serious** that has a value of 'Y' if there was disability, death, emergency room visit, or hospitalization, and a value of 'N' otherwise.

- Close the **MYLIB.VAEREXT** window.
- Select the project name **Vaccine Adverse Events** to display the project Properties Panel.
- Click  for the **Project Start Code** property.

The Project Start Code dialog box appears.

- Enter the following code on the **Code** tab after the LIBNAME statement you previously added to create the mylib library:

```
data mylib.vaerext_serious;
  set mylib.vaerext;
  if DISABLE='Y' or DIED='Y' or ER_VISIT='Y' or HOSPITAL='Y' then serious='Y';
  else serious='N';
run;
```

This code creates a new SAS file, **vaerext_serious** from the **vaerext** file in the **mylib** library, adds a variable **serious**, and assigns it a value of **Y** or **N**, depending on the value of the **DISABLE**, **DIED**, **ER_VISIT**, and **HOSPITAL** variables.

9. Click **Run Now**.
10. Click **OK**.
11. Select **View** ⇒ **Explorer** from the main menu.

The Explorer window appears.

12. Select **Mylib** in the SAS Libraries tree.

Notice that the Mylib library now contains a new entry for the file **Vaerext_serious**.

13. Double-click **Vaerext_serious**.

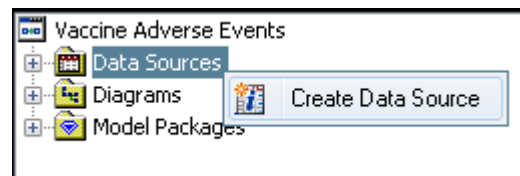
The contents of the **Vaerext_serious** file appears in a new window.

14. Scroll all the way to the right to see the new column **serious**.
15. Close the **MYLIB.VAEREXT_SERIOUS** window.
16. Close the Explorer window.

Create a Data Source

To create a data source:

1. Right-click the Data Sources folder in the Project Panel and select **Create Data Source**.



The Data Source wizard appears.

2. Select **SAS Table** in the Source drop-down menu.
3. Click **Next**.

The Select a SAS Table window appears.

4. Click **Browse**.
5. Click the SAS library named **Mylib** in the library tree.

The Mylib library folder contents are displayed on the Select a SAS Table dialog box.

Note: If you do not see SAS data files in the Mylib folder, click **Refresh**.

6. Select the **Vaerext_serious** table.
7. Click **OK**.

The two-level name **MYLIB.VAEREXT_SERIOUS** is displayed in the **Table** field.

8. Click **Next**.

The Table Information page appears.

9. Click **Next**.

The Metadata Advisor Options page appears.

10. Select **Advanced**.11. Click **Next**.

The Column Metadata page appears.

12. Select the following variable roles by clicking the role value for each variable value and selecting the indicated value from the drop-down list.

- Set the role for **V_ADMINBY** to **Input**.
- Set the role for **V_FUNDBY** to **Input**.
- Set the role for **serious** to **Target**.

13. Click **Next**.

The Decision Configuration page appears.

14. Click **Next**.

The Create Sample page appears.

15. Click **Next**.

The Data Source Attributes page appears.

16. Click **Next**.

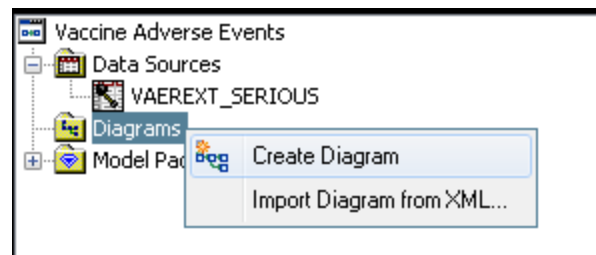
The Summary page appears.

17. Click **Finish**.

The VAEREXT_SERIOUS table is added to the Data Sources folder in the Project Panel.

Create a Diagram

To create a diagram, complete the following steps:

1. Right-click the Diagrams folder in the Project Panel and select **Create Diagram**.

The Create New Diagram dialog box appears.

2. Type *VAERS Example* in the **Diagram Name** field.3. Click **OK**.

The empty VAERS Example diagram opens in the diagram workspace.

Chapter 4

Analyzing the SYMPTOM_TEXT Variable

About the Tasks That You Will Perform	15
Identify Input Data	15
Partition Input Data	16
Parse Data	17
Filter Data	18
Cluster Data	19
View Results	20
Examine Data Segments	24

About the Tasks That You Will Perform

The SYMPTOM_TEXT variable contains the text of an adverse event as it was reported. This chapter explains how you can analyze the SYMPTOM_TEXT variable by performing the following tasks:

1. Identify the VAERS_SERIOUS data source with an **Input Data** node.
 2. Partition the input data using the **Data Partition** node.
 3. Parse the document collection using the **Text Parsing** node.
 4. Reduce the total number of parsed terms using the **Text Filter** node.
 5. Cluster documents using the **Text Cluster** node.
 6. View the results.
 7. Examine data segments using the **Segment Profile** node.
-

Identify Input Data

To identify input data:

1. Select the **VAEREXT_SERIOUS** data source from the Data Sources folder in the Project Panel.

2. Drag **VAEREXT_SERIOUS** into the diagram workspace to create an **Input Data** node.

Partition Input Data

The **Data Partition** node enables you to partition your input data into one of the following data sets:

- **Train** — used for preliminary model fitting. The analyst attempts to find the best model weights by using this data set.
- **Validation** — used to assess the adequacy of the model in the Model Comparison node. The validation data set is also used for model fine-tuning in the **Decision Tree** model node to create the best subtree.
- **Test** — used to obtain a final, unbiased estimate of the generalization error of the model.

For more information about the **Data Partition** node, see the SAS Enterprise Miner Help.

Perform the following steps to add a **Data Partition** node to the analysis:

1. Select the **Sample** tab on the node toolbar and drag a **Data Partition** node into the diagram workspace.
2. Connect the **VAEREXT_SERIOUS** input data node to the **Data Partition** node.

Note: To connect one node to another node in the default horizontal view, position the mouse pointer at the right edge of a node. A pencil icon appears. Hold the left mouse button down, and drag the line to the left edge of the node that you want to connect to, and then release the left mouse button. To change your view of connected nodes to a vertical view, right-click in the diagram workspace, and select **Layout** ⇒ **Vertically** in the menu that appears.



3. Select the **Data Partition** node to view its properties.
Details about the node appear in the Properties Panel.
4. Set the Data Set Allocations properties as follows:
 - Set the **Training** property to **60.0**.
 - Set the **Validation** property to **20.0**.
 - Set the **Test** property to **20.0**.

These data partition settings ensure adequate data when you build prediction models with the **VAEREXT_SERIOUS** data.

Parse Data

The **Text Parsing** node enables you to parse a document collection in order to quantify information about the terms that are contained therein. You can use the **Text Parsing** node with volumes of textual data such as e-mail messages, news articles, Web pages, research papers, and surveys. For more information about the **Text Parsing** node, see the SAS Text Miner Help.

Perform the following steps to add a **Text Parsing** node to the analysis:

1. Select the **Text Mining** tab on the node toolbar, and drag a **Text Parsing** node into the diagram workspace.
2. Connect the **Data Partition** node to the **Text Parsing** node.




3. Select the **Text Parsing** node.

The properties for the **Text Parsing** node appear in the Properties Panel.

4. Set the **Different Parts of Speech** property value to **No**.

For the VAERS data, this setting offers a more compact set of terms.

5. Click the  for the **Synonyms** property.

A dialog box appears.

6. Click **Import**.

The Select a SAS Table dialog box appears.

7. Select **No data set to be specified**.

8. Click **OK** to exit the Select a SAS Table dialog box.

9. Click **OK** to exit the Synonyms dialog box.

10. Click the  for the **Ignore Parts of Speech** property.

The Ignore Parts of Speech dialog box appears.

11. Select the following items, which represent parts of speech:

- Aux
- Conj
- Det

- Interj
- Part
- Prep
- Pron
- Num

Note: Hold down the CTRL key to select more than one.

Any terms with the parts of speech that you select in the Ignore Parts of Speech dialog box are ignored during parsing. The selections indicated here ensure that the analysis ignores low-content words such as prepositions and determiners.

12. Click **OK**.

Filter Data

The **Text Filter** node can be used to reduce the total number of parsed terms or documents that are analyzed. Therefore, you can eliminate extraneous information so that only the most valuable and relevant information is considered. For example, the **Text Filter** node can be used to remove unwanted terms and to keep only documents that discuss a particular issue. This reduced data set can be orders of magnitude smaller than the one that represents the original collection, which might contain hundreds of thousands of documents and hundreds of thousands of distinct terms. For more information about the **Text Filter** node, see the SAS Text Miner help.

To filter the data:

1. Select the **Text Mining** tab on the node toolbar, and drag a **Text Filter** node into the diagram workspace.
2. Connect the **Text Parsing** node to the **Text Filter** node.



3. Select the **Text Filter** node.
4. Set the value of the **Term Weight** property to **Mutual Information**.

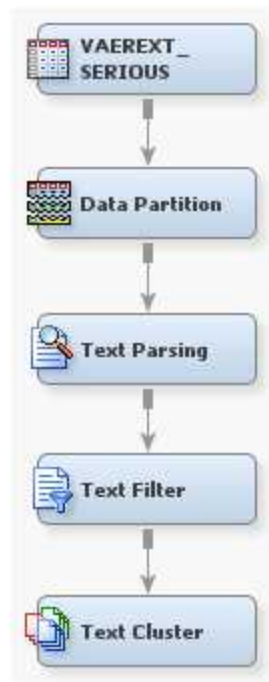
This causes the terms to be differentially weighted when they correspond to serious reactions.

Cluster Data

The **Text Cluster** node clusters documents into disjointed sets of documents and reports on the descriptive terms for those clusters. Two algorithms are available. The Expectation Maximization algorithm clusters documents with a flat representation, and the Hierarchical clustering algorithm groups clusters into a tree hierarchy. Both approaches rely on the singular value decomposition (SVD) to transform the original weighted, term-document frequency matrix into a dense but low dimensional representation. For more information about the **Text Cluster** node, see the SAS Text Miner help.

To cluster the data:

1. Select the **Text Mining** tab on the node toolbar, and drag a **Text Cluster** node into the diagram workspace.
2. Connect the **Text Filter** node to the **Text Cluster** node.



3. Select the **Text Cluster** node.
4. Set the **Descriptive Terms** to **12** to ease cluster labeling.
5. Right-click the **Text Cluster** node in the diagram workspace, and select **Run**.
6. Click **Yes** in the Confirmation dialog box when you are asked whether you want to run the path.
7. Click **OK** in the Run Status dialog box that appears after the **Text Cluster** node has finished running.

View Results

After the process flow diagram has completed running, you can view the results that were obtained by each node.

1. Select the **Text Parsing** node.

The Properties for the **Text Parsing** node appear in the Properties Panel.

Notice that the **Text Parsing** node's **Parse Variable** property has been populated with the SYMPTOM_TEXT variable. This is because the SYMPTOM_TEXT variable was the longest variable with a role of **Text** in the VAEREXT_SERIOUS input data source.

2. Right-click the **Text Parsing** node and select **Results**.

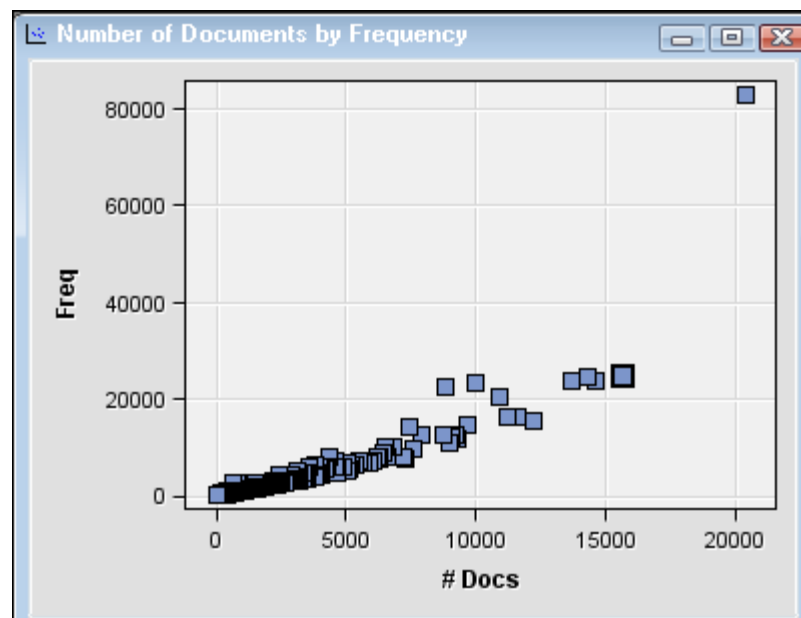
The Results window appears for the **Text Parsing** node.

3. Select the Terms window.

4. Click the **Freq** column heading to sort the terms by frequency.

Scroll through the list of terms. Notice that for each term, the Terms window provides the number of documents the term appeared in, the frequency of the term, and whether the term was kept.

5. Select a term. Notice that the point corresponding to this term is selected in the ZIPF Plot and the Number of Documents by Frequency plot.



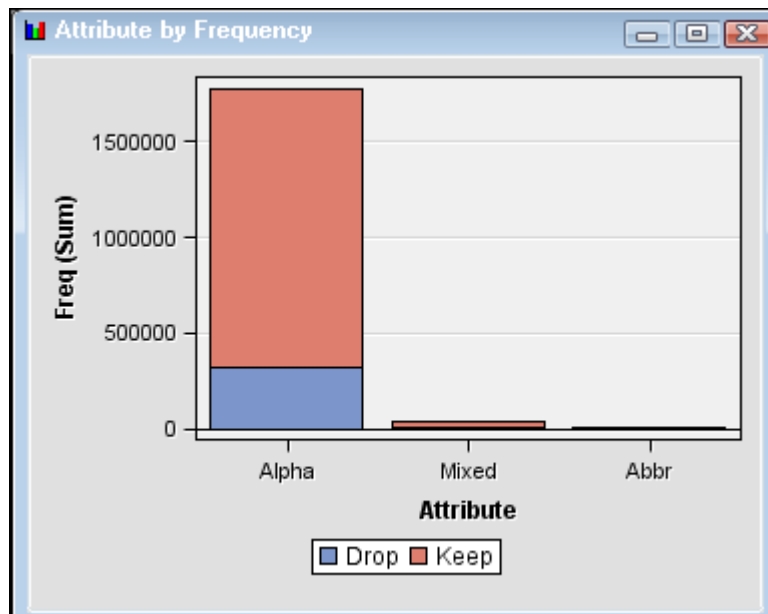
6. Close the Results window.

7. Select the **Text Filter** node.

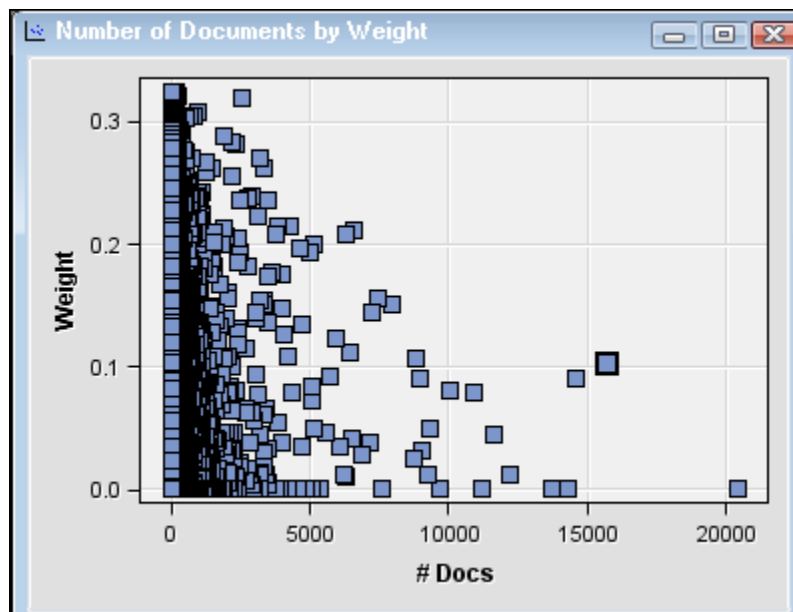
8. Right-click the **Text Filter** node and select **Results**.


The Results window appears for the **Text Filter** node.

Notice that the Attribute by Frequency window and the Role by Freq window now show the number of terms in each category that were dropped or kept.




The Number of Documents by Weight plot shows the number of documents in which each term appears relative to each term's weight.



9. Close the Results window.
10. Click the  for the **Filter Viewer** property.
The Interactive Filter Viewer window appears.
11. View the terms in the Terms window. The terms are sorted first by their keep status and then by the number of documents that they appear in.
Note: You can change the sorted order by clicking a column heading.
12. View the documents in the Documents window.
13. Right-click a cell in the **SYMPTOM_TEXT** column, and then select **Toggle Show Full Text** to see the full text contained in SYMPTOM_TEXT.

14. Select a term that is related to an adverse reaction that you want to investigate further. For example, select **fever** under the TERM column of the Terms window. Right-click on the term and select **Add Term to Search Expression**.

Terms				
	TERM	FREQ	# DOCS	KEEP ▼
+	receive	24917	15632	<input checked="" type="checkbox"/>
+	vaccine	23886	14602	<input checked="" type="checkbox"/>
+	swell	15564	12209	<input checked="" type="checkbox"/>
+	day	16467	11647	<input checked="" type="checkbox"/>
+	arm	16366	11209	<input checked="" type="checkbox"/>
+	report	20322	10932	<input checked="" type="checkbox"/>
	pt	23389	10037	<input checked="" type="checkbox"/>
+	fever	11752	9324	<input checked="" type="checkbox"/>
+	leave	Add Term to Search Expression		
+	site			
			Treat as Synonyms	

15. Click **Apply**.

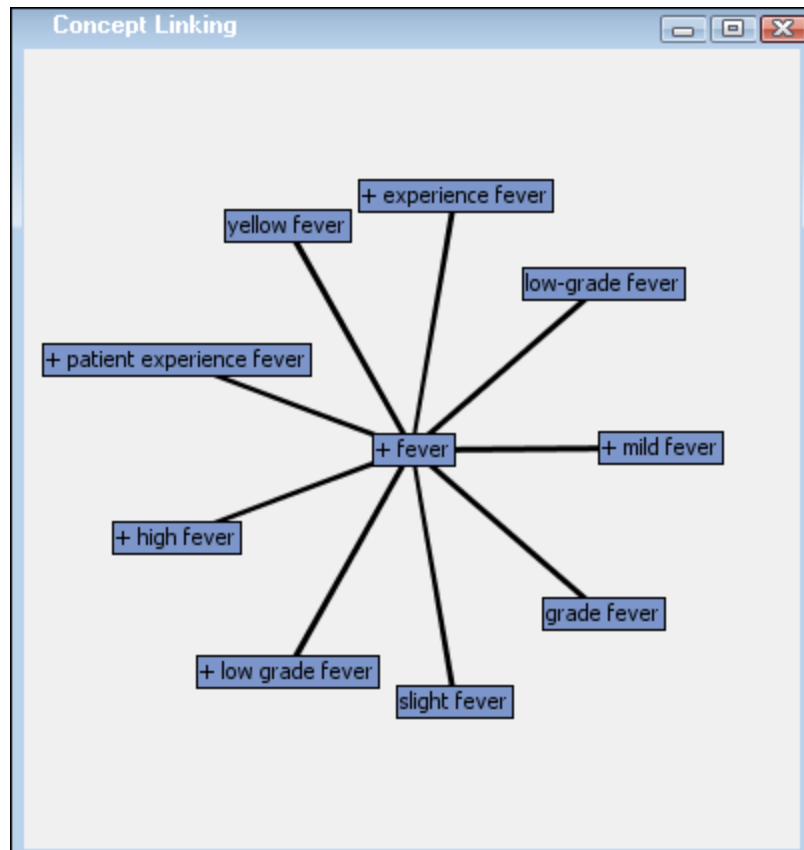
Notice that the Documents window updates to only include entries that contain the term **fever**.

16. Click **Clear**, and then click **Apply**.

The terms in the Terms window resets.

17. Select the term **fever** in the Terms window, and then right-click it and select **View Concept Links**.

The Concept Linking window appears. Concept linking is a way to find and display the terms that are highly associated with the selected term in the Terms table. The selected term is surrounded by the terms that correlate the strongest with it. The Concept Linking window shows a hyperbolic tree graph with **fever** in the center of the tree structure. It shows you the other terms that are strongly associated with the term **fever**.



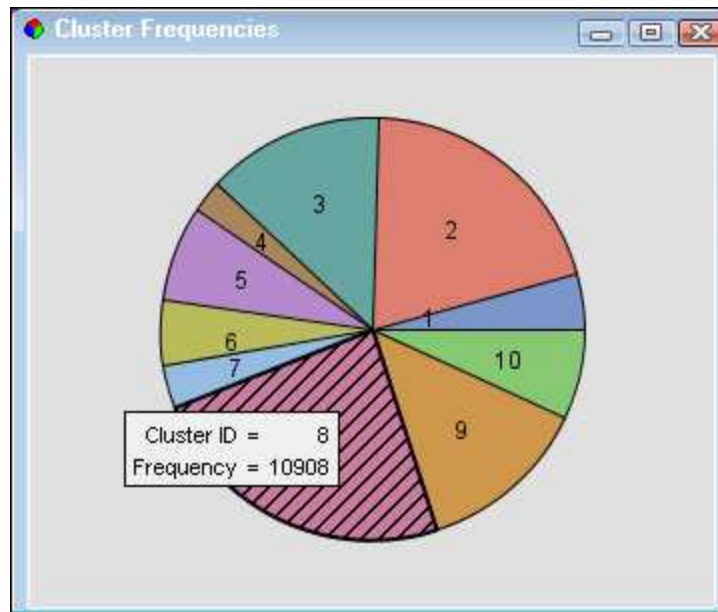
To expand the Concept Linking view, select a term that is not in the center of the graph, right-click it, and select **Expand Links**.

18. Close the Results window.
19. Select the **Text Cluster** node.
20. Right-click the **Text Cluster** node and select **Results**.

The Results window appears.

21. View the clusters in the Clusters window. Select a cluster.

Notice how the corresponding cluster is selected in the Cluster Frequencies chart, the Cluster Frequency by RMS plot, and the Distance Between Clusters plot.



22. Close the Results window.


Examine Data Segments

In this section, you will examine segmented or clustered data using the **Segment Profile** node. A segment is a cluster number that you derive analytically by using SAS Text Miner clustering techniques. The **Segment Profile** node enables you to get a better idea of what makes each segment unique or at least different from the population. The node generates various reports that aid in exploring and comparing the distribution of these factors within the segments and population. For more information about the **Segment Profile** node, see the SAS Enterprise Miner Help.

To examine data segments, complete the following steps:

1. Select the **Assess** tab on the node toolbar, and drag a **Segment Profile** node into the diagram workspace.
2. Connect the **Text Cluster** node to the **Segment Profile** node.

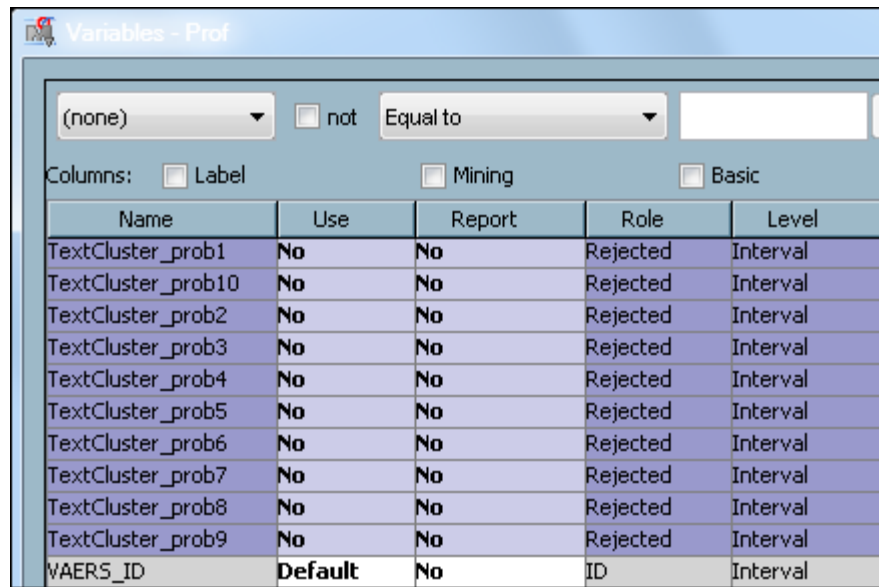


3. Select the **Segment Profile** node.
4. Click the  for the **Variables** property.

The Variables window appears.

5. Select all the “_prob” variables and set their Use value to **No**.

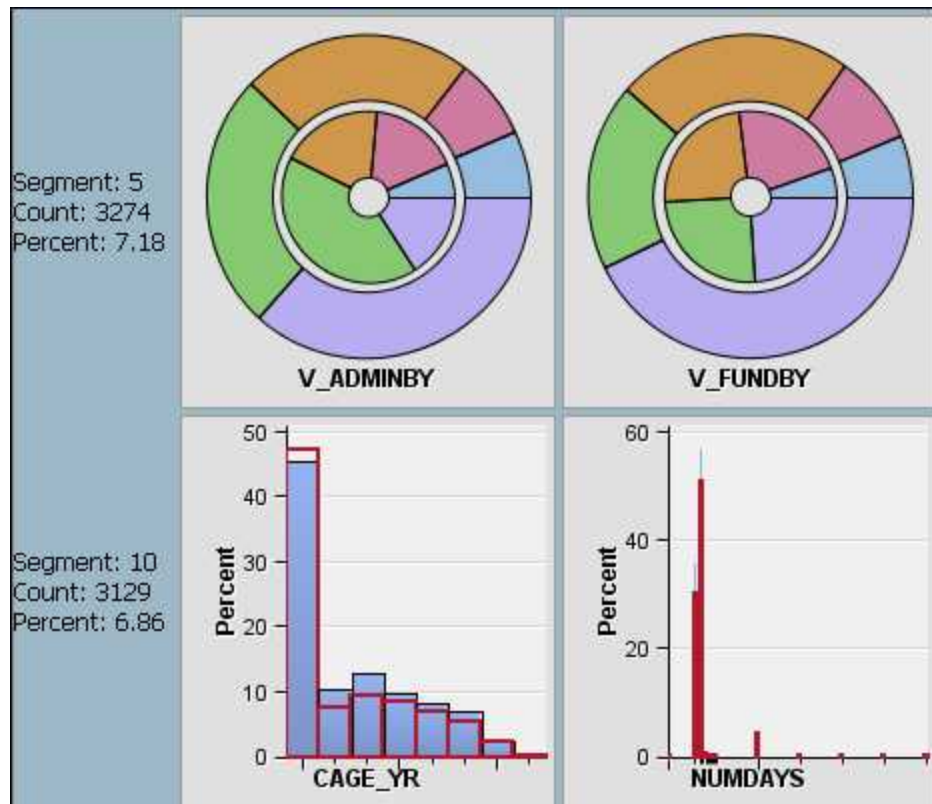
Note: You can hold down Shift and select all the “_prob” variables by clicking on the first “_prob” variable and dragging the pointer to select all “_prob” variables. After all “_prob” variables are selected, you can change the Use value of each selected “_prob” variable by changing the Use value of one of the “_prob” variables. This will change the other “_prob” Use values to the selected value as well.



Name	Use	Report	Role	Level
TextCluster_prob1	No	No	Rejected	Interval
TextCluster_prob10	No	No	Rejected	Interval
TextCluster_prob2	No	No	Rejected	Interval
TextCluster_prob3	No	No	Rejected	Interval
TextCluster_prob4	No	No	Rejected	Interval
TextCluster_prob5	No	No	Rejected	Interval
TextCluster_prob6	No	No	Rejected	Interval
TextCluster_prob7	No	No	Rejected	Interval
TextCluster_prob8	No	No	Rejected	Interval
TextCluster_prob9	No	No	Rejected	Interval
VAERS_ID	Default	No	ID	Interval

6. Select all the “_SVD” variables and set their **Use** value to **No**.
7. Click **OK**.
8. Select the **Segment Profile** node in the diagram workspace.
9. Enter *0.0010* as the value for the **Minimum Worth** property.
10. Right-click the **Segment Profile** node, and select **Run**.
11. Click **Yes** in the Confirmation dialog box when you are asked whether you want to run the path.
12. After the node finishes running, click **Results** in the Run Status dialog box.
13. Maximize the Profile window.

The following shows a portion of this window.

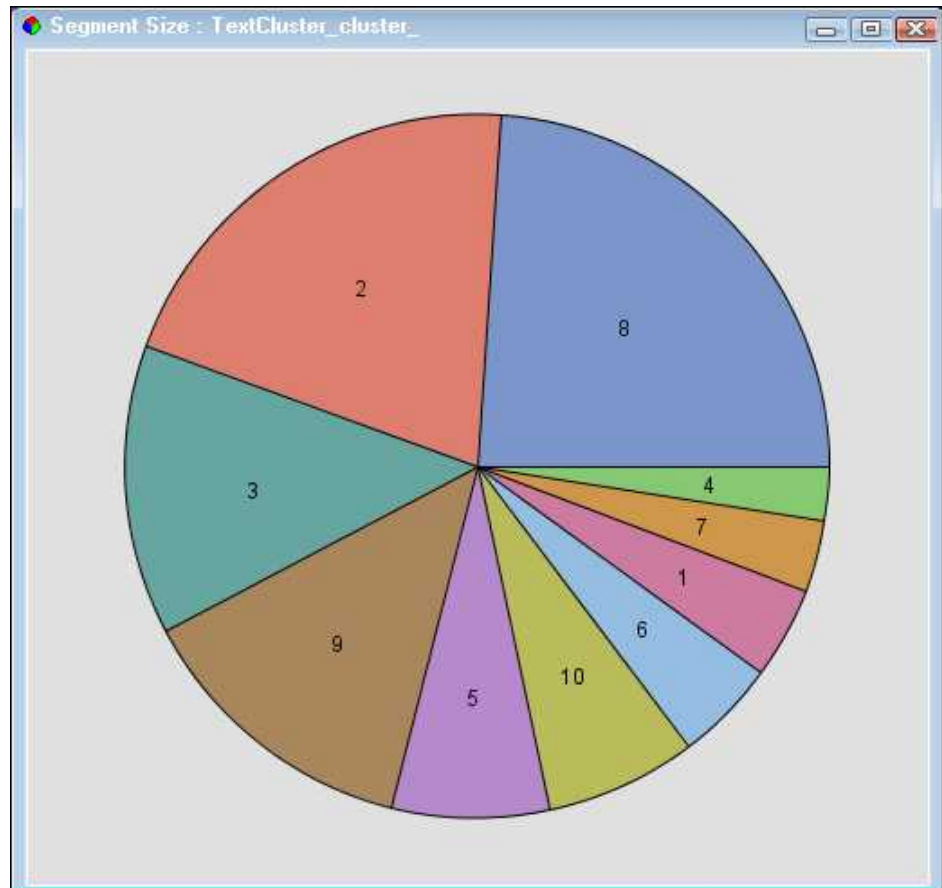


The Profile window displays a lattice, or grid, of plots that compare the distribution for the identified and report variables for both the segment and the population. The graphs shown in this window illustrate variables that have been identified as factors that distinguish the segment from the population that it represents. Each row represents a single segment. The far-left margin identifies the segment, its count, and the percentage of the total population.

The columns are organized from left to right according to their ability to discriminate that segment from the population. Report variables, if specified, appear on the right in alphabetical order after the selected inputs. The lattice graph has the following features:

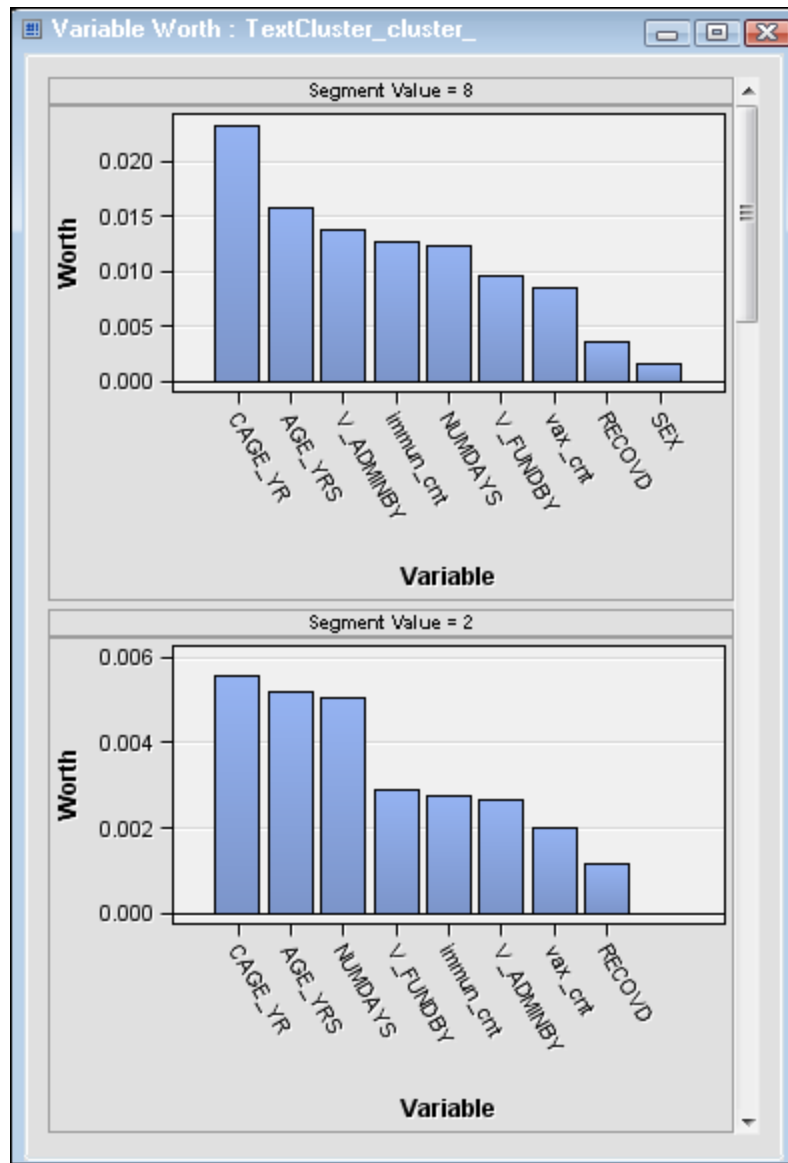
- **Class variable** — is displayed as two nested pie charts that consist of two concentric rings. The inner ring represents the distribution of the total population. The outer ring represents the distribution for the given segment.
- **Interval variable** — is displayed as a histogram. The blue shaded region represents the within-segment distribution. The red outline represents the population distribution. The height of the histogram bars can be scaled by count or by percentage of the segment population. When you are using the percentage, the view shows the relative difference between the segment and the population. When you are using the count, the view shows the absolute difference between the segment and the population.

14. Maximize the Segment Size chart.



15. Maximize the Variable Worth window.

The following shows a portion of this window.



16. Close the **Results** window.

Chapter 5

Cleaning Up Text

About the Tasks That You Will Perform	31
Use a Synonym Data Set	32
Create a New Synonym Data Set	34
Use Merged Synonym Data Sets	37

About the Tasks That You Will Perform

As demonstrated in the previous chapter, SAS Text Miner does a good job of finding themes that are clear in the data. But, when the data needs cleaning, SAS Text Miner can be less effective at uncovering useful themes. In this chapter, you will encounter manually edited data that contains many misspellings and abbreviations, and you will work on cleaning the data to get better results.

The README.TXT file provided in the Getting Started with SAS Text Miner 12.1 zip file contains a list of abbreviations that are commonly used in the adverse event reports. SAS Text Miner enables you to specify a synonym list. A VAER_ABBREV synonym list is provided for you in the Getting Started with SAS Text Miner 12.1 zip file. So that you can create such a synonym list, the abbreviations list from README.TXT was copied into a Microsoft Excel file. The list was manually edited in the Microsoft Excel file and then imported into a SAS data set. For example, CT was marked as equivalent to computerized axial tomography.

For more information about importing data into a SAS data set, see the following documentation resource:

<http://support.sas.com/documentation/>

You will perform the following tasks to clean the text and examine the results:

1. Use a synonym data set from the Getting Started with SAS Text Miner 12.1 zip file.
2. Create a new synonym data set by using the **SAS Code** node and the %TEXTSYN macro. The %TEXTSYN macro will run through all the terms, automatically identify which ones are misspellings, and create synonyms that map correctly spelled terms to the misspelled terms.
3. Examine results using merged synonym data sets.

Use a Synonym Data Set

To use a synonym data set:

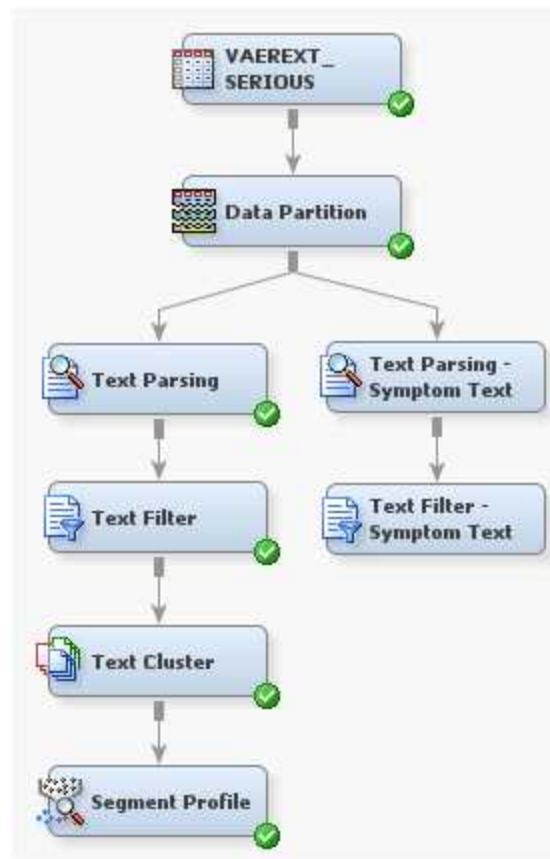
1. Right-click the **Text Parsing** node, and select **Copy**.



For this example, it is important to copy the node instead of creating a new **Text Parsing** node because the settings you previously specified in the **Text Parsing** node properties panel will be used.

2. Right-click in the empty diagram workspace, and select **Paste**.
3. To distinguish this newly pasted **Text Parsing** node from the first node, right-click it, and select **Rename**.
4. Enter *Text Parsing — Symptom Text* in the **Node Name** field, and then click **OK**.
5. Right-click the **Text Filter** node, and select **Copy**.

For this example, it is important to copy the node instead of creating a new **Text Filter** node because the settings you previously specified in the **Text Filter** node properties panel will be used.

6. Right-click in the empty diagram workspace, and select **Paste**.
7. To distinguish this newly pasted **Text Filter** node from the first node, right-click it, and select **Rename**.
8. Enter *Text Filter — Symptom Text* in the **Node Name** field, and then click **OK**.
9. Connect the **Data Partition** node to the **Text Parsing — Symptom Text** node.
10. Connect the **Text Parsing — Symptom Text** node to the **Text Filter — Symptom Text** node.



11. Select the **Text Parsing — Symptom Text** node.
12. Select the  for the **Synonyms** property.
A dialog box appears.
13. Click **Import**.
The Select a SAS Table dialog box appears.
14. Select the **Mylib** library in the folder tree.
The contents of the **Mylib** library appear.
15. Select **Vaer_abbrev**, and then click **OK**.
The contents of the **Vaer_abbrev** data source appear in a dialog box.
16. Click **OK**.
Leave all other settings the same as in the original **Text Parsing** node.
17. Right-click the **Text Filter — Symptom Text** node in the diagram workspace, and select **Run**. Click **Yes** in the Confirmation dialog box.
18. Click **OK** in the Run Status dialog box when the node has finished running.
19. Click the  for the **Filter Viewer** property of the **Text Filter — Symptom Text** node.
The Interactive Filter Viewer window appears.
20. Click the **TERM** column heading to sort the Terms table.
21. Select **abdomen** under the **TERM** column in the Terms window.

You might need to scroll down to see the term. In the Terms window, there should be a plus (+) sign next to **abdomen**. Click on the plus sign to expand the term. This shows all synonyms and stems that are mapped to that term. A stem is the root form of a term. The child term **abd** is included. Both **abdomen** and **abd** will be treated the same.

Terms						
	TERM ▲	FREQ	# DOCS	KEEP	WEIGHT	ROLE
	abdome	1	1	<input type="checkbox"/>	0.0	
<input checked="" type="checkbox"/>	abdomen	661	600	<input checked="" type="checkbox"/>	0.041	
<input type="checkbox"/>	abdomen	561	513			
<input type="checkbox"/>	abd	100	93			
	abdomen area	1	1	<input type="checkbox"/>	0.0	Noun Group

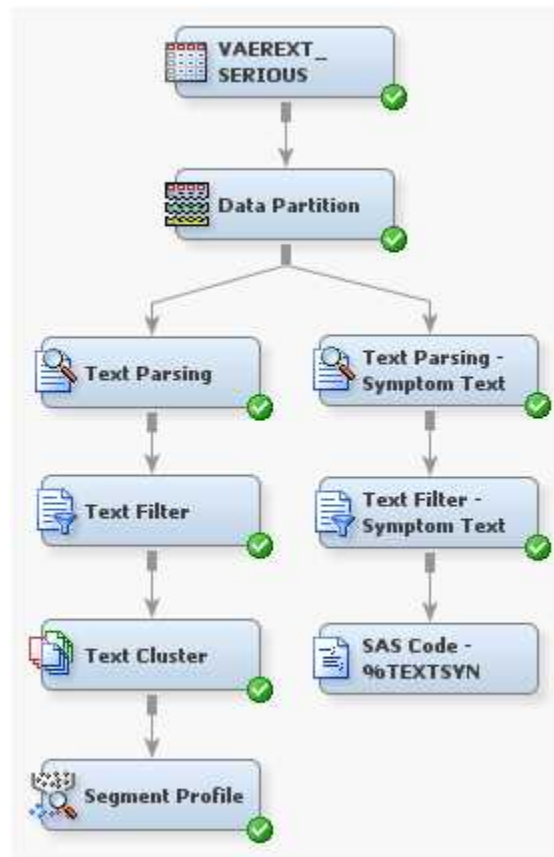
22. Close the Interactive Filter Viewer window.


Create a New Synonym Data Set

You can use the SAS Text Miner %TEXTSYN macro to create a new synonym data set. The %TEXTSYN macro evaluates all the terms, automatically identifies which terms are misspellings, and creates synonyms that map correctly spelled terms to misspelled terms.

To create a new synonym data set:

1. Select the **Utility** tab on the node toolbar and drag a **SAS Code** node into the diagram workspace.
2. Right-click the **SAS Code** node, and select **Rename**.
3. Enter *SAS Code — %TEXTSYN* in the **Node Name** field, and then click **OK**.
4. Connect the **Text Filter — Symptom Text** node to the **SAS Code — %TEXTSYN** node.



5. Select the **SAS Code — %TEXTSYN** node, and then click the  for the **Code Editor** property in the Properties Panel.

The Code Editor window appears.

6. Enter the following code in the Code Editor:



```

%textsyn( termds=<libref>.<nodeID>_terms
          , docds=&em_import_data
          , outds=&em_import_transaction
          , textvar=symptom_text
          , mnpardoc=8
          , mxchddoc=10
          , synds=mylib.vaerextsyms
          , dict=mylib.engdict
          , maxsped=15
          ) ;

```

Note: You will need to replace **<libref>** and **<nodeID>** in the first line in the above code with the correct library name and node ID. To determine what these values are, close the Code Editor window, and then select the arrow that connects the **Text Filter — Symptom Text** node to the **SAS Code — %TEXTSYN** node. The value for **<libref>** will be the first part of the table name that appears in the Properties panel, such as **emws**, **emws2**, and so on. The node ID will appear after the value for **<libref>**, and will be **TextFilter**, **TextFilter2**, and so on. After you determine the value for **<libref>** and **<nodeID>**, a possible first line might be **termds=emws2.textfilter2_terms**. Your libref and node ID values could differ depending on how many **Text Filter** nodes and diagrams have been created in your workspace.

For details about the %TEXTSYN macro, see SAS Text Miner Help documentation.

7. After you have added the %TEXTSYN macro code to the Code Editor window, and modified it to add values for <libref> and <nodeID>, click the  to save the changes.
8. Click the  to run the **SAS Code — %TEXTSYN** node.
9. Click **Yes** in the Confirmation dialog box.
10. Click **OK** in the dialog box that indicates that the node has finished running.
11. Close the **Code Editor** window.
12. Select **View** ⇒ **Explorer** from the main menu.

The Explorer window appears.

13. Click **Mylib** in the SAS Libraries tree, and then select **Vaerextsyms**.

Note: If the **Mylib** library is already selected and you do not see the Vaerextsyms data set, you might need to click **Show Project Data** or refresh the Explorer window to see the **Vaerextsyms** data set.

14. Double-click **Vaerextsyms** to see its contents.

	example1	example2	Term	parent
116	... 4 days following. Heat, !!anti-inflammato ries!! 11/18/2002, Muscle relaxant...	... over the counter non steroidal !!anti-inflam...	anti-inflammatori es	anti-inflammat ory
117	... TO TOUCH. WAS GIVEN !!ANTIBIOTICS!! IF SWELLING INCREASES.	... where they gave her !!antibiotics!!/steroids.	antibiotics	antibiotics
118	... erythematous papules over B/L !!anticubi tal!! spaces, elbows RLQABD, Axilla,...	... hives on face, neck, ! !anticubital!! and right fo ot and...	anticubital	antecubital

Here is a list of what the Vaerextsyms columns provide:

- **Term** is the misspelled word.
- **parent** is a guess at the word that was meant.
- **example1** and **example2** are two examples of the term in a document.
- **childndocs** is the number of documents that contained that term.
- **numdocs** is the number of documents that contained the parent.
- **minsped** is an indication of how close the terms are.
- **dict** indicates whether the term is a legitimate English word. Legitimate words can still be deemed misspellings, but only if they occur rarely and are very close in spelling to a frequent target term.

For example, Observation 117 shows **antibiotics** to be a misspelling of **antibiotics**. Four documents contain **antibiotics**, and 745 documents contain the parent. Note that double exclamation marks (!!) both precede and follow the child term in the example text so that you can see the term in context.

15. Examine the Vaerextsyms table to see whether you disagree with some of the choices made. For this example, however, assume that the %TEXTSYN macro has done a good enough job of detecting misspellings.

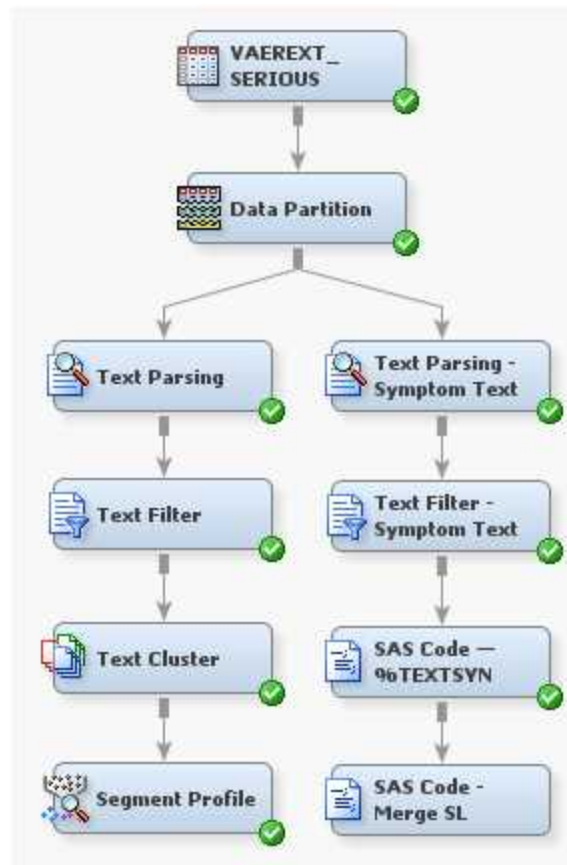
Note: The Vaerextsyms table can be edited using any SAS table editor. You cannot edit this table in the SAS Enterprise Miner GUI. You can change a parent for any misspellings that appear incorrect or delete a row if the Term column contains a valid term.


16. Close the Mylib.Vaerextsyms table and the Explorer window.

Use Merged Synonym Data Sets

In this set of tasks, you will create a new data set that contains all the observations from both the Mylib.Vaerextsyms and Mylib.Vaer_abbrev data sets. You will examine the results by using the merged synonym data set. Complete the following steps:


1. Select the **Utility** tab on the node toolbar and drag a **SAS Code** node into the diagram workspace.
2. Right-click the **SAS Code** node, and select **Rename**.
3. Enter *SAS Code — Merge SL* in the **Node Name** field, and then click **OK**.
SL stands for *Synonym Lists*.
4. Connect the **SAS Code — %TEXTSYN** node to the **SAS Code — Merge SL** node.



5. Select the **SAS Code — Merge SL** node.
6. Click the  for the **Code Editor** property.
The Code Editor appears.
7. Enter the following code in the Code Editor:

```
data mylib.vaerextsyms_new;
    set mylib.vaerextsyms mylib.vaer_abbrev;
run;
```

This code merges the resulting synonyms data set from the first **SAS Code** — **%TEXTSYN** node with the abbreviations data set.

8. Click .
9. Close the Code Editor window.
10. Right-click the **SAS Code** — **Merge SL** node, and select **Run**. Click **Yes** in the Confirmation dialog box.
11. Click **Results** in the Run Status dialog box when the node has finished running.
12. From the Results window, select **View** ⇒ **SAS Results** ⇒ **Log** to see the SAS code where the new data set is created.

Close the Results window.

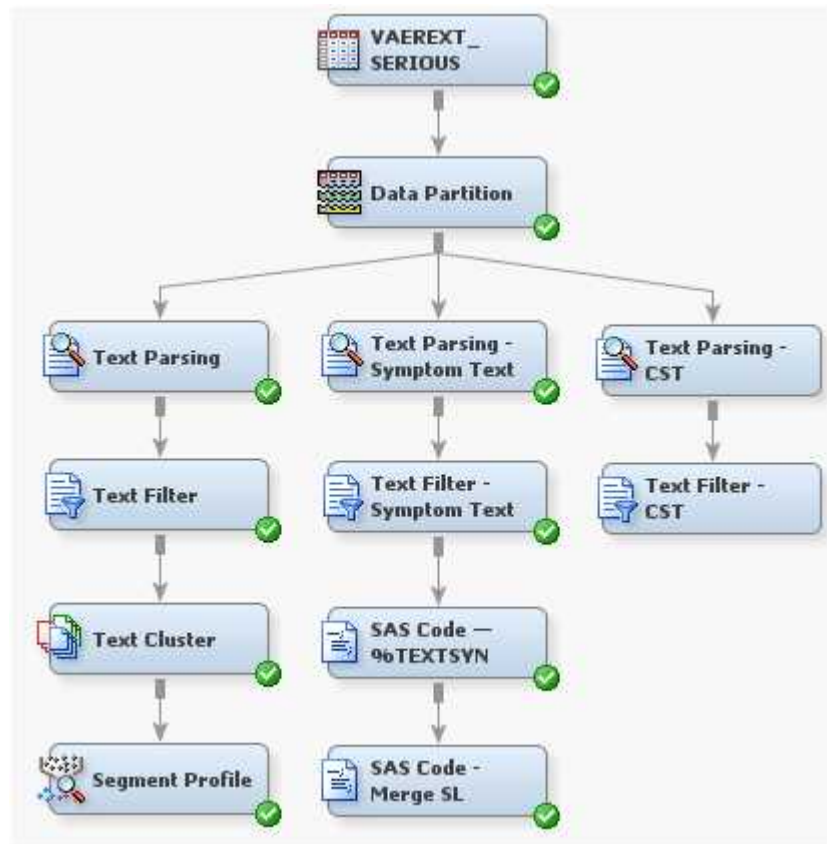
13. Right-click the **Text Parsing** — **Symptom Text** node, and then select **Copy**.



Note: It is important to copy the **Text Parsing** — **Symptom Text** node instead of creating a new Text Miner node. You do this in order to keep the same property settings you previously configured for the **Text Parsing** — **Symptom Text** node.

14. Right-click an empty space in the diagram workspace and select **Paste**.
15. Right-click the new **Text Parsing** node, and select **Rename**.
16. Enter *Text Parsing* — *CST* in the **Node Name** field.

CST stands for *Cleaned Symptom Text*.

17. Click **OK**.
18. Right-click the **Text Filter** — **Symptom Text** node, and then select **Copy**.
19. Right-click an empty space in the diagram workspace and select **Paste**.
20. Right-click the new **Text Filter** node, and select **Rename**.
21. Enter *Text Filter* — *CST* in the **Node Name** field.
22. Click **OK**.
23. Connect the **Data Partition** node to the **Text Parsing** — **CST** node.
24. Connect the **Text Parsing** — **CST** node to the **Text Filter** — **CST** node.



25. Select the **Text Parsing — CST** node.
26. Click the  for the **Synonyms** property.
27. Click **Import**.
28. Click **Mylib** in the SAS Libraries tree.
The contents of the Mylib library appear.
29. Select **Mylib.Vaerextsyms_new**.
30. Click **OK**.
The contents of the data set appear in the dialog box.
31. Click **OK**.
32. Right-click the **Text Filter — CST** node, and select **Run**. Click **Yes** in the Confirmation dialog box.
33. Click **OK** in the Run Status dialog box when the node has finished running.
34. Select the **Text Filter — CST** node.
35. Click the  for the **Filter Viewer** property.
The Interactive Filter Viewer window appears.
36. Select the plus sign (+) next to **patient** in the Terms table. Note that the misspellings **patien**, **patietn**, and **paitent** are included as child terms.

Terms			
	TERM	FREQ	# DOCS
<input checked="" type="checkbox"/>	patient	45164	18019
<input type="checkbox"/>	patien	7	7
<input type="checkbox"/>	patients	14	14
<input type="checkbox"/>	pts	278	215
<input type="checkbox"/>	patietn	5	5
<input type="checkbox"/>	paitent	5	5
<input type="checkbox"/>	pt	23367	10036
<input type="checkbox"/>	patient	21484	8653
<input type="checkbox"/>	ptsd	4	4

37. Close the Interactive Filter Viewer.

Chapter 6

Create Topics and Rules

About the Tasks That You Will Perform	41
Create Topics	41
Create Rules	43

About the Tasks That You Will Perform

This chapter shows how you can create topics and rules from filtered terms using the **Text Topic** node, and the **Text Rule Builder** node.

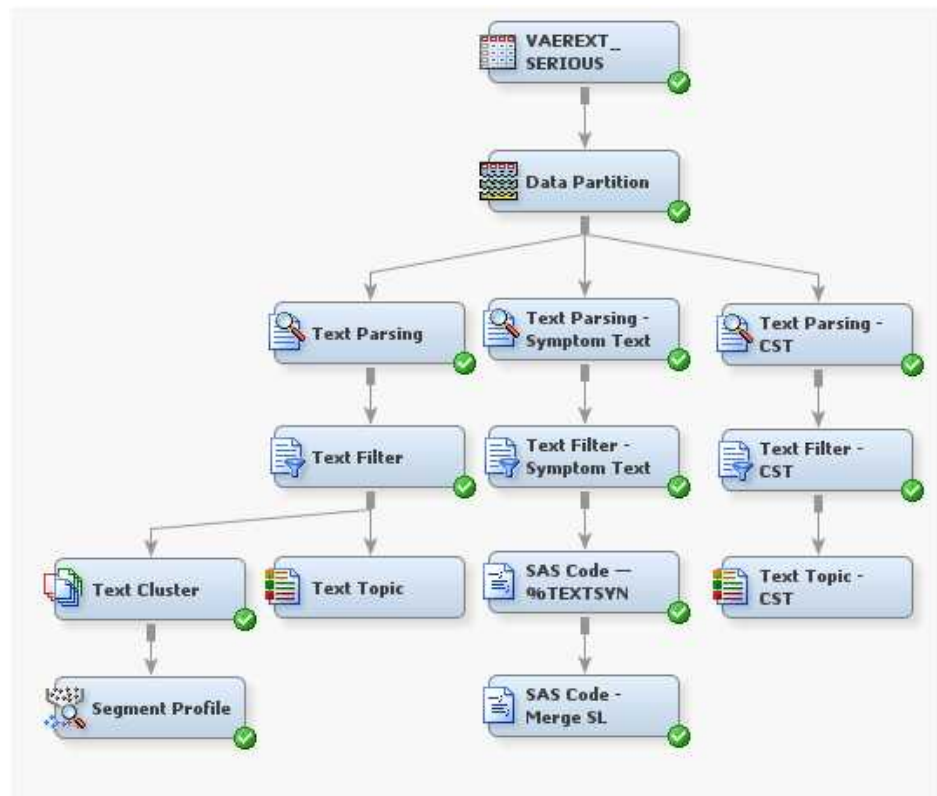
The **Text Topic** node enables you to explore the document collection by automatically associating terms and documents according to both discovered and user-defined topics. Topics are collections of terms that describe and characterize a main theme or idea. The goal in creating a list of topics is to establish combinations of words that you are interested in analyzing. The ability to combine individual terms into topics can improve your text mining analysis. Through combining, you can narrow the amount of text that is subject to analysis to specific groupings of words that you are interested in. For more information about the **Text Topic** node, see the SAS Text Miner Help.

The **Text Rule Builder** node generates an ordered set of rules from small subsets of terms that together are useful in describing and predicting a target variable. Each rule in the set is associated with a specific target category that consists of a conjunction that indicates the presence or absence of one or a small subset of terms (for example, “term1” AND “term2” AND (NOT “term3”)). A particular document matches this rule if and only if it contains at least one occurrence of term1 and of term2 but no occurrences of term3. This set of derived rules creates a model that is both descriptive and predictive. When categorizing a new document, it proceeds through the ordered set and chooses the target that is associated with the first rule that matches that document. The rules are provided in the syntax that can be used within SAS Content Categorization Studio, and can be deployed there. For more information about the **Text Rule Builder** node, see the SAS Text Miner help.

Create Topics

After filtering text, you can use the **Text Topic** node to create topics. Perform the following steps to use **Text Topic** nodes in the analysis:

1. Select the **Text Mining** tab on the node toolbar and drag a **Text Topic** node into the diagram workspace.
2. Connect the **Text Filter** node to the **Text Topic** node.
3. Select the **Text Mining** tab on the node toolbar and drag a **Text Topic** node into the diagram workspace.
4. Right-click the **Text Topic** node, and select **Rename**.
5. Enter *Text Topic — CST* in the **Node Name** field, and then click **OK**.
6. Select the **Text Topic — CST** node.
7. Enter **50** as the value for the **Number of Multi-term Topics** property.
8. Connect the **Text Filter — CST** node to the **Text Topic — CST** node.



9. Right-click the **Text Topic** node in the diagram workspace, and select **Run**. Click **Yes** in the Confirmation dialog box.
10. Click **Results** in the Run Status dialog box when the node has finished running.
11. Review the topics in the Topics table to see which terms make up each topic.

Topics						
Category	Topic ID ▲	Document Cutoff	Term Cutoff	Topic	Number of Terms	# Docs
Multiple	1	16.42	0.414	quality,+product quality complaint,+unsp...	85	5145
Multiple	2	12.34	0.322	information,+vaccinate,pneumococcal,...	157	6736
Multiple	3	10.40	0.286	varicella,+varicella virus vaccine,+virus,...	151	5976
Multiple	4	5.181	0.240	fluvirin,vaccinee,+receive,insufficient inf...	76	4985

12. Close the Results window.
13. Right-click the **Text Topic — CST** node in the diagram workspace, and select **Run**. Click **Yes** in the Confirmation dialog box.

14. Click **Results** in the Run Status dialog box when the node has finished running.
15. Review the topics in the Topics table to see which terms make up each topic.

Notice that the topics that were generated by running the **Text Topic — CST** are different from those that were generated by running the **Text Topic** node.

Topics						
Category	Topic ID	Document Cutoff	Term Cutoff	Topic	Number of Terms	# Docs
Multiple	1	10.56	0.241	+virus,+live,+varicella virus vaccine...	248	5231
Multiple	2	5.455	0.150	received,+patient,vaccinated,infor...	285	6731
Multiple	3	5.312	0.146	23v,23v polysaccharide vaccine,+i...	283	6519
Multiple	4	4.028	0.138	+seek,+unspecified medical attent...	163	3959

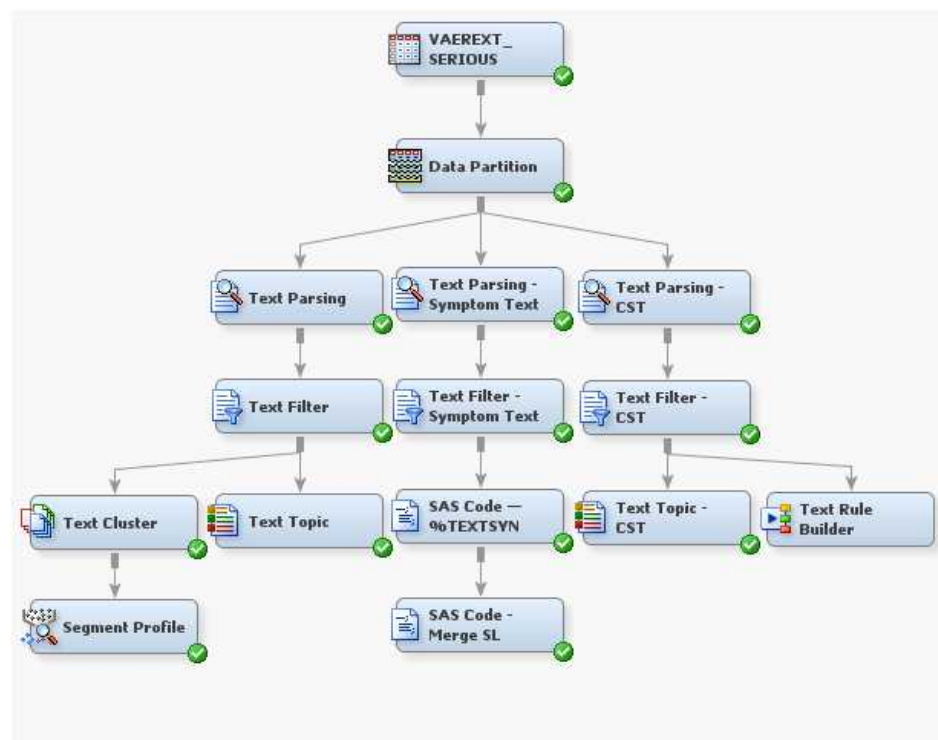
The differences follow from previously performed text cleaning activities, and the number of multi-term topics.

16. Close the Results window.

Create Rules

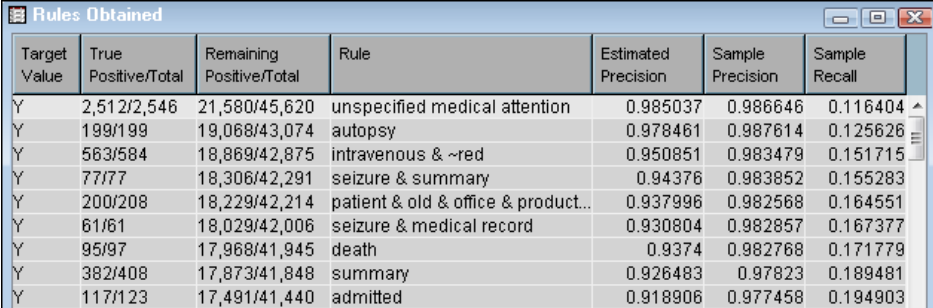
After filtering text, you can use the **Text Rule Builder** node to create rules. Perform the following steps to use a **Text Rule Builder** node in the analysis:

1. Select the **Text Mining** tab on the node toolbar and drag a **Text Rule Builder** node into the diagram workspace.
2. Connect the **Text Filter — CST** node to the **Text Rule Builder** node.



3. Right-click the **Text Rule Builder** node in the diagram workspace, and select **Run**. Click **Yes** in the Confirmation dialog box.

4. Click **Results** in the Run Status dialog box when the node has finished running.
5. Select the Rules Obtained window.



Target Value	True Positive/Total	Remaining Positive/Total	Rule	Estimated Precision	Sample Precision	Sample Recall
Y	2,512/2,546	21,580/45,620	unspecified medical attention	0.985037	0.986646	0.116404
Y	199/199	19,068/43,074	autopsy	0.978461	0.987614	0.125626
Y	563/584	18,869/42,875	intravenous & ~red	0.950851	0.983479	0.151715
Y	77/77	18,306/42,291	seizure & summary	0.94376	0.983852	0.155283
Y	200/208	18,229/42,214	patient & old & office & product...	0.937996	0.982568	0.164551
Y	61/61	18,029/42,006	seizure & medical record	0.930804	0.982857	0.167377
Y	95/97	17,968/41,945	death	0.9374	0.982768	0.171779
Y	382/408	17,873/41,848	summary	0.926483	0.97823	0.189481
Y	117/123	17,491/41,440	admitted	0.918906	0.977458	0.194903

In the second column above, the True Positive (the first number) is the number of documents that were correctly assigned to the rule. The Total (the second number) is the total positive.

In the third column above, the Remaining Positive (the first number) is the total number of remaining documents in the category. The Total (the second number) is the total number of documents remaining.

For more information about the Rules Obtained window or the **Text Rule Builder** node, see [“The Text Rule Builder Node” on page 77](#), and the SAS Text Miner Help.

6. Close the Results window.

Chapter 7

Create Models and Compare Them

About the Tasks That You Will Perform	45
Create Models	45
Compare the Models	47

About the Tasks That You Will Perform

This section shows how to use **Decision Tree** nodes to create models, and compare them with a **Model Comparison** node.

A **Decision Tree** node can be used to classify observations based on the values of nominal, binary, or ordinal targets. It can also predict outcomes for interval targets or the appropriate decision when you specify decision alternatives. An empirical tree represents a segmentation of the data that is created by applying a series of simple rules. Each rule assigns an observation to a segment based on the value of one input.

One rule is applied after another, resulting in a hierarchy of segments within segments. The hierarchy is called a tree, and each segment is called a node. The original segment contains the entire data set and is called the root node of the tree. A node with all its successors forms a branch of the node that created it.

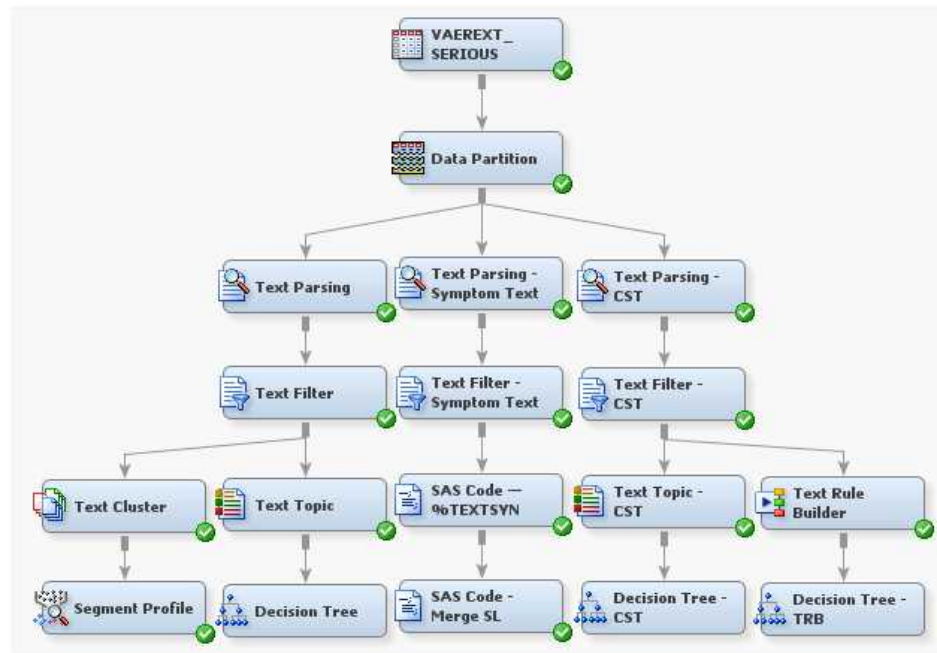
The final nodes are called leaves. For each leaf, a decision is made and applied to all observations in the leaf. The type of decision depends on the context. In predictive modeling, the decision is the predicted value. For more information about **Decision Tree** nodes, see the SAS Enterprise Miner help.

Create Models

To create models by using **Decision Tree** nodes:

1. Select the **Model** tab on the node toolbar and drag a **Decision Tree** node into the diagram workspace.
2. Connect the **Text Topic** node to the **Decision Tree** node.
3. Select the **Model** tab on the node toolbar and drag a **Decision Tree** node into the diagram workspace.
4. Right-click the **Decision Tree** node, and select **Rename**.

5. Enter *Decision Tree* — CST in the **Node Name** field, and then click **OK**.
6. Connect the **Text Topic** — CST node to the **Decision Tree** — CST node.
7. Select the **Model** tab on the node toolbar and drag a **Decision Tree** node into the diagram workspace.
8. Right-click the **Decision Tree** node, and select **Rename**.
9. Enter *Decision Tree* — TRB in the **Node Name** field, and then click **OK**.
10. Connect the **Text Rule Builder** node to the **Decision Tree** — TRB node.

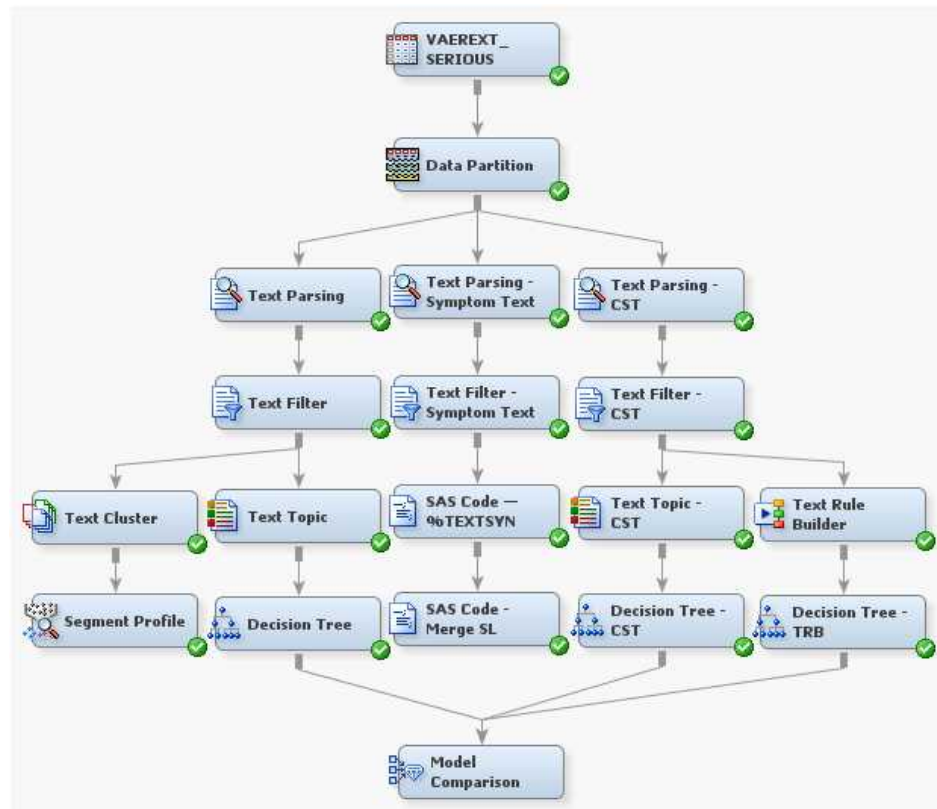


11. Right-click the **Decision Tree** node in the diagram workspace, and select **Run**. Click **Yes** in the Confirmation dialog box.
12. Click **Results** in the Run Status dialog box when the node has finished running.
13. Select the Tree window, and explore the tree that was obtained.
14. Close the Results window.
15. Right-click the **Decision Tree** — CST node in the diagram workspace, and select **Run**. Click **Yes** in the Confirmation dialog box.
16. Click **Results** in the Run Status dialog box when the node has finished running.
17. Select the Tree window, and explore the tree that was obtained. How does this tree differ from the previous tree? The primary difference is that different topics were used for each decision point.
18. Close the Results window.
19. Right-click the **Decision Tree** — TRB node in the diagram workspace, and select **Run**. Click **Yes** in the Confirmation dialog box.
20. Click **Results** in the Run Status dialog box when the node has finished running.
21. Select the Tree window, and explore the tree that was obtained. How does this tree differ from the previous two trees? The primary difference is that instead of topics, single or multi-term rules were used for each decision point.
22. Close the Results window.

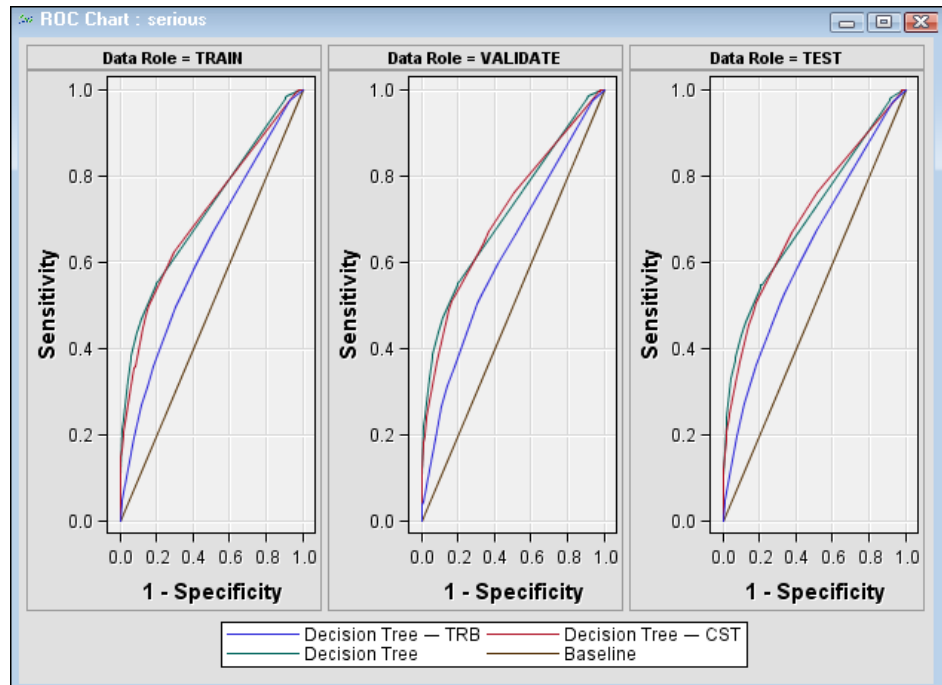
Compare the Models

To compare models using a **Model Comparison** node:

1. Select the **Assess** tab on the node toolbar and drag a **Model Comparison** node into the diagram workspace.
2. Connect the **Decision Tree** node, the **Decision Tree — CST** node, and the **Decision Tree — TRB** node to the **Model Comparison** node.



3. Right-click the **Model Comparison** node in the diagram workspace, and select **Run**. Click **Yes** in the Confirmation dialog box.
4. Click **Results** in the Run Status dialog box when the node has finished running.
The Results window appears.
5. Select the ROC Chart.



The greater the area under the curve, the better the model. The brown line represents the baseline to compare the models by. The blue line represents how well the **Decision Tree — TRB** model did at predicting the target SERIOUS. This model used input from the **Text Rule Builder** node. The green line represents how well the **Decision Tree** model did at predicting the target SERIOUS. The red line represents how well the **Decision Tree — CST** model did at predicting the target SERIOUS.

Both the **Decision Tree** and the **Decision Tree — CST** models performed better than the **Decision Tree — TRB** model. The **Decision Tree** and the **Decision Tree — CST** models performed similarly with respect to each other.

As an additional exercise, you could try modifying the number of multiple or single term topics in the **Text Topic** or **Text Topic — CST** nodes. Then rerun the **Decision Tree** and the **Decision Tree — CST** nodes to see whether the models have been improved.

Chapter 8

The Text Import Node

About the Text Import Node	49
Using the Text Import Node	49
Contents	49
Import Documents from a Directory	50
Import Documents from the Web	51

About the Text Import Node

The **Text Import** node serves as a replacement for an **Input Data** node. It enables you to create data sets dynamically from files contained in a directory or from the Web. The **Text Import** node takes an import directory that contains text files in potentially proprietary formats such as MS Word and PDF files as input. The tool traverses this directory and filters or extracts the text from the files, places a copy of the text in a plain text file, and places a snippet (or possibly even all) of the text in a SAS data set.

If a URL is specified, the node crawls Web sites, retrieves files from the Web, and puts them in an import directory before doing this filtering process. The output of a **Text Import** node is a data set that can be imported into the **Text Parsing** node. In addition to filtering the text, the **Text Import** node can also identify the language that the document is in and take care of transcoding documents to the session encoding.

For more information about the **Text Import** node, see the SAS Text Miner Help.

The rest of this chapter presents two examples of how you can use the **Text Import** node.

Using the Text Import Node

Contents


The following examples show you how you can use the **Text Import** node to import documents from a directory or the Web. These examples assume that SAS Enterprise Miner is running, the SAS Document Conversion server is running, and a diagram workspace has been opened in a project. For information about creating a project and a diagram, see [“Setting Up Your Project” on page 9](#).

- [“Import Documents from a Directory” on page 50](#)


- “Import Documents from the Web” on page 51

Import Documents from a Directory

To import documents from a directory:

1. Select the **Text Mining** tab, and drag a **Text Import** node into the diagram workspace.
2. Click the  for the **Import File Directory** property of the **Text Import** node.
A Select Server Directory dialog box appears.
3. Navigate to a folder that contains documents that you want to create a data set from, select it, and then click **OK**.

Note: To see the file types that you want to select, you might need to select **All Files** in the type drop-down menu.

4. Click the  for the **Language** property.

The Languages dialog box appears.

5. Select one or more licensed languages in which to require the language identifier to assign each document’s language, and then click **OK**.
6. (Optional) Specify the file types to process for the **Extensions** property. For example, if you want to look at only documents with a .txt and a .pdf extension, specify **.txt .pdf** for the **Extensions** property, and click **Enter**.

Note: If you do not specify file types to process, the **Text Import** node processes all file types in the specified import file directory.

7. Right-click the **Text Import** node, and select **Run**.
8. Click **Yes** in the Confirmation dialog box.
9. Click **Results** in the Run Status dialog box when the node has finished running.
The Results window appears.
10. Examine results from the documents that you imported.


You can now use the **Text Import** node as an input data source for your text mining analysis.

11. Select the **Text Mining** tab, and drag a **Text Parsing** node into the diagram workspace.
12. Connect the **Text Import** node to the **Text Parsing** node.
13. Right-click the **Text Parsing** node, and select **Run**.
14. Click **Yes** in the Confirmation dialog box.
15. Click **OK** in the Run Status dialog box when the node has finished running.

Import Documents from the Web

To import documents from the Web:

Note: Web crawling is supported only on Windows operating systems.

1. Select the **Text Mining** tab, and drag a **Text Import** node into the diagram workspace.
2. Click the  for the **Import File Directory** property of the **Text Import** node.

A Select Server Directory dialog box appears.

3. Navigate to a folder, select it, and then click **OK**.

The documents are first written to the **Import File Directory** location. The files are processed from the **Import File Directory** location, and then are written to the **Destination Directory** location.

4. Enter the uniform resource locator (URL) of a Web page that you want to crawl in the **URL** property of the **Text Import** node. For example, try *www.sas.com*.
5. Type *1* as the number of levels to crawl in the **Depth** property.
6. Set the **Domain** property to **Unrestricted**.

Note: If you want to crawl a password-protected Web site, set the **Domain** property to **Restricted**, and provide a user name for the **User Name** property and a password for the **Password** property.

7. Right-click the **Text Import** node and select **Run**.
8. Click **Yes** in the Confirmation dialog box.
9. Click **Results** in the Run Status dialog box when the node has finished running.
10. Examine results from the Web site.

You can now use the **Text Import** node as an input data source for your text mining analysis.

11. Select the **Text Mining** tab, and drag a **Text Parsing** node into the diagram workspace.
12. Connect the **Text Import** node to the **Text Parsing** node.
13. Right-click the **Text Parsing** node, and select **Run**.
14. Click **Yes** in the Confirmation dialog box.
15. Click **OK** in the Run Status dialog box when the node has finished running.

Chapter 9

The Text Parsing Node

About the Text Parsing Node	53
Using the Text Parsing Node	53

About the Text Parsing Node

The **Text Parsing** node enables you to parse a document collection in order to quantify information about the terms that are contained therein. You can use the **Text Parsing** node with volumes of textual data such as e-mail messages, news articles, Web pages, research papers, and surveys. For more information about the **Text Parsing** node, see the SAS Text Miner Help.

The rest of this chapter presents an example of how you can use the **Text Parsing** node.

Using the Text Parsing Node

This example shows you how to use the **Text Parsing** node to identify terms and their instances in a data set that contains text. This example assumes that SAS Enterprise Miner is running, and that a diagram workspace has been opened in a project. For information about creating a project and a diagram, see [“Setting Up Your Project” on page 9](#).

1. The SAS data set SAMPSIO.ABSTRACT contains the titles and text of abstracts from conferences. Create the ABSTRACT data source and add it to your diagram workspace. Set the Role value of the TEXT and TITLE variables to **Text**.
2. Select the **Text Mining** tab on the toolbar, and drag a **Text Parsing** node into the diagram workspace.
3. Connect the ABSTRACT data source to the **Text Parsing** node.



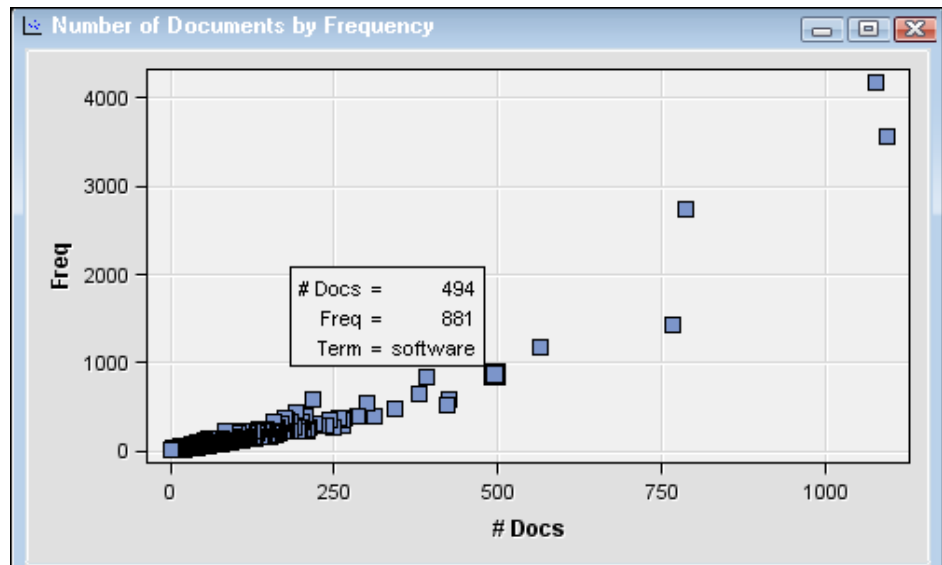
4. In the diagram workspace, right-click the **Text Parsing** node and select **Run**. Click **Yes** in the Confirmation dialog box that appears.

- Click **Results** in the Run Status dialog box when the node finishes running. The Results window displays a variety of tabular and graphical output to help you analyze the terms and their instances in the ABSTRACT data source.
- Sort the terms in the Terms table by frequency, and then select the term “software.” As the Terms table illustrates, the term “software” is a noun that occurs in 494 documents in the ABSTRACT data source, and appears a total number of 881 times.

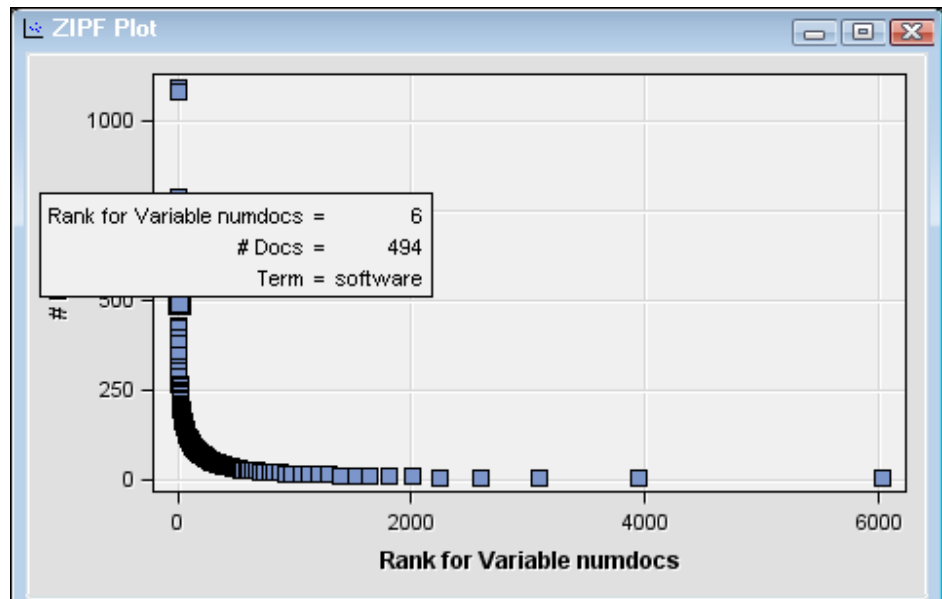
Term	Role	Attribute	Freq ▼	# Docs	Keep	Parent/Child Status	Parent ID	Rank for Variable numdocs
+ sas in...	Comp...	Entity	4187	1077Y	+		23263	2
+ be ...	Verb	Alpha	3571	1093N	+		141	1
data ...	Noun	Alpha	2747	786Y			16	3
+ use ...	Verb	Alpha	1429	766N	+		468	4
+ syste...	Noun	Alpha	1164	565Y	+		64	5
software...	Noun	Alpha	881	494Y			20	6
+ applic...	Noun	Alpha	844	392Y	+		33	9
+ user ...	Noun	Alpha	645	379Y	+		122	10
+ have ...	Verb	Alpha	582	425N	+		190	7

When you select a term in the Terms table, the point corresponding to that term in the Text Parsing Results plots is highlighted.

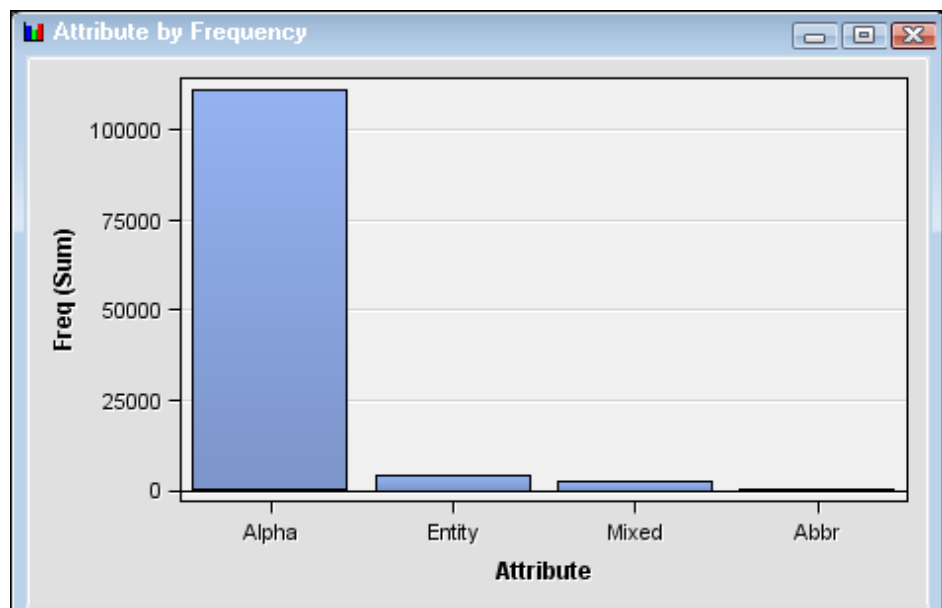
- Select the Number of Documents by Frequency plot, and position the cursor over the highlighted point for information about the term “software.”



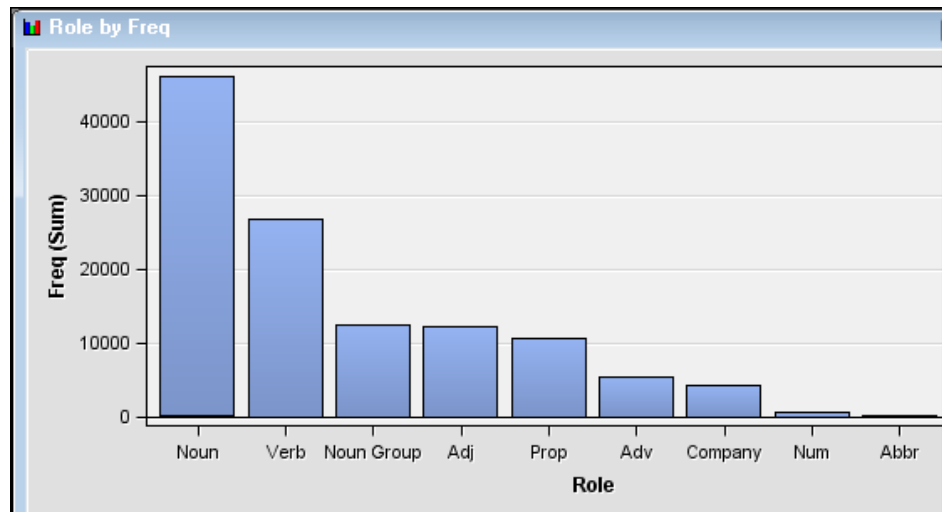
Similar information is also presented in a ZIPF plot.



The Attribute by Frequency chart shows that **Alpha** has the highest frequency among attributes in the document collection.




The Role by Freq chart illustrates that **Noun** has the highest frequency among roles in the document collection.

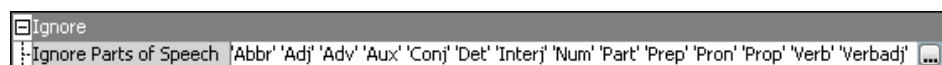


8. Return to the Terms table, and notice that the term “software” is kept in the text parsing analysis. This is illustrated by the value of **Y** in the Keep column. Notice that not all terms are kept when you run the **Text Parsing** node with default settings.

Term	Role	Attribute	Freq ▼	# Docs	Keep	Parent/Child Status	Parent ID	Rank for Variable numdocs
+ sas in...	Comp...	Entity	4187	1077	Y	+	23263	2
+ be	Verb	Alpha	3571	1093	N	+	141	1
data	Noun	Alpha	2747	786	Y		16	3
+ use	Verb	Alpha	1429	766	N	+	468	4
+ syste...	Noun	Alpha	1164	565	Y	+	64	5
software...	Noun	Alpha	881	494	Y		20	6
+ applic...	Noun	Alpha	844	392	Y	+	33	9
+ user	Noun	Alpha	645	379	Y	+	122	10
+ have	Verb	Alpha	582	425	N	+	190	7

The **Text Parsing** node not only enables you to gather statistical data about the terms in a document collection. It also enables you to modify your output set of parsed terms by dropping terms that are a certain part of speech, type of entity, or attribute. Scroll down the list of terms in the Terms table and notice that many of the terms with a role other than **Noun** are kept. Assume that you want to limit your text parsing results to terms with a role of **Noun**.

9. Close the Results window.
10. Select the **Text Parsing** node, and then select the  for the **Ignore Parts of Speech** property.
11. In the Ignore Parts of Speech dialog box, select all parts of speech except for **Noun** by holding down CTRL and clicking on each option. Click **OK**. Notice that the value for the **Ignore Parts of Speech** property is updated with your selection.

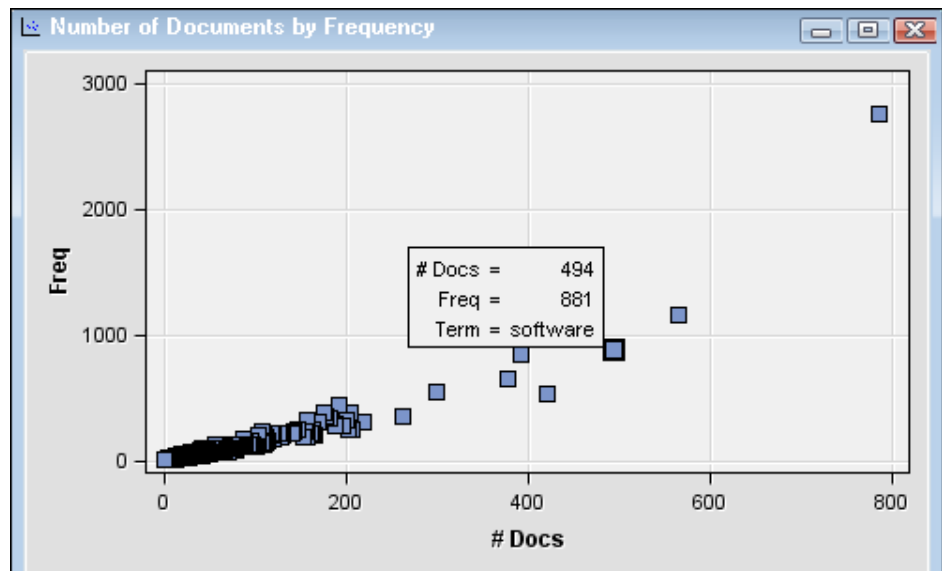


12. In addition to nouns, also keep noun groups. Set the **Noun Groups** property to **Yes**.

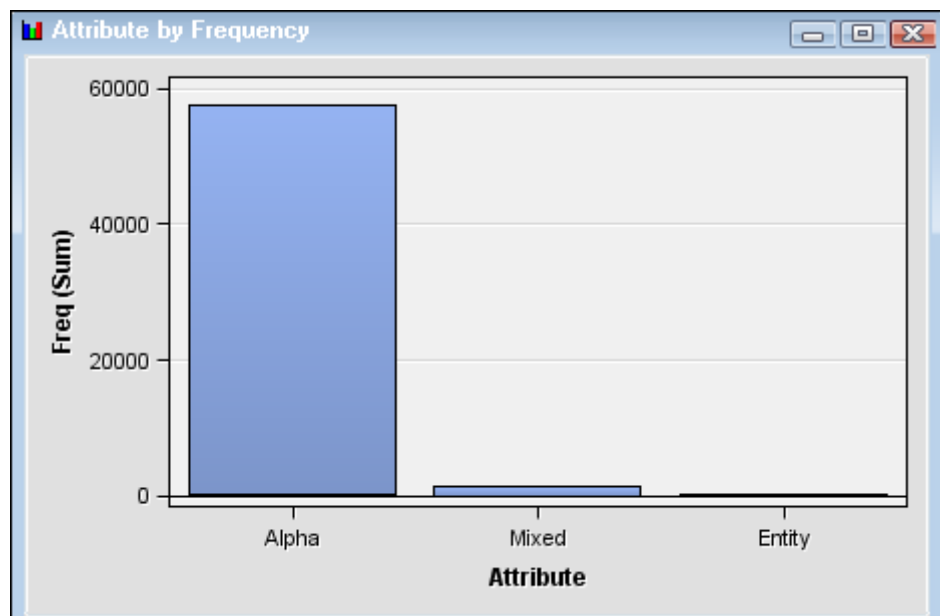
13. Right-click the **Text Parsing** node and select **Run**. Click **Yes** in the Confirmation dialog box that appears. Select **Results** in the Run Status dialog box when the node has finished running. Notice that the term “software” has a higher rank among terms with a role of just “noun” or “noun group” than it did when other roles were included. If you scroll down in the Terms table, you can see that just terms with a **Noun** or **Noun Group** role are included.

Terms								
Term	Role	Attribute	Freq	# Docs	Keep	Parent/Child Status	Parent ID	Rank for Variable numdocs
data	...Noun	Alpha	2747	786Y			14	1
+ system	...Noun	Alpha	1164	565Y	+		31	2
software	...Noun	Alpha	881	494Y			16	3
+ paper	...Noun	Alpha	524	422Y	+		62	4

As we would expect, there are fewer terms plotted in the Number of Documents by Frequency plot:



Similarly, the total number of terms in the output results with an attribute of **Alpha** has decreased, as can be seen in the Attribute by Frequency chart:



Chapter 10

The Text Filter Node

About the Text Filter Node	59
Using the Text Filter Node	59

About the Text Filter Node

The **Text Filter** node can be used to reduce the total number of parsed terms or documents that are analyzed. Therefore, you can eliminate extraneous information so that only the most valuable and relevant information is considered. For example, the **Text Filter** node can be used to remove unwanted terms and to keep only documents that discuss a particular issue. This reduced data set can be orders of magnitude smaller than the one that represents the original collection, which might contain hundreds of thousands of documents and hundreds of thousands of distinct terms.

For more information about the **Text Filter** node, see the SAS Text Miner Help.

The rest of this chapter presents an example of how you can use the **Text Filter** node.

Using the Text Filter Node

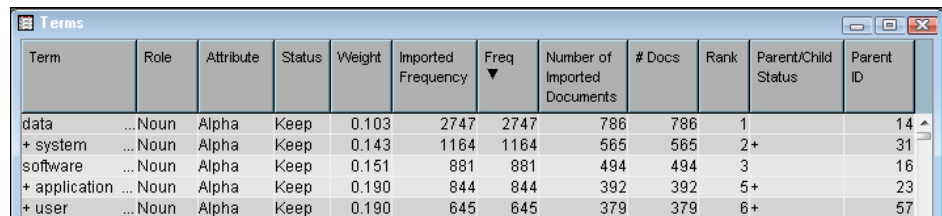
This example shows you how to filter out terms using the **Text Filter** node. This example assumes that you have performed “[Using the Text Parsing Node](#)” on page 53, and builds off the process flow diagram created there.

1. Select the **Text Mining** tab on the toolbar, and drag a **Text Filter** node into the diagram workspace.
2. Connect the **Text Parsing** node to the **Text Filter** node.




3. In the diagram workspace, right-click the **Text Filter** node and select **Run**. Click **Yes** in the Confirmation dialog box.
4. Click **Results** in the Run Status dialog box when the node finishes running.

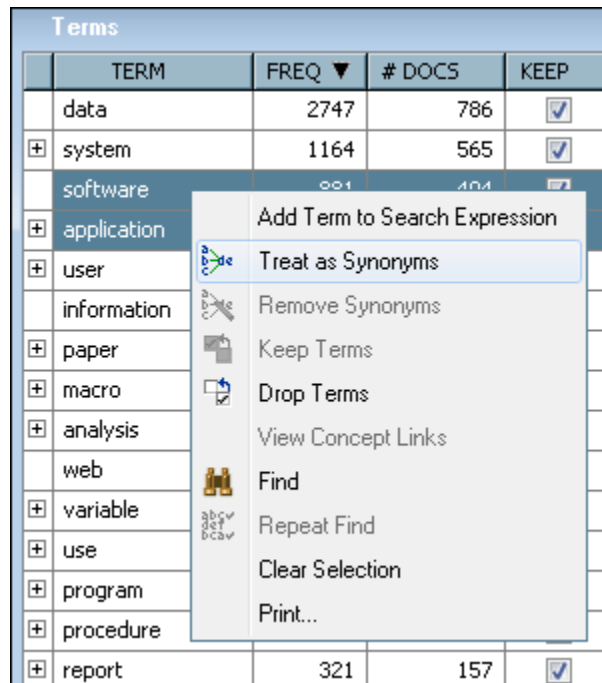
5. Select the Terms table. Sort the terms by frequency by clicking the Freq column heading.



Term	Role	Attribute	Status	Weight	Imported Frequency	Freq	Number of Imported Documents	# Docs	Rank	Parent/Child Status	Parent ID
data	... Noun	Alpha	Keep	0.103	2747	2747	786	786	1		14
+ system	... Noun	Alpha	Keep	0.143	1164	1164	565	565	2+		31
software	... Noun	Alpha	Keep	0.151	881	881	494	494	3		16
+ application	... Noun	Alpha	Keep	0.190	844	844	392	392	5+		23
+ user	... Noun	Alpha	Keep	0.190	645	645	379	379	6+		57

Assume that for this text mining analysis, you know that “software” and “application” are really used as synonyms in the documents that you want to analyze, and that you want to treat them as the same term.

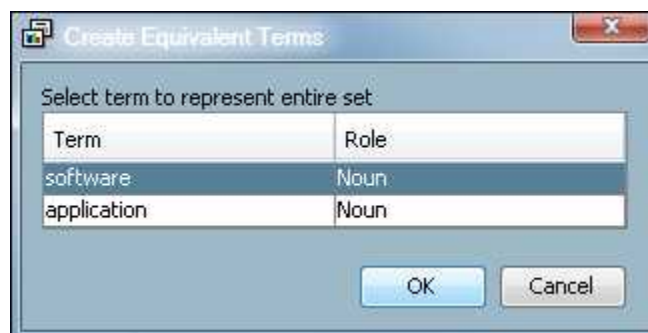
6. Close the Results window. Select the **Text Filter** node, and then click the  for the **Filter Viewer** property.
7. In the Interactive Filter Viewer sort the terms in the Terms table by frequency. Hold Ctrl down on your keyboard, select “software” and “application”, and then right-click “software” and select **Treat as Synonyms** from the drop-down menu.



	TERM	FREQ ▼	# DOCS	KEEP
	data	2747	786	<input checked="" type="checkbox"/>
+	system	1164	565	<input checked="" type="checkbox"/>
	software	881	494	<input checked="" type="checkbox"/>
+	application			<input checked="" type="checkbox"/>
+	user			<input checked="" type="checkbox"/>
	information			<input checked="" type="checkbox"/>
+	paper			<input checked="" type="checkbox"/>
+	macro			<input checked="" type="checkbox"/>
+	analysis			<input checked="" type="checkbox"/>
	web			<input checked="" type="checkbox"/>
+	variable			<input checked="" type="checkbox"/>
+	use			<input checked="" type="checkbox"/>
+	program			<input checked="" type="checkbox"/>
+	procedure			<input checked="" type="checkbox"/>
+	report	321	157	<input checked="" type="checkbox"/>

Add Term to Search Expression
 Treat as Synonyms
 Remove Synonyms
 Keep Terms
 Drop Terms
 View Concept Links
 Find
 Repeat Find
 Clear Selection
 Print...

8. In the Create Equivalent Terms dialog box, select **software** as the term to represent both terms in the Terms table.



Create Equivalent Terms

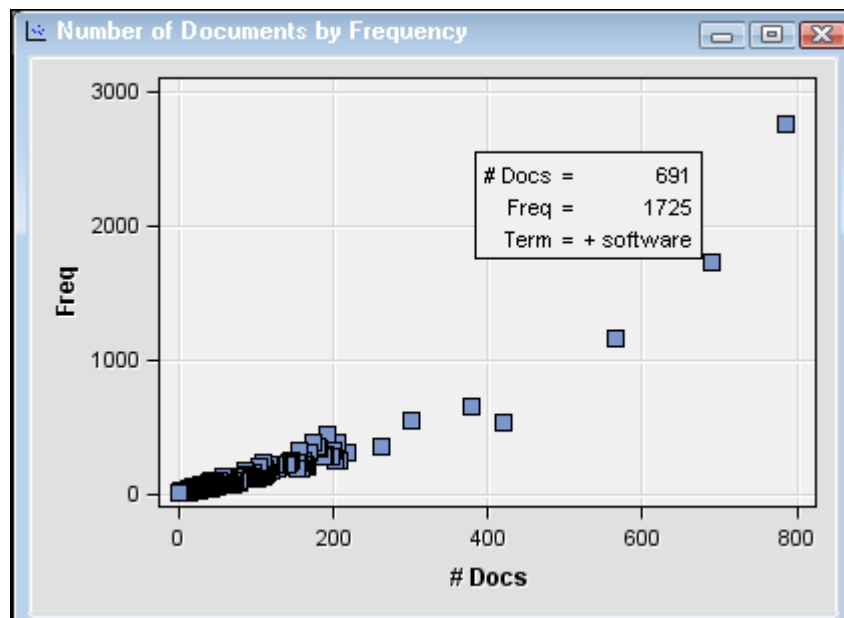
Select term to represent entire set:

Term	Role
software	Noun
application	Noun

9. Click **OK** in the Create Equivalent Terms dialog box. Notice that the term “software” now represents both terms in the Terms table. Expand the term “software”.

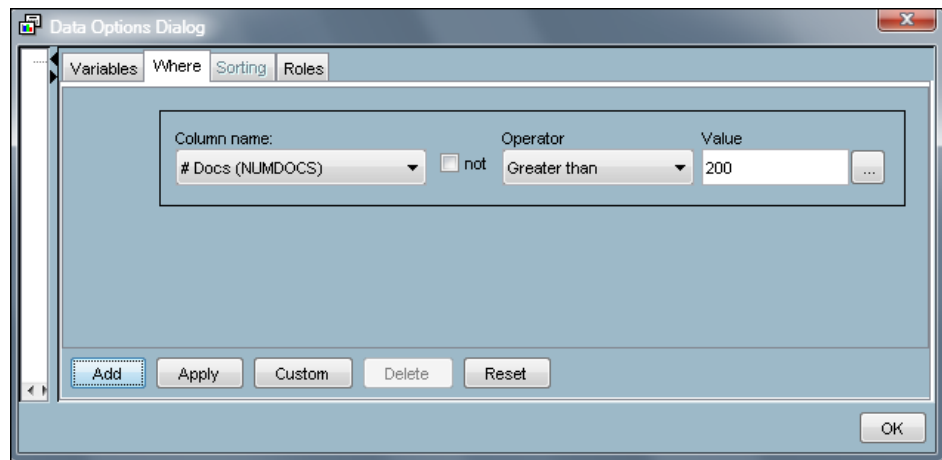
Terms							
	TERM	FREQ ▼	# DOCS	KEEP	WEIGHT	ROLE	ATTRIBUTE
	data	2747	786	<input checked="" type="checkbox"/>	0.103	Noun	Alpha
<input type="checkbox"/>	software	1725	691	<input checked="" type="checkbox"/>	0.114	Noun	Alpha
	application	499	264			Noun	Alpha
	software	881	494			Noun	Alpha
	applications	345	219			Noun	Alpha
<input checked="" type="checkbox"/>	system	1164	565	<input checked="" type="checkbox"/>	0.143	Noun	Alpha
<input checked="" type="checkbox"/>	user	645	379	<input checked="" type="checkbox"/>	0.19	Noun	Alpha

10. Close the Interactive Filter Viewer. When prompted whether you would like to save your changes, select **Yes**.
11. Right-click the **Text Filter** node, and select **Run**. Select **Yes** in the Confirmation dialog box. Select **Results** in the Run Status dialog box when the node has finished running.
12. Select the Number of Documents by Frequency plot to see how both terms are now treated as the same.

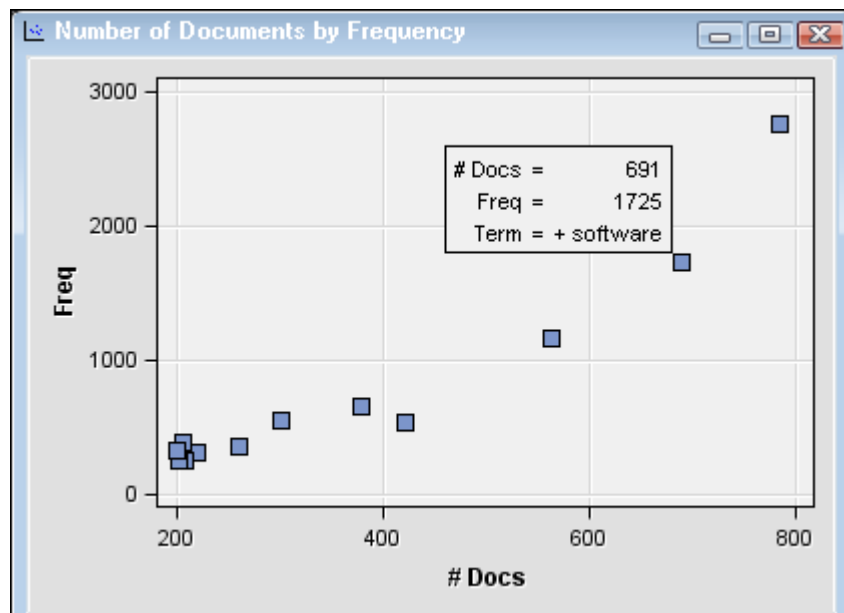



You can also use options to change your view or specify a subset of results to appear in a plot. For example, consider that you want to refine this plot to only show terms that appear in more than 200 documents.

13. Right-click the Number of Documents by Frequency plot, and select **Data Options**.
14. Select the **Where** tab in the Data Options Dialog box. Select **# Docs** from the **Column name** drop-down menu. Select **Greater than** from the **Operator** drop-down menu. Type **200** in the **Value** text box.



15. Click **Apply**, and then click **OK**. The Number of Documents by Frequency plot resizes and includes only terms that occur in more than 200 documents.





16. Close the Results window. In addition to resizing or subsetting a plot to help focus your analysis, you can also directly search for terms using the Interactive Filter Viewer.
17. Select the **Text Filter** node, and then click the  for the **Filter Viewer** property. In the Interactive Filter Viewer, type *software* in the **Search** text box, and click **Apply**.

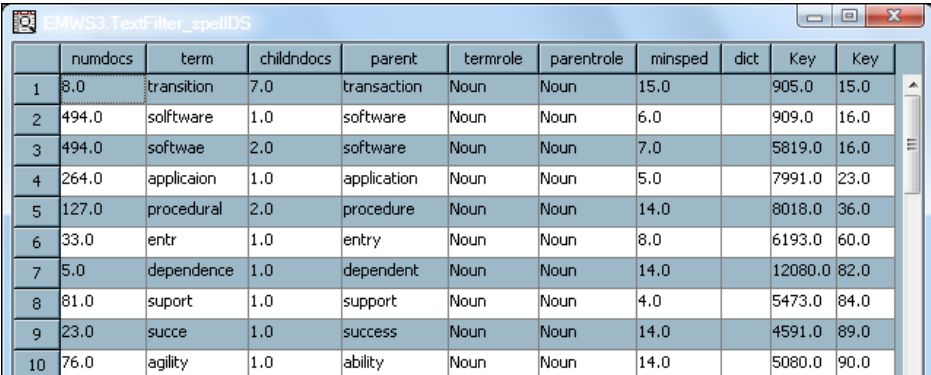
TEXT	TEXTFILTER_SNIPPET	TEXTFILTER_RELEVANCE	TITLE
New Features in SAS/ACCESS Software What is	... SAS / ACCESS Software What is	1.0	New Features in SAS/ACCES...
Extending the Power of Your SAS System Applications	... with Enterprise Reporter	1.0	Extending the Power of Your...
Collecting Data Via the Internet with SAS/IntrNet and	... SAS / SHARE Software The Texas	1.0	Collecting Data Via the Inter...
Producing Structured Clinical Trial Reports using SAS	... Reports using SAS Software : A	1.0	Producing Structured Clinical ...
A Table Production System That Meets the Challenges	... SAS / AF Software and the	0.857	A Table Production System T...
Data Warehousing on a Shoestring Perhaps the	... cost of additional software ... '	0.857	Data Warehousing on a Sho...
Delivering South Carolina Health and Demographic	... SAS / IntrNet Software This paper	0.714	Delivering South Carolina He...

The Documents table provides a snippet of text that includes the term that you are searching for. You can use information in the Documents table to help you understand the context in which a term is being used. To do so, examine the snippet result in addition to the full text and title of the document. For more information about the Interactive Filter Viewer, see the Interactive Filter Viewer topic in the SAS Text Miner Help.


Searching for a term in the Interactive Filter Viewer raises an interesting problem. As shown above, a search for “software” is case insensitive. However, what if there are instances of a term that you want to find that are misspelled in the document collection? You can also check for spelling when filtering terms using a dictionary data set.

18. Close the Interactive Filter Viewer, and select **No** when prompted for whether you want to save changes.
19. (Optional) Select the **Text Filter** node, and set the **Check Spelling** property to **Yes**. When you rerun the **Text Filter** node, terms will be checked for misspellings. You can also specify a data set to use in spell-checking by clicking the  for the **Dictionary** property and selecting a data set. For information about creating a dictionary data set, see the How to Create a Dictionary Data Set topic in the SAS Text Miner help.

Right-click the **Text Filter** node, and select **Run**. Select **Yes** in the Confirmation dialog box. When the node finishes running, select **OK** in the Run Status dialog box. Click the  for the **Spell-Checking Results** property to access a window in which you can view the data set that contains spelling corrections that were generated during spell-checking. For example, the term "softwae" is identified as a misspelling of the term "software."



	numdocs	term	childndocs	parent	termrole	parentrole	minsped	dict	Key	Key
1	8.0	transition	7.0	transaction	Noun	Noun	15.0		905.0	15.0
2	494.0	software	1.0	software	Noun	Noun	6.0		909.0	16.0
3	494.0	softwae	2.0	software	Noun	Noun	7.0		5819.0	16.0
4	264.0	applicaiion	1.0	application	Noun	Noun	5.0		7991.0	23.0
5	127.0	procedural	2.0	procedure	Noun	Noun	14.0		8018.0	36.0
6	33.0	entr	1.0	entry	Noun	Noun	8.0		6193.0	60.0
7	5.0	dependence	1.0	dependent	Noun	Noun	14.0		12080.0	82.0
8	81.0	suport	1.0	support	Noun	Noun	4.0		5473.0	84.0
9	23.0	succe	1.0	success	Noun	Noun	14.0		4591.0	89.0
10	76.0	agility	1.0	ability	Noun	Noun	14.0		5080.0	90.0

You can see this relationship in the Terms table in the Interactive Filter Viewer. Click the  for the **Filter Viewer** property. Expand the term "software" in the Terms table to view its synonyms. The synonyms include "softwae," which was identified as a misspelled term during spell-checking.

Terms							
	TERM	FREQ ▼	# DOCS	KEEP	WEIGHT	ROLE	ATTRIBUTE
	data	2747	786	<input checked="" type="checkbox"/>	0.103	Noun	Alpha
[-]	software	1729	691	<input checked="" type="checkbox"/>	0.114	Noun	Alpha
	applications	345	219			Noun	Alpha
	applicaion	1	1			Noun	Alpha
	application	499	264			Noun	Alpha
	softwae	2	2			Noun	Alpha
	software	881	494			Noun	Alpha
	software	1	1			Noun	Alpha
[+]	system	1164	565	<input checked="" type="checkbox"/>	0.143	Noun	Alpha
[+]	user	645	379	<input checked="" type="checkbox"/>	0.19	Noun	Alpha

Notice that the synonyms also include "application," which was created in steps 7-10 of this example, and "applicaion," which was identified during spell-checking as a misspelling of "application."

Chapter 11

The Text Topic Node

About the Text Topic Node	65
Using the Text Topic Node	65

About the Text Topic Node

The **Text Topic** node enables you to explore the document collection by automatically associating terms and documents according to both discovered and user-defined topics. Topics are collections of terms that describe and characterize a main theme or idea. The goal in creating a list of topics is to establish combinations of words that you are interested in analyzing. The ability to combine individual terms into topics can improve your text mining analysis. Through combining, you can narrow the amount of text that is subject to analysis to specific groupings of words that you are interested in.

For example, you might be interested in mining articles that discuss the activities of a "company president." One way to approach this task is to look at all articles that have the term "company," and all articles that have the term "president." The **Text Topic** node enables you to combine the terms "company" and "president" into the topic "company president." The approach is different from clustering. Clustering assigns each document to a unique group, while the **Text Topic** node assigns a score for each document and term to each topic. Then thresholds are used to determine whether the association is strong enough to consider that the document or term belongs to the topic. As a result, documents and terms can belong to more than one topic or to none at all. The number of topics that you request should be directly related to the size of the document collection (for example, a large number for a large collection).

For more information about the **Text Topic** node, see the SAS Text Miner Help.

Using the Text Topic Node

This example shows you how to create topics using the **Text Topic** node. This example assumes that you have performed [“Using the Text Filter Node” on page 59](#), and builds off the process flow diagram created there.

1. Select the **Text Mining** tab on the toolbar, and drag a **Text Topic** node into the diagram workspace.
2. Connect the **Text Filter** node to the **Text Topic** node.

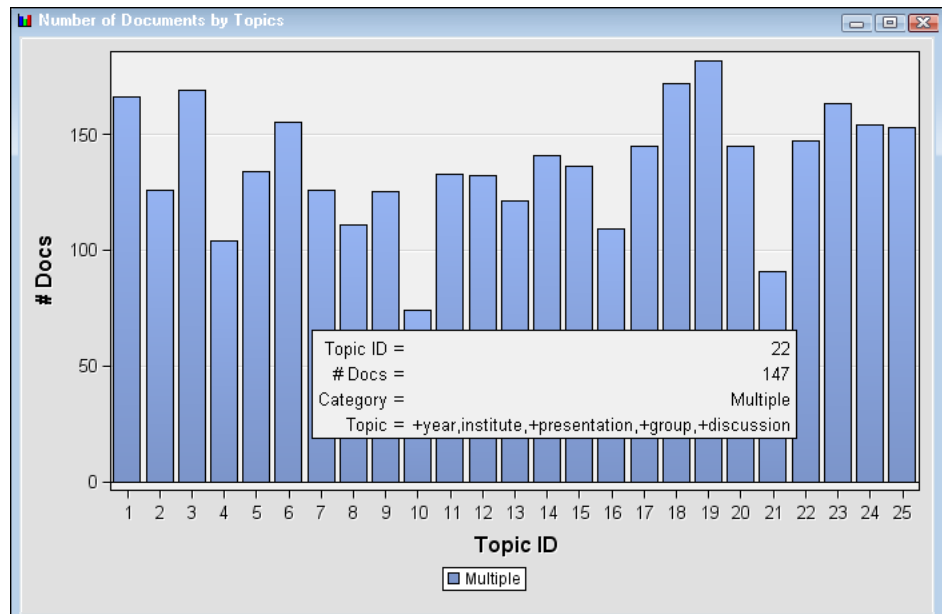


3. In the diagram workspace, right-click the **Text Topic** node and select **Run**. Click **Yes** in the Confirmation dialog box that appears. Click **Results** in the Run Status dialog box when the node finishes running.
4. Select the Topics table to view the topics that have been created with a default run of the **Text Topic** node.

Category	Topic ID	Document Cutoff	Term Cutoff	Topic	Number of Terms	# Docs
Multiple	1	0.537	0.124	+data set,+variable,+valu...	87	166
Multiple	2	0.485	0.118	+warehouse,+data ware...	79	126
Multiple	3	0.485	0.119	+server,+mainframe,+cli...	85	169
Multiple	4	0.443	0.111	+macro,+macro variable,...	47	104
Multiple	5	0.461	0.115	+customer,+business,+s...	99	134
Multiple	6	0.457	0.115	web,+page,+browser,inf...	85	155
Multiple	7	0.457	0.117	+sample,+analysis,+test,...	110	126
Multiple	8	0.407	0.109	+graph,graph software,gr...	83	111
Multiple	9	0.414	0.109	+output,delivery,output d...	87	125
Multiple	10	0.416	0.105	+statement,+data set,+s...	53	74
Multiple	11	0.410	0.107	java,+client,+component,...	82	133
Multiple	12	0.418	0.108	+decision,+support,infor...	99	132
Multiple	13	0.377	0.106	+performance,+server,sc...	93	121
Multiple	14	0.411	0.105	+report,print,+macro,+pr...	92	141
Multiple	15	0.331	0.105	+model,regression,+vari...	105	136
Multiple	16	0.339	0.102	+group,+treatment,+trial,...	94	109
Multiple	17	0.362	0.104	health,information,+care,...	116	145
Multiple	18	0.386	0.104	+analysis,+interface,data...	114	172
Multiple	19	0.390	0.104	+function,+macro,+progr...	109	182
Multiple	20	0.356	0.103	+entry,+frame,+develope...	108	145
Multiple	21	0.341	0.097	+entry,+catalog,+catalog ...	77	91
Multiple	22	0.315	0.098	+year,institute,+presentat...	121	147
Multiple	23	0.264	0.095	+programmer,+statemen...	98	163
Multiple	24	0.259	0.095	+program,+development,...	121	154
Multiple	25	0.255	0.092	+version,integration,+ent...	97	153


Note: If you ran the optional spell-checking step in “Using the Text Filter Node” on page 59, then the Topics shown here and those represented in subsequent steps might be different from what you see.

5. Select the Number of Documents by Topics chart to see a topic by the number of documents that it contains.



Note: You might need to resize the default graph to see the topic ID values.

In addition to multi-term topics, you can use the **Text Topic** node to create single-term topics or to create your own topics.

6. Close the Results window, and select the **Text Topic** node.
7. Select the **Number of Single-term Topics** property, type *10*, and press **Enter** on your keyboard.
8. Click the  for the **User Topics** property.
9. In the User Topics dialog box, click the **Add** button twice to create two rows. Enter the terms *company* and *president*, each with a weight of *0.5*, and specify the topic *company and president* for both.

EMWS3.TextTopic_INITTOPICS

Topic	Term	Role	Weight
company and president	company		0.5
company and president	president		0.5

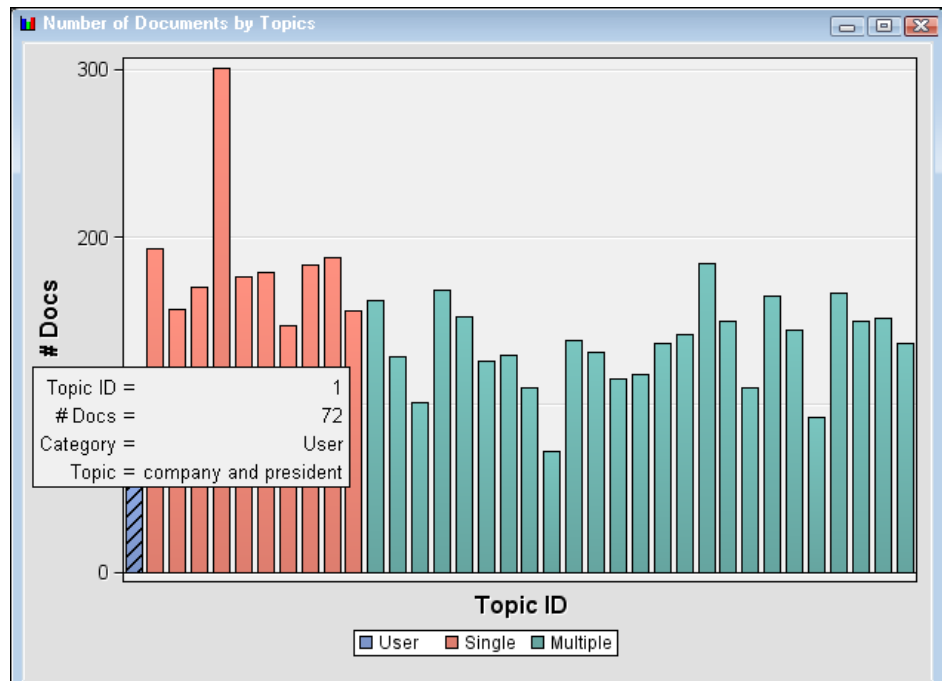
Import OK Cancel

10. Click **OK**.
11. Right-click the **Text Topic** node and select **Run**. Select **Yes** in the Confirmation dialog box, and then **Results** in the Run Status dialog box when the node finishes running.


12. Select the Topics table. Notice that 10 new single-term topics have been created along with the topic that you specified in the User Topics dialog box.

Topics						
Category ▼	Topic ID	Document Cutoff	Term Cutoff	Topic	Number of Terms	# Docs
User	1	0.001	0.001	company and president	1	72
Single	2	0.001	0.001	+macro	1	193
Single	3	0.001	0.001	+report	1	157
Single	4	0.001	0.001	+data set	1	170
Single	5	0.001	0.001	information	1	301
Single	6	0.001	0.001	web	1	176
Single	7	0.001	0.001	+variable	1	179
Single	8	0.001	0.001	+server	1	147
Single	9	0.001	0.001	+program	1	183
Single	10	0.001	0.001	+technique	1	188
Single	11	0.001	0.001	access	1	156
Multiple	12	0.515	0.122	+data set,+variable,+va...	81	162
Multiple	13	0.481	0.117	+warehouse,+data war...	79	128
Multiple	14	0.438	0.111	+macro,+macro variabl...	45	101
Multiple	15	0.485	0.117	+server,+mainframe,+c...	89	168
Multiple	16	0.457	0.115	web,+page,+browser,i...	85	152
Multiple	17	0.455	0.117	+sample,+analysis,+te...	109	126
Multiple	18	0.211	0.092	+program,+date,+time,...	106	129
Multiple	19	0.403	0.109	+graph,graph software,...	84	110
Multiple	20	0.418	0.104	+statement,+data set,+	54	72
Multiple	21	0.413	0.108	java,+client,+server,+co...	83	138
Multiple	22	0.423	0.108	+decision,+support,inf...	98	131
Multiple	23	0.404	0.105	+output,delivery,output ...	89	115
Multiple	24	0.378	0.106	+performance,+server,...	98	118
Multiple	25	0.409	0.105	+report,print,+macro,+p...	90	136
Multiple	26	0.327	0.105	+model,regression,+va...	107	142
Multiple	27	0.390	0.105	+programmer,+functio...	103	184
Multiple	28	0.370	0.105	health,information,+car...	113	150
Multiple	29	0.336	0.102	+group,+treatment,+tria...	96	110
Multiple	30	0.380	0.104	+analysis,+tool,+proce...	115	165
Multiple	31	0.352	0.102	+entry,+frame,+databa...	107	144
Multiple	32	0.338	0.097	+entry,+catalog,+catalo...	76	92
Multiple	33	0.321	0.098	+version,+enhanceme...	97	166
Multiple	34	0.306	0.098	+year,institute,+progra...	115	150
Multiple	35	0.314	0.097	+development,+test,+p...	120	151
Multiple	36	0.223	0.091	+table,+procedure,+for...	108	136

13. Select the Number of Documents by Topics window to see the multi-term, single-term, and user-created topics by the number of documents that they contain.




You can use the Interactive Topic Viewer to view and modify topic properties.

- Close the Results window, and select the **Text Topic** node. Click the  for the **Topic Viewer** property.

Interactive Topic Viewer

File Edit



Topics

Recalculate

Topic	Category	Term Cutoff	Document Cutoff	Number of Terms	# Docs
company and president	User	0.001	0.001	1	72
+macro	Single	0.001	0.001	1	193
+report	Single	0.001	0.001	1	157
+data set	Single	0.001	0.001	1	170
information	Single	0.001	0.001	1	301

Terms

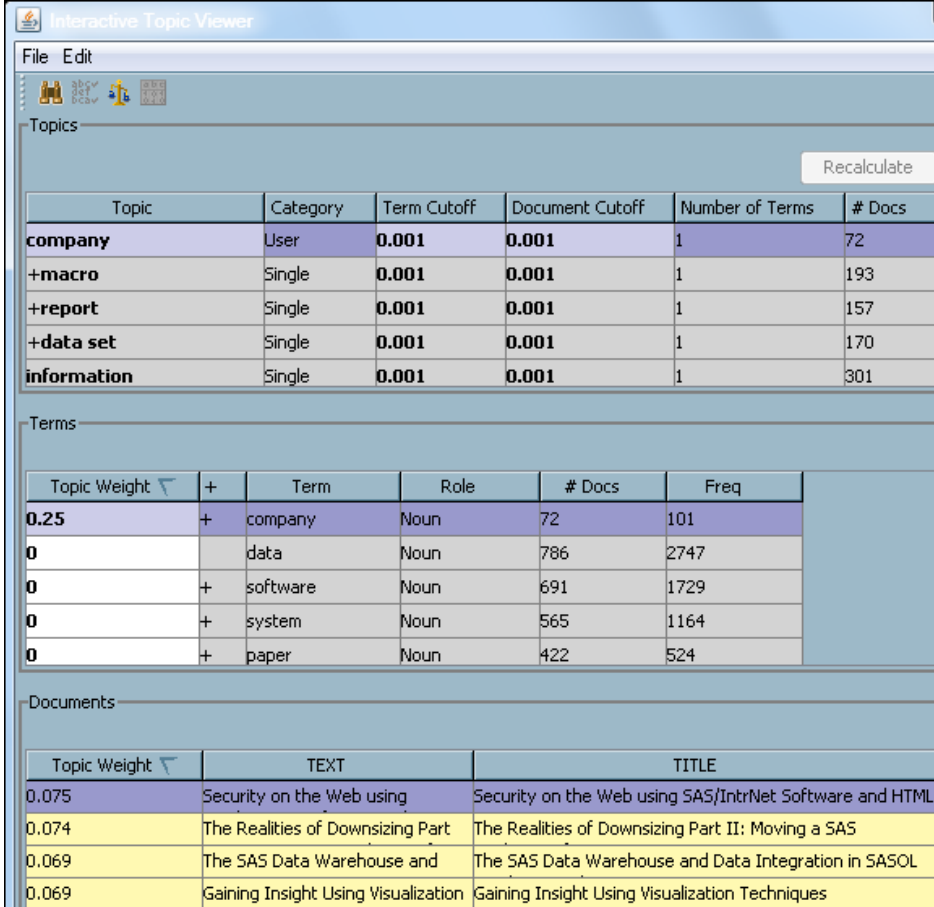
Topic Weight		Term	Role	# Docs	Freq
0.5	+	company	Noun	72	101
0		data	Noun	786	2747
0	+	software	Noun	691	1729
0	+	system	Noun	565	1164
0	+	paper	Noun	422	524

Documents

Topic Weight	TEXT	TITLE
0.149	The Realities of Downsizing Part	The Realities of Downsizing Part II: Moving a SAS
0.149	Security on the Web using	Security on the Web using SAS/IntrNet Software and HTML
0.138	The SAS Data Warehouse and	The SAS Data Warehouse and Data Integration in SASOL
0.138	Gaining Insight Using Visualization	Gaining Insight Using Visualization Techniques

In the Interactive Topic Viewer, you can change the topic name, term and document cutoff values, and the topic weight.

15. Select the topic value “company and president” in the Topics table and rename the topic to *company*. Select the topic weight for the term “company” in the Terms table, and change it to 0.25. Click **Recalculate**.



The screenshot shows the Interactive Topic Viewer application. It has a menu bar with 'File' and 'Edit', and a toolbar with icons for file operations and calculations. The main area is divided into three sections: Topics, Terms, and Documents.

Topics Table:

Topic	Category	Term Cutoff	Document Cutoff	Number of Terms	# Docs
company	User	0.001	0.001	1	72
+macro	Single	0.001	0.001	1	193
+report	Single	0.001	0.001	1	157
+data set	Single	0.001	0.001	1	170
information	Single	0.001	0.001	1	301

A 'Recalculate' button is located to the right of the Topics table.

Terms Table:

Topic Weight	+	Term	Role	# Docs	Freq
0.25	+	company	Noun	72	101
0		data	Noun	786	2747
0	+	software	Noun	691	1729
0	+	system	Noun	565	1164
0	+	paper	Noun	422	524

Documents Table:

Topic Weight	TEXT	TITLE
0.075	Security on the Web using	Security on the Web using SAS/IntrNet Software and HTML
0.074	The Realities of Downsizing Part	The Realities of Downsizing Part II: Moving a SAS
0.069	The SAS Data Warehouse and	The SAS Data Warehouse and Data Integration in SASOL
0.069	Gaining Insight Using Visualization	Gaining Insight Using Visualization Techniques

16. Close the Interactive Topic Viewer, and select **No** when prompted to save your changes. For more information about the Interactive Topic Viewer, see the Interactive Topic Viewer topic in the SAS Text Miner help.

Chapter 12

The Text Cluster Node

About the Text Cluster Node	71
Using the Text Cluster Node	71

About the Text Cluster Node

The **Text Cluster** node clusters documents into disjointed sets of documents and reports on the descriptive terms for those clusters. Two algorithms are available. The Expectation Maximization algorithm clusters documents with a flat representation, and the Hierarchical clustering algorithm groups clusters into a tree hierarchy. Both approaches rely on the singular value decomposition (SVD) to transform the original weighted, term-document frequency matrix into a dense but low dimensional representation.

For more information about the **Text Cluster** node, see the SAS Text Miner Help.

The rest of this chapter presents an example of how you can use the **Text Cluster** node.

Using the Text Cluster Node

This example uses the **Text Cluster** node to cluster SAS Users Group International (SUGI) abstracts. This example assumes that SAS Enterprise Miner is running, and that a diagram workspace has been opened in a project. For information about creating a project and a diagram, see [“Setting Up Your Project” on page 9](#).



Note: SAS Users Group International is now SAS Global Forum.

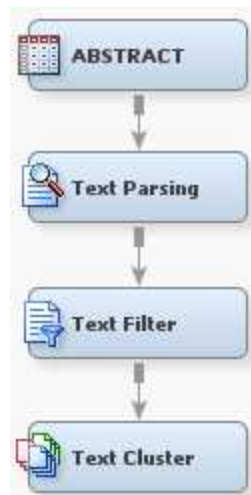
Perform the following steps:

1. Create a data source for SAMPSIO.ABSTRACT. Change the Role of the variable TITLE to **ID**.

Note: The SAMPSIO.ABSTRACT data set contains information about 1,238 papers that were prepared for meetings of SUGI from 1998 through 2001 (SUGI 23 through 26). The variable TITLE is the title of the SUGI paper. The variable TEXT contains the abstract of the SUGI paper.

2. Add the SAMPSIO.ABSTRACT data source to the diagram workspace.

3. Select the **Text Mining** tab on the Toolbar, and drag a **Text Parsing** node into the diagram workspace.
4. Connect the **Input Data** node to the **Text Parsing** node.
5. Select the **Text Parsing** node, and then click the  for the **Stop List** property.
6. Click the **Import** button, browse to select SAMPSIO.SUGISTOP as the stop list, and then click **OK**. Click **OK** to exit the dialog box for the **Stop List** property.
7. Set the **Find Entities** property to **Standard**.
8. Click the  for the **Ignore Types of Entities** property to open the Ignore Types of Entities dialog box.
9. Select all entity types except for: **Location**, **Organization**, **Person**, and **Product**. Click **OK**.
10. Select the **Text Mining** tab, and drag a **Text Filter** node into the diagram workspace.
11. Connect the **Text Parsing** node to the **Text Filter** node.
12. Select the **Text Mining** tab, and drag a **Text Cluster** node into the diagram workspace.
13. Connect the **Text Filter** node to the **Text Cluster** node. Your process flow diagram should resemble the following:

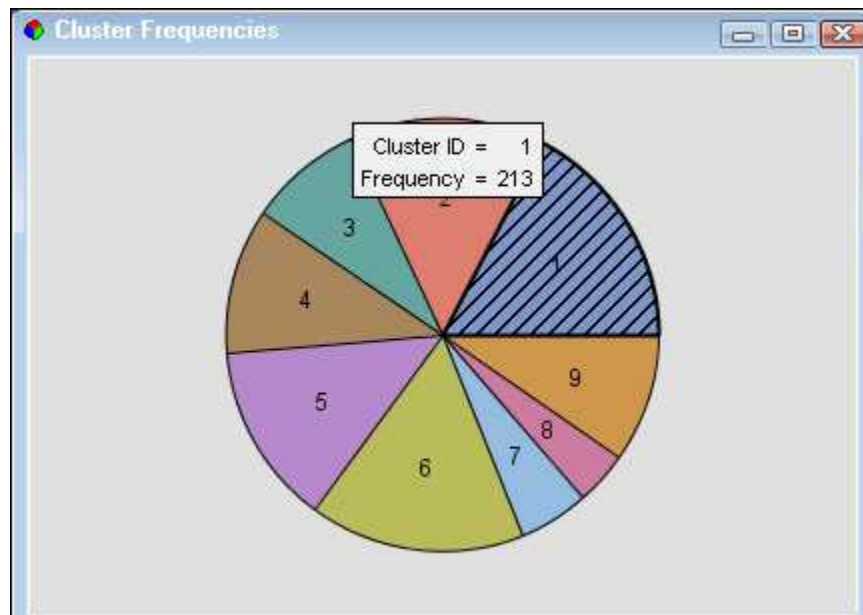


14. Right-click the **Text Cluster** node and select **Run**. Click **Yes** in the Confirmation dialog box.
15. Click **Results** in the Run Status dialog box when the node has finished running.
16. Select the Clusters table.

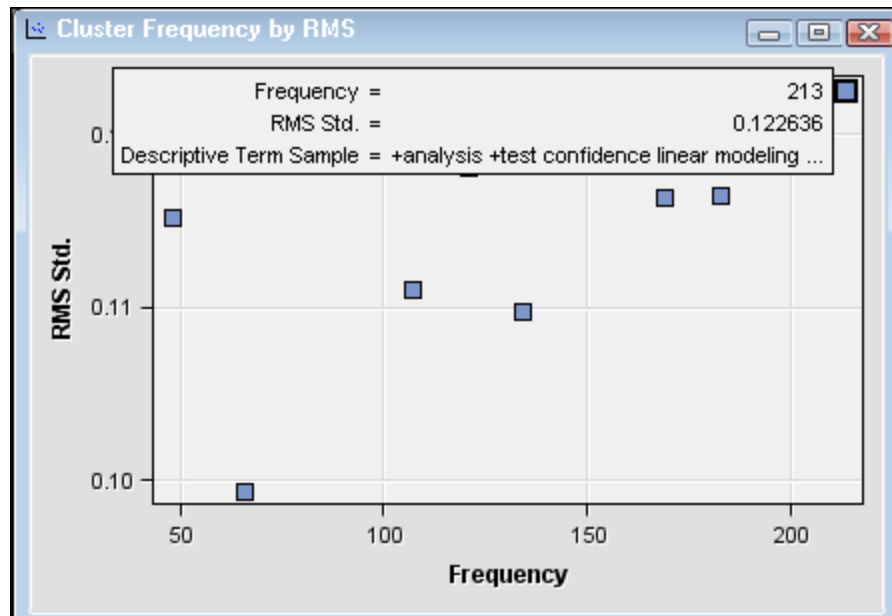
The Clusters table contains an ID for each cluster, the descriptive terms that make up that cluster, and statistics for each cluster.

Clusters			
Cluster ID	Descriptive Terms	Frequency	Percentage
1	+analysis +test confidence linear modeling models regression...	213	17%
2	'data set' +macro +report +set +statement +step variables +pr...	183	15%
3	+client +performance +server scalable tuning +version integrati...	107	9%
4	'output delivery system' +browser +output delivery html intrnet o...	134	11%
5	'data warehousing' +data warehouse' +business +customer +...	169	14%
6	programs +function windows operating functions +program file...	197	16%
7	'graph software' +graph charts graphs graphics +annotate +pro...	66	5%
8	'data entry' +entry +screen dates +date +frame +find +change w...	48	4%
9	+object af objects developers +development +database +fram...	121	10%

17. Select the first cluster in the Clusters table.
18. Select the Cluster Frequencies window to see a pie chart of the clusters by frequency. Position the mouse pointer over a section to see the frequency for that cluster in a tooltip.



19. Select the Cluster Frequency by RMS window, and then position the mouse pointer over the highlighted cluster.



The frequency of the first cluster is the highest, but how does it compare to the other clusters in terms of distance?

20. Select the Distance Between Clusters window. Then position the mouse pointer over the highlighted cluster to see the position of the first cluster in an X and Y coordinate grid.



Position the mouse pointer over other clusters to compare distances.

21. Close the Results window.

Now compare the clustering results that were obtained with the Expectation-Maximization clustering algorithm with using a Hierarchical clustering algorithm.

22. Select the **Text Cluster** node.

23. Select **Exact** for the **Exact or Maximum Number** property.

24. Specify *10* for the **Number of Clusters** property.
25. Select **Hierarchical** for the **Cluster Algorithm** property.
26. Right-click the **Text Cluster** node and select **Run**. Click **Yes** in the Confirmation dialog box.
27. Click **Results** in the Run Status dialog box when the node has finished running.
28. Select the Clusters table.

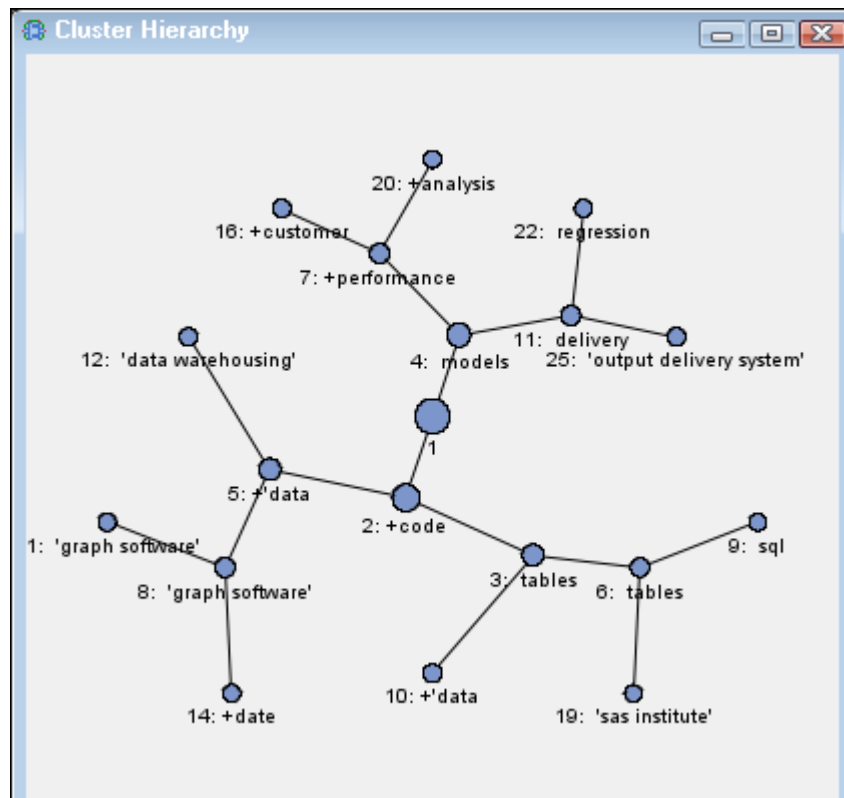
Cluster ID	Descriptive Terms	Frequency	Percentage
9	sql +database statistics +group +proce...	101	8%
10	'data set' +set +program sets +langua...	163	13%
12	'data warehousing' +data warehouse' +...	128	10%
14	+date functions +variable +find +functio...	121	10%
16	+customer +performance +server custo...	142	11%
19	'sas institute' institute windows +develo...	98	8%
20	+analysis clinical mixed models statisti...	108	9%
21	'graph software' +graph charts graphics...	128	10%
22	regression confidence models tests +p...	96	8%
25	'output delivery system' +browser +outp...	153	12%

Notice that while there are 10 clusters in the table, the Cluster IDs do not range from 1 to 10.

29. Select the Hierarchy Data table for more information about the clusters that appear in the Clusters table.

Hierarchy Level ▲	Cluster ID	Parent	Descriptive Terms	Frequency	Graph Description
1	1	.		12381	
2	2	1	+code +program tables windo...	7392	+code
2	4	1	models web output +perform...	4994	models
3	3	2	tables windows +'data set' +p...	3623	tables
3	5	2	'data warehouse' +graph +w...	3775	'data
3	7	4	+performance models +analy...	2507	+perfor...
3	11	4	delivery html intrnet output pa...	24911	delivery
4	6	3	tables windows +table sql +gr...	1996	tables
4	10	3	'data set' +set +program set...	16310	'data
4	8	5	'graph software' +graph functi...	2498	'graph ...
4	12	5	'data warehousing' +'data war...	12812	'data ...

30. Select the Cluster Hierarchy table for a hierarchical graphical representation of the clusters.



31. Close the Results window.

Chapter 13

The Text Rule Builder Node

About the Text Rule Builder Node	77
Using the Text Rule Builder Node	77

About the Text Rule Builder Node

The **Text Rule Builder** node generates an ordered set of rules from small subsets of terms that together are useful in describing and predicting a target variable. Each rule in the set is associated with a specific target category. Each target category consists of a conjunction that indicates the presence or absence of one or a small subset of terms (for example, “term1” AND “term2” AND (NOT “term3”)). A particular document matches this rule if and only if it contains at least one occurrence of term1 and of term2 but no occurrences of term3.

This set of derived rules creates a model that is both descriptive and predictive. When categorizing a new document, the model will proceed through the ordered set and choose the target that is associated with the first rule that matches that document. The rules are provided in the syntax that can be used within SAS Content Categorization Studio, and can be deployed there.

For more information about the **Text Rule Builder** node, see the SAS Text Miner Help.

The rest of this chapter presents an example of how you can use the **Text Rule Builder** node.

Using the Text Rule Builder Node

This example uses the SAMPSIO.NEWS data set to show you how to predict a categorical target variable with the **Text Rule Builder** node. The results will also show that the model is highly interpretable and useful for explanatory and summary purposes as well. This example assumes that SAS Enterprise Miner is running, and that a diagram workspace has been opened in a project. For information about creating a project and a diagram, see [“Setting Up Your Project” on page 9](#).

The SAMPSIO.NEWS data set consists of 600 brief news articles. Most of the news articles fall into one of these categories: computer graphics, hockey, and medical issues.

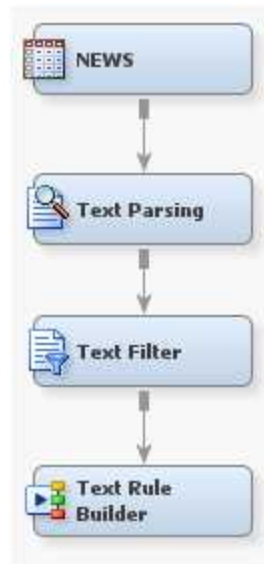
The SAMPSIO.NEWS data set contains 600 observations and the following variables:

- **TEXT** is a nominal variable that contains the text of the news article.
- **graphics** is a binary variable that indicates whether the document belongs to the computer graphics category (1=yes, 0=no).
- **hockey** is a binary variable that indicates whether the document belongs to the hockey category (1=yes, 0=no).
- **medical** is a binary variable that indicates whether the document is related to medical issues (1=yes, 0=no).
- **newsgroup** is a nominal variable that contains the group that a news article fits into.

To use the **Text Rule Builder** node to predict the categorical target variable, **newsgroup**, in the SAMPSIO.NEWS data set:

1. Use the Data Source Wizard to define a data source for the data set SAMPSIO.NEWS.
 - a. Set the measurement levels of the variables **graphics**, **hockey**, and **medical** to **Binary**.
 - b. Set the model role of the variable **newsgroup** to **Target** and leave the roles of **graphics**, **hockey**, and **medical** as **Input**.
 - c. Set the variable **text** to have a role of **Text**.
 - d. Select **No** in the Data Source Wizard — Decision Configuration dialog box.
 - e. Use the default target profile for the target **newsgroup**.
2. After you create the **NEWS** data source, drag it to the diagram workspace.
The **Text Rule Builder** node must be preceded by **Text Parsing** and **Text Filter** nodes.
3. Select the **Text Mining** tab on the toolbar, and drag a **Text Parsing** node into the diagram workspace.
4. Connect the **NEWS** data source to the **Text Parsing** node.
5. Select the **Text Mining** tab on the toolbar, and drag a **Text Filter** node into the diagram workspace.
6. Connect the **Text Parsing** node to the **Text Filter** node.
7. Select the **Text Mining** tab on the toolbar, and drag a **Text Rule Builder** node into the diagram workspace.
8. Connect the **Text Filter** node to the **Text Rule Builder** node.

Your process flow diagram should resemble the following:



9. Select the **Text Rule Builder** node in the process flow diagram.
10. Click the value for the **Generalization Error** property, and select **Very Low**.
11. Click the value for the **Purity of Rules** property, and select **Very Low**.
12. Click the value for the **Exhaustiveness** property, and select **Very Low**.
13. In the diagram workspace, right-click the **Text Rule Builder** node and select **Run**. Click **Yes** in the Confirmation dialog box that appears.
14. Click **Results** in the Run Status dialog box when the node finishes running.
15. Select the Rules Obtained table to see information about the rules that were obtained.
The words in the Rule column have the corresponding estimated precision at implying the target, **newsgroup**.

Rules Obtained						
Target Value	True Positive/Total	Remaining Positive/Total	Rule	Estimated Precision	Sample Precision	Sample Recall
MEDICAL	58/58	200/600	gordon	0.977778	1	0.29
MEDICAL	17/17	142/542	msg	0.922315	1	0.375
MEDICAL	14/14	125/525	treat	0.904762	1	0.445
MEDICAL	11/11	111/511	medicine	0.879572	1	0.5
MEDICAL	10/10	100/500	pain	0.866667	1	0.55
MEDICAL	10/10	90/490	merrill	0.863946	1	0.6
MEDICAL	7/7	80/480	health	0.814815	1	0.635
MEDICAL	7/7	73/473	treatment	0.812074	1	0.67
MEDICAL	5/5	66/466	symptom	0.754752	1	0.695
MEDICAL	5/5	61/461	study	0.752092	1	0.72
MEDICAL	4/4	56/456	infection	0.707602	1	0.74
MEDICAL	5/6	52/452	normal	0.653761	0.993506	0.765
MEDICAL	3/3	47/446	diet	0.642152	0.993631	0.78
MEDICAL	7/10	44/443	drug	0.599887	0.976048	0.815
MEDICAL	4/5	37/433	russell	0.595843	0.97093	0.835
MEDICAL	4/4	33/428	amount & ~team	0.544977	0.971591	0.855
MEDICAL	2/2	29/424	antibiotic	0.534198	0.97191	0.865
MEDICAL	2/2	27/422	med	0.531991	0.972222	0.875
MEDICAL	2/2	25/420	kekule	0.529762	0.972527	0.885
MEDICAL	2/3	23/418	disease	0.42201	0.967568	0.895

In the second column above, the True Positive (the first number) is the number of documents that were correctly assigned to the rule. The Total (the second number) is the total positive.

In the third column above, the Remaining Positive (the first number) is the total number of remaining documents in the category. The Total (the second number) is the total number of documents remaining.

In the above example, in the first row, 200 documents have been assigned to the MEDICAL newsgroup, and 600 total documents exist in the data set. Fifty-eight of the documents were assigned to the rule “gordon” (58 were correctly assigned). This means that if a document contains the word “gordon,” and you assign all those documents to the MEDICAL newsgroup, 58 out of 58 will be assigned correctly. In the next row, there are $200 - 58 = 142$ MEDICAL newsgroup documents left that can be evaluated for rule assignment, out of a total of $600 - 58 = 542$ documents. In this second row, 17 documents are correctly assigned to the rule “msg.” This means that if a document contains the term “msg,” and you assign all those documents to the MEDICAL newsgroup, 17 out of 17 will be assigned correctly.

Most of the rules are single term rules because the NEWS data set is limited in size. However, there is one multiple term rule above. In the 16th row, the rule “amount & ~team” means that if a document contains the word “amount” and does not contain the word “team,” then 4 of the remaining documents will be correctly assigned to the MEDICAL newsgroup.

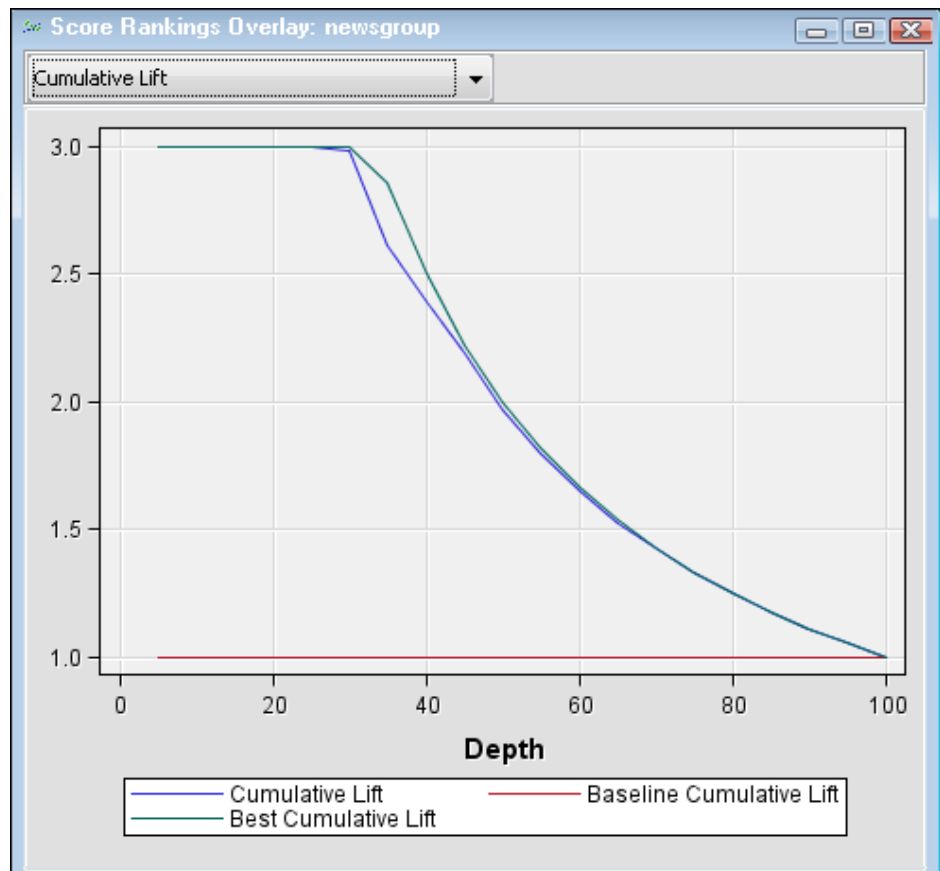
Note: ~ means logical not.

16. Select the Score Rankings Overlay graph to view the following types of information about the target variable:

- Cumulative Lift
- Lift
- Gain
- % Response
- Cumulative % Response

- % Captured Response
- Cumulative % Captured Response


Note: To change the statistic, select one of the above options from the drop-down menu.



17. Select the Fit Statistics window for statistical information about the target variable, **newsgroup**.

Fit Statistics					
Target	Fit Statistics	Statistics Label	Train	Validation	Test
newsgroup	_ASE_	Average Squared Error	0.033461	.	.
newsgroup	_DIV_	Divisor for ASE	1800	.	.
newsgroup	_MAX_	Maximum Absolute Error	1	.	.
newsgroup	_NOBS_	Sum of Frequencies	600	.	.
newsgroup	_RASE_	Root Average Squared Error	0.182924	.	.
newsgroup	_SSE_	Sum of Squared Errors	60.23006	.	.
newsgroup	_DISF_	Frequency of Classified Cases	600	.	.
newsgroup	_MISC_	Misclassification Rate	0.07	.	.
newsgroup	_WRONG_	Number of Wrong Classifications	42	.	.

18. Close the Results window.
19. Click the value for the **Generalization Error** property, and select **Medium**.
20. Click the value for the **Purity of Rules** property, and select **Medium**.
21. Click the value for the **Exhaustiveness** property, and select **Medium**.

22. Select the **NEWS** data source.
23. Click the  for the **Variables** property.
24. Change the role of the **HOCKEY** variable to **Target**, and change the role of the **NEWSGROUP** variable to **Input**.
25. Click **OK**.
26. In the diagram workspace, right-click the **Text Rule Builder** node and select **Run**. Click **Yes** in the Confirmation dialog box that appears.
27. Click **Results** in the Run Status dialog box when the node finishes running.
28. Select the Rules Obtained table to see information about the rules that predicted the target — the **HOCKEY** newsgroup.

The words in the Rule column have the corresponding estimated precision at implying the hockey target.

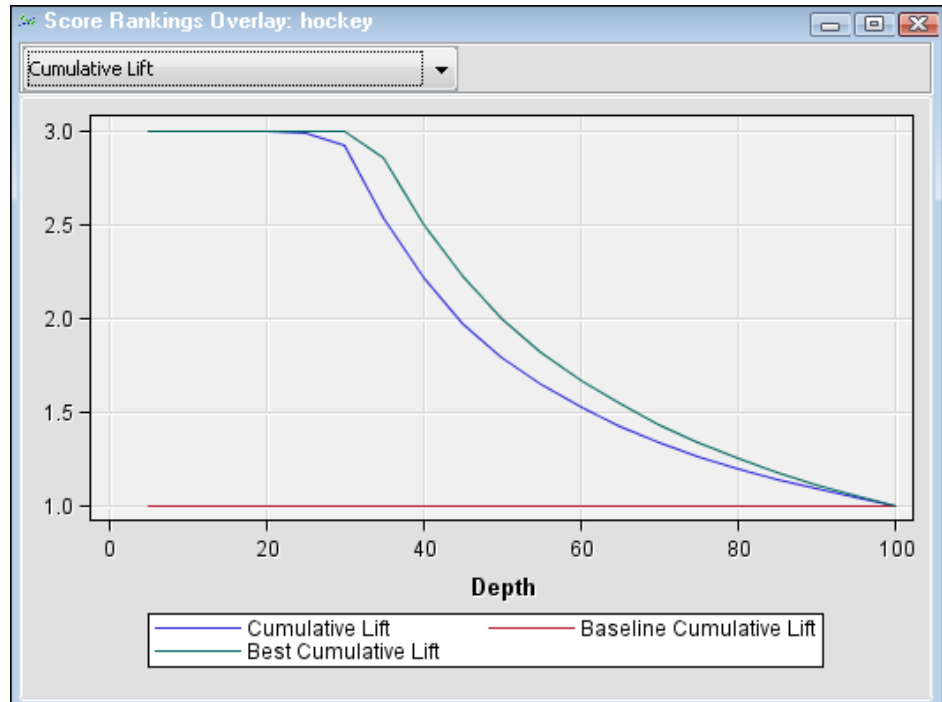
Target Value	True Positive/Total	Remaining Positive/Total	Rule	Estimated Precision	Sample Precision	Sample Recall
1	69/70	200/600	team	0.918803	0.985714	0.345
1	23/23	131/530	hockey	0.805721	0.989247	0.46
1	13/13	108/507	cup	0.700197	0.990566	0.525
1	11/11	95/494	playoff	0.659919	0.991453	0.58
1	10/10	84/483	lemieux	0.63285	0.992126	0.63
1	9/9	74/473	sfu	0.603034	0.992647	0.675
1	8/8	65/464	uwaterloo	0.570043	0.993056	0.715
1	9/10	57/456	fan	0.555556	0.987013	0.76
1	6/6	48/446	montreal	0.49007	0.9875	0.79
1	5/7	42/440	player	0.384242	0.976048	0.815
1	3/3	37/433	ranger	0.334873	0.976471	0.83
1	3/4	34/430	goal	0.302713	0.971264	0.845
1	2/2	31/426	laurentian	0.258216	0.971591	0.855
1	2/2	29/424	ucs	0.254717	0.97191	0.865
1	2/2	27/422	belfour	0.251185	0.972222	0.875
1	2/2	25/420	gerald	0.247619	0.972527	0.885
0	96/96	395/418	know	0.995767	1	0.243038

In the above example, in the first row, 200 documents have been assigned to the **HOCKEY** newsgroup, and 600 total documents exist in the data set. The target value is **1**, instead of “**HOCKEY**,” because you set the **hockey** variable to be the target instead of the **newsgroup** variable. 70 of the documents were assigned to the rule “team” (69 were correctly assigned). This means that if a document contains the word “team,” and you assign all those documents to the **HOCKEY** newsgroup, 69 out of 70 will be assigned correctly. In the next row, there are $200 - 69 = 131$ **HOCKEY** documents left that can be evaluated for rule assignment, out of a total of $600 - 70 = 530$ documents. In this second row, 23 documents are correctly assigned to the rule “hockey.” This means that if a document contains the word “hockey,” and you assign all those documents to the **HOCKEY** newsgroup, 23 out of 23 will be assigned correctly.

29. Select the Score Rankings Overlay graph to view the following types of information about the target variable:
 - Cumulative Lift


- Lift
- Gain
- % Response
- Cumulative % Response
- % Captured Response
- Cumulative % Captured Response

Note: To change the statistic, select one of the above options from the drop-down menu.



30. Select the Fit Statistics table for statistical information about the hockey target variable.

Fit Statistics					
Target	Fit Statistics	Statistics Label	Train	Validation	Test
hockey	_ASE_	Average Squared Error	0.006483	.	.
hockey	_DIV_	Divisor for ASE	1200	.	.
hockey	_MAX_	Maximum Absolute Error	0.383761	.	.
hockey	_NOBS_	Sum of Frequencies	600	.	.
hockey	_RASE_	Root Average Squared Error	0.080515	.	.
hockey	_SSE_	Sum of Squared Errors	7.779142	.	.
hockey	_DISF_	Frequency of Classified Cases	600	.	.
hockey	_MISC_	Misclassification Rate	0.046667	.	.
hockey	_WRONG_	Number of Wrong Classificati...	28	.	.

31. Close the Results window.
32. Click the  for the **Content Categorization Code** property.

The Content Categorization Code window appears. The code provided in this window is the code that is output for SAS Content Categorization and is ready for compilation.

33. Click **Cancel**.

34. Click the  for the **Change Target Values** property.

The Change Target Values window appears.

You can use the Change Target Values window to improve the model.

35. Select one or more cells in the **Assigned Target** column, and select a new target value.

36. Click **OK**.

37. Rerun the **Text Rule Builder** node, and then check whether your model has been improved.

Chapter 14

Tips for Text Mining

Processing a Large Collection of Documents	85
Dealing with Long Documents	85
Processing Documents from an Unsupported Language or Encoding	86

Processing a Large Collection of Documents

Using SAS Text Miner nodes to process a large collection of documents can require a lot of computing time and resources. If you have limited resources, it might be necessary to take one or more of the following actions:

- Use a sample of the document collection.
 - Set some of the parse properties to **No** or **None**, such as **Noun Groups** or **Find Entities**.
 - Reduce the number of SVD dimensions or roll-up terms. If you are running into memory problems with the SVD approach, you can roll up a certain number of terms, and then the remaining terms are automatically dropped.
 - Limit parsing to high information words by turning off all parts of speech other than nouns, proper nouns, noun groups, and verbs.
 - Structure sentences properly for best results, including correct grammar, punctuation, and capitalization. Entity extraction does not always generate reasonable results.
-

Dealing with Long Documents

SAS Text Miner uses the "bag-of-words" approach to represent documents. That means that documents are represented with a vector that contains the frequency with which each term occurs in each document. In addition, word order is ignored. This approach is very effective for short, paragraph-sized documents, but it can cause a harmful loss of information with longer documents. You might want to consider preprocessing your long documents in order to isolate the content that is really of use in your model. For example, if you are analyzing journal papers, you might find that analyzing only the abstract gives the best results. Consider using the SAS DATA step or an alternative programming language such as Perl to extract the relevant content from long documents.

Processing Documents from an Unsupported Language or Encoding

If you have a collection of documents from an unsupported language or encoding, you might still be able to successfully process the text and get useful results. Follow these steps:

1. Set the language to **English**.
2. Turn off these parse properties:
 - **Different Parts of Speech**
 - **Noun Groups**
 - **Find Entities**
 - **Stem Terms**
3. Run the **Text Parsing** node.

Chapter 15

Next Steps: A Quick Look at Additional Features

The %TEXTSYN Macro	87
The %TMFILTER Macro	87

The %TEXTSYN Macro

The %TEXTSYN macro is provided with SAS Text Miner. You can use this macro after a **Text Parsing** node has been successfully run to find and correct misspellings that appear in the input data source. It is not supported for use with the Chinese language.

The macro creates a synonym data set, which you can use in SAS Text Miner, that contains misspelled terms and candidate parents (correctly spelled terms). The data set includes the variables “term,” “parent,” and “category.” Using optional arguments, you can also specify that the synonym data set include example usages (from up to two documents) of the misspelled terms.

See the SAS Text Miner Help for more information about the %TEXTSYN macro.

The %TMFILTER Macro

The %TMFILTER macro is a SAS macro that enables you to convert files into SAS data sets. The %TMFILTER macro is provided with SAS Text Miner. It is supported in all operating systems for filtering and on Windows for crawling. The %TMFILTER macro relies on the SAS Document Conversion Server that is installed and running on a Windows machine. See SAS Document Conversion server for more information. You can use the macro to perform the following tasks:

- filter a collection of documents that is saved in any supported file format and output a SAS data set that can be used to create a SAS Text Miner data source.
- Web crawl and output a SAS data set that can be used to create a SAS Text Miner data source. Web crawling retrieves the text of a starting Web page, extracts the URL links within that page, and then repeats the process within the linked pages recursively. You can restrict a crawl to the domain of the starting URL, or you can let a crawl process any linked pages that are not in the domain of the starting URL. The crawl continues until a specified number of levels of drill-down is reached or until all the Web pages that satisfy the domain constraint are found. Web crawling is supported only on Windows operating systems.

- identify the languages of all documents in a collection.

See the SAS Text Miner Help for more information about the %TMFILTER macro.

Glossary

catalog directory

a part of a SAS catalog that stores and maintains information about the name, type, description, and update status of each member of the catalog.

clustering

the process of dividing a data set into mutually exclusive groups so that the observations for each group are as close as possible to one another and different groups are as far as possible from one another. In SAS Text Miner, clustering involves discovering groups of documents that are more similar to each other than they are to the rest of the documents in the collection. When the clusters are determined, examining the words that occur in the cluster reveals the focus of the cluster. Forming clusters within the document collection can help you understand and summarize the collection without reading every document. The clusters can reveal the central themes and key concepts that are emphasized by the collection.

concept linking

finding and displaying the terms that are highly associated with the selected term in the Terms table.

data source

a data object that represents a SAS data set in the Java-based Enterprise Miner GUI. A data source contains all the metadata for a SAS data set that Enterprise Miner needs in order to use the data set in a data mining process flow diagram. The SAS data set metadata that is required to create an SAS Enterprise data source includes the name and location of the data set; the SAS code that is used to define its library path; and the variable roles, measurement levels, and associated attributes that are used in the data mining process.

diagram

See process flow diagram.

entity

any of several types of information that SAS Text Miner is able to distinguish from general text. For example, SAS Text Miner can identify names (of people, places, companies, or products, for example), addresses (including street addresses, post office addresses, e-mail addresses, and URLs), dates, measurements, currency amounts, and many other types of entities.

libref

a name that is temporarily associated with a SAS library. The complete name of a SAS file consists of two words, separated by a period. The libref, which is the first

word, indicates the library. The second word is the name of the specific SAS file. For example, in VLIB.NEWBDAY, the libref VLIB tells SAS which library contains the file NEWBDAY. You assign a libref with a LIBNAME statement or with an operating system command.

model

a formula or algorithm that computes outputs from inputs. A data mining model includes information about the conditional distribution of the target variables, given the input variables.

node

(1) in the SAS Enterprise Miner user interface, a graphical object that represents a data mining task in a process flow diagram. The statistical tools that perform the data mining tasks are called nodes when they are placed on a data mining process flow diagram. Each node performs a mathematical or graphical operation as a component of an analytical and predictive data model. (2) in a neural network, a linear or nonlinear computing element that accepts one or more inputs, computes a function of the inputs, and optionally directs the result to one or more other neurons. Nodes are also known as neurons or units. (3) a leaf in a tree diagram. The terms leaf, node, and segment are closely related and sometimes refer to the same part of a tree.

parsing

to analyze text for the purpose of separating it into its constituent words, phrases, multiword terms, punctuation marks, or other types of information.

partitioning

to divide available data into training, validation, and test data sets.

process flow diagram

a graphical representation of the various data mining tasks that are performed by individual Enterprise Miner nodes during a data mining analysis. A process flow diagram consists of two or more individual nodes that are connected in the order in which the data miner wants the corresponding statistical operations to be performed. Short form: PFD.

roll-up terms

the highest-weighted terms in the document collection.

SAS data set

a file whose contents are in one of the native SAS file formats. There are two types of SAS data sets: SAS data files and SAS data views. SAS data files contain data values in addition to descriptor information that is associated with the data. SAS data views contain only the descriptor information plus other information that is required for retrieving data values from other SAS data sets or from files that are stored in other software vendors' file formats.

scoring

the process of applying a model to new data in order to compute output. Scoring is the last process that is performed in data mining.

segmentation

the process of dividing a population into sub-populations of similar individuals. Segmentation can be done in a supervisory mode (using a target variable and various techniques, including decision trees) or without supervision (using clustering or a Kohonen network).

singular value decomposition

a technique through which high-dimensional data is transformed into lower-dimensional data.

source-level debugger

an interactive environment in SAS that enables you to detect and resolve logical errors in programs that are being developed. The debugger consists of windows and a group of commands.

stemming

the process of finding and returning the root form of a word. For example, the root form of grind, grinds, grinding, and ground is grind.

stop list

a SAS data set that contains a simple collection of low-information or extraneous words that you want to remove from text mining analysis.

test data

currently available data that contains input values and target values that are not used during training, but which instead are used for generalization and model comparisons.

training data

currently available data that contains input values and target values that are used for model training.

validation data

data that is used to validate the suitability of a data model that was developed using training data. Both training data sets and validation data sets contain target variable values. Target variable values in the training data are used to train the model. Target variable values in the validation data set are used to compare the training model's predictions to the known target values, assessing the model's fit before using the model to score new data.

variable

a column in a SAS data set or in a SAS data view. The data values for each variable describe a single characteristic for all observations. Each SAS variable can have the following attributes: name, data type (character or numeric), length, format, informat, and label.

Index

A

accessibility features [4](#)

C

cleaning data [31](#)
 creating a synonym data set [34](#)
 examining results using merged
 synonym data sets [37](#)
 using a synonym data set [32](#)
 compatibility [4](#)
 converting files into data sets [87](#)

D

data cleaning
 See [cleaning data](#)
 Data Partition node [16](#)
 data segments [24](#)
 data sets
 converting files into [87](#)
 importing [31](#)
 merged synonym data sets [37](#)
 synonym data sets [32, 34](#)
 data source
 creating for projects [12](#)
 descriptive mining [1](#)
 diagrams
 creating [13](#)
 document analysis [3](#)
 document requirements [1](#)
 documents
 from unsupported language or encoding
 [86](#)
 large collection of [85](#)
 long [85](#)

E

encoding
 unsupported [86](#)

F

file preprocessing [3](#)
 files
 converting into data sets [87](#)

H

Help [7](#)

I

importing data sets [31](#)
 input data
 identifying [15](#)
 partitioning [16](#)

L

languages
 unsupported [86](#)
 large collection of documents [85](#)
 library
 create [10](#)
 long documents [85](#)

M

macros
 %TMFILTER [87](#)
 merged synonym data sets [37](#)
 misspelled terms [34](#)

P

partitioning input data [16](#)
 path for projects [10](#)
 predictive mining [1](#)
 predictive modeling [1](#)
 projects
 creating [9](#)
 creating data source [12](#)

- creating diagrams 13
- path for 10
- setting up 9

R

- results
 - examining with merged synonym data sets 37

S

- SAS Enterprise Miner 12.1 2
- SAS Text Miner 12.1 2
 - accessibility features 4
 - Help 7
- Section 508 standards 4
- segments 24
- stems 33
- SYMPTOM_TEXT variable analysis
 - examining data segments 24
 - identifying input data 15
 - partitioning input data 16
 - setting node properties 16
- synonym data sets
 - creating 34
 - merged 37

T

- text cleaning
 - See [cleaning data](#)
- text mining
 - descriptive mining 1
 - document requirements for 1
 - general order for 3
 - large collection of documents 85
 - long documents 85
 - predictive mining 1
 - process 3
 - tips for 85
 - unsupported language or encoding 86
- text parsing 3
- tips for text mining 85
- transformation 3

U

- unsupported languages or encoding 86

V

- Vaccine Adverse Event Reporting System 5