

# SAS® Contextual Analysis 14.3: User's Guide

The correct bibliographic citation for this manual is as follows: SAS Institute Inc. 2017. SAS® Contextual Analysis 14.3: User's Guide. Cary, NC: SAS Institute Inc.

#### SAS® Contextual Analysis 14.3: User's Guide

Copyright © 2017, SAS Institute Inc., Cary, NC, USA

All Rights Reserved. Produced in the United States of America.

**For a hard copy 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.

The scanning, uploading, and distribution of this book via the Internet or any other means without the permission of the publisher is illegal and punishable by law. Please purchase only authorized electronic editions and do not participate in or encourage electronic piracy of copyrighted materials. Your support of others' rights is appreciated.

**U.S. Government License Rights; Restricted Rights:** The Software and its documentation is commercial computer software developed at private expense and is provided with RESTRICTED RIGHTS to the United States Government. Use, duplication, or disclosure of the Software by the United States Government is subject to the license terms of this Agreement pursuant to, as applicable, FAR 12.212, DFAR 227.7202-1(a), DFAR 227.7202-3(a), and DFAR 227.7202-4, and, to the extent required under U.S. federal law, the minimum restricted rights as set out in FAR 52.227-19 (DEC 2007). If FAR 52.227-19 is applicable, this provision serves as notice under clause (c) thereof and no other notice is required to be affixed to the Software or documentation. The Government's rights in Software and documentation shall be only those set forth in this Agreement.

SAS Institute Inc., SAS Campus Drive, Cary, NC 27513-2414

September 2017

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 trademarks of their respective companies.

14.3-P1:utagsug

### **Contents**

L	Ising This Book	
V	Vhat's New in SAS Contextual Analysis 14.3	
	Accessibility	
Chapter 1 / Introd	duction to SAS Contextual Analysis	
V	Vhat Is SAS Contextual Analysis?	1
	low Does SAS Contextual Analysis Work?	
S	Supported Languages	3
L	Jsing Taxonomies	4
L	Ising the Interface	5
Chapter 2 / Proje	cts in SAS Contextual Analysis	<b>7</b>
C	Overview of a Project	8
F	Preparing the Document Collection	9
	mporting Projects	
C	Creating a New Project	12
Ų	Ising the Properties Page	19
	Sharing Projects	
S	Scoring an External Data Set	25
Α	bout Sentiment Analysis	26
Chapter 3 / Perfo	rming the Analysis Tasks	29
C	Overview of the Analysis Tasks	30
	Ising the Analysis Task Pages	
	Vriting Concept Rules: Basic LITI Syntax	
V	Vriting Category Rules: Boolean Rules	83
Appendix 1 / Par	t-of-Speech Tags (for Languages Other Than English)	93
Ir	troduction to Part-of-Speech and Other Tags	94
Р	art-of-Speech Tags	94

#### iv Contents

Appendix 2 / P	redefined Concept Priorities (for Languages Other	
Than English	n)	137
	Using Priority Values in Predefined Concepts	138
	Priority Values for Predefined Concepts	138
	Recommended Reading	163
	Glossary	165

## **Using This Book**

#### **Audience**

This book is designed for users of SAS Contextual Analysis. It describes the terminology used in SAS Contextual Analysis and provides instructions for tasks.

SAS Contextual Analysis can currently process 31 languages including English. This document is not tailored for any specific language, but English is used in the example text. A subset of predefined concepts is provided for most of the supported languages. For a full list of languages, see "Supported Languages" on page 3.

### **What's New**

# What's New in SAS Contextual Analysis 14.3

#### **Overview**

SAS Contextual Analysis 14.3 includes the following new and enhanced features:

- Additional language support.
- Feature-level sentiment in score code

#### **Details**

#### Language Support and Documentation Enhancements

SAS Contextual Analysis now supports 31 languages, including English. For a list of all languages supported, see "Supported Languages" on page 3.

Included in this document are lists of predefined concepts in all supported languages and their priority values. For more information, see "Concepts" on page 30 and

viii What's New

Appendix 2, "Predefined Concept Priorities (for Languages Other Than English)," on page 137.

#### **Feature-Level Sentiment in Score Code**

Feature-level sentiment score code is now available. To download score code, see "Viewing and Downloading Code" on page 22.

## **Accessibility**

For information about the accessibility of this product, see Accessibility Features of SAS Contextual Analysis 14.3 at support.sas.com.

# Introduction to SAS Contextual Analysis

What Is SAS Contextual Analysis?	1
How Does SAS Contextual Analysis Work?	2
Supported Languages	3
Using Taxonomies	4
Using the Interface	5

### **What Is SAS Contextual Analysis?**

SAS Contextual Analysis is a web-based text analytics application that uses contextual analysis to provide a comprehensive solution to the challenge of identifying and categorizing key textual data. Using this application, you can build models (based on training documents) that automatically analyze and categorize a set of documents. You can then customize your models in order to realize the value of your text-based data.

SAS Contextual Analysis combines the machine-learning capabilities of SAS Text Miner with the rules-based linguistic methods of categorization and extraction in SAS Enterprise Content Categorization. These capabilities, along with document-level sentiment scoring, are combined in a single user interface.

Using SAS Contextual Analysis, you can identify key textual data in your document collections, categorize those data, build concept models, and remove meaningless textual data.

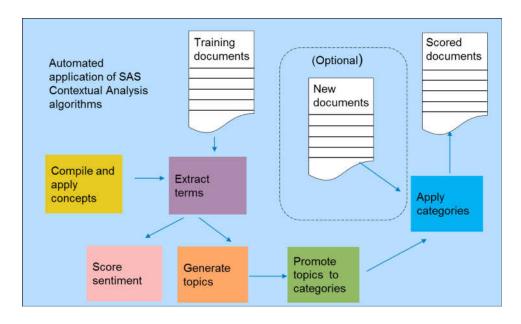
By default, words that provide little or no value are excluded from analysis. Examples of these words include the articles *a*, *an*, and *the* and conjunctions such as *and*, *or*, and *but*. Other terms that are specific to your document collection but provide little or no value are also identified and excluded.

SAS Contextual Analysis is designed for users who have no SAS programming or SAS macro language experience.

# **How Does SAS Contextual Analysis Work?**

Figure 1.1 provides an overview of the SAS Contextual Analysis processes.

Figure 1.1 Process Overview



SAS Contextual Analysis enables you to extract pre-defined concepts or create additional custom concepts that you want to discover in a document or set of documents. For more information about concepts, see "Concepts" on page 30.

The SAS Contextual Analysis algorithms group similar documents in a collection into topics. The documents in each topic often contain similar subject matter, such as motorcycle accidents, computer graphics, or weather patterns. Automatic topic identification enables you to easily categorize each document in your collection.

You can create categories using these methods:

- import categories from SAS Enterprise Content Categorization
- specify category variables in your training documents
- create new categories
- promote topics to categories

Preliminary rules are generated when you promote a topic to a category or when you specify category variables in your training documents. These rules can be edited and refined.

Whether you use the automated processes and rules that are available for extracting terms and subject matter or you customize the processes and rules, context sensitivity is an essential component of your model. To enhance context sensitivity, you can modify the preliminary rules. You can add or modify Boolean operators, characters, and other selections to make the rule matching more context-sensitive.

Finally, you deploy your model to automate the process of classifying a set of input documents.

### **Supported Languages**

Table 1.1 shows a full list of supported languages. See your SAS sales representative for information about licensing additional languages.

 Table 1.1
 SAS Contextual Analysis 14.3 Supported Languages

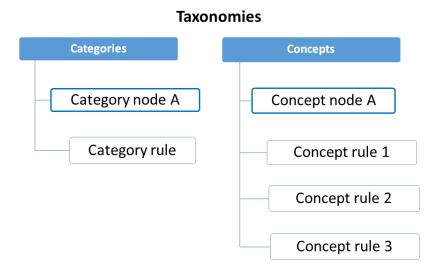
Arabic	Chinese (Simp./Trad.)
Croatian	Czech
Danish	Dutch
English	Farsi
Finnish	French
German	Greek
Hebrew	Hindi
Hungarian	Indonesian
Italian	Japanese
Korean	Norwegian (Bok./Nyn.)
Polish	Portuguese
Romanian	Russian
Slovak	Slovene
Spanish	Swedish
Thai	Turkish
Vietnamese	

### **Using Taxonomies**

In SAS Contextual Analysis, you can create category and concept rules, which are displayed in a taxonomic structure. Each taxonomy consists of a tree of *nodes*. Each

node is a container for rules. By contrast, under a concept node, there can exist multiple rules. Figure 1.2 on page 5 demonstrates how category and concept taxonomies differ.

Taxonomies in SAS Contextual Analysis Figure 1.2

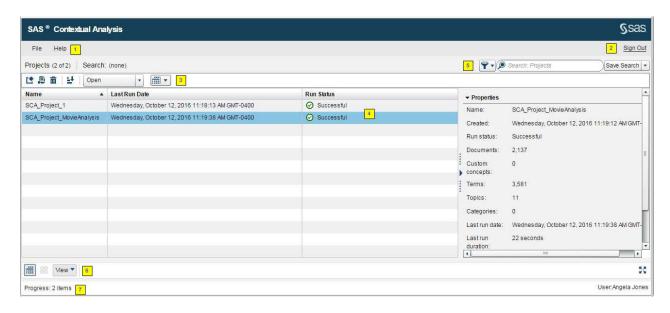


### **Using the Interface**

The main components of the user interface are shown in Figure 1.3.

6 Chapter 1 / Introduction to SAS Contextual Analysis

Figure 1.3 SAS Contextual Analysis Interface



- 1 Application menu
- 2 Sign out
- 3 Application toolbar
- 4 Project list
- 5 Search options
- 6 View options
- 7 Progress panel (click to open)

## Projects in SAS Contextual Analysis

Overview of a Project	8
Preparing the Document Collection	9
Importing Projects Importing an Existing SAS Contextual Analysis Project Model Importing an Existing SAS Enterprise Content Categorization Project Special Considerations When Importing Projects	10
Creating a New Project  Using the Create Project Wizard  Step 1: Identify Project Files, Servers, and Other Properties  Step 2: Specify Term and Synonym Lists  Step 3: Choose Predefined Concepts  Step 4: Identify a Data Source  Step 5: Run the Project	13 14 15
Using the Properties Page Checking Project Status Editing Project Information Viewing and Downloading Code Exporting a Project Model	19
Sharing Projects	. 24
Scoring an External Data Set	25

About Sentiment Analysis	. 26
Introduction to Document Scoring	26
Using SAS Sentiment Analysis Models in SAS	
Contextual Analysis	27

#### **Overview of a Project**

In SAS Contextual Analysis, you create projects, which are basically containers for your data and analysis. A project contains the input data, text mining options, and analysis tasks (working with concepts, terms, topics, and categories). SAS Contextual Analysis is designed so that you can create and run multiple projects simultaneously. Text analysis is performed in the background so that you can open one project while performing analysis on a different project. Projects can be shared among users and updated collaboratively.

When you build a model, you choose input data that contain the document collection that you want to use as a training data set. It is important to ensure that your training data are representative of the data to which this model will be applied. Topics and categories are built based on the terms in this document collection.

Next, you can choose to specify either a start list or a stop list. You can also specify whether to use a synonym list. Before you can run your project on a SAS data set, you must specify the text field that you want to analyze. You can also specify one or more category variables for the analysis.

After the project runs, you can view the terms and automatically discovered topics that were created during the initial text mining. Then, you use the topics to create categories. Categories are groups of documents that contain similar terms. SAS Contextual Analysis builds a set of rules for each category.

#### **Preparing the Document Collection**

Before you create a project in SAS Contextual Analysis, you need to prepare your document collection for analysis. SAS Contextual Analysis enables you to analyze a document collection that is stored as a SAS data set or in text-based file formats such as MS Office, OpenDocument (OpenOffice), PDF, XML, HTML, and others. You can select a SAS data set and then identify the text variables and category variables to be analyzed. Or you can specify a directory that contains the files that you want to use as training data.

When you prepare the input document collection, you should select a set of documents that is representative of the documents that you want to categorize later. The terms that exist in the input document collection are used to build the topics and categories.

There are no standard rules for creating an input document collection. However, the following guidelines can help you prepare your input document collection:

- Include at least 15 to 20 documents for each category that you want to discover.
- Be familiar with the contents of the documents in order to anticipate term discovery and rule creation.
- Do not store SAS data sets in the same directory where you store a collection of Microsoft Word or Adobe PDF documents.

Note: When you use a SAS data set, you must register that data set with the SAS Metadata Server before it is available in SAS Contextual Analysis. You can use SAS Management Console and SAS Enterprise Guide to register data sets. When you use a collection of documents in a folder (rather than a SAS data set), you must locate the folder on the server where the SAS Contextual Analysis workspace server is installed. For more information, see SAS Contextual Analysis: Administrator's Guide.

#### **Importing Projects**

# Importing an Existing SAS Contextual Analysis Project Model

During project creation, you can import a SAS Contextual Analysis project so that you can reuse existing category and concept rules. Imported projects include only category and concept rules; other project components such as topics, terms, data, project settings, and so on, are not imported. The file that you import is in JSON (JavaScript Object Notation) format.

For information about exporting a SAS Contextual Analysis project, see "Exporting a Project Model" on page 23.

# Importing an Existing SAS Enterprise Content Categorization Project

During project creation, you can import an existing SAS Enterprise Content Categorization project for analysis. If you plan to import a project, note the following:

- Concepts that were defined by using the LITI (language interpretation and text interpretation) syntax in an imported SAS Enterprise Content Categorization project can be used in your SAS Contextual Analysis project.
- Categories that were defined using Boolean rules (MCAT syntax) in an imported SAS Enterprise Content Categorization project can be used in your SAS Contextual Analysis project.
- Categories that were created using linguistic rules in SAS Enterprise Content Categorization are not supported.

**Note:** In order for the LITI concepts to be parsed correctly in SAS Contextual Analysis, the parsing priority for disabled concepts must be honored. To ensure this, open your existing project in SAS Enterprise Content Categorization. For any child concept that

was disabled, modify its parent concept so that the parent has a higher priority than the child. Save the project before you import it into SAS Contextual Analysis.

#### **Special Considerations When Importing Projects**

When you import categories and concepts from a SAS Contextual Analysis or SAS Enterprise Content Categorization project into your new project, note that duplicate names might occur. Here are some items to consider when duplicate category or concept names occur during the importing process:

- Check messages to see how SAS Contextual Analysis handled the duplicate name. Click son the tool bar to see messages.
- Plan for possible duplicate names if you are using predefined concepts in your project. For example, suppose a LITI concept that you are importing has the same name as a predefined concept in the current project. In that case, the imported LITI concept's rules are added to the predefined concept's rules.
- When importing concepts from an existing project, the source of the concept determines which concept name takes precedence when the system resolves a naming conflict. Predefined concepts in the current project take first priority, followed by concepts imported from SAS Contextual Analysis, and finally, concepts imported from SAS Enterprise Content Categorization.
- When importing categories from an existing project, the source of the category determines which category name takes precedence when the system resolves a naming conflict. Categories imported from SAS Contextual Analysis take first priority, followed by external categories in the current project, and finally, concepts imported from SAS Enterprise Content Categorization.

### **Creating a New Project**

#### **Using the Create Project Wizard**

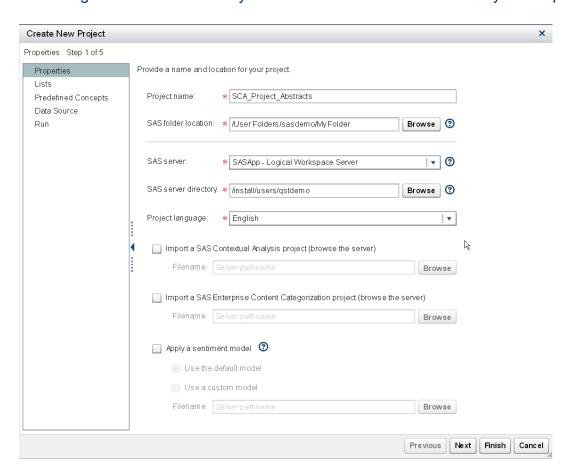
The first time you log on to SAS Contextual Analysis, you must create a project before you can do anything else. To create a new project, click the icon near the upper left corner of the main window. The Create New Project wizard appears, where you can enter all the specifics for your project.

TIP Click the Help icon ② in the Create New Project wizard for information about a specific field or page.

# **Step 1: Identify Project Files, Servers, and Other Properties**

- 1 Enter the project name, and specify where the project folder can be accessed in SAS metadata. If you want the project to be shared by others, select **Browse** and select the **Products** folder, and then select the **SAS Contextual Analysis** folder.
- Indicate the SAS server and SAS server directory. (Be aware that if you want to share the project with other users, those users must have Read and Write access to the SAS server directory.)
- 3 Select a language for the project data. This is the language that will be processed in your data files.
- **4** (Optional) Import a SAS Contextual Analysis project. For more information, see "Importing an Existing SAS Contextual Analysis Project Model" on page 10.
- 5 (Optional) Import a SAS Enterprise Content Categorization project. For more information, see "Importing an Existing SAS Enterprise Content Categorization Project" on page 10.

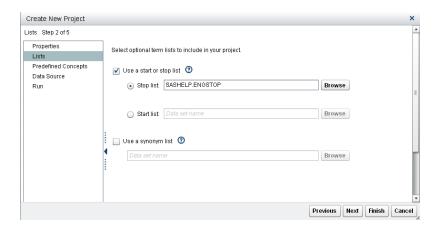
6 (Optional) Apply a sentiment model. For more information about sentiment, see "Using SAS Sentiment Analysis Models in SAS Contextual Analysis" on page 27.



#### **Step 2: Specify Term and Synonym Lists**

- 1 (Optional) Identify a start list or a stop list (but not both) to control which terms to include or exclude during text mining analysis. For more information, see "Start Lists and Stop Lists" on page 34. A default stop list is selected by default.
- 2 (Optional) Specify a synonym list to identify pairs of words that should be treated as single terms for analysis. For more information, see "Terms and Synonyms" on page 33.

#### 14 Chapter 2 / Projects in SAS Contextual Analysis

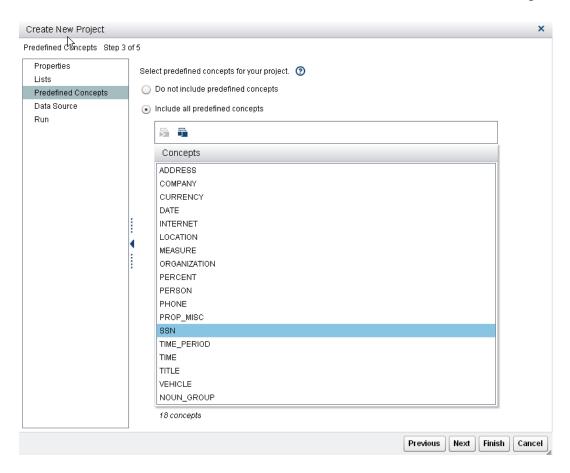


#### **Step 3: Choose Predefined Concepts**

SAS Contextual Analysis provides *predefined concepts*, which are concepts whose rules are already written. Predefined concepts save time by providing you with commonly used concepts and their definitions, such as **COMPANY** or **TITLE**. You can choose to include them or not. If you do not include predefined concepts, they cannot be added later. If you include the predefined concepts, you can disable one or more predefined concepts by selecting one or more predefined concepts and then clicking

 $\blacksquare$  . Disabled concepts are ignored during data processing. You can re-enable any predefined concept by selecting a concept and then clicking  $\blacksquare$  .

For information about predefined concepts, see "Concepts" on page 30.



#### **Step 4: Identify a Data Source**

Here are the options for selecting a data source:

- You can choose a data source now or later. If you choose a data source later, you can still enter more information for your project in the Edit Project wizard.
- You can select analysis variables from within a SAS data set. If you choose this option, you must specify the data set library and name and specify the text variable that you want to analyze. Rather than including the full text of a document in a SAS variable, you can enter a file reference, which identifies the location of a file. Using a file reference is the only way to analyze documents that are longer than 32,767 characters.

**Note:** The referenced file must be in plain text format (TXT).

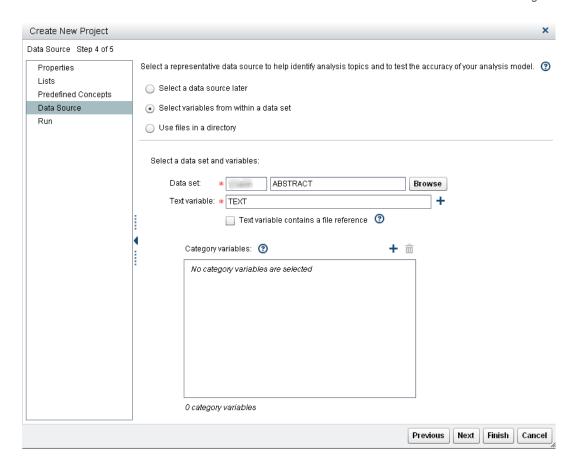
In addition, you can specify one or more category variables to indicate how you want the documents to be grouped. For example, suppose you are analyzing customer comments from hotel stays. The data column Hotels includes names of hotels where customers stayed. If you specify Hotels as a category variable, then category rules are automatically generated. Subcategory rules are also generated for each hotel that appears frequently in the data.

**Note:** SAS Contextual Analysis reserves variables names that begin with an underscore (\_\_). Therefore, if you select a data set that includes variable names that begin with an underscore (\_\_), you could encounter an error. If an error occurs, rename the variables in the data set and try again.

You can specify a document collection that is stored in text-based file formats such as MS Office, OpenDocument (OpenOffice), PDF, XML, or HTML. The files must be located in a folder. You can define categories later.

**Note:** Files that have no content are not imported—they are ignored.

In the following example, the text variable *TEXT* is selected from the data set *ABSTRACT*. There is no category variable.



#### **Step 5: Run the Project**

You can choose to run the entire project now or later. Select **None** on the Run page to run the project later. See the Help for more information about when to run the project.

If you choose to run the entire project, the following events occur for data sources that are provided:

- Parsing takes place.
- Topics are generated.
- Rules are generated and run for any category variables that you specified.
- Sentiment is applied (if you specified a sentiment model).

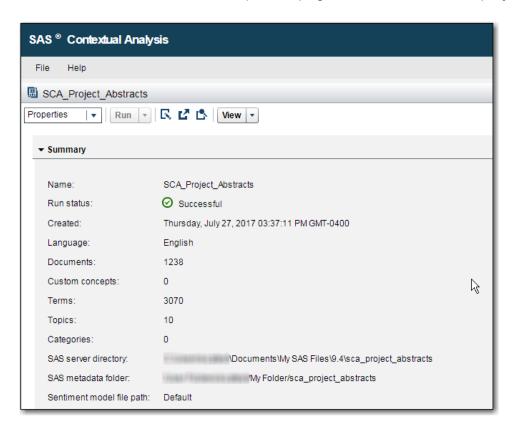
Concepts are applied (if you included predefined concepts in the project).

If you import a SAS Enterprise Content Categorization project and then run the project, the following events occur:

- Imported concepts are compiled and applied to the data source that is provided.
- Imported categories are compiled and applied to the data source provided.

**Note:** Subcategories are imported only if their parent categories were successfully imported.

The **Run status** field on the Properties page indicates whether a project is running.



**Note:** You can also check the **Progress** panel, which is located in the lower left corner of the main window. Click the word **Progress** to see which projects are running and which projects have finished running.



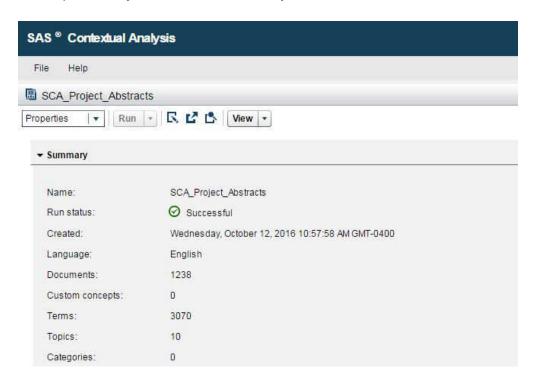
The **Run** menu enables you to run analysis tasks individually or run only the tasks that are out of date (and their dependent tasks, if any).



### **Using the Properties Page**

#### **Checking Project Status**

The Properties page indicates whether the project ran successfully and provides basic information about the data that were analyzed.



The resulting status of each analysis task that was run is displayed in the following fields in the **Status** section:



#### Task

The name of the analysis task

#### Task Up-to-Date

Indicates whether information in the task has changed since the last time the task was run. If no information has changed, the value is Yes and no further action is required. If information has changed in the task since the last run, then the value is No and the task should be rerun.

#### Last Run Date and Last Run Time

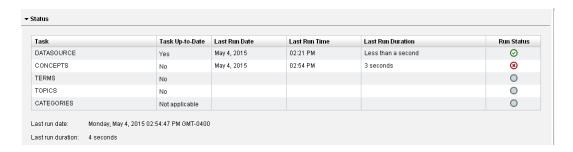
The last date and time the task was run

#### **Last Run Duration**

The duration of the task's last run

#### **Last Run Status**

Indicates whether the task's run was successful  $\bigcirc$ , the run failed  $\bigcirc$ , the task did not run  $\bigcirc$ , or the task was unable to run (because of missing or insufficient information) or warnings occurred  $\triangle$ . Click **View messages** on the toolbar for specific information. A status of **Not run**  $\bigcirc$  is displayed if the failure of a task prevents dependent tasks from running. In the following example, the failure of the CONCEPTS task prevented the TERMS and TOPICS tasks from running. You must correct the error and rerun the project until it runs successfully.



Click on the toolbar to see messages about the status of each task. Here is an example of the Messages window for a task that has an error. The **Message Type** column indicates the corresponding analysis task for each message.



The data source information is displayed at the bottom of the Properties page.

▼ Data Source		
Library:	scasio	
Data set:	ABSTRACT	
Column:	TEXT	

#### **Editing Project Information**

Click to edit basic information for your project (such as project name). The Edit Project wizard appears. Items that you cannot edit appear in gray.

Note: You must rerun the project to see the effects of your changes.

#### **Viewing and Downloading Code**

You can view and download SAS score code that is created. Score code enables you to apply the text analytic models in your project (concepts, categories, and sentiment) to other data.

On the Properties page, click View. Select Concept code, Sentiment code, or

Categories code. Click to copy the code for use in other programs.

TIP There are comments embedded in the code that give details about the generated code. It is recommended that you read the comments about the code and the default settings that are included.

When working with score code, please note the following:

When a category is created from either a topic promotion or a category variable, it contains a list of automatically generated subcategories whose rules, when processed together, define the category. You can remove these automatically generated subcategories from the output (and thereby see the results for the parent categories only) by setting the following flag to ▼:

%let drop auto generated

- If any predefined concept is enabled or disabled in SAS Contextual Analysis after the code is generated, you must regenerate the score code before copying and pasting it into another program (or downloading the code).
- If you have run concepts while predefined concepts were enabled, the generated scored data displays matches that do not include a full path.

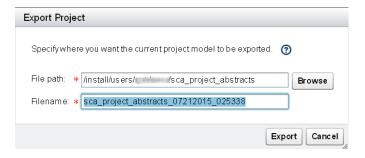
If you want to download the code to a file, click 🔼 and follow the prompts to specify a location for the file.

#### **Exporting a Project Model**

You can export a SAS Contextual Analysis project model so that you can reuse its rules in a new project. When you export a project, only the category and concept rules are exported. Any predefined concepts that you have selected are preserved. The other project components such as topics, terms, data, project settings, and so on, are not exported. The file that you export is in JSON (JavaScript Object Notation) format.

To export a project:

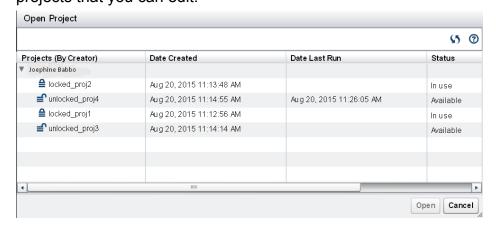
- Open the project that you want to export. Click not not to olbar.
- In the **File path** field, select a server location for the file that you are exporting. Note that the project name appears in the file path in all lowercase letters, and underscores ( ) might be inserted.
- In the **Filename** field, review the automatically generated name for the exported project and edit the name if desired. The generated filename consists of the project name plus the current date (in MMDDYY format) and time (in MMHHSS format). respectively.



#### **Sharing Projects**

Project to be shared by multiple users must be saved in a specific project folder named **Products/SAS Contextual Analysis**. When creating a project, save it in the appropriate folder by selecting **Browse** and select the **Products** folder, and then select the **SAS Contextual Analysis** folder.

You can open and edit a project that has been created by another user. From the main window, click . The Open Project window displays all of the shared projects and their owners. Projects that are in use are locked and cannot be opened. The locked icon appears next to projects that cannot be opened. The unlocked icon appears next to projects that you can edit.



**Note:** To keep your shared projects unlocked for other users, close open projects that are not in use and sign out of your SAS Contextual Analysis session (rather than closing your web browser).

After you open a project, you can edit the project information such as synonym lists. For information about editing project information, see "Editing Project Information" on page 22. For information about editing concept or category rules, terms, or other analysis tasks, see Chapter 3, "Performing the Analysis Tasks," on page 29.

#### **Scoring an External Data Set**

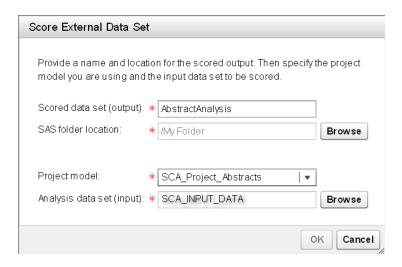
You can use the model that you built in your SAS Contextual Analysis project to score an external data set. When you score an external data set, the category model is applied to the external data set (called the target data set). The categorization information for the document collection is then output into a scored data set.

Note: The data set must be stored in a file outside a folder. If your project uses a folder as a data source, you cannot score data sets within the same folder.

To score an external data set:

- Select File ->Score External Data Set from the application's main menu.
- In the **Scored data set (output)** field, enter the name for the scored data set that is to be generated.
- 3 Enter the SAS folder location for saving the data.
- In the Project Model field, select the name of the project that contains the analysis model that you are using.
- In the **Analysis data set (input)** field, provide the data set to be scored. The analysis data set must have the same text variable as the selected project model's data source.

Note: To be eligible for scoring, a project must have compiled a category binary file. A category binary file is generated when you run a project that contains categories.



After the scoring begins, the project's run status changes to **Running**. When the scoring is complete, the scored data set is placed in the library folder where the project that you used as your project model is stored.

#### **About Sentiment Analysis**

#### **Introduction to Document Scoring**

Sentiment analysis is the process of identifying the author's tone or attitude (positive, negative, or neutral) expressed in a document. SAS Contextual Analysis uses a set of proprietary rules that identify and analyze terms, phrases, and character strings that imply sentiment. A sentiment score is then assigned, based on that analysis. Using these rules, the software is able to provide repeatable, high quality results.

The assignment of sentiment to a document is based on the attitude that is associated with the document as a whole. For example, the following document would have a positive sentiment: Had an awesome time yesterday. Glad I brought my tent from Store XYZ.

Because documents can be associated with multiple words or terms that imply sentiment, SAS Contextual Analysis uses a scoring system to assign a final sentiment score. The following list provides basic information about how sentiment scoring works. (The information has been simplified to illustrate key concepts.)

- Each positive term or phrase is worth a single (positive) point.
- Each negative term or phrase is worth a negative point. If there are more positive terms or phrases than negative, the final sentiment score is positive.
- If there are more negative terms or phrases, the final sentiment score is negative.
- If there are an equal number of positive and negative terms or phrases, the sentiment score is neutral.

# **Using SAS Sentiment Analysis Models in SAS Contextual Analysis**

Rules that are generated using SAS Sentiment Analysis are stored in a .sam binary file. When you create a project in SAS Contextual Analysis, you can use a .sam binary file that you have created to your specifications, or you can use the default file that is available for your project's language.

Note: Not all languages have default sentiment models available for use.

For more information about the sentiment analysis and scoring, see SAS Sentiment Analysis 12.2: User's Guide.

# Performing the Analysis Tasks

Overview of the Analysis Tasks	<b>30</b>
Introduction	30
Concepts	30
Terms and Synonyms	33
Start Lists and Stop Lists	34
Topics	35
Categories	35
Using the Analysis Task Pages	36
Concepts Page	36
Terms Page	41
Topics Page	44
Categories Page	49
Writing Concept Rules: Basic LITI Syntax	56
Introduction to Concept Rules	56
Concepts versus Facts	57
Which Rule Type Should I Use?	58
Using Punctuation	62
Adding Rule Modifiers	63
Using Boolean Operators for Extracting Concept	
Rules and Facts	65
Using the Coreference Operator	70
Using the Export Feature	71
Using Part-of-Speech and Other Tags	72

Using Regular Expressions (Regex)	75
Using Morphological Expansion Symbols	79
Adding Comments	80
Concept Rule Types: Examples	80
Writing Category Rules: Boolean Rules	83
Introduction to Category Rules	83
Boolean and Proximity Operators for Category Rules	85
Using Symbols in Boolean Rules	90
Using _tmac for Referencing Categories	92

# **Overview of the Analysis Tasks**

#### Introduction

When you run a project, the following analysis tasks are performed (if data are present):

- concept extraction
- term identification (including synonyms)
- topic discovery
- category analysis

The following sections describe each task.

# Concepts

A *concept* is a property such as a book title, last name, city, gender, and so on. Concepts are useful for analyzing information in context. You can write rules for recognizing concepts that are important to you, thereby creating custom concepts. For example, you can specify that the concept *kitchen* is identified when the terms *refrigerator*, *sink*, and *countertop* are encountered in text.

SAS Contextual Analysis provides *predefined concepts*, which are concepts whose rules are already written. Predefined concepts save time by providing you with

commonly used concepts and their definitions, such as COMPANY or TITLE. You cannot rename predefined concepts, nor can you view or edit their base definitions. You can provide additional rules in the Edit tab for processing.

**Note:** If an imported concept has the same name as a predefined concept, the imported concept's rules are added to the predefined concept's rules.

For custom concepts, you can prioritize which matches are returned when overlapping matches occur (for example, a concept node that matches New York and another concept node that matches New York City). You do this by setting a priority value. When setting priority values, it is helpful to know the preset values of predefined concepts so that you can set a custom concept's priority at a higher value. For more information about setting priorities, see "Concepts Page" on page 36

Table 3.1 on page 31 shows a list of the predefined concepts for English that are included with SAS Contextual Analysis, along with their priority values. See Appendix 2, "Predefined Concept Priorities (for Languages Other Than English)," on page 137 for a complete list of predefined concepts and their priority values for supported languages other than English.

**Note:** Some languages use a subset of the predefined concepts listed here.

You can disable or enable any of the concepts during project creation (or in the Concepts task window).

Table 2.1	Dradafinad	Concento and	Driorition	for English
rabie 3. i	Predelilled	Concepts and	Priorities	ioi Erigiisti

Predefined Concept	Description	Priority Value
ADDRESS	Postal address or number and street name	20
COMPANY	Company name	25
	Note: SAS Contextual Analysis uses a fixed dictionary of company and organization names in order to identify this concept. This concept is frequently associated with a parent. For example, if IBM appears in the text, it is returned with the predefined parent International Business Machines. Typically, the longest and most precise version of a name is used as the parent form.	

CURRENCY	Currency or currency expression. Examples: \$300, 300 million dollars	18
DATE	Date expression including date, day, month, or year	18
INTERNET	Digital location, including URL, path, filename, or email address	18
LOCATION	City, country, state, geographical place or region, or political place or region	30*
MEASURE	Measurement or measurement expression. Examples: 500kg, 2300 sq f	20
NOUN_GROUP	Multiple words that act as a single term (for example, "for sale" or "clinical trial")	17
ORGANIZATION	Government, legal, or service agency. (See the note associated with COMPANY.)	25
PERCENT	Percentage or percentage expression. Examples: 97%, 12 percentage points	18
PERSON		20
	Examples: 97%, 12 percentage points	-
PERSON	Examples: 97%, 12 percentage points  Person's name	20
PERSON PHONE	Examples: 97%, 12 percentage points  Person's name  Phone number  Proper noun with an ambiguous (miscellaneous) classification that could be in another category, such as product, book,	20
PERSON PHONE PROP_MISC	Examples: 97%, 12 percentage points  Person's name  Phone number  Proper noun with an ambiguous (miscellaneous) classification that could be in another category, such as product, book, person, and so on. Example: Fargo	20 18 5
PERSON PHONE PROP_MISC	Examples: 97%, 12 percentage points  Person's name  Phone number  Proper noun with an ambiguous (miscellaneous) classification that could be in another category, such as product, book, person, and so on. Example: Fargo  Social Security number  Time or time expression. Example: 0800,	20 18 5
PERSON PHONE PROP_MISC  SSN TIME	Examples: 97%, 12 percentage points  Person's name  Phone number  Proper noun with an ambiguous (miscellaneous) classification that could be in another category, such as product, book, person, and so on. Example: Fargo  Social Security number  Time or time expression. Example: 0800, 1:15, 6pm	20 18 5 18

VEHICLE Motor vehicle including color, year, make, and 2 model	20
--	----

Highest priority value for English.

A *custom concept* is a concept whose rules you must write.

**Note:** If you have imported a SAS Enterprise Content Categorization project, the concepts that were created using LITI rules appear in your project as custom concepts. You can edit them further by using the rules editor.

For more information about writing concept rules, see "Writing Concept Rules: Basic LITI Syntax" on page 56. For information about writing Boolean rules, see "Writing Category Rules: Boolean Rules" on page 83.

# **Terms and Synonyms**

A term is defined as a label for a group of strings or patterns that represent a single concept (an idea) as defined by underlying rules or algorithms. In SAS Contextual Analysis, a term is the basic building block for topics, term maps, and category rules. Each term has an associated role that either is blank or identifies the term's part of speech. A term reflects one or more surface forms. A surface form is a variant of a term that is located in a matched subset of text. Surface forms can include inflected forms. synonyms, misspellings, and other ways of referring to a term. SAS Contextual Analysis can identify and classify misspellings of terms based on similarity and frequency. Because misspellings actually refer to another term, they are treated as synonyms during analysis.

A synonym list is a SAS data set that identifies pairs of words that should be treated as single terms for the purposes of analysis. Synonyms are applied at the parent level. You can specify a synonym list in the Create New Project wizard and in the Edit Project wizard. Synonym lists are stored in data sets and have a required format. You must include the following variables:

- TERM, which contains a term to treat as a synonym of the PARENT.
- PARENT, which contains the representative term to which the TERM should be assigned.

You can also include the following variables:

- TERMROLE, which enables you to specify that the synonym is assigned only when the TERM occurs in the role specified in this variable. A *term role* is a function performed by a term in a particular context; term roles include part-of-speech roles, entity roles, and user-defined roles.
- PARENTROLE, which enables you to specify the role of the PARENT.

TIP If you want to use any of the terms in your synonym list to extract concepts, you must create custom concepts for the PARENTROLE that match the PARENT terms (or verify that the concepts exist). After the concepts are established, rerun the terms. For example, suppose the parent term Luke Skywalker specifies the parent role JEDI\_MASTER. You must create a custom concept called JEDI\_MASTER that includes a rule that matches Luke Skywalker and then rerun the terms.

For more information about roles, see the section "Term Roles and Attributes" in SAS Text Miner: Reference Help.

**Note:** If a synonym list includes multiple entries that assign the same terms to different parents, then the parsing results reflect only the first entry.

### **Start Lists and Stop Lists**

You use start lists and stop lists to control which terms are or are not used in a text mining analysis. A *start list* is a data set that contains a list of terms to include in the parsing results. If you use a start list, then only terms that are included in that list appear in parsing results. A *stop list* is a data set that contains a list of terms to exclude from the parsing results. You can use stop lists to exclude terms that contain little information or that are extraneous to your text mining tasks. A default stop list is provided for English (Sashelp.EngStop).

Start lists and stop lists have the same required format. You must include the variable TERM, which contains the terms to include (start) or exclude (stop). You can also include the variable ROLE, which contains an associated role. If you specify a ROLE

variable, then terms are kept (for a start list) or dropped (for a stop list) only if their role is the one that is specified in the ROLE variable.

# **Topics**

Topics are derived from natural groupings of important terms that occur in your documents. In SAS Contextual Analysis, topics are automatically generated and assigned to documents. A single document can contain more than one topic.

The Topics page displays all the topics that SAS Contextual Analysis identified. The default name of a topic is the top five terms that appear frequently in the topic. These terms are sorted in descending order based on their weight.

# **Categories**

A *category* identifies a group of documents that share a common characteristic.

For example, you could use categories to identify the following:

- areas of complaints for hotel stays
- themes in abstracts of published articles
- recurring problems in a warranty call center

You create categories by promoting a topic to a category, specifying a category variable in the New Project wizard, or creating a new category. You can also import categories from SAS Enterprise Content Categorization. You can edit the rules that are automatically generated for category variables and for topics that are promoted to categories.

**Note:** The category rules are in the format that SAS Enterprise Content Categorization uses (MCAT), rather than in LITI format. You can refer to LITI concepts from within categories.

For more information about writing concept rules, see "Writing Concept Rules: Basic LITI Syntax" on page 56. For information about writing Boolean rules, see "Writing Category Rules: Boolean Rules" on page 83.

# **Using the Analysis Task Pages**

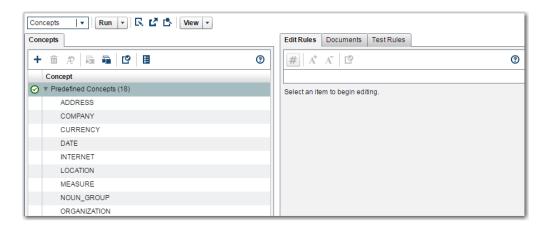
# **Concepts Page**

The Concepts page enables you to view predefined and imported concepts, add custom concepts, test concept rules, edit concept properties, and view the documents that contain matches.

TIP Customize your view of the items that are associated with a concept node by dragging the **Edit Rules**, **Documents**, and **Test Rules** tabs from one pane to another.

Expand the list of predefined and custom concept nodes to see what is included in your analysis.

Note: If you chose to exclude predefined concepts during project creation, you cannot access predefined concepts on the Concepts page.



Click the toolbar buttons to disable • or enable • a concept node.

Note: Any terms that are associated with a disabled concept are removed from the terms list and disregarded during parsing.

Here are other important actions that you can take on the Concepts page:

#### Add a custom concept

Click + to add a custom concept node for which you create your own rules.

TIP When you create a custom concept node, follow these naming guidelines:

- Use valid characters numbers, letters, and underscores ( ). (See the Note below regarding the use of underscores).
- Concept names are case-sensitive.
- Create names that are not regular words; using mixed case is recommended to help with readability. For example, MyConcept or myConcept are good names. Do not use names for custom concepts that are also words (for example, Problem or Mechanics ) that could be matched in your text. Instead, use names that cannot be interpreted as words, such as MyNewConcept.

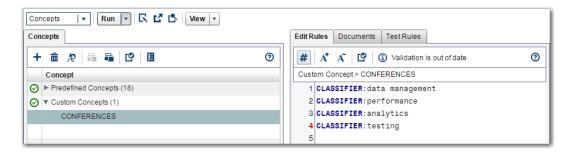
#### Note:

You must follow these guidelines when using underscores ( ) in concept names, or your concept rules will not work as expected.

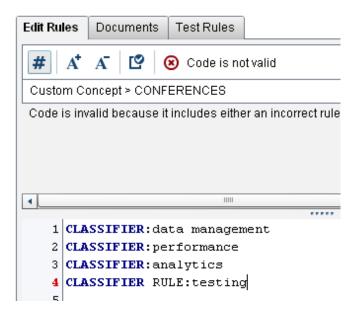
- If you use underscores at either end of the concept name, be sure there is a matched pair at both ends. For example, Domestic is permitted, but Domestic is not permitted.
- Do not include Q, a character combination reserved by the application, anywhere in a concept name.
- If a concept name begins with an underscore, the next character must be a letter. For example, the concept name 25anniv is not permitted.

TIP Use mixed case to enhance the readability of concept names. For example, truckMechanicalIssues is easier to read than truckmechanicalissues.

On the Edit Rules tab, enter the LITI rules for the concept node. You must validate the rules before the concept node can be used in the analysis.

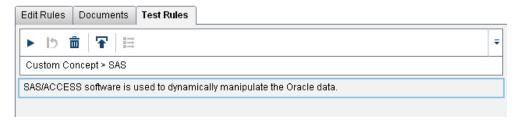


Click on the **Edit Rules** tab to validate each rule individually, or click in the **Concepts** to validate all the rules. Errors and other messages are displayed on the **Edit Rules** tab.



#### Test concept node rules

Select a concept node in the **Concepts** tab and then click the **Test Rules** tab. Upload a file to test by clicking ¬, or simply type test text for the rule that you have selected. Click > to test the rule.



In the following sample screen, the matched strings are highlighted for the concept node **SAS**.

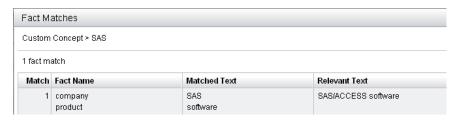


**Note:** The matched strings (highlighted text) can contain multiple matches in concept rule testing.

Note: Concept nodes that are named in categories might return more matches than concepts that are run outside of categories. In categories, matches on concepts are based on an "all matches" method, which returns all matches found in the text. By contrast, in concepts, matches are based on a "best match" method. The best match method detects when text that matches one concept overlaps text that matches another concept (for example, a concept that matches New York and another concept that matches New York City). When concept matches overlap and the best match method is used, only the concept that is assigned the highest number for the priority is returned (1 is the lowest). When two or more concepts have the same priority assigned, SAS Contextual Analysis selects a match by using the "longest match" method.

You can view matches for facts (related pieces of information in text that are located and matched together) in a separate window so that you can examine the matched strings in your test text. Click to open the window and view the fact matches. The following sample screen shows the matches for testing the text SAS/ACCESS software is used to dynamically manipulate the Oracle data. With the rule

```
\label{eq:predicate_rule:(company, product):(AND, "\_company{SAS}", "\_product{software}")} \\
```



For more information about facts, see "Concepts versus Facts" on page 57.

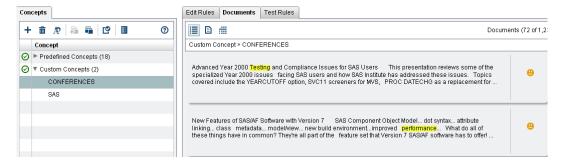
View and explore matching documents

To view the training documents that contain matches, click the **Documents** tab. Click one of the icons to switch between document views. Suppose

you created a concept node conferences, which contains the rules

CLASSIFIER:testing CLASSIFIER:performance

Matches within the documents are highlighted, as shown in the following sample screen:



**Note:** Sentiment values are displayed only if you applied a sentiment model when you created the project.

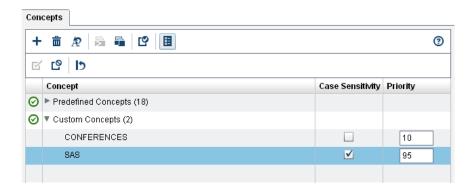
Edit custom concept node properties to refine concept matches You can edit certain properties to help refine the matches from your custom concept rules. Click to view the properties.

Select the **Case Sensitivity** check box to ensure that matches occur for the cases that are specified in the rule.

You can prioritize which matches are returned when overlapping matches occur (for example, a concept node that matches New York and another concept node that

matches New York City). In the case of overlapping matches, the concept node with the highest number entered in the **Priority** column is returned. The value must be a positive number (1 is the lowest priority). There is no limit to the highest priority value. The default value is 10.

Note: If you want to make sure that your concept matches do not overlap with predefined concept values, use a priority value that is higher than that of the predefined concepts that you have included. For more information, see "Concepts" on page 30 or Appendix 2, "Predefined Concept Priorities (for Languages Other Than English)," on page 137.



### **Terms Page**

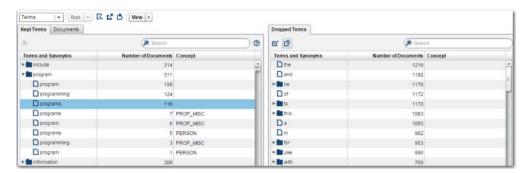
After a project is successfully run, open the **Terms** page to view the terms that were discovered in your document collection. The default view shows the Kept Terms on the left and the **Dropped Terms** on the right. Use the icons and to switch views within a tab.

TIP To customize your view, drag the **Document**, **Kept Terms**, or **Dropped Terms** tabs from one pane to another.

Here are other important tasks that you can complete in the Terms page:

#### View Terms

The **Kept Terms** displays all the terms in the document collection that were kept. The **Number of Documents** column displays the number of training documents that contain the selected term. The **Concept** column displays each term's role, if one can be determined. To view the surface forms that were assigned to a term, click the triangle that appears next to that term.



#### Drag terms from one tab to another

By default, the lists of terms are sorted in descending order of the number of documents in which each term appears. You can drag parent terms from the **Kept** tab to the **Dropped** tab, and back again.

**Note:** If you make changes to the terms and you want to see the effects of your changes, you must click **Run** to rerun the project.

**CAUTION!** If concept rules are out-of-date when you rerun any tasks (all out-of-date tasks or topics only), any changes that you made to terms are overwritten with the original terms list.

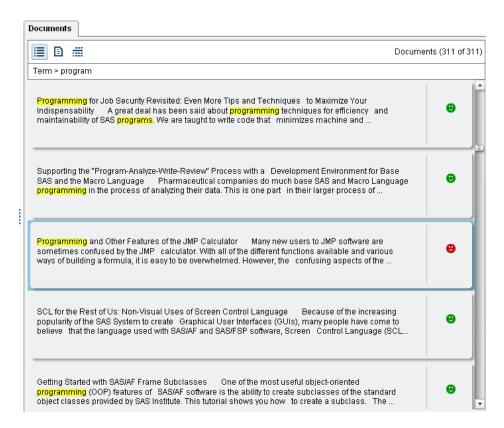
### View and explore matching documents

To view the training documents that contain matches, click the **Documents** tab.

Click one of the icons 

to switch between document views. Matches are

highlighted, as shown in the following sample screen:



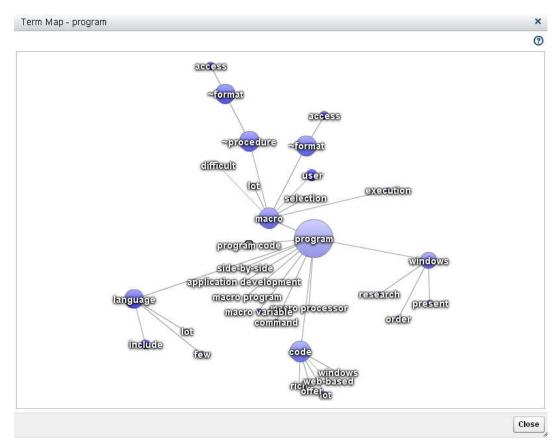
**Note:** Sentiment values are displayed only if you applied a sentiment model when you created the project.

#### View a term map

To view a **Term Map** for a term, select that term in the **Kept Terms** and click **...**.



TIP Sometimes a term map is not generated because there are too few documents existing in the corpus to find significant relationships with the term. You can choose another term or, if necessary, increase the size of your document collection.



The Term Map window displays a term map for the selected term. In the preceding sample screen, the selected term is *program*, and it is represented by the largest circle in the map. For more information about reading the map, click ② above the term map.

# **Topics Page**

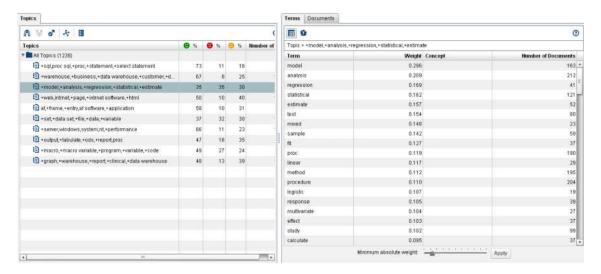
To analyze a topic, select that topic on the **Topics** tab. The selected topic is identified by its five most important terms. Here are the tasks that you can perform on the Topics page:

### View terms that comprise the topic

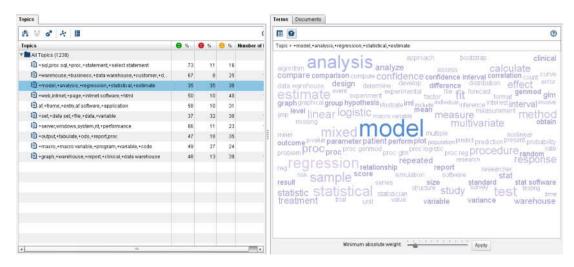
In the following sample screen, the topic is identified by the terms **model**, **analysis**, **regression**, **statistical**, and **estimate**.

Note: The percentage of the documents in the topic that have a sentiment score of positive, negative, and neutral appears with each topic, provided that you included a sentiment model when you created the project.

When you select a topic, the view in the right is updated. Use the icons in the right 's toolbar to switch views. For information about each view, click ? in the right. Click in the right to view the terms list as a table. The terms table displays every term in the topic, its calculated weight, its assigned role (concept), and the number of documents that contain that term.



In the following sample screen, the word-cloud view icon \_ was selected. A slider at the bottom of the word-cloud pane enables you to adjust the minimum absolute weight necessary for a term to be included in this topic. The word cloud is updated as you move the slider to the right. Click **Apply** to finalize the changes that you make.

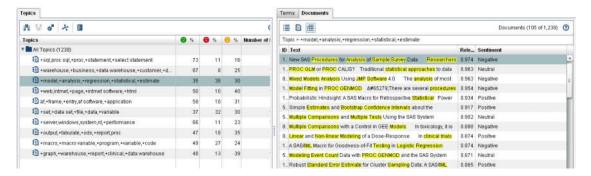


### View documents associated with the topic

To view the training documents that are associated with a topic, click the **Documents** tab.

Select one of the document view icons to see the selected topic's

terms. In the following sample screen, all of the terms that mark the documents as being a part of this topic are highlighted.



#### Split and merge topics

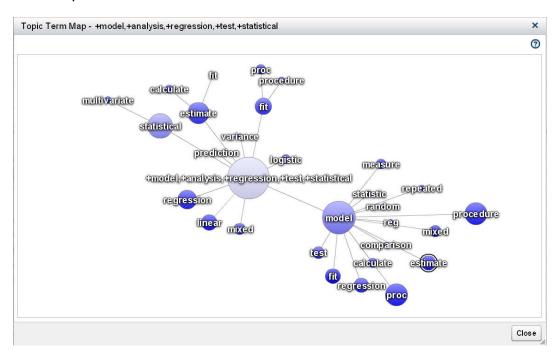
If you consider a topic too broad for your purposes, you can split it by selecting it in the **Topics** and then clicking . This action splits the selected topic into two new topics.

If you see two topics that seem related, you can merge them by selecting them and clicking . This action combines all the selected topics into the same topic.

**Note:** If you want to revert to the unmerged topics after you merge them, you can do so by rerunning topics. Your changes to terms and topics up to that point will be lost.

#### View a topic term map

You can view a topic term map from the **Topics** pane by selecting a topic and clicking . In a topic term map, the tilde at the beginning of a term is treated as a NOT operator. For more information about reading the map, click nabove the topic term map.



### Promote topics to categories

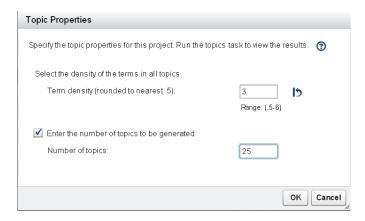
A key step in the analysis is to identify which topics you want to promote to categories. To promote a topic to a category, select that topic in the Topics pane and to the **Categories** page. You can promote multiple topics to categories at one time.



### Edit topic properties

You can edit the properties that affect all topics. Term density refers to how topics are populated with terms; it is defined by a number between 0.5 and 6 (the default value is 2). When term density is closer to 0.5, topics are more densely populated by terms. When term density is closer to 6, topics are less densely populated by terms. This value affects the number of documents that belong to a topic (for example, having fewer terms in a topic captures fewer documents). Values that you enter are rounded to the nearest integer or half-integer.

You can also designate a maximum number of topics that you want generated for the project. If you leave the setting blank, the software uses default methods to generate topics from important terms.



**Note:** You must run the topics to see the results of your changes.

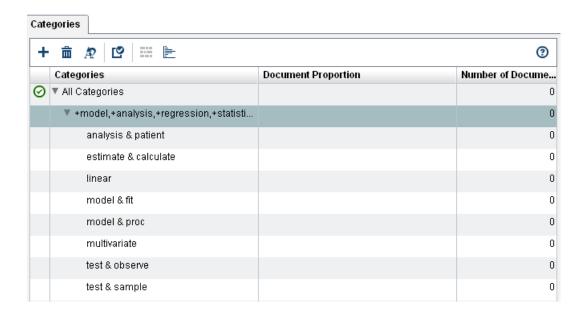
# **Categories Page**

After you create a category from a topic in the **Topics** page, the category appears on the **Categories** page. In the **Edit Rules** tab, you see the rules that were generated for that category.

**Note:** Rules are generated automatically for a maximum of 25 categories. When less than eight categories exist, rules are grouped into a category marked "OTHER."

Note: Subcategories are created when three or more values exist for an external category's variable. For example, a category Flavors that contains variables blueberry, cherry, and peach creates three subcategories that contain the generated rules. Subcategories are not created for external categories where binary variable values exist. For example, a category In stock with values Yes and No does not create subcategories; rather, the generated rule is listed within the category In stock.

The **Documents** tab is not populated until you run the category.



TIP Customize your view of the items associated with a category by dragging the Edit Rules, Documents, and Test Rules tabs from one to another.

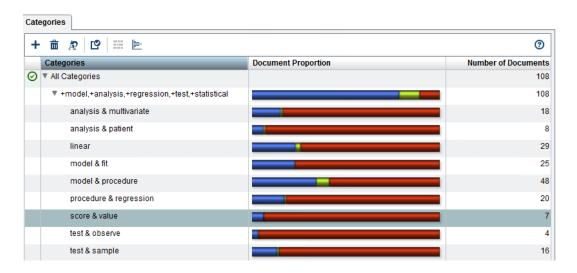
Here are the important tasks that you can perform on the **Categories** page:

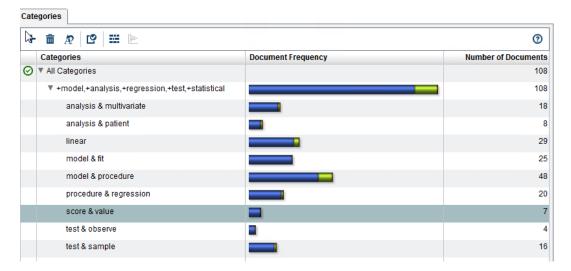
#### Run the categories

Use the **Run** menu to compile your topics into categories. The rules for each category are executed against the input data set.

Note: The least-occurring value is used as the target for analysis in categories where binary variable values exist. For example, suppose the category In stock has a value Yes that occurs 1200 times and a value No that occurs 940 times. In this case, No is used as the target.

This action updates the **Document Proportion** and **Document Frequency** columns. Use  $\models$  and  $\rightleftharpoons$  to switch between the column views.





The following colors are used in the **Document Proportion** and **Document** Frequency columns:

#### **Blue - True Positive**

The category rules captured documents that you intended to capture. You want to maximize this number.

#### **Green - False Positive**

The category rules captured documents that you did not intend to capture because the documents do not fit the category.

### **Red - False Negative**

The category rules missed documents that were found in the topic.

### **Gray - Matches**

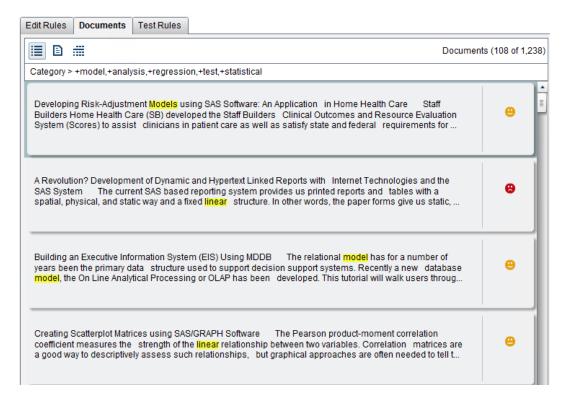
Gray bars show the statistics for matches on categories that are custom or imported.

**Note:** If you run categories that contain no rules or subcategory rules, the results show as false negatives (for categories generated from promoted topics or external category variables only).

To view only the documents that are true positives, click the blue portion of the bar and click the **Documents** tab. To view only the false positives, click the green portion of the bar and click the **Documents** tab, and so on.

### View document matches for categories

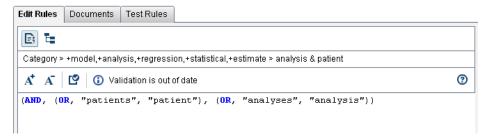
The **Documents** tab is updated to display only the documents that meet your selection. Use the icons to switch between views. Highlighted terms were used to determine the document's membership in the category.



**Note:** The sentiment score for each document is displayed only if you specified a sentiment model when you created the project.

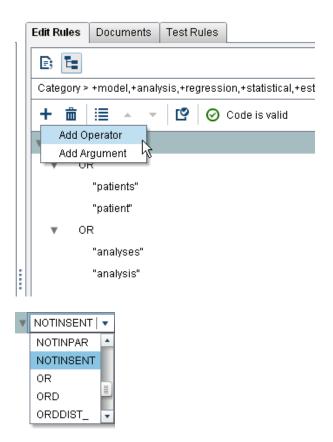
### Edit category rules

To begin editing, select a rule and then click the **Edit Rules** tab. Use the rules code icon and rules tree icon to switch between editing modes. The following sample screen shows the rule code view.

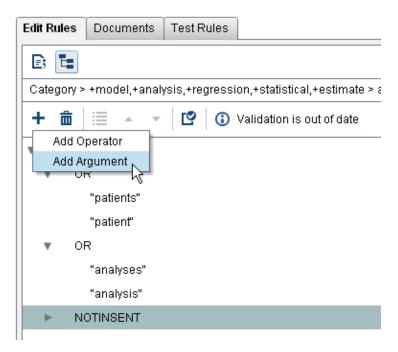


Edit rule code by entering a rule and then validating the code.

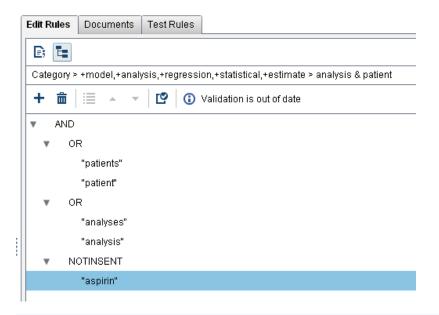
The rule tree enables you to build rules by choosing operators and adding arguments in a visual display. Click + to add an operator or an argument to the rule. The following sample screen show adding an operator to a rule in the rule tree view.



To add arguments, click + and select **Add Argument**.



Enter the argument in the space provided.

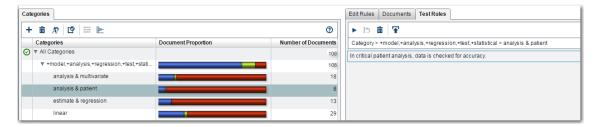


Click on the Edit Rules tab to validate each rule individually, or click on the Categories tab to validate all the rules. Errors and other messages are displayed on the Edit Rules tab.

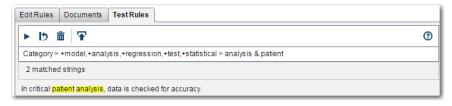
For information about writing category rules, see "Writing Category Rules: Boolean Rules" on page 83.

#### Test category rules

To test category rules, select a rule and then click the **Test Rules** tab. Upload a file to test by clicking  $\mathbf{r}$ , or simply type (or copy and paste) test text for the rule that you have selected. Click > to test the rule.



In the following sample screen, the matched strings are highlighted for the rule Analysis & patient.



Clear the highlighting by clicking 15.

# Writing Concept Rules: Basic LITI **Syntax**

### **Introduction to Concept Rules**

Concept rules are written using LITI (language interpretation and text interpretation) syntax. Concept rules recognize items in context so that you can extract only the pieces of the document that match the rule. For example, you can create a custom concept node named LaGuardia Airport Comments, and then write a rule that extracts all documents in your document set that contain the word LGA. In other words, all of the documents displayed for the concept node LaGuardia Airport Comments would contain LGA.

Each document is evaluated separately for matches; matches do not span documents.

Note: In concepts, matches are based on a "best match" method. The best match method detects when text that matches one concept overlaps text that matches another concept (for example, a concept that matches New York and another concept that matches New York City). When concept matches overlap and the best match method is used, only the concept that is assigned the highest number for the priority is returned (1 is the lowest). When two or more concepts have the same priority assigned, SAS Contextual Analysis selects a match by using the "longest match" method.

For information about editing rules by using the interface and by using properties settings, see "Concepts Page" on page 36. For a list of rule types, see "Which Rule Type Should I Use?" on page 58.

The following list provides basic guidelines for using LITI syntax to write concept rules. The syntax is flexible, and therefore the syntax elements can be combined in numerous ways.

A rule consists of a rule type (which is written in uppercase letters), followed by a colon, then by arguments. For example, in the rule CLASSIFIER: LGA, CLASSIFIER is the rule type, LGA is the argument, and they are separated by a colon. Rule modifiers can be used to further refine the set of matches. The rule syntax varies greatly

depending on the rule type; the basic syntax is included in the description of each rule in Table 3.2 on page 59 and Table 3.3 on page 61. For a list of rule modifiers, see "Adding Rule Modifiers" on page 63.

- Use descriptive concept rule names that cannot be used as single words (for example, BASEBALLSCORE). You can also include the type of rule as a prefix (for example, CONCEPT BASEBALLSCORE).
- A single concept rule can reference one or more other concepts nodes. You can also write rules that recognize key words or elements within a specific context. For example, you can extract documents that contain the string LGA only if it appears before the word Airport.
- Use part-of-speech tags in rules to identify linguistic structures. For more information, see "Using Part-of-Speech and Other Tags" on page 72.
- Use Boolean and proximity operators to enhance the precision of your rules. For more information, see "Using Boolean Operators for Extracting Concept Rules and Facts" on page 65.
- Use morphological expansion operators to return inflected forms of a word.
- Use coreference operators to resolve pronouns. For example, if the pronoun he were used to refer to Walt Disney, you can write a rule that specifies the canonical form (full form) and returns it in the concept. For more information, see "Using the Coreference Operator" on page 70.

# **Concepts versus Facts**

Facts (also called predicates) are related pieces of information in text that are located and matched together.

Facts can be identified within a custom concept. For example, suppose you want to identify US universities that were named after presidents. You could write a rule that identifies George Washington as a US president (US President Names) and also identifies George Washington University as a university named for him (UNIVERSITY).

So, in the sentence There are countless active student organizations at George Washington University, the String George Washington would match the CONCEPT US\_President\_Names and George Washington University Would match UNIVERSITY.

You can use the following special types of concept rules to locate facts:

■ A predicate rule (PREDICATE\_RULE) uses Boolean and proximity operators to help locate facts. For example, you can use Boolean and proximity operators to specify terms that you want to occur within a certain number of terms of each other. The following rule identifies occurrences of the term America (denoted as country) that occurs within three terms of flag, emblem, or crest:

You can use a sequence rule (SEQUENCE) when the order of the items in the fact is important. A sequence rule can detect a structure so that each term in the fact matches in the order that you specify.

**CAUTION!** Although you can create and test fact rules (SEQUENCE and PREDICATE\_RULE) in SAS Contextual Analysis, they are not applied to the project's documents when the project is run. As a result, you will not see fact matches within document views, topics, and auto-generated rules. To obtain fact rule matches, you can choose one of the following options: (1) Use the project's concept score code feature. For more information, see "Viewing and Downloading Code" on page 22. (2) Deploy the project's LITI binary file (which includes the fact rules) for use with SAS Enterprise Content Categorization Server. For more information, see SAS Notes for SAS Contextual Analysis, available at http://support.sas.com/notes/index.html.

# Which Rule Type Should I Use?

There are several distinct types of rules for extracting concepts and facts. You can specify more than one rule in each custom concept or fact. It is important to understand the rule types so that you can select those that efficiently generate the most matches for your purposes.

**Note:** For the concept rule syntax listed in the following tables, < > denotes an optional syntax element. Items in *italics* denote values that you must supply, such as strings and concept node names.

Table 3.2 lists the types of rules that are used for extracting concepts. Included is a brief description of how each rule type is used, along with basic syntax. For examples of concept rule syntax, see "Concept Rule Types: Examples" on page 80.

 Table 3.2
 Overview of Rules for Extracting Concepts

ntifies single terms or strings that you want matched in context. For mple, in a concept definition, you can create CLASSIFIER rules that tain specific airport codes. The portions of text that contain the airport es are considered matches to the CLASSIFIER rules.
mple, in a concept definition, you can create CLASSIFIER rules that tain specific airport codes. The portions of text that contain the airport
return a canonical (full) form for the matched string, usie the optional ument $<$ , $information>$ . Enter the canonical form in the ument, after the comma (,).
ASSIFIER:string <, information>
ntifies related information by referencing other concepts. For example, apture documents that contain certain US airport names and ations, you can create a CONCEPT rule type in the definition. The NCEPT rule type can reference a CLASSIFIER rule type by its name, reby accessing a list of airport codes.
NCEPT is a rule type. It is not to be confused with a "concept" in the eral sense.
The concept that you are referencing in the rule is also matched a string. For example, in the rule CONCEPT: SCORE, the string SCORE is sched. Therefore, it is recommended that you use concept names that not be used as single words (for example, BASEBALLSCORE).
NCEPT: <priority=<i>n&gt;:<i>argument-1</i>&lt;<i>argument-n</i>&gt;ere <i>argument</i> can be a concept name, rule modifier, or string.</priority=<i>
urns matches that occur in the specified context only. For example, to ract matches that include names of university professors, you could rate a C_CONCEPT rule that identifies matches on a concept eviously defined) that identifies last names only when the matched nes are preceded by the word Professor.
e: This rule type requires the _c modifier.
CONCEPT: <argument> _c{argument}<argument> ere argument can be a concept name, rule modifier, or string.</argument></argument>

#### CONCEPT RULE

Uses Boolean and proximity operators to determine matches. For a list of Boolean operators, see "Boolean and Proximity Operators for Category Rules" on page 85.

**Note:** This rule type requires the c modifier. Quotation marks (") must surround the strings that you want to match. The  $c\{\}$  can surround only one argument (unless it is within an OR operator). The argument is highlighted when matches are returned. The other arguments that appear in quotation marks provide context for the match and must be present for a match to occur.

CONCEPT\_RULE:(<Boolean-rule-1>...<Boolean-rule-n> where *Boolean-rule* can be nested *n* times and is written as: Boolean-operator " c{argument-1}",<"argument-2">...<"argument-n">)

#### NO\_BREAK

Prevents partial matches by ensuring that a match occurs only if the entire string is located. For example, suppose you want to capture text that includes the item National Gallery of Art. You can create a rule that ensures that the entire string National Gallery of Art is matched and not Gallery and Art as separate items.

**Note:** This rule type requires the c modifier.

**Note:** NO BREAK applies across the entire taxonomy regardless of where the rule appears or whether the rule is enabled or disabled.

**Note:** Do not insert NO BREAK rules into the same concept that they reference with the c modifier; this can cause unpredictable match behavior.

**Tip:** When you are writing NO\_BREAK rules, it is helpful to insert them all in one concept. That is, create a concept that contains globally implemented rules only (such as NO BREAK or REMOVE ITEM). Having such rules all in one place aids in troubleshooting the matching behavior across your taxonomy.

**NO BREAK**: c{argument} where argument can be a concept name or string.

#### REGEX

Identifies recurring patterns of textual information that can be expressed in numbers and characters, such as telephone numbers, license plates number and character combinations, or word pairings such as merrygo-round. For example, you could write a REGEX rule that matches the expression 32,768. For more information, see "Using Regular" Expressions (Regex)" on page 75.

**REGEX**:regular-expression

#### REMOVE ITEM

Ensures that a correct match is made when one word is a unique identifier for more than one concept. For example, you can write a rule that distinguishes between the Arizona Cardinals football team and the St. Louis Cardinals baseball team. The context of each match is used to eliminate incorrect matches.

Note: This rule type requires the c modifier and the ALIGNED operator. Quotation marks (") must surround the strings that you want to match.

REMOVE\_ITEM:(ALIGNED, "\_c{concept name}", <"argument"> where argument can be a concept name or string.

Table 3.3 lists the rules used for extracting facts. Included is a brief description of how each rule type is used, along with basic syntax.

 Table 3.3
 Overview of the Rules for Extracting Facts

Rule Type	Description and Basic Syntax
PREDICATE_ RULE	Helps you define facts that you want identified in text. For information about facts, see "Concepts versus Facts" on page 57.
	PREDICATE_RULE:(argument-name-1 <argument-name-n>): (Boolean-rule-1<boolean-rule-n>) where argument-name refers to a name you specify for fact matching, and where Boolean-rule can be nested n times and is written as: (Boolean-operator, "_argument-name {argument}", "&lt;_argument-name&gt;{<argument>}")</argument></boolean-rule-n></argument-name-n>
	The PREDICATE_RULE rule type requires arguments. It is more flexible than the SEQUENCE rule type because it does not always specify order.
SEQUENCE	Identifies facts in documents if the facts appear in the order specified. For information about facts, see "Concepts versus Facts" on page 57.  SEQUENCE:(argument-name-1
	<argument-name-n>):_argument-name-1{argument} &lt;_argument-name-n argument&gt; where argument-name refers to a name you specify for fact matching, and where argument refers to the name of one or more concept nodes.</argument-name-n>
	<b>Note:</b> This syntax is written in its simplest form. Additional modifiers and arguments for concept rule matching can be inserted.
	The SEQUENCE rule type requires arguments. The number of argument-names specified must match the number of _argument-names.

# **Using Punctuation**

Use punctuation to qualify the matches for all rule types except CLASSIFIER and CONCEPT.

#### Colon:

Separates rule types and tags. When to use a colon:

- After a concept rule type (for example, CLASSIFIER:)
- Between the arguments list and the rules list in a SEQUENCE or PREDICATE\_RULE definition.
- Before a part-of-speech tag (for example, :Prep).

#### Comma,

Separates elements (such as arguments) in a concept rule definition. Add a space after the comma and before the next element. Also used to separate logical operators in a PREDICATE\_RULE definition.

#### Single space

Separates strings, concept node names, part-of-speech tags, and rule modifiers in CONCEPT, CONCEPT\_RULE, and C\_CONCEPT rule types.

#### Quotation marks ""

Encloses concept node names and strings in CONCEPT\_RULE, REMOVE\_ITEM, and PREDICATE\_RULE rule types.

#### Parentheses ()

Groups the elements in CONCEPT\_RULE, REMOVE\_ITEM, SEQUENCE, and PREDICATE\_RULE rule types.

### Square braces []

Groups elements in the REGEX rule type.

#### Curly braces { }

Delimits information that is returned as a match. Braces {} can be used in combination with parentheses () in some rule types.

## **Adding Rule Modifiers**

Several types of concept rule modifiers can enhance the matching ability of a rule. Table 3.4 and Table 3.5 list the type of rule modifiers available and denote which rule types they can be used in.

 Table 3.4
 Concept Rule Modifiers and Associated Rule Types

Modifier	CLASSIFIER	CONCEPT	C_CONCEPT	CONCEPT_ RULE
Comments	X	X	X	X
Context (_c)			X (Required)	X (Required)
Word (_w)		X	X	X
Word with initial capital letter (_cap)		X	X	X
Multiple matches symbol (>)			Х	Х
Morphological expansion symbols (@, @A, @N, and @V)		X	X	X
Boolean andproximity operators				X
Part-of-speech tags		X	Х	Х
Export feature	X			
Coreference symbols (_ref, _P, and _F)		X	X	X

## **64** Chapter 3 / Performing the Analysis Tasks

Regular expressions (Regex)			
Predefined concepts	X	Х	Х

 Table 3.5
 Concept Rule Modifiers and Associated Rule Types, Continued

Modifier	REMOVE_ ITEM	NO_BREAK	SEQUENCE	PREDICATE _RULE	REGEX
Comments	X	X	X	X	
Context (_c)	X (Required)	X (Required)			
Word (_w)	X	X	X	X	
Word with initial capital letter (_cap)	Х	X	Х	X	
> symbol					
Morphological expansion symbols (@, @A, @N, and @V)	X	X	X	X	
Boolean and proximity operators				X	
Part-of-speech tags	X	X	X	X	
Export feature					
Coreference symbols (_ref, _P, and _F)					
Regular expressions (Regex)					X (Required)

Predefined concepts	Х	X	Х	Х	
•					

## **Using Boolean Operators for Extracting Concept Rules and Facts**

Table 3.6 lists Boolean operators that you can use when you write concept rules and identify facts.

 Table 3.6
 Boolean Operators for Extracting Concept Rules and Facts

Operator	Description
ALIGNED	Takes two arguments. Returns a match when both arguments are present (aligned) in a document. Used with the REMOVE_ITEM rule type. For example, the following rule specifies that if a match on rules in the LOC concept node also matches rules in the PERSON concept node, then the match on LOC should be removed:  REMOVE_ITEM:ALIGNED, ("_c{LOC}", "PERSON")
AND	Takes one or more arguments. Matches if all arguments occur in the document, in any order. For example, the following rule returns a match on King Louis XIV if it occurs in the document with France:  CONCEPT_RULE: (AND, "_c{King Louis XIV}", "France")
DIST_n	(Distance) Takes a value for $n$ and two or more arguments. Matches if all arguments occur within $n$ (or fewer) tokens of each other, regardless of their order. For example, the following rule returns a match in the phrase the picture with the best lighting:
	CONCEPT_RULE: (DIST_5, "best", "_c{picture}")
	<b>Note:</b> For calculation purposes, the distance between tokens is not inclusive. For example, the distance between best and show in the phrase best in show is two tokens. Words that include hyphens are counted as one token (for example, merry-go-round is one token).

#### NOT

Takes one argument. Matches if the argument does not occur in the document. Must be used with the AND operator. For example, the following rule returns a match if cinema, theater, or theatre occur in the document, but Broadway does not:

```
CONCEPT RULE: (AND, (OR, "c{cinema}", "c{theater}", "c{theater}"), (NOT, "Broadway"))
```

**Note:** The NOT operator applies across the entire document. If you specify the OR operator in addition to the AND operator, you must enclose the OR arguments in parentheses.

#### OR

Takes one or more arguments. Matches if at least one argument occurs in the document. For example, the following rule returns a match if one or more of the tokens U.S., US, or United States appear in the document:

```
CONCEPT_RULE: (OR, "_c{U.S.}", "_c{US} ", "_c{United States}")
```

Note: Rules that are generated by SAS Contextual Analysis nest the OR operator within the AND operator. However, the OR operator can stand alone.

#### ORD

(Order) Takes one or more arguments. Matches if all of the arguments occur in the order specified in the rule. For example, the following rule returns a match in the sentence The warranty claim for the washing machine was denied.:

```
(ORD, "warranty", "claim", "denied")
```

## ORDDIST n

(Order and distance) Takes a value for *n* and two or more arguments. Matches if all arguments occur in the same order that is specified in the rule and if all arguments are within *n* tokens of each other. For example, the following rule returns a match in the phrase the teacher introduced elementary statistics because the arguments appear in the correct order and within five tokens of each other:

```
CONCEPT RULE: (ORDDIST 5, "elementary", " c{statistics}")
```

**Note:** For calculation purposes, the distance between tokens is not inclusive. For example, the distance between the best and show in the phrase best in show is two tokens. Words that include hyphens are counted as one token (for example, merry-go-round is one token).

#### PARA

(Paragraph) Matches if all the arguments occur in a single paragraph, in any order. For example, the following rule returns a match if the paragraph contains the term Manhattan and also includes the token apartment. (Only Manhattan is highlighted.)

```
CONCEPT_RULE: (PARA, "_c{Manhattan}", "apartment")
```

**Note:** PARA rules work properly only when they are applied to data sets that contain paragraph delimiters \n\n (newline), \t\t (tab), or <P> (paragraph). PARA cannot be applied on the **Test Rules** tab. PARA also cannot be applied to data that is contained in folders.

#### SENT

(Sentence) Takes two or more arguments. Matches if all the arguments occur in the same sentence, in any order. For example, the following rule returns a match when Amazon and river occur within the same sentence:

```
CONCEPT RULE:(SENT, " c{Amazon}", "river")
```

Delimiters are used for sentence tokenization, which is a process that breaks up sentences into words, phrases, symbols, or other meaningful elements (tokens). Note that a period ( . ) does not necessarily indicate an end of sentence (for example, Mr. Quackenbush or Boston, Mass. could occur in the middle of a sentence). Here is a list of sentence delimiters:

\r\n\r\n	Two consecutive carriage returns and new lines (for documents created in Windows)
\r\n \r\n	Two consecutive carriage returns and new lines, separated by a space
. <space></space>	Period (.) followed by an ASCII space
.\n	Period (.) followed by a new line
.\r	Period (.) followed by a carriage return
!	Exclamation point
!\n	Exclamation point followed by a new line
!\r	Exclamation point followed by a carriage return
?	Question mark
?\n	Question mark followed by a newline
?\r	Question mark followed by a carriage return
.)	Period followed by a closing parenthesis
!)	Exclamation point followed by a closing parenthesis
?)	Question mark followed by a closing parenthesis
,"	Period followed by double quotation marks.

## SENT\_n

(Multiple sentences) Takes a value for n and two or more arguments. Returns matches within n sentences. For example, the following rule returns a match for the concept node GENDER and the term he within two sentences. Suppose the GENDER concept node contains the following rule:

CLASSIFIER: male

You can then write this rule:

CONCEPT RULE: (SENT 2, " c{GENDER}", "he")

For more information, see the SENT operator.

#### SENTEND\_n

(End of sentence) Takes a value for *n* and one or more arguments. Returns matches within *n* tokens from the end of the sentence. For example, suppose the GENDER concept node contains the following rule:

CLASSIFIER: female

Then the following rule returns a match for the concept node GENDER and the term she within five tokens from the end of a sentence:

```
CONCEPT RULE: (SENTEND 5, " c{GENDER}", "she")
```

For more information, see the SENT operator.

**Note:** When you specify the value of *n*, consider that the end of the sentence is 0. Words that include hyphens are counted as one token (for example, merry-go-round is one token).

## SENTSTART\_n

(Start of sentence) Takes a value for n and one or more arguments. Returns matches within n tokens from the beginning of the sentence. For example, the following rule locates matches for the sentence The patient experienced breathing difficulty.:

```
CONCEPT RULE: (SENTSTART 5, " c{patient}", "breathing", "difficulty")
```

For more information, see the SENT operator.

**Note:** When you specify the value of n, consider that the beginning of the sentence is 0. Words that include hyphens are counted as one token (for example, merry-go-round is one token).

#### UNLESS

Takes two arguments. Restricts certain matches within the parameters that you specify when both arguments are matched in the same document. Used in rule types PREDICATE RULE and CONCEPT RULE only. Specify only these Boolean operators with the UNLESS operator: AND, SENT, DIST, ORD, and ORDDIST. The Boolean operators should appear with the second argument, as shown in the example.

For example, the following rule does not include the token river in its matches; in addition, the rule returns matches for Mississippi the state and not Mississippi the river:

```
CONCEPT RULE: (UNLESS, "river", (SENT, " c{Mississippi}", "United States"))
```

The rule ensures that river does not appear between Mississippi and United States in the matches.

Note: When you specify a concept node in a rule that uses the UNLESS operator, specify a concept nodes that contains only CLASSIFIER or REGEX rules.

## **Using the Coreference Operator**

Use the coreference operator (ref) when you want to link pronouns and other words with the canonical form (full form) of the terms that they reference.

Suppose you have a concept node **LEADERS** that includes this rule:

```
CLASSIFIER: Congressional leaders
```

You can create a concept node THEY SAID that enables they to reference its canonical form, Congressional leaders. Both forms are matched in the document.

```
C CONCEPT: c{LEADERS} said ref{they}
```

You can use the following symbols with the coreference operator (ref). Place the symbol after the ref{concept} operator.

## > (Multiple matches)

Locates multiple instances of a match that is specified by the coreference operator ( ref). For example, you might want to return the canonical form of the name Ms. Geraldine Jones each time the nickname Geri is encountered. The > symbol enables this match to occur after the first time the canonical form of the name is located.

```
C CONCEPT: c{Ms. Geraldine Jones} ref{Geri}>
```

F (Forward)

Returns only matches that occur after the coreference rule match. Sample syntax:

```
C CONCEPT: c{PERSON} as ref{TITLE} F
```

P (Preceding)

Returns only matches that occur before the coreference rule match. Sample syntax:

```
C CONCEPT: c{MILITARY BRANCH} as ref{HONOR} P
```

## **Using the Export Feature**

The Export feature enables you to find matching occurrences of terms or phrases found in CLASSIFIER rules and then export them to one or more concepts. This feature is useful for conditional matching of terms or phrases. You can export matches from multiple concepts to one concept, or to more than one concept.

**Note:** The Export feature can be used only with CLASSIFIER rules.

For example, suppose you want to find all the occurrences of the term accounts receivable that occur together with the name Sokolov, and export those matches to the concept AR. You could write the following rule in a concept node named ACCOUNT HOLDER:

```
CLASSIFIER: [export=AR:accounts receivable]:Sokolov
```

The rule first matches the term Sokolov. If that match is found, the rule checks the documents for any occurrences of the term accounts receivable and assigns any matches to the concept AR. In the list of matches for ACCOUNT HOLDER, the term Sokolov would be highlighted. In the list of matches for AR, the term accounts receivable would be highlighted. Note that in order for the rule to work, the primary term (in the example, Sokolov) needs to be present anywhere in the document before accounts receivable can be returned as a match for concept node AR.

Concepts that you are exporting to (such as AR in the example) must exist in the list of concepts and can contain additional rules (or be empty).

The following example illustrates how to export two sets of terms to the same concept.

```
CLASSIFIER: [export=text2]:text1
```

If text1 and text2 appear in a document, return text1 and text2 as separate matches for the concept where this line is located.

For example, suppose you have written the following rule:

```
CLASSIFIER: [export=SAS]:institute
```

The string SAS institute returns SAS and institute as matches to the concept where this line is located. The string institute (occurring alone) is a match, but not SAS occurring alone.

## **Using Part-of-Speech and Other Tags**

Part-of-speech tags enable you to locate matches by the part of speech that the searched item belongs to, rather than locating a specific term. These tags are useful when you know the syntax but not the exact wording of an item that you are seeking. Also included are other tags that are not considered parts of speech (such as punctuation).

Because the parts of speech are sensitive to the context in which they appear, the same word might be tagged differently, depending on the surrounding text. For example, the word will could be tagged as modal (she will be a big star someday) or noun (a last will and testament).

Part-of-speech tags are preceded by a colon (:). The tags are case-sensitive. For example, suppose you want to match an attribution for a quotation in a news article. You know that the syntax for the match will appear as Senator from state or Senator of state but you do not know the name of the senator. You can use the following rule:

```
C_CONCEPT:SENATE_TITLE _c{ cap _cap} :Prep STATE
```

The rule assumes that there is a concept SENATE\_TITLE that contains words such as majority leader, senator, and senators, and a concept STATE that includes names of states. The :Prep tag indicates a preposition (for example, from or of). A match on the C\_CONCEPT rule would occur on the text Senator Phineas Craymoor from North Carolina took the floor. However, the following text would not produce a match because the word and is not a preposition: Senators Phineas Craymoor and Garrett Garcia from North Carolina pushed the bill through.

Table 3.7 lists the part-of-speech tags in English. For tags in other languages, see Appendix 1, "Part-of-Speech Tags (for Languages Other Than English)," on page 93.

 Table 3.7
 Part-of-Speech Tags (for English)

Part-of-Speech Tag	Definition	Examples
:ABBREV	Abbreviation	etc., Ms, cm
:Acomp	Comparative adjective	cooler, luckier, worse
:Adv	Adverb	lyrically, physically
:Asup	Superlative adjective	mellowest, merriest, best
:C	Conjunction	when, yet, after, except
:date	Date	2000-02-21, 04/03/2012
:digit	Sequence of numbers	2345, 234.22, 21/234
:Det	Determiner	the, an, every
:F	Foreign	facto, klieg, modus
:inc	Unknown word	slaster, lijer
:Int	Interjection	hah, hello, tallyho
:Md	Modal	can, should, will
:N	Noun	cake, love, shoe
:Npl	Plural noun	peas, sheep, shoes
:Num	Number	one, twenty, hundred
:PN	Proper noun	SAS, Cary, Goodnight
:PossDet	Possessive determiner	our, his, my

Chapter 3 / Performing the Analysis Tasks

:PossPro	Possessive pronoun	mine, yours, hers
:PreDet	Pre-determiner	quite, such, all
:Prefix	Prefix	cross, ex, multi
:Prep	Preposition	on, under, across
:Pro	Pronoun Relative pronoun	he, one, somebody, me myself, oneself, themselves
:Ptl	Particle	away, forward, in
:sep	Separator and punctuation	;,/
:time	Time	7AM, 10:00 pm
:url	Filenames, pathnames, URL	A:/mydir/file.txt, www.sas.com
:V	Undeclined <i>be</i> , <i>do</i> , or <i>have</i> auxiliary Undeclined verb First person singular verb	be, do, have go, see, love am
:V3sg	Third person singular <i>be</i> , <i>do</i> , or <i>have</i> auxiliary Third person singular verb	is, does, has goes, sees, loves
:Ving	Present participle <i>be</i> , <i>do</i> , or <i>have</i> auxiliary  Present participle	being, doing, having bucketing, climbing
:Vpp	Past participle <i>be</i> , <i>do</i> , or <i>have</i> auxiliary Past participle	been, done, had dashed, factored, gone
:Vpt	Past tense <i>be</i> , <i>do</i> , or <i>have</i> auxiliary Past tense verb	was, were, did, have dashed, factored, went
:WAdv	Adverbial <i>wh</i>	how, when, whereby

:Wdet	Demonstrative determiner wh	which, what, whatever
:WPossPro	Possessive determiner wh	whose
:WPro	Nominal wh	whose, what, whoever

## **Using Regular Expressions (Regex)**

Use regular expressions (Regex syntax) to identify regularly occurring patterns in the text that include numbers and characters. You can use regular expressions to match patterns such as license plate numbers (example: ABX-0444), part numbers for manufacturing components (example: TMS1T3B1M5R-23), hyphenated words (example: fifty-nine), and so on.

The following guidelines apply to Regex syntax:

Include one or more characters inside square brackets ([ ]) to match the specified characters. This provides flexibility in character matching. For example, the following rule matches c, r, a s, or h:

```
REGEX: [crash]
```

If you add a plus sign (+) as follows, the rule matches the characters specified in any combination, such as rash, cash, ash, and crass (but not crashpad or crashdummy):

```
REGEX: [crash] +
```

Characters are matched within a string in sequence when represented without square brackets ([]). For example, the following rule matches only the word any (anyone or anything would not be matched):

```
REGEX: any
```

To match words that contain any, you can modify the rule to use asterisks (\*) to match other character occurrences (or none) surrounding any. For example, the following rule matches any, anyone, anything, and Many:

```
REGEX: [A-Za-z]*any[A-Za-z]*
```

You can specify a range of characters to be matched. For example, the following rule matches lowercase characters between a and f, inclusively:

```
REGEX: [a-f]
```

To add uppercase characters, use the following rule:

```
REGEX:[A-Fa-f]
```

You can specify characters that should not be matched (negated characters) by inserting a caret (^) before a set of characters. For example, the following rule matches all characters, numbers, and symbols in text except a, e, i, o, and u:

```
REGEX: [^aeiou]
```

Characters that are reserved for special meaning (metacharacters) must be escaped with a backward slash (\) to be literally matched in a regular expression. The metacharacters are: [, ], (, ), ?, \*, +, ., -, \, and |

For example, [\?] matches a question mark? in text.

Numbers are matched as-is within a string when represented without square brackets ([ ] ). For example, the following rule matches part numbers that begin with 0125- and end with a letter:

```
REGEX: 0125\-[A-Za-z]
```

Numbers are matched by specifying ranges when enclosed in square brackets ([ ] ). For example, the following rule returns a match on a number between 0 and 9:

```
REGEX: [0-9]
```

The special characters used for matching in Regex syntax can be used in combination and are shown in Table 3.8 on page 76.

 Table 3.8
 Special Characters (Metacharacters) Used in Regular Expressions

Character or Expression	Description
1	(Alternative) Indicates that matches occur when either regular expression $a$ or $b$ is matched. Example: $a \mid b$

()	Grouping mechanism (non for clarity. Example: (ababa	n-remembering). Used in expressions ab)   b)	
	(Wildcard) Matches any ch	naracter.	
%	Matches % or percent		
?	Matches 0 or 1 occurrence	es .	
*	Matches 0 or more occurre	ences	
+	Matches 1 or more occurre	ences	
{}	Indicates repetition:		
	<pre>{n} matches exactly n occurrences</pre>	{n,m} matches at least n occurrences but no more than m occurrences	
	{n,} matches at least n occurrences		
\a	Alarm (beep)		
\n	New line		
\r	Carriage return		
\t	Tab		
\f	Form feed		
\e	Escape		
\d	Digit (same as [0-9])		
\D	Not a digit (same as [^0-	9])	
\w	Word character (same as	[a-zA-Z_0-9])	
\W	Non-word character (same as [^a-zA-Z_0-9])		
\s	Whitespace character (sar	me as [ \t\n\r\f]])	

IS	Non-white-space character (same as [^ \t\n\r\f]])
\xh	Hexadecimal number, where h is a hexadecimal character
\xhh	Hexadecimal number, where h is a hexadecimal character
/00	Octal number, where o is an octal digit
\000	Octal number, where o is an octal digit

## The following restrictions apply to Regex syntax:

- Regex syntax works similarly to regular expressions in Perl; however, the two are not identical.
- Character matching for characters, numbers, or symbols that are specified inside square brackets ([]) does not occur at the word level. For example, the following rule matches the isolated letters x, y, and z, but no matching occurs for the words xylitol, yes, or recognize:

```
REGEX: [xyz]
```

If you add a plus sign (+) to match multiple occurrences (or one occurrence) as follows, the rule matches any combination of the characters that are specified, such as xzx, yz, and zyzy:

```
REGEX: [xyz] +
```

However, because word-level matching does not occur, there is no matching for words xx1, syzygy, or diy.

- You cannot refer to concepts in a Regex expression.
- Backward references to matches in the text are not supported.
- Parentheses ( ) as a grouping mechanism where matches are remembered are not supported. Parentheses are used merely for clarifying matching rules.

## **Using Morphological Expansion Symbols**

You can use morphological expansion in all rule types except CLASSIFIER and REGEX. For example, to expand the word breathe to all verb forms, which include breathes and breathing, use the following syntax for the argument: "breathe@V".

 Table 3.9
 Morphological Expansion Symbols in Concept Rules

Symbol	Description
@	Expands the concept rule to match all inflectional forms of the word in the argument. For example, the argument "wonder@" returns the matches wonder, wonders, wondered, wondering, and so on.
	<b>Note:</b> If you apply @ to a word that SAS Contextual Analysis does not recognize, no expansion occurs. Only the exact string specified before the @ is matched. For example, "grath" would not expand. Only the string grath would return a match in the rule.
@A	Expands the concept rule to match inflected comparative and superlative adjective forms of the word in the argument. For example, the argument "happy@A" returns the matches happier and happiest.  Note: If you apply @A to a word that is not an adjective, no expansion occurs.
@N	Expands the concept rule to match all inflected noun forms of the word in the argument. For example, the argument "quality@N" returns the matches quality and qualities.
	<b>Note:</b> If you apply @N to a word that is not a noun, no expansion occurs.
@V	Expands the concept rule to match all inflected verb forms of the word in the argument. For example, the argument "transfer@V" returns the matches transfer, transfers, transferred, and transferring.
	Note: If you apply @V to a word that is not a verb, no expansion occurs.

## **Adding Comments**

You can insert comments into rule definitions that have separate rules appearing on successive lines, such as CLASSIFIER rules. The comment continues until the end of the line. Comments are written as

```
# comment text
```

Note: The pound character (#) denotes a comment. If you want to match # in a rule definition, you must use a backward slash (\) as an escape character before the #. (Example: The expression 99\# attempts to match the string 99#.)

TIP You can comment out a rule by inserting a pound character (#) at the beginning of a line that contains a rule.

## **Concept Rule Types: Examples**

Examine the syntax in the examples to understand how to write different types of concept rules.

#### CLASSIFIER

Example: To extract documents that contain US airport codes, you can create a concept node named US AIRPORTS that includes these CLASSIFIER rules:

CLASSIFIER: BUF CLASSIFIER: BUR CLASSIFIER: BVK

So, documents that include a match on one or more of the airport codes BUF, BUR, or BVK, return a match for US AIRPORTS.

#### CONCEPT

Example: To extract documents that contain flight arrival information, create a concept node on time arrivals. The rule definition for on time arrivals contains the CONCEPT rule type. The CONCEPT rule type can reference the concept node US AIRPORTS, which enables airport codes to be detected. The rule definition for the

concept node on TIME ARRIVALS is as follows: CONCEPT: at US AIRPORTS on time (where US AIRPORTS includes CLASSIFIER rules that identify US airport codes).

## C CONCEPT

Example: To extract documents that include names of university professors, create a C CONCEPT rule named PROFESSORS whose definition includes this rule: C CONCEPT: Professor c{FIRSTNAME LASTNAME}. The rule indicates that matches are returned when FIRSTNAME and LASTNAME (previously defined) are found, but only when they are preceded by the word Professor. Provide the context for the match by using the modifier c and enclosing the argument that you want to match in braces ({}).

The rule modifier c indicates that the match occurs within the context of the specified concept nodes.

## NO BREAK

Example: Suppose you want to extract National Gallery of Art. You defined a concept node US ART GALLERIES that includes the CLASSIFIER rule National Gallery of Art. There also exists a concept node called CLASS TYPES that includes the CLASSIFIER rule Art. You can create the following rule that prevents a partial match on CLASS TYPES and ensures that the entire string National Gallery of Art is matched: NO BREAK: c{US ART GALLERIES}

The rule modifier c indicates that the match occurs within the context of another concept node.

Note: NO BREAK applies across the entire taxonomy regardless of where the rule appears or whether the rule is enabled or disabled.

## REMOVE ITEM

Example: Suppose you want to extract the baseball team St. Louis Cardinals, but not the football team Arizona Cardinals. You have a concept node named FOOTBALL that includes the rule CLASSIFIER: Cardinals. You have another concept node named BASEBALL that includes the rule CLASSIFIER: Cardinals. The following rule returns matches for the baseball team only:

```
REMOVE ITEM: (ALIGNED, " c{FOOTBALL}", "BASEBALL")
```

**Note:** The REMOVE\_ITEM rule type could influence matches outside of the concept node in which it is used. In this case, the rule could influence matches in the FOOTBALL rule because the rule specifies that items be removed.

#### REGEX

Example: To extract whole numbers in text (such as 1, 23, 456, and so on), use the rule

```
REGEX: [0-9] +
```

This rule requires that one or more consecutive digits occur and are without decimals.

Example: To extract a number that uses decimal notation, such as 392.55, 45.25, and 0,987654321, use the following rule:

```
REGEX: [0-9] + [, \] [0-9] +
```

This rule returns a match on any digit 0 to 9 followed by any number of digits, commas, or periods (in any combination), and then ending in a digit.

For more information about writing Regex rules, see "Using Regular Expressions (Regex)" on page 75.

## CONCEPT\_RULE

Example: Suppose you want to extract Amazon the company, not Amazon the river. You could use this rule, which would return a company name within three words of company, but not if there were nature-related words in the document.

```
CONCEPT_RULE: (AND, (DIST_3, "_c{COMPANY}", "company"), (NOT, "NATURE"))
```

#### SEQUENCE

Example: Suppose you want to extract first and last names only from a list of first, middle, and last names. You can use a SEQUENCE rule to define the arguments first and last. By using these arguments, matches are made on the concept nodes FIRST\_NAME, MIDDLE\_NAME, and LAST\_NAME, but matches are returned on only FIRST\_NAME and LAST\_NAME.

```
SEQUENCE:(first, last): _first{FIRST_NAME} MIDDLE_NAME _last{LAST_NAME}
```

## PREDICATE RULE

Example: Suppose you want to match a company to its products. You could use the following PREDICATE RULE, which assumes that the concept node COMPANY includes CLASSIFIER rules that list company names and the concept node PRODUCTS contains CLASSIFIER rules that list products. Items must appear in the same sentence.

```
PREDICATE RULE: (company, product): (SENT, " company{COMPANY}",
"produces", " product{PRODUCTS}")
```

## **Writing Category Rules: Boolean Rules**

## **Introduction to Category Rules**

Category rules resolve to true or false. "True" results in a match. Boolean rules use Boolean and proximity operators, arguments, and modifiers to define the conditions that are necessary for category matches. Category rules are simpler to write than LITI rules and are recommended when there is no need to extract specific information from the data. For a list of operators, see Table 3.10 on page 85.

Use the following syntax for a Boolean rule:

```
(OPERATOR, <argument1>, <argument2>, ...)
where arguments can be terms, strings, or nested rules.
```

## General rules for syntax:

- Boolean and proximity operators are enclosed in parentheses and separated with commas. Strings within arguments are included in quotation marks (" "). Example: (AND, "holiday", "vacation")
- Rules can be nested. Example: (AND, (OR, "courage", "courageous"), (OR, "brave", "bravery"))

- Reference a category from another category by using special syntax called tmac syntax (\_tmac). For more information, see "Using \_tmac for Referencing Categories" on page 92.
- Concept node names can be referenced in category rules. If you reference a concept node name, all concept matches will also match in the category. Concept node names must be enclosed in braces ([]) and quotation marks ("") For example, to reference the concept node GAME\_SHOWS in a category rule, you could write the rule (OR, "[GAME\_SHOWS]").

Note: Concept nodes that are named in categories might return more matches than concepts that are run outside of categories. In categories, matches on concepts are based on an "all matches" method, which returns all matches found in the text. By contrast, in concepts, matches are based on a "best match" method. The best match method detects when text that matches one concept overlaps text that matches another concept (for example, a concept that matches New York and another concept that matches New York City). When concept matches overlap and the best match method is used, only the concept that is assigned the highest number for the priority is returned (1 is the lowest). When two or more concepts have the same priority assigned, SAS Contextual Analysis selects a match by using the "longest match" method.

- The enabled or disabled status of concepts that are named in categories is ignored during category matching. As a result, the concepts are processed as if they were all enabled, regardless of whether they were previously disabled.
- Special symbols can be used to modify the rules to include, wildcards, case sensitivity, and so on. For a list of symbols, see Table 3.11 on page 90.

**Note:** XPath expressions are not supported.

## **Boolean and Proximity Operators for Category Rules**

Table 3.10 shows a list of Boolean and proximity operators that you can use to write category rules.

 Table 3.10
 Boolean and Proximity Operators for Category Rules

Operator	Description
AND	Takes one or more arguments. Matches if all arguments occur in the document, in any order. For example, the rule (AND, "King", "Louis", "XIV") returns a match if King, Louis, and XIV all occur in the document.
DIST_n	(Distance) Takes a value for $n$ and two or more arguments. Matches if all arguments occur within $n$ (or fewer) tokens of each other, regardless of their order. For example, the rule (DIST_5, "best", "picture") returns a match in the phrase the picture with the best lighting.
	<b>Note:</b> For calculation purposes, the distance between tokens is not inclusive. For example, the distance between the tokens best and show in the phrase best in show is two tokens. Words that include hyphens are counted as one token (for example, merry-go-round is one token).
END_n	(From the end of the document) Takes a value for $n$ and one or more arguments. Matches if the argument occurs within $n$ tokens from the end of the document. For example, the rule (END_35, "conclusion") returns a match if conclusion is found within 35 tokens from the last token in the document.
	<b>Note:</b> Words that include hyphens are counted as one token (for example, merry-go-round is one token).
MIN_n	(Minimum) Takes a value for $n$ and one or more arguments. Matches if the document contains at least $n$ of the arguments specified (in any order). For example, the rule (MIN_2, "Hollywood", "tinseltown", "movies") returns a match if Hollywood and movies occur in the document. However, there is no match if Hollywood occurs twice and no other arguments occur.

## MINOC\_n

(Minimum occurrence) Takes a value for n and one or more arguments. Matches if the document contains at least n occurrences of the arguments specified (in any order or combination). For example, the rule (MINOC\_2, "Hollywood", "tinseltown", "movies") returns a match if Hollywood and movies occur in the document. There is also a match if Hollywood occurs twice and no other arguments occur.

## MAXOC\_n

(Maximum occurrence) Takes a value for n and one or more arguments. Matches if the document contains n or fewer occurrences of the arguments (in any order or combination). Useful for filtering spam documents. For example, the rule (MAXOC\_8, "savings", "offer", "best") returns a match if savings occurs in the document six times. There is also a match if offer occurs in the document six times and best occurs twice.

#### MAXPAR\_n

(Maximum paragraph) Takes a value for n and one or more arguments. Matches if all arguments occur within the first n (or fewer) paragraphs of the document, in any order. For example, the rule (MAXPAR\_4, "seasonal", "herbs", "native") returns a match if seasonal occurs in paragraph 4, herbs occurs in paragraph 2, and native occurs in paragraph 2.

**Note:** MAXPAR rules work properly only when applied to data sets that contain paragraph delimiters (\n\n). MAXPAR cannot be applied on the **Test Rules** tab. MAXPAR also cannot be applied in the **Categories** tab to data that is contained in folders.

#### MAXSENT n

(Maximum sentence) Takes a value for n and one or more arguments. Matches if all arguments occur within the first n sentences of the document, in any order. For example, the rule (MAXSENT\_4, "weight loss", "plan") returns a match if weight loss and plan occur in sentence 3 of the document. For a list of sentence delimiters, see the SENT operator.

#### NOT

Takes one argument. Matches if the argument does not occur in the document. Must be used with the AND operator. For example, the rule (AND, (OR, "cinema", "theater", "theatre"), (NOT, "Broadway")) returns a match if cinema, theater, or theatre occur in the document and Broadway does not.

**Note:** The NOT operator applies across the entire document. If you specify the OR operator in addition to the AND operator, you must enclose the OR arguments in parentheses.

## NOTIN

(Not in) Takes two arguments and matches if the first argument does not appear within the second argument. For example, the rule (NOTIN, "butter", "peanut butter") identifies butter when it does not appear within the noun phrase peanut butter. This sentence returns a match: Early American colonists churned their own butter.

## NOTINDIST n

(Not in distance) Takes a value for *n* and two arguments. Matches if the arguments do not occur within *n* tokens of each other, or if the first argument listed in the rule occurs in the document and the second argument does not. For example, the rule (NOTINDIST 3 "orange", "green") returns a match if orange and green do not occur within three tokens of each other, or if only orange appears in the document. The following sentence returns a match because the tokens that are specified in the rule are more than three tokens apart: How green is my valley, how orange is the sunset?

**Note:** For calculation purposes, the distance between tokens is not inclusive. For example, the distance between the tokens best and show in the phrase best in show is two tokens. Words that include hyphens are counted as one token (for example, merry-go-round is one token).

#### NOTINPAR

(Not in paragraph) Takes two or more arguments and matches if all arguments occur within the document but appear in separate paragraphs. For example, the rule (NOTINPAR, "China", "export") returns a match if China and export occur in separate paragraphs (without the other argument present).

**Note:** NOTINPAR rules work properly only when applied to data sets that contain paragraph delimiters (\n\n). NOTINPAR cannot be applied on the **Test Rules** tab. NOTINPAR also cannot be applied in the **Categories** tab to data that is contained in folders.

#### NOTINSENT

(Not in sentence) Takes two or more arguments and matches only if all arguments occur within the document but appear in separate sentences. For example, the rule (NOTINSENT, "China", "trade") returns a match if China and trade occur in separate sentences (without the other argument present), as in China is our biggest partner. The trade it generates is huge. For a list of sentence delimiters, see the SENT operator.

#### OR

Takes one or more arguments. Matches if at least one argument occurs in the document. For example, the rule (OR, "U.S.", "US ", "United States") returns a match if one or more of the items U.S., US, or United States appear in the document.

**Note:** Rules that are generated by SAS Contextual Analysis nest the OR operator within the AND operator. However, the OR operator can stand alone.

## ORD

(Order) Takes one or more arguments. Matches if all of the arguments occur in the order that is specified in the rule. It cannot be used with SENT (or any other operator that limits the scope of matches). For example, the rule (ORD, "warranty", "claim", "denied") returns a match in the sentence The warranty claim for the washing machine was denied.

## ORDDIST\_n

(Order and distance) Takes a value for n and two or more arguments. Matches if both arguments occur in the same order that is specified in the rule and if both arguments are within n tokens of each other. For example, the rule (ORDDIST\_5, "elementary", "statistics") returns a match in the phrase the teacher introduced elementary statistics.

**Note:** For calculation purposes, the distance between tokens is not inclusive. For example, the distance between the tokens best and show in the phrase best in show is two tokens. Words that include hyphens are counted as one token (for example, merry-go-round is one token).

#### PAR

(Paragraph) Takes one or more arguments. Matches if all the arguments occur in a single paragraph, in any order. For example, the rule (PAR, "director", "budget") returns a match if the paragraph includes both director and budget.

**Note:** PAR rules work properly only when applied to data sets that contain paragraph delimiters (\n\n). PAR cannot be applied on the **Test Rules** tab. PAR also cannot be applied in the **Categories** tab to data that is contained in folders.

## PARPOS n

(Paragraph position) Takes a value for *n* and one or more arguments. Matches if all arguments occur within the  $n^{th}$  paragraph, in any order. For example, the rule (PARAPOS\_2, "journalists", "detained", "overseas") returns a match if journalists, detained, and overseas occur within paragraph 2 of the document.

**Note:** PARPOS rules work properly only when applied to data sets that contain paragraph delimiters (\n\n). PARPOS cannot be applied on the **Test Rules** tab. PARPOS also cannot be applied in the **Categories** tab to data that is contained in folders.

#### SENT

(Sentence) Takes two or more arguments. Matches if all the arguments occur in the same sentence, in any order. For example, the rule (SENT, "growth", "hormone") returns a match in the sentence Patients who take a growth hormone might experience side effects. Sentence delimiters are as follows:

deficited definitions are as follows.			
\r\n\r\n	Two consecutive carriage returns and new lines (for documents created in Windows)		
\r\n \r\n	Two consecutive carriage returns and new lines, separated by a space		
. <space></space>	Period (.) followed by an ASCII space		
.\n	Period (.) followed by a new line		
.\r	Period (.) followed by a carriage return		
!	Exclamation point		
!\n	Exclamation point followed by a new line		
!\r	Exclamation point followed by a carriage return		
?	Question mark		
?\n	Question mark followed by a newline		
?\r	Question mark followed by a carriage return		
.)	Period followed by a closing parenthesis		
!)	Exclamation point followed by a closing parenthesis		
?)	Question mark followed by a closing parenthesis		
."	Period followed by double quotation marks		

START_n	(From the start of the document) Takes a value for $n$ and one or more arguments. Matches if the argument occurs within $n$ words from the start of the document. For example, the rule (START_22, "infection") returns a match if infection occurs within 22 words of the first word in the document.
	<b>Note:</b> Words that include hyphens are counted as one token (for example, merry-go-round is one word).

## **Using Symbols in Boolean Rules**

You can use the symbols listed in Table 3.11 to modify your Boolean rules for category matching. Symbols are written as suffixes to strings in arguments. For example, to expand the word breathe to all inflected verb forms, which include breathes and breathing, use the following syntax for the argument: "breathe@V".

 Table 3.11
 Special Symbols Used in Boolean Rules

Symbol	Description
*	(Wildcard matching) Matches any characters that occur at the beginning or end of the word. For example, the argument "travel*" returns the matches travels, traveled, traveler, traveling, and so on. The argument "*room" matches bedroom, cloakroom, ballroom, room, and so on.
٨	(Beginning of sentence) Starts searching at the beginning of the sentence to find a match. For example, the argument "^Independent" returns a match in this sentence: Independent research was conducted.
	Note: Tokens (words, phrases, symbols, or other meaningful elements) need to be entered specifically to be considered for matching. For example, if you are searching for **In this case, use the argument "^\*\In this case". Also note that backward slashes (\) are used as escape characters for the asterisks (*) so that the asterisks are not treated as wildcards.

\$	(End of sentence) Starts searching at the end of the sentence to find a match. For example, the argument "deleted.\$" returns a match on the following sentence: All the files were hastily deleted.
	Note: Tokens (words, phrases, symbols, or other meaningful elements) need to be entered specifically to be considered for matching. For example, the argument "deleted\$" would not produce a match on the following sentence: All the files were hastily deleted. because the ending period (.) was not specified.
@	(Morphological expansion) Expands the category rule to match all inflectional forms of the word in the argument. For example, the argument "wonder@" returns the matches wonder, wonders, wondered, wondering, and so on (but does not return a match on wonderful).
	<b>Note:</b> If you apply @ to a word that SAS Contextual Analysis does not recognize, no expansion occurs. Only the exact string specified before the @ is returned. For example, "grath" would not expand. Only the string grath would return a match in the rule.
@A	(Morphological expansion for adjectives) Expands the category rule to
	match inflected comparative and superlative adjective forms of the word in the argument. For example, the argument "happy@A" returns the matches happier and happiest.
	in the argument. For example, the argument "happy@A" returns the
@N	in the argument. For example, the argument "happy@A" returns the matches happier and happiest.  Note: If you apply @A to a word that is not an adjective, no expansion
@N	in the argument. For example, the argument "happy@A" returns the matches happier and happiest.  Note: If you apply @A to a word that is not an adjective, no expansion occurs.  (Morphological expansion for nouns) Expands the category rule to match all noun forms of the word in the argument. For example, the argument
@N @V	in the argument. For example, the argument "happy@A" returns the matches happier and happiest.  Note: If you apply @A to a word that is not an adjective, no expansion occurs.  (Morphological expansion for nouns) Expands the category rule to match all noun forms of the word in the argument. For example, the argument "quality@N" returns the matches quality and qualities.  Note: If you apply @N to a word that is not a noun, no expansion

_L	(Literal matching) Matches a literal string. Useful when you want to match a string that includes symbols. For example, the argument "\$USD_L" returns the match \$USD.
	<b>Note:</b> Tokens (words, phrases, symbols, or other meaningful elements) need to be entered specifically to be considered for matching.
_C	(Case matching) Specifies case-sensitive matching. For example, the argument "Iris_C" returns the match Iris, but not iris.

## **Using \_tmac for Referencing Categories**

Referencing a category enables you to use the rules in an existing category without having to duplicate the rules. Use tmac syntax (\_tmac) to reference an existing category in a category rule. The definition of the existing rule is processed in the category that references it.

To reference a category, you must identify its path. All category paths begin with @Top/. From there, you can specify the path by following the category hierarchy.

For example, suppose you have the following category structure under **All Categories**:

NAME

**FIRST** 

LAST

You would reference the category FIRST as @Top/NAME/FIRST.

You can use the tmac syntax with Boolean operators. For example, suppose you want to reference the category **first** from a category called **first\_name**. You could add this rule in the **first\_name** definition:

```
(OR, tmac: "@Top/NAME/FIRST")
```

To enforce a first name followed by last name (FIRST LAST), you could add this rule in a category called COMPLETE\_NAME::

```
(ORD,_tmac:"@Top/NAME/FIRST",_tmac:"@Top/NAME/LAST")
```

The definitions written in FIRST and LAST are automatically processed.

## **Appendix 1**

# Part-of-Speech Tags (for Languages Other Than English)

Introduction to Part-of-Speech and Other Tags	94
Part-of-Speech Tags	94
Arabic	94
Chinese	96
Croatian	97
Czech	98
Danish	100
Dutch	101
Farsi	103
Finnish	104
French	105
German	107
Greek	108
Hebrew	109
Hindi	
Hungarian	111
Indonesian	112
Italian	
Japanese	114
Korean	118
Norwegian	
Polish	121

Portuguese	122
Romanian	123
Russian	125
Slovak	127
Slovene	129
Spanish	130
Swedish	131
Thai	132
Turkish	134
Vietnamese	135

# Introduction to Part-of-Speech and Other Tags

The part-of-speech tags for languages other than English are listed in the following tables. Also included are other tags that are not considered parts of speech (such as punctuation). All tags are case-sensitive and are preceded by a colon (:) in concept rules. For more information, including English tags, see "Using Part-of-Speech and Other Tags" on page 72.

## **Part-of-Speech Tags**

## **Arabic**

Table A1.1Part-of-Speech Tags for Arabic

Part-of-Speech Tag	Description	Examples
:ADJ	Adjective	أبدي, أثري
:ADV	Adverb	أيضا, ربما

:CONJ	Conjunction	بل, حتى
:DET	Determiner	ال
:DIALECT	Dialect	آسم, أثول
:FUT	Future particle	س, سوف
:INTERJ	Interjection	أجل, لا
:INTERROG	Interrogative	أين, عمّا
:NEGPART	Negative particle	لم
:NOUN	Noun	تفاحة, شجرة
:NUM	Number	آلاف, أربعة
:PART	Particle	قد, لقد
:PREP	Preposition	إلا, على
:PRON	Pronoun	أنا, أنت
:PROP	Proper noun	أمريكا
:PUNC	Punctuation	٠, ٩
:CV	Imperative verb	انتيا, العبان
:IV	Present verb	تأتون, تلعبا
:PV	Past verb	أتتا, لعبت
:ASCII	English word	memory, tablets
:DEFAULT	Unknown word	اعتياديًا, وشيئً
:NUMBER	Number	1.8, 200
:URL	URL	http://www.sas.com

## Chinese

 Table A1.2
 Part-of-Speech Tags for Chinese

Part-of-Speech Tag	Description	Examples
:A	Adjective	俊俏, 开心, 兇險, 凌亂
:ASCII	ASCII characters	sas, do, happy, day2136456
:C	Conjunction	或, 与, 雖然
:D	Adverb	非常, 偏偏, 稍微, 永遠
:digit	Number	1051, 1.9
:E	Interjection	咦, 呸, 哦喲
:F	Location / direction	中間,下边,南侧
:G	Other morpheme	馨, 慚
:H	Other prefix	亚, 非
:K	Other suffix	们, 者, 們
:L	Idiom (chengyu)	囫囵吞枣, 博古通今, 一廂情願
:M	Quantifier	十, 卅, 成千上万, 上萬, 1 0 5 1
:N	Noun	人, 桌子, 香蕉, 枷鎖
:NR	Proper noun, name	习近平, 梁振英, 奥巴马
:NS	Proper noun, geographic	中国,美國,山東
:NT	Proper noun, organization	北京大学, 上汽集團
:NZ	Proper noun, miscellaneous	潘婷, 劍南春

:O	Onomatopoeia	吱呀, 叽叽喳喳, 劈裏啪啦
:P	Preposition	依照, 对于
:Q	Classifier	个, 斤, 艘, 加侖
:R	Pronoun	我, 他們, 这
:S	Subcountry location (general; specifics only within sinosphere)	地上, 上空, 高处, 內廳
:T	Temporal phrase	今天, 夜间, 十月, 去歲
:U	Particle	的, 了, 着
:UNKNOWN	Unknown word	婳, 繟
:V	Verb	看, 认为, 彈奏, 徵納
:W	Punctuation or symbols	!, o, \$, ¥
:Y	Interjectional particle	吧, 吗, 麽

## Croatian

 Table A1.3
 Part-of-Speech Tags for Croation

Part-of-Speech Tag	Description	Examples
:Adj	Adjective	svaki, hrvatskim, koje
:Abbr	Abbreviation	dr, itd.
:Adv	Adverb	uistinu, tamo
:Conj	Conjunction	a, ali, kad
:Interj	Interjection	hej, hajde, oh

:Noun	Noun	dan, april, http:// www.sas.com
:Part	Particle	ne, bilo (as in "bilo koje")
:Prep	Preposition	sa, bez, o
:Pron	Pronoun	ja, me, ih, nas, vam, njihovoj, svašta
:Aux	Auxiliary verb	bih, je (before a main verb)
:Verb	Main verb	voli, došao, pozvala, dođite
:Num	Number	2, dva, sedmi, 1.23.2015
:time	Time	23:30:01
:Punct	Separator or punctuation	, .
:Prop	Proper noun	Aleksandar, Jelenu, Gorenje, Zagreb

## Czech

 Table A1.4
 Part-of-Speech Tags for Czech

Part-of-Speech Tag	Description	Examples
:A	Adjective	duchovní, celý
:ADJIND	Indefinite adjective	všechny, čertvíjaký
:ADJINT	Interrogative adjective	která, jakém
:ADV	Adverb	například, dál, zároveň, někam, ne
:CONJ	Conjunction	a, nebo

:DEM	Demonstrative	tomto, tím
:INTJ	Interjection	ahoj, fuj
:N	Noun	autorů, lidem
:NADJ	Negative adjective	žádnej
:NPRO	Negative pronoun	nikoho, nic
:NUM	Numeral	tři, dvoje
:NUMO	Ordinal numeral	šestatřicáté
:digit	Number	33, 1844, 14.3.2014
:POSS	Possessive adjective	její, mou
:PREP	Preposition	V, Z
:PRO	Pronoun	kdo, sobě, nás
:V	Verb	nebyl, jdou
:sep	Separator or punctuation	.,:
:PN	Proper noun	Pavel, Valenta, Chotěbořským
:inc	Unknown foreign word	mp3, larger
:time	Time	23:30:01
:url	URL	www.sas.com, http:// www.sas.com

#### Danish

 Table A1.5
 Part-of-Speech Tags for Danish

Part-of-Speech Tag	Description	Example
:A	Adjective	socialest, udartendere
:ABB	Abbreviation	DVS, FL
:ADV	Adverb	sydsydøst
:CONJ	Conjunction	Såsom
:DET	Determiner	dens, hans
:INT	Interjection	joh, pøj
:INV	Invariable	ibm, netscape
:N	Noun	thyboernes, centerer
:NUM	Spelled out number	tyvefem, tredive
:PN	Proper noun	Egholm, Puccini
:PNF	First name	Franck, Carlos
:PNG	Geographical name	Mallorca
:PNH	Last name	Groth, Leth
:PNO	Organization	Renault, Corel
:PREP	Preposition	fra, trods
:PRO	Pronoun	dens
:PROP	Personal pronoun	jerselv, sigselv

:TIM	Month or day	tirsdag
:VB	Infinitive	opofre, læse
:VE	Present participle	læsende
:VH	Past participle	tredjebehandlet
:VI	Passive	fuldkommengøredes, bemyndiges
:VJ	Preterite	læste
:VP	Present	anvender, bliver
:VY	Imperative	tilvirk
:date	Date	23-12-2012, 12/12/2012
:time	Time	:23:50, 09:23
:digit	Digit	2012, 12.23
:url	Internet address	http://www.sas.com
:sep	Separator or punctuation	.,;
:inc	Unknown word	bl, erne

#### **Dutch**

 Table A1.6
 Part-of-Speech Tags for Dutch

Part-of-Speech Tag	Definition	Examples
:A	Adjective	betrouwbaar, gelukkig, mooi
:ABB	Abbreviation	enz, kg, zgn

**102** Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:ADV	Adverb	eenmaal, hier, nu
:CONJ	Conjunction	als, doch, hoe
:DET	Determiner	de, der, een
:digit	Number	21
:DNUM	Determiner, number	acht, elf, miljard, duizend
:inc	Unknown word	XIXX
:N	Noun	geluk , schoonheid
:PFX	Prefix	anti
:PN	Proper noun	Amerika, Nederland
:PREP	Preposition	met, per, te, van
.DDEDDET	Decembration and determines	t t
:PREPDET	Preposition and determiner contraction	ten, ter
:PREPUE I		alles, beide, hetgeen
	contraction	· 
:PRO	Pronoun	alles, beide, hetgeen
:PRO :sep	Pronoun Separator or punctuation	alles, beide, hetgeen
:PRO :sep :url	Pronoun Separator or punctuation URL	alles, beide, hetgeen , http://www.sas.com
:PRO :sep :url :V	contraction  Pronoun  Separator or punctuation  URL  Verb	alles, beide, hetgeen , http://www.sas.com helpt, vernieuwt
:PRO :sep :url :V :VB	contraction  Pronoun  Separator or punctuation  URL  Verb  Infinitive	alles, beide, hetgeen , http://www.sas.com helpt, vernieuwt helpen, vernieuwen
:PRO :sep :url :V :VB	contraction  Pronoun  Separator or punctuation  URL  Verb  Infinitive  Present progressive	alles, beide, hetgeen  ,  http://www.sas.com  helpt, vernieuwt  helpen, vernieuwen  helpende, vernieuwende

## **Farsi**

 Table A1.7
 Part-of-Speech Tags in Farsi

Part-of-Speech Tag	Description	Examples
:A	Adjective	خوشگل, خوشحال
:Acomp	Comparative adjective	خوشگلتر, خوشحالتر
:Asup	Superlative adjective	خوشگلترين, خوشحالترين
:Appl	Participle used as adjective	آسایانیده, آبانانده
:ADV	Adverb	هنوز, آنگه, ابتدائا:
:CLASS	Classifier	باب, تخته, رأس
:CONJ	Conjunction	اگر, تااینکه
:DET	Determiner	اون, این
:INTJ	Interjection	آه, آفرین, ای
:N	Noun	آذوقه, آرنج, چشم
:Npl	Plural noun	آرنجها, چشمها
:NUM	Numeral	دو, صد, میلیون
:NUMord	Ordinal numeral	دومین, سوم, صدمین
:PN	Proper noun	اسر ائيل, آتوسا
:PPOS	Preposition	از, الا, چون
:PRO	Pronoun	ن, او, شما
:PUNC	Punctuation or symbol	"( ? %

**104** Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:Vinf	Infinitive (usage similar to English gerund)	خواندن ,خوردن
:V	Verb	بخوان, بخوانم, خواندم
:ASCII	ASCII characters and digits	happy, 2017, love123
:DEFAULT	Unknown word	بخوانبخوان

## **Finnish**

 Table A1.8
 Part-of-Speech Tags for Finnish

Part-of-Speech Tag	Definition	Examples
:A	Adjective	loistava, korkea
:ADV	Adverb	ohitse, juuri
:CLX	Clitic	kinko, pas
:CONJC	Coordinating conjunction	ja, vaan
:CONJS	Subordinating conjunction	ellei, jotta
:date	Date	12-14, 2001-12-02
:digit	Number	1234, 7
:inc	Unknown word	auttonkkan, eggs
:N	Noun	siltoineen, postiksi
:PN	Proper noun	Pertti, Fazer
:PREP	Preposition	pitkin, kanssaan
:PRO	Pronoun	noihin, muussa, ketkä

:sep	Separator or punctuation	; / +
:time	Time	12:00:00, 7PM
:PROP	Personal pronoun	sinun, heissä, me
:url	URL	http://www.sas.com
:VB	Infinitive verb	heilahtamassa, heilauttaen, olla
:VC	Potential present verb	lähennemme, luvannette
:VE	Present participle verb	kumarrettava, ilmaisevaa
:VH	Past participle verb	jaettu, ilmaistu
:VJ	Indicative preterit verb	meditoitpa, matkattu
:VP	Indicative present verb	ihastele, hörähdä
:VS	Conditional present verb	omistautuisi, hehkuisikaan
:VY	Imperative verb	parannuttako, pakkaa

## French

 Table A1.9
 Part-of-Speech Tags for French

Part-of-Speech Tag	Definition	Examples
:A	Adjective	comparable, compassionnelle, intraduisibles
:ADV	Adverb	plutôt, individuellement
:CONJC	Coordinating conjunction	et, ou

Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:CONJS	Subordinating conjunction	lorsque, puisque
:DET	Determiner	sa, tes
:digit	Number	123, 12.3, 12.3.2003, 12/3/2003
:inc	Unknown word	analytics
:INTJ	Interjection	tralala, zzz
:N	Noun	zèbre, encyclopédie
:PN	Proper noun	Eurotunnel, Égypte
:PFX	Prefix	anglo, éco
:PREP	Preposition	après, jusque
:PREPDET	Preposition and determiner contraction	aux, du
:PRO	Pronoun	lui, ce
:sep		
.эер	Separator or punctuation	, . !
:url	Separator or punctuation  URL	http://www.sas.com
	•	
:url	URL	http://www.sas.com
:url	URL Verb	http://www.sas.com vais, obligez
:url :V :Vpp	Verb Past participle	http://www.sas.com vais, obligez mangé, relaxée, travaillées
:url :V :Vpp :VB	URL Verb Past participle Infinitive	http://www.sas.com  vais, obligez  mangé, relaxée, travaillées  traduire, rompre

#### German

 Table A1.10
 Part-of-Speech Tags for German

Part-of-Speech Tag	Definition	Examples
:A	Adjective	schön, zuverlässig
:ADV	Adverb	gern, sehr
:CONJ	Conjunction	und, oder
:CPO	Compounding (prefix only)	Lustigkeits
:DET	Determiner	der, eine
:digit	Number	21
:DNUM	Determiner, number	fünf, zwölf
:EMP	Emphatic/intensifier	ganz
:inc	Unknown word	XIXX
:N	Noun	Schönheit, Zuverlässigkeit
:PFX	Prefix	Irr, lob
:PN	Proper noun	Mozart
:PN.gen	Proper noun, genitive	Nirvanas
:PNG.dat	Proper noun, geographic, dative	Niederlanden
:PREDET	Predeterminer	manch
:PREP	Preposition	kontra, ober
:PRO	Pronoun	er, sie

**108** Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:PXPRO	Pronominal adverb	heraus
:sep	Separator or punctuation	,
:url	URL	http://www.sas.com
:V	Verb	ging, half
:VI	Infinite (infinitives and participles)	gehen, helfen

#### Greek

 Table A1.11
 Part-of-Speech Tags for Greek

Part-of-Speech Tag	Description	Examples
:A	Adjective	ενορμητικός, άβαθος
:ADV	:Adverb	πολύ, επίσης
:CONJ	Conjunction	και, αλλά
:DET	Determiner	ένας, ο
:INTJ	Interjection	χαίρε, όπα
:N	Noun	μήλο, δέντρο
:N :PART	Noun Particle	μήλο, δέντρο πάρα
:PART	Particle	πάρα
:PART :PREP	Particle Preposition	πάρα άχρι, διά
:PART :PREP :PRO	Particle Preposition Pronoun	πάρα άχρι, διά εσύ, αυτός

:VI	Infinite verb	έρχονταν, παίζαμε
:VP	Present verb	έρχεσαι, παίζουμε
:VPP	Present participle verb	b παίζοντας, παίρνοντάς
:VPS	Present subjunctive verb	αιμοδοτήσουν, κατασκευαστώ
:VY	Imperative verb	έλα
:url	URL	http://www.sas.com
:date	Date	2015-12
:digit	Number	1, 20
:sep	Separator or punctuation	.,»
:inc	Unknown word	Χγh
:time	Time	23:59
:PN	Proper noun	Μάντσεστερ

#### **Hebrew**

 Table A1.12
 Part-of-Speech Tags for Hebrew

Part-of-Speech	Descriptions	Examples
:A	Adjective	יפה, אדיר
:ADV	Adverb	באמת, בבטחה
:CONJ	Conjunction	או, בגלל
:INT	Interjection	אוף, אהה
:N	Noun	רחוב, ברחוב, אבזור, אבטחה

:PN	Proper noun	ישראל, אבוג'ה, אדוארד
:PREP	Preposition	אודות, אצל
:PRO	אנחנו, באתה, ה"הן	אנחנו, באתה, ה"הן
:QUANT	Quantifier	אחד, ביליון, שתיהן
:V	Verb	שמח, אבטח, אהבו
:WPro	Interrogative pronoun	מהיכן
:date	Date	12/31/2016, 2016-12-31
:digit	Number	100, 6,666, 6.000
:inc	Unknown word	happy, happy123, בוויטנאם
:sep	Separator or punctuation	.,!-
:time	Time	14:30:30
:url	:URL	http://www.sas.com

#### Hindi

 Table A1.13
 Part-of-Speech Tags for Hindi

Part-of-Speech Tag	Definition	Examples
:A	Adjective	ज्ञात, ज्ञानी
:PRO	Pronoun	तेरा, मेरा
:N	Noun	मेयर, मैग्नोलिया
:ADV	Adverb	यथायोग्य, यथोचित
:CONJ	Conjunction	यदि, यद्यपि

:DET	Determiner	ऎसा, इसी
:INTJ	Interjection	आह, अहा
:NUM	Number	अस्सी, अड़तालीस
:PN	Proper noun	अग्नीवो
:POST	Particles	का, का
:V	Verb	खरीदना, गुजर
:PUNC	Separator or punctuation	Ι, ΙΙ
:sep	Separator or punctuation	,.)
:inc	Unknown words	आिद, २२५
:digit	Number	0, 3

# Hungarian

 Table A1.14
 Part-of-Speech Tags for Hungarian

Part-of-Speech Tag	Description	Examples
:A	Adjective	bámulatos, isteni
:ABR	Abbreviation	stb.
:ADV	Adverb	amerről, hova
:CONJ	Conjunction	és
:DET	Determiner	az, egy
:DNUM	Number	ötvenhetedik, hat, 13.2, 9
:INT	Interjection	Teringettét

**112** Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:N	Noun	szállodánkba, ablakok
:PART	Particle	nélkül, múlva
:PN	Proper noun	szaúdiakat, Zsoltinak
:PRO	Pronoun	őneki, ti
:V	Verb	megtudhattuk, elmentek
:VPFX	Verbal prefix	meg, ki
:date	Date	12/21/2013, 03/04/2011
:time	Time	23:50, 09:21
:url	URL	http://www.sas.com
:inc	Unknown word	txt, doc

#### Indonesian

 Table A1.15
 Part-of-Speech Tags for Indonesian

Part-of-Speech Tag	Definition	Examples
:A	Adjective	lonjong, menjengkelkan
:N	Noun	kosmologiku, lotengnya
:ADV	Adverb	mingguan, perlahan
:CONJ	Conjunction	sambil, biarpun
:V	Verb	biaskanlah, membuntutiku
:PREP	Preposition	dari
:ABBREV	Abbreviations	dpa

:DET	Determiners	sebuah
:NUM	number words	empat, delapan
:INTJ	Interjections	hai, hoi
:PRO	Pronoun	dikau, engkau
:PN	Proper noun	irlandia, filipina
:PHR	Phrasal; the word can be combined with another word to form a phrase	sebiru, secantik
:sep	Separator or punctuation	"(,
:inc	Unknown words	jpg, png
:digit	Number	22, 490
:url	URL	www.jakarta.go.id
:date	Date	12/31/2016

## Italian

 Table A1.16
 Part-of-Speech Tags for Italian

Part-of-Speech Tag	Definition	Examples
:A	Adjective	affidabile, bellissimo, felice
:AVV	Adverb	felicemente, rapidamente
:CONG	Conjunction	ma, oppure, sebbene
:DET	Determiner	il, la, uno
:digit	Number	21

:ESC	Interjection	ah, ahimè
:inc	Unknown word	Xrxx
:N	Noun	affidabilità, bellezza, felicità
:PN	Proper noun	Roma, Italia
:PRON	Pronoun	io, ne, tu
:PREFIX	Prefix	anti, ri
:PREP	Preposition	con, in, per
:sep	Separator or punctuation	,
•	'	•
:SUFFIX	Suffix	anza, issimo
·		
:SUFFIX	Suffix	anza, issimo
:SUFFIX	Suffix	anza, issimo  http://www.sas.com
:SUFFIX :url	Suffix URL Verb	anza, issimo  http://www.sas.com  andare, vedono
:SUFFIX :url :V :VGerund	Suffix URL Verb Gerund	anza, issimo  http://www.sas.com  andare, vedono  andando, vedendo

#### **Japanese**

 Table A1.17
 Part-of-Speech Tags for Japanese

Part-of-Speech Tag	Description	Examples
:AJ	Adjective	長い, 忙しい,便利だ
:AV	Adverb	いかが, やはり
:AVC	Adverbs of form or condition	直に, ぐっすり

:AVD	Adverb of degree	とっても, 大して
:AVE	Adverb of evaluation	たまたま, 無論
:AVF	Adverb of frequency	あくまで, しばしば
:AVO	Adverb of opinion	いわば, 概して
:AVQ	Adverb of quantity	大方, いくら
:AVS	Adverb of statement	いかに, あたかも
:AVT	Adverb of tense or aspect	急遽, 直ぐ
:AX	Auxiliary verbs	べきだ, らしい, ようだ
:CN	Conjunction	並びに, 但し, だけど
:CP	Copula	だ, なんだ
:DA	Adverbial demonstrative	こう, そう, あのように
:DM	Prenominal demonstrative	この, あの, そん な
:DN	Pronoun	あれ, こちら, あそこ
:MD	Prenominal modifier	小さな, 主たる, 色ん な
:IT	Interjection	あれれ, あ~, え えと
:NA	Adverbial noun	おおむね, なにぶん
:NC	Common noun	風, 学校, 雑誌
:NK	Content noun	の, もの, こと
:NT	Noun of time	長年, 夏, 先月
:NV	Verbal noun	請求, 弁解, 勉強
:NP	Proper noun	WTO繊維協定, 米州

Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:NH	Proper noun of Person	中川秀直, 中川浩明, 中川勝
:NHM	Proper noun of Given name	奈江子, 太郎, 那惠子
:NHS	Proper noun of Family name	鈴木, 佐藤, 田中
:NPO	Proper noun of Organization	米軍,米国,米国際貿易委員会
:NL	Proper noun of Place	米国,越南,奈央島
:NN	Numeral	千, 零, 6
:PC	Particles of case marker	を, で, の, へ
:PE	Particles that appear at the end of the sentence	っけ, な, なぁ
:PN	Particles that combine nominals	ないし, ないしは, 並びに
:PP	Particles that combine clauses	ながら, なら, のに
:PQ	Particles of quotation	て, と, っと
:PS	Particles that mean <i>only</i> or <i>too</i>	も, のみ, くらい
:PRJ1	Prefixes to i-adjective	か, こ, 真
:PRJ2	Prefixes to na-adjective	無, 不, 非
:PRN	Prefixes to nominals	高, 前, 全
:PRV	Prefixes to predicates	相, 猛, 最
:SJN	Suffixes to nouns and configure adjectives	っぽい, くさい
:SJV	Suffixes to verbs and configure adjectives	たい, づらい

:SNA	Suffixes to adjectives and configure nouns	さ
:SNC	Suffixes to classifiers and configure nouns	せんち, ペーじ
:SNN	Suffixes to nouns	っ子, 中, 所
:SNV	Suffixes to verbs and configure nouns	かた, っぷり
:SV	Suffixes to verbs	せる, れる, 上げる
:V1	Ichidan Verb	治せる, 泣ける, 叫べる
:V5	Godan Verb	直す, 長びく, 産む
:VK	Kuru Verb	来る
:VS1	Suru Verb	する
:VS2	Suru Verb d	賀する, 刑する, 御する
:VSN	Suru Verb	きりきり, 毅然と
:VZ	Zuru verb	準ずる, 同ずる
:SC	Special category-comma	` ,
:SCP	Special category-closed parentheses	) » ]
:SOP	Special category-opened parentheses	( 《 [
:SK	Special category-other symbols	? ~
:SP	Special category-period	o ·
:SS	Special category-space	
:digit	Number	1.0, 10

**118** Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:sep	Separator or punctuation	٠,
:KATAKANA	Unknown word in katakana	ポータブルオプション, オブ ザベーション
:HIRAGANA	Unknown word in hirakana	きんぽうげ
:UNKNOWN	Unknown word	嘘, 甦
:ASCII	English word	Display, Momente

#### Korean

 Table A1.18
 Part-of-Speech Tags for Korean

Part-of-Speech Tag	Description	Examples
:AD	Adverb	매우, 정말, 빨리
:AJ	Adjective	예쁘다, 귀엽다, 차분하다
:ASCII	Foreign	Korean, iPhone, SK
:DATE	Date	2015-04-28, 20150428
:DEFAULT	Unknown word	하페즈, 샤리프, 쿠레쉬
:GAC	Case grammatical affix	가, 를, 로
:GAD	Determinative grammatical affix	은, 을, 는
:GAH	Change grammatical affix	이다, 기, 음
:GAJ	Conjunctive grammatical affix	는데, 는지, 느라고
:GAP	Predicate grammatical affix	다, 습니다, 더구만
:GAR	Respect grammatical affix	시, 으시, 옵

:GAT	Time grammatical affix	겠, 었, 였었
:GAX	Auxiliary grammatical affix	도, 만, 까지
:IJ	Interjection	아, 네, 그래
:NN	Noun	하늘, 산, 바다
:NNB	Bound noun	것, 수, 개
:NNP	Proper noun	서울, 이순신, 국립국어원
:NUMBER	Numeral	하나, 둘, 셋
:PF	Prefix	제-, 햇-, 명-
:PN	Prenoun	각, 첫,기초적
:PR	Pronoun	이것, 언제, 이분
:PUNC	Punctuation	.?!()
:SF	Suffix	-꾼, 꾸러기, -감
:TIME	Time	23:59:59
:URL	URL	http://www.sas.com
:VB	Verb	웃다, 뛰다, 날다

# Norwegian

 Table A1.19
 Part-of-Speech Tags for Norwegian

Part-of-Speech Tag	Description	Examples
:A	Adjective	leket
:ABV	Abbreviation	mfl, mht

Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:ADV	Adverbr	alltid, framover
:CONJ	Conjunction	som
:DET	Determiner	dens
:IMRK	Auxiliary	å
:INV	Invariable (foreign word)	ruskursus, ørknen
:N	Noun	anordningen, tydeets
:NUM	Spelled out number	tusen, seks
:PN	Proper noun	Egholm, Puccini
:PNF	First name	Kristine, Curtis
:PNG	Geograhical name	Tertnes, Sandefjord
:PNH	Last nameg	Høyem, Lundberg
:PNO	Organization	Braathens, Santana
:PREP	Preposition	fra
:PREPDET	Preposition+ determiner	idette, idenne
:PRO	Pronoun	jeg, det
:PROP	Personal pronoun	sjølve
:VB	Verb: Infinitive	trikse
:VE	Verb: Present participle	brukende, krislende
:VH	erb: Past participle	brukt
:VI	Verb: Passive	bemyndiges, fyltes
:VJ	Verb: Preterite	Preterite brukte

:VP	Verb: Present	gasjerer
:VY	Verb: Imperative	slepp
:date	Date	12/23/2012, 23/12/2012
:url	URL	http://www.sas.com
:digit	Number	12, 23, 23.4
:inc	Unknown word	txt, sms
:sep	Punctuation	,.!

#### **Polish**

 Table A1.20
 Part-of-Speech Tags for Polish

Part-of-Speech Tag	Description	Examples
:A	Adjective	własne, każda, głównych
:ABBREV	Abbreviation	ang., tzw.
:ADV	Adverb	więcej, tylko
:CONJ	Conjunction	i, czyli
:INTER	Interjection	ej, fuj, amen
:N	Noun	teorie, miejscach, Wojciech
:NUM	Numeral	siedmiu, tysięcy
:PART	Particle	też
:PREP	:Preposition	za, z, na, do
:PRON	Pronoun	się, sami, go, tobie

**122** Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:V	Verb	wiedzieć, dotarł
:date	Date	:01/01/2012, 12/12/17, 12-23-2001, 23-12-01
:time	Time	23:30:01
:digit	Number	12, -5, 23,45
:sep	Separator or punctuation	.,-
:url	URL	http://www.sas.com
:PN	Unknown or foreign proper noun	Achitophel, Trzciński, LP-vinyl
:inc	Unknown or foreign word	sapiens, ela544

## **Portuguese**

 Table A1.21
 Part-of-Speech Tags for Portuguese

Part-of-Speech Tag	Definition	Examples
:A	Adjective	confiável, belo, feliz
:ADV	Adverb	belamente, felizmente
:CONJ	Conjunction	e, que
:DET	Determiner	alguns, cada, os
:digit	Number	21
:DNUM	Numeric determiner	bilionésimo, cinco
:inc	Unknown word	XIXX
:INTJ	Interjection	caramba, eh

:N	Noun	beleza, felicidade
:PFX	Prefix	anti, circum
:PN	Proper noun	Brasil, Portugal
:PREP	Preposition	com, de, em
:PREPDET	Preposition and determiner contraction	dessas, dum
:PRO	Pronoun	me, nós, quem
:sep	Separator or punctuation	,
:url	URL	http://www.sas.com
:V	Verb	agradecem, garanto
:VB	Infinitive	agradecer, garantir
:VG	Gerund	agradecendo, garantindo
:VH	Past historic	agradecido, garantido
:XL	Foreign word	cf, ibid, sic

## Romanian

 Table A1.22
 Part-of-Speech Tags for Romanian

Part-of-Speech Tag	Description	Examples
:A	Adjective	înalt
:ABBR	Abbreviation	etc
:ADV	Adverb	taman

**124** Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:CONJ	Conjunction	şi
:DEMDET	Demonstrativer determiner	acesta
:DET	Determiner	un
:DOPRO	direct object pronoun	vă
:INTDET	interrogative determiner	câtora
:INTJ	Interjection	hei
:INTPRO	Interrogative pronoun	се
:IOPRO	Indirect object pronoun	îmi
:N	Noun	carte
:NUM	Number	trei
:POLPOSSPRO	Polite possessive pronoun	dumneavoastră
:POSSADJ	Possessive adjective	voastre
:POSSADJ :POSSDET	Possessive adjective Possessive determiner	voastre
	-	
:POSSDET	Possessive determiner	alui
:POSSDET :PREP	Possessive determiner Preposition	alui
:POSSDET :PREP :PRO	Possessive determiner Preposition Iimpersonal pronoun	alui pro vreuna
:POSSDET :PREP :PRO :PRONADJ	Possessive determiner Preposition limpersonal pronoun Pronominal adjective	alui pro vreuna întâia
:POSSDET :PREP :PRO :PRONADJ :ROPRO	Possessive determiner  Preposition  limpersonal pronoun  Pronominal adjective  Reflexive pronoun	alui pro vreuna întâia însele
:POSSDET :PREP :PRO :PRONADJ :ROPRO :SPRO	Possessive determiner  Preposition  limpersonal pronoun  Pronominal adjective  Reflexive pronoun  Personal pronoun	alui  pro  vreuna  întâia  însele  eu
:POSSDET :PREP :PRO :PRONADJ :ROPRO :SPRO :VIMPERAT	Possessive determiner  Preposition  Iimpersonal pronoun  Pronominal adjective  Reflexive pronoun  Personal pronoun  Imperative verb	alui  pro  vreuna  întâia  însele  eu  ziceţi

:VPASTPART	Ppast participle	zis
:VPLUPERFIND	Pluperfect indicative verb	zisesem
:VPRESIND	Present indicative verb	zic
:VPRESPART	Present participle	zicând
:VPRESSUBJ	Present subjunctive verb t	zică
:VPRETIND	Preterite indicative verb	ziseră
:inc	Unknown word or nonword	aasdqwert
:PN	Proper noun	Elena
:time	Time	23:30:00, 8:45
:date	Date	12-23-2012, 12/12/2012, 01.12.2012
:digit	Number	100, 999
:url	URL	http://www.sas.com
:sep	Punctuation or separator	.,!

## Russian

Table A1.23 Part-of-Speech Tags for Russian

Part-of-Speech Tag	Definition	Examples
:A	Adjective	духовитый, красивая, лучших
:ABBREV	Abbreviation	др, км
:adv	Comparative adverb	дальше

Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:adverbial	Adverb	хорошо, сколько-нибудь
:conj	Conjunction	если, и
:digit	Number	123, 12.3, 12.3.2003, 12/3/2013
:inc	Unknown word	геминг, analytics
:idet	Interrogative determiner	который
:INT	Interjection	ax
:intadv	Interrogative adverb	где
:intquant	Interrogative quantifier	сколькие, почём
:N	Noun	велосипед, история, малолетство
:NONDECL	Nondeclinable word	мартини, маэстро
:NONDECL-ADJ	Nondeclinable adjective	баскервиллей
:NONDECL-PN	Nondeclinable proper noun	Шевроле, Айдахо
:NONDECL-PRO	Nondeclinable pronoun	всяко
:num	Number	один, десятью
:particle	Particle	бы, же
:PN	Proper noun	Миа, Тузла
:PNA	Proper adjective	Роханский, Сашина
:PNN	Proper noun, name	Свердловск, Мария, Давыдович
:prep	Preposition	до, вроде
:pron	Pronoun	я, её

:sep	Separator or punctuation	,.!
:url	URL	http://www.sas.com
:VB	Infinitive	автоматизировать, менять, кончить
:V	Verb	нажимает, кладите, плавала
:VG	Gerund	адаптировав, вальсируя

## Slovak

 Table A1.24
 Part-of-Speech Tags for Slovak

Description	Examples
Adjective	všeobecné , verejnej
Abbreviatione	ul, Dr
Adverb	pravidelne, vyslovene vakrail.sk 23:30:00 23/12/2012, 23-12-2012
Conjunction	ak , iba
Interjection	oj, stop
Noun	doručení, partnerov
Numeral	štyritisíc, prvom
Particle	by, tiež wslettri http:// www.sas.com, info@slo
Possessive adjective	vaše, jeho
Preposition	o, v, pre
	Adjective Abbreviatione Adverb  Conjunction Interjection Noun Numeral Particle  Possessive adjective

Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:PREPPRO	Preposition+Pronoun	uňho
:PRO	Pronoun	si, Vám, to
:PROR	Relative pronoun	ktoré, akékoľvek
:V	Verb	prinášame, budú
:VN	Negative verb	nespráva, neviete
:VB	Infinitive verb	využívať, stiahnuť
:VBN	Negative infinitive verbe	nezaostávať e
:VG	Past participle verbim	prešli, chutili
:VGN	Negative past participle T	nemali
:VY	Imperative I	zažite, pozrite
:VYN	Negative imperative	nevybrali
:digit	Numbermai	1.4, -10, +421
:sep	Separator or punctuation	.,1
:PN	Proper noun e	Oetker, KEP
:inc	Unknown or foreign word	newslettri
:url	URL or email	http://www.sas.com, info@slovakrail.sk
:time	Time	23:30:00
:date	Date	23/12/2012, 23-12-2012

## Slovene

Part-of-Speech Tag	Description	Example
:Adj	Adjective	prvi, črna
:Abbr	Abbreviation	itd.
:Aux	Auxiliary verb	je (in front of a main verb)
:Adv	Adverb	hmalu, daleč
:Conj	Conjunction	ali, in
:Interj	Interjection	bravo, ah
:Noun	Noun	dni, dogodka
:Num	Numeral	dva, šest
:digit	Number	20.3, 123
:Part	Particle	pa, ne, spet
:Prep	Preposition	v, za
:Pron	Pronoun	te, mi, vsak, kdo
:Verb	Verb	sta, uporablja, suspendirali, pozabite
:sep	Separator or punctuation	.∶, «
:Prop	Proper noun	Maribor, Roglič
:date	Date	23/12/2012, 23-12-2012
:time	Time	23:30:00

:url	URL	http://www.sas.com, info@sas.com

# **Spanish**

Table A1.25 Part-of-Speech Tags for Spanish

Part-of-Speech Tag	Definition	Examples	
:A	Adjective	confiable, feliz, hermoso	
:ABBREV	Abbreviation	km, pág, Sra	
:Adv	Adverb	ahora, felizmente	
:CONJ	Conjunction	ni, pero, y	
:DET	Determiner	el, las, mi, nuestro	
:digit	Number	21	
:inc	Unknown word	XIXX	
:INTJ	Interjection	hola	
:N	Noun	belleza, felicidad	
:PN	Proper noun	Chile, España	
:PREP	Preposition	con, de, en, por	
:PREPDET	Preposition and determiner contraction	al, del	
:PRON	Pronoun	alguien, ellos, me	
:sep	Separator or punctuation	,	
:url	URL	http://www.sas.com	

:V	Verb	ayudan, pide	
:VB	Infinitive	ayudar, pedir	
:VE	Present progressive	ayudando, pidiendo	
:VH	Past participle	ayudado, pedido	

#### **Swedish**

 Table A1.26
 Part-of-Speech Tags for Swedish

Part-of-Speech Tag	Definition	Examples	
:A	Adjective	fört	
:ABB	Abbreviation	st.	
:ADV	Adverb	väl	
:CONJ	Conjunction	samt	
:DET	Determiner	ens	
:DNUM	Number	två	
:INT	Interjection	hej	
:INV	Invariant	morse	
:N	Noun	bok	
:PN	Proper noun	Øsel	
:PNF	Proper noun - first name	Tove	
:PNG	Proper noun - geographic	Östmark	
:PNH	Proper noun - last name	Viklund	

:PNO	Proper noun - organization	Toshiba
:PREP	Preposition	till
:PRO	Impersonal pronoun	somlig
:PROP	Personal pronoun	du
:VB	Infinitive verb	vara
:VD	Active supine verb	varit
:VE	Present participle	varande
:VF	Passive supine verb	varats
:VH	Perfect participle	sedd
:VI	Passive present	ses
:VJ	Active preterite	såg
:VK	Passive preterite	sågs
:VP	Active present verb	varar
:VY	Imperative	vara

## Thai

 Table A1.27
 Part-of-Speech Tags for Thai

Part-of-Speech Tag	Description	Examples	
:ADJ	Adjective	กตัญญู, กตัญญูกตเวที	
:ADV	Adverb	กระง่องกระแง่ง, กระดิบๆ	
:AUXVERB	Auxiliary verbs	ควรจะ, ต้อง	

:CLAS	Classifiers	กก., กม.	
:CONJ	Conjunction	ก่อน, จน	
:DET	Determiner	ทั้ง, ทุก	
:END	Particle used at the end of a question, command or entreaty	ล่ะ, เหรอ	
:INTERJ	Interjection	ชะชะ, ดูกร	
:NEG	Negation	มิใช่, ไม่	
:NOUN	Noun	กงพัด, กฎหมายบ้านเมือง	
:NUMBER	Number	สอง, เก้า	
:PREF	Prefix	ปรา, อน	
:PREP	Preposition	กว่า, ก่อนหน้า	
:PRON	Pronoun	คนอื่นๆ, คนใด	
:PROPLOC	Proper noun, location	กมลา, กรีซ	
:PROPMISC	Proper noun, others	กุชชี่, คลีนิกซ์	
:PROPNAME	Proper noun, person names	กปิลกาญจน์, กตัญญุตานนท์	
:PROPORG	Proper noun, organizations	กรุงเทพธุรกิจ, กระทรวงมหหาด ไทย	
:PUNC	Separator or punctuation	"()	
:SUFF	Suffix s	สิ, เอย	
:VERB	Verb	กทรรป, กรมเกรียม	
:DEFAULT	Unknown wordz	Josephson, microbridge	

#### **Turkish**

Table A1.28 Part-of-Speech Tags for Turkish

Part-of-Speech Tag	Description	Examples
:A	Adjective	iyi, zor
:ADV	Adverb	yine, zaten
:CONJ	Conjunction	veya, hem
:date	Date	12/30/2000, 12/30/00, 2000-30-12
:digit	Number	12.302.000, 5
:inc	Unknown word	wug
:N	Noun	kitap, insan
:NUMERAL	Numeral	dokuz, onbir
:PARTICLE	Particle	beri, diye
:PN	Proper noun	Ayşe, Türkçe
:PRONOUN	Impersonal pronoun	bunlar, kendi
:PROP	Personal pronoun	onlar, sen
:QUANT	Quantifier	çok, her
:sep	Separator or punctuation	!.,
:time	Time	12:30:00
:url	URL	sas.com, www.sas.com, http:// www.sas.com

:V_GER_STEM	Untagged verb	Not applicable	
	Y (imperative)	bil	
	O (subjunctive)	bile	
	P (continuous, including archaic forms)	biliyor, bilmekte	
	N (necessitative)	bilmeli	
	J (past)	bildi	
	I (indefinite/inferential)	bilmiş	
	F (future)	bilecek	
	C (conditional)	bilse	
	B (infinitive)	bilmek	
	A (habitual/nomic)	bilir	
:V	Verb (can be followed by any number and combination of any of the following):		

### **Vietnamese**

 Table A1.29
 Part-of-Speech Tags for Vietnamese

Part-of-Speech Tag	Description	Examples
:A	Adjective	an toàn, bận rộn, lịch sự
:ABBREV	Abbreviation	APEC, ANÐT, ÐTNN
:Adv	Adverb	bỗng chốc, chưa chừng
:Aux	Particle	chính
:C	Conjugation	dù rằng, hoặc là
:F	Foreign word	cà-rem, Ampe, ăng ten
:Int	Interjection	hỡi, ái chà, ô hay
:N	Noun	áo quần, cừu, cương vị

Appendix 1 / Part-of-Speech Tags (for Languages Other Than English)

:Num	Numeral	2007, bảy, mươi n
:PreDet	Determiner	một số
:Prep	Preposition	cho, vào
:PN	Proper noun	Việt Nam, Trung Quốc
:Pro	Pronoun	tôi, chúng mày, chúng nó
:PUNCT	Punctuation or symbol	!:()@
:RelPro	Relative pronoun	ai nấy
:V	Verb	ngưỡng mộ, lưu nghiệm
:DEFAULT	Unrecognized character	,

## **Appendix 2**

# Predefined Concept Priorities (for Languages Other Than English)

Using Priority Values in Predefined Concepts	138
Priority Values for Predefined Concepts	138
Arabic	138
Chinese	139
Croatian	140
Czech	141
Danish	142
Dutch	143
Farsi	144
Finnish	144
French	145
German	
Greek	147
Hebrew	148
Hindi	148
Hungarian	149
Indonesian	149
Italian	150
Japanese	151
Korean	152
Norwegian	152
Polish	153

Portuguese	154
Romanian	155
Russian	155
Slovak	156
Slovene	
Spanish	157
Swedish	158
Thai	159
Turkish	160
Vietnamese	160

# Using Priority Values in Predefined Concepts

To accurately set priorities for matching custom concepts in your language, see the topic "Priority Values for Predefined Concepts". For information about setting priorities, see "Concepts Page" on page 36. For priority values in English, see Table 3.1 on page 31.

**Note:** Use the highest priority value per language to ensure that there are no conflicts with custom concepts during document processing. The highest priority value for each language is marked in the tables in the following section with a footnote.

### **Priority Values for Predefined Concepts**

### **Arabic**

For Arabic, there are no specific priority values for predefined concepts. The default value of 10 is used for all of the predefined concepts listed below.

Predefined Concept
CURRENCY
DATE
INTERNET
LOCATION
MEASURE
NOUN_GROUP
NUMBER
ORGANIZATION
PERCENT
PERSON
PHONE
TIME
TIME_PERIOD
TITLE

### Chinese

 Table A2.2
 Predefined Concept Priorities for Chinese

Predefined Concept	Priority Value
ADDRESS	30*

**140** Appendix 2 / Predefined Concept Priorities (for Languages Other Than English)

CURRENCY	20
DATE	30*
INTERNET	20
LOCATION	20
MEASURE	20
ORGANIZATION	25
PERCENT	20
PERSON	20
PHONE	20
PROP_MISC	20
TIME	20
TITLE	20

<sup>\*</sup> Highest value for this language.

### Croatian

 Table A2.3
 Predefined Concept Priorities for Croatian

Predefined Concept	Priority Value
LOCATION	12*
ORGANIZATION	10
CURRENCY	10
PERSON	11

PHONE	10
TITLE	10
MEASURE	10
NOUN_GROUP	10
DATE	10
TIME	10
PERCENT	10

Highest value for this language.

### Czech

 Table A2.4
 Predefined Concept Priorities for Czech

Predefined Concept	Priority Value
NOUN_GROUP	9
PERSON	10*
ORGANIZATION	10*
DATE	10*
TIME	10*
PERCENT	10*
CURRENCY	10*
LOCATION	10*

Highest value for this language.

### **Danish**

 Table A2.5
 Predefined Concept Priorities for Danish

Predefined Concept	Priority Value
PERSON*	20*
ORGANIZATION*	20*
LOCATION*	20*
COMPANY	19
TITLE	18
PHONE	18
DATE	18
TIME	18
INTERNET	18
MEASURE	18
NOUN_GROUP	15
PERCENT	18
SSN	18
CURRENCY	18
TIME_PERIOD	18
PROP_MISC	13
VEHICLE	15

ADDRESS	20*

Highest value for this language.

### **Dutch**

 Table A2.6
 Predefined Concept Priorities for Dutch

Predefined Concept	Priority Value
ADDRESS	20*
COMPANY	19
CURRENCY	18
DATE	18
INTERNET	18
LOCATION	20*
MEASURE	18
NOUN_GROUP	15
ORGANIZATION	20*
PERCENT	18
PERSON	20*
PHONE	18
PROP_MISC	13
TIME	18
TIME_PERIOD	18

Highest value for this language.

### **Farsi**

For Farsi, there are no specific priority values for predefined concepts. The default value of 10 is used for all of the predefined concepts listed below.

 Table A2.7
 Predefined Concept Priorities for Farsi

Predefined Concept	
CURRENCY	
PERCENT	
DATE	
TIME	
LOCATION	
PERSON	
ORGANIZATION	

### **Finnish**

 Table A2.8
 Predefined Concept Priorities for Finnish

Predefined Concept	Priority Value
NOUN_GROUP	15
LOCATION	25*
PERSON	20
ORGANIZATION	10

COMPANY	25*
DATE	10
TIME	10
CURRENCY	10
INTERNET	10

Highest value for this language.

### **French**

 Table A2.9
 Predefined Concept Priorities for French

Predefined Concept	Priority Value
ADDRESS	20*
COMPANY	19
CURRENCY	18
DATE	18
INTERNET	18
LOCATION	20*
MEASURE	18
NOUN_GROUP	15
ORGANIZATION	20*
PERCENT	18
PERSON	20*

**146** Appendix 2 / Predefined Concept Priorities (for Languages Other Than English)

PHONE	18
PROP_MISC	13
SSN	18
TIME	18
TIME_PERIOD	18
TITLE	18

<sup>\*</sup> Highest value for this language.

### German

 Table A2.10
 Predefined Concept Priorities for German

Priority Value
25
25
25
18
18
40
18
15
25
18

PERSON	60*
PHONE	18
PROP_MISC	30
TIME	18
TIME_PERIOD	18

Highest value for this language.

### Greek

 Table A2.11
 Predefined Concept Priorities for Greek

Predefined Concept	Priority Value
DATE	18
LOCATION	25*
MEASURE	18
NOUN_GROUP	15
COMPANY	25*
PERCENT	18
PERSON	20
PHONE	18
TIME	18
INTERNET	18
CURRENCY	18

**148** Appendix 2 / Predefined Concept Priorities (for Languages Other Than English)

TIME_PERIOD	18
ADDRESS	25*
ORGANIZATION	20

<sup>\*</sup> Highest value for this language.

### Hebrew

For Hebrew, there are no specific priority values for predefined concepts. The default value of 10 is used for all of the predefined concepts listed below.

 Table A2.12
 Predefined Concept Priorities for Hebrew

Predefined Concept	
NOUN_GROUP	

### Hindi

 Table A2.13
 Predefined Concept Priorities for Hindi

Predefined Concept	Priority Value
NOUN_GROUP	10
CURRENCY	10
DATE	10
LOCATION	10
ORGANIZATION	10
PERCENT	10

PERSON	40*
TIME	10

Highest value for this language.

### Hungarian

 Table A2.14
 Predefined Concept Priorities for Hungarian

Predefined Concept	Priority Value
PERSON	20*
DATE	15
TIME	15
LOCATION	15
ADDRESS	10
CURRENCY	15
PERCENT	15
ORGANIZATION	15
NOUN_GROUP	10

Highest value for this language.

### **Indonesian**

For Indonesian, there are no specific priority values for predefined concepts. The default value of 10 is used for all of the predefined concepts listed below.

 Table A2.15
 Predefined Concept Priorities for Indonesian

Predefined Concept

NOUN\_GROUP

COMPANY

### Italian

 Table A2.16
 Predefined Concept Priorities for Italian

Predefined Concept	Priority Value
ADDRESS	25*
COMPANY	25*
CURRENCY	18
DATE	18
INTERNET	18
LOCATION	25*
MEASURE	18
NOUN_GROUP	15
PERCENT	18
PERSON	20
PHONE	18
PROP_MISC	13
SSN	18

TIME	18
TIME_PERIOD	18
TITLE	18

Highest value for this language.

### **Japanese**

For Japanese, there are no specific priority values for predefined concepts. The default value of 10 is used for all of the predefined concepts listed below.

 Table A2.17
 Predefined Concept Priorities for Japanese

Predefined Concept	
ADDRESS	
CURRENCY	
DATE	
LOCATION	
MEASURE	
ORGANIZATION	
PERCENT	
PERSON	
PHONE	
TIME	

### Korean

 Table A2.18
 Predefined Concept Priorities for Korean

Predefined Concept	Priority Value
CURRENCY	19
DATE	20
INTERNET	18
LOCATION	22*
MEASURE	19
NUMBER	18
ORDNUMBER	18
ORGANIZATION	21
PERCENT	18
PERSON	20
PHONE	18
TIME	20
TITLE	18

<sup>\*</sup> Highest value for this language.

### Norwegian

For Norwegian, there are no specific priority values for predefined concepts. The default value of 10 is used for all of the predefined concepts listed below.

 Table A2.19
 Predefined Concept Priorities for Norwegian

**Predefined Concept** 

NOUN\_GROUP

### **Polish**

 Table A2.20
 Predefined Concept Priorities for Polish

Predefined Concept	Priority Value
INTERNET	18
ADDRESS	20
PROP_MISC	13
CURRENCY	18
DATE	18
LOCATION	20
MEASURE	18
NOUN_GROUP	15
ORGANIZATION	21*
COMPANY	19
PERCENT	18
PERSON	20
PHONE	18
TIME	18

**154** Appendix 2 / Predefined Concept Priorities (for Languages Other Than English)

TIME_PERIOD	18
TITLE	18

<sup>\*</sup> Highest value for this language.

### **Portuguese**

 Table A2.21
 Predefined Concept Priorities for Portuguese

Predefined Concept	Priority Value
ADDRESS	25*
COMPANY	25*
CURRENCY	18
DATE	18
INTERNET	18
LOCATION	25*
MEASURE	18
NOUN_GROUP	15
PERCENT	18
PERSON	20
PHONE	18
PROP_MISC	13
SSN	18
TIME	18

TIME_PERIOD	18
TITLE	18

Highest value for this language.

### Romanian

For Romanian, there are no specific priority values for predefined concepts. The default value of 10 is used for all of the predefined concepts listed below.

 Table A2.22
 Predefined Concept Priorities for Romanian

Predefined Concept	
NOUN_GROUP	
COMPANY	

### Russian

For Russian, there are no specific priority values for predefined concepts. The default value of 10 is used.

 Table A2.23
 Predefined Concept Priorities for Russian

Predefined Concept	
NOUN_GROUP	
LOCATION	
ORGANIZATION	
PERSON	
CURRENCY	

**156** Appendix 2 / Predefined Concept Priorities (for Languages Other Than English)

DATE	
VEHICLE	
TIME	
PERCENT	
INTERNET	

### Slovak

 Table A2.24
 Predefined Concept Priorities for Slovak

Predefined Concept	Priority Value
LOCATION	8
PERSON	7
ORGANIZATION	10*
INTERNET	10*
ADDRESS	10*
NOUN_GROUP	6
CURRENCY	10*
DATE	10*
TIME	10*
PERCENT	10*

<sup>\*</sup> Highest value for this language.

### Slovene

For Slovene, there are no specific priority values for predefined concepts. The default value of 10 is used.

 Table A2.25
 Predefined Concept Priorities for Slovene

Predefined Concept
LOCATION
ORGANIZATION
COMPANY
PERSON
MEASURE
PRODUCT
PROP_MISC
VEHICLE
NOUN_GROUP

### **Spanish**

 Table A2.26
 Predefined Concept Priorities for Spanish

Predefined Concept	Priority Value
ADDRESS	25*
COMPANY	25*

**158** Appendix 2 / Predefined Concept Priorities (for Languages Other Than English)

CURRENCY	18
DATE	18
INTERNET	18
LOCATION	25*
MEASURE	18
NOUN_GROUP	15
ORGANIZATION	20
PERCENT	18
PERSON	20
PHONE	18
PROP_MISC	13
SSN	18
TIME	18
TIME_PERIOD	18
TITLE	18

<sup>\*</sup> Highest value for this language.

### **Swedish**

 Table A2.27
 Predefined Concept Priorities for Swedish

Predefined Concept	Priority Value
ADDRESS	20*

COMPANY	19
CURRENCY	18
DATE	18
INTERNET	18
LOCATION	20*
MEASURE	18
NOUN_GROUP	15
ORGANIZATION	20*
PERCENT	18
PERSON	20*
PHONE	19
PROP_MISC	13
SSN	18
TIME	18
TIME_PERIOD	18
VEHICLE	15

Highest value for this language.

### Thai

 Table A2.28
 Predefined Concept Priorities for Thai

	Predefined Concept	Predefined Concept
--	--------------------	--------------------

**160** Appendix 2 / Predefined Concept Priorities (for Languages Other Than English)

PERSON	20
ORGANIZATION	30*
LOCATION	30*

<sup>\*</sup> Highest value for this language.

### **Turkish**

 Table A2.29
 Predefined Concept Priorities for Turkish

Predefined Concept	Predefined Concept
NOUN_GROUP	10
ORGANIZATION	11*
COMPANY	10
PERSON	10
LOCATION	10
PERCENT	10
CURRENCY	10
DATE	10
TIME	10

<sup>\*</sup> Highest value for this language.

### **Vietnamese**

For Vietnamese, there are no specific priority values for predefined concepts. The default value of 10 is used.

 Table A2.30
 Predefined Concept Priorities for Turkish

Predefined Concept
DATE
INTERNET
PROP_MISC
TIME
TIME_PERIOD

### **Recommended Reading**

Here is the recommended reading list for this title:

- SAS Contextual Analysis: Administrator's Guide
- SAS Content Categorization Single User Servers: Administrator's Guide
- SAS Text Miner: Reference Help
- SAS Encoding: Understanding the Details
- SAS Enterprise Content Categorization Studio: User's Guide
- Text Mining and Analysis: Practical Methods, Examples, and Case Studies Using SAS

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

SAS Books

SAS Campus Drive

Cary, NC 27513-2414

Phone: 1-800-727-0025 Fax: 1-919-677-4444

Email: sasbook@sas.com

Web address: sas.com/store/books

### **Glossary**

### category

a classification for documents that is based on a common characteristic. Category membership is indicated as a binary property. In order to determine when a document is likely to be a member of a category, one or more Boolean rules comprising the category text definition must be satisfied.

### concept

an abstract class of meanings. In order to determine when a concept is likely to be referenced in a subset of text, the rules comprising the concept text definition must be satisfied.

### model scoring

the process of applying a model to new data in order to compute outputs.

### parse

to analyze text, such as a SAS statement, for the purpose of separating it into its constituent words, phrases, punctuation marks, values, or other types of information. The information can then be analyzed according to a definition or set of rules.

### relevancy score

a score that indicates how well a document satisfies a rule or model. The best match has a score of 1 and reflects a perfect (100%) match.

### scoring

See model scoring.

#### sentiment

an attitude that is expressed about an item that is being analyzed, which can be a segment of text, a grouping of text segments, or a specific subject of interest.

### sentiment analysis

the use of natural language processing, computational linguistics, and text analytics to determine the attitude of a speaker or writer with respect to a topic, document, or other item of analysis. Sentiment analysis results in a positive, negative, or neutral score on the target of analysis.

### 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.

### string

See text string.

#### subset of text

the matched text for a concept text definition; this consists of one or more strings that are contained in a document.

### surface form

a variant of a term that is contained in a matched subset of text in one or more documents. These forms include stems, synonyms, misspellings, and alternate ways of referring to the same entity.

### taxonomy

a hierarchical relationship of parent and child category nodes. In a true taxonomy, whenever a category is detected, it is implied that all parents are also represented. For example, if something is identified as human, it must also be a primate, mammal, animal, and so on.

#### term

a representation of a single concept in one or more textual forms, as defined by rules or algorithms.

### term map

a node-arc graph that centers around an "object of interest," which could be a category, concept, topic, or term. Corresponding nodes in the graph indicate rules that are predictive of the object of interest. Better rules are shown as larger nodes. The arcs represent the addition or exclusion of terms that are used to build up the rules.

#### term role

a function that is performed by a term in a particular context. A term can function as a part of speech, entity type, or other purpose that is user-defined.

#### term table

a list of every term in a collection of documents including the representative text form for each term, its role, and all of its surface forms that appear within that collection.

### text string

a subset of text that consists of adjacent characters of any type. Depending on the specified options, strings can be either case-sensitive or case-insensitive.

#### token

in the SAS programming language, a collection of characters that communicates a meaning to SAS and that cannot be divided into smaller functional units. A token such as a variable name might look like an English word, but can also be a mathematical operator, or even an individual character such as a semicolon. A token can contain a maximum of 32,767 characters.

### topic

a machine-generated category, the purpose of which is to indicate what documents are about. A topic identifies groupings of important terms in a document collection. A single document can contain one or more topics, or no topics.

### topic document weight

See topic-specific document weight

### topic term weight

See topic-specific term weight

### topic-specific document weight

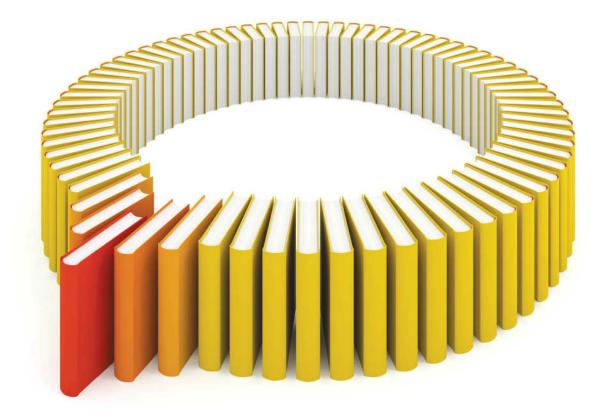
an indicator of the importance of a topic to a document. A value that is above a specified cutoff value indicates that a document contains that topic.

### topic-specific term weight

an indicator of the relative importance of a term in a topic as compared to other terms. A term with a value above a specified cutoff value contributes to the assignment of a document to the topic.

### weight

a numeric indicator that is assigned to an item and that indicates the relative importance of the item in a frequency distribution or population.



# Gain Greater Insight into Your SAS® Software with SAS Books.

Discover all that you need on your journey to knowledge and empowerment.



