# TEXT MINING WITH RAPIDMINER

*G. Ertek, D. Tapucu, and I. Arın*
*Sabancı University, Istanbul, Turkey*

The goal of this chapter is to introduce the text mining capabilities of RAPIDMINER through a use case. The use case involves mining reviews for hotels at TripAdvisor.com, a popular web portal. We will be demonstrating basic text mining in RAPIDMINER using the text mining extension. We will present two different RAPIDMINER processes, namely Process01 andProcess02, which respectively describe how text mining can be combined with association mining and cluster modeling. While it is possible to construct each of these processes from scratch by inserting the appropriate operators into the process view, we will instead import these two processes readily from existing model files. Throughout the chapter, we will at times deliberately instruct the reader to take erroneous steps that result in undesired outcomes. We believe that this is a very realistic way of learning to use RAPIDMINER, since in practice, the modeling process frequently involves such steps that are later corrected.

# USE CASES WITH RAPIDMINER

# USE CASES WITH RAPIDMINER
## Working Title

**Dr. Markus Hofmann**
Institute of Technology Blanchardstown, Ireland

**Ralf Klinkenberg**
Rapid-i

*to be added*

# CONTRIBUTORS

GURDAL ERTEK,    Sabancı University, Istanbul, Turkey
DILEK TAPUCU,    Sabancı University, Istanbul, Turkey
INANC ARIN,    Sabancı University, Istanbul, Turkey

# CONTENTS IN BRIEF

# CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

**CHAPTER 1**

# TEXT MINING WITH RAPIDMINER

G. Ertek, D. Tapucu, and I. Arin

Sabancı University, Istanbul, Turkey

## 1.1 INTRODUCTION

The goal of this chapter is to introduce the text mining capabilities of RAPIDMINER through a use case. The use case involves mining reviews for hotels at TripAdvisor.com, a popular web portal. We will be demonstrating basic text mining in RAPIDMINER using the text mining extension. We will present two different RAPIDMINER processes, namely Process01 and Process02, which respectively describe how text mining can be combined with association mining and cluster modeling. While it is possible to construct each of these processes from scratch by inserting the appropriate operators into the process view, we will instead import these two processes readily from existing model files.

Throughout the chapter, we will at times deliberately instruct the reader to take erroneous steps that result in undesired outcomes. We believe that this is a very realistic way of learning to use RAPIDMINER, since in practice, the modeling process frequently involves such steps that are later corrected.

**1**

### 1.1.1   Text Mining

*Text mining* (also referred to as *text data mining* or *knowledge discovery from textual databases*), refers to the process of discovering interesting and non-trivial knowledge from text documents. The common practice in text mining is the analysis of the information extracted through text processing to form new facts and new hypotheses, that can be explored further with other data mining algorithms. Text mining applications typically deal with large and complex data sets of textual documents that contain significant amount of irrelevant and noisy information. Feature selection aims to remove this irrelevant and noisy information by focusing only on relevant and informative data for use in text mining. Some of the topics within text mining include feature extraction, text categorization, clustering, trends analysis, association mining and visualization.

### 1.1.2   Data Description

The files required for this chapter, including the data and pre-built processes reside within a folder titled `LocalRepository`. The data used in this chapter comes from TripAdvisor.com, a popular web portal in the hospitality industry, and is shown in Figure 1.1. This publicly available data set contains the reviews and ratings (1 through 5) of clients or customers for 1850 hotels. The original data was extracted by The Database and Systems Information Laboratory at the University of Illinois at Urbana-Champaign,and is available under `http://sifaka.cs.uiuc.edu/~wang296/Data/`. There are 1850 text documents in the original data set,corresponding to the reviews of 1850 hotels. Each document contains all the reviews for that hotel. While it is possible to run text mining processes with the original data, we will be using a subset of the data containing only the first 100 text documents. The data set used in this chapter may be downloaded from `http://people.sabanciuniv.edu/ertekg/papers/supp/09.zip` or the short web link `http://bit.ly/Kzdn5x`, as well as http://research.sabanciuniv.edu .

### 1.1.3   Running RAPIDMINER

When running RAPIDMINER, it is strongly recommended to right click the mouse button on the start menu and `Run as administrator`, as shown in Figure 1.2, rather than simply clicking or double clicking the executable or shortcut icon. By running as administrator, we are granted the permissions to install extension packages. Running the software without the administrator rights may cause errors when trying to install updates, including the extension packages. Figure 1.3 shows the splash screen displayed when RAPIDMINER is run for the first time or when it is run without any extensions installed. In Figure 1.3, the highlighted region is currently blank but it will later contain the icons for the installed extensions. When RAPIDMINER is loaded, the user starts with the *welcome perspective* (a *perspective* is a particular combination of *views*), as shown in Figure 1.4.

### 1.1.4   RapidMiner Text Processing Extension Package

RAPIDMINER is the most popular open source software in the world for data mining, and strongly supports text mining and other data mining techniques that are applied in combination with text mining. The power and flexibility of RAPIDMINER is due to the GUI-based IDE (integrated development environment) it provides for rapid prototyping

**Figure 1.1**   The TripAdvisor data set



**Figure 1.2**   Running RAPIDMINER as administrator

and development of data mining models, as well as its strong support for scripting based on XML (extensible mark-up language). The visual modeling in the RAPIDMINER IDE is based on the defining of the data mining process in terms of operators and the flow of

**Figure 1.3**    RAPIDMINER splash screen when no extension packages are installed

process through these operators. Users specify the expected inputs, the delivered outputs, the mandatory and optional parameters, and the core functionalities of the operators, and the complete process is automatically executed by RAPIDMINER. Many packages are available for RAPIDMINER, such as text processing, Weka extension, parallel processing, web mining, reporting extension, series processing, PMML, community, and R extension packages. The package that is needed and used for text mining is the *Text Processing* package, which can be installed and updated through the Update RapidMiner menu item under the Help menu.

### 1.1.5   Installing Text Mining Extensions

We will initiate our text mining analysis by importing the two previously built processes. However, even before that we have to check and make sure that the extensions required for text mining are installed within RAPIDMINER. To manage the extensions, select the Help menu and then  manage Extensions menu item, as shown in Figure 1.4. The dialog box that comes up, as shown in Figure 1.5, does not list any extensions. Thus, the extensions required have to be installed.



**Figure 1.4**    Managing RAPIDMINER extension packages

**Figure 1.5**    Dialog box stating that no extension packages are yet installed



**Figure 1.6**    Updating RAPIDMINER



**Figure 1.7**    Installing extension packages for text mining

Click the Close button and select Help menu and then Update RapidMiner menu item as shown in Figure 1.6. RAPIDMINER will connect to the internet and fetch the list of available updates, eventually displaying all the available updates, as in Figure 1.7. In this window, the selection of an update for installation can be made only by double clicking the check box on the left hand side of that update's name. When an update is selected for installation, a small green check sign appears on the check box, as in Figure 1.7. The extension packages (available as updates) needed for text mining are Text Processing, Web Mining, and Wordnet Extension. Select these updates as shown in Figure 1.7. Then, click the install button. When the Terms of Use window appears, click the radio button for I accept the terms of use, and click Ok . During the installation of updates, the installation process can be stopped by first clicking inside the Progress window, and clicking Stop , but don't do this now. When the Update Complete dialog box appears, click Yes to restart RAPIDMINER with the newly installed updates available . When the splash screen for RAPIDMINER is displayed, you will notice that the icons for the installed extensions appear in the previously blank region (Figure 1.8).



**Figure 1.8**    RAPIDMINER splash screen with the three extension packages installed for text mining

## 1.2   ASSOCIATION MINING OF TEXT DOCUMENT COLLECTION (Process01)

### 1.2.1   Importing Process01

We will now initiate the text mining analysis by importing the processes supplied with this chapter. For this, select Files and then the Import Process menu item as in Figure 1.9. Go to the LocalRepository folder (supplied with this chapter), then to Processes folder, click on Process01.rmp, and click the Open button, as in Figure 1.10. The RAPIDMINER process Process01 is now displayed in the *design perspective*, as shown in Figure 1.11.



**Figure 1.9**   Importing an existing process

**Figure 1.10**    Selecting the process to import



**Figure 1.11**    Operators for Process01 and the Parameters for Process Documents from Files operator
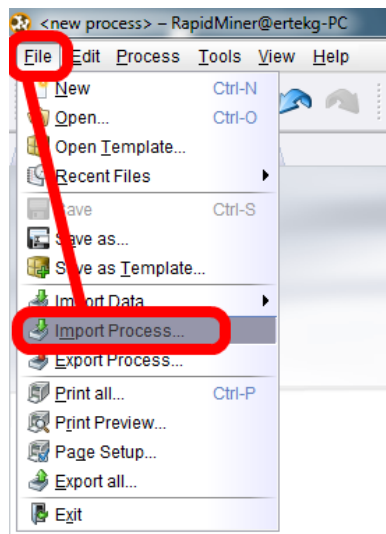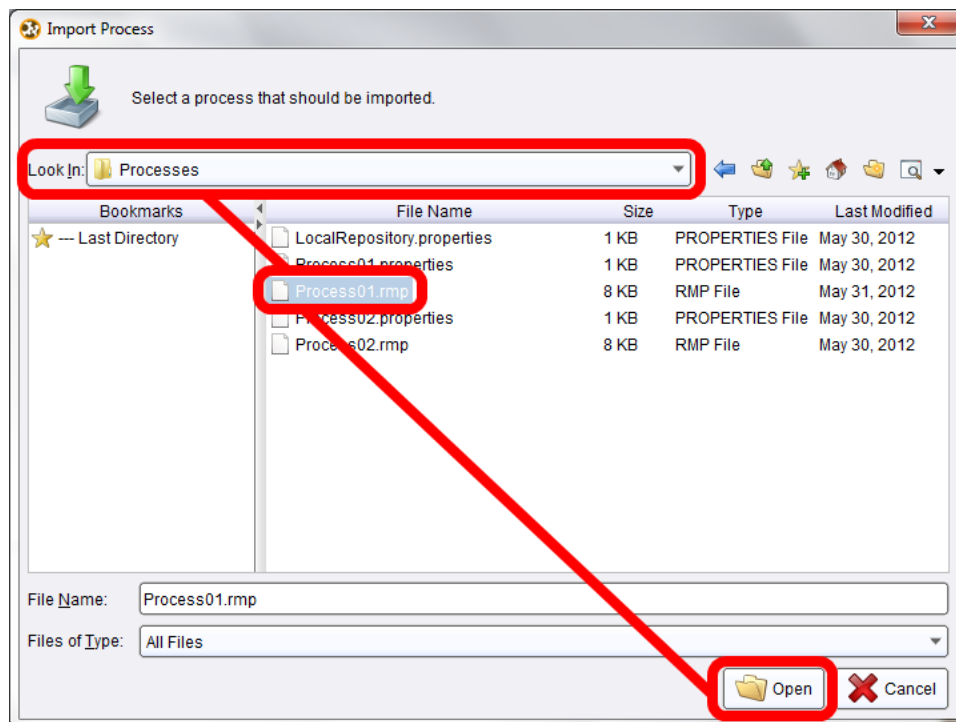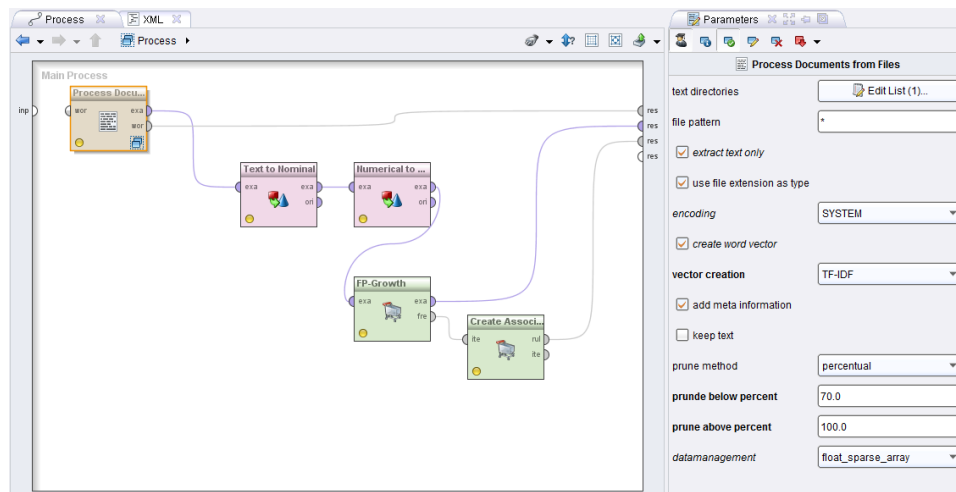
### 1.2.2  Operators in Process01

The parameters for the operators in this process are given on the right-hand side of the process, listed under the Parameters tab. For the first process, the parameter text directories specifies where to read the text data from. A very important parameter is the vector creation used. In Process01, the selected vector creation method is TF-IDF *TF-IDF* (*Term Frequency-Inverse Document Frequency*) is a term weighting method. It gives higher weight to terms that appear frequently in the document but not many times in other documents. However, this may yield too many words, even up to tens of thousands of them. Too many words would make it prohibitive to carry out the successive data mining steps due to lengthy running times for the data mining algorithms. Hence, it is a very good idea to prune the resulting word set using a prune method, by selecting the method and its parameters within the Parameters view.

Click on the Process Documents from Files operator, as in Figure 1.11. In Process01, words that appear in less than 70.0 % of the documents are pruned, as can be seen from the value of 70.0 for the prune below percent parameter. It is also possible to prune the words that appear in too many documents, but this was not done in this example, as can be seen from the value of 100.0 for the prune above percent parameter. In association mining, we are interested in the items (words in text mining context) that appear in 100.0% of the transactions (documents in our text mining context), since they can form interesting frequent itemsets (word lists) and association rules, that provide actionable insights. Thus, it is appropriate to set prune above percent to 100.0 in Process01 (Figure 1.11), including the items (words) that appear in every document.

Process01 consists of five operators (Figure 1.11). Firstly, the Process Documents from Files operator performs text processing which involves preparing the text data for the application of conventional data mining techniques. Process Documents from Files operator reads data from a collection of text files and manipulates this data using text processing algorithms. This is a nested operator, meaning that it can contain a sub-process consisting of a multitude of operators. Indeed, in Process01, this nested operator contains other operators inside. Double click on this operator, and you will see the sub-process inside it, as outlined in Figure 1.12. This sub-process consists of six operators that are serially linked (Figure 1.12):

- Tokenize Non-letters (Tokenize)

- Tokenize Linguistic (Tokenize)

- Filter Stopwords (English)

- Filter Tokens (by Length)

- Stem (Porter)

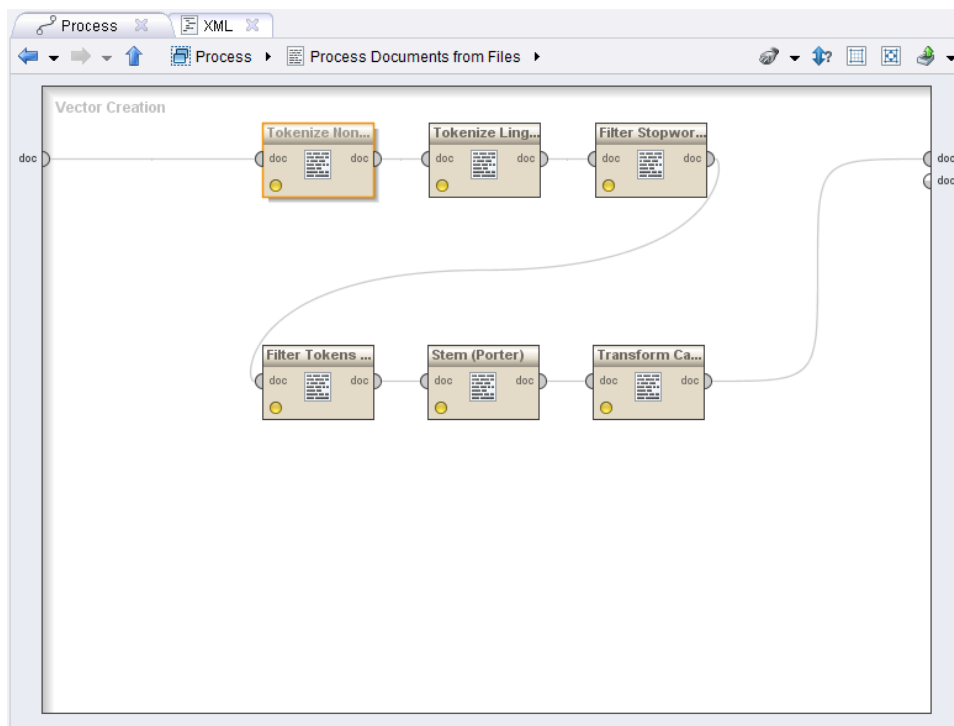- Transform Cases (2) (Transform Cases)

**Figure 1.12** Operators within the Process Documents from Files nested operator
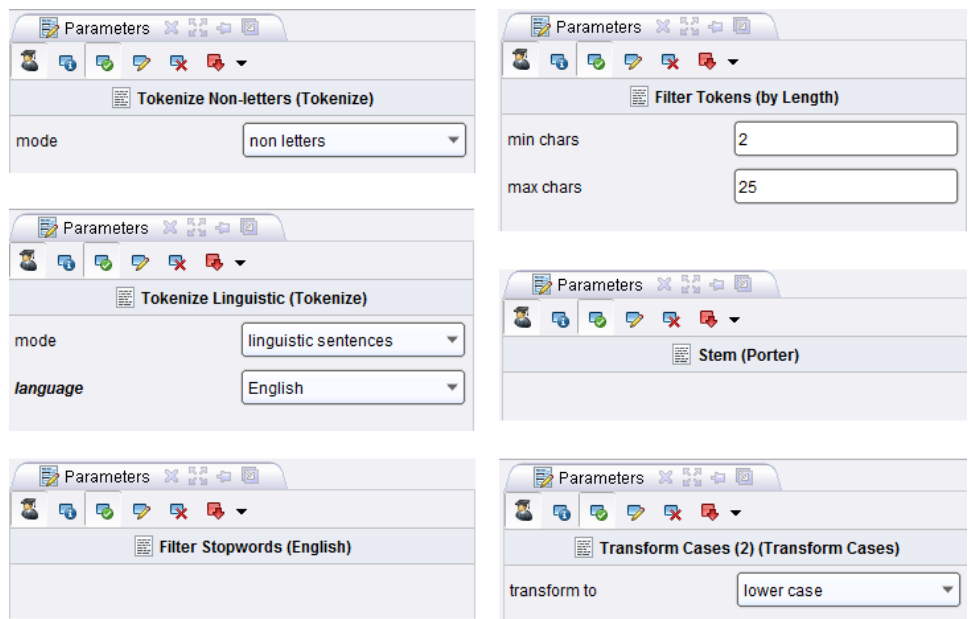
**Figure 1.13**    Parameters for the operators within the Process Documents from Files operator

The sub-process basically transforms the text data into a format that can be easily analyzed using conventional data mining techniques such as association mining and cluster modeling. The parameters for each of the operators in this sub-process (within the Process Documents from Files operator) are displayed in Figure 1.13. Notice that Figure 1.13 was created manually by combining multiple snapshots.

In this sub-process, the Tokenize Non-letters (Tokenize) and Tokenize Linguistic (Tokenize) operators are both created by selecting the Tokenize operator, but with different parameter selections. The former operater tokenizes based on non letters whereas the latter operater tokenizes based on the linguistic sentences within the English language. The Filter Stopwords (English) operator removes the stop words in the English language from the text data set. The Filter Tokens (by Length) operator removes all the words composed of less than min chars characters and more than max chars characters. In this example, words that have less than 2 characters or more than 25 characters are removed from the data set. The Stem (Porter) operator performs stemming and the Transform Cases(2) (Transform Cases) operator transforms all the characters in the text into lower case. It should be noted that the name of this last operator is not a good name, and we, as the modelers have forgotten to rename the operator after its name was automatically assigned by RAPIDMINER. This mistake should be avoided in constructing the processes.
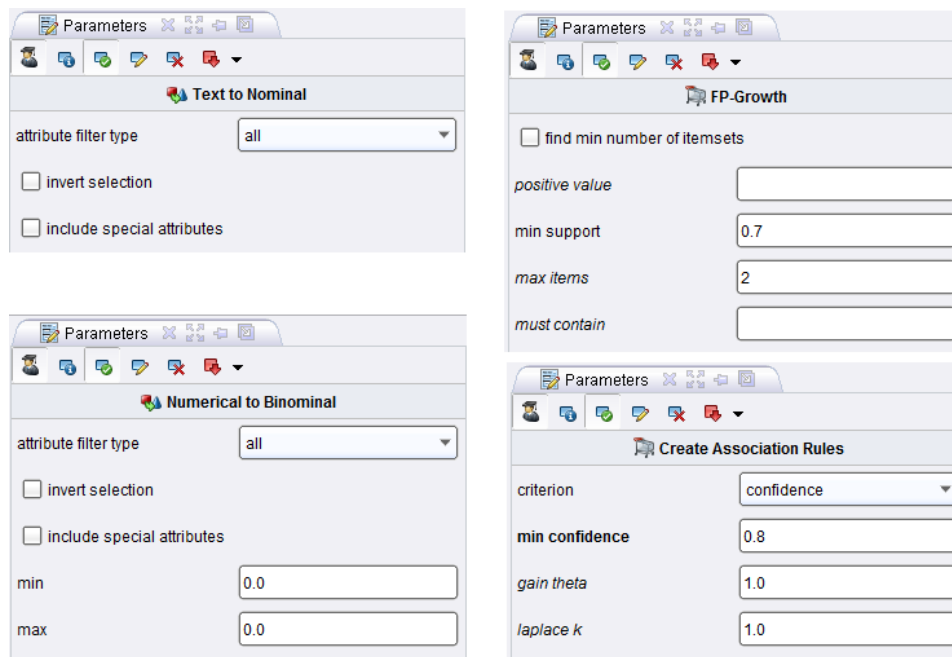
**Figure 1.14**   Parameters for the operators in Process01

The parameters for each of the operators in Process01 (Figure 1.11) are displayed in Figure 1.14. Notice that Figure 1.14 was created manually by combining multiple snapshots. The Text to Nominal operator transforms the text data into nominal (categorical) data. The Numerical to Binomial operator then transforms the data into binominal form. This means that each row represents a document, a few columns provide meta-data about that document and the remaining columns represent the words appearing in all the documents, with the cell contents telling (true or false) whether that word exist in that document or not. FP-Growth algorithm is used for identifying the frequent item sets. In this example, the min support parameter is 0.7 (Figure 1.14), meaning that the operator generates a list of the frequent sets of words (itemsets) that appear in at least 70 % of the documents. Notice that, it will be computationally efficient to select the min support value in the FP-Growth operator to be equal to prune below percent value for the Process Documents from File operator (Figure 1.11) divided by 100. Also, the max items parameter is 2, meaning that the generated list is limited to pairs of words (2-itemsets), and the list will not contain frequent word sets (itemsets) with 3 or more words in them. The final operator in Process01, namely Create Association Rules, receives the list of frequent word sets from the FP-Growth operator, and computes the rules that satisfy the specified constraints on selected association mining criteria. In this example, the association rules are computed according to the the criterion of confidence, as well as gain theta and laplace k. The specified minimal values for these 3 criteria are 0.8, 1.0 and 1.0, respectively.

### 1.2.3  Saving Process01

So far Process01 has been imported in to the workspace of RAPIDMINER. Now it is a good time to keep it in the LocalRepository so that we will not have to import it again next time we work with RAPIDMINER. Click on the Repositories view, right-click on the LocalRepository and select Configure Repository as shown in Figure 1.15. This is actually an initialization step before using RAPIDMINER for the first time, but we will still go through this step to ensure that we are saving everything to a preferred folder directory in our computer. In this chapter, we will be saving everything under the Root directory of C:\RapidMiner, as in Figure 1.16. Next, we can save Process01 in to the LocalRepository. Right-click on the LocalRepository text and select Store Process Here, as in Figure 1.17. When the Store Process dialog window appears, click Ok. Our importing and saving of Process01 is now completed.

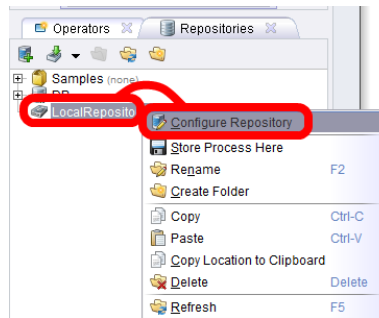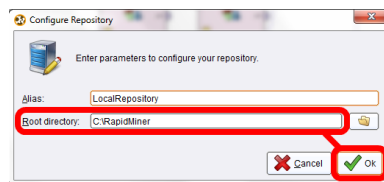**Figure 1.15**    Configuring LocalRepository



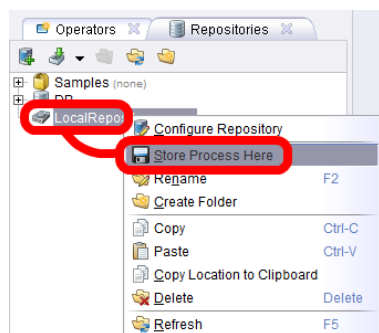**Figure 1.16**    Specifying the Root directory for LocalRepository



**Figure 1.17**    Storing Process01 in LocalRepository

## 1.3  CLUSTERING TEXT DOCUMENTS (Process02)

The second text mining process that we will introduce in this chapter is Process02, and involves the clustering of the 100 documents in the text collection.

### 1.3.1  Importing Process02

To import this process, select File menu, Import Process menu item, select Process02.rmp within the Import Process window, and click Open.

### 1.3.2  Operators in Process02

Process02 is given in Figure 1.18. Similar to Process01, Process02 begins with the Process Documents from Files operator, whose parameters are given in Figure 1.18. In Process02, the selected vector creation method is different; it is Term Frequency. The impact of this new selection will be illustrated later. For now, it should be noted that this selection results in the computation of the relative frequencies of each of the words in each of the documents in the data set. For example, if a word appears 5 times within a document that consist of 200 words, then the relative frequency of that word will be $5/200 = 0.025$. This value of 0.025 will appear under the column for that word, at the row for that document. In Process02, the prune method is again percentual, just as in Process01 (Figure 1.11). However, the value for the prune below percent parameter is now different.
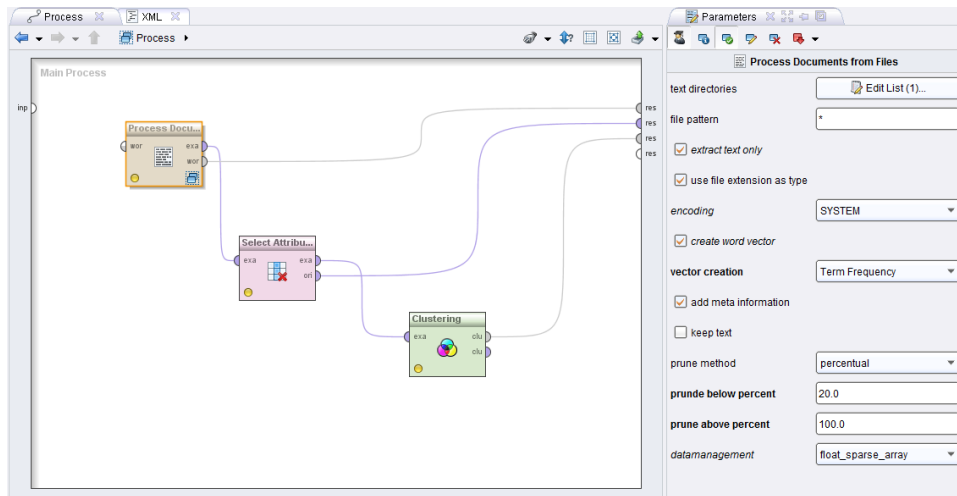


**Figure 1.18**  Operators in Process02 and the Parameters for Process Documents from Files operator
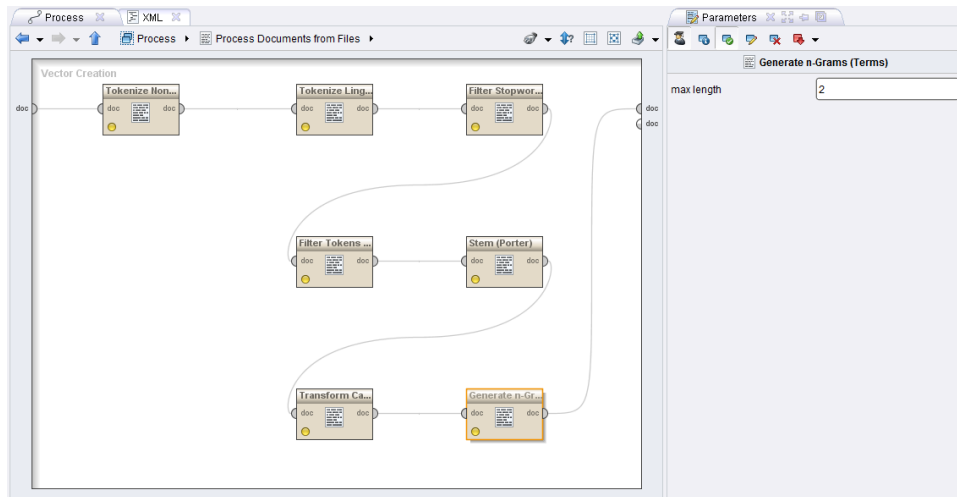
**Figure 1.19** Operators within the Process Documents from Files nested operator and the Parameters for the Generate n-Grams (Terms) operator
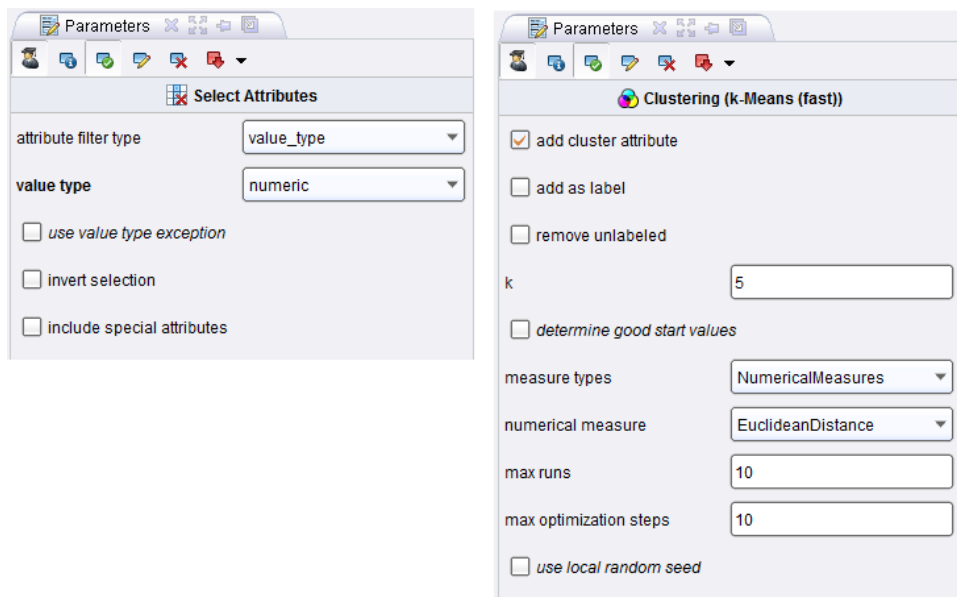


**Figure 1.20** Parameters for the operators in Process02

The prune below percent parameter in Process01 was 70.0, whereas it is now 20.0 in Process02. The reason for this change is due to the fact that the applied data mining technique in Process01 is association mining, whereas in Process02 it is clustering.

The former technique is computationally much more expensive then the latter, meaning that running association mining algorithms on a data set takes much longer running time compared to the k-Means clustering algorithm on the same data set. Therefore, with the same amount of time available for our data mining process, we have to work with much smaller data sets if we are carrying out association mining, as in Process01. The values displayed in Figures 1.11 and 1.18 have been determined through trial and error, such that the running time for the processes do not exceed 30 seconds on a laptop with Intel i5 processor and 4GB RAM.

Figure 1.19 shows the contents of this operator. While there were six operators within this nested operator in Process01, there are now seven operators in Process02. The newly-added operator is Generate n-Grams (Terms). The only parameter for this operator is max length which is set equal to 2 in our example (Figure 1.19).

The parameters for the Select Attributes and Clustering (k-Means(fast)) operators within Process02 are displayed in Figure 1.20. The Select Attributes operator takes the complete data set and transforms it into a new one by selecting only the columns with numeric values, i.e. columns corresponding to the words. This transformation is required for the next operator, which performs clustering based on numerical values. Clustering (k-Means (fast)) operator carries out the k-Means clustering algorithm on the numerical data set. Since each row of the data (each document in the text collection) is characterized in terms of the occurrence frequency of words in it, this operator will place together the documents that have a similar distribution of the word frequencies. The k value, which denotes the number of clusters to be constructed, is set to 5 in our example.

### 1.3.3  Saving Process02

Now, let us save Process02 into LocalRepository: Click on Repositories view, right click on LocalRepository view and select Store Process Here. When the store process dialog box appears click Ok. As shown in Figure 1.21, both processes are now saved under local repository. Since we had earlier defined the directory root for local repository as C:\RapidMiner, the saved processes will appear as .rpm (RAPIDMINER process) files under that directory.

Unfortunately, there is a problem with the naming of the files: the file names also contain the extension .rpm, in addition to the .rpm extension itself. For example, the correct name for Process01rmp.rmp under C:\RapidMiner should be Process01.rmp. Right click on Process01.rmp text within RAPIDMINER, and select rename as in Figure 1.22 then, as in Figure 1.23 remove the .rmp extension from the file name and click OK. Do the same for renaming Process02 and you will obtain the results in Figure 1.24.

**Figure 1.21**    Saved processes within LocalRepository

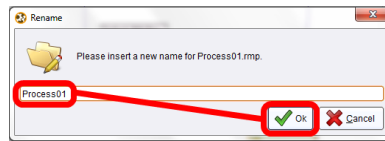

**Figure 1.22**    Renaming a process

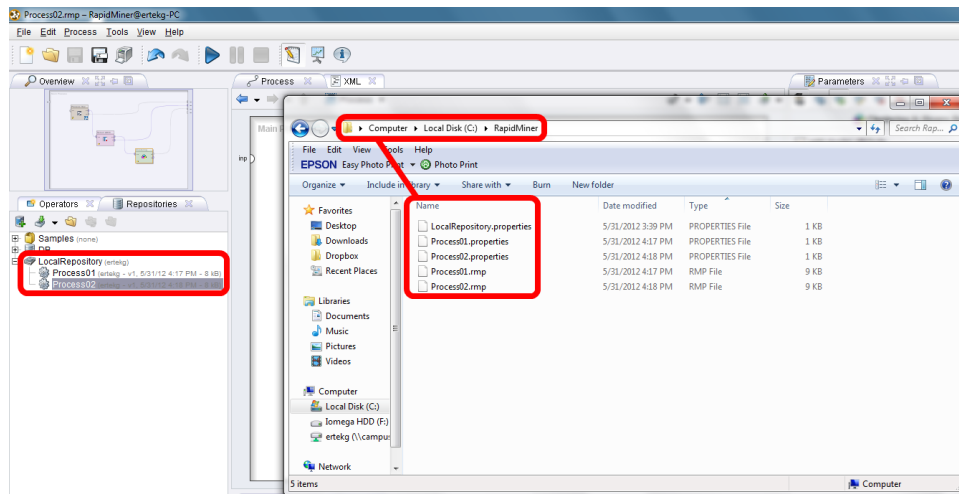**Figure 1.23**   Renaming Process01 correctly



**Figure 1.24**   Corrected names for Process01 and Process02 in LocalRepository and the corresponding files

## 1.4   RUNNING Process01 AND ANALYZING THE RESULTS

Having imported and saved the processes into local repository, now we will run these processes.

### 1.4.1   Running Process01

For running Process01, double click on Process01.rmp under the Repositories view. Process01 will be loaded into the Process view. Click the run button in the tool bar, as in Figure 1.25. You may be prompted and asked whether you would like to close all results before starting process (Figure 1.25). Click Yes. When the process run is completed, you will be prompted again and asked whether you would like to switch to the result perspective (Figure 1.26). Click Yes and you will have the screen in Figure 1.27. This is indeed the *result perspective*, displaying the result obtained through your process. Until now, you worked on your process in the *design perspective*, where we could change your process by adding/removing/editing operators. Now, in the result perspective, we are not able to modify our process. In the tool bar, we can switch to the design perspective or the result perspective by clicking their respective buttons.



**Figure 1.25**   Opening and running Process01



**Figure 1.26**   Dialog box alerting the switch to the result perspective

**Figure 1.27** Result Overview for Process01 results

### 1.4.2  Empty Results for Process01

In result perspective, we can see the latest results at the bottom of the result overview. To see the details click inside the bar with the most recent date and time and you will see an expanded view as in Figure 1.27. In the extended view you will note that there are three types of results obtained as a result of running Proces01. These three results correspond precisely to the three result ports for Process01 which are represented by the three half-circles on the right of the process view in Figure 1.11 and are labeled with res. The first result is the wordlist (wor) arising from the Process Documents From Files operator (Figure 1.11). The second result is the example set exa, which arises from the FP-Growth operator. The third and final result is the list of association rules (rul) which arises from the Create Association Rules.

In Figure 1.27 the word list contains 0 entries and the number of examples is 0. But why? Was there something wrong with our processes? We will answer this question next.

### 1.4.3  Specifying the Source Data for Process01

Running Process01 resulted in no tangible results and the reason is very simple. We have not yet specified the text data set that will be processed and mined in Process01. To specify the data source, as in Figure 1.28 switch to the design perspective by clicking the notepad button, click on the Process Documents from Files operator and click on the box to the right of the parameter text directories.

In the dialog box click on the directory button, (see Figure 1.29) select the source data folder, double click TripAdvisor_First_100_Hotels folder, and click the open button (Figure 1.30). Now the directory for the source text data appears in the dialog box, as in Figure 1.31. Now click Ok.

**Figure 1.28**    Specifying the data source text directories for Process01



**Figure 1.29**    Specifying the text directories

**Figure 1.30** Specifying the text directories

**Figure 1.31** Specified text directories

### 1.4.4   Re-Running Process01

Having specified the data source, we can run Process01 again. Click on the run button (see Figure 1.32) and click Yes when you are asked. You should click Save Process Before Start?, and click Yes when you are asked whether to Close old results before starting process. When the process is running you can observe which operator is currently being used. For example, in Figure 1.33 the FP-Growth algorithm is running, as can be understood from the small green triangle on its lower left corner. From Figure 1.33 we can also read (from the lower left corner of the RAPIDMINER window) that Process01 has been running for 14 seconds (14 s) and the FP-Growth operator has been running for 7 seconds (7 s).



**Figure 1.32**   Running Process01 again



**Figure 1.33**   Running of Process01

### 1.4.5    Process01 Results

When the new results are created, click Yes to the result perspective. In the result perspective you will now see a see a blue bar. Click inside the bottommost blue bar and you will see an overview of the results, as in Figure 1.34. This time the world list is not empty, it Contains 631 entries. The Number of examples in the example set is equal to 100, and there are 250 attributes (Figure 1.34).



**Figure 1.34**    Result Overview for Process01 results



**Figure 1.35**    WordList generated by Process01

**Figure 1.36**    Meta Data View for the ExampleSet generated by Process01



**Figure 1.37**    Data View for the ExampleSet generated by Process01

Click on the world list view to display the words found by the Process Documents from Files operator, as in Figure 1.35.

Click on the ExampleSet (Numerical to Binomial) view to display the example set. Figure 1.36 presents the Meta Data View, which gives the metadata (information on the data attributes) of the ExampleSet. Now click on the Data View button to display the data itself (Figure 1.37). Each row corresponds to a text document and the first column is an automatically generated key attribute. The extra four columns contain the label, file name, file location and the file creation data. The remaining columns correspond to the world list. The cells under the word attributes take the binominal value of true or false, denoting whether the word exists in that document.

The final result is the set of association rules. Click on the AssociationRules(Create Association Rules) tab to display the association rules (Figure 1.38). In the Table View, table grid presents the generated association rules, with one rule in each row. For example, the second row states "IF worth THEN holidai" with a Support level of 0.760 and Confidence level of 0.800. This rule means that in 76 of the 100 documents, words with stem worth and holidai appear together. Furthermore in 80% of the documents where a word derived from the stem worth appears, at least one word derived from the stem holidai is observed. On the left hand side of the grid we can select and Show rules matching a particular set of words (Figure 1.38).



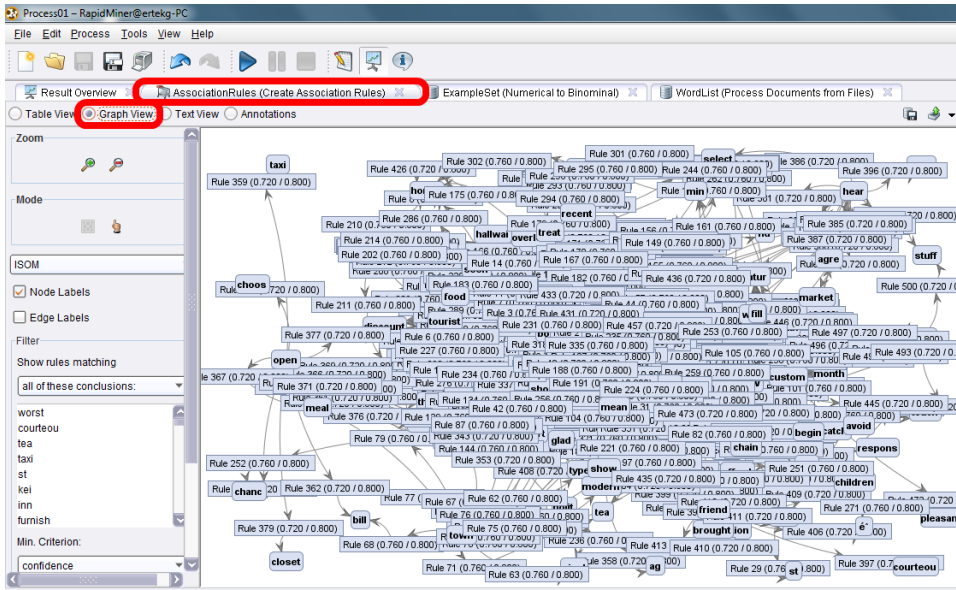**Figure 1.38**   Table View for the AssociationRules generated by Process01

**Figure 1.39** Graph View for the AssociationRules generated by Process01
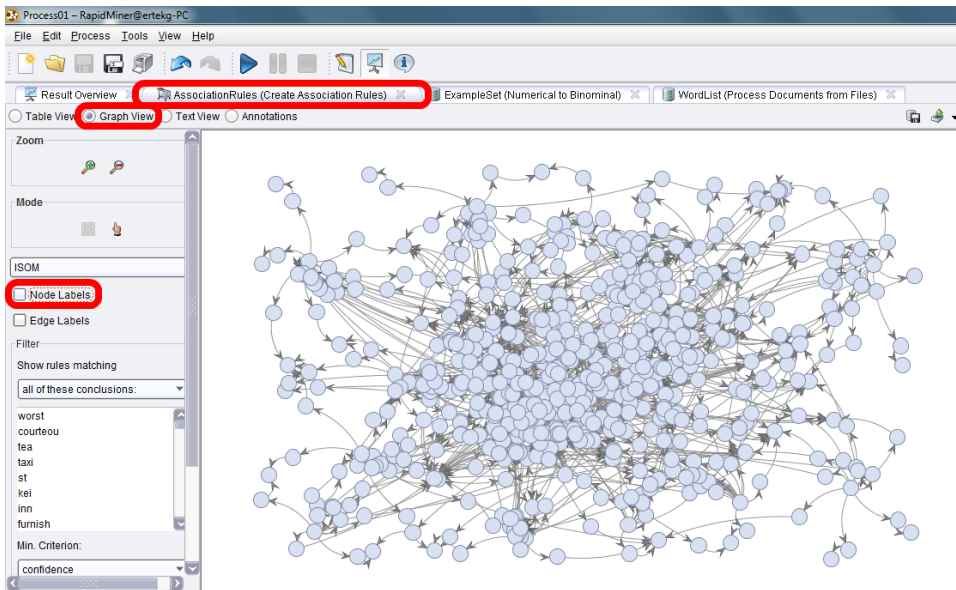


**Figure 1.40** Graph View for the AssociationRules, without the node labels
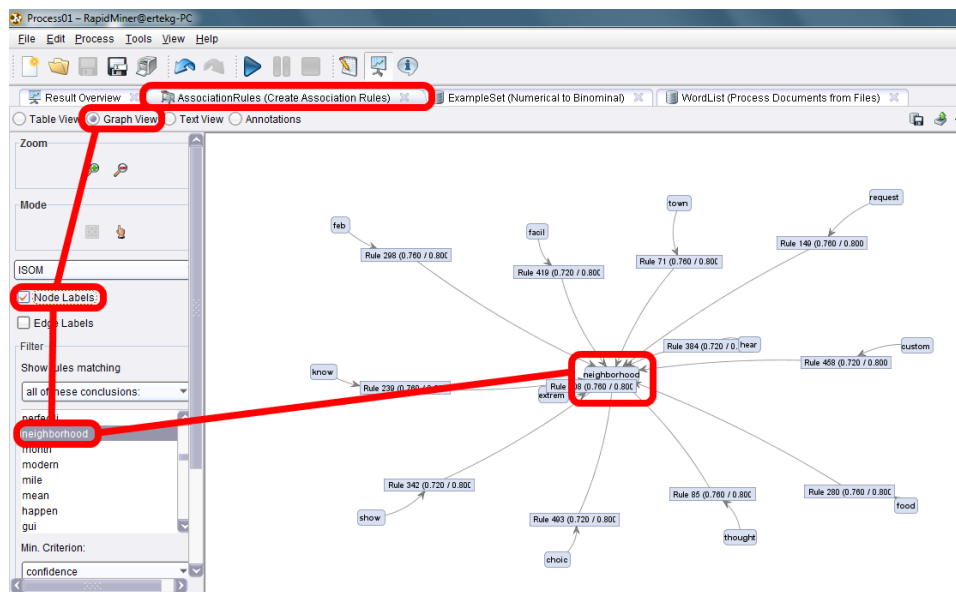
**Figure 1.41** Filtering rules in the Graph View for the AssociationRules

These association rules can also be displayed in the form of a graph. For this, click on the Graph View radio button as in Figure 1.39. The initial graph visualization may be cluttered due to the Node Labels. As Click on the check box Node Labels to uncheck that option and eliminate or reduce the clutter in (Figure 1.40). In this window, filtering can again be carried out by selecting words from the list box on the left hand side. Once a filtered word is selected you will want to again display the Node Labels as in Figure 1.41. At this point, select a word from the list box and check the Node Label check box (Figure 1.41).

The final interesting result that we would like to share is displayed in Figure 1.42. Click on WordList(Process Documents from Files) tab to return back to the word list. Click on the attribute label Document Occurences to obtain a descending order. Then move the vertical scroll bar to see the border line where the value changes from 100 to 99. From Figure 1.42 we can observe that the word work appears in all 100 documents, whereas the word big appears in 99 documents.



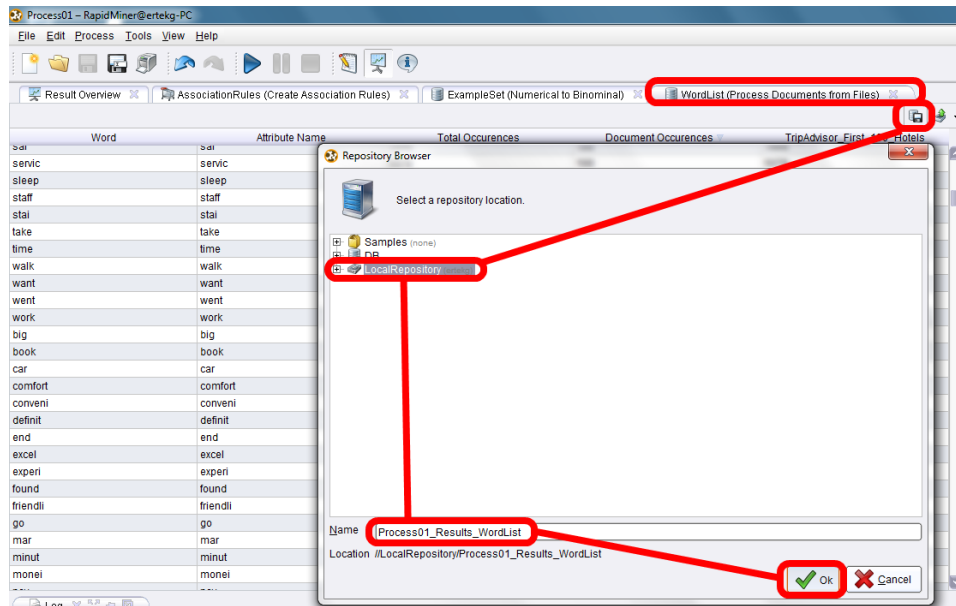**Figure 1.42**   Document Occurences of the words in the WordList

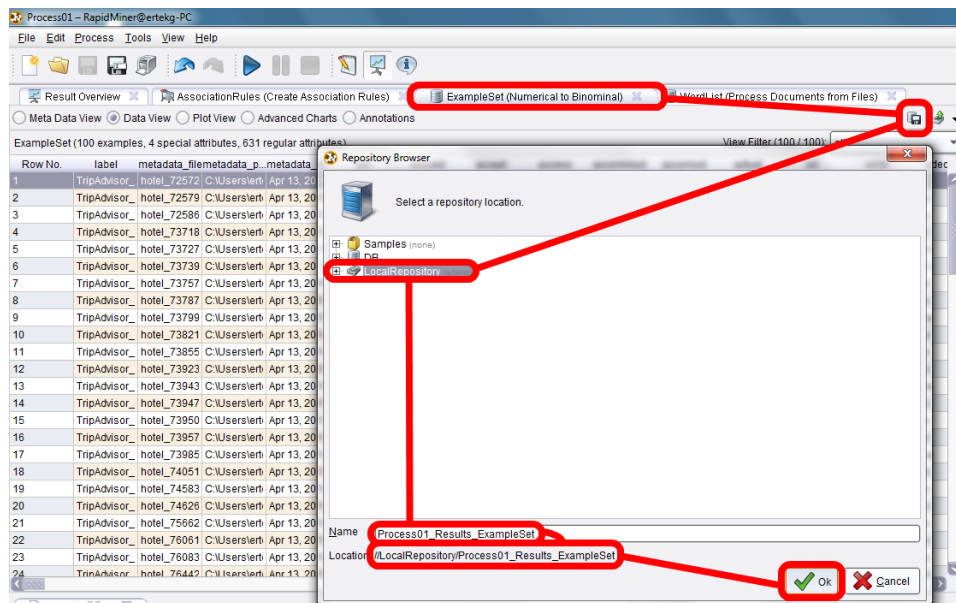**Figure 1.43** Saving the WordList for Process01



**Figure 1.44** Saving the ExampleSet for Process01

### 1.4.6 Saving Process01 Results

Having generated results of the Process01 we will save these results into LocalRepository, so that we can save time when we want to analyze the results next time. The three types of results, namely the Word List, ExampleSet, AssociationRules have to be saved individually. First, click on the save button in the upper right corner of the results, click on LocalRepository in the appearing window, and type in an appropriate name inside the Name text box as in Figure 1.43, then click Ok.

Next, click on the ExampleSet (Numerical to Binominal) tab, and repeat the same steps, as in Figure 1.44.

Finally, click on the AssociationRules, as in Figure 1.45. For this last result we would like to show you a common possible mistake. After clicking the save button, omit clicking on LocalRepository and click Ok as Figure 1.45. Please note that in this case we did not specify a location for the generated results. This will cause an error . When you see the error dialog box , click Close. Then, follow the correct steps; Click on the save button, click on LocalRepository, give an appropriate unique name and click Ok as in Figure 1.46.
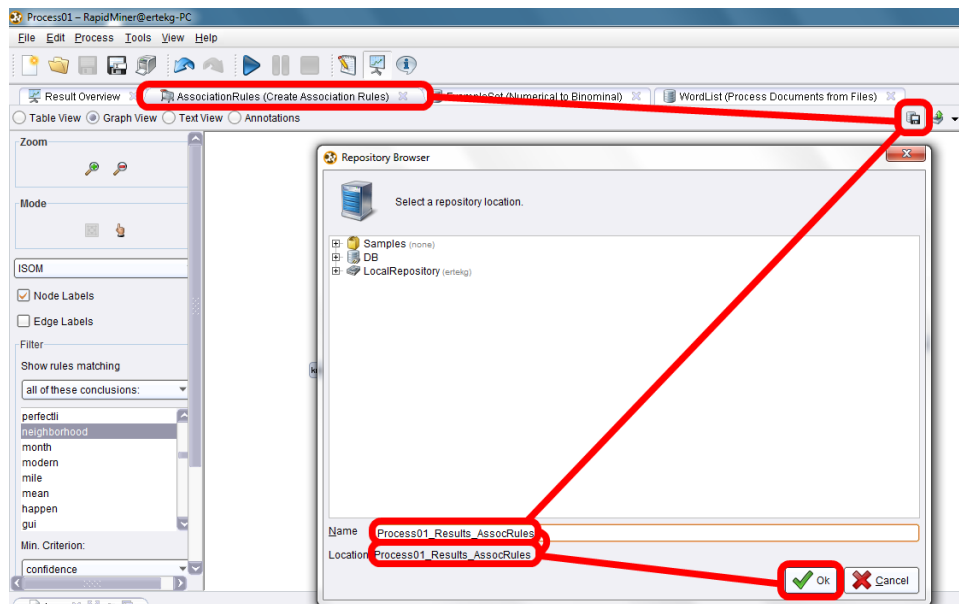


**Figure 1.45** Saving the AssociationRules incorrectly, without selecting LocalRepository or another directory
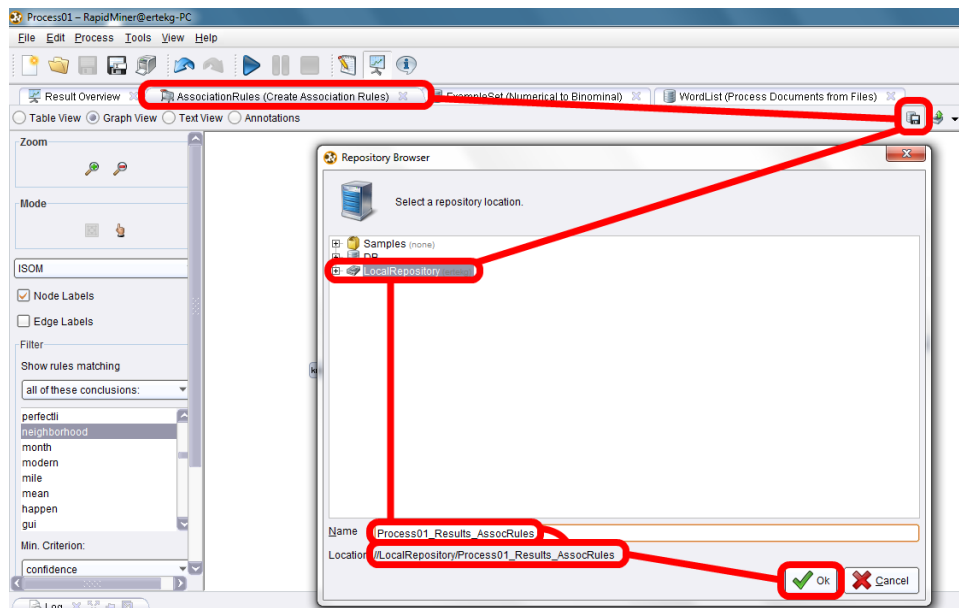
**Figure 1.46**   Saving the AssociationRules correctly, selecting the directory to be saved into

Now that all the data for the generated results are saved let us also see how we can save (export) the data visualizations. Click on the export button to the right of save button, as in Figure 1.47. Then specify the directory where the visualization will be saved (exported) and select the image file type as in Figure 1.48. In this example, select the .png image format, and then click Ok. You can check that the visualization is now saved as an image file, as specified.



**Figure 1.47**   Exporting the graph visualization as an image file



**Figure 1.48**   Specifying the directory to be saved into, the file name, and the file type

## 1.5 RUNNING Process02 AND ANALYZING THE RESULTS

In this section we will run Process02 and analyze the generated results.

### 1.5.1 Running Process02

To switch to the design view click on the design view button (Notepad icon) as in Figure 1.49. Then double click on Process02 inside the Repositories view, and click the run button as in Figure 1.50. You will be prompted with an error dialog box, as in Figure 1.51.



**Figure 1.49**    Switching to the design view



**Figure 1.50**    Opening and running Process02

### 1.5.2 Specifying the Source Data for Process02

The reason for this error is that we have not specified the data set for Process02 to solve this problem. To resolve this, first click on Close (Figure 1.51). Then, as in Figure 1.52, click on Process Documents from Files operator and click on the button next to the text directories parameter inside the Parameters view. Then, select the same data set as you have selected for Process01, as shown in Figures 1.28, 1.29, 1.30 and 1.31. Now the data source has been specified and we can run Process02 to obtain tangible results. Click on the run button as in Figure 1.50, and answer Yes to the questions that are prompted.



**Figure 1.51**    Error message due to not specifying the source data for Process02



**Figure 1.52**    Specifying the text directories for the source data of Process02

### 1.5.3 Process02 Results

The results of Process02 are shown in Figure 1.53. There are three types of results generated after running Process02:



**Figure 1.53** Result Overview for Process02 results

Firstly, the Word List contains the filtered words from the documents in our example, the Word List generated Contains 5175 entries.

The second result is the Example Set, which contains data requiring the word composition of each of the documents in the database.

The third result of Process02 is Centroid Cluster Model. In our example five clusters (Cluster 0, Cluster 1, Cluster2, Cluster 3, Cluster 4) have been generated, where each cluster contains a subset of 100 documents.

The Word List is very similar to the word list of Process01 word list. However, there is a difference due to the additional Generate n-Grams (Terms) operator (Figure 1.19). This additional operator extracts n-grams (sequences of words), as well as single words. To see these generated n-grams as a part of the word list you can click on the WordList view or the Example view. Click on the Example view, as in Figure 1.54 and you will see several of these n-grams. For example, besides the word absolut the n-gram words absolut_love, absolut_perfect, absolut_stai are part of the WordList, and appear as attributes in the ExampleSet (Figure 1.54). Now, to see the data itself, click on Data View radio button as in Figure 1.55. In this grid, each number tells the relative frequency of a word in a document. For example, the word absolut (the word that have absolut as their stem) constitutes 0.030 proportion (3%) of all the words in document number 13, that contains the reviews for hotel_73943.



**Figure 1.54**    Meta Data View for the ExampleSet generated by Process02, including the n-Grams

The main contribution of Process02 is the cluster model, which clusters documents, and thus the hotels, according to the similarity of the frequency distribution of the words contained in their TripAdvisor reviews. Now click on Cluster Model (clustering) view to view the cluster model results, as in Figure 1.56. Then, to view the content of each of the clusters, click on Folder View radio button, as in Figure 1.57. In our example the first cluster of hotels is cluster_0, which contains 28 hotels, including hotels with row numbers 47, 51, 55, 95, 100.



**Figure 1.55** Data View for the ExampleSet generated by Process02, and the relative occurence frequency of the word absolut in `hotel_73943.txt`



**Figure 1.56** Text View for the Cluster Model generated by Process02, displaying the number of examples in each cluster
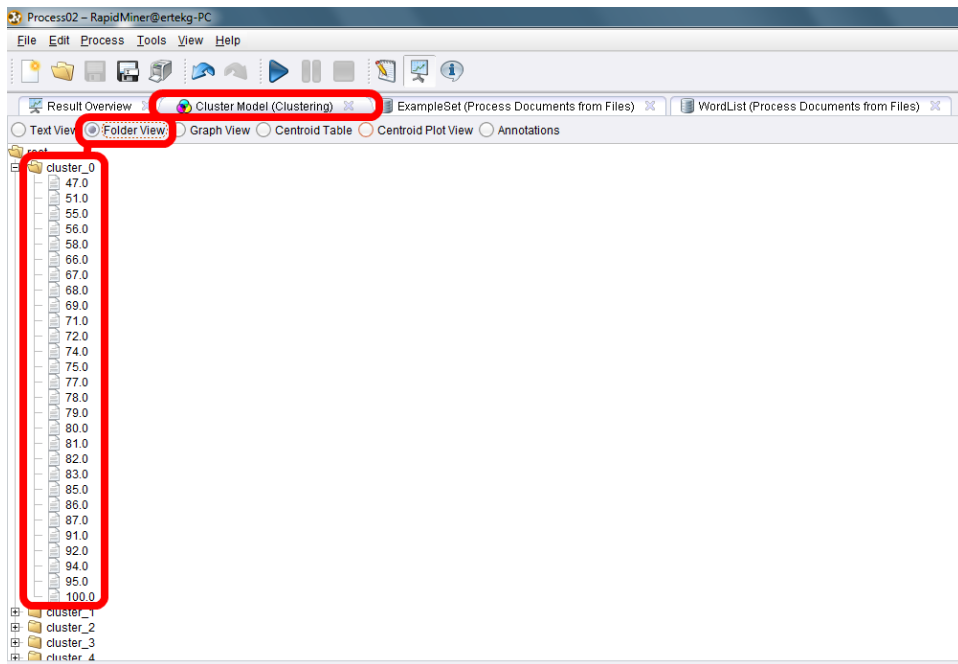
**Figure 1.57**     Folder View for the Cluster Model generated by Process02, displaying the examples in each cluster

How do these clusters differ from each other? To get an answer to this question, click on Centroid Table radio button (Figure 1.58). In this Centroid table we can observe the average frequency of each word in the documents of each cluster. For example the word absolut appears much more frequently in cluster_3, compared to the other clusters: In cluster_3, the average frequency for absolut is 0.013, whereas its at most 0.008 in the others.



**Figure 1.58** Centroid Table for the Cluster Model generated by Process02, displaying the average frequency of each word in each cluster

Finally, save the results of Process02, just as you did for Process01. The final picture in the Repositories view will be as in Figure 1.59.
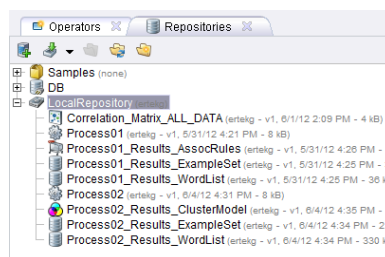


**Figure 1.59** Final view of LocalRepository with the processes and their results saved

## 1.6   CONCLUSIONS

In this chapter, we have shown the basic RAPIDMINER operators for text mining, as well as how they can be used in conjunction with other operators that implement popular conventional data mining techniques. Throughout the chapter, we have discussed two processes, Process01 and Process02, which complement text processing with association mining and cluster modeling, respectively. We have shown how these text mining processes, which are formed as combinations of text processing and data mining techniques, are modeled and run in RAPIDMINER, and discussed how their results are analyzed. RAPIDMINER has an extensive set of operators available for text processing and text mining, which can also be used for extracting data from the web, and perform a multitude of other types of analysis. While most of these operators were not discussed in this chapter due to space limitations, the reader is strongly encouraged not to suffice with the contents of this chapter, and explore and experiment other text processing and mining operators and capabilities within RAPIDMINER, its official extensions, and the processes posted by the RAPIDMINER user community.

### Acknowledgement

# GLOSSARY

**Decision Tree**    Decision Trees are....

**Data Exploration**    Data Exploration is of particular importance when...