

# User Manual: Using RPM with TSAT

June 12, 2018

## Contents

<b>1</b>	<b>Introduction to TSAT</b>	<b>2</b>
1.1	Interface Layout . . . . .	2
1.2	Tools Introduction . . . . .	2
1.3	Overview . . . . .	2
<b>2</b>	<b>Time Series Classification using RPM</b>	<b>2</b>
2.1	File formats . . . . .	2
2.2	Training . . . . .	4
2.3	Testing . . . . .	9
2.4	Testing Unlabeled Data . . . . .	11
2.5	Saving a Trained RPM Model . . . . .	12
2.6	Loading an RPM Model . . . . .	15
2.7	Settings . . . . .	18
2.7.1	Dynamic Time Warping . . . . .	18
2.7.2	Iterations . . . . .	23
<b>3</b>	<b>Motif Discovery</b>	<b>24</b>
3.1	File format . . . . .	24
3.2	Guide to Motif Discovery . . . . .	24
<b>4</b>	<b>Anomaly Detection</b>	<b>24</b>
4.1	File format . . . . .	24
4.2	Guide to Anomaly Detection . . . . .	24
<b>5</b>	<b>Error Messages</b>	<b>24</b>
5.1	File Errors . . . . .	24
5.2	RPM Errors . . . . .	24
<b>6</b>	<b>FAQs</b>	<b>24</b>

# 1 Introduction to TSAT

## 1.1 Interface Layout

## 1.2 Tools Introduction

Representative Pattern Mining - RPM

Motif Discovery

Anomaly Detection

## 1.3 Overview

# 2 Time Series Classification using RPM

TSAT implements Representative Pattern Mining or RPM (see Section 1.2 for more details) to perform time series classification. In order to perform time series classification you will need a training and a test dataset containing time series data.

The standard method to train a supervised learning classifier is to take the labeled dataset and split it into two datasets, training and testing data. One common way to split the data is to have 80% training and 20% testing.

**Training Data** Training data is the primary data and will be used to create a model that can identify similar patterns in new, unlabeled, data. This data must have a label for each time series so that RPM can learn what the labels can look like. This is where the bulk of the data should be set aside for as RPM will need many samples to find representative patterns.

**Testing Data** Testing data is a small subset of the data usually from the same source as the training data, but not found in the training data. This set of data will be used to test the model that RPM made for accuracy or to predict labels for unlabeled test data.

Splitting data into a training and a test set is beyond the scope of this manual and is not done by TSAT. The goal of this section is to first detail the proper file formats for training and testing data in Section 2.1. Then the proper procedures to train (Section 2.2), test (Section 2.3), and review the results (Section ??) are presented step by step.

## 2.1 File formats

File formatting is very important in TSAT and especially when using RPM. If the file is not in the correct format TSAT will not be able to read the file and may produce unexpected results or error messages. The data may be formatted by column, row, or following the ARFF file format. Additionally, the labels for the time series may be any string excluding white space and “?” as this is reserved for unknown values in test data.

Figure 1: Examples of RPM Data

# 1 1 1 2 2 2 2 2 2	# 1.0000000e+000 1.0000000e+000 1.0000000e+000
0 0 0 122880 122880 122880 122880 0 0	-4.6427649e-001 -8.9697208e-001 -4.6469596e-001
0 0 0 0 0 0 0 0 0	-5.5504787e-001 -6.8568553e-001 -5.6773891e-001
0 0 0 0 0 0 0 0 0	-8.4284310e-001 -1.3513818e+000 -3.2022764e-002
0 0 0 0 0 0 0 0 0	-8.6589548e-001 -1.4586668e+000 -6.3504562e-001
0 0 0 0 0 0 0 0 0	-9.3639631e-001 -1.1653456e+000 -6.0282554e-001
0 0 0 0 0 0 0 0 0	-8.1726995e-001 -1.4039293e+000 -2.6685628e-001
0 0 0 0 0 0 0 0 0	-2.6361216e-001 -1.8217996e+000 -2.6706128e-001
0 0 0 0 0 0 0 0 0	-1.2580483e+000 -8.3160109e-001 -9.3104230e-001
0 0 0 0 0 0 0 0 0	-1.2503934e+000 -1.0163124e+000 -4.4938186e-001
0 0 0 0 0 0 0 0 0	-9.1830825e-001 -8.0353040e-001 -7.2134200e-001
0 0 0 0 0 0 0 0 0	-9.2210226e-001 -1.2595048e+000 -3.9727192e-001
0 0 0 0 122880 122880 0 0 0	-9.8448828e-001 -1.1392341e+000 -9.6212589e-001
	-1.2880511e+000 -8.7865203e-001 -1.4206669e+000

(a) Example 1

(b) Example 2

**Column Formatted Data** The data files are simple text files that store the time series data with one entry per column, with a space delimiter, with each row representing a time step in the time series data. With RPM compatible data the first row in the file starts with a “#” with rest of the row containing the label for each time series rather than the time series values. If the file is missing this row RPM will not be enabled in TSAT. Examples of column formatted RPM compatible data can be seen in figure 1. Another thing to keep in mind is that in this format the time series must all be the same length.

**Row Formatted Data** Another acceptable format is the row format. This format is especially useful when the time series are not all the same length as each row or time series may have its own length. In this format the first line of the file is a “#” followed by a new line. Starting on the second line, each line starts with the label followed by the corresponding time series (each value separated by a space). There should be no empty lines. For example,

```
#
1 -5.3 -23 5 ...
1 23 1 5 3 1 ...
two 23 3 4 200 ...
two 42 3 4 102 ...
...
```

In this example the labels are “1” and “two” and the time series follow after the labels.

**ARFF Formatted data** A standard format for many public time series datasets is the ARFF file format. For example, <http://timeseriesclassification.com/dataset.php> has a number of time series in ARFF format. ARFF files are more complicated than both the column and row formats, but is more widely used outside TSAT. Here is an abbreviated example ARFF file:

```
@relation Adiac
```

```
@attribute att0 numeric
@attribute att1 numeric
...
@attribute target 1,2,...
```

```
@data
1.3749,1.2894,1.2043,1.1194,1.0347, ... 1
1.7257,1.7001,1.6611,1.6089,1.5319, ... 2
...
```

The ARFF file begins with the name of the dataset Adiac by using the ARFF formatting by putting it after the `@relation` element. After the name of the dataset each timestep is listed as an attribute `@attribute <timestepName> numeric` where you can choose what to name each timestep. After listing the timesteps as attributes the labels are listed as the target attribute `@attribute target {1, 2, ...}` where these are the labels for the time series. Finally, the time series data is in comma separated value (CSV) format following the `@data` line. Each value in a time series is separated by a comma on a single line and the last value on the line is the label for the time series.

**Unknown Test Data** In column, row, or ARFF format when predicting unlabeled test data, the test data must be labeled as “?” (note that there must only be test data that is labeled with a “?”). For example, a row formatted test dataset might be:

```
#
? -5.3 -23 5 ...
? 23 1 5 3 1 ...
? 23 3 4 200 ...
...
```

As can be seen the label is “?” and the time series follows after the label. When training there must always be more than one example from each class label and there must be more than one label.

## 2.2 Training

Once you have the data in the proper format and TSAT open training RPM can begin.

**Step 1** First click on the “Browse” button under the “Data Sources” section of the window, as seen in figure 2.

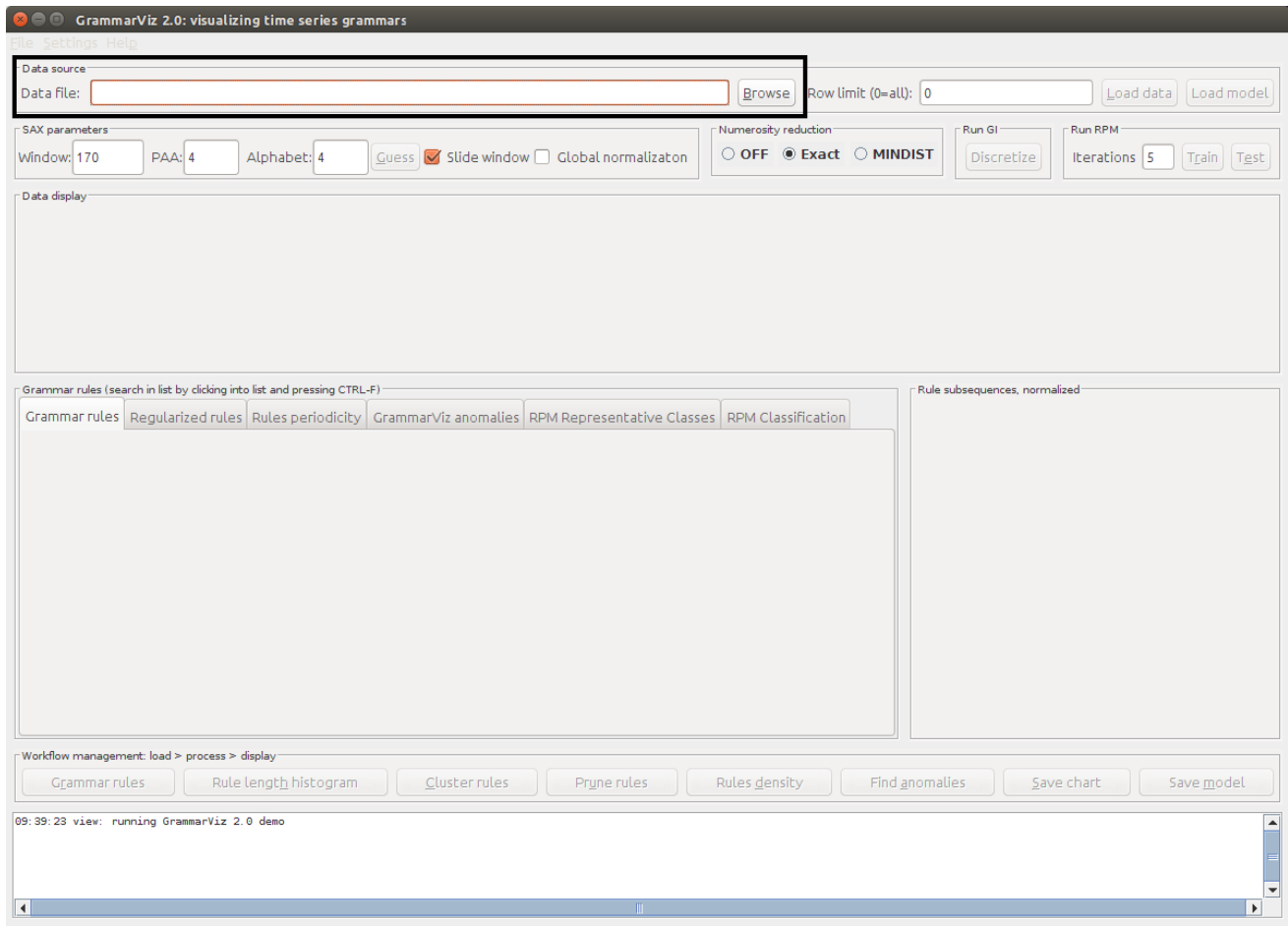


Figure 2: Open TSAT

**Step 2** This should bring up the file browser prompt in figure 3. Using this prompt select the file containing the training set in the RPM compatible format, figure 4.

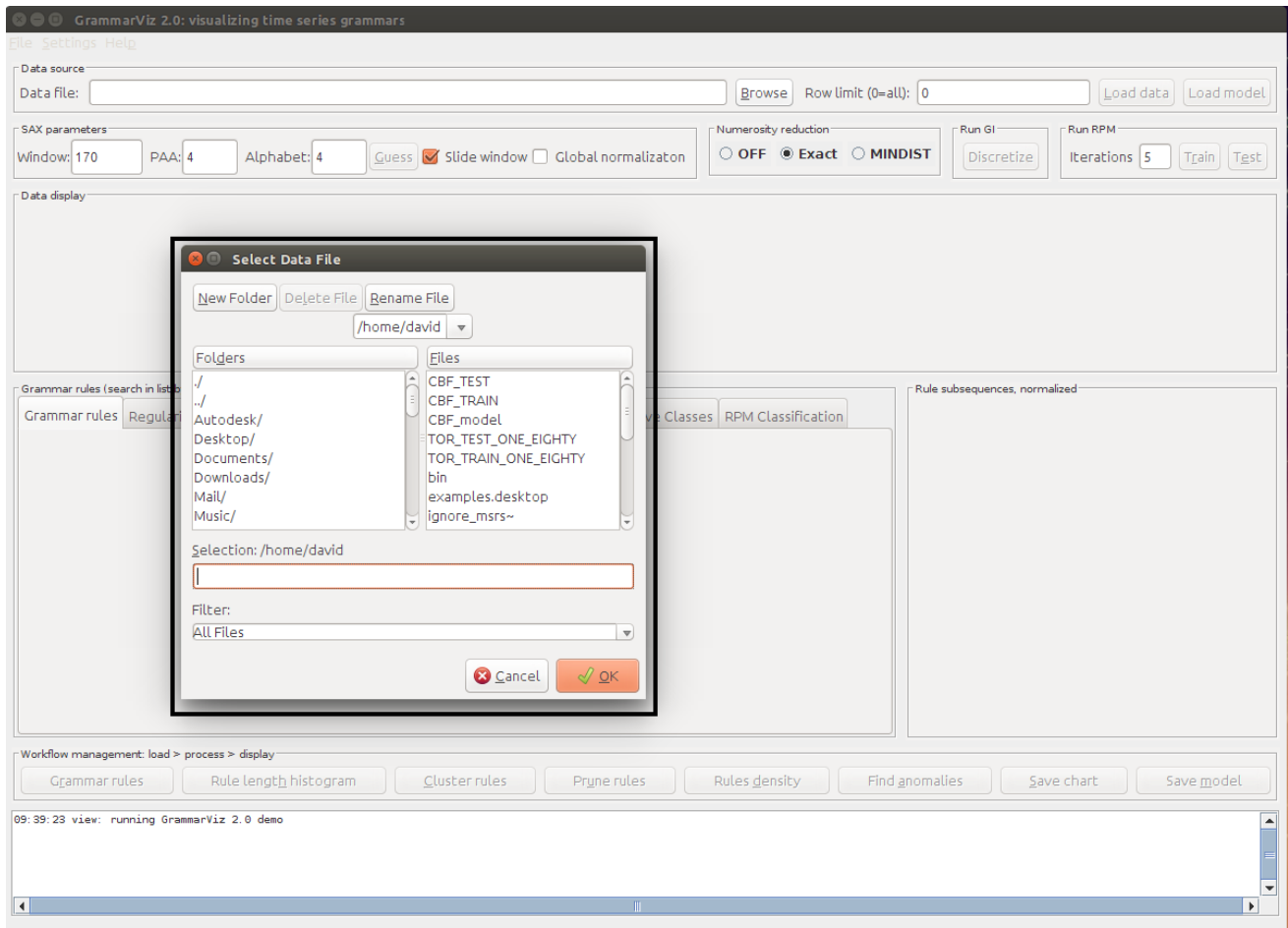


Figure 3: Open the file browser prompt

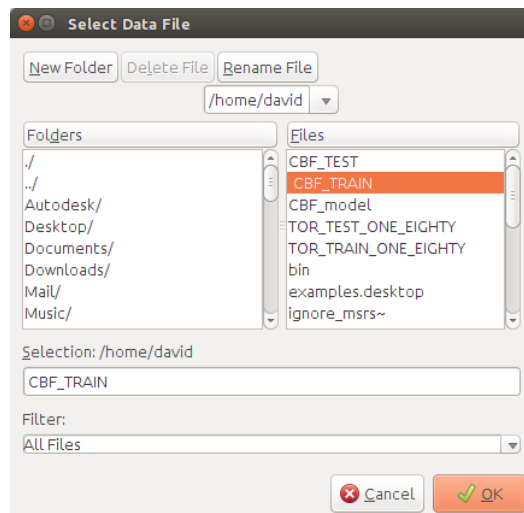


Figure 4: Browser prompt

**Step 3** After selecting the file press the button labeled “Load Data” and TSAT will load the data and the graphs will be populated, and if the data is found to be RPM compatible data then the “Train” button should become available. The text field labeled “Row Limit” allows the user to limit the number of rows that are read in from file, for example if the file contains 100 rows the user could limit it to the first 50.

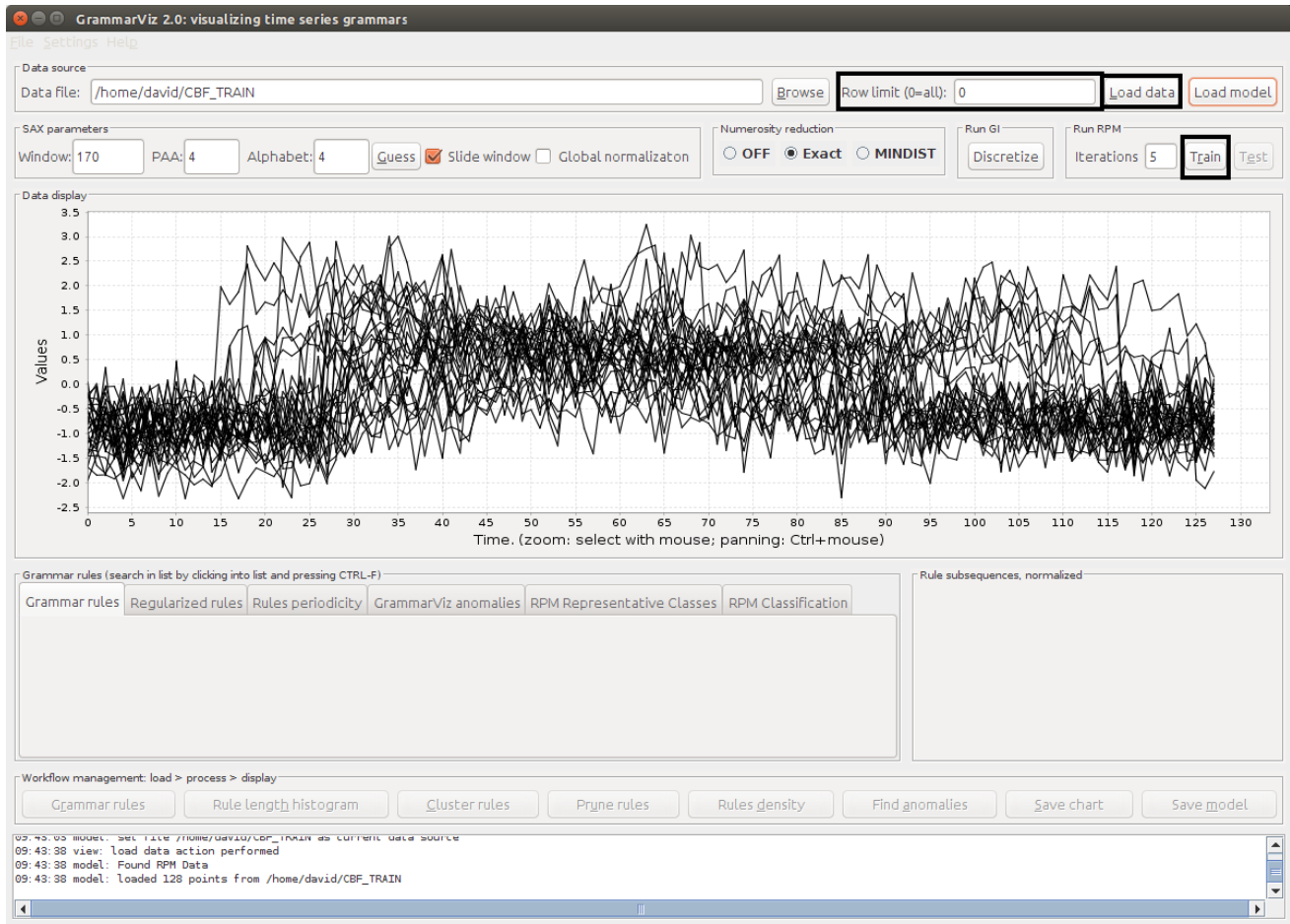


Figure 5: Loaded data

Hitting this button will begin the training phase of RPM, this can take some time depending on the data and the number of iterations RPM will run. The text field labeled “Iterations” sets the maximum number of iterations RPM will go, this prevents RPM from running for too long trying to refine the model. Once the training is complete the tab “RPM Representative Classes” will become populated with patterns RPM thinks represent the labels given. The fields “Window”, “PAA”, and “Alphabet” will also be populated with the values RPM believes are the best fit for the data to aid in further analysis.

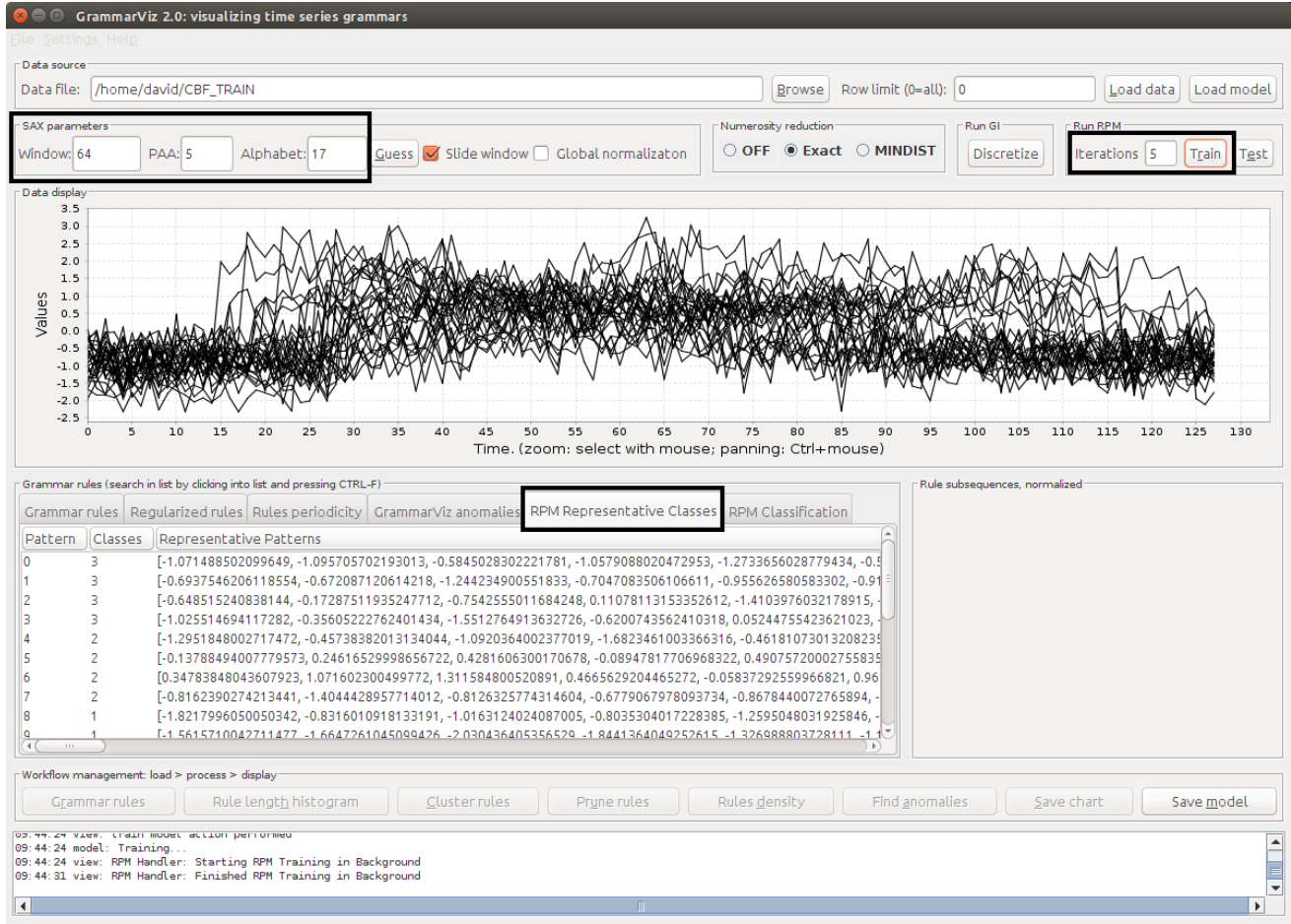


Figure 6: Representative Classes after Training



Selecting the patterns will display their graph on the right hand side of the window, multiple patterns can be selected.

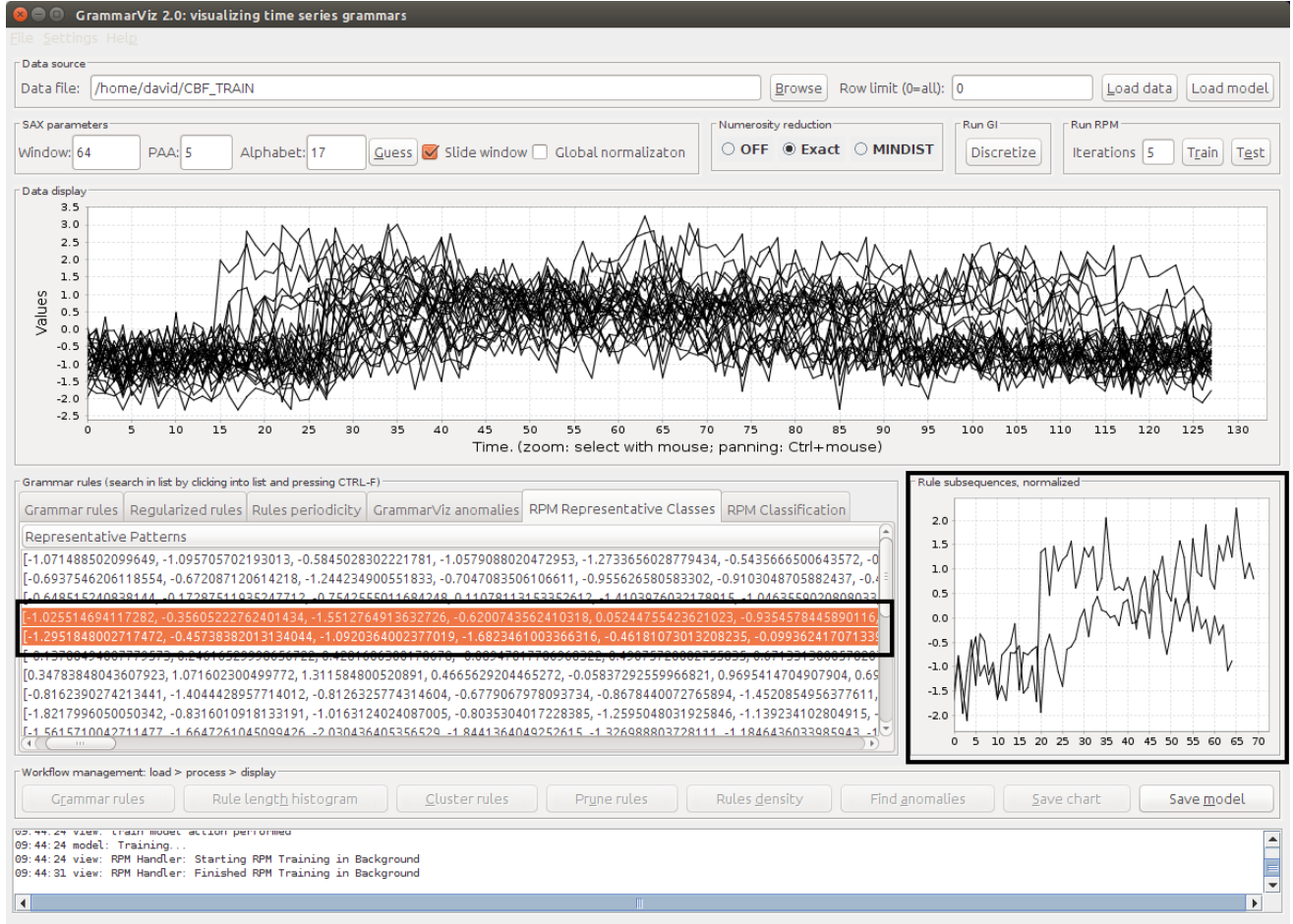


Figure 7: Representative pattern preview

## 2.3 Testing

Once the model has been trained it should be tested for accuracy, this will use a smaller dataset in the RPM compatible format to measure how well the model does.

**Step 1** Click the “Test” button and a file browser prompt will appear, depending on how large the dataset is this may take a moment.

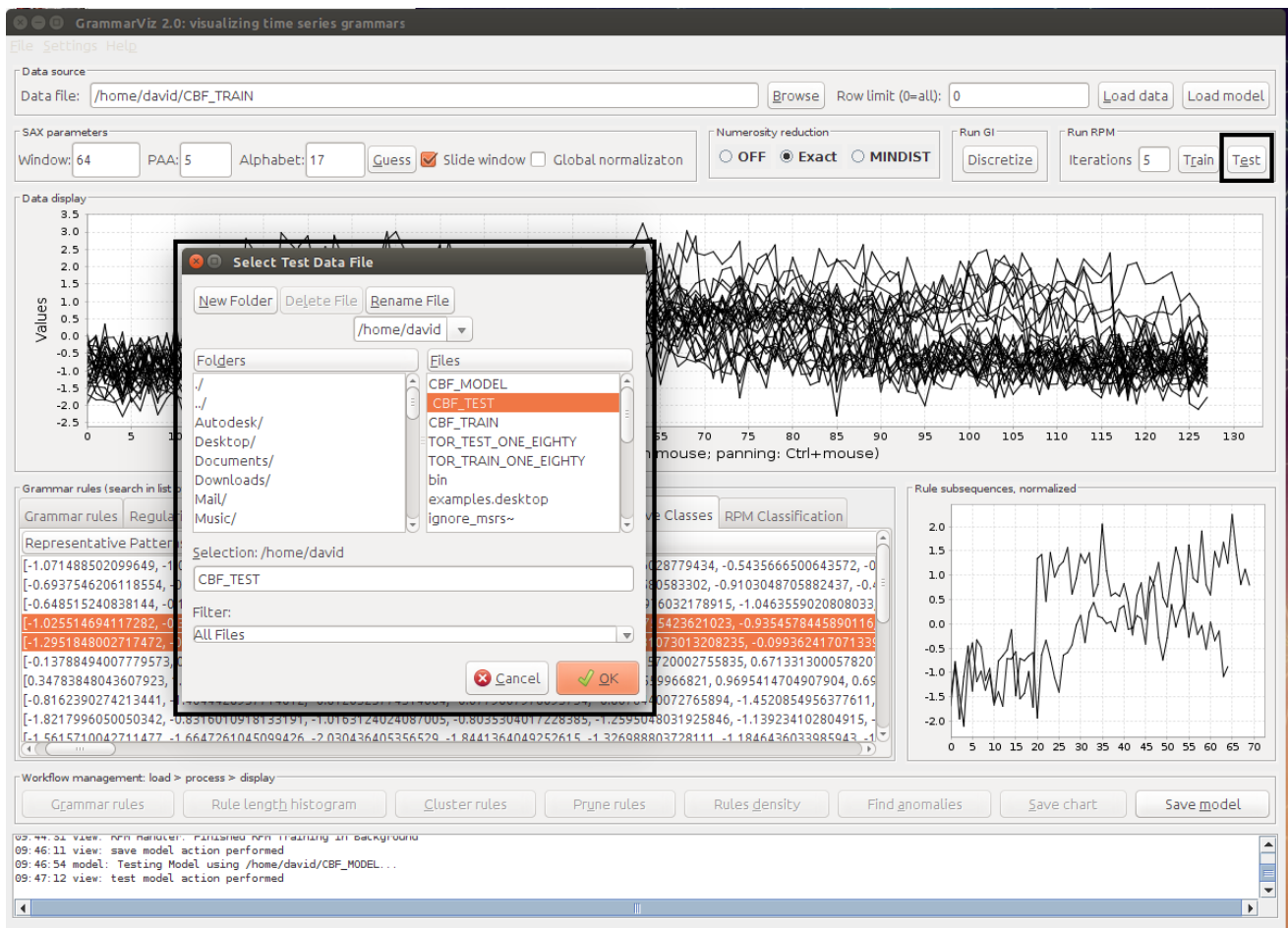


Figure 8: Testing the RPM model

Once the testing is complete the tab labeled “RPM Classification” will be populated. This provides statistics on the effectiveness of the model by reporting the number of samples that were incorrectly labeled by the model.

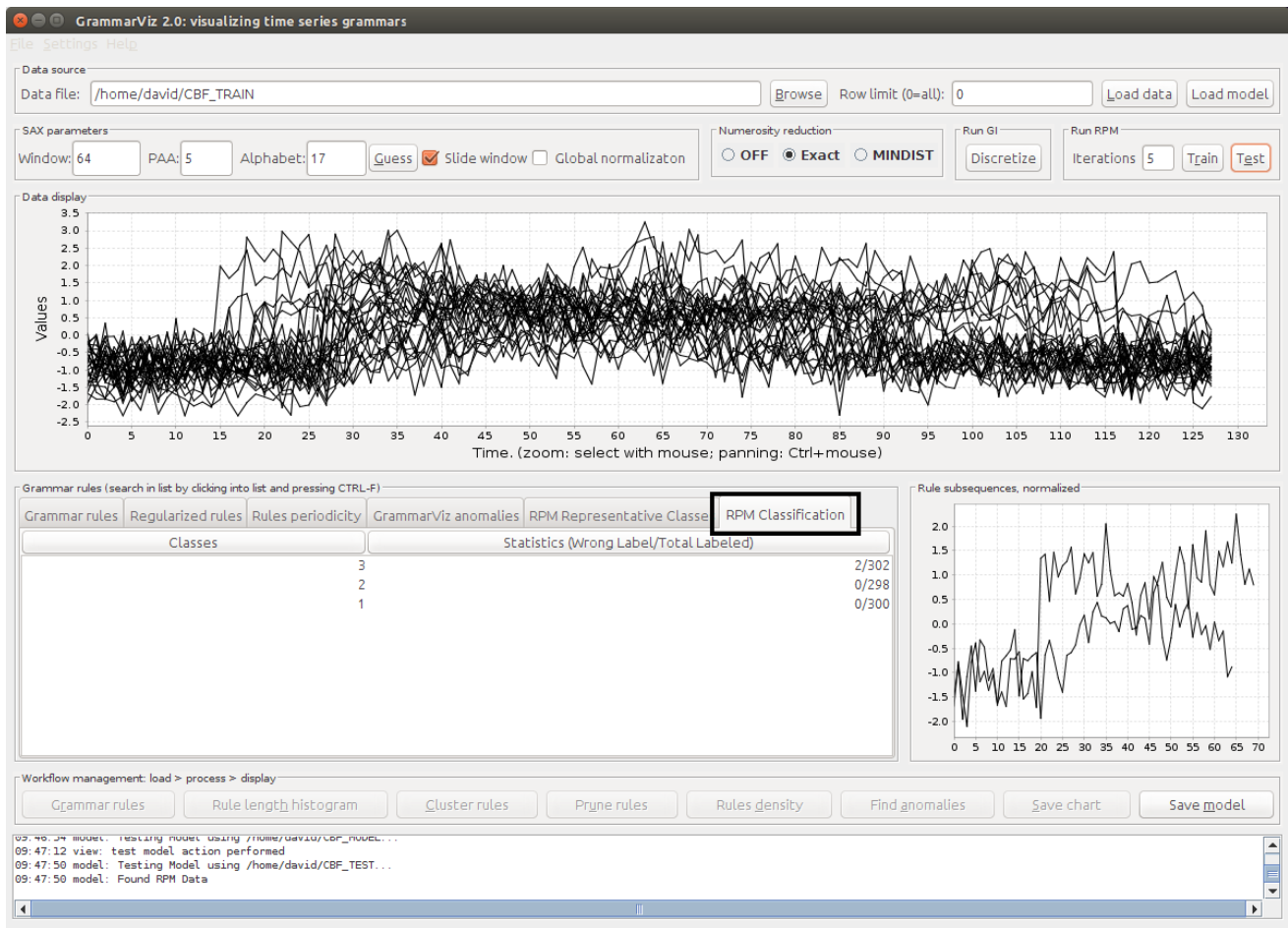


Figure 9: The results from the testing

## 2.4 Testing Unlabeled Data

Using the same method for loading the test data when the data is labeled we can see the results for unlabeled data. Here the test data labels are all question marks so the results will consist of the probability that the test example is in each of the different training classes and the predicted label. For example, in figure 10 the solid box has the label probability for each class and dashed box has the predicted class label.

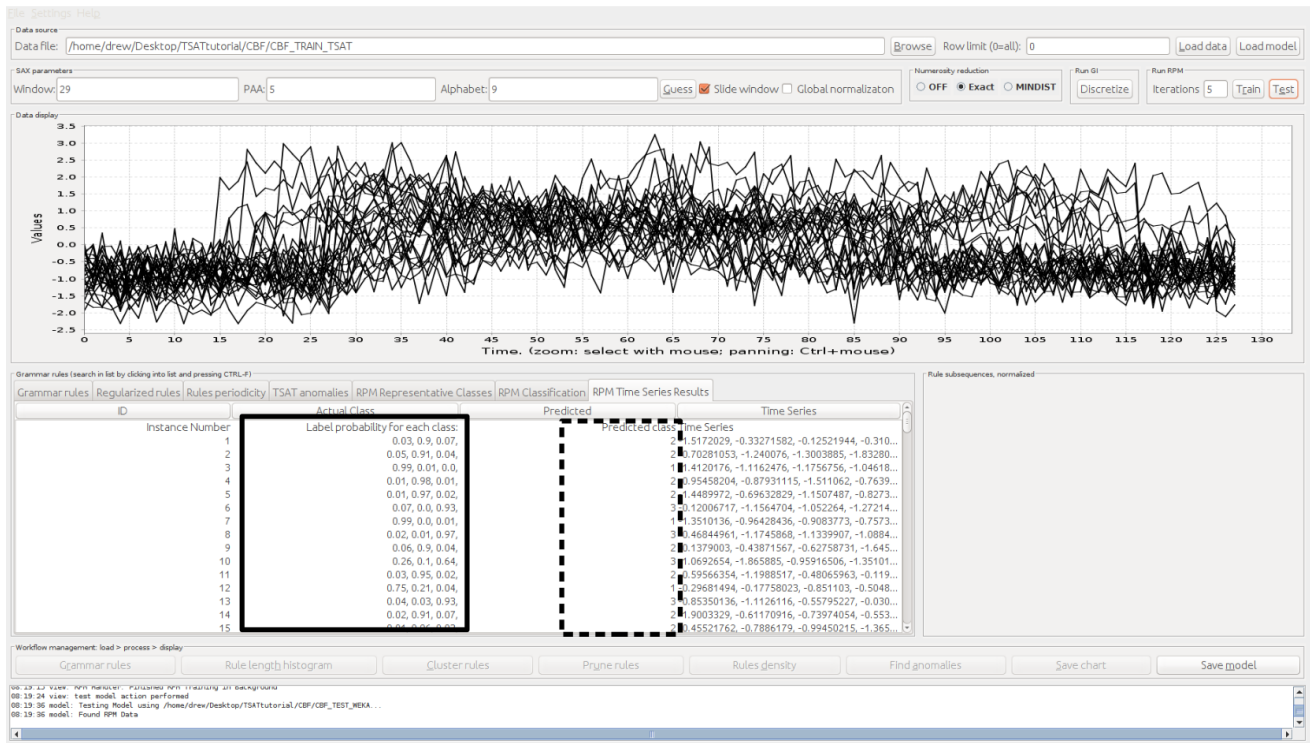


Figure 10: Solid box highlights the probability that the time series was in each of the different class labels and the dashed box highlights the predicted label.

## 2.5 Saving a Trained RPM Model

Creating a model can take some time and there for being able to save the model for later uses is a useful feature. Saving the RPM model will generate a file that can be loaded in later for further testing. One thing to note is that the saved model does not contain the training data however the training data is still needed when doing testing there for a copy of the training data must be retained.

**Step 1** Once a model has been trained up clicking the save model button, as in figure 11, a file browser prompt will appear.

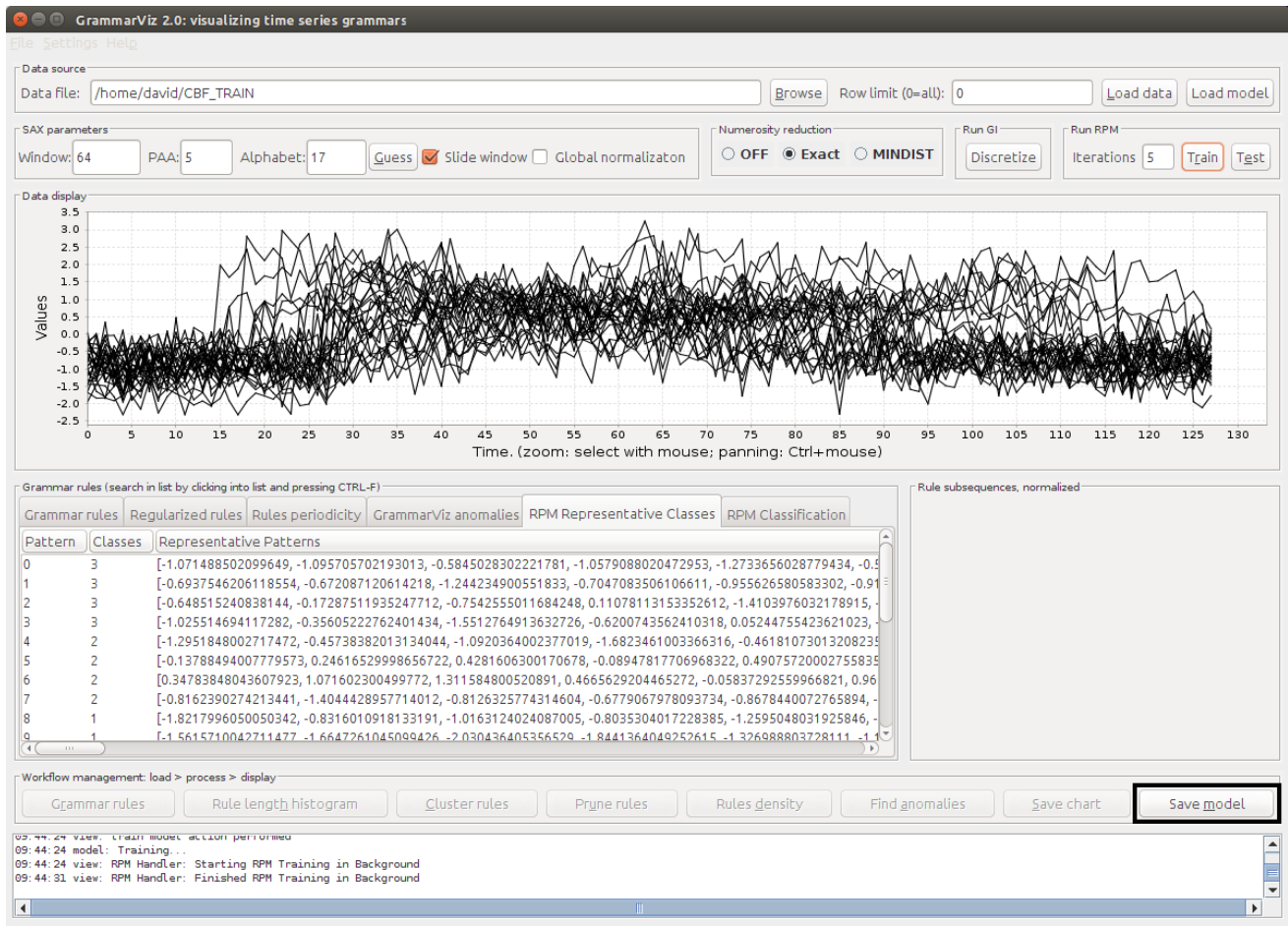


Figure 11: Saving the RPM model



**Step 2** With the file browser prompt select a location to save the model and give it a name, then click the “OK” button to save the model.

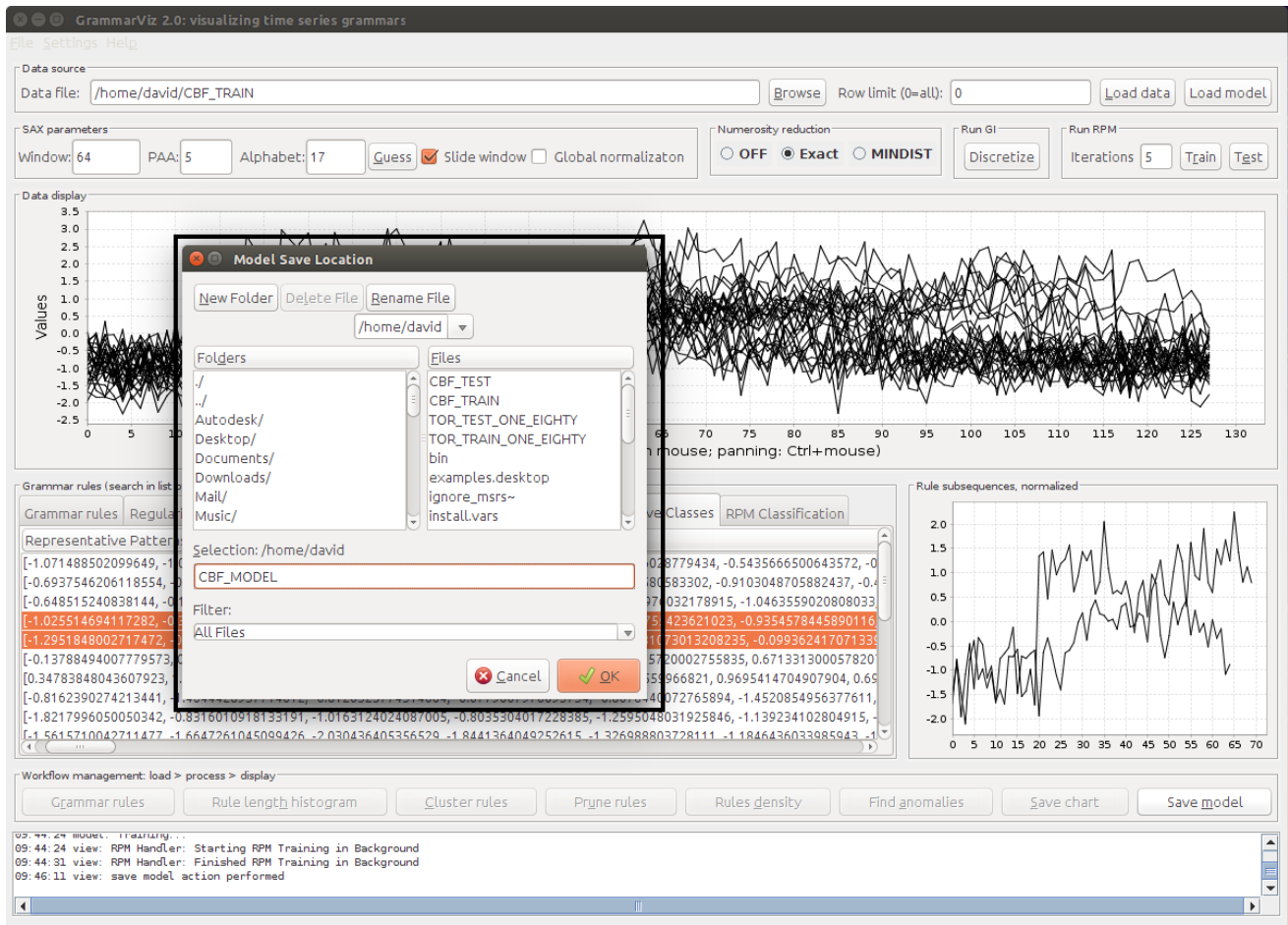


Figure 12: Saving the RPM model to file

## 2.6 Loading an RPM Model

When a model has already been saved, simply loading the will allow for further testing. When loading a model the software will look for the original training data from where it was when it was originally trained. If the data is not there then the software will ask for the location of the data.

**Step 1** First click on the “Browse” button under the “Data Sources” section of the window, as seen in figure 13.

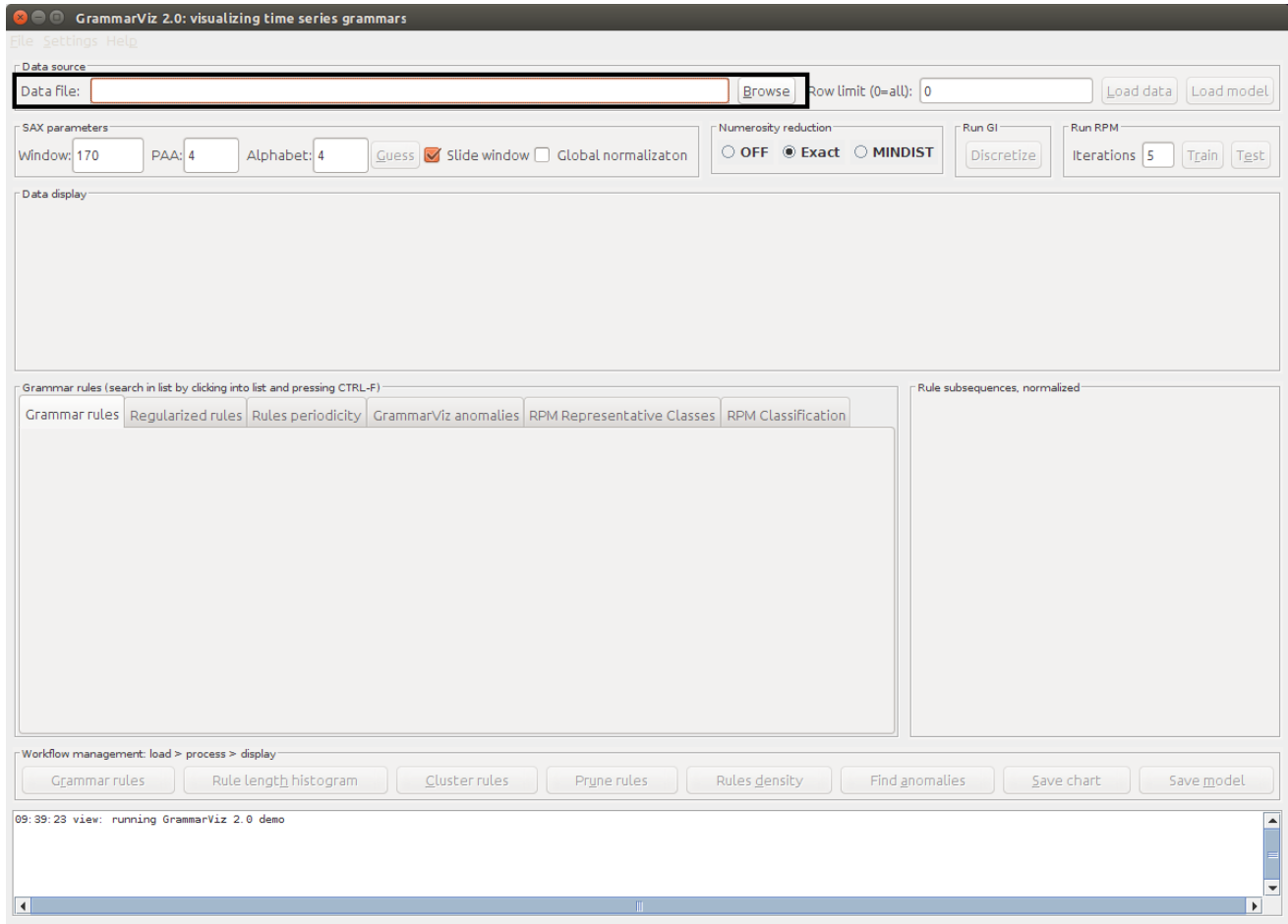


Figure 13: Loading a model

**Step 2** This should bring up the file browser prompt in figure 14. Using this prompt select the previously saved model.

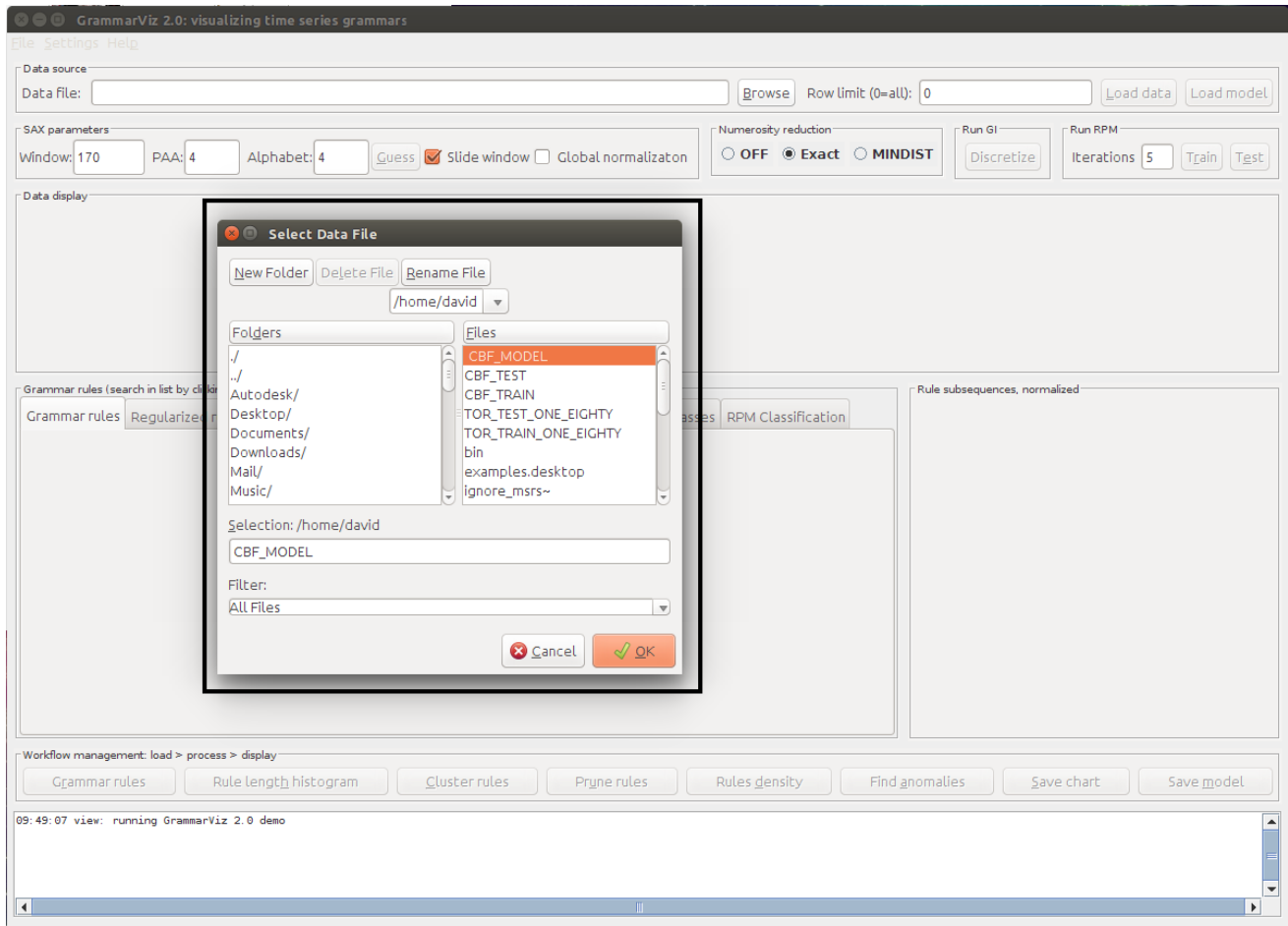


Figure 14: Open the file browser prompt



**Step 3** Once the model has been selected, click the “Load Model” button and the model will be loaded into TSAT. If the data is not found during the loading step TSAT will ask for the location of the data using a file browser prompt, like in figure 16, simply provide the data and the model will finish loading.

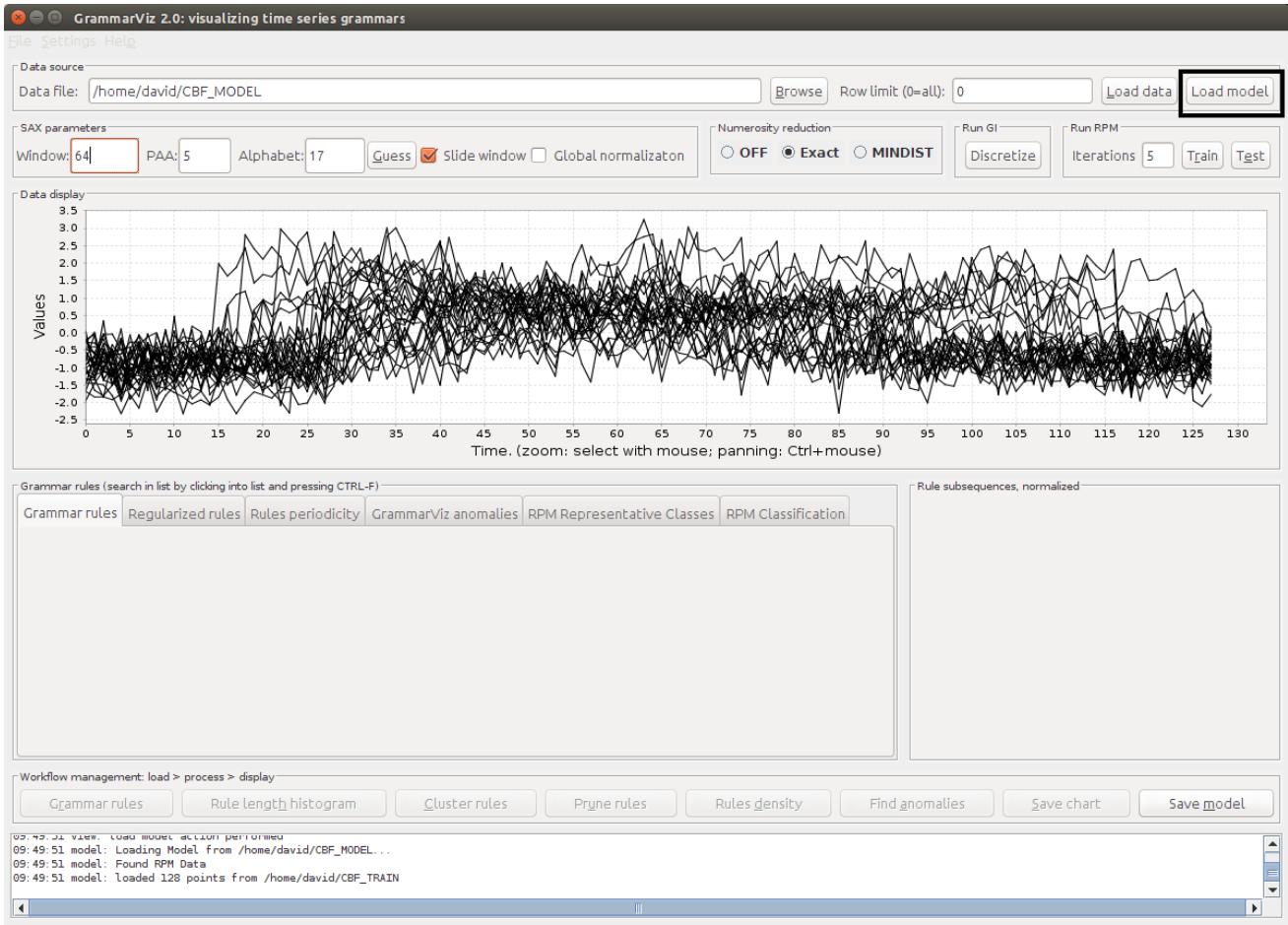


Figure 15: Model loaded

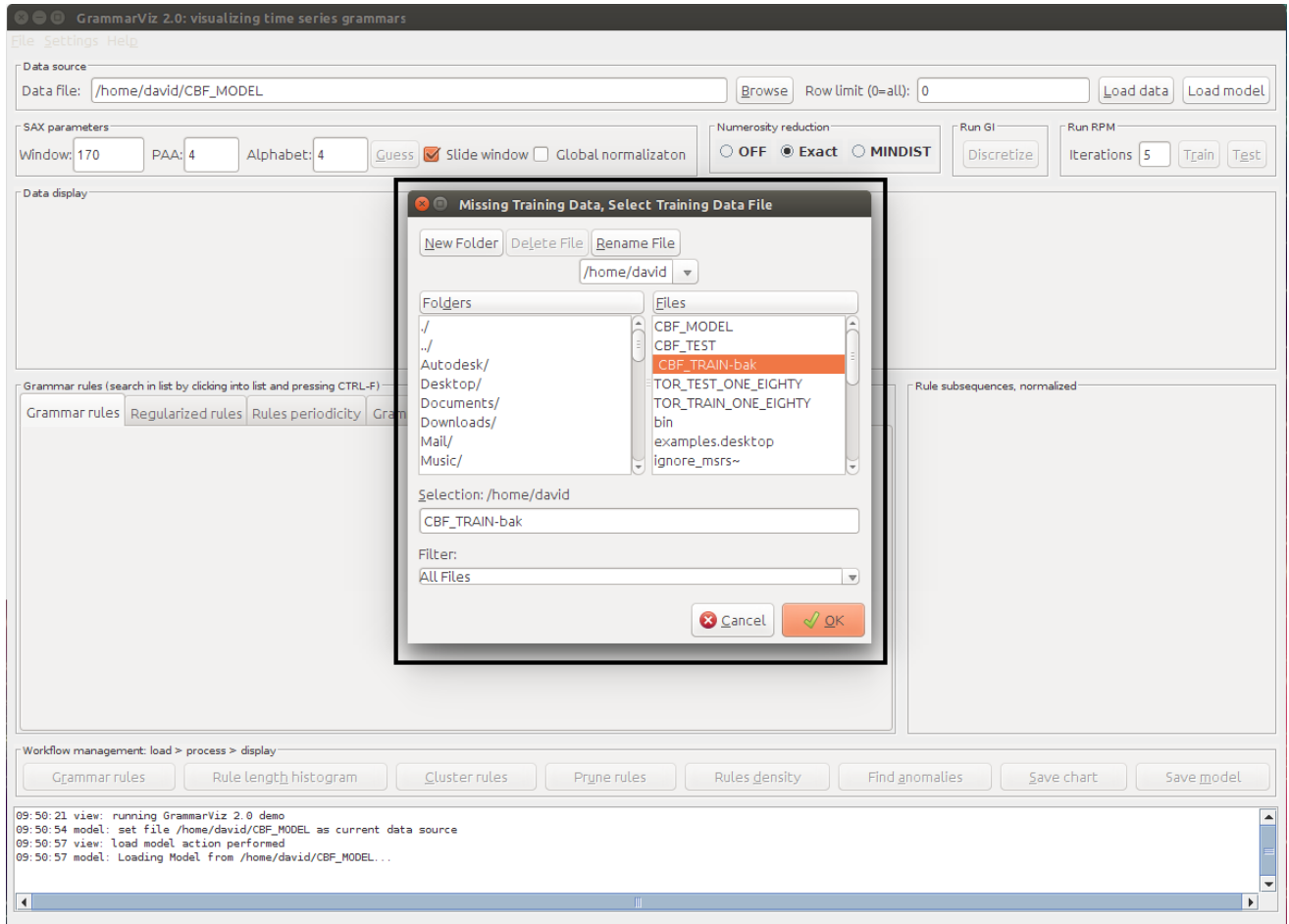


Figure 16: Missing training data file browser prompt

## 2.7 Settings

There are a few options that can be changed when using RPM in TSAT, some of them have already been mentioned and will be covered again.

### 2.7.1 Dynamic Time Warping

Dynamic Time Warping, or DTW, is a method of measuring distance between two time series, this means how similar or different they are to each other. By default RPM uses Euclidean distance which is a simple and fast measurement, however it does not do well when the similar patterns between time series occur at different positions. This is where DTW comes in, it can handle temporal shifts in patterns and, depending on the data, can vastly improve the accuracy of the model. There is a cost however, DTW is a much slower operation and is very expensive to run so it is left as an option for the user.

DTW also has another parameter called “Window” which can have dramatic effects on DTW both in how long it takes to run and its accuracy. The window size basically limits how far DTW will go to try to accurately try to match the two time series. A smaller window will stop DTW from trying to over match them and will take less time to compute. A larger window will take much longer to compute but can allow DTW to match patterns that are farther apart.

Choosing a good window size can be highly dependent on the data and what is being compared, and therefore some experimentation may be needed to find a good window size. There are a few good rules when choosing a window size, for one a window size greater than 10 will usually give bad results so 10 is considered a good starting point. Often for the more common types of data a 3-5 window size can be much better option with significant speed ups. Note DTW's window should not be confused with the Window size in the SAX parameters section of the main window, these are two different and distinct uses of the word window.

**Step 1** To change between Euclidean distance and DTW first open the settings menu:  
“Settings” → “TSAT options” or press Ctrl+p. This will bring up the settings menu in figure 17.

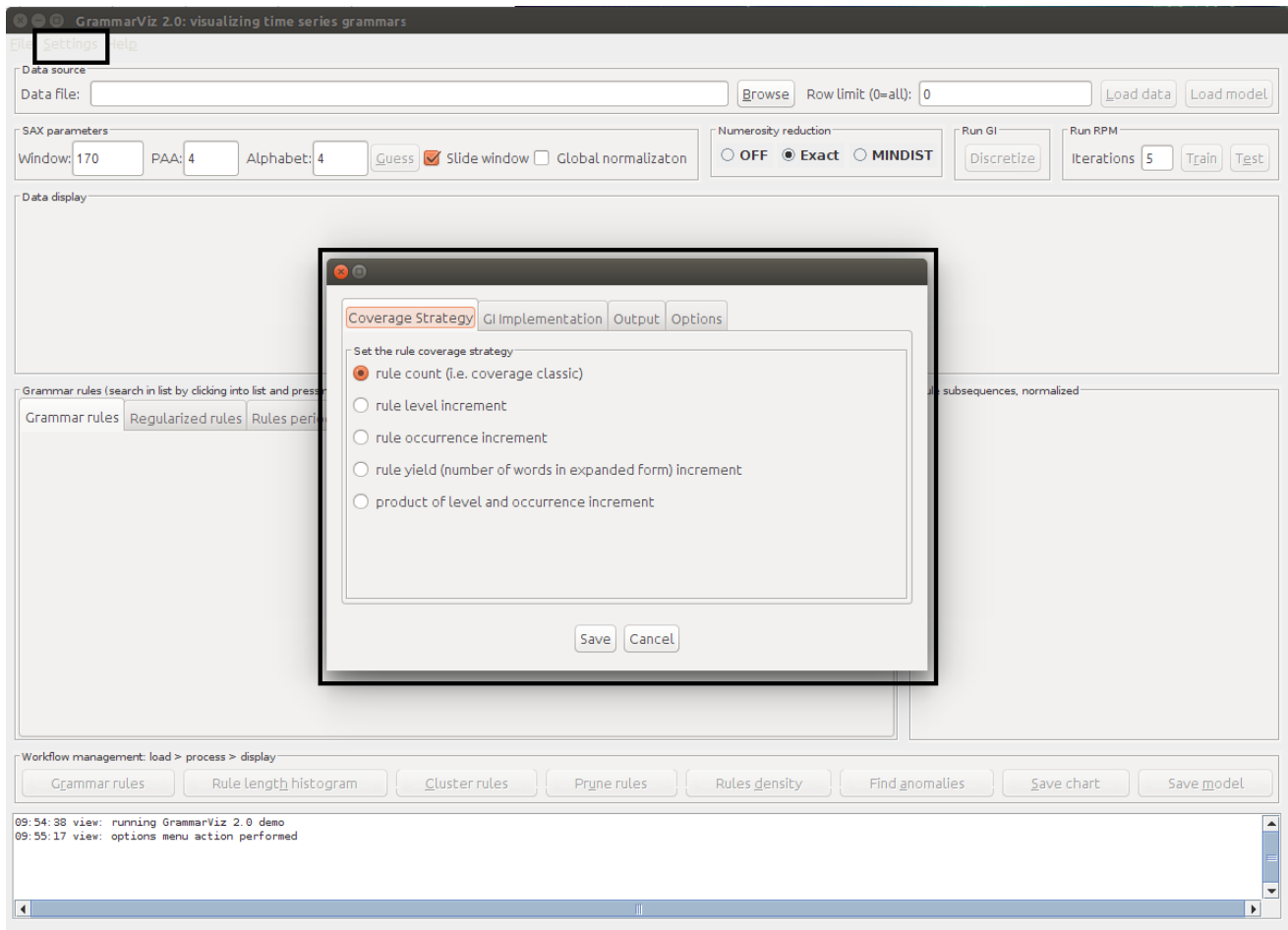


Figure 17: TSAT Settings Dialog

**Step 2** Now click on the “Options” tab.

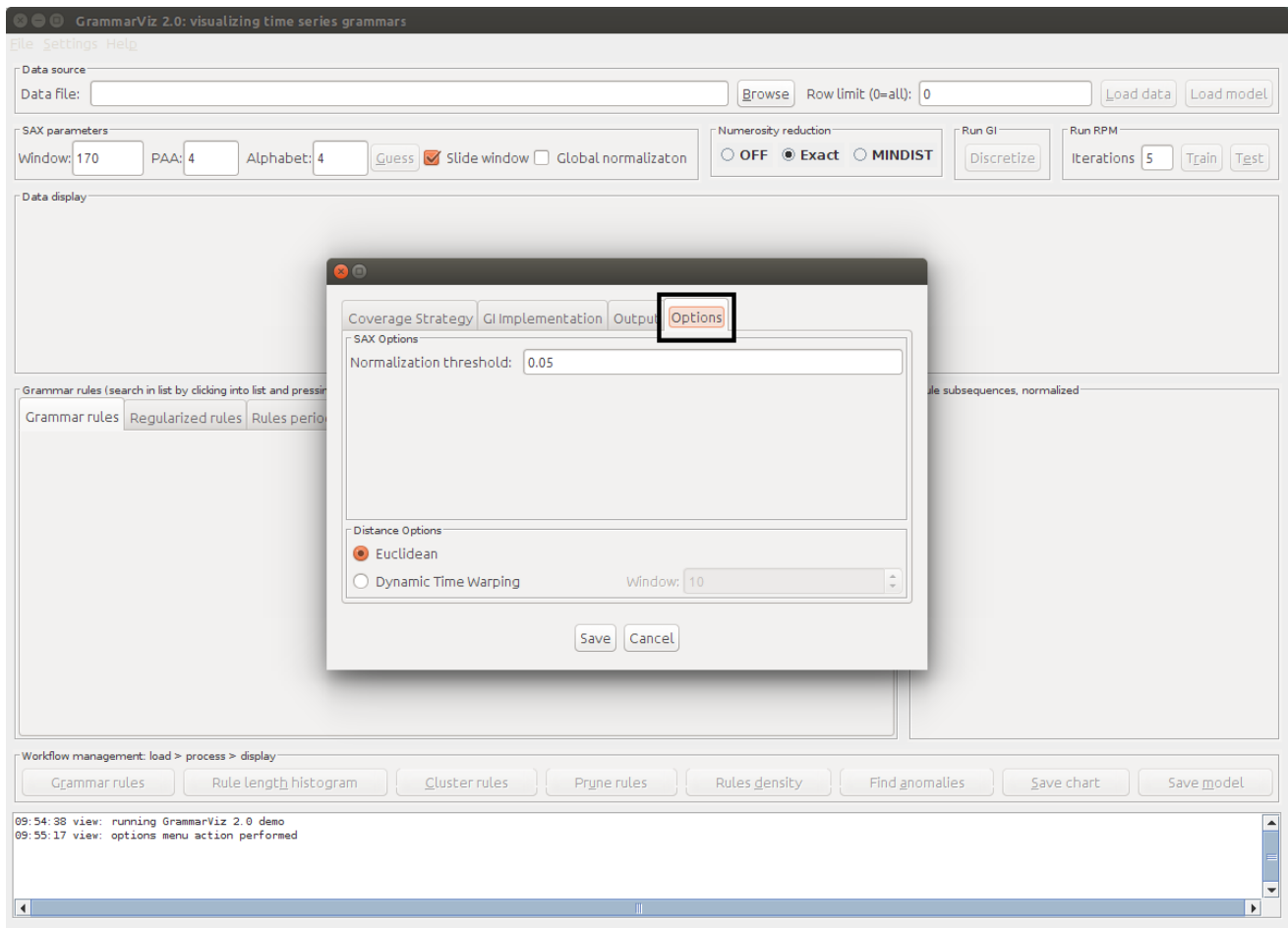


Figure 18: TSAT Settings Dialog Options

**Step 3** Now select the “Dynamic Time Warping” option and the desired “Window” then click save.

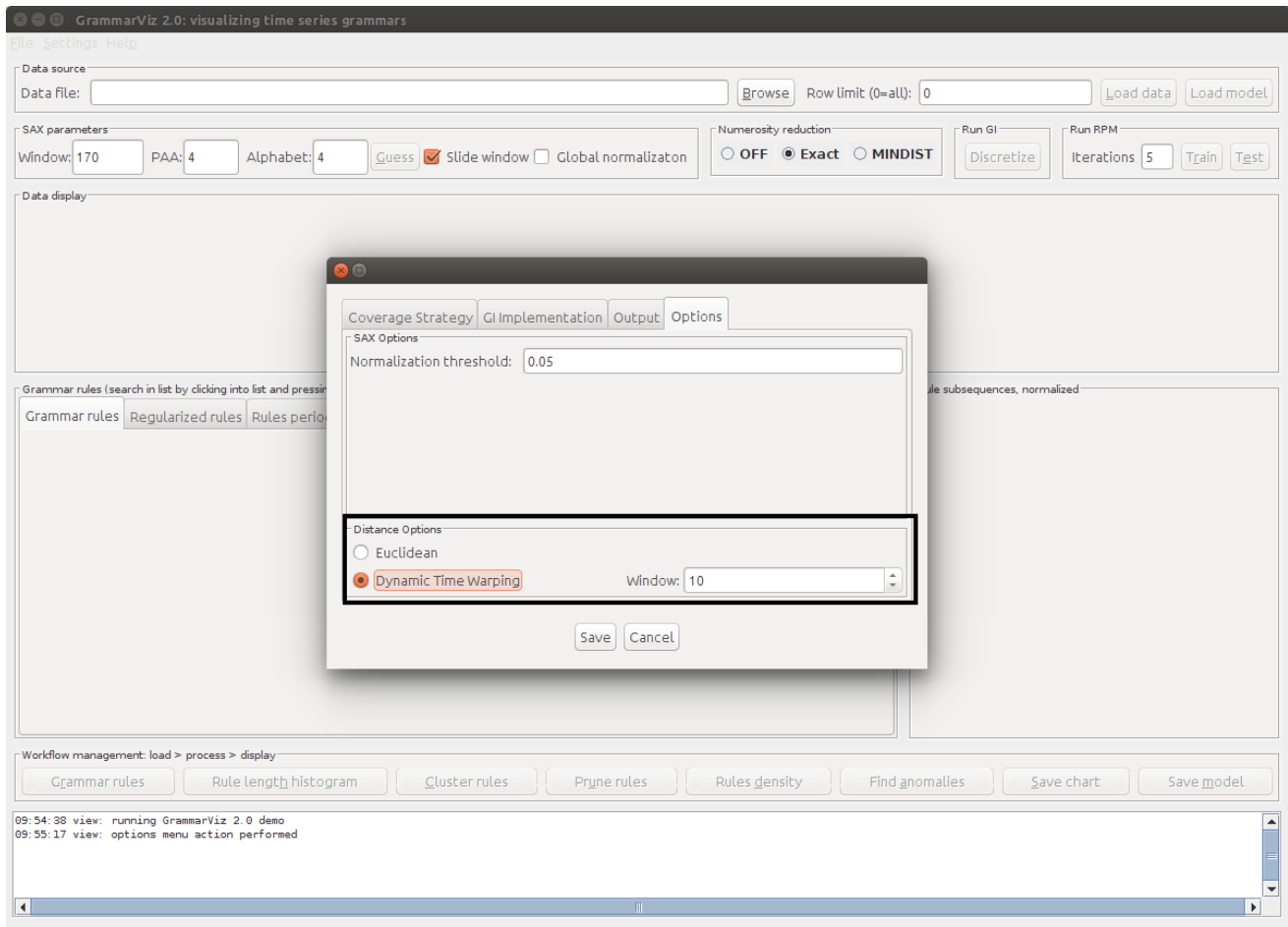


Figure 19: TSAT Settings Dialog Options DTW

## 2.7.2 Iterations

During the operation of RPM it goes through a step that gets repeated many times. This step only stops under two conditions, a minimum threshold is met or if the maximum number of iterations are reached. The iterations setting found under the “Run RPM” section of the main window in TSAT is how the user can control the maximum number of iterations, figure 20. The number of iterations can have an effect on how accurate the model can get, however the more iterations RPM runs through the longer it will take to complete. This becomes a balance between the quality of the model and the how long the training phase will take. It should also be noted that RPM can stop before the maximum number of iterations is met if the model has reached an ideal state. However, this does not mean that all models will or even can reach an ideal state before the maximum number of iterations is reached, indeed some data sets may never return a model that meets the requirements. As RPM runs through the iterations the model should get better but the amount it gets better by can become increasingly insignificant and therefore adding another 10 iterations may not add any significant results to the model. The only way to know if adding more iterations will improve the model is by experimentation which would involve training multiple times, increasing the maximum number of iterations every run until the testing results return no significant improvements.

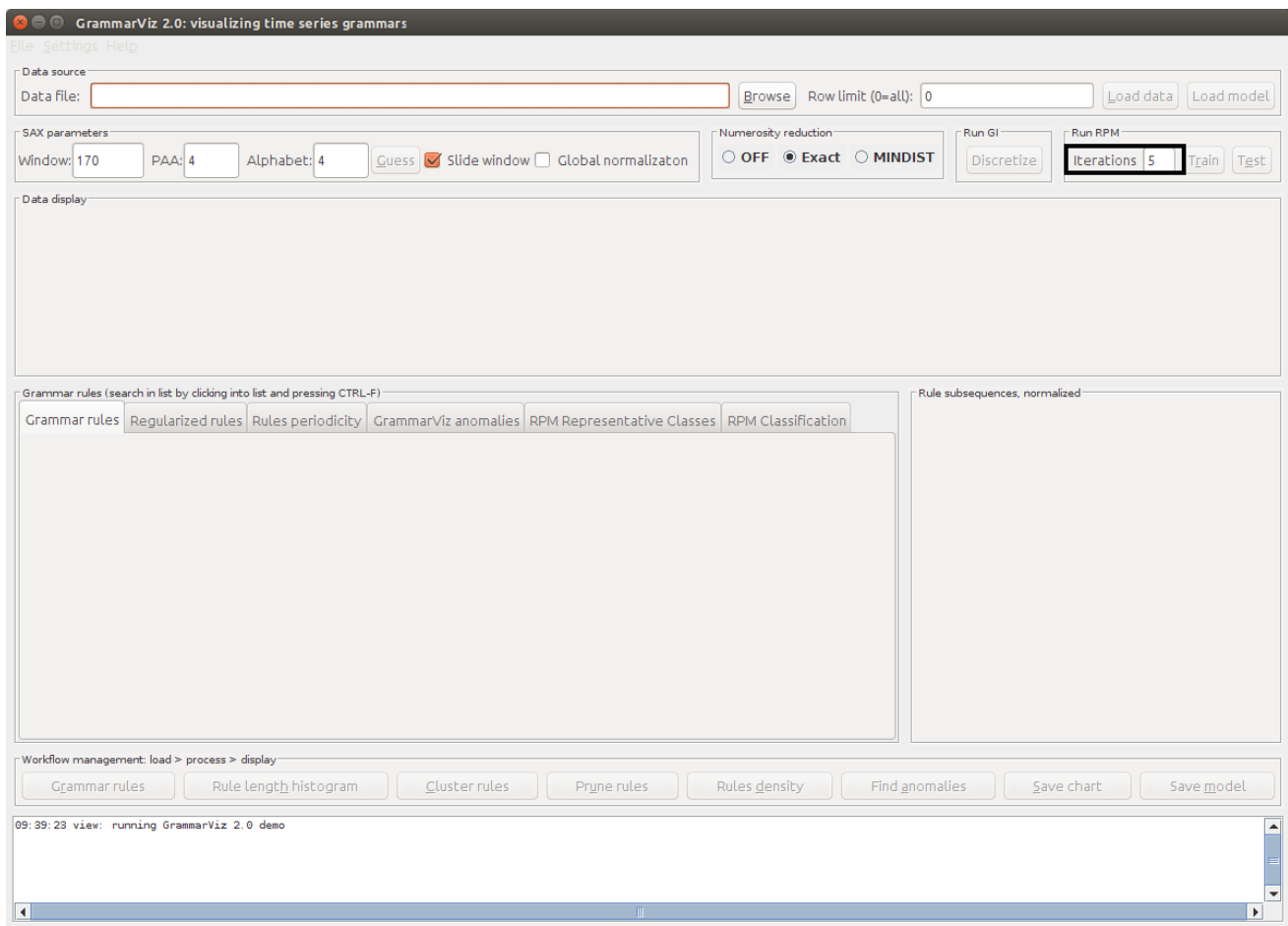


Figure 20: TSAT RPM Iteration Setting

## 3 Motif Discovery

### 3.1 File format

### 3.2 Guide to Motif Discovery

## 4 Anomaly Detection

### 4.1 File format

### 4.2 Guide to Anomaly Detection

## 5 Error Messages

### 5.1 File Errors

### 5.2 RPM Errors

## 6 FAQs

When training there must always be more than one example from each class label and there must be more than one label.

**Installation** This tutorial assumes that you are running Ubuntu 16.04 with Java 1.8 or greater installed.

```
git clone https://github.com/dwicke/TSAT.git
cd TSAT
mvn package -Psingle
```

This will create `tsat-0.0.1-SNAPSHOT-jar-with-dependencies.jar` in the target directory. You can execute the jar and run the GUI by double clicking on it after changing its permissions:

```
chmod +x tsat-0.0.1-SNAPSHOT-jar-with-dependencies.jar
```

To run the GUI from a shell you can do:

```
$ java -Xmx2g -jar target/tsat-0.0.1-SNAPSHOT-jar-with-dependencies.jar
```

The `-Xmx2g` allocates max of 2Gb of memory for the software.