# Data Preprocessing Applying Filters Example

# WEKA — Feature selection

When a database contains a large number of attributes, there will be several attributes which do not become significant in the analysis that you are currently seeking.

Thus, removing the unwanted attributes from the dataset becomes an important task in developing a good machine learning model.
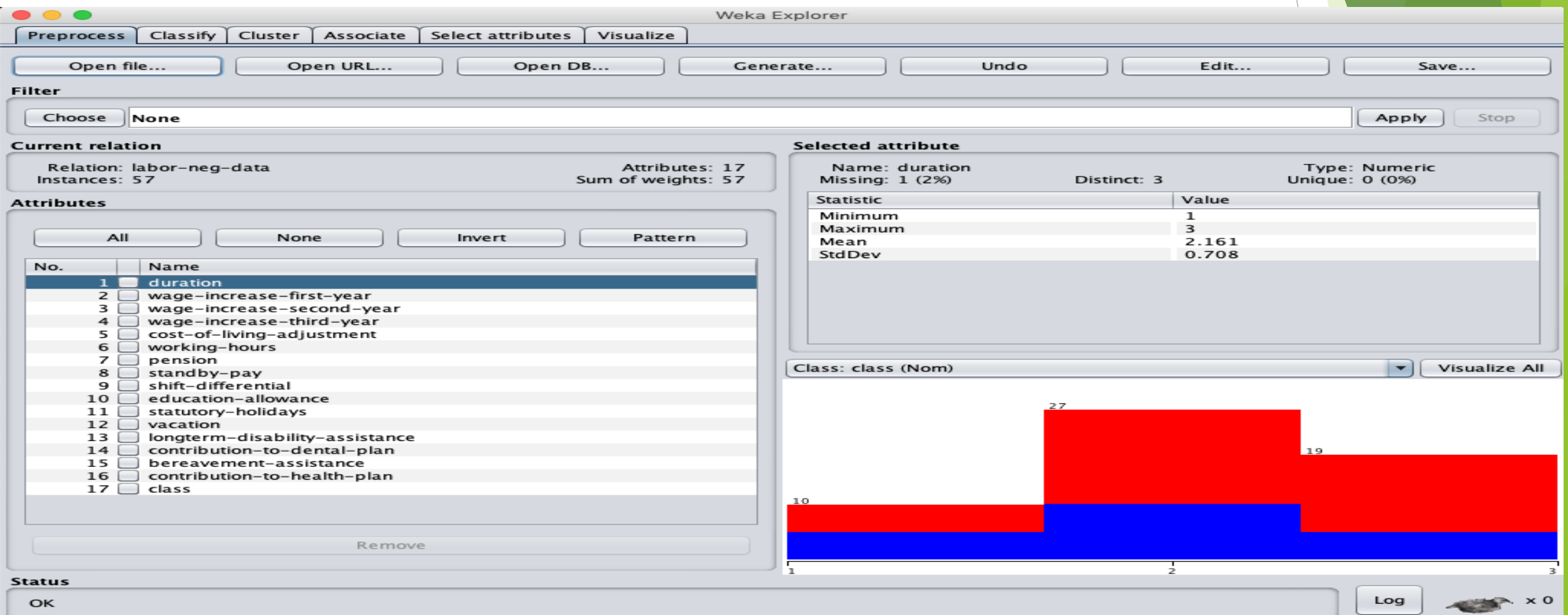
# WEKA — Feature selection

You may examine the entire dataset visually and decide on the irrelevant attributes. This could be a huge task for databases containing a large number of attributes like the supermarket case that you saw in an earlier discussion.

WEKA provides an automated tool for feature selection. To demonstrate this feature on a database containing a large number of attributes.

In the Preprocess tag of the WEKA explorer, select the labor.arff file for loading into the system.

# WEKA – Feature selection

When you load the data, you will see the following screen:Notice that there are 17 attributes. Our task is to create a reduced dataset by eliminating some of the attributes which are irrelevant to our analysis.

# WEKA — Feature selection

Click on the **Select attributes** TAB. You will see the following screen:

# WEKA — Feature selection

Under the **Attribute Evaluator** and **Search Method**, you will find several options. We will just use the defaults here. In the **Attribute Selection Mode**, use full training set option. Click on the **Start** button to process the dataset. You will see the following output:
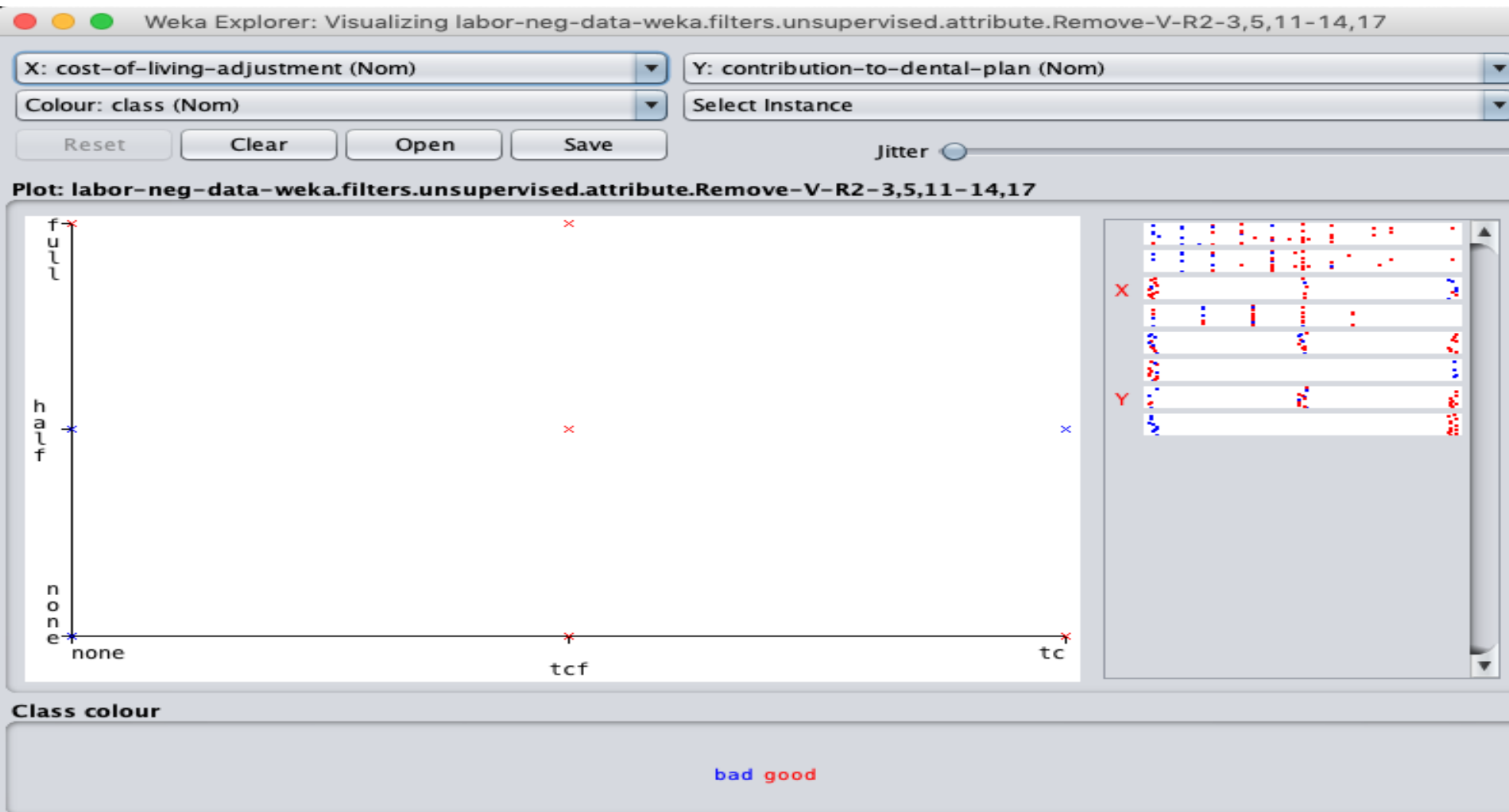
# WEKA — Feature selection

At the bottom of the result window, you will get the list of **Selected** attributes. To get the visual representation, right click on the result in the **Result** list.
The output is shown in the following screenshot: This is similar to the ones we have seen in the earlier chapters.

# WEKA — Feature selection

Clicking on any of the squares will give you the data plot for your further analysis. A typical data plot is shown below: This is similar to the ones we have seen in the earlier slides in previous lab sesssions. Play around with the different options available to analyze the results.
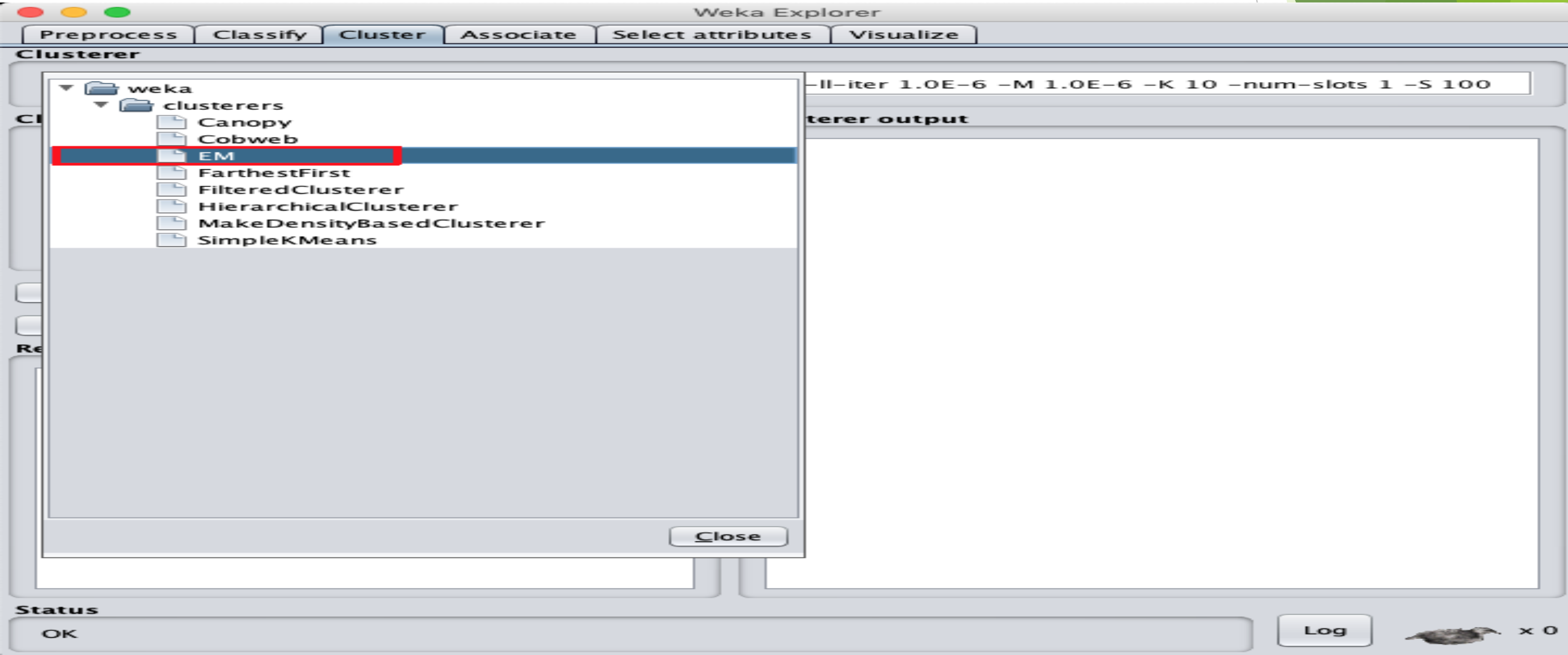
# WEKA — Clustering

You can observe that there are 150 instances and 5 attributes. The names of attributes are listed as **sepallength**, **sepalwidth**, **petallength**, **petalwidth** and **class**.

The first four attributes are of numeric type while the **class** is a nominal type with 3 distinct values.

Examine each attribute to understand the features of the database. We will not do any preprocessing on this data and straight-away proceed to model building.
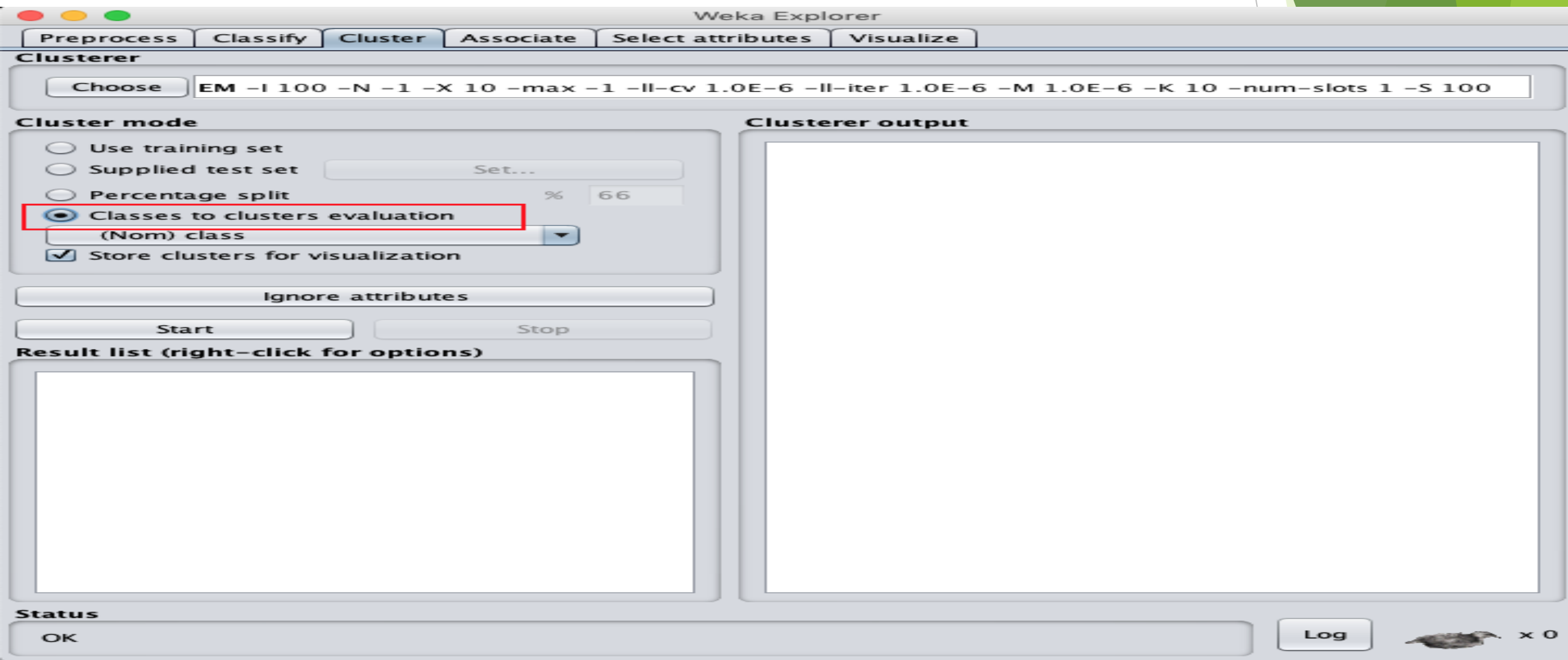
# WEKA — Clustering

**Clustering -** Click on the **Cluster** TAB to apply the clustering algorithms to our loaded data. Click on the **Choose** button. You will see the following screen:
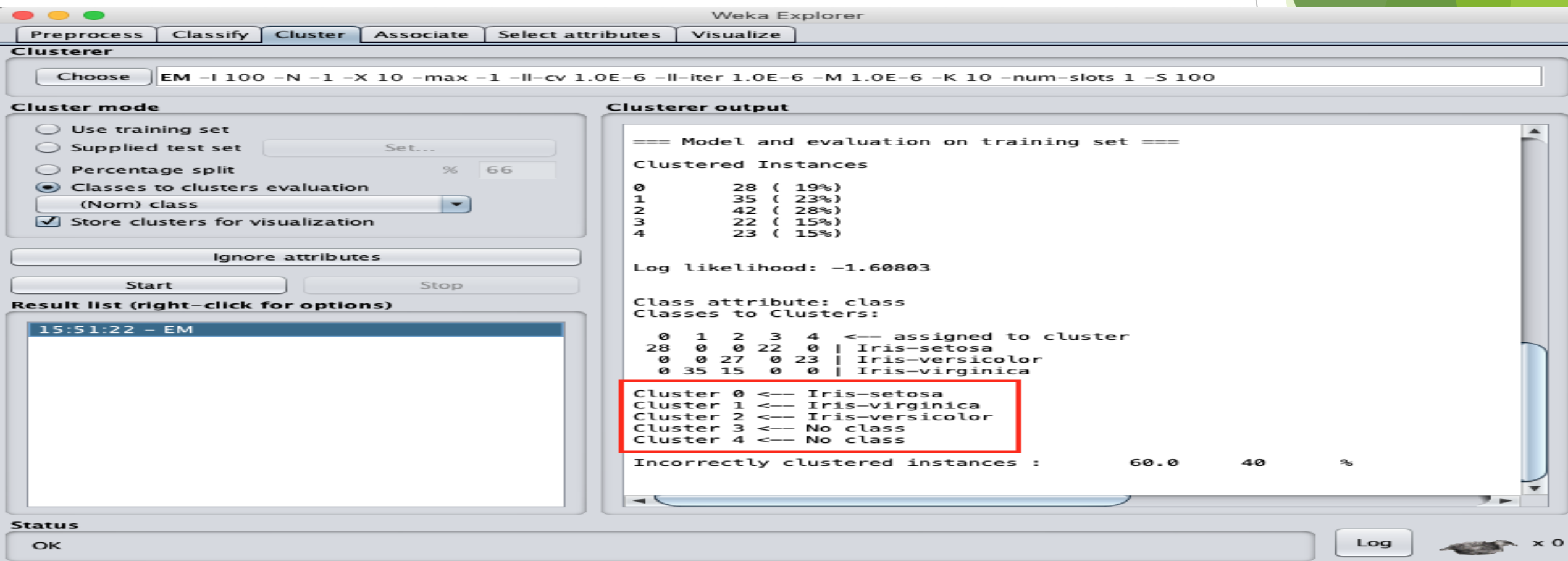
# WEKA — Clustering

Now, select **EM** as the clustering algorithm. In the **Cluster mode** sub window, select the **Classes to clusters evaluation** option as shown in the screenshot below:

# WEKA – Clustering

Click on the **Start** button to process the data. After a while, the results will be presented on the screen. Next, let us study the results. The output of the data processing is shown in the screen below:



From the output screen, you can observe that: ☐ There are 5 clustered instances detected in the database.

☐ The **Cluster 0** represents setosa, **Cluster 1** represents virginica, **Cluster 2** represents versicolor, while the last two clusters do not have any class associated with them.

# WEKA — Clustering

If you scroll up the output window, you will also see some statistics that gives the mean and standard deviation for each of the attributes in the various detected clusters. This is shown in the screenshot given below:



Next, we will look at the visual representation of the clusters.

# WEKA — Clustering

To visualize the clusters, right click on the **EM** result in the **Result list.** You will see the following options:

# WEKA — Clustering

Select **Visualize cluster assignments**. You will see the following output:



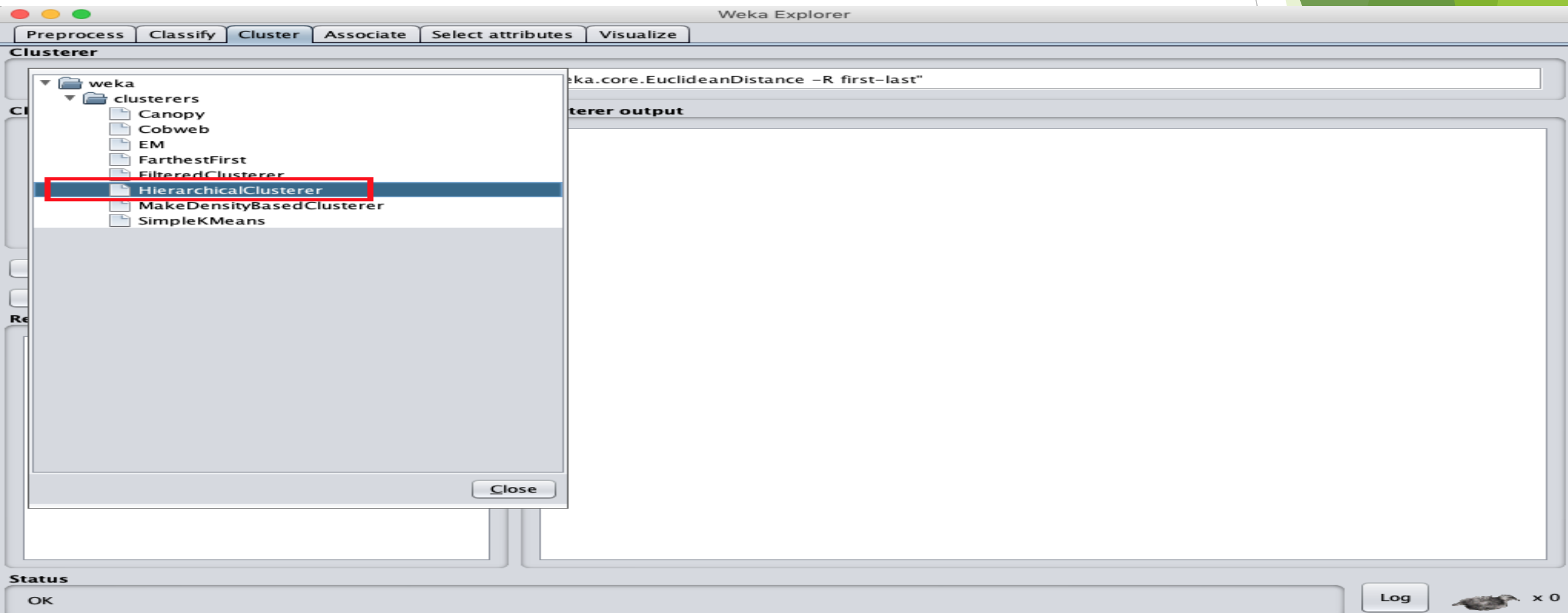As in the case of classification, you will notice the distinction between the correctly and incorrectly identified instances. You can play around by changing the X and Y axes to analyze the results. You may use jittering as in the case of classification to find out the concentration of correctly identified instances. The operations in visualization plot are similar to the one you studied in the case of classification.

# WEKA — Clustering

**Applying Hierarchical Clusterer**

To demonstrate the power of WEKA, let us now look into an application of another clustering algorithm. In the WEKA explorer, select the **HierarchicalClusterer** as your ML algorithm as shown in the screenshot shown below:

# WEKA — Clustering

Choose the **Cluster mode** selection to **Classes to cluster evaluation**, and click on the **Start** button. You will see the following output:



Notice that in the **Result list,** there are two results listed: the first one is the EM result and the second one is the current Hierarchical. Likewise, you can apply multiple ML algorithms to the same dataset and quickly compare their results.

# WEKA — Clustering

If you examine the tree produced by this algorithm, you will see the following output: