

Assignment 1

Classification and Clustering using Weka Explorer

Hafara Firdausi 05111950010040

Description

- Find dataset in UCI machine learning, Kaggle, etc, to be processed.
- Do **classification** using the following methods (supervised learning):
 - Naive Bayes
 - K-Nearest Neighbor (KNN)
 - Support Vector Machine (SVM)
- Do **clustering** using the following methods (unsupervised learning):
 - K-Means

Tools

- Weka Explorer 3.8.2

Dataset Description

- Title :

Forest Covertype data

- Description :

Forest Cover Type (FC) data contains **tree observations** from four wilderness areas located in the **Roosevelt National Forest of northern Colorado**. All observations are cartographic variables (no remote sensing) from 30 meter x 30 meter sections of forest. This dataset includes information on tree type, shadow coverage, distance to nearby landmarks (roads etcetera), soil type, and local topography.

This dataset is part of the **UCI Machine Learning Repository** ([here](#)), but covertype dataset that I use comes from **Kaggle** that can be found [here](#). The original database owners are **Jock A. Blackard, Dr. Denis J. Dean, and Dr. Charles W. Anderson** of the Remote Sensing and GIS Program at Colorado State University.

- Details :

| | |
|----------------------------------|----------------------|
| Data Set Characteristics | Multivariate |
| Attribute Characteristics | Categorical, Integer |
| Associated Tasks | Classification |
| Number of Instances | 581012 |
| Number of Attributes | 54 |
| Missing Values? | No |

| Date Donated | 1998-08-01 |
|--------------|------------|
|--------------|------------|

- Attributes:

| Column Name | Number of Columns |
|------------------------------------|-------------------|
| Elevation | 1 |
| Aspect | 1 |
| Slope | 1 |
| Horizontal_Distance_To_Hydrology | 1 |
| Vertical_Distance_To_Hydrology | 1 |
| Hillshade_9am | 1 |
| Hillshade_Noon | 1 |
| Hillshade_3pm | 1 |
| Horizontal_Distance_To_Fire_Points | 1 |
| Wilderness_Area | 4 |
| Soil_Type | 40 |
| Cover_Type | 1 |

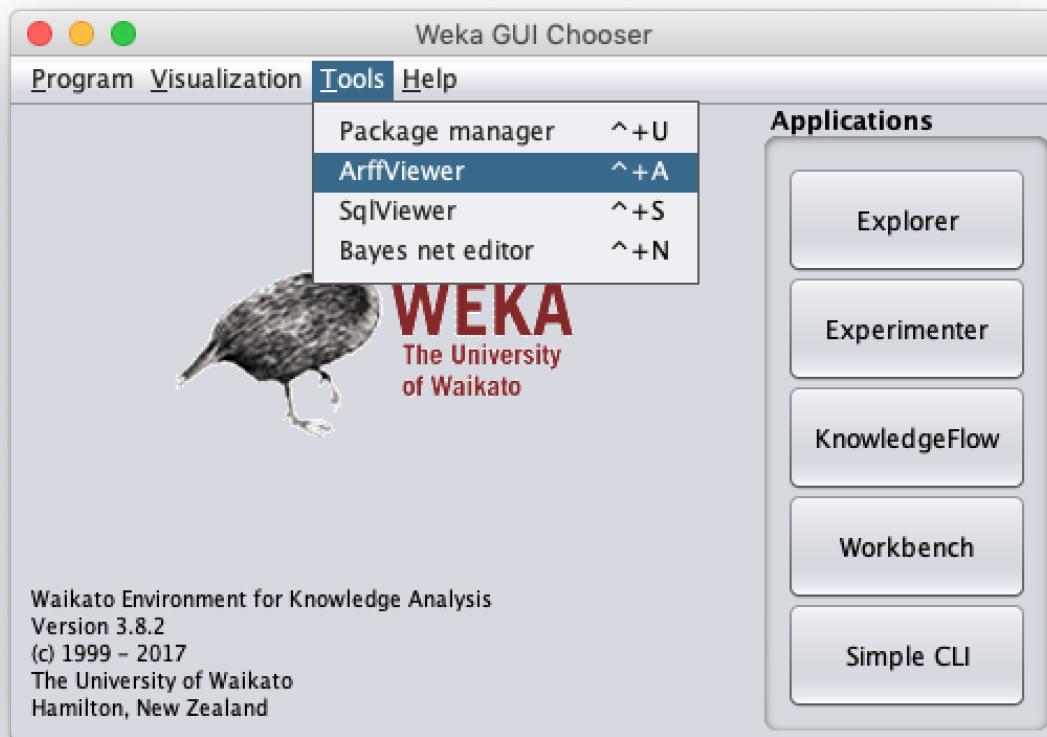
- There are 7 classes (covertype) here, that is:

1. Spruce/Fir
2. Lodgepole Pine
3. Ponderosa Pine
4. Cottonwood/Willow
5. Aspen
6. Douglas-fir
7. Krummholz

Report

Data Preparation

1. Convert dataset file **covertype.csv** into **covertype.arff** using **ArffViewer** (it can be found in the Weka's main menu > Tools > ArffViewer) because Weka Explorer just open data with **.arff** extension (it's like usual **.csv** file with header information).



2. Reduce the number of data become **10000 rows** as a sample because **too much data** cause some processes in Weka to get stuck (source: [here](#) | data reduction process: [here](#)).

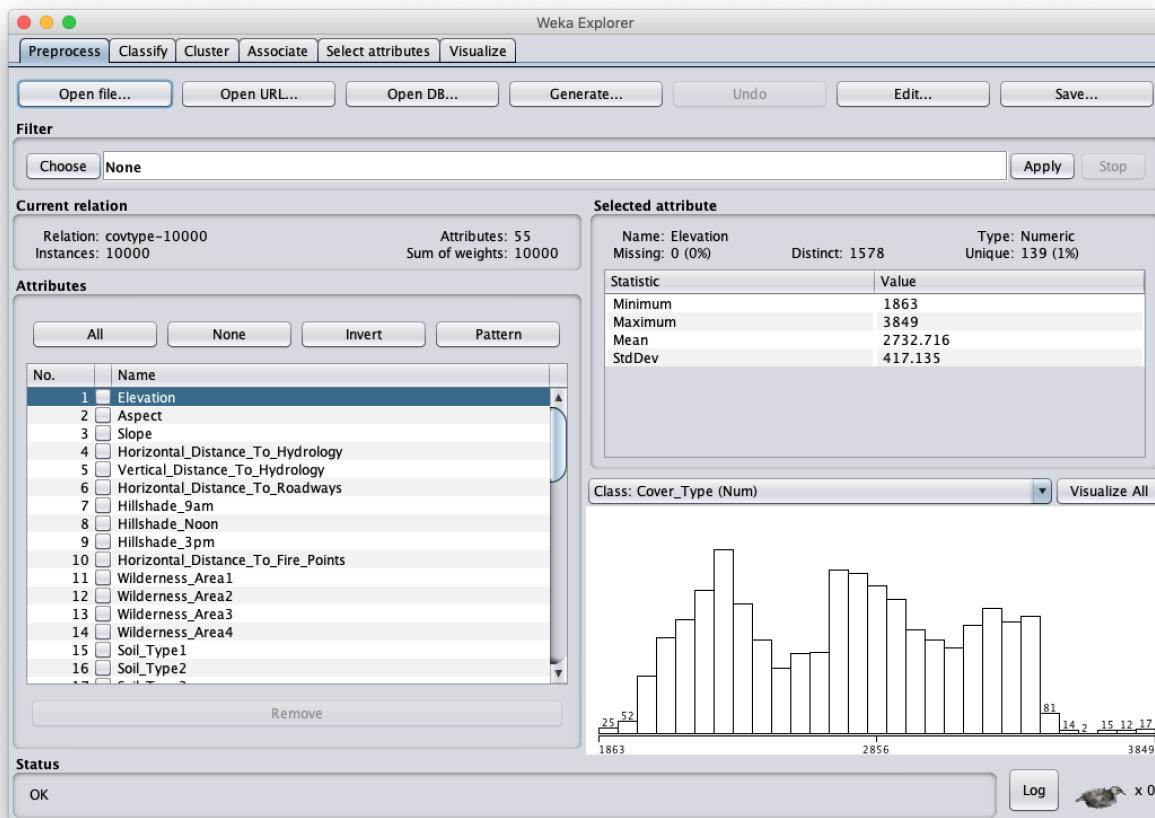
- Here's the **arff** file.

```
@attribute Elevation numeric
@attribute Aspect numeric
@attribute Slope numeric
@attribute Horizontal_Distance_To_Hydrology numeric
@attribute Vertical_Distance_To_Hydrology numeric
@attribute Horizontal_Distance_To_Roadways numeric
@attribute Hillshade_9am numeric
@attribute Hillshade_Noon numeric
@attribute Hillshade_3pm numeric
@attribute Horizontal_Distance_To_Fire_Points numeric
@attribute Wilderness_Area1 numeric
@attribute Wilderness_Area2 numeric
@attribute Wilderness_Area3 numeric
@attribute Wilderness_Area4 numeric
@attribute Soil_Type1 numeric
@attribute Soil_Type2 numeric
@attribute Soil_Type3 numeric
@attribute Soil_Type4 numeric
@attribute Soil_Type5 numeric
@attribute Soil_Type6 numeric
@attribute Soil_Type7 numeric
```

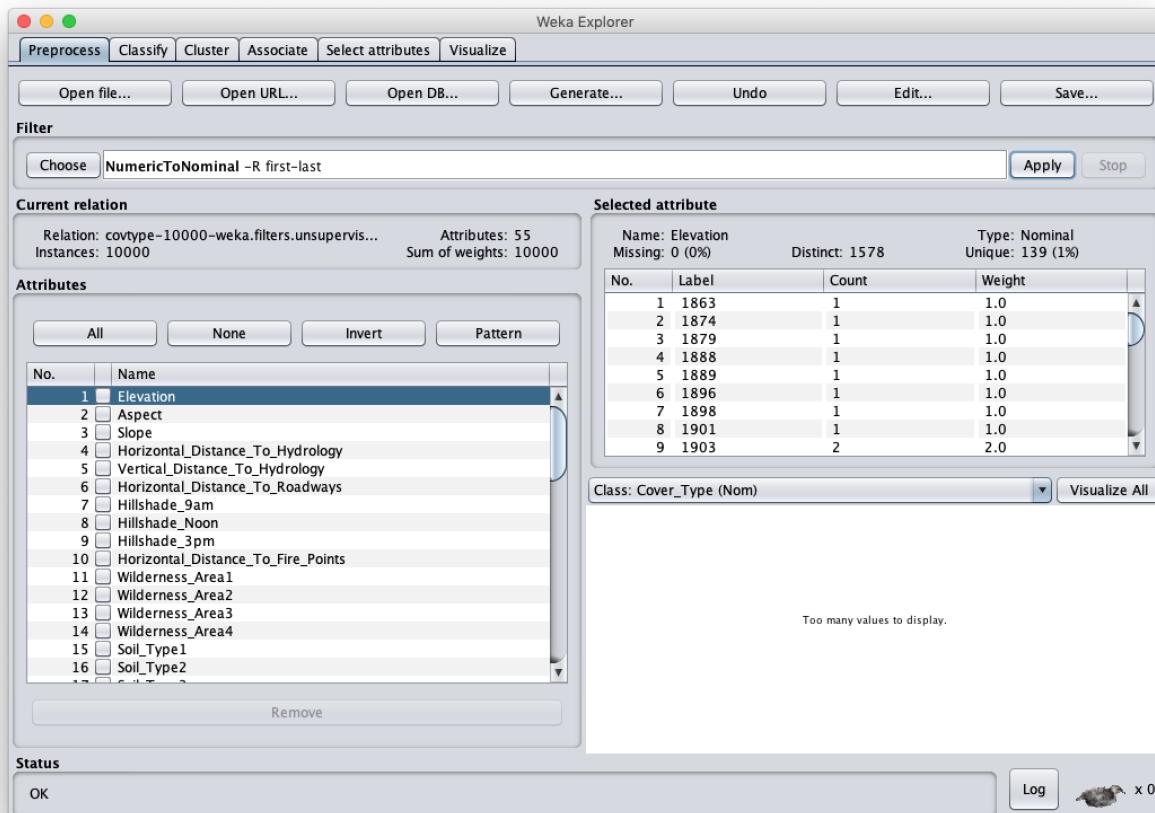
```
@attribute Soil_Type8 numeric  
@attribute Soil_Type9 numeric  
@attribute Soil_Type10 numeric  
@attribute Soil_Type11 numeric  
@attribute Soil_Type12 numeric  
@attribute Soil_Type13 numeric  
@attribute Soil_Type14 numeric  
@attribute Soil_Type15 numeric  
@attribute Soil_Type16 numeric  
@attribute Soil_Type17 numeric  
@attribute Soil_Type18 numeric  
@attribute Soil_Type19 numeric  
@attribute Soil_Type20 numeric  
@attribute Soil_Type21 numeric  
@attribute Soil_Type22 numeric  
@attribute Soil_Type23 numeric  
@attribute Soil_Type24 numeric  
@attribute Soil_Type25 numeric  
@attribute Soil_Type26 numeric  
@attribute Soil_Type27 numeric  
@attribute Soil_Type28 numeric  
@attribute Soil_Type29 numeric  
@attribute Soil_Type30 numeric  
@attribute Soil_Type31 numeric  
@attribute Soil_Type32 numeric  
@attribute Soil_Type33 numeric  
@attribute Soil_Type34 numeric  
@attribute Soil_Type35 numeric  
@attribute Soil_Type36 numeric  
@attribute Soil_Type37 numeric  
@attribute Soil_Type38 numeric  
@attribute Soil_Type39 numeric  
@attribute Soil_Type40 numeric  
@attribute Cover_Type numeric
```

Data Preprocessing

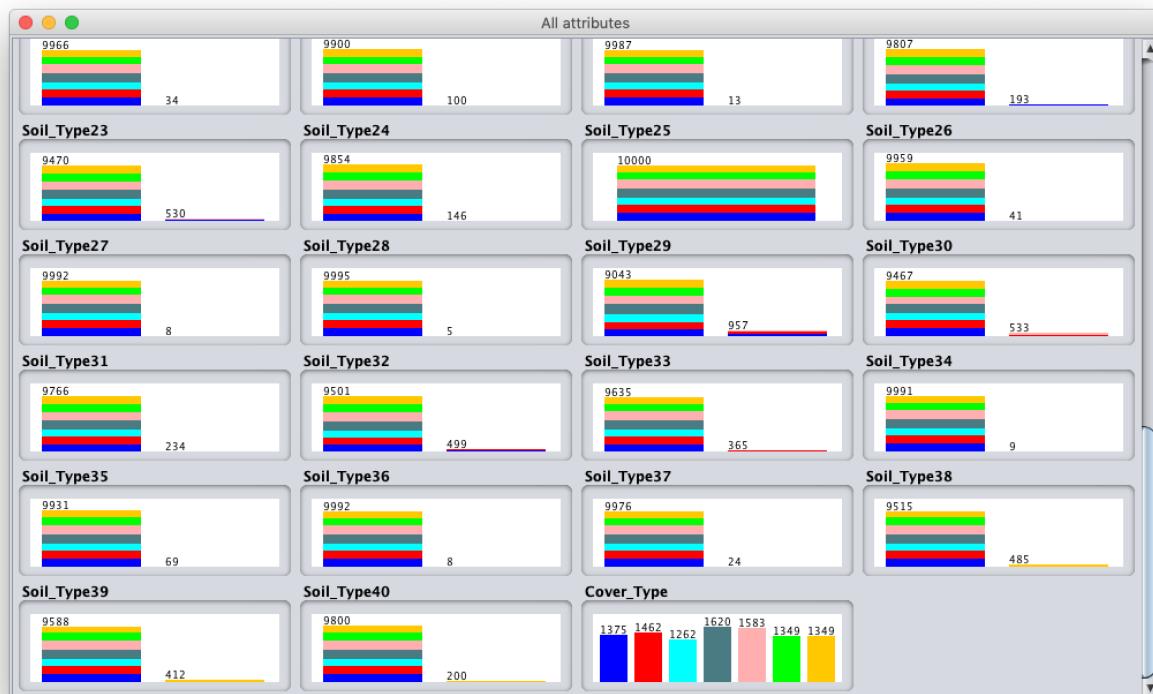
1. Open the dataset with **.arff** extension in the Weka.



2. Some algorithms must use data with **Nominal** type. So, we must preprocess the data first by using filter **NumericToNominal**.

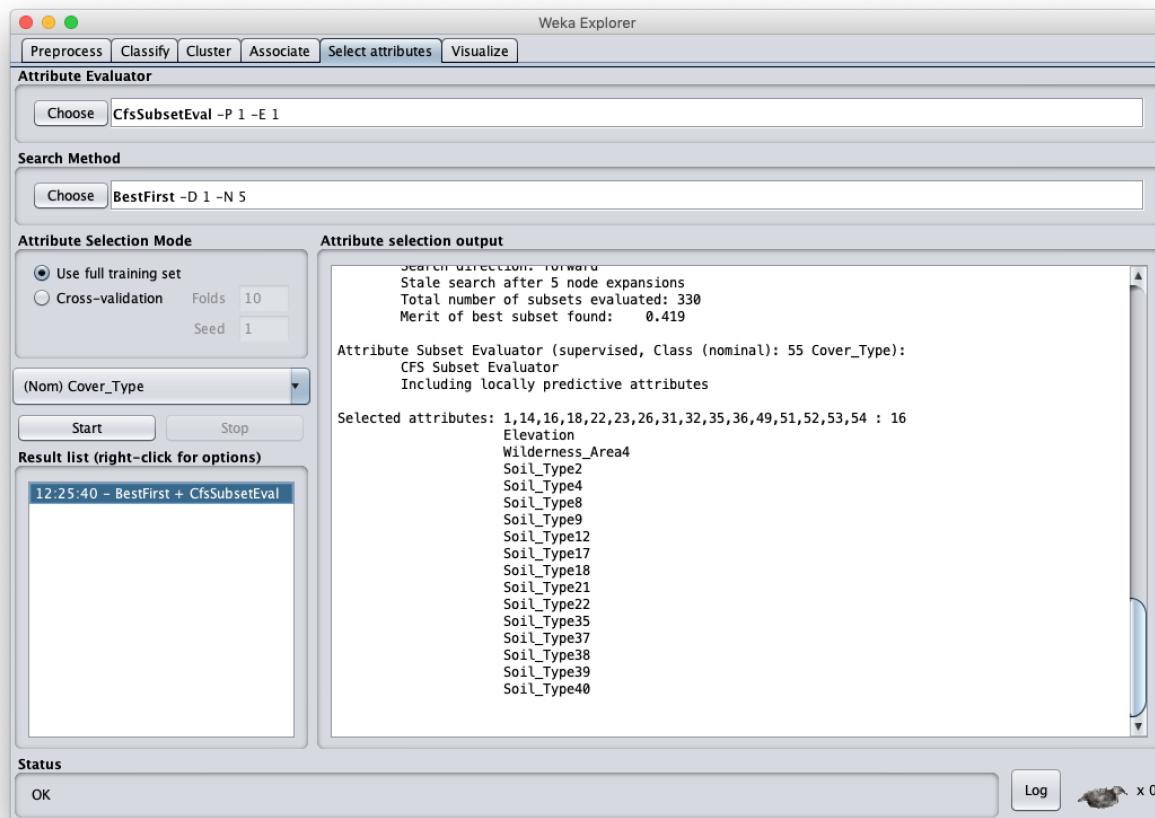


3. Click **Visualize all** to visualize all attributes.

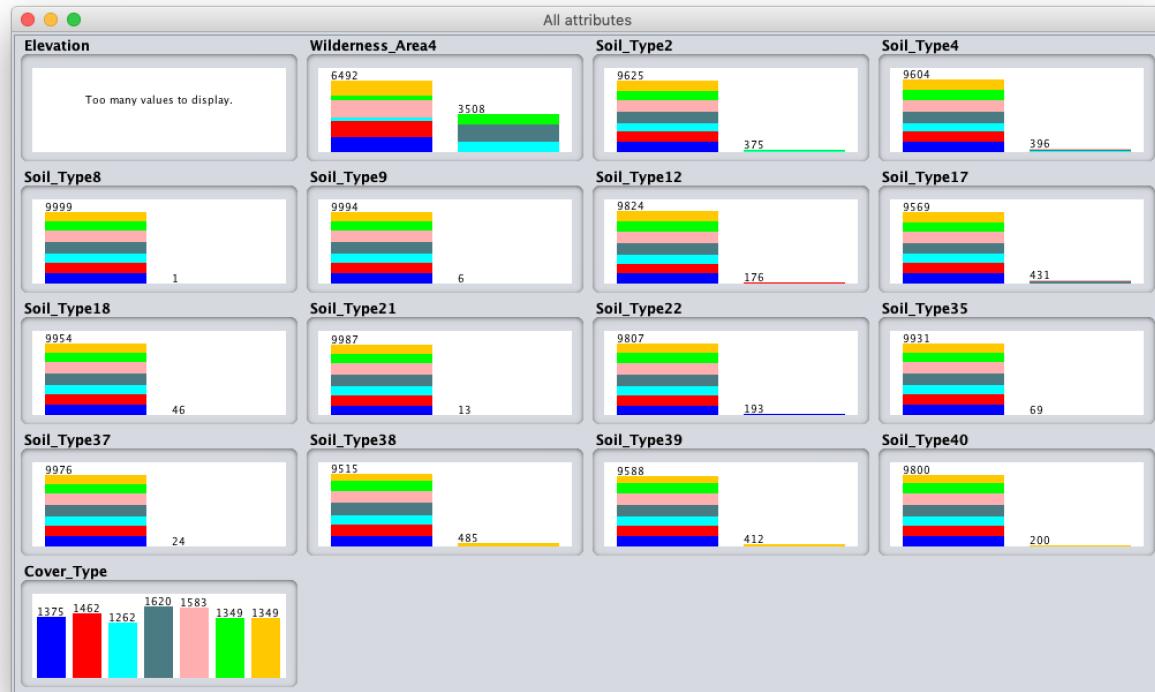


Select Attributes

1. Use default attribute evaluator (**CfsSubsetEval**) and search method (**BestFirst**).
2. Here's the selected attributes.



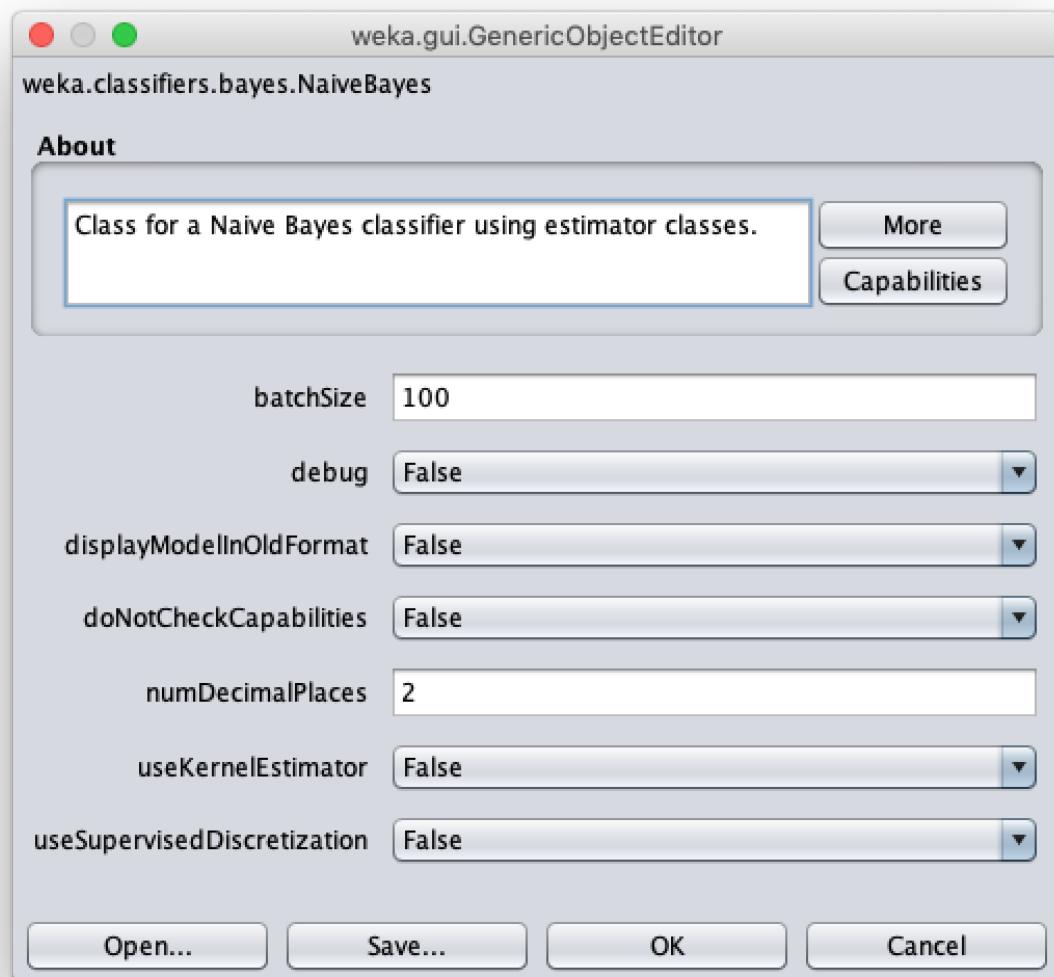
3. Visualize attributes.



Classification - Naive Bayes

1. Click **Classify** tab.

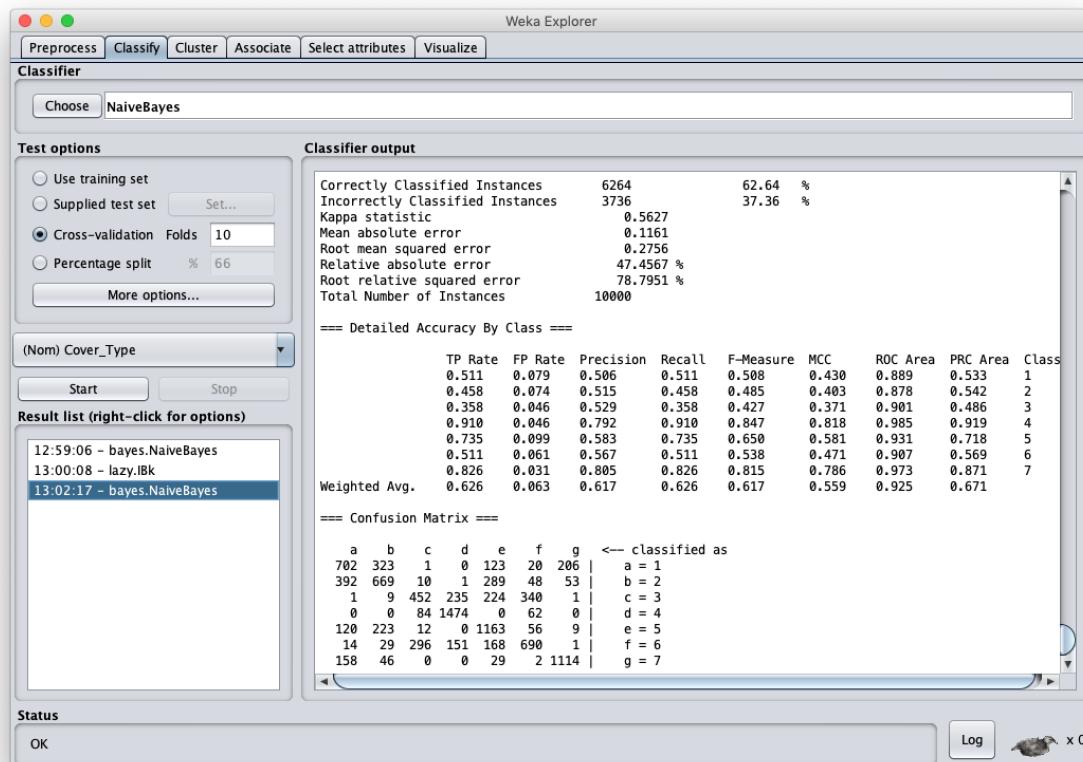
2. Click **Choose** button and select **NaiveBayes** under the **bayes** group.



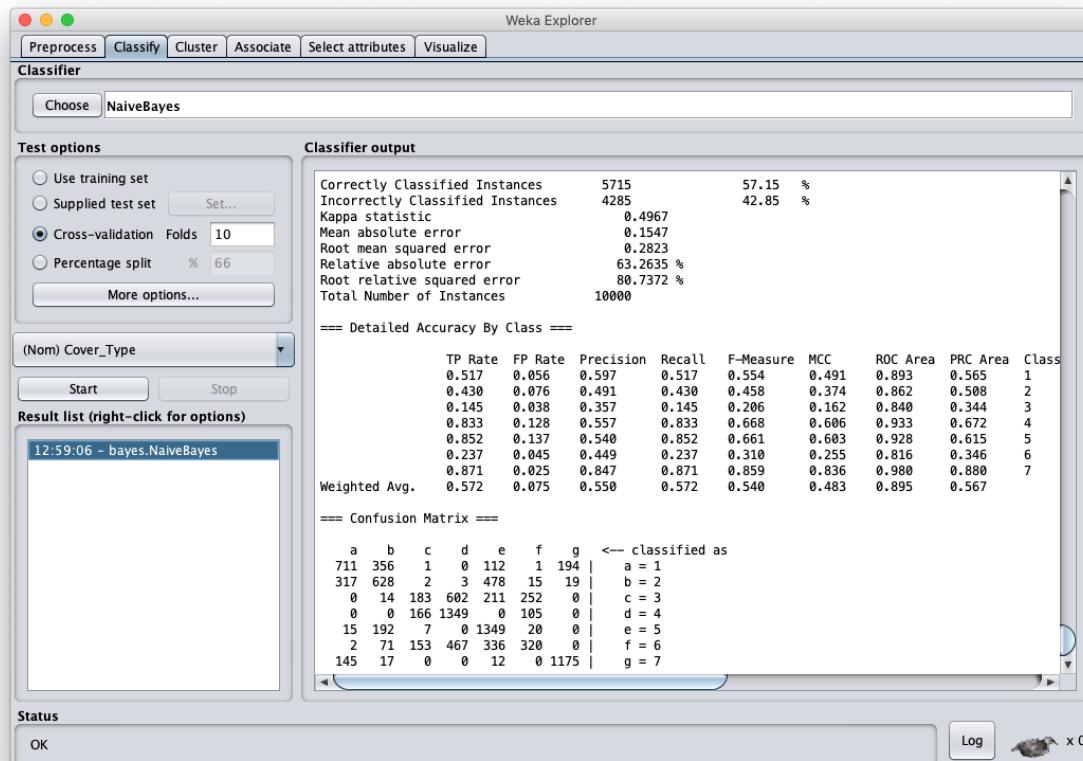
3. Click the **Start** button to run the algorithm on the Covertype dataset.

4. Here's the result.

- Origin Dataset



- After Attributes Selection

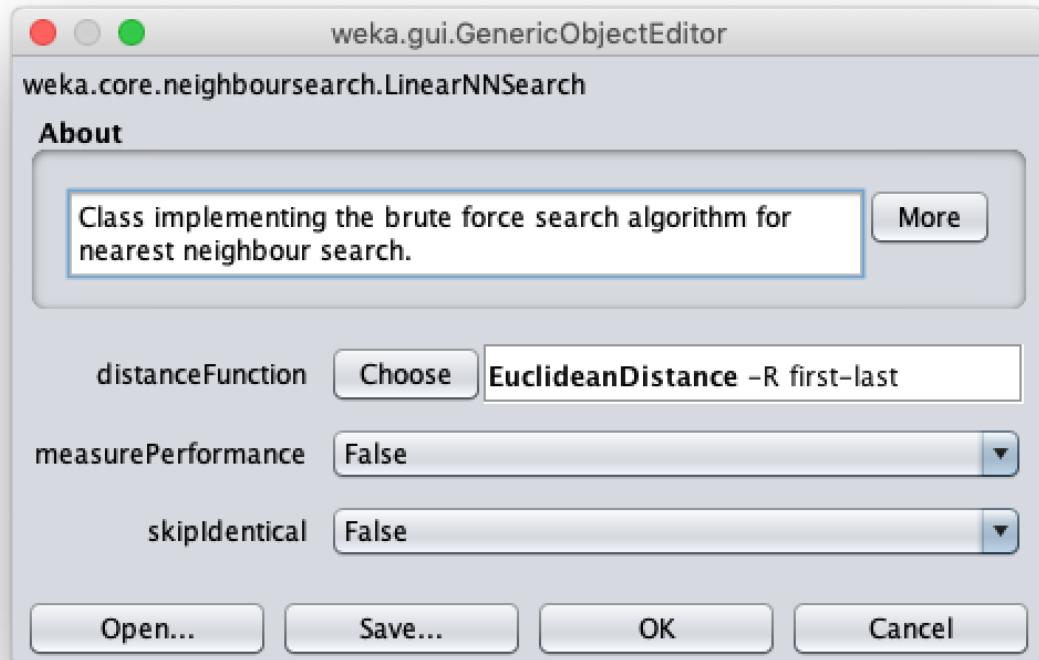


Based on the result above, **origin dataset has better accuracy** than reduced dataset after attribute selection.

Classification - K-Nearest Neighbor (KNN)

1. Still in the **Classify** tab.
2. Click the **Choose** button and select **IBk** under the **lazy** group.

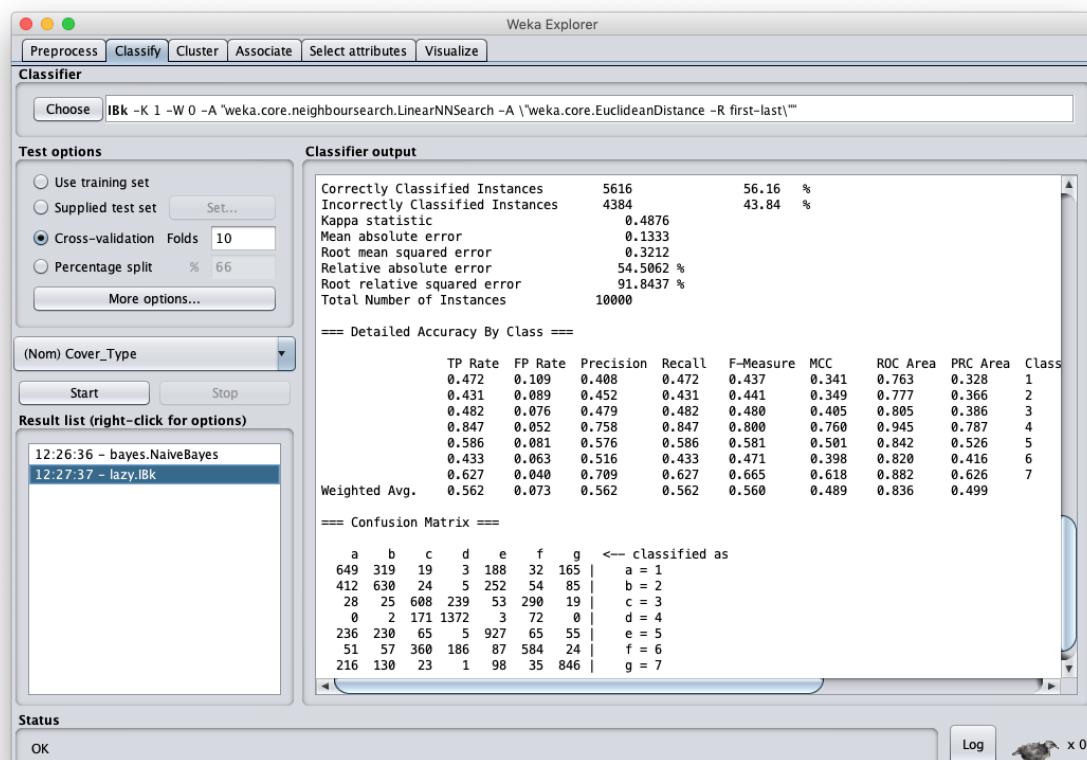




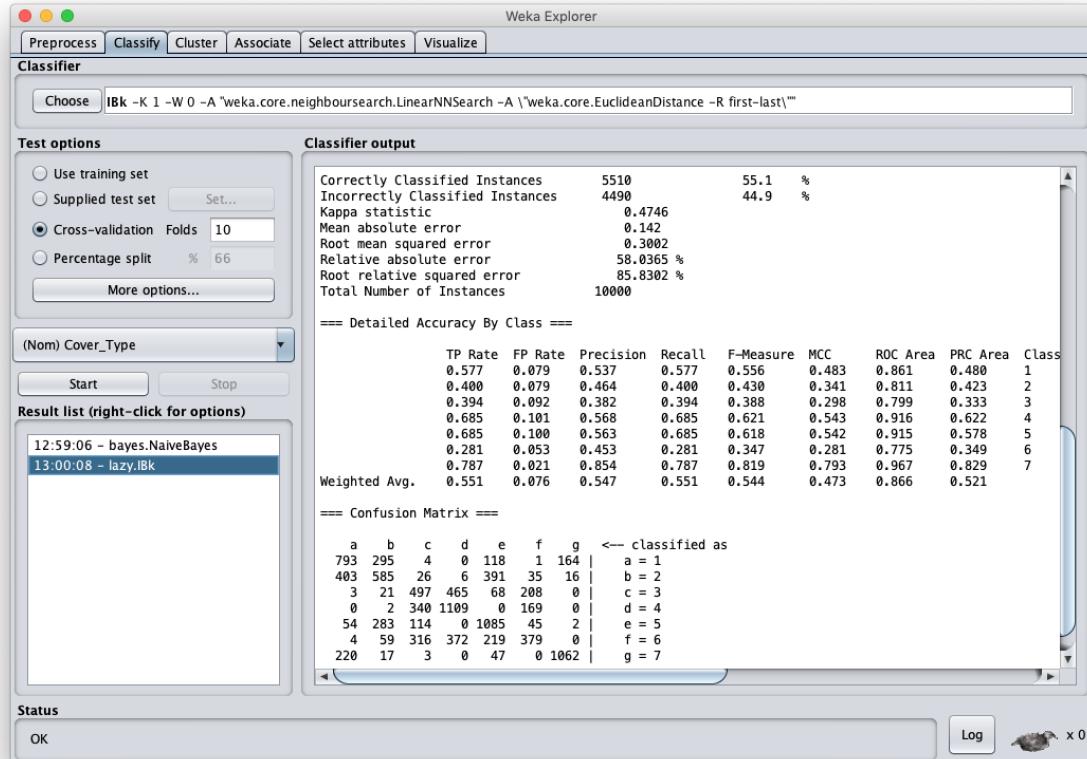
3. Click the **Start** button to run the algorithm on the Covertype dataset.

4. Here's the result using **k = 1**

- Origin Dataset

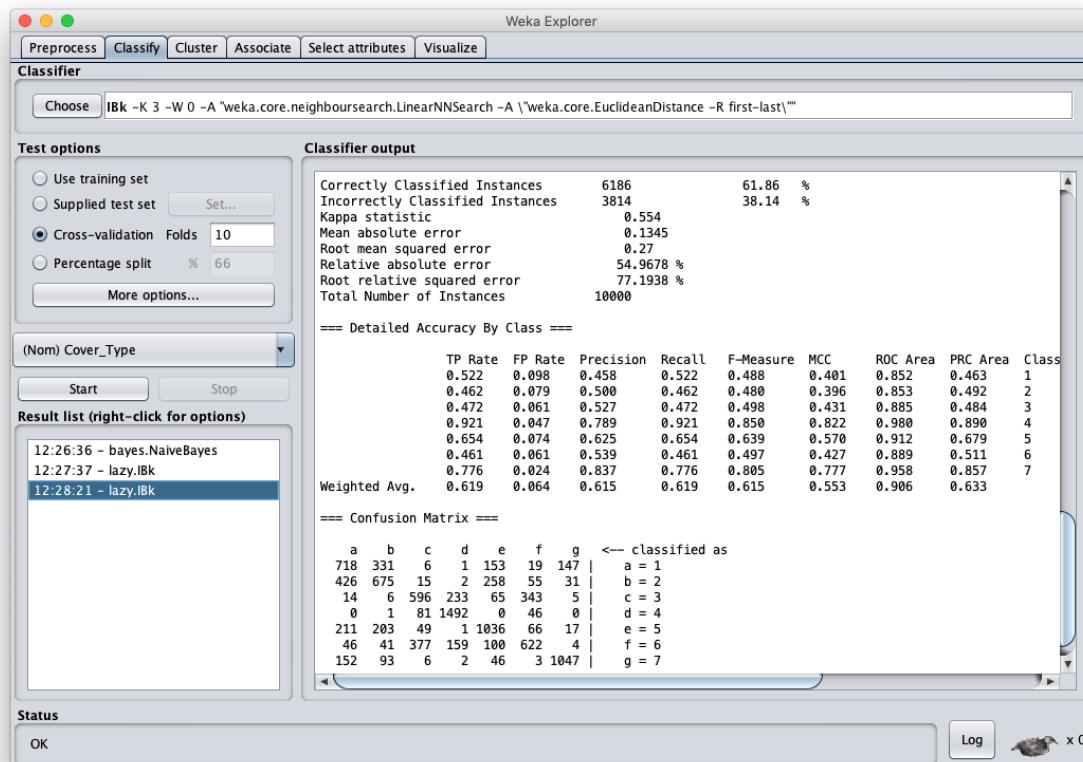


- o After Attributes Selection

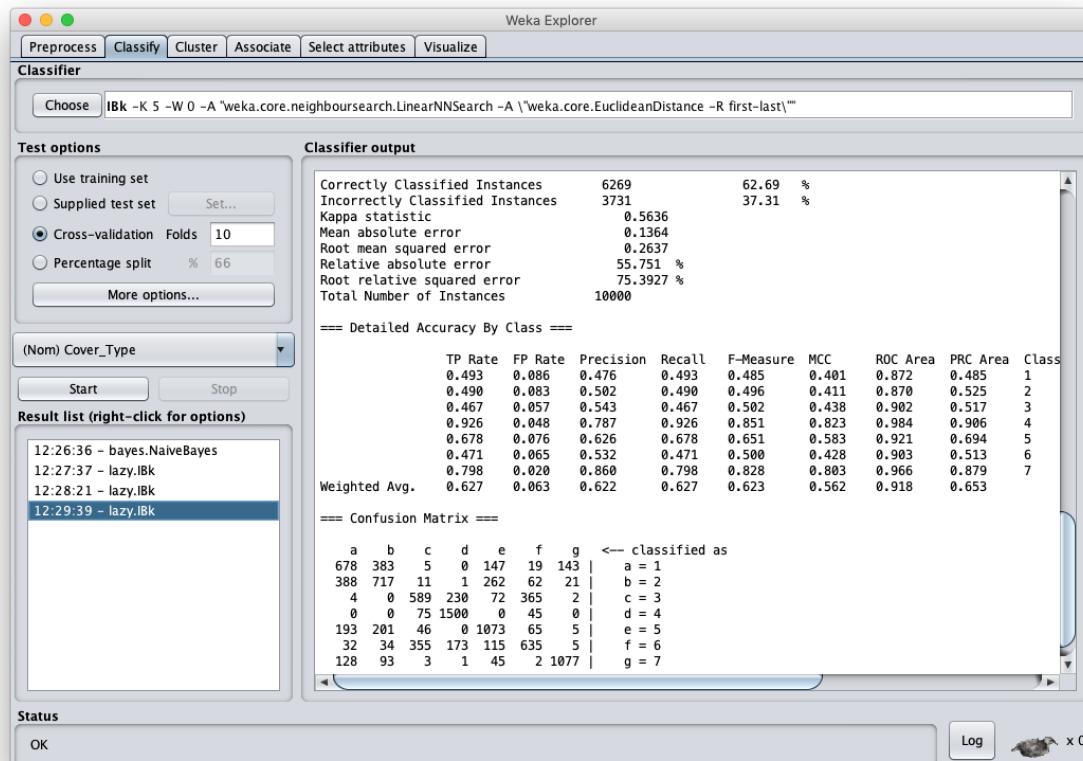


Based on the result above, **origin dataset still has better accuracy** than reduced dataset after attribute selection. So, the reduced dataset will not be used for tests anymore.

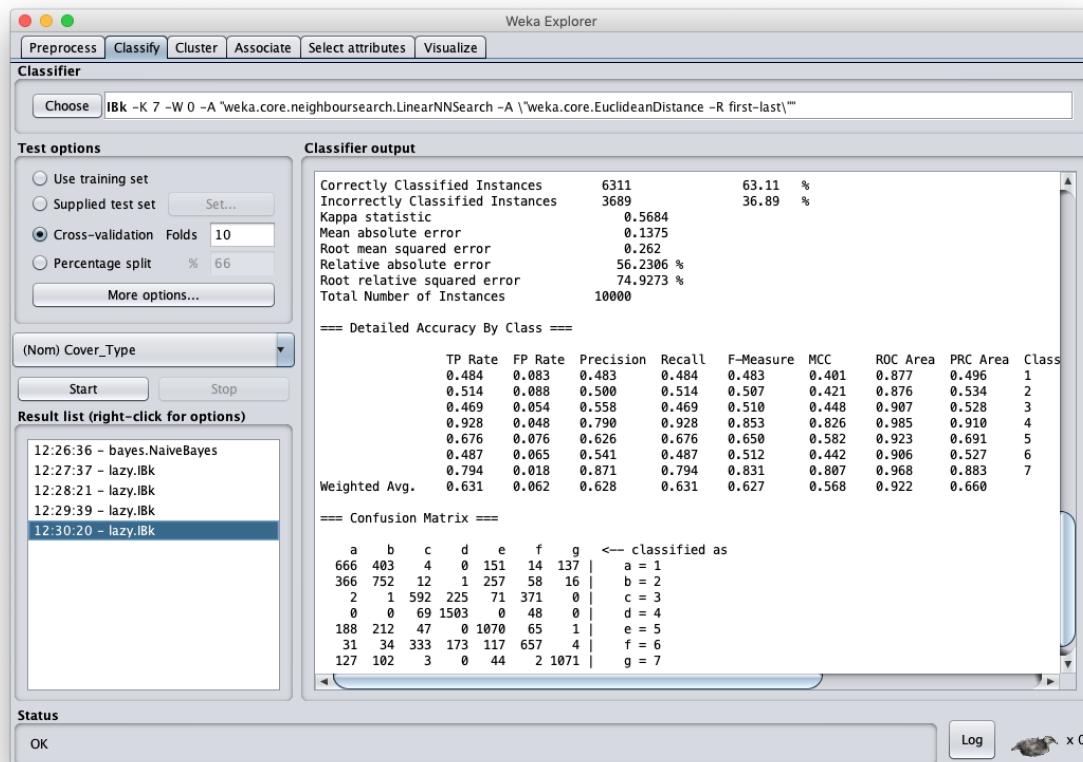
- o k = 3



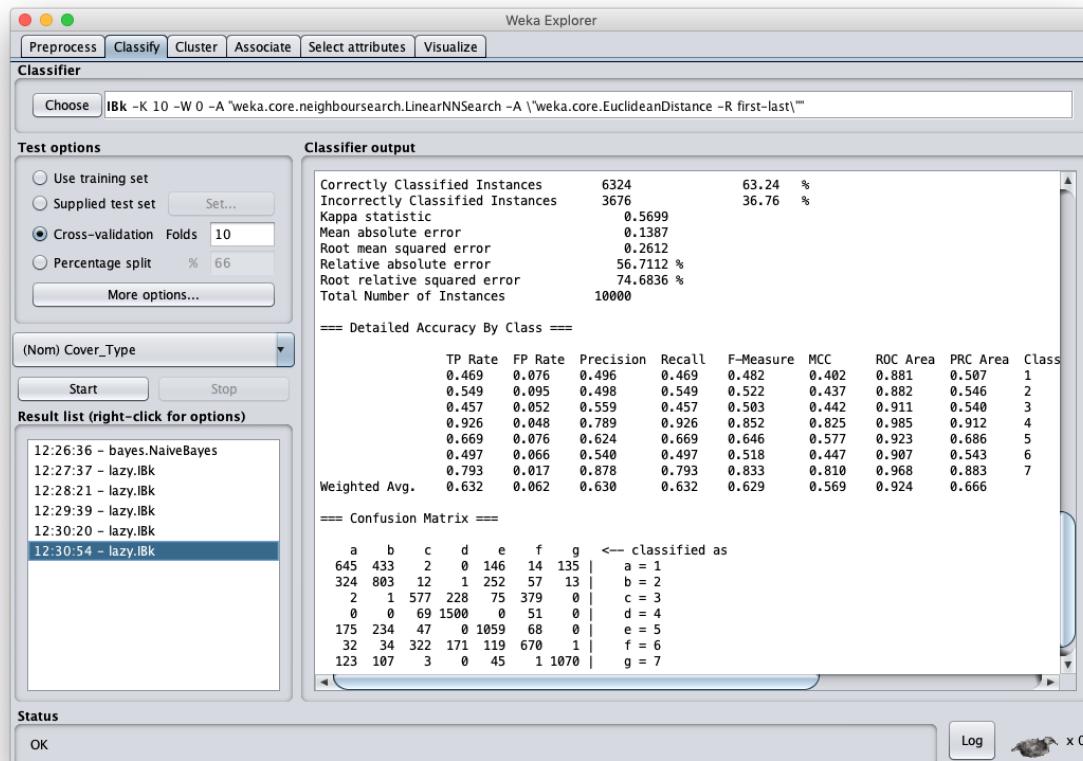
- o $k = 5$



- o $k = 7$



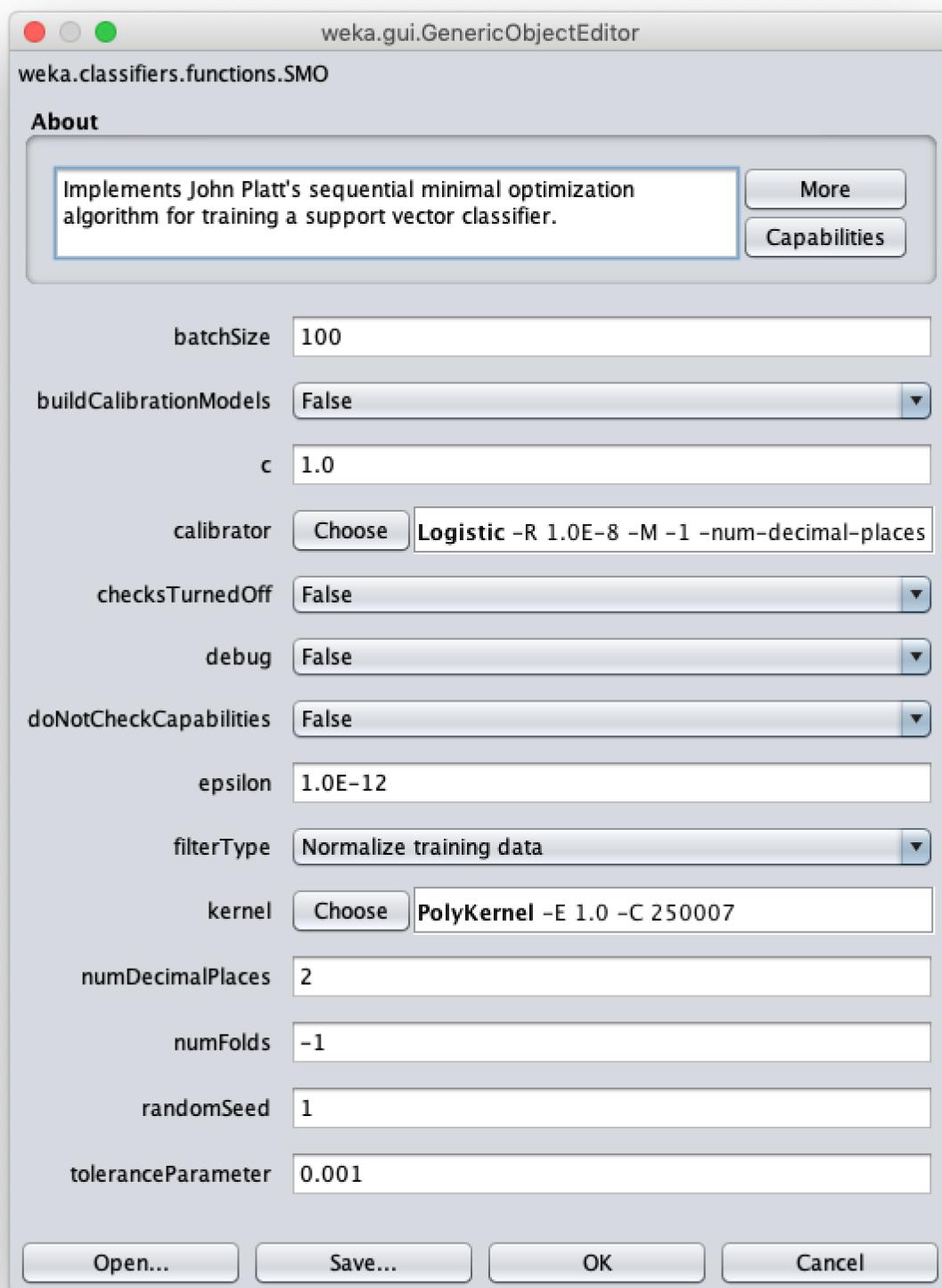
- o $k = 10$



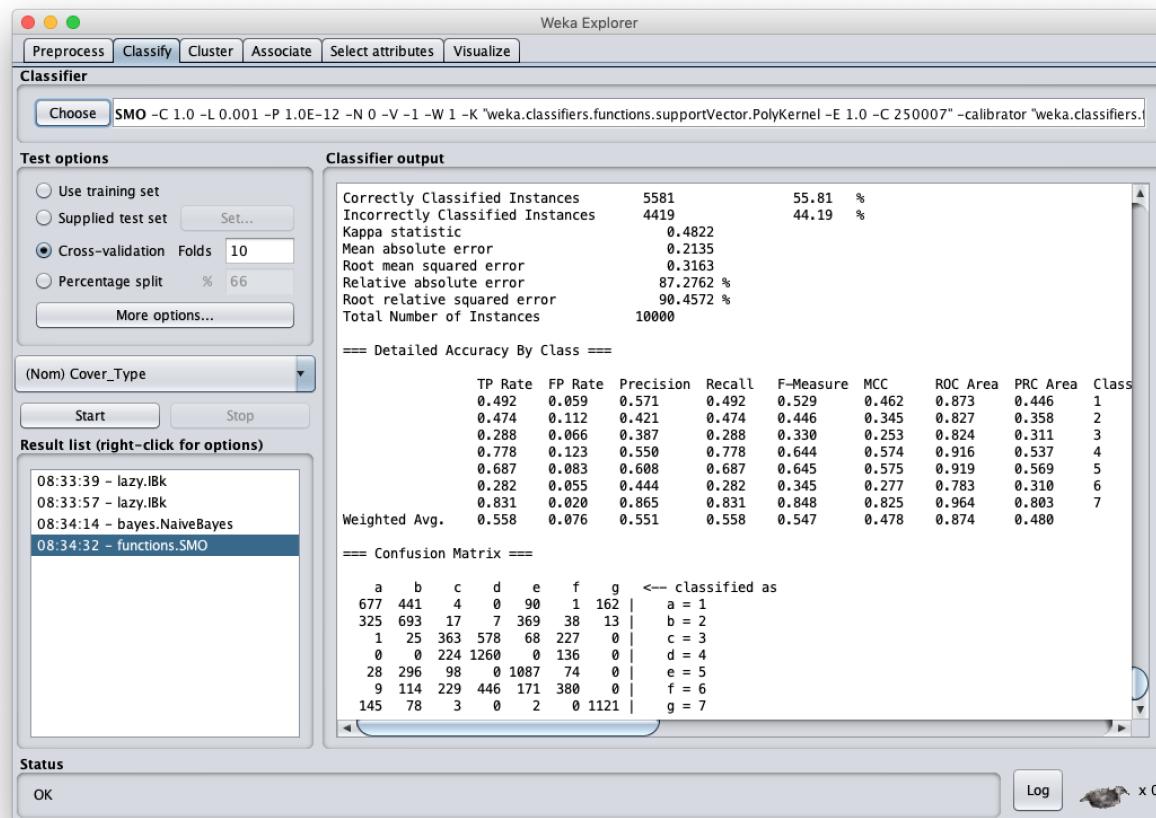
Based on the result above, **the greater k used, the better accuracy result obtained.**

Classification - Support Vector Machine (SVM)

1. Still in the **Classify** tab.
2. Click the **Choose** button and select **SMO** under the **function** group.

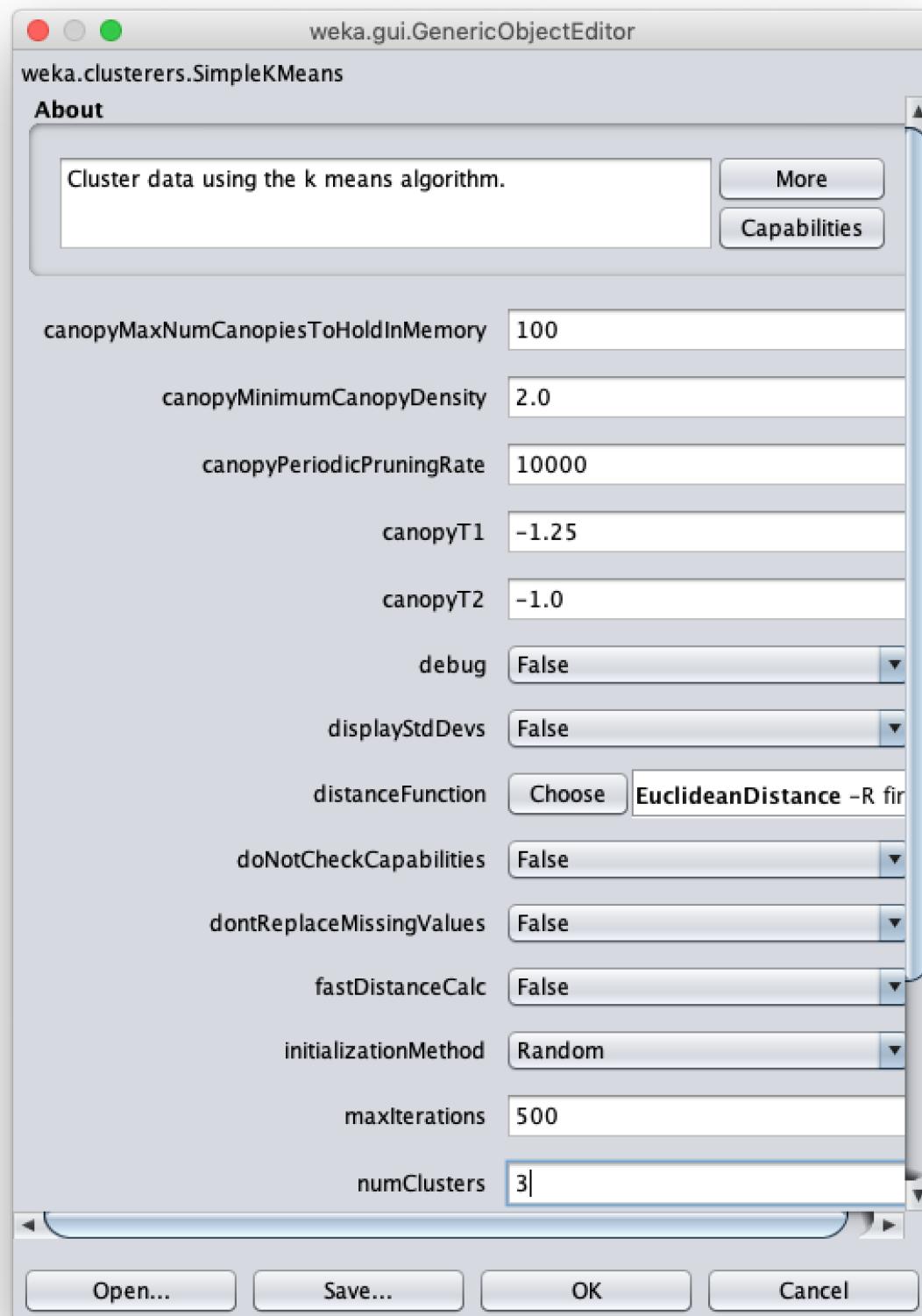


3. Click the **Start** button to run the algorithm on the Covertype dataset.
4. Here's the result.

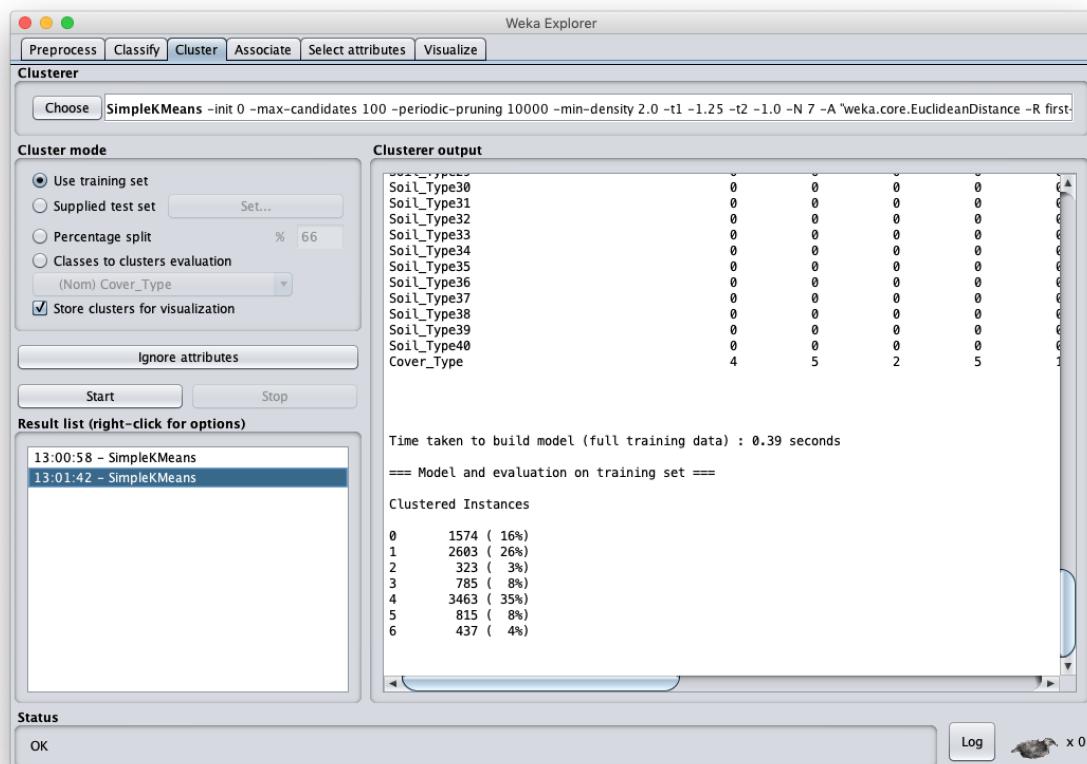


Clustering - K-Means

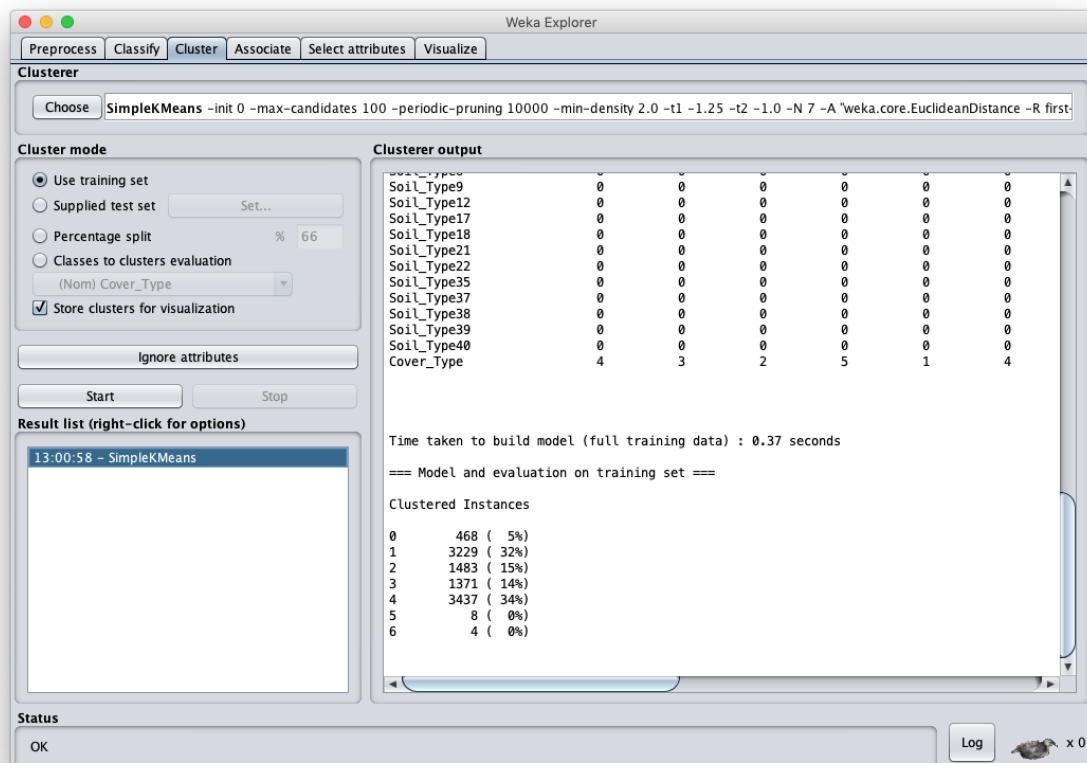
1. Click **Cluster** tab.
2. Click the Clusterer “Choose” button and select .
3. Click the **Choose** button and select **SimpleKMeans** under the **clusterers** group.



4. Click the **Start** button to run the algorithm on the Covertype dataset.
5. Here's the result using **k = 7** because there are 7 classes in this dataset.
 - Origin Dataset



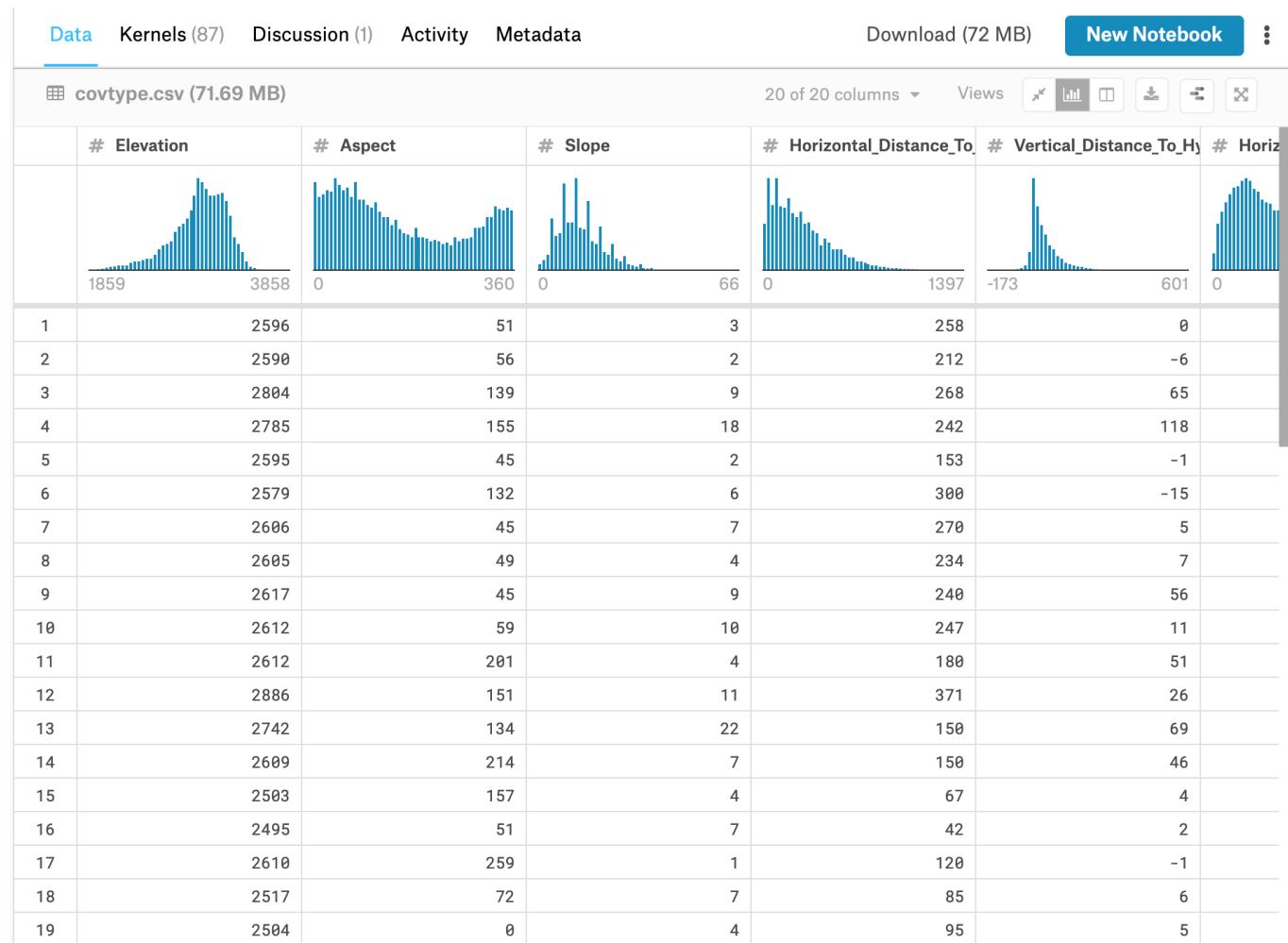
- After Attributes Selection



Conclusion

The best classification method for this dataset in this experiment is **K-Nearest Neighbour** with **$k = 10$** because it has the highest accuracy result. However if it's using **$k = 1$** it will get worse accuracy than **Naive Bayes**. The classification method that got the worst accuracy results is **SVM**, so this method is not suitable to classify this dataset.

In fact, this data processing is still bad, because the range of values for each feature in this dataset is too different. So, I think this data must be normalized first.



References

- <https://archive.ics.uci.edu/ml/datasets/covtype>
- <https://www.kaggle.com/uciml/forest-cover-type-dataset/>