# Clustering using WEKA

Dataset

- 3D_spatiao_network.arff
- au1.csv, au2.csv, and au3.csv
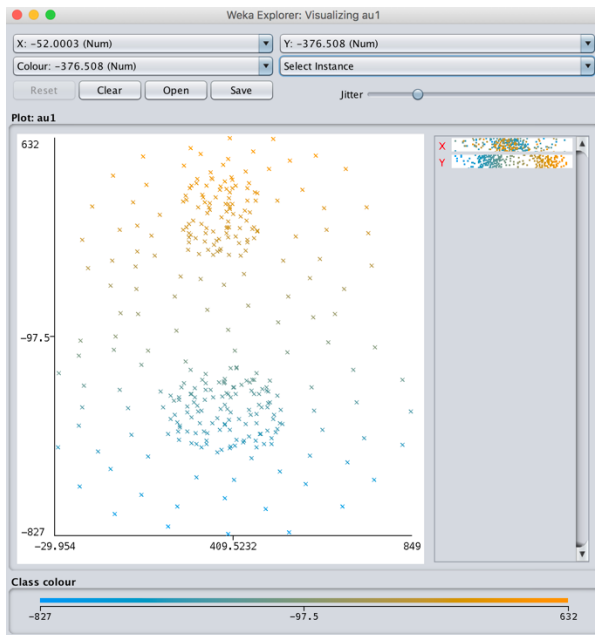- BlackFriday.csv
- highdimensional.csv

Task

1. Clustering with DBSCAN (au1.csv and au2.csv)
2. Clustering result comparison (au1.csv)
3. DBSCAN with high dimensional data (highdimensional.csv)
4. Hierarchical clustering
5. More hierarchical clustering practice

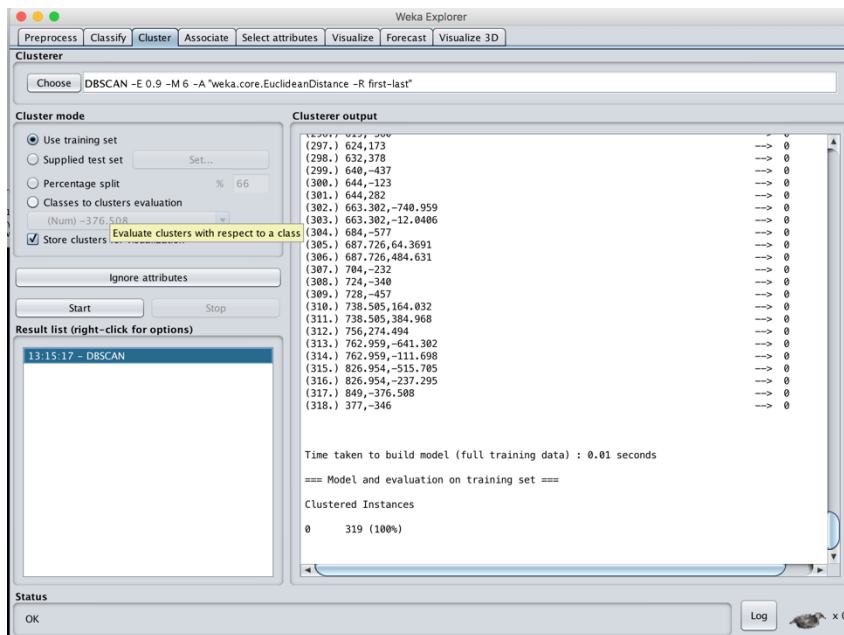**Task 1: Clustering with DBSCAN**

Step 1: if you do not have DBSCAN and OPTICS installed in your Weka, then please go to **Weka GUI Choose->Tools->Package Manager** and search and install **optics_dbScan**

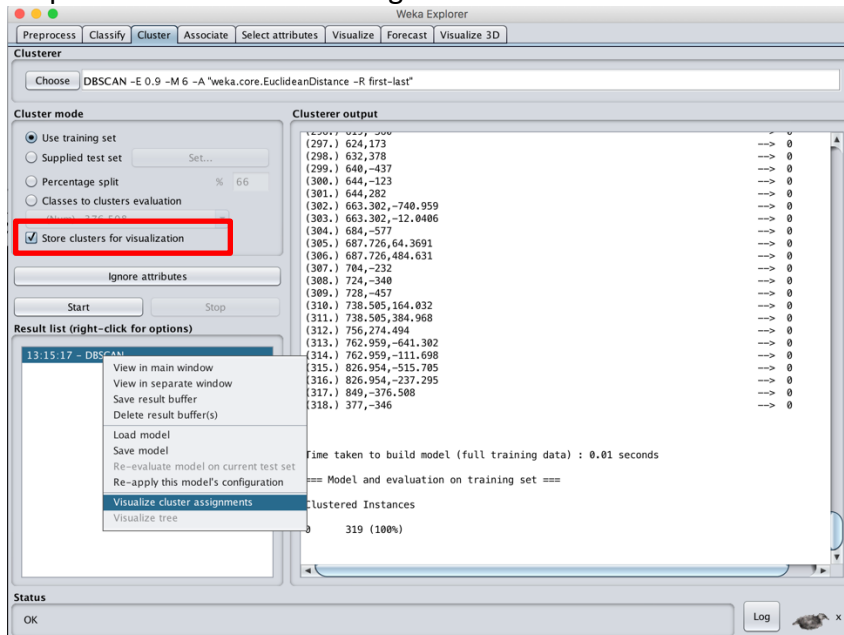Step 2: Load au1.csv data into WEKA.
The data has two dense clusters and a sparse cluster (could be seen as noise) around these two dense clusters as below.
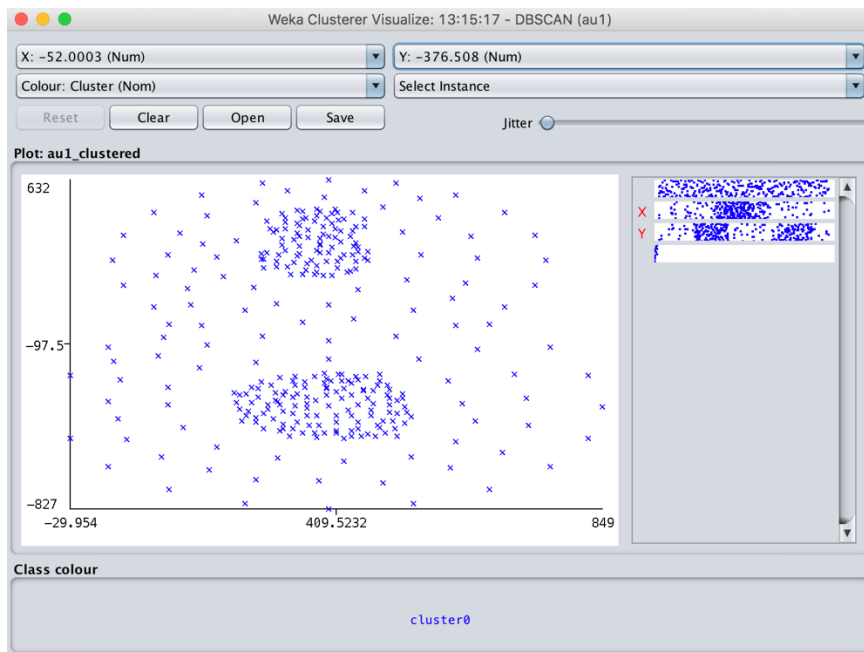


Step 3: Run DBSCAN with a default value. Check 'Store clusters for visualization' to see the clustering result.
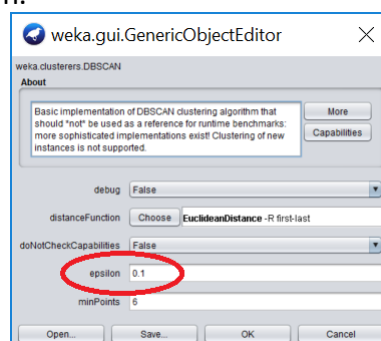
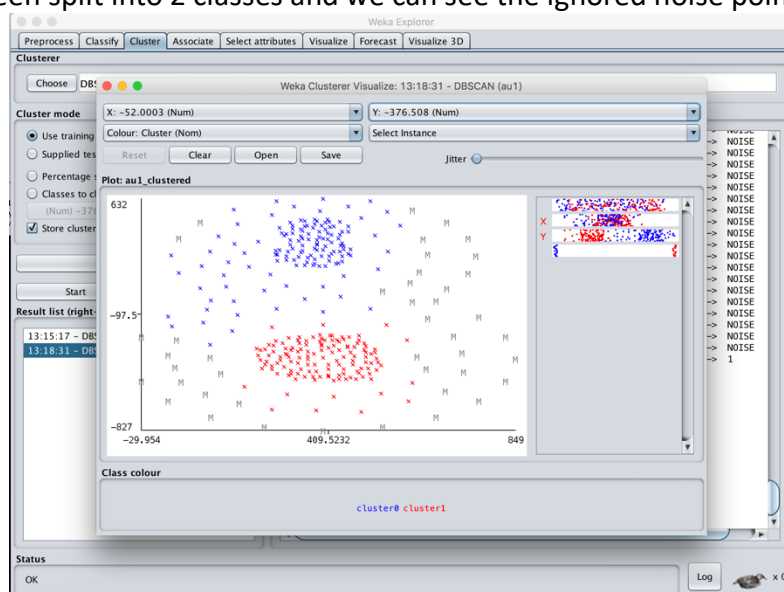Step 4: Visualise the clustering result.



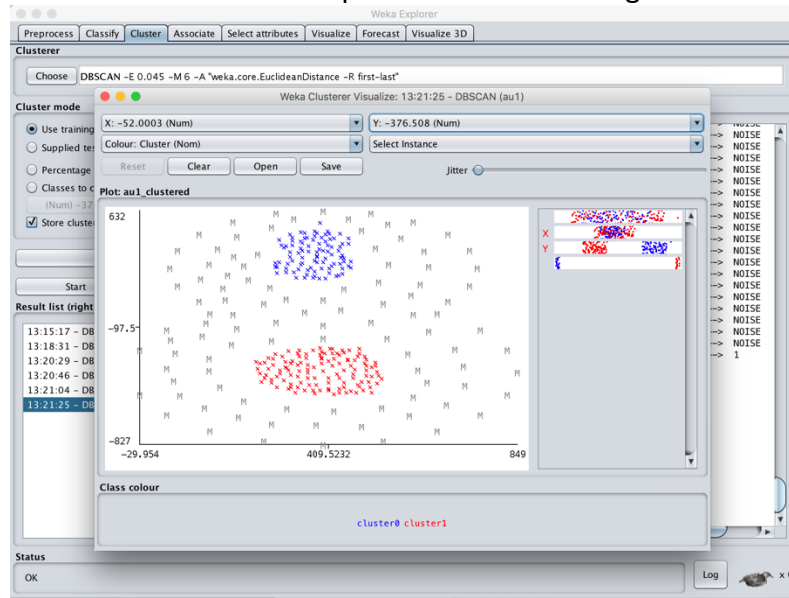Visualisation shows one cluster with everything as below.

DBScan has clustered everything together into one cluster due to epsilon being too large. Change epsilon distance to 0.1 and run again.
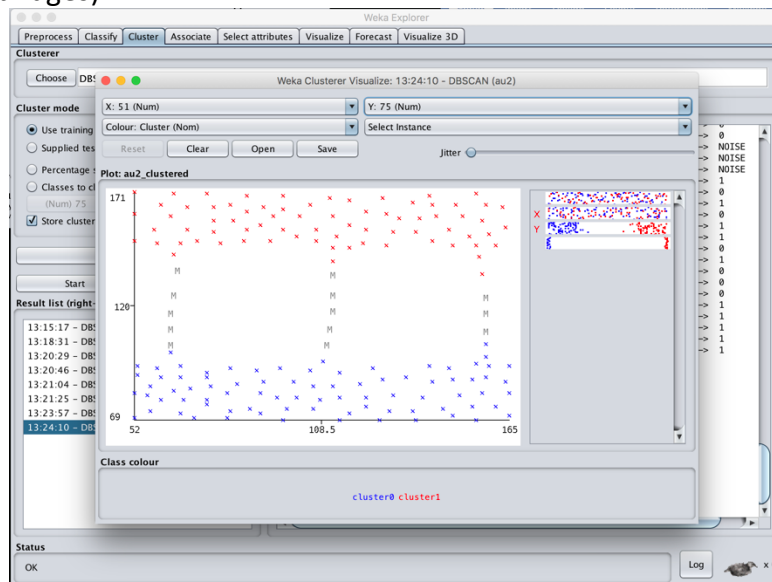


Now the data has been split into 2 classes and we can see the ignored noise points marked as M.

Step 5: But still the clustering is not perfect. Try tweaking the epsilon parameter to detect better clusters. Epsilon value = 0.045 and minPts = 6 produces the following result.
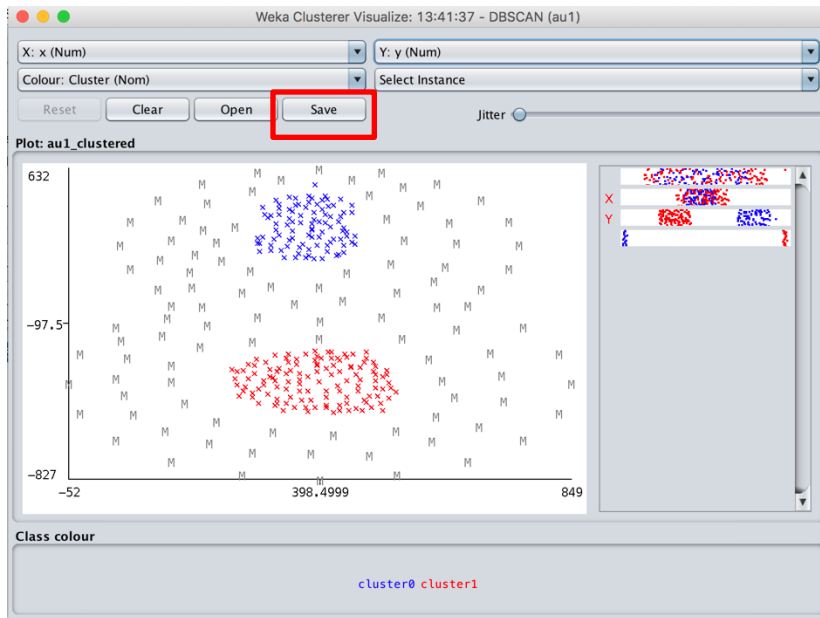


Step 6: load au2.csv, and try to tweak the epsilon parameter to find the following two clusters connected by links (bridges).



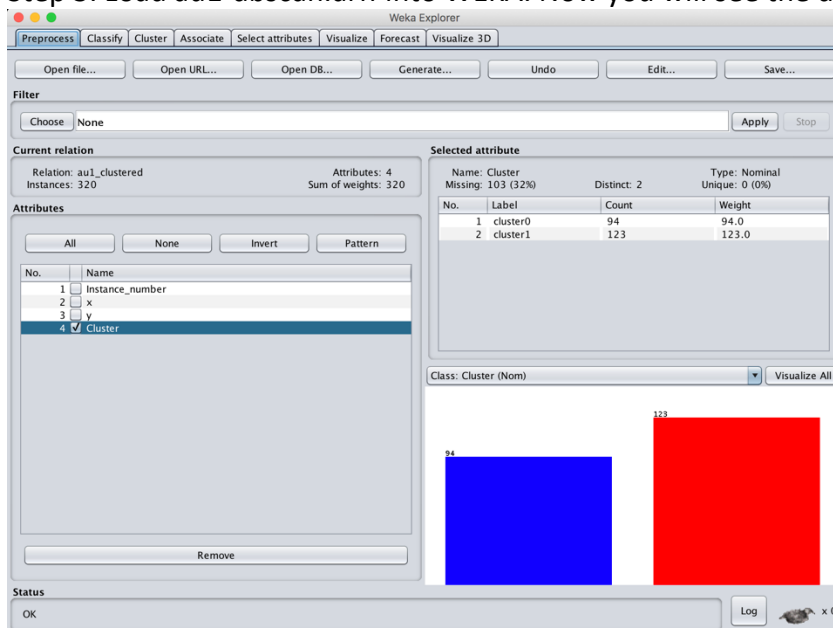## Task 2: Clustering result comparison
You may wish to store the clustering result to the file.
Step 1: You can save it when the clustering result is visualised.
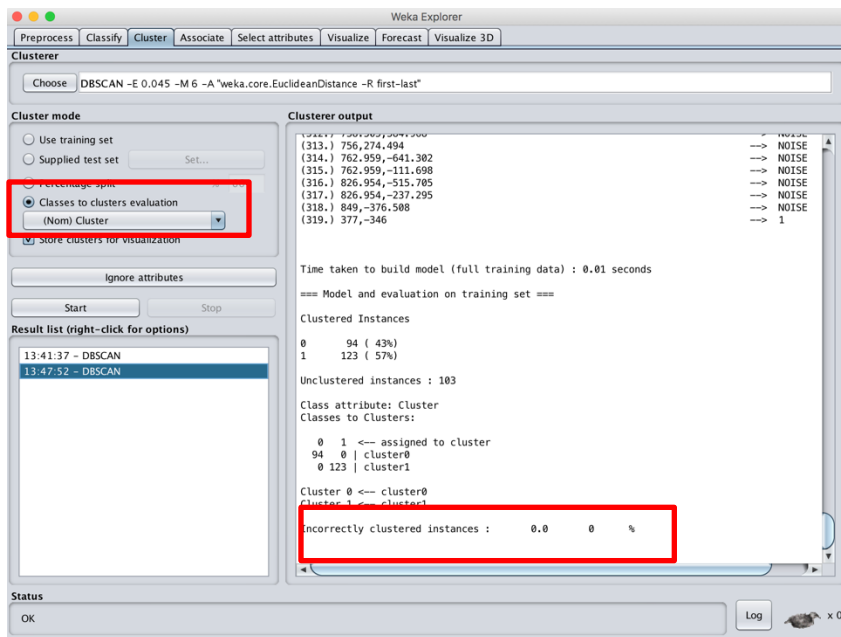
Step 2: Name the file au1-dbscan.arff.

Step 3: Load au1-dbscan.arff into WEKA. Now you will see the data with 4 attributes (ID, x, y, Cluster).
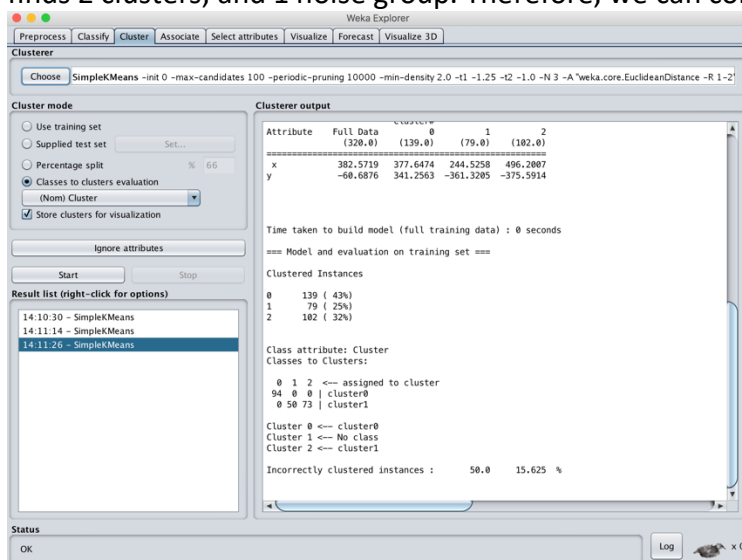


Step 5: remove 'Instance_number' attribute.

Step 6: Cluster using DBSCAN (epison = 0.045 and minPts = 6), and compare this clustering result with 'Cluster' attribute in the file. Then not surprisingly, you will see it is 100% accurate.

Step 7: now compare this DBSCAN result with k-means clustering with k=3. Please note that DBSCAN finds 2 clusters, and 1 noise group. Therefore, we can compare it with k-means k=3.
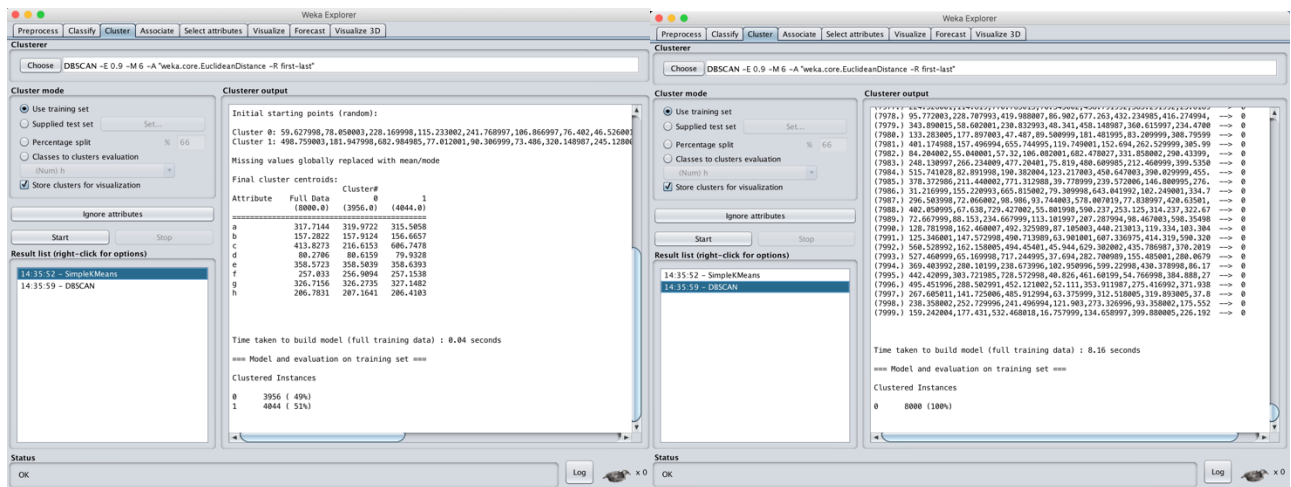


You can see that there is around 15% misclassification. It appears that all those noise points are misclassified in k-means.

## Task 3: DBSCAN with high dimensional data

k-means clustering is fast and works well for high dimensional data, however DBSCAN is able to detect quality clusters (arbitrary shapes) but computationally very expensive (slow) for high dimensional data.

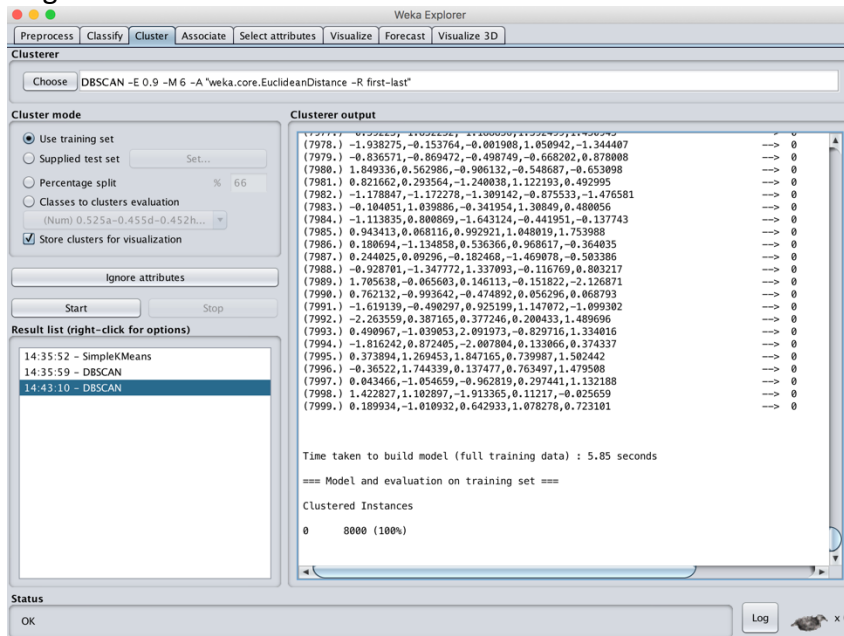Step 1: open highdimensional.csv. The data set has 8000 instances with 8 dimensions.
Step 2: Run k-means and DBSCAN with default values and measure time taken to calculate. As you can see below, it takes 0.04 sec and 8.16 sec in my machine (this would vary with your machine).

Step 3: Once possible way to handle high dimensional data with DBSCAN, is to use Feature Selection or Feature Reduction (Dimension Reduction) approach. Let's use PCA for this data.

Step 4: with 60% of variable captured, PCA produces 5 principle components.

Step 5: Now run DBSCAN with the same default value. Now it takes 5.85 sec which is faster than the original 8 attributes.
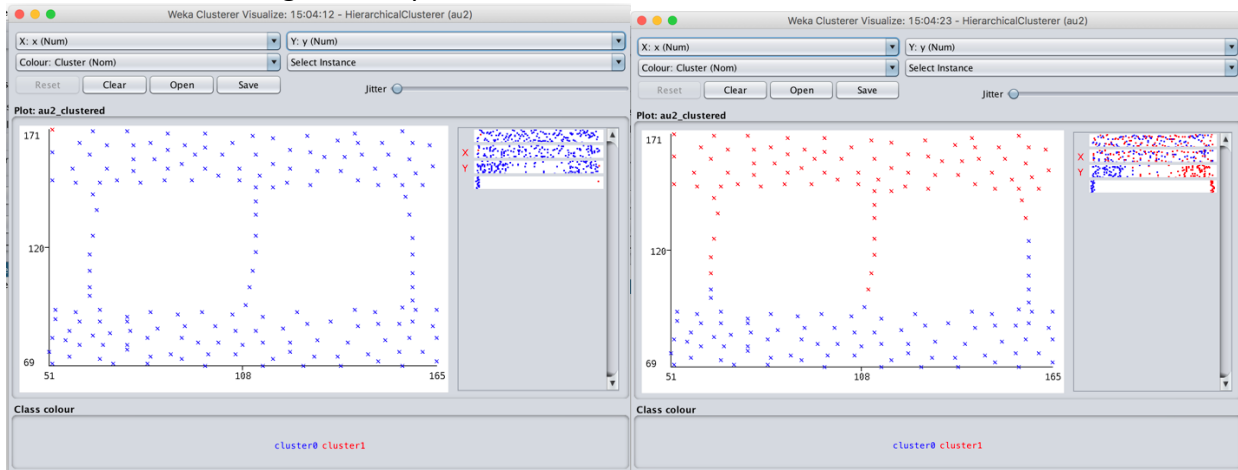
**Task 4: Hierarchical clustering with WEKA**

Bottom-up, agglomerative, hierarchical clustering is a popular approach to generate a clustering hierarchy from the dataset. There are many different types of merge techniques including single-linkage, complete-linkage, mean-linkage, and average-linkage. As you can easily notice, hierarchical clustering is more computationally expensive (slower) than partitioning clustering, so it would take longer time with larger data.

Step 1: load au2.csv into WEKA.

Step 2: Run Hierarchical clustering with the number of clusters = 2. Try to observe the result with various different merge techniques. Which one is able to detect the two clusters?



## Task 5: More Practice with DBSCAN and Hierarchical clustering

More clustering practice with 3D_spatial_network and BlackFriday.csv data with DBSCAN and hierarchical clustering. You might need to do preprocessing before you are able to use them.

This is the end of this week's prac.