# Associations and Clustering

*PROF. K. H. LIM*

50.038 Computational data science

# Pre-processing

o Various pre-processing steps to handle feature

o Look at Filter→Unsupervised→Attribute

o Which will you need for association rule mining?

# Association Rule Mining

o Weka provides an implementation of the Apriori algorithm

# Association Rule Mining

o Weka provides an implementation of the Apriori algorithm

 ◦ Options to set various parameters, e.g., minSup and minConf thresholds

# Exercise

1. Load in the ./weka-3.8/data/credit-g.arff dataset

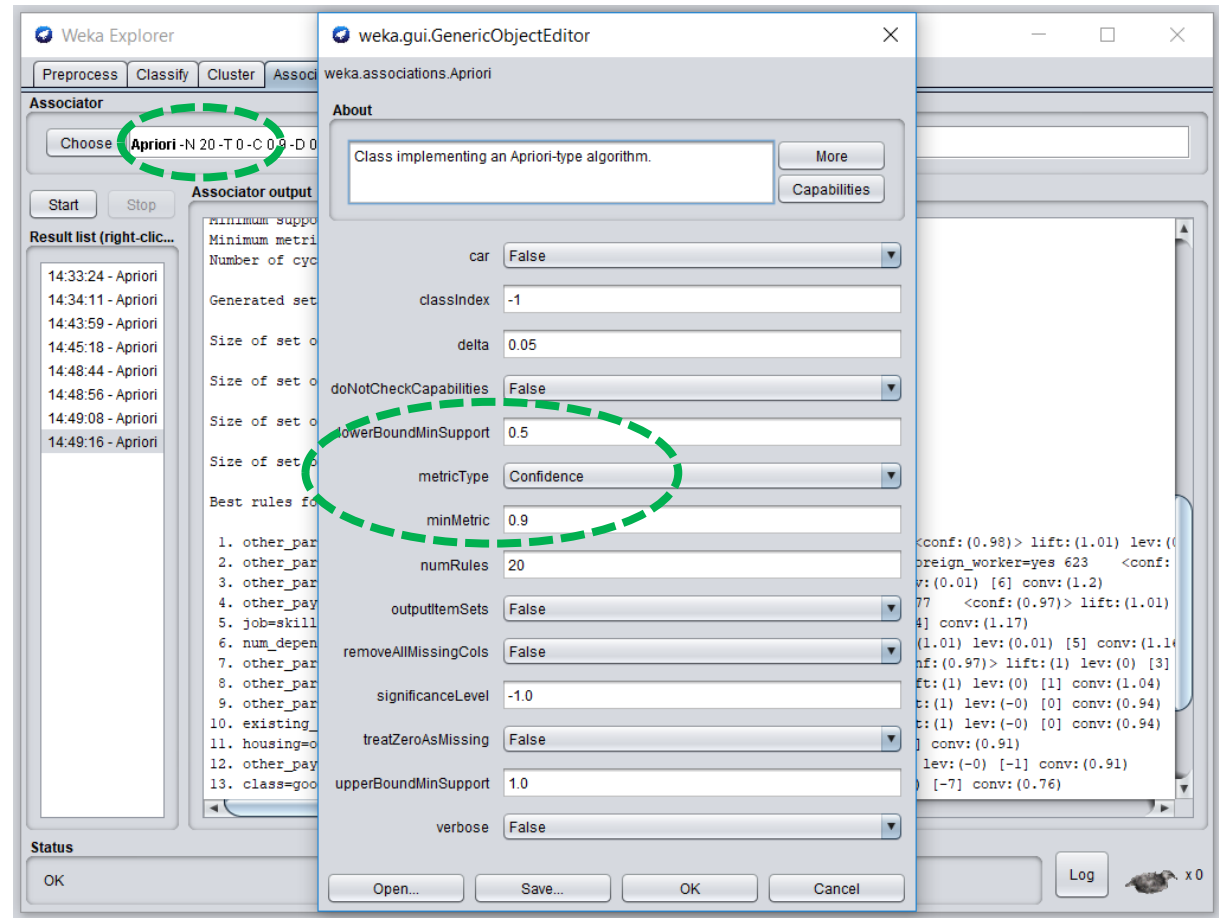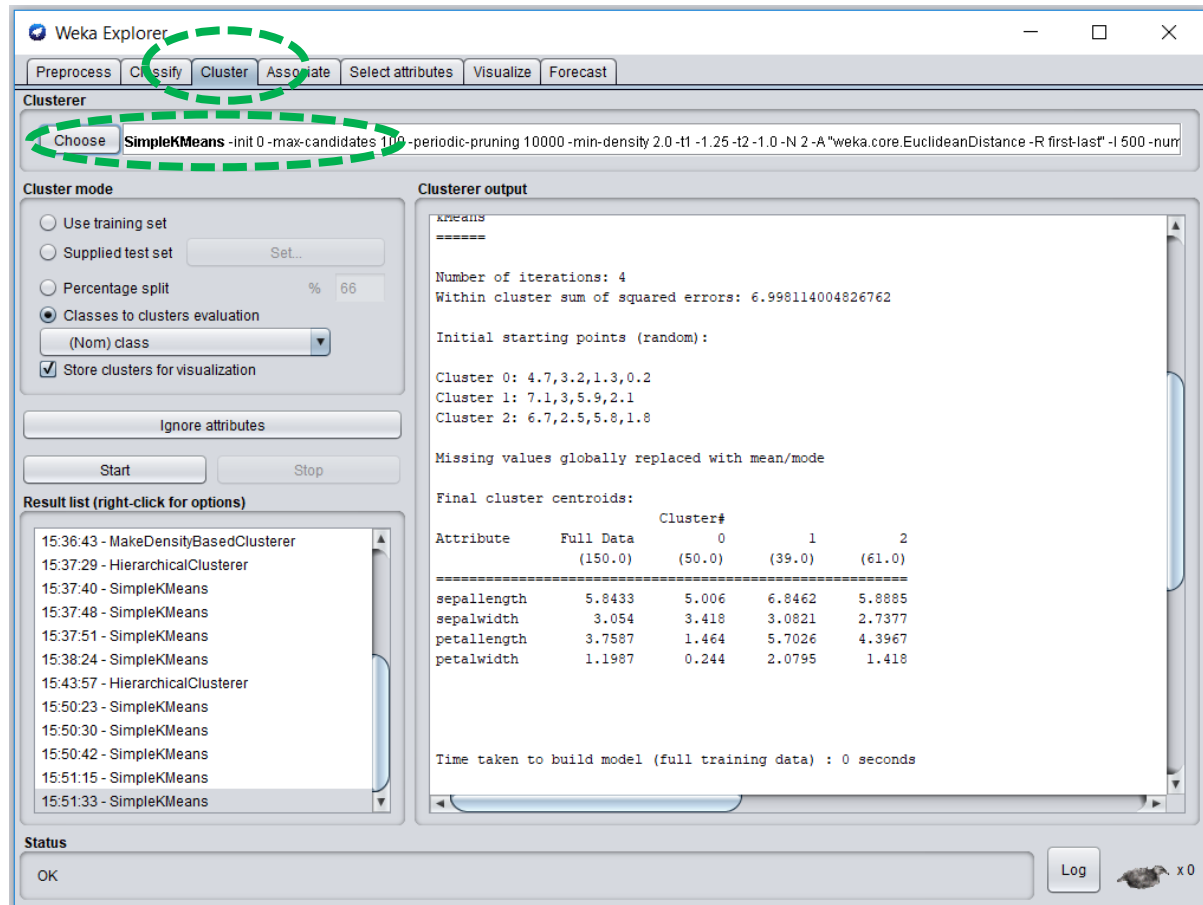2. What types of features are in the dataset?
   norminal, numeric

3. How should you pre-process the dataset before applying association rule mining?   discretize

4. With a minSup=0.8 threshold, identify the top 10 association rules (based on confidence scores)?
   - What do you observe? How can you obtain the top 10 rules?

5. Now load in the ./weka-3.8/data/supermarket.arff dataset
   - With a minSup=0.5 threshold, what are the frequent 2-itemsets?
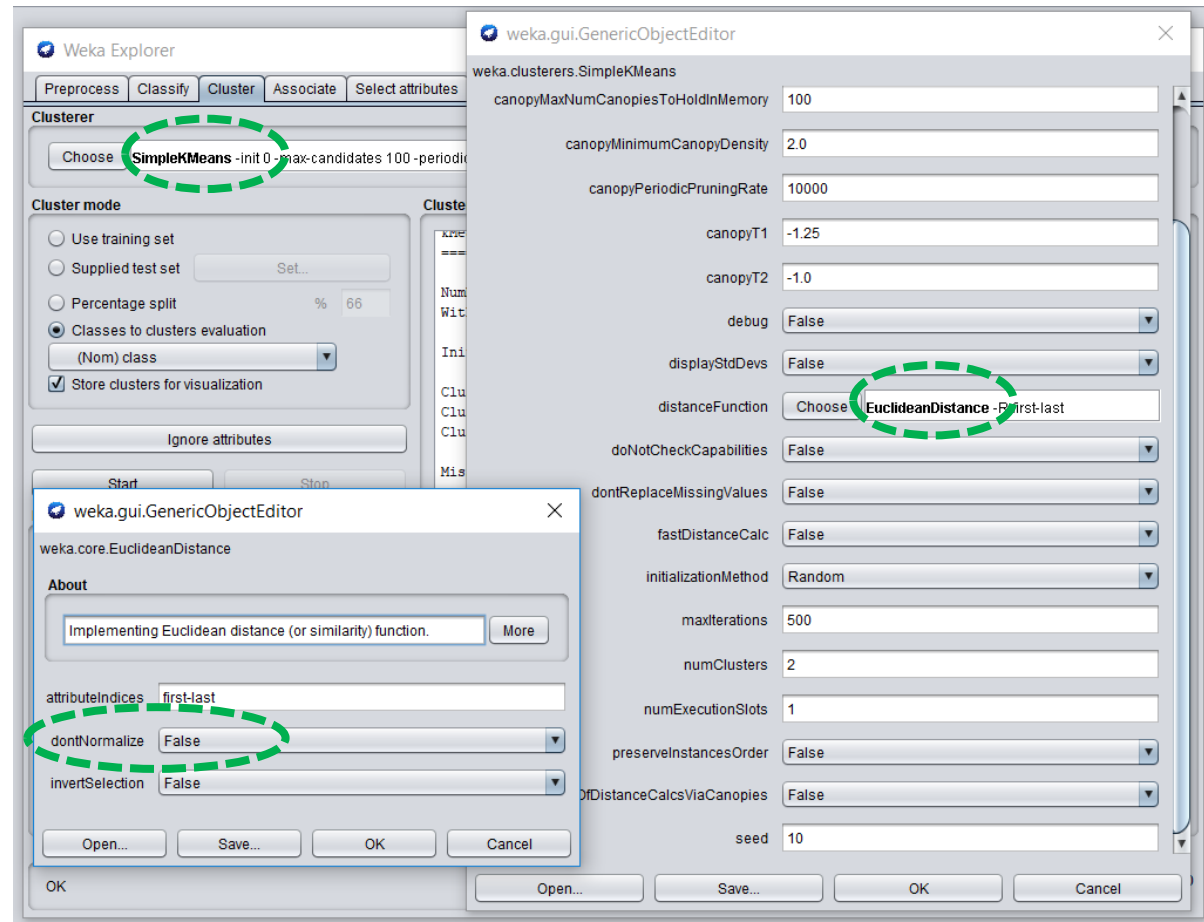   - What are the top 3 association rules based on confidence?

# Clustering

o Weka provides implementations of various clustering algorithms, including k-means and hierarchical

50.038 COMPUTATIONAL DATA SCIENCE – ASSOCIATION, CLUSTERING AND COMMUNITY DETECTION
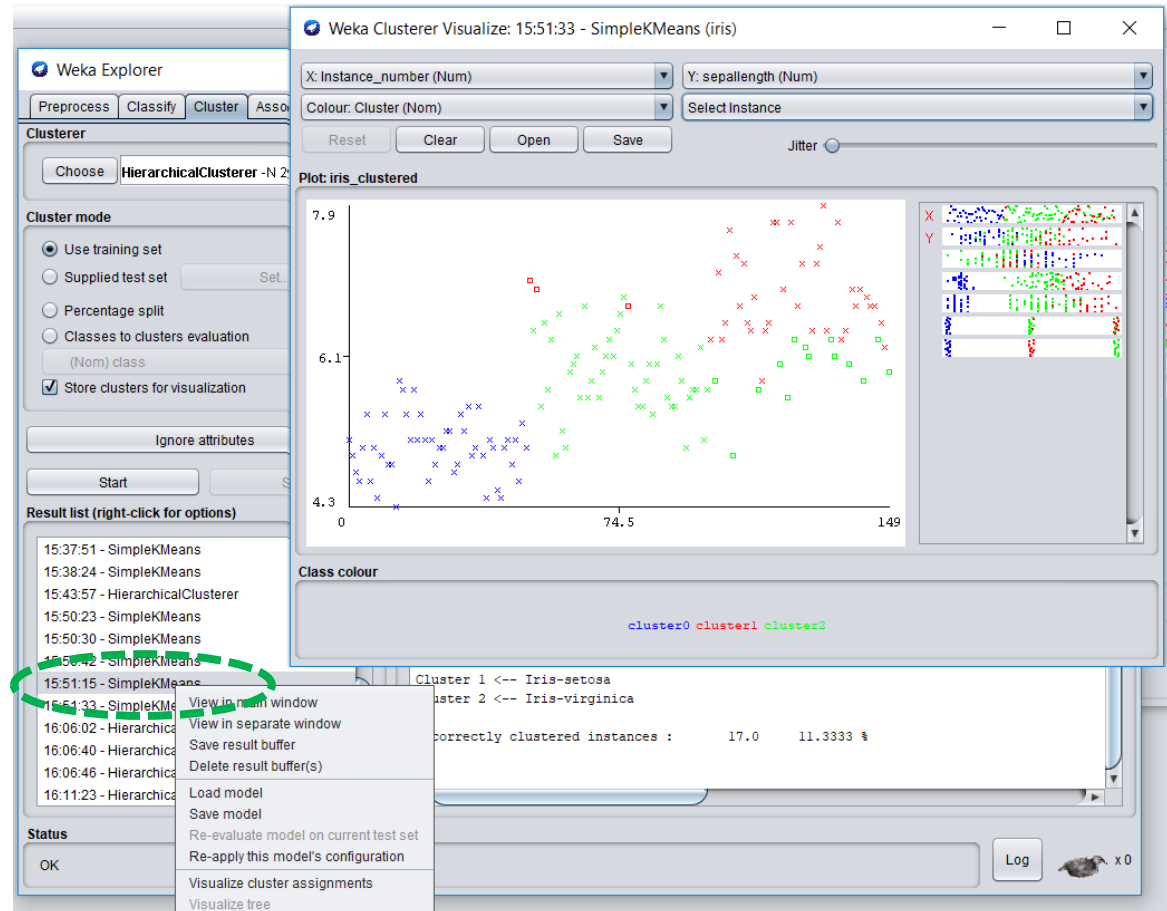
# Clustering

o Weka provides implementations of various clustering algorithms, including k-means and hierarchical

  ◦ Able to fine-tune k-means in various ways, e.g., select distance measure, set seeds, feature normalization, set k-value, etc

# Clustering

o Able to visualize cluster assignments based on different features

◦ Right-click on "result list", and select "visualize cluster assignment"

# Exercise

1. Load in the ./weka-3.8/data/iris.arff dataset

2. Run the k-means (SimpleKMeans) algorithm multiple times with k=3 and observe the sum of squared errors (SSE) values.
   - K-means typically return different clusters with each run, why do you observe in terms of SSE and why is this so?

3. Run k-means again, with feature normalization and without.
   - What do you observe now in terms of SSE?

# Project

o Use the rest of the lab to work on your projects

o Presentation during Week 8
  ◦ Allocated time of 10min per group
  ◦ Details on presentation slots available from next week
  ◦ Please ***sign up for a presentation slot***