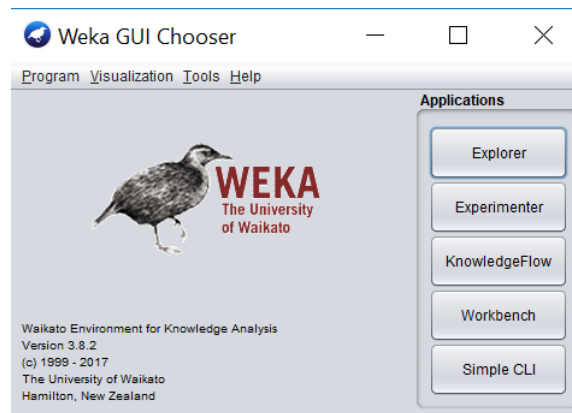# Documentation: Attribute Importance Analysis of Bank Marketing Data
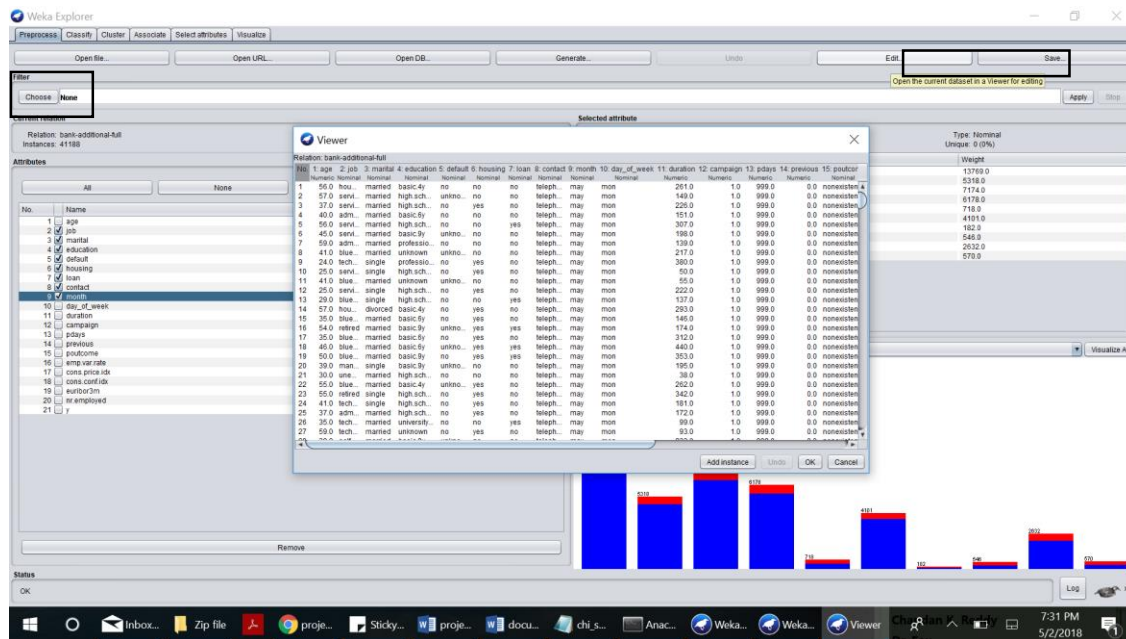
Meghana Gaopande    Pranavi Rambhakta
meghanag@vt.edu        pranavi7@vt.edu

In the project, Weka software is used for data-preprocessing, model building and model evaluation of the data while Matlab and Python (Jupyter Notebook web application) is used for exploratory analysis and visualization of the data.

The Weka explorer can be accessed from the weka homepage by clicking on the "Explorer" option.



There are direct navigation options available for Preprocess, Classify, Cluster, Associate, Select attributes and Visualize in which Preprocess page is the default. As shown below in the screenshot, the dataset can be loaded from the computer using the 'open file' option at the top left corner. The current dataset can be viewed for editing using the "Edit" option. The attributes are all listed and they can be selectively chosen individually and filters can be applied only on the attributes selected. The "filter" option is on the left top corner, just below the open file option. Additionally, the distribution of class labels for each attribute can be visualized in the form of histograms in the right bottom corner.
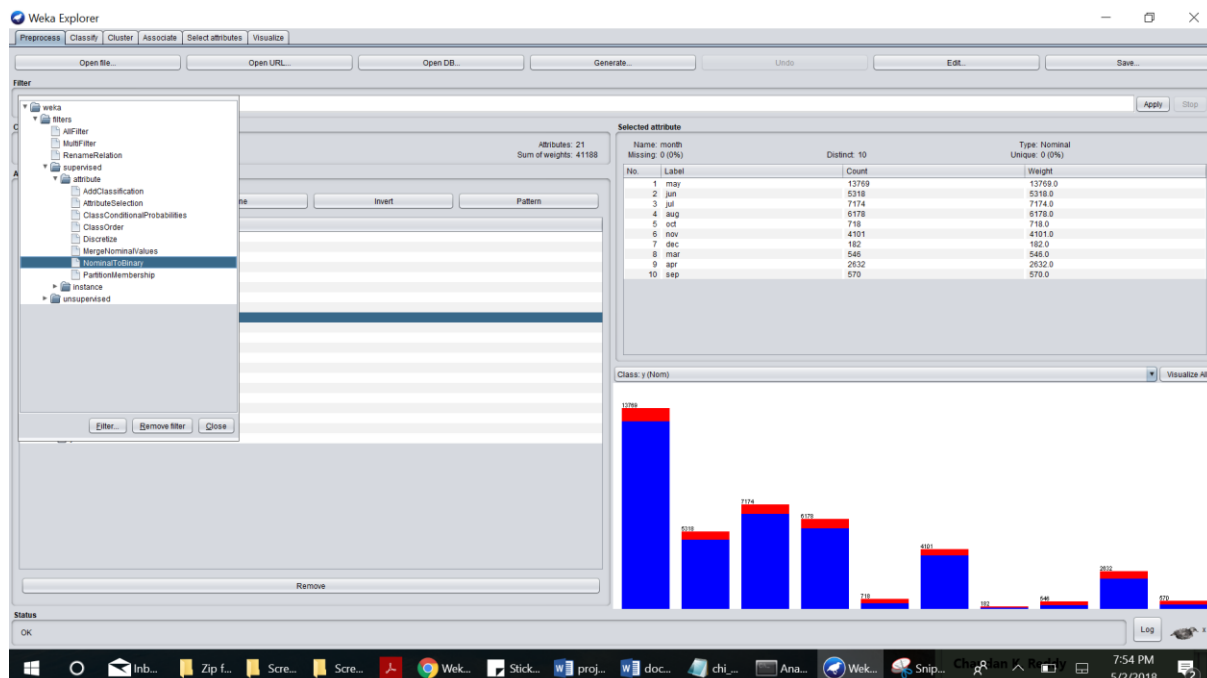
## Data Preprocessing:

Unordered categorical attributes can be selected and to perform One Hot Encoding on them, we used the filter: weka->supervised->NominalToBinary

Integer Encoding on ordinal categorical data can be performed by using the filter: weka->unsupervised->OrdinalToNumeric

Though weka is used for One Hot Encoding and Integer encoding, it is also coded in python (The code for is available in "data_exploration.py"- code folder)

The percentage of the missing values for each attribute in the data is shown in the form of a bar plot (The code for is available in "data_exploration.py"- code folder)
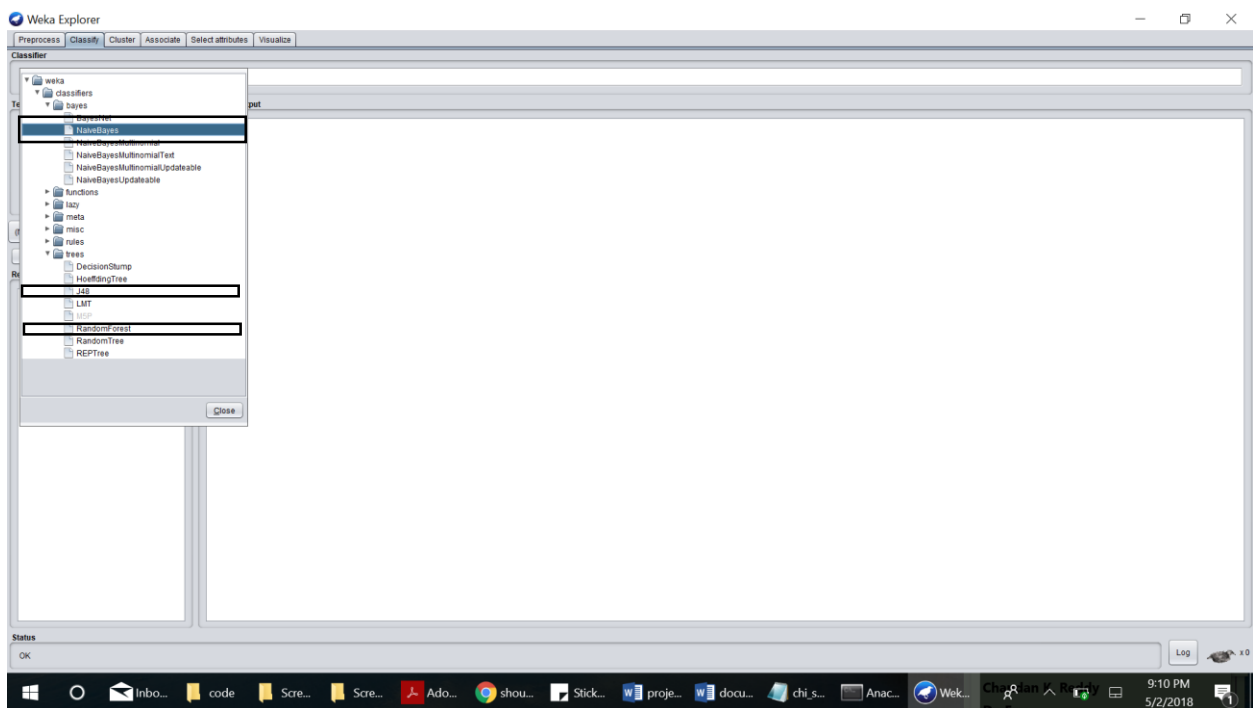
## Data Exploration:

The exploratory analysis and visualization of numeric data is shown in the form of correlation matrices and boxplots (The code is available in "box_and_corr_additional.m" and "box_and_corr_full.m"- code folder).
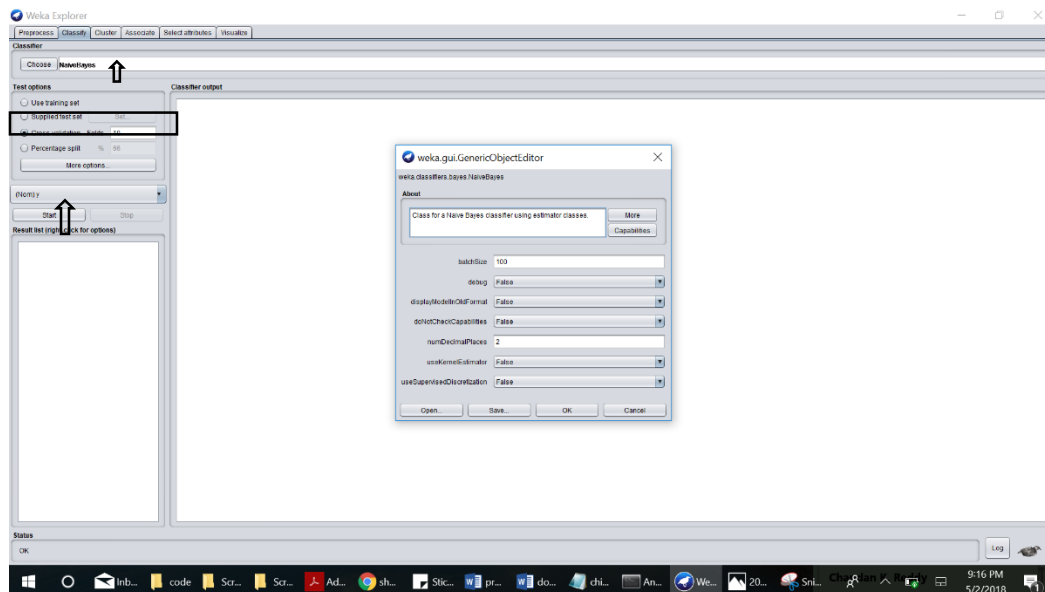
Additionally, the relation between nominal attributes is shown in the form of scatters plots and correlation plots. Chi-squared test is also performed to further analyze their correlation. (The code is available in "data_exploration.py"- code folder)
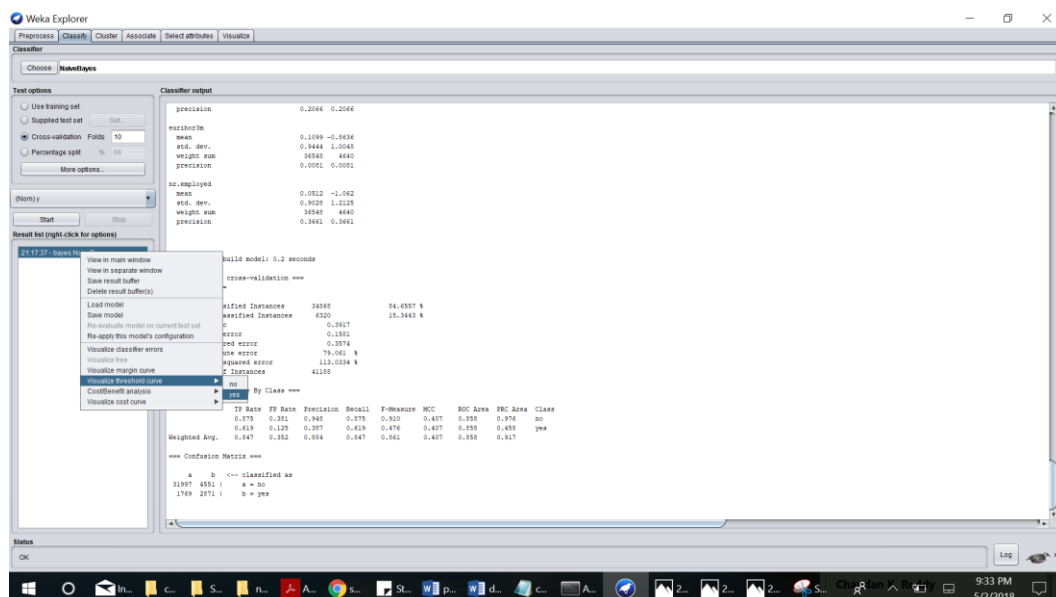
## Model Building:

Classification models can be chosen from the option "choose" in classify page. The drop down has all the model from which one can be selected. We have used three classification models: NaïveBayes, J48 and RandomForest. Decision tree is J48 in weka.

The parameters of the classification model can be specified or modified by clicking on the chosen model as shown below. Further, datasets for training and testing can also specified. We have used the default parameters and have chosen 10-fold cross validation. On clicking on the start, the classifier model starts running to produce output which will be displayed in the classifier output section. The result buffers are stored in the result list.



Apart from the output, results can be visualized by right-clicking on the respective result buffer and choosing from the options as shown below.

The results have Kappa statistic, confusion matrix and many other evaluation metrics listed. We have used Kappa statistic as our evaluation metric.