# Assignment 4
**Submission Date: Thursday, October 4, 2012**

**Learning Objectives:**

1. Gain insight into the concept of confidence intervals
2. Learn how to use Weka's Remove Folds filter and the Experimenter

(1) We used the SpeakerIDforLect data set for this assignment.

The 5-fold info is as follows.
Fold 1
Success rate: 19.8094%
Number of instances: 3988

Fold 2
Success rate: 22.1721%
Number of instances: 3987

Fold 3
Success rate: 23.3258%
Number of instances: 3987

Fold 4
Success rate: 27.4893%
Number of instances: 3987

Fold 5
Success rate: 23.702%
Number of instances: 3987

We compute the confidence interval on fold 1, with $Z = 1.65$, we have
Confidence interval: [0.1878878, 0.2087121]

With doubled data set, the confidence interval will be smaller, because with more data, we are more certain about the point estimate.

(2-3) For a multiple classification problem like this one, an instance could be incorrectly classified into any class other than the actual class label. For each actual class value in this question, we focus on the two cells in the confusion matrix with the most errors, because that is the main source of errors.

For each value of the class label, we list the features that possibly caused misclassification, because the average values of those features are very different between the two cells with most errors and other cells in the confusion matrix.

For class value irf03, feature 13, 14

For class value irf04, feature 1

For class value irf06, feature 16, 18, 19, 17, 23

For class value irf07, feature 23, 17, 20, 19, 18, 16
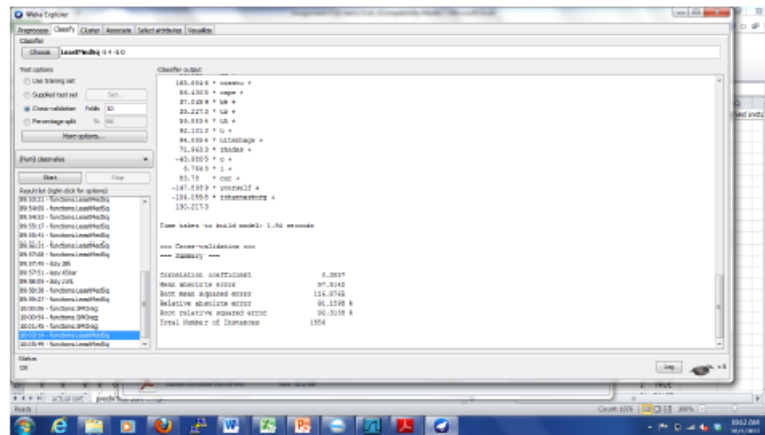
For class value irm02, feature 19, 20, 21, 23

For class value irm05, feature 19

For class value irm06, no apparent indicator since most of the instances were classified correctly
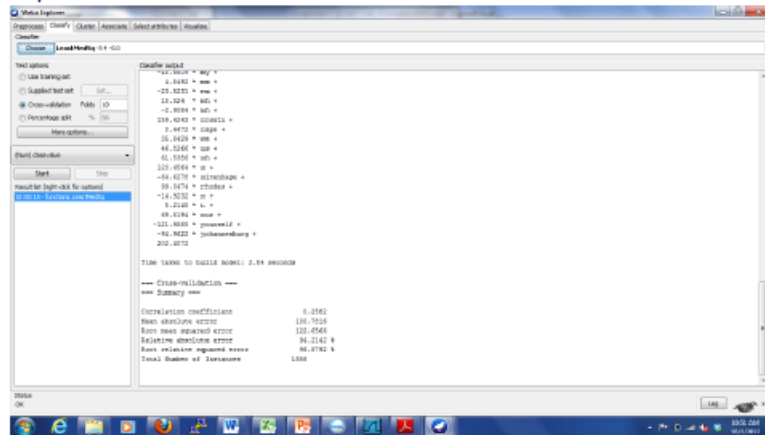
For class value irm07, feature 20

(4) Here is one example solution:

Using the normal data set, replacement preprocessing and the LinearMedSq classifier the correlation coefficient = .3857.



Simple data set

(1) Is the algorithm you picked significantly different from SimpleLinearRegression on the simple version of the dataset? The LinearMedSq algorithm does not perform significantly different from the SimpleLinearRegression algorithm. This is because the performance values fall within the confidence interval of the baseline. (2) And for the algorithm you picked, does its performance differ significantly between the two versions of the dataset? The performance on the simple data set performs more poorly.

These pages contain the comparison to two other algorithms. These results do not indicate a significant improvement over the Multilayer preceptor or the SimpleLinearRegression algorithm.

The results from the comparison between the two data sets does not show significant improvement between the two data sets since the confidence intervals overlap by .01.