# Solution to Laboratory Week 3

**1. The data files for this lab can be found at**
**/KDrive/SEH/SCSIT/Students/Courses/COSC2111/DataMining/data**

**2. Load the file arff/UCI/iris.arff.**

**(a) Go to the classifiers screen and select sepallength as the class attribute.**
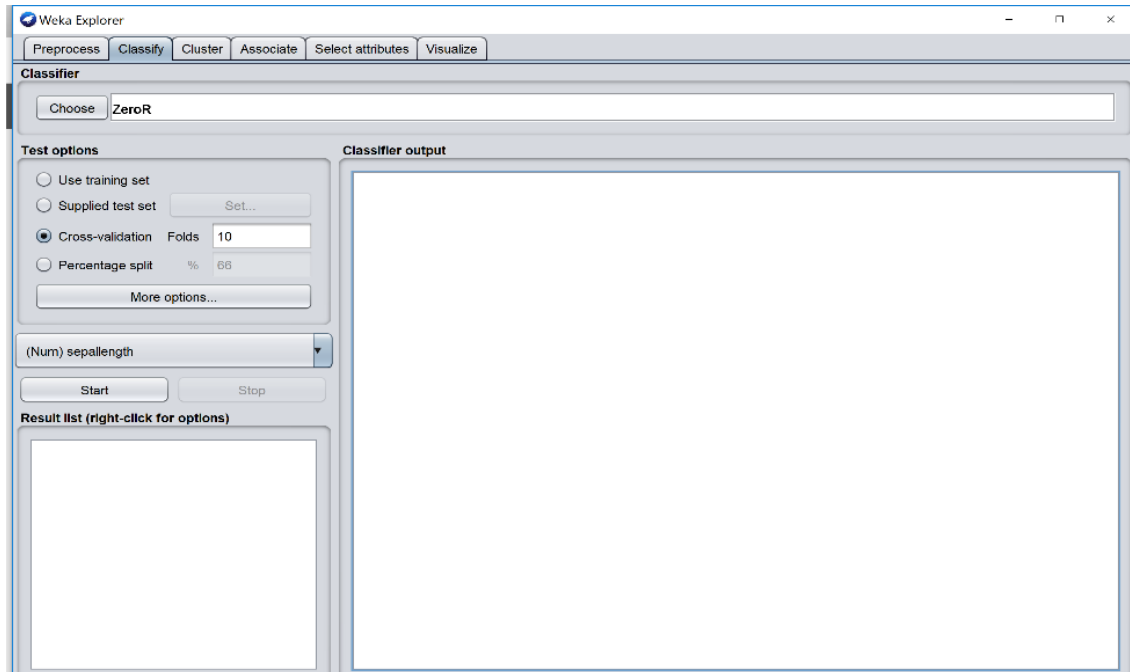


Figure 1. Screenshot of Weka on select sepallength as the class attribute.

**(b) Select 'More Options' and 'Output Predictions'.**


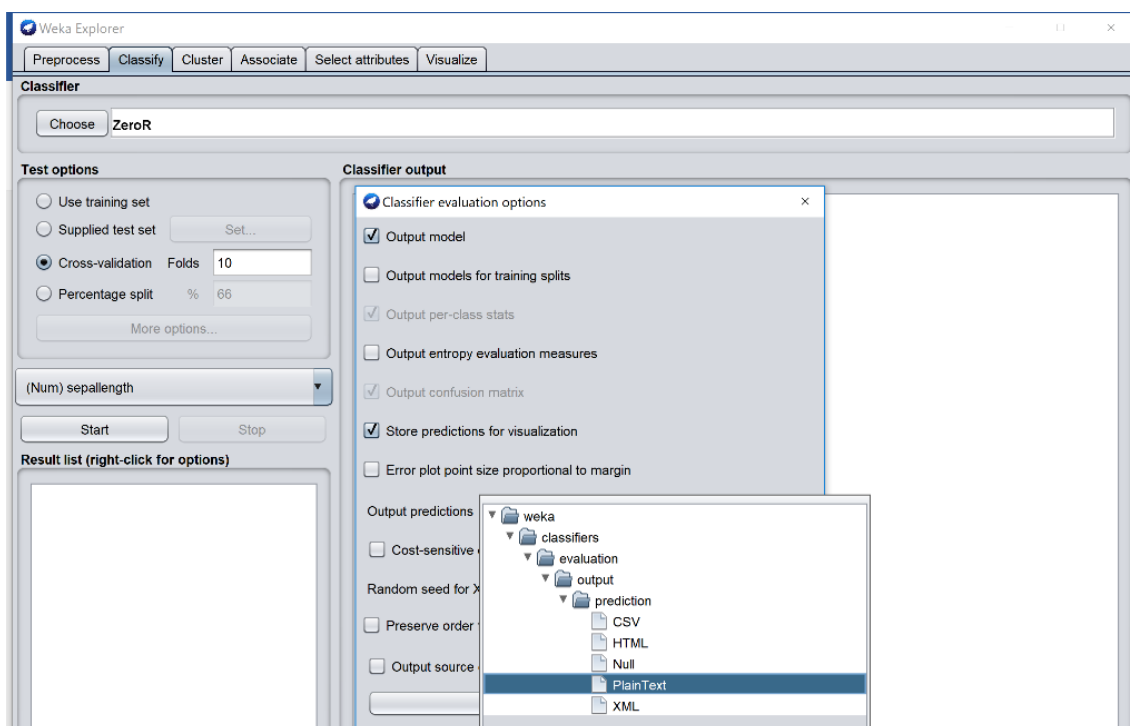
Figure 2. Screenshot of Weka on how to print 'output Predictions'.

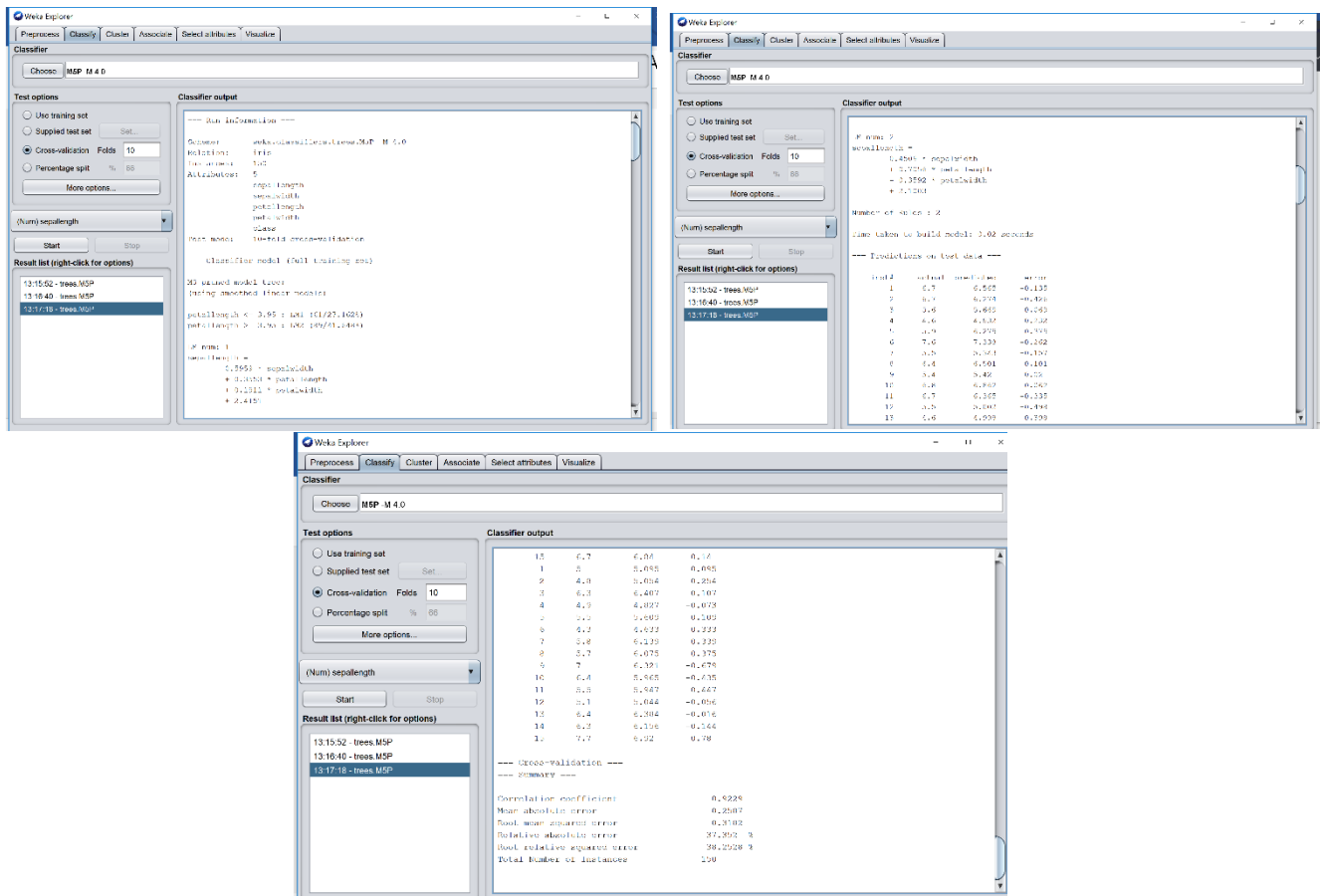## (c) Run the M5P classifier with default parameters.



Figure 3. Screenshots of a run of M5P classifier with default parameters.

## (d) Examine the output. What do you think of the accuracy of the predictions?

The error rate of M5-P classifier(classifier/ numeric predictor for numeric data set)  has been displayed in weka as :

Correlation coefficient            0.9229

Mean absolute error              0.2587

Root mean squared error         0.3182

Relative absolute error            37.352%

Root relative squared error       38.2528%

Total Number of Instances        150

A correlation coefficient above 9 indicates a very strong correlation.  In this case the 0.92 correlation between predicted and actual indicates good agreement which, indicates high accuracy.  Also the

values of petal length range from 4.3 to 7.9 with an average of about 6.  An average error of +/- 0.25 when the target is 6 looks pretty good to me.   The relative error looks surprisingly high, but we'll see late that Weka doesn't calculate it the way you might expect.

### (e) Experiment with different values for the parameters. What is the effect on accuracy?

By left-click on the name of classifier as shown in figure you would get access to a menu of all parameters of the classifier where these parameters can be edited to see their effects on outputs of the classifier.
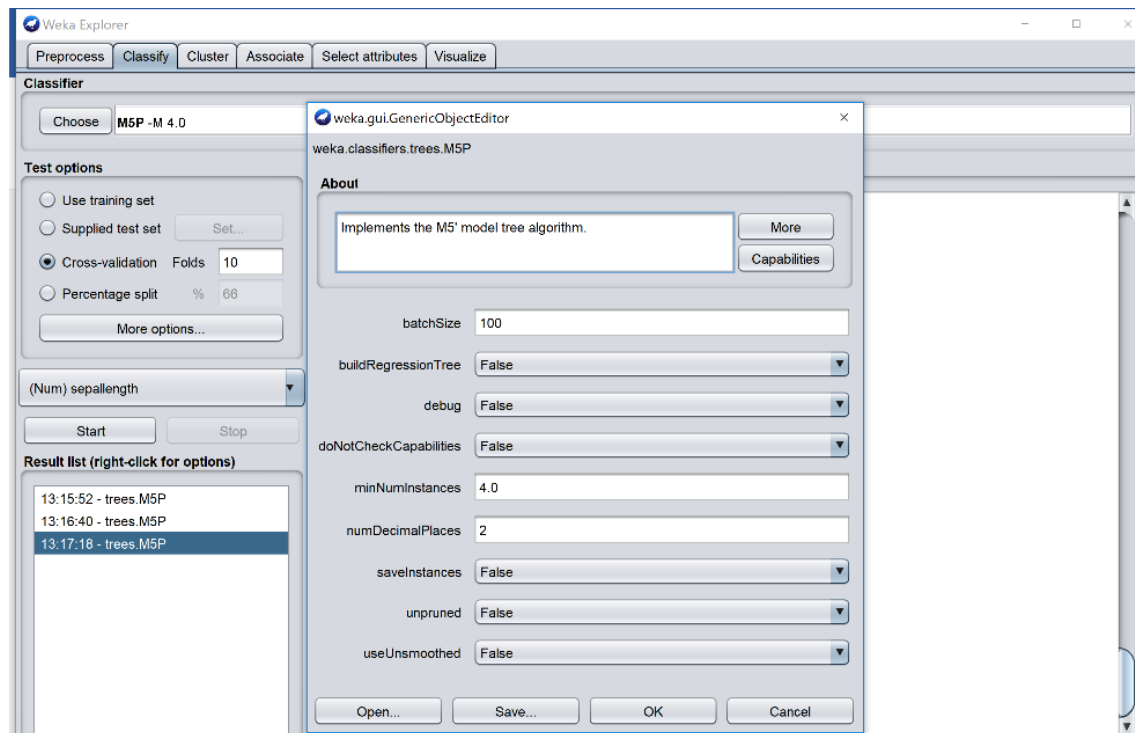


Figure 4. Screenshot of Weka on how to change parameters of M5-P.

Table1. Effect of varying batch-size on running time and error rate of M5-P classifier.

| Batchsize | Time taken to build model (second) | Error rate (relative absolute error) |
|---|---|---|
| 1000 | 0.01 | 37.352 |
| 500 | 0 | 37.352 |
| 200 | 0.01 | 37.352 |
| 100 | 0.02 | 37.352 |
| 50 | 0 | 37.352 |
| 30 | 0.01 | 37.352 |
| 20 | 0 | 37.352 |
| 10 | 0 | 37.352 |
| 5 | 0 | 37.352 |
| 2 | 0.01 | 37.352 |
| 1 | 0 | 37.352 |

As it can be seen from table 1, changing batchzise does not influence error rate of the classifier while slightly changes time taken on building model.

Table 2. Effect of varying the "M" parameter on running time and error rate of M5-P classifier.

| "M" parameter | Time taken to build model (seconds) | Error rate (relative absolute error) |
|---|---|---|
| 1 | 0.01 | 37.352 |
| 2 | 0.01 | 37.352 |
| 3 | 0.01 | 37.352 |
| 4 | 0.01 | 37.352 |
| 5 | 0.01 | 37.352 |
| 6 | >0.01 | 37.352 |
| 7 | >0.01 | 37.352 |
| 8 | >0.01 | 37.352 |
| 10 | >0.01 | 37.352 |
| 15 | >0.01 | 37.6953 |
| 20 | >0.01 | 37.8761 |
| 30 | >0.01 | 38.504 |

Table 2 displays results from changing the "M" parameter (ie. Altering the minimum number of instances in the tree's leaves). The effect on model build time is minor yet the larger the "M" parameter the faster the model can be built. The relative absolute error remains stable at around 37.352 until the "M" parameter increases beyond 10, whereby we see the error rate increase.

**(f) Experiment with ZeroR and IBK and their various parameters.**

Table 3. Effect of varying batchsize and K on accuracy and running time of zeroR and IBK.

| Classifier Model | ZeroR | | IBK | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | K=1 | | K=5 | | K=50 | | K=100 | |
| Batch size | Time taken (second) | Error rate (%relative absolute e) | Time taken (second) | Error rate (%relative absolute e) | Time taken (second) | Error rate (%relative absolute e) | Time taken (second) | Error rate (%relative absolute e) | Time taken (second) | Error rate (%relative absolute e) |
| 1 | 0 | 100 | 0 | 45.5808 | 0 | 42.6264 | 0 | 57.6746 | 0 | 74.0479 |
| 5 | 0 | 100 | 0 | 45.5808 | 0 | 42.6264 | 0 | 57.6746 | 0 | 74.0479 |
| 10 | 0 | 100 | 0 | 45.5808 | 0 | 42.6264 | 0 | 57.6746 | 0 | 74.0479 |
| 20 | 0 | 100 | 0 | 45.5808 | 0 | 42.6264 | 0 | 57.6746 | 0 | 74.0479 |
| 50 | 0 | 100 | 0 | 45.5808 | 0 | 42.6264 | 0 | 57.6746 | 0 | 74.0479 |
| 100 | 0 | 100 | 0 | 45.5808 | 0 | 42.6264 | 0 | 57.6746 | 0 | 74.0479 |
| 200 | 0 | 100 | 0 | 45.5808 | 0 | 42.6264 | 0 | 57.6746 | 0 | 74.0479 |
| 500 | 0 | 100 | 0 | 45.5808 | 0 | 42.6264 | 0 | 57.6746 | 0 | 74.0479 |

As it can be seen from table 2, changing batchsize does not have any affect on accuracy and running time of IBK and zeroR classifiers. However, accuracy of IBK is still significantly higher than zeroR. Increasing value of K in IBK first leads to improvement in accuracy. However if K is too big the accuracy suffers.

**(g) Build a table of classifier, parameter values and error. What combination gives the most accurate predictions?**

Looking at table 1 and comparing these limited number of runs based on varying batchsize and K parameters, IBK with k=5 gives the most accurate prediction. However we don't know yet whether there is a better value for k between 1 and 50.

**(h) Can you explain the differences in errors?**

The more complex model is more accurate.

**3. Repeat the previous exercise with cpu.with.vendor.arff**

**4. Load the file soybean.arff**

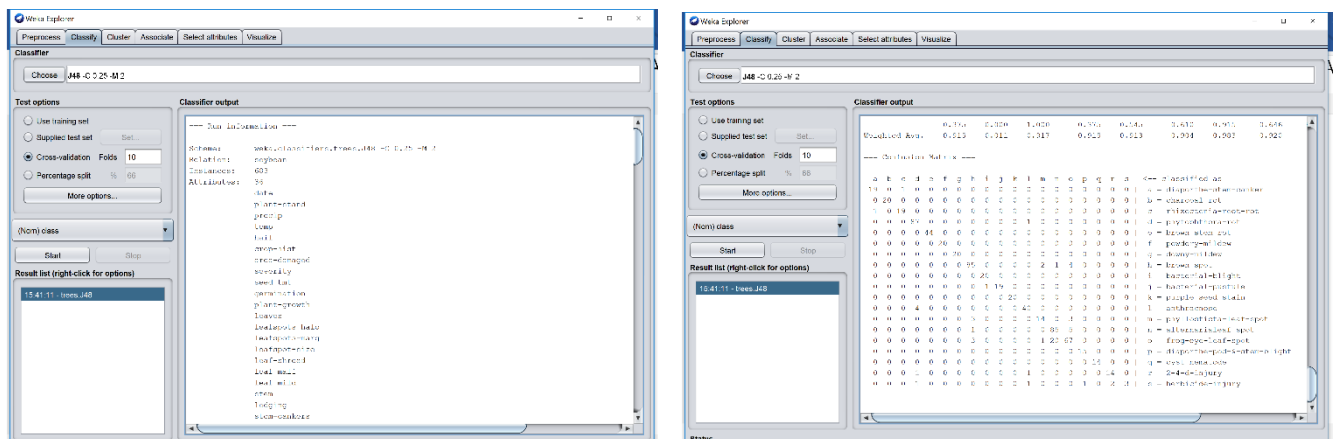**(a) Run the J48 classifier with default parameters.**



Figure5. Screenshots of running J48 classifier on soybean data set.

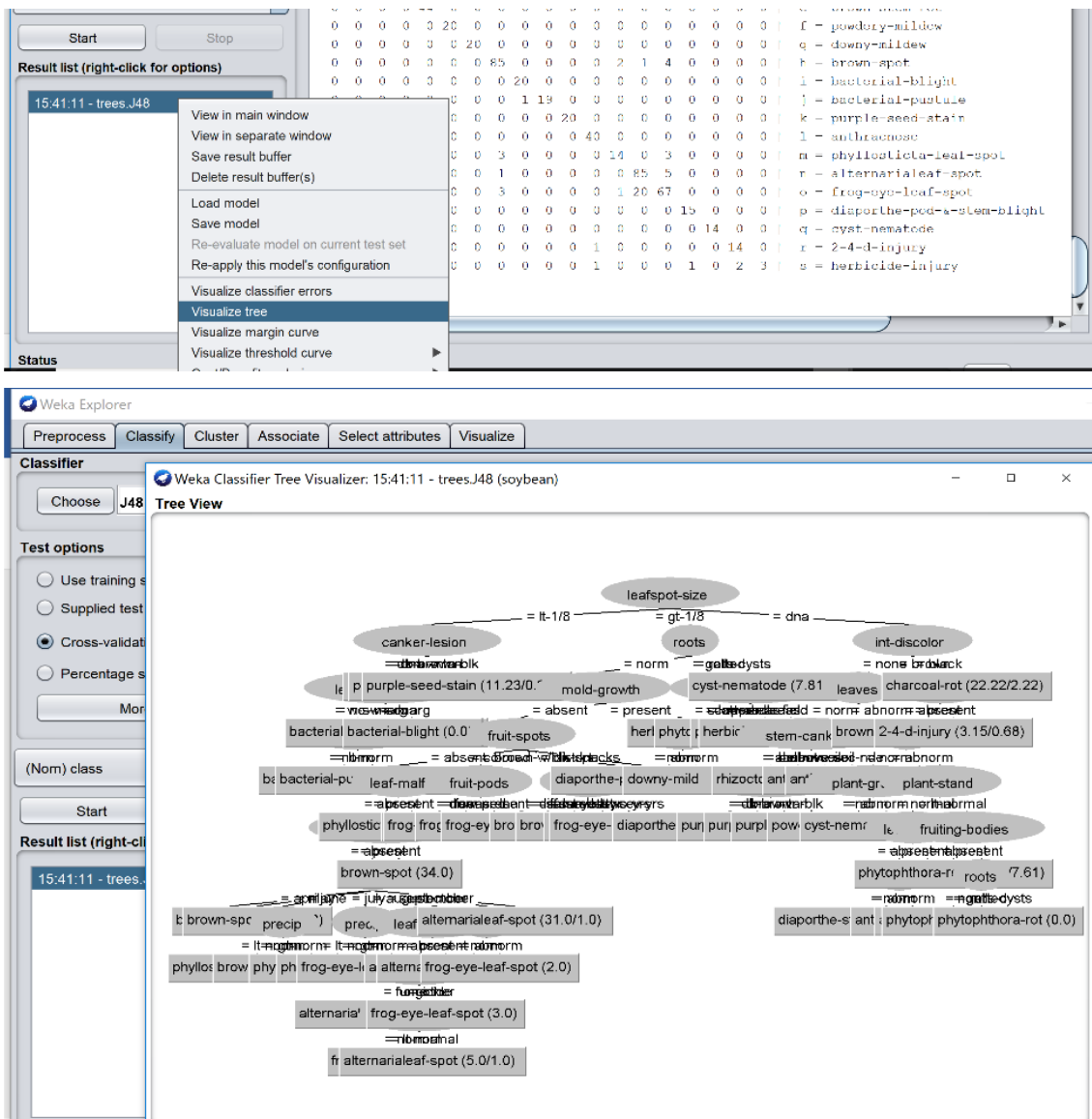**(b) Make sure you can visualise the tree.**



Figure 6. Screenshots of visualising J48 tree.

**(c) Experiment with different values of the C and M parameters.**

Table 4. Effects of varying M and C parameters on accuracy and running time of J48.

| C | M | Time Taken | Accuracy |
|---|---|---|---|
| 0.25 | 2 | 0.03 | 91.5081 |
| | 3 | 0 | 89.3119 |
| | 5 | 0 | 89.6047 |
| | 10 | 0 | 84.4802 |
| | 20 | 0 | 78.4773 |
| | 50 | 0 | 54.3192 |
| 1 | 2 | 0.63 | 91.5081 |
| | 3 | 0.56 | 89.6047 |
| | 5 | 0.44 | 88.1406 % |
| | 10 | 0.28 | 84.7731 |
| | 20 | 0.19 | 78.4773 |
| | 50 | 0.05 | 54.3192 |

| | 2 | 0.72 | 91.5081 |
|---|---|---|---|
| | 3 | 0.6 | 89.6047 |
| 10 | 10 | 0.43 | 88.1406 |
| | 20 | 0.13 | 78.4773 |
| | 50 | 0.02 | 54.3192 |

As shown in table 3, accuracy of J48 is decreasing by increasing M, while increasing confidence C does not have any affect on accuracy of classifier. Therefore, the best combination of c and M is when M = 2 and c = 0.25.

**(d) Using percentage split build a table of training and test errors**

| # | Percentage Split | Test Accuracy | Train Accuracy |
|---|---|---|---|
| 1 | %66 | 90.5172 | 96.3397 |
| 2 | %40 | 85.3659 | 96.3397 |
| 3 | %80 | 90.5109 | 96.3397 |

**(e) Is there any overfitting?**

In all runs, accuracy of training is greater than test error therefore there is overfitting.

**(f) What would you say is the best combination of parameter values**

The best combination of c and M is when M = 2 and c = 0.25.

**5. Repeat the previous exercise with the glass.arff**

**6. Experiment with other classifiers and data files.**