# Solution to Laboratory Week 6.

**1. You will need to have access to the WEKA package.**

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

**3. Load the file /arff/UCI/ionosphere.arff**

**(a) Get the classification accuracy with J48.**

The classification accuracy of J48, using 10-fold cross-validation on ionosphere data set is as follows:

Total number of instances: 351

Correctly Classified Instances     321        91.453 %

Incorrectly Classified Instances    30        8.547 %

**(b) Apply attribute selection with the default settings, ie CfsSubsetEval and BestFirst. Go back to Preprocess, remove all but the selected attributes and rerun J48. What is the accuracy?**
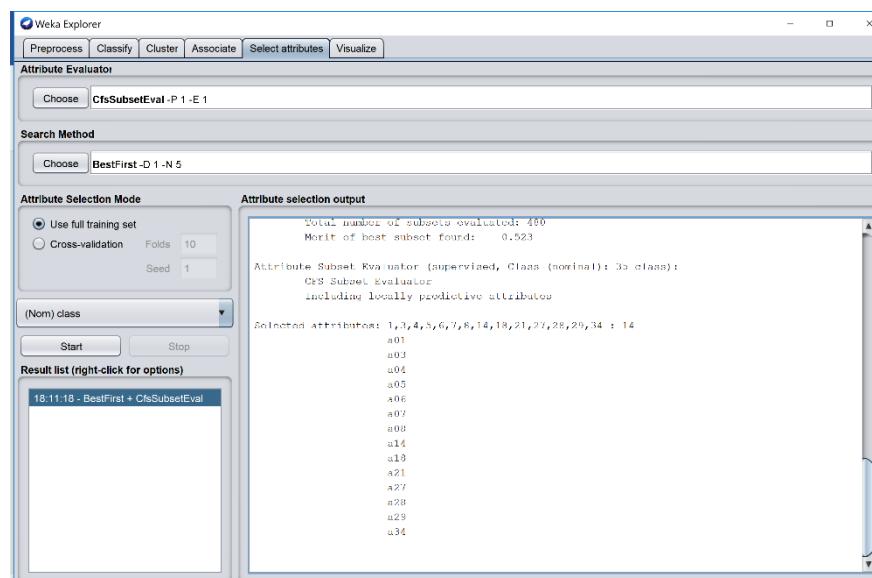


Figure 1. List of selected attributes by running CfsSubsetEval and BestFirst attribute selection algorithms.
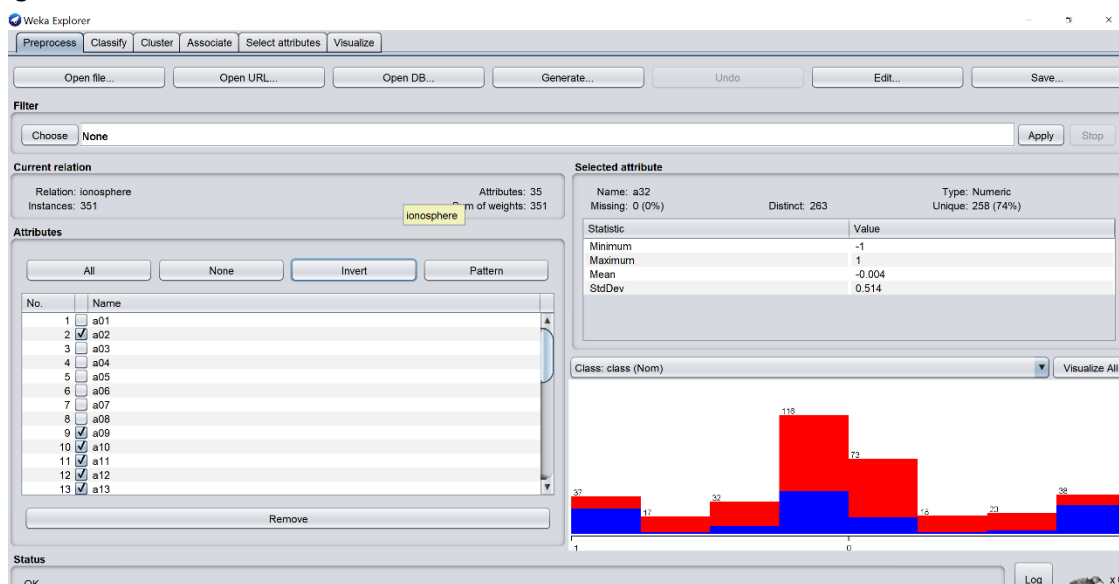
Figure 2.  Remove all but the selected attributes from preprocess tab by clicking on Remove button.

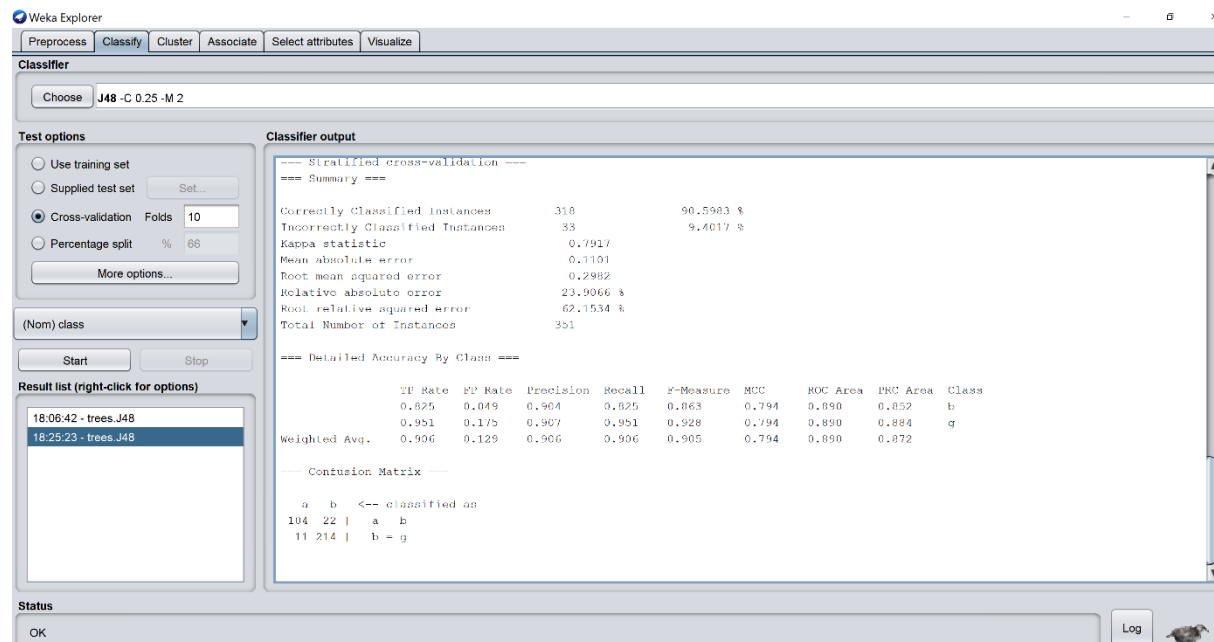Be careful do not remove class label.



Figure 3. Output of running J48 after attribute selection.

As it can be seen from figure 3, accuracy of J48 is reducing to 90.593% after attribute selection.

[Useful hint] Go to Preprocess -->Filters--> Supervised -Attribute and select the AttributeSelection filter. Using this filter (with Undo) will save the tedious task of manually selecting the results of attribute selection.

**(c) Reload ionosphere.arff Apply attribute selection with WrapperSubsetEval, BestFirst and J48 as the classifier in WrapperSubsetEval. Go back to Preprocess, remove all but the selected attributes and rerun J48. What is the accuracy?**
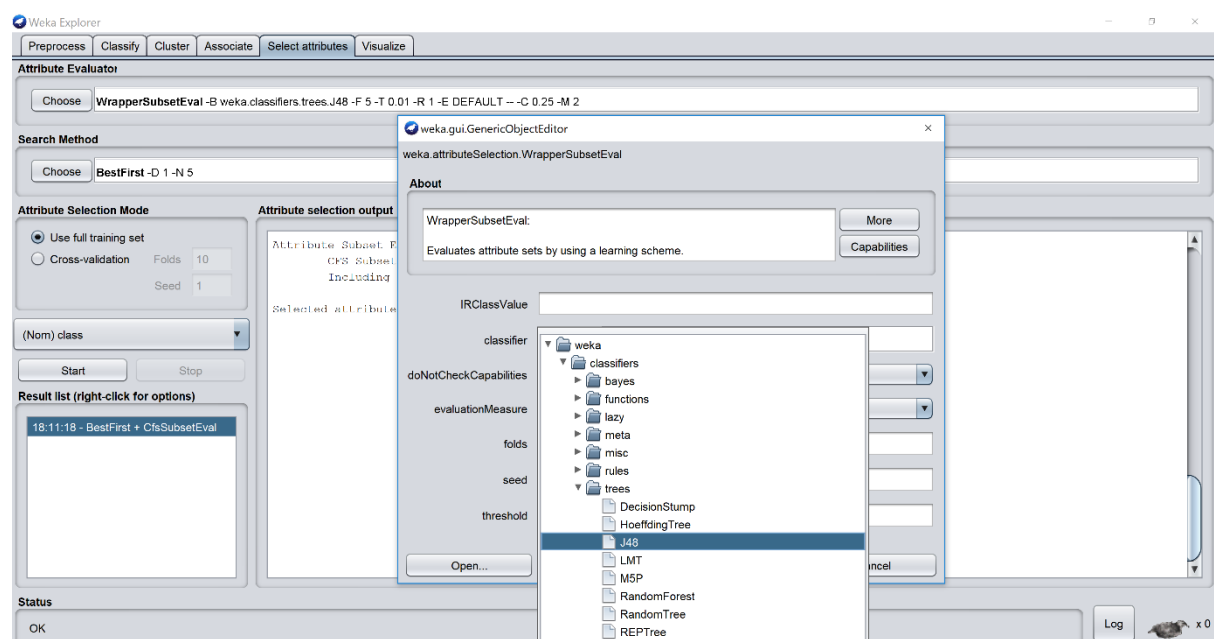
Figure 4. Snapshot of setting attribute selection algorithm to WrapperSubsetEval and setting its classifier parameter to J48.
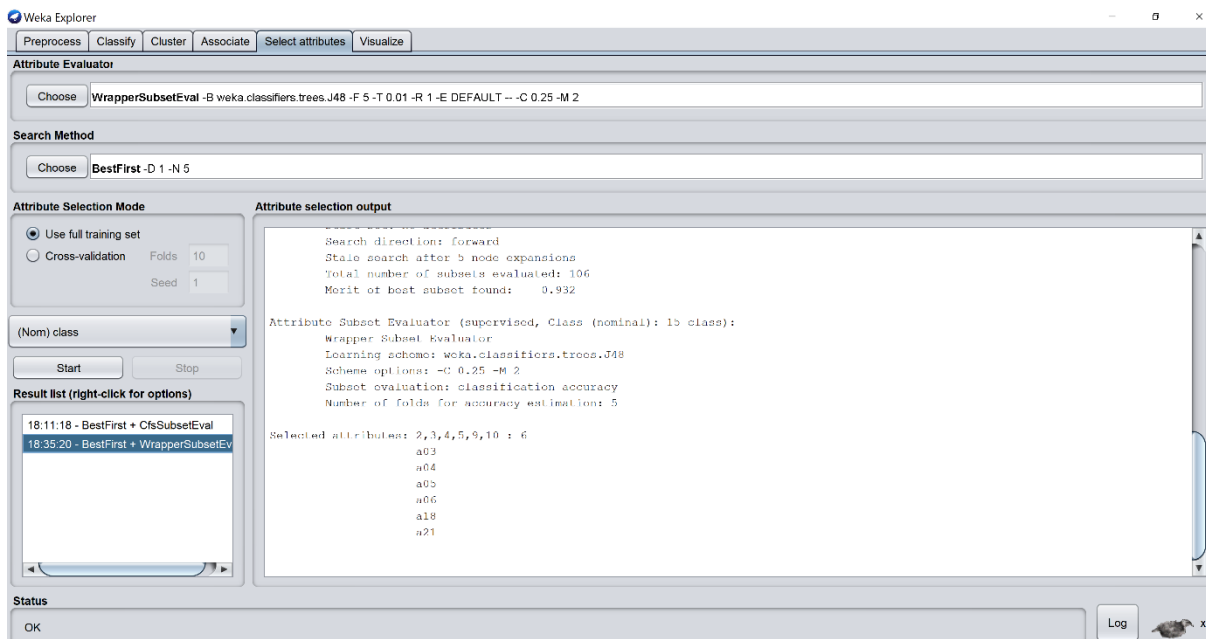


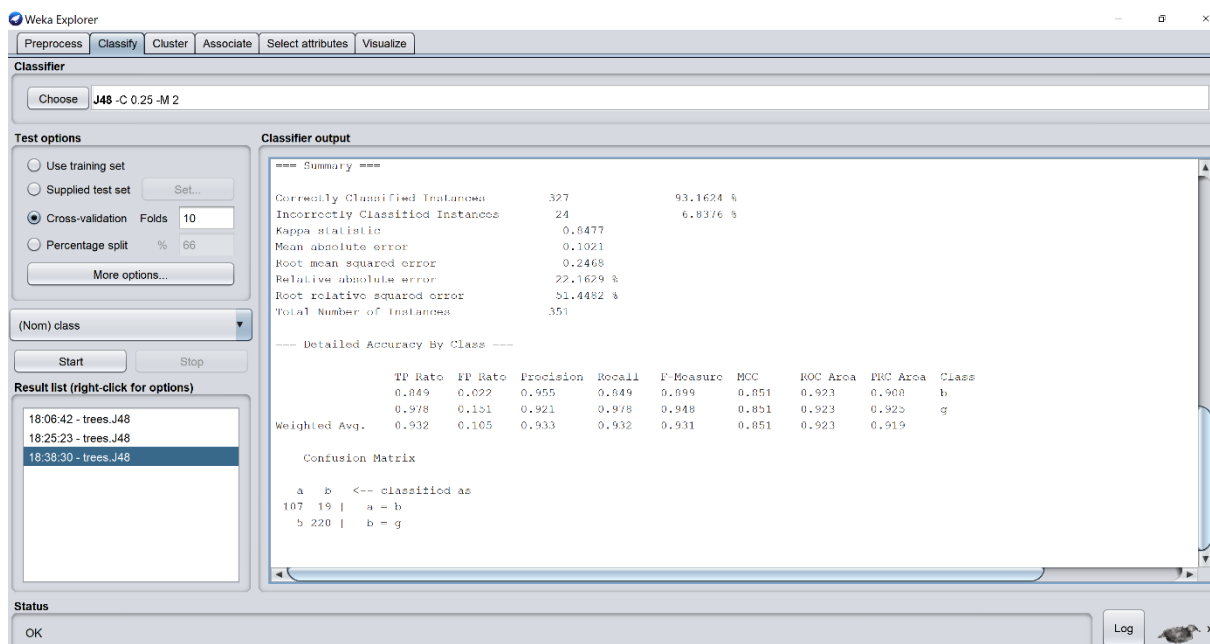Figure 5. Output of attribute selection from WrapperSubsetEval.



Figure 6. Output of running J48 after wrapper attribute selection from figure 5.

As it can be seen from figure 6, the accuracy of classifier is increased to 93.1624% by using wrapper attribute selection algorithm.

**(d) What do you conclude about the value of attribute selection?**

Using attribute selection algorithm could improve accuracy of classifier or make it worse. In this case, using the CfsSubsetEval filter results in fewer attributes, but lower accuracy. Using the wrapper method results in fewer attributes and higher accuracy.

**(e) Explore other combinations of evaluator and search method. Can you find anything better?**

**4. The file /arff/UCI/isolet.arff has 618 attributes. Explore a variety of attribute selection techniques to reduce the number of attributes without reducing accuracy. What is your best result?**
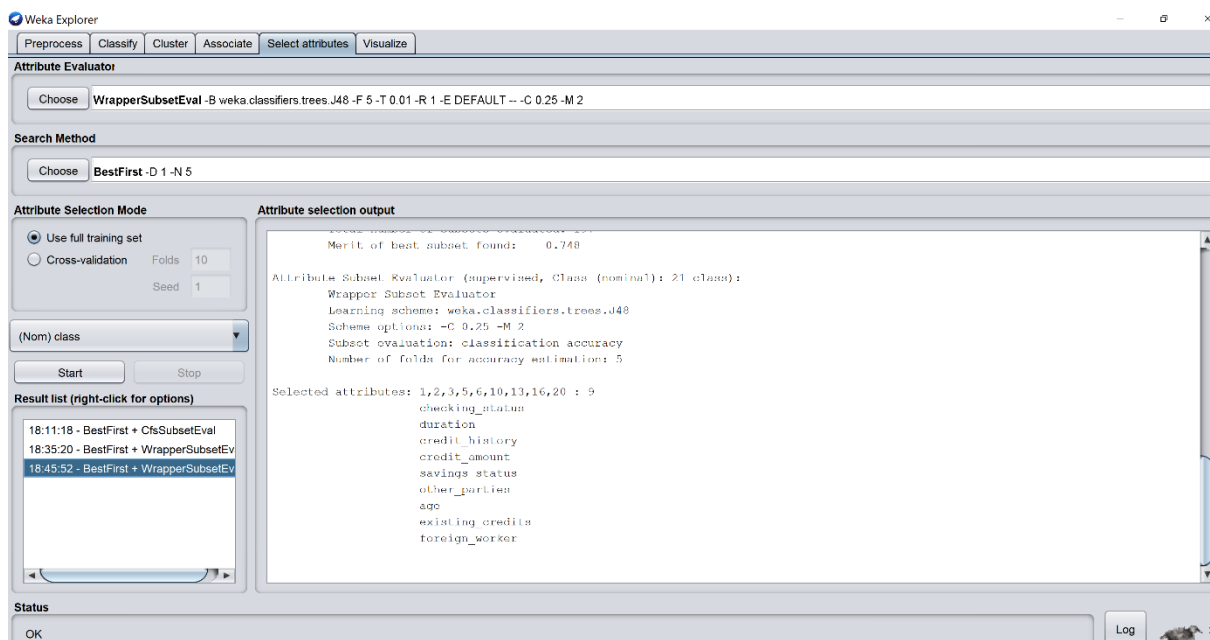
Using isolet1+2+3+4.arff

```
Evaluator                   Search        Num Features      Accuracy
Full Feature set                               617          82.27
CfssubsetEval               Best First         190          83.39
CfssubsetEval               Best First         191          83.42
CorrelationAttrEval         Ranker             190          80.15
GainRatio                   Ranker             200          81.65
InfoGain                    Ranker             200          78.83
OneR                        Ranker             200          79.25
Wrapper J48, M=10           Best First         ???          ??.??
```

Wrapper has run for 24 hours and didn't finish.

Results so far indicate that slightly better accuracy can be achieved with a reduced number of attributes. There is a fair variation in the performance of the attribute selection methods.

**5. Load the file /arff/UCI/credit-g.arff**

**(a) Repeat (3) on this file.**

## Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

**Classifier**

Choose | J48 -C 0.25 -M 2

**Test options**

- Use training set
- Supplied test set | Set...
- Cross-validation | Folds | 10
- Percentage split | % | 66

More options...

(Nom) class

Start | Stop

**Result list (right-click for options)**

18:06:42 - trees.J48
18:25:23 - trees.J48
18:38:30 - trees.J48
18:49:56 - trees.J48

**Classifier output**

```
--- Stratified cross-validation ---
--- Summary ---

Correctly Classified Instances         705              70.5   %
Incorrectly Classified Instances       295              29.5   %
Kappa statistic                          0.2467
Mean absolute error                      0.3467
Root mean squared error                  0.4796
Relative absolute error                 82.5233 %
Root relative squared error            104.6565 %
Total Number of Instances             1000

--- Detailed Accuracy By Class ---

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
                 0.840    0.610    0.763      0.840   0.799      0.251  0.639     0.746     good
                 0.390    0.160    0.511      0.390   0.442      0.251  0.639     0.449     bad
Weighted Avg.    0.705    0.475    0.687      0.705   0.692      0.251  0.639     0.657

--- Confusion Matrix ---

   a    b   <-- classified as
 588  112 |   a = good
 183  117 |   b = bad
```

**Status**

OK

Log | x 0

---

## Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

**Attribute Evaluator**

Choose | CfsSubsetEval -P 1 -E 1

**Search Method**

Choose | BestFirst -D 1 -N 5

**Attribute Selection Mode**

- Use full training set
- Cross-validation | Folds | 10
- Seed | 1

(Nom) class

Start | Stop

**Result list (right-click for options)**

18:11:18 - BestFirst + CfsSubsetEval
18:35:20 - BestFirst + WrapperSubsetEv
18:45:52 - BestFirst + WrapperSubsetEv
18:52:31 - BestFirst + CfsSubsetEval

**Attribute selection output**

```
=== Attribute Selection on all input data ===

Search Method:
        Best first.
        Start set: no attributes
        Search direction: forward
        Stale search after 5 node expansions
        Total number of subsets evaluated: 134
        Merit of best subset found:    0.076

Attribute Subset Evaluator (supervised, Class (nominal): 21 class):
        CFS Subset Evaluator
        Including locally predictive attributes

Selected attributes: 1,2,3 : 3
                     checking_status
                     duration
                     credit_history
```

**Status**

OK

Log | x

Using credit-g.arff

| Evaluator | Search | Num Features | Accuracy |
|---|---|---|---|
| Full Feature set | | 20 | 70.5 |
| CfssubsetEval | Best First | 3 | 70.5 |
| Wrapper J48 | Best First | 9 | 70.5 |

The same accuracy is achieved with 3, 9 and 20 features. Looking at the file below, the 9 features look reasonable. Not sure what to make of the 3 features. Maybe an artefact of the data.

**(b) Open the file with a text editor and read the descriptions of the attributes. Does the set of selected attributes make sense from what you know about credit ratings?**