

COMP30120 Tutorial

Dimension Reduction and Feature Selection

Derek Greene

**School of Computer Science and Informatics
Autumn 2015**



Tutorial Q1(a)

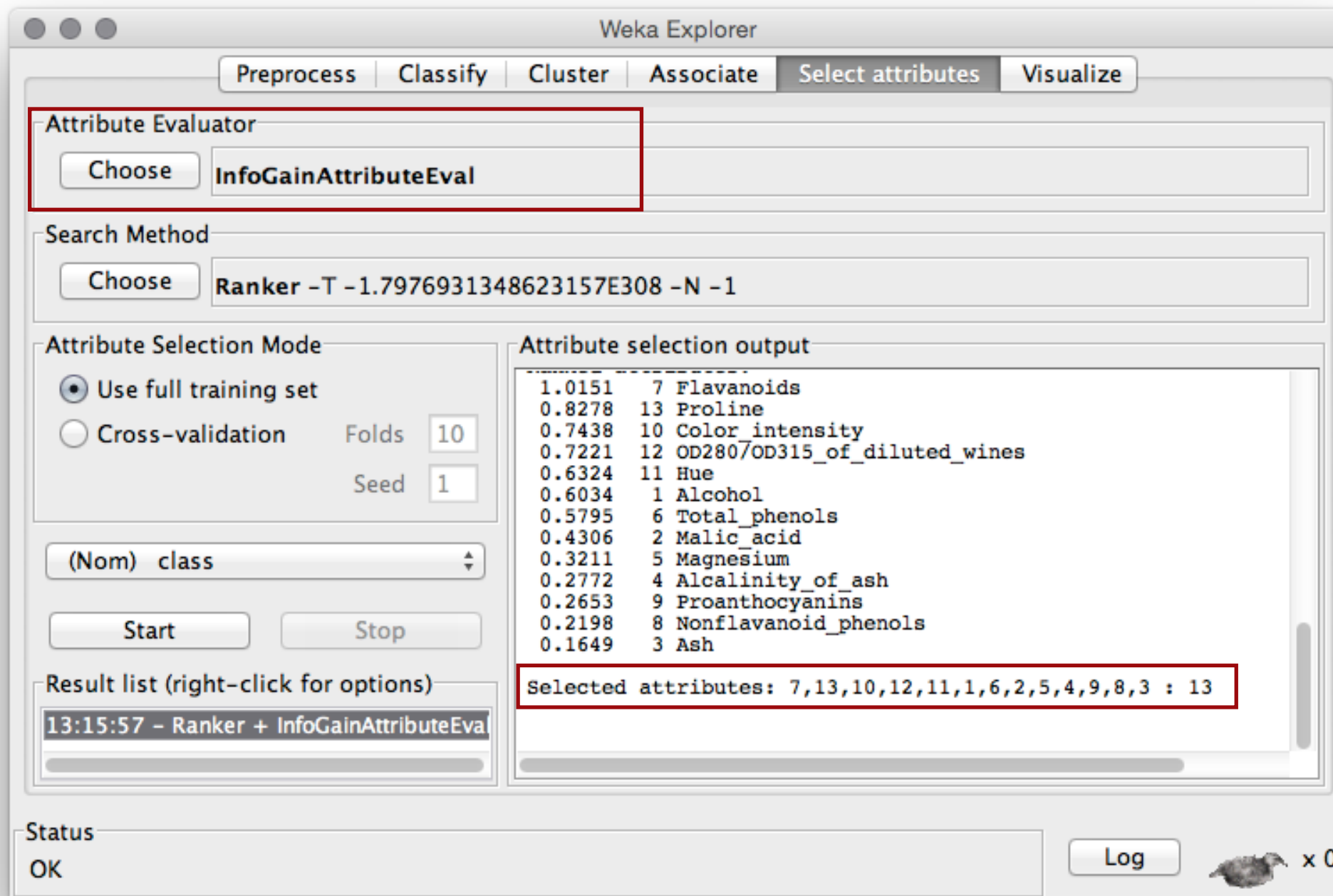
In Weka, use filter-based feature selection with Information Gain to identify the 3 most discriminating and 3 least discriminating features in the *Wine* data set in the ARFF file provided.

Assess the accuracy of a 1-nearest neighbour classifier with:

- (i) only the 3 most discriminating features included.
- (ii) only the 3 least discriminating features included.

Tutorial Q1(a)

In Weka *Select attributes* tab, choose *InfoGainAttributeEval* as the evaluator, *Ranker* as the method.



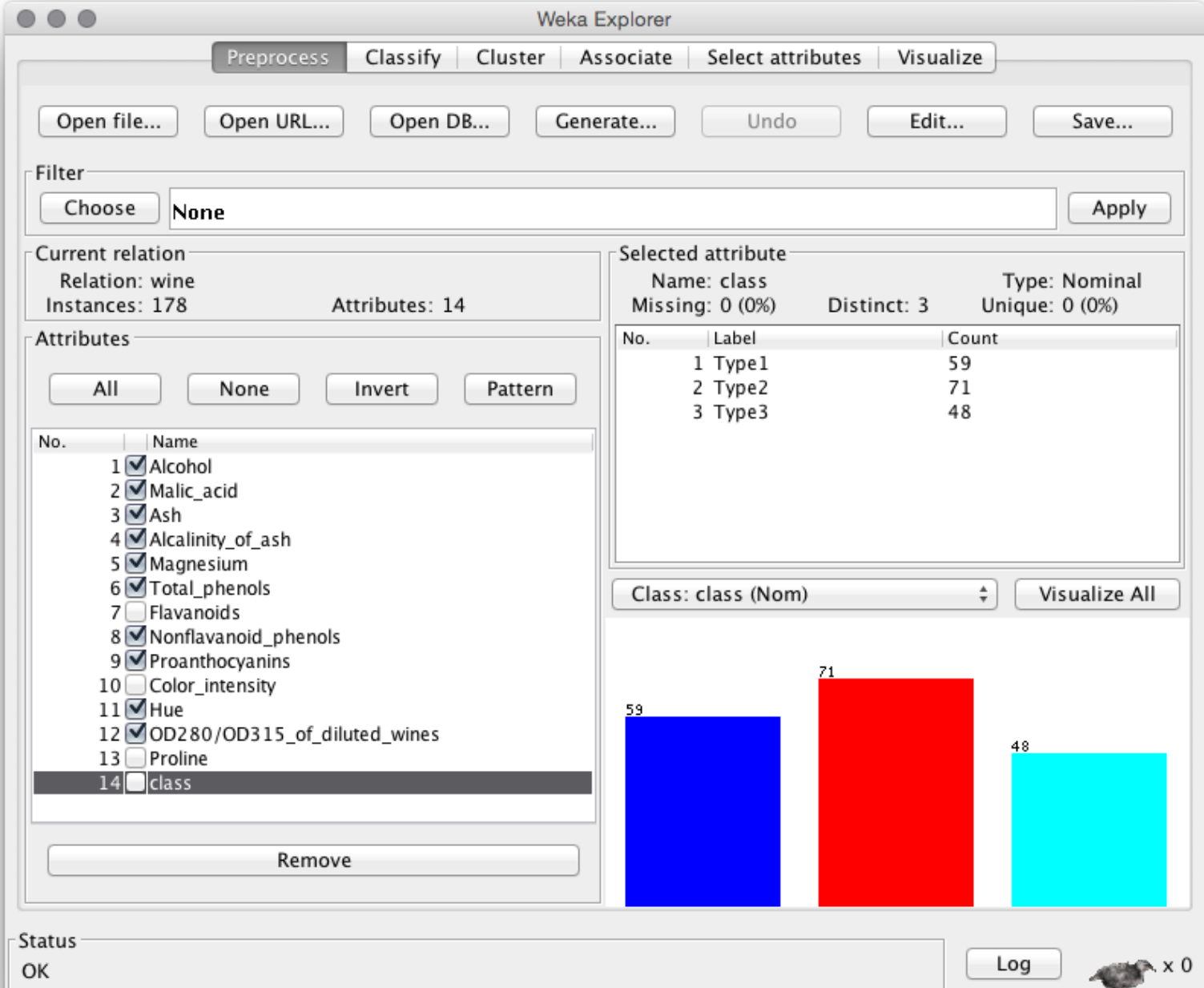
Tutorial Q1(a)

Assess the accuracy of a 1-nearest neighbour classifier with only the 3 most discriminating features included.

Most discriminating:
7, 13, 10

In the *Preprocess* tab,
remove the unwanted
features.

NB: Keep the “class”
feature!



The screenshot shows the Weka Explorer interface with the 'Preprocess' tab selected. The 'Current relation' is 'wine' with 178 instances and 14 attributes. The 'Attributes' list shows 14 attributes, with 'class' selected. The 'Selected attribute' section shows 'class' with 3 distinct values: Type1 (59), Type2 (71), and Type3 (48). A bar chart at the bottom right visualizes these counts.

No.	Label	Count
1	Type1	59
2	Type2	71
3	Type3	48

Class: class (Nom) Visualize All

Tutorial Q1(a)

Assess the accuracy of a 1-nearest neighbour classifier with only the 3 most discriminating features included.

Re-run the 1NN classifier with the new feature subset.

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier: Choose IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"weka.core.EuclideanDistance -R first-las

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)

- 13:45:05 - lazy.IBk

Classifier output

Time taken to build model: 0 seconds

=== Stratified cross-validation ===

=== Summary ===

Correctly Classified Instances	170	95.5056 %
Incorrectly Classified Instances	8	4.4944 %
Kappa statistic	0.9319	
Mean absolute error	0.0376	
Root mean squared error	0.1717	
Relative absolute error	8.5601 %	
Root relative squared error	36.651 %	
Total Number of Instances	178	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	RC
	0.983	0.034	0.935	0.983	0.959	
	0.915	0.019	0.97	0.915	0.942	
	0.979	0.015	0.959	0.979	0.969	
Weighted Avg.	0.955	0.023	0.956	0.955	0.955	

=== Confusion Matrix ===

a	b	c	<-- classified as
58	1	0	a = Type1
4	65	2	b = Type2
0	1	47	c = Type3

Status: OK

Log x 0

Tutorial Q1(a)

Assess the accuracy of a 1-nearest neighbour classifier with only the 3 least discriminating features included.

Least discriminating:
3, 8, 9

Load the ARFF file
again, remove the
unwanted features.

Re-run the 1NN
classifier with the new
feature subset.

Weka Explorer

Preprocess **Classify** Cluster Associate Select attributes Visualize

Classifier: Choose IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"weka.core.EuclideanDistance -R first-last\""

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):

- 13:48:00 - lazy.IBk

Classifier output:

using 1 nearest neighbour(s) for classification

Time taken to build model: 0 seconds

=== Stratified cross-validation ===

=== Summary ===

Correctly Classified Instances	110	61.7978 %
Incorrectly Classified Instances	68	38.2022 %
Kappa statistic	0.422	
Mean absolute error	0.2582	
Root mean squared error	0.5001	
Relative absolute error	58.8025 %	
Root relative squared error	106.7339 %	
Total Number of Instances	178	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area
	0.644	0.235	0.576	0.644	0.608	0.71
	0.521	0.252	0.578	0.521	0.548	0.638
	0.729	0.1	0.729	0.729	0.729	0.803
Weighted Avg.	0.618	0.206	0.618	0.618	0.617	0.706

=== Confusion Matrix ===

a	b	c	<-- classified as
38	20	1	a = Type1
22	37	12	b = Type2
6	7	35	c = Type3

Status: OK

Log x 0

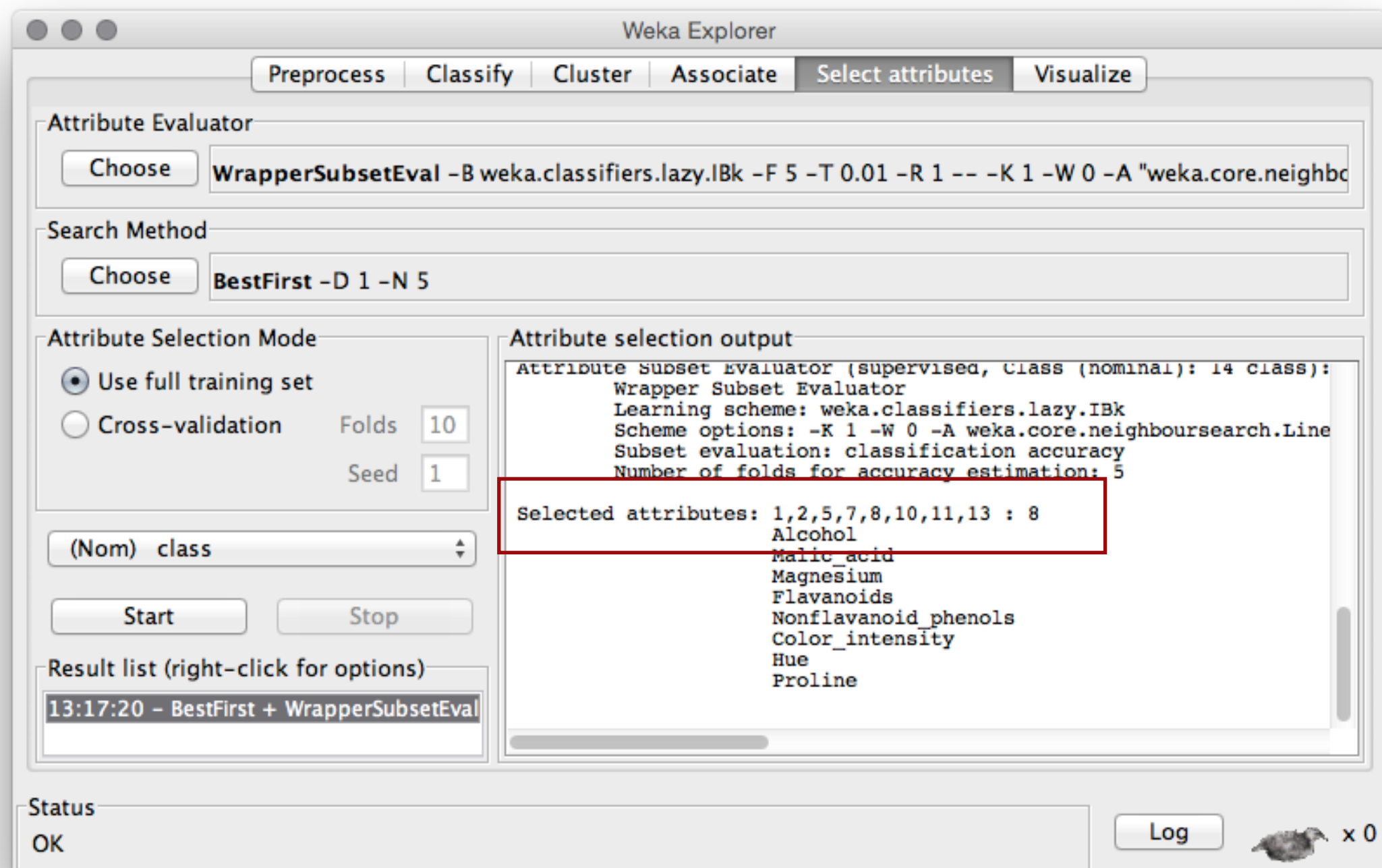
Tutorial Q1(b)

In Weka, apply wrapper-based feature selection to the *Wine* data set using a 1-nearest neighbour classifier and the following search strategies:

- (i) forward sequential search
- (ii) backward elimination

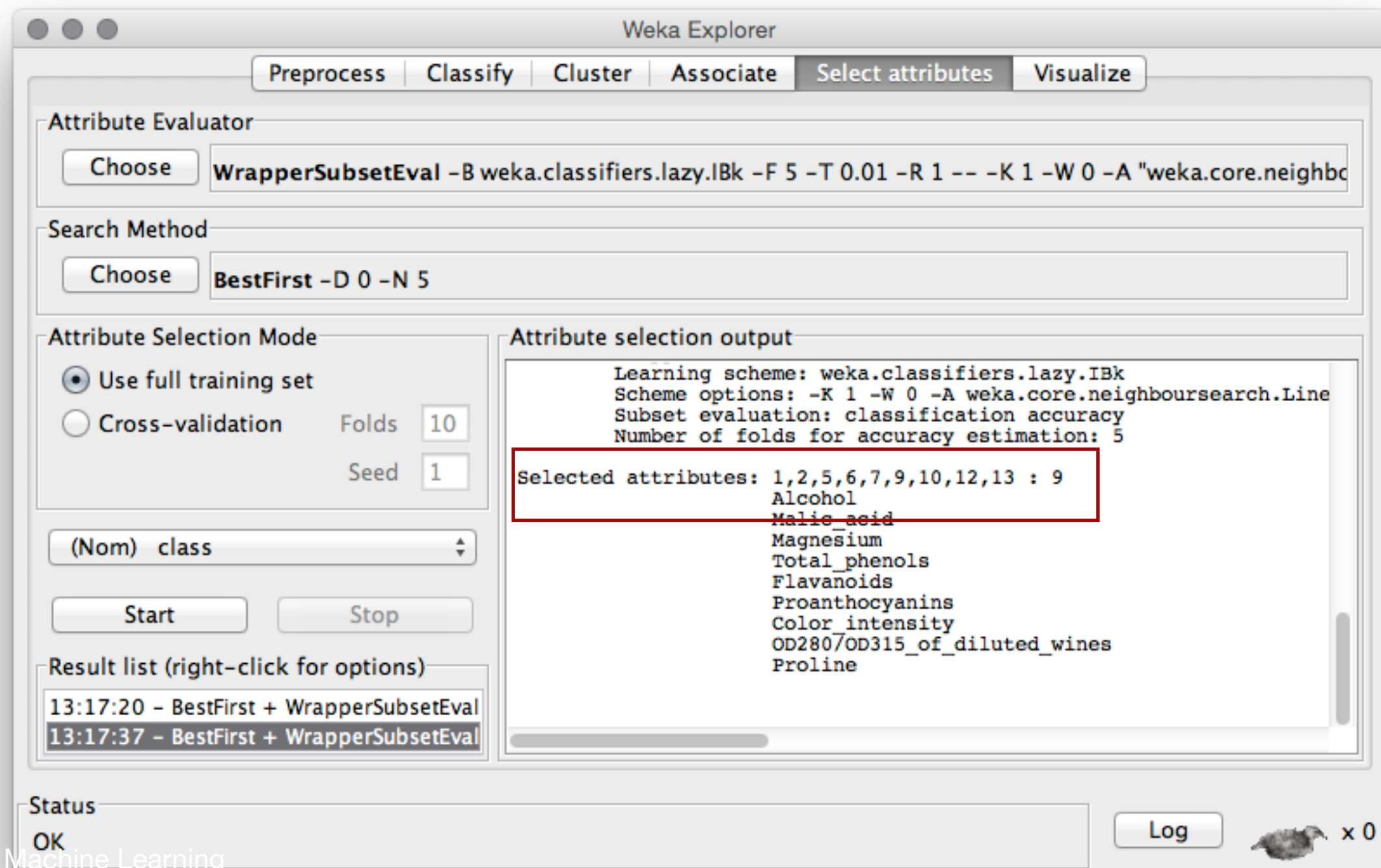
Tutorial Q1(b)

For Forward Sequential Search: In *Select attributes* tab, choose *WrapperSubsetEval* as the evaluator, *BestFirst* as the search method. Change options for search to *Forward* (1).



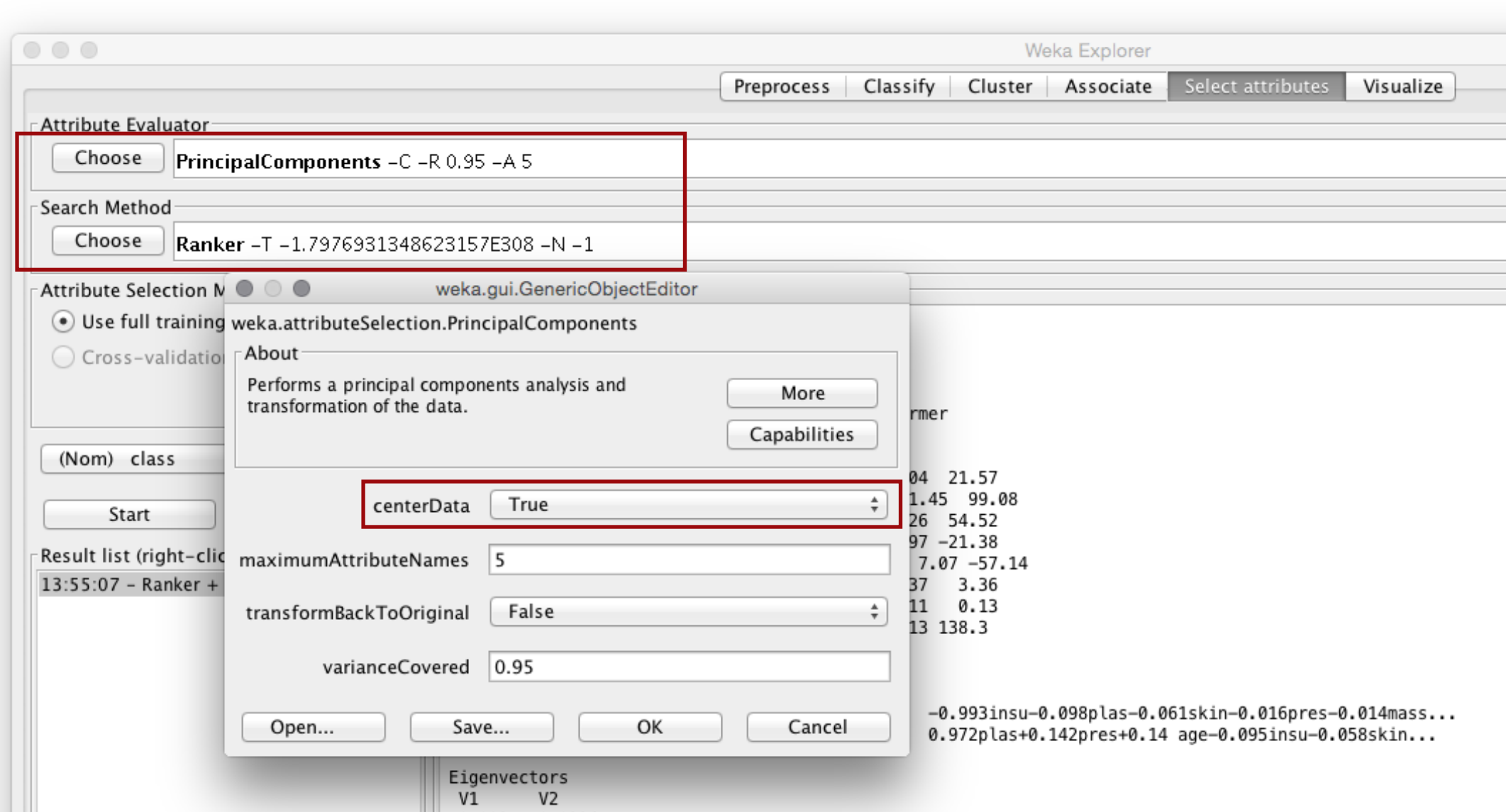
Tutorial Q1(b)

For Backward Elimination: In *Select attributes* tab, choose *WrapperSubsetEval* as the evaluator, *BestFirst* as the search method. Change options for search to *Backward* (0).



Tutorial Q2

In *Select attributes* tab, choose *PrincipalComponents* as the evaluator, and *Ranker* as the search method. Change options for PCA, set *centerData* to True to use the standard covariance matrix approach.



Tutorial Q2

PCA produces a new 2D feature space that accounts for $> 95\%$ of the variance of the data.

The screenshot shows the Weka Explorer interface with the 'Select attributes' tab selected. The 'Attribute Evaluator' is set to 'PrincipalComponents -C -R 0.95 -A 5'. The 'Search Method' is set to 'Ranker -T -1.7976931348623157E308 -N -1'. The 'Attribute Selection Mode' is set to 'Use full training set' with 'Folds' set to 10 and 'Seed' set to 1. The 'Result list' shows '13:55:07 - Ranker + PrincipalComponents'. The 'Attribute selection output' pane displays the following information:

Attribute Evaluator (unsupervised):
Principal Components Attribute Transformer

Covariance matrix

11.35	13.95	9.21	-4.39	-28.56	0.47	-0.04	21.57
13.95	1022.25	94.43	29.24	1220.94	55.73	1.45	99.08
9.21	94.43	374.65	64.03	198.38	43	0.26	54.52
-4.39	29.24	64.03	254.47	802.98	49.37	0.97	-21.38
-28.56	1220.94	198.38	802.98	13281.18	179.78	7.07	-57.14
0.47	55.73	43	49.37	179.78	62.16	0.37	3.36
-0.04	1.45	0.26	0.97	7.07	0.37	0.11	0.13
21.57	99.08	54.52	-21.38	-57.14	3.36	0.13	138.3

Eigenvalue proportion cumulative

13456.57298	0.88855	0.88855	-0.993insu-0.098plas-0.061skin-0.016pres-0.014mass...
932.76013	0.06159	0.95014	0.972plas+0.142pres+0.14 age-0.095insu-0.058skin...

Eigenvectors

V1	V2
0.002	0.0226
-0.0978	0.9722
-0.0161	0.1419
-0.0608	-0.0579
-0.9931	-0.0946
-0.014	0.047
-0.0005	0.0008
0.0036	0.1402

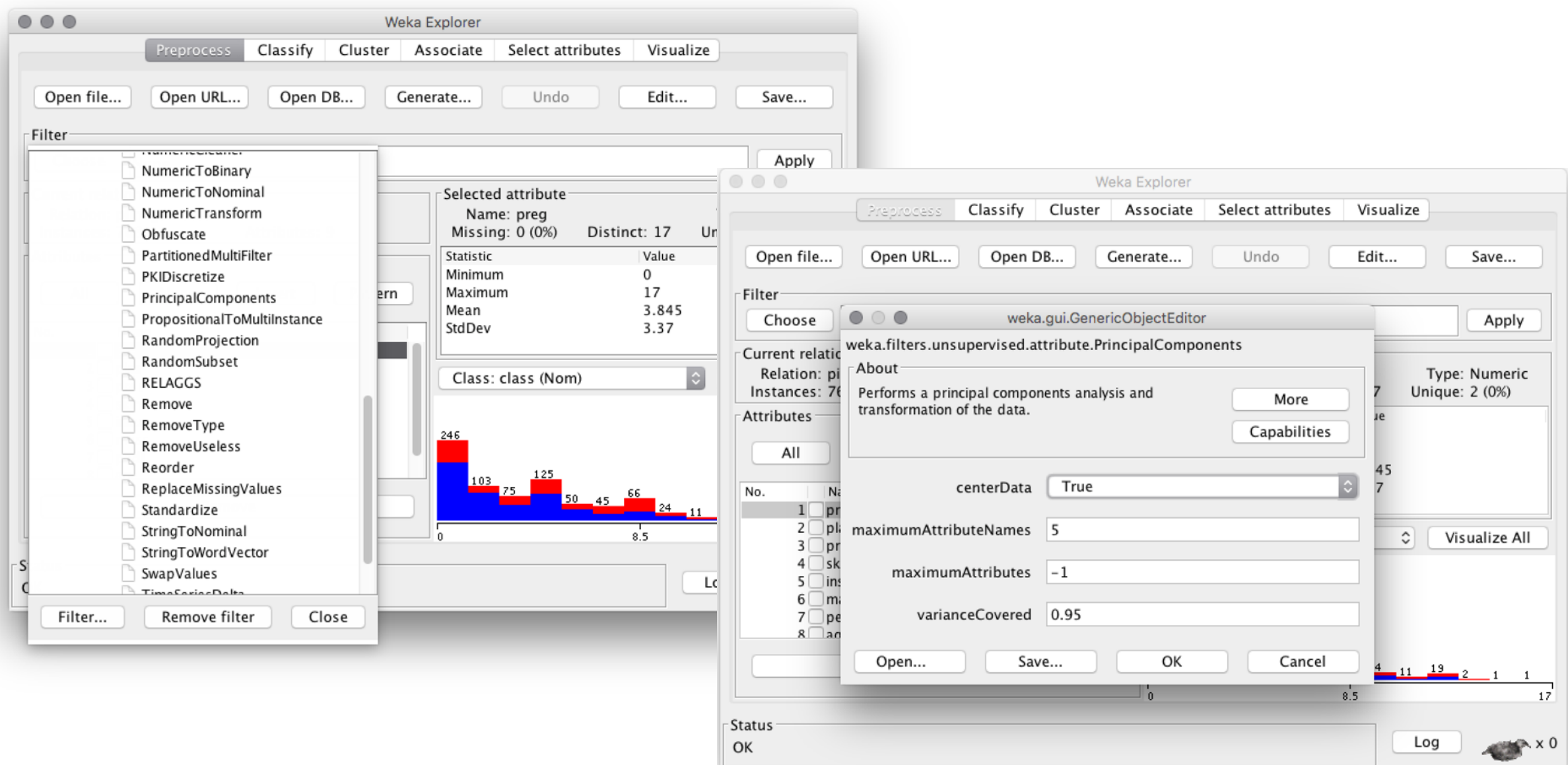
Ranked attributes:

0.1115	1	-0.993insu-0.098plas-0.061skin-0.016pres-0.014mass...
0.0499	2	0.972plas+0.142pres+0.14 age-0.095insu-0.058skin...

Selected attributes: 1,2 : 2

Tutorial Q2

- Can also apply PCA in the Weka *Preprocess* tab.
- Click *Filter*, then *Filters - Unsupervised - Attribute - Principal Components*. In Options, choose *Centre Data*, click *Apply*.



Tutorial Q3

- (a) Explain why the feature subset selection problem with a k-Nearest Neighbour classifier is an exponential search problem.
- (b) Describe in outline a Genetic Algorithm solution to this search problem.
- (c) Describe crossover and mutation techniques for feature subset selection.
- (d) Why is over-fitting a potential risk in wrapper feature subset selection?

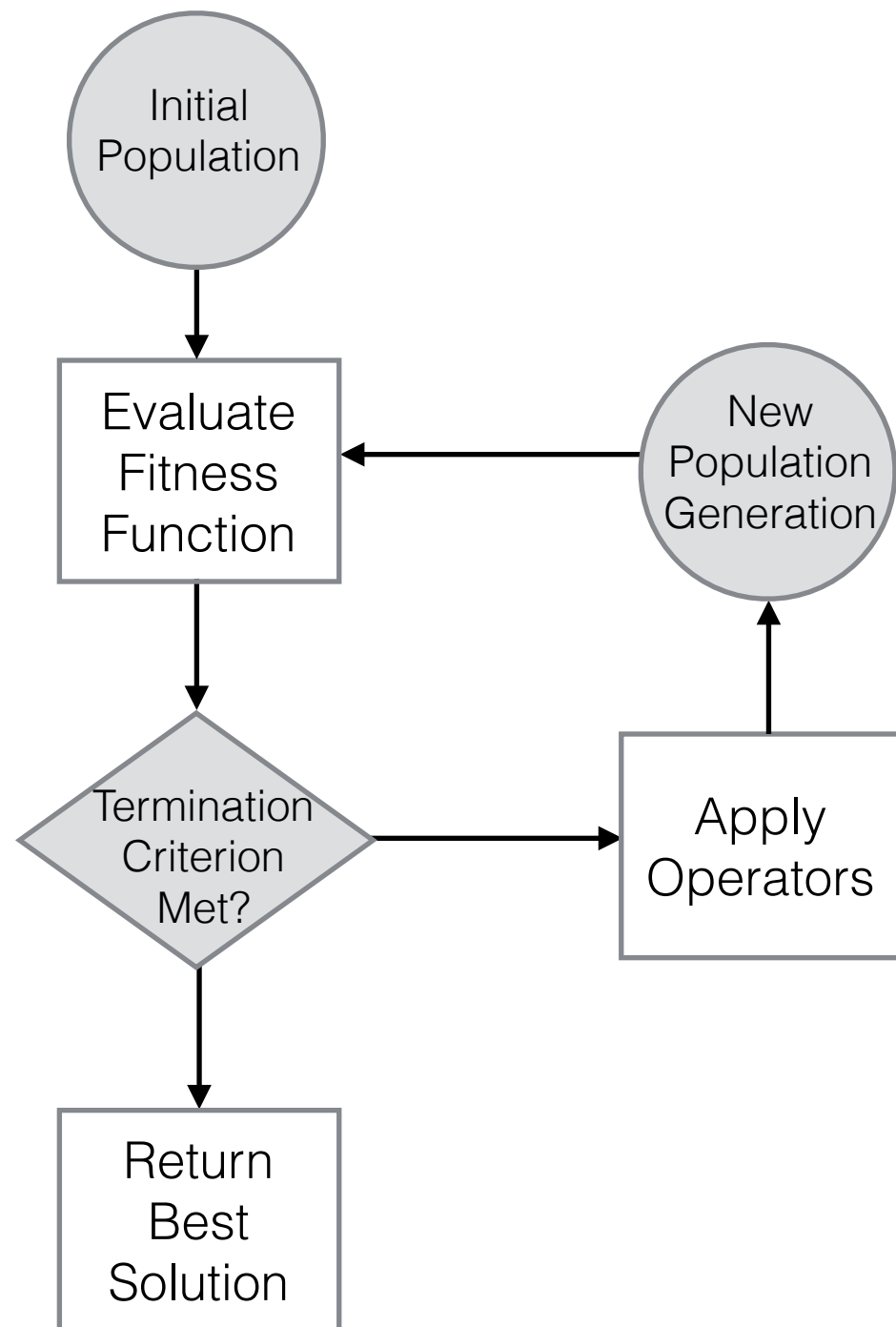
Tutorial Q3(a)

Explain why the feature selection problem with k -Nearest Neighbour Classifiers is an exponential search problem.

- Feature selection is in general NP-hard.
- Brute force evaluation of all feature subsets involves $\binom{d}{k}$ combinations if k is fixed, or 2^d subsets if not fixed.
- If we apply k -NN as a classifier in the context of a wrapper feature selection strategy, we need to evaluate accuracy.
- To measure generalisation accuracy, will typically want to apply k -fold cross validation
e.g. 10-fold cross validation \Rightarrow 10 k -NN runs for every subset

Tutorial Q3(b)

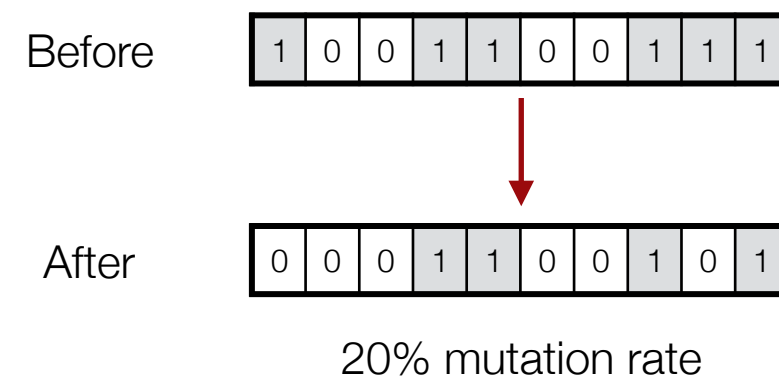
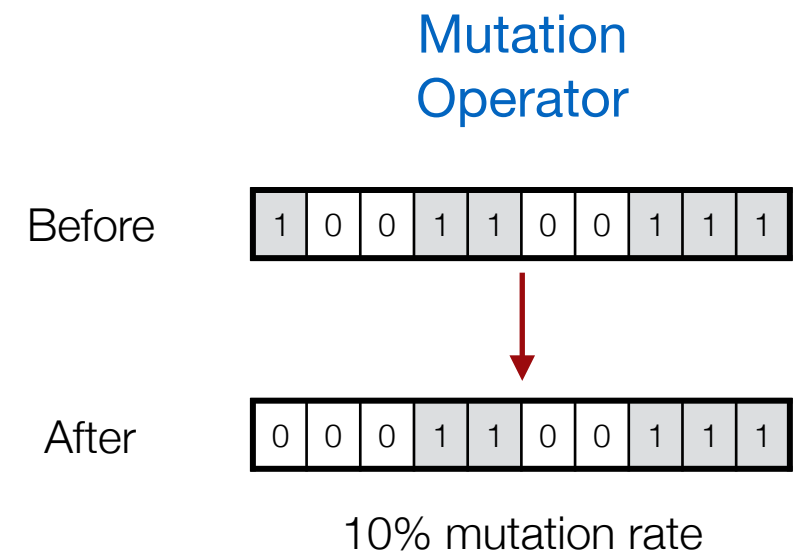
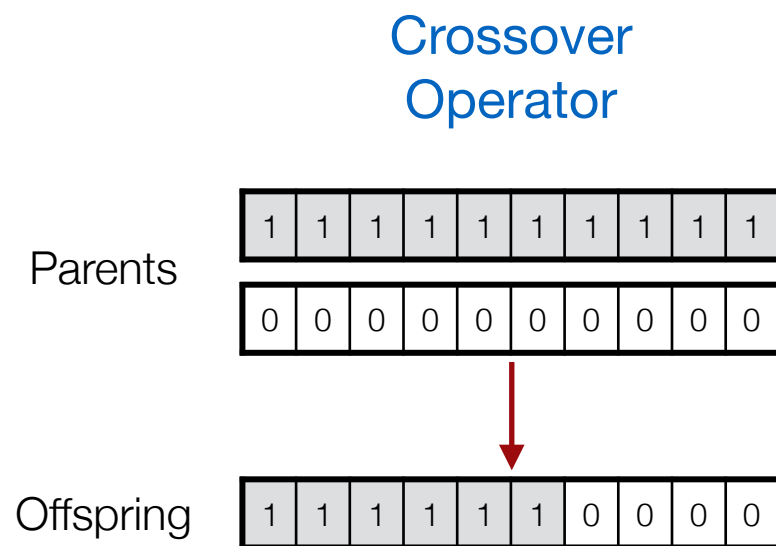
Outline of a Genetic Algorithm solution...



- Involves applying creating “generations” of feature subset solutions, which are iteratively improved.
- *Population*: Different candidate feature subsets.
- *Fitness function*: Classification accuracy using a given classifier on each subset, based on 10 fold cross-validation

Tutorial Q3(c)

Describe crossover and mutation techniques for a Genetic Algorithm solution to the same problem.



Tutorial Q3(d)

- Wrappers can be prone to overfitting when the chosen feature subsets are used to build classifiers on unseen data. We simply found the best feature subset for the training data!
- Using 10-fold cross validation and/or training-test splits to evaluate accuracy of the subsets can reduce overfitting to some extent.
- However, intensive search methods (e.g. Genetic Algorithms) can still identify feature subsets that “on average” favour the characteristics of the training data, and perform poorly on real unseen data that was never used in the feature selection process.

“Overfitting in feature selection appears to be exacerbated by the intensity of the search, since the more feature subsets that are visited the more likely the search is to find a subset that overfits”

- Loughrey & Cunningham, 2004

Tutorial Q3(d)

Journal of Machine Learning Research 3 (2003) 1371-1382

Submitted 5/02; Published 3/03

Overfitting in Making Comparisons Between Variable Selection Methods

Juha Reunanen

*ABB, Web Imaging Systems
P.O. Box 94, 00381 Helsinki, Finland*

JUHA.REUNANEN@FL.ABB.COM

Editors: Isabelle Guyon and André Elisseeff

Abstract

This paper addresses a common methodological flaw in the comparison of variable selection methods. A practical approach to guide the search or the selection process is to compute cross-validation performance estimates of the different variable subsets. Used with computationally intensive algorithms, these estimates may overfit and yield biased predictions. Therefore, they cannot reliably be used to compare two selection methods, as is shown by the empirical results of this paper. Instead, like in other instances of the model selection problem, independent test sets should be used for determining the final performance. The claims made in the literature about the superiority of more exhaustive search algorithms over simpler ones are also revisited, and some of them are refuted.

Keywords: Variable selection; Algorithm comparison; Overfitting; Cross-validation; k -nearest neighbors

Overfitting in Wrapper-Based Feature Subset Selection: The Harder You Try the Worse it Gets

John Loughrey, Pádraig Cunningham

Trinity College Dublin, College Green, Dublin 2, Ireland.
{John.Loughrey, Padraig.Cunningham}@cs.tcd.ie

Abstract. In Wrapper based feature selection, the more states that are visited during the search phase of the algorithm the greater the likelihood of finding a feature subset that has a high internal accuracy while generalizing poorly. When this occurs, we say that the algorithm has overfitted to the training data. We outline a set of experiments to show this and we introduce a modified genetic algorithm to address this overfitting problem by stopping the search before overfitting occurs. This new algorithm called GAWES (Genetic Algorithm With Early Stopping) reduces the level of overfitting and yields feature subsets that have a better generalization accuracy.