# Finding a gene profile classifier for high-grade tumors

David Bozarth
Engineering Science Department, Sonoma State University, Rohnert Park, CA 94928

**Using tumor gene expression signatures to distinguish high-grade tumors (which tend to be more aggressive) from low-grade tumors, performance measures were compared for 45 supervised classifier variations. For 69 genes known to be correlated with high-grade tumors in more than one of 6 different tissue types, published profile data was mined for ability of classifier methods to predict the grade range (high or low) of cancerous tissue samples. Several classifiers yielded accuracy in excess of 70%, with false-negatives (as proportion of all incorrect predictions) well below 0.4.**

## Introduction

DNA microarray analysis focuses data mining and bioinformatics technologies on biomedical problems [1]. Recent studies have used microarray analysis to obtain specific gene transcription profile information on conditions of clinical interest and public health need [2][3][4][5].

Previous work has shown that gene expression signatures from high-grade cancers tend to differ from those of the host tissue, and to exhibit similarity across various host tissue types [3][5].

Data [5][6] published jointly by the University of Michigan Medical School and the Institute of Bioinformatics (Bangalore, India) compile normalized results of DNA profiling studies on cancer, and identify those genes that appear frequently with distinguishing power among host tissue types and cancer types. These data include a set of genes that show significant differential expression between low-grade and high-grade cancers. This data set was downloaded and applied to each of a series of supervised classifiers provided in Weka 3.4, a publicly-available data-mining package [7] associated with the University of Waikato (New Zealand) and our course textbook [8] authors.

## Methods

The data were downloaded as individual sets, each of which gives normalized and comparable [5] gene expression levels for one gene over a number of samples from the same tissue type. The range of tissue types was {bladder, brain, breast, lung, ovarian, prostate}. The web site provides two complete sets of breast cancer data. I chose to use one of these – the one that more closely matched the number of data points in the sets provided for other tissue types.

Of the 69 genes in the data set, all were correlated with more than one of the 6 tissue types, but no gene was correlated with all 6 tissues. For each gene-tissue correlation, a table was downloaded, and all the tables were concatenated into one master table. Irrelevant columns were stripped, leaving just 4 attributes: a gene identifier, the host tissue type, the normalized gene expression level, and the actual classification of the tumor (high-grade or low-grade). The value TRUE was selected to represent high-grade diagnoses and predictions, with FALSE representing

low-grade diagnoses and predictions [*Table 1*]. The number of rows (data points) in the master table was 12,279.

## Table 1. Sample of data table

| tissue | gene | value | class |
|--------|------|-------|-------|
| prostate | CCT6A | -0.46866 | FALSE |
| brain | CKS2 | -0.36081 | TRUE |
| ovarian | CXCL9 | -0.20535 | TRUE |
| bladder | TOP2A | 0.09554 | FALSE |
| brain | DLG7 | 0.19716 | TRUE |
| ovarian | TRIP13 | 0.26510 | TRUE |
| breast | KIF23 | 0.33564 | TRUE |
| prostate | SLC7A5 | 0.39893 | TRUE |
| bladder | DLG7 | 0.45542 | TRUE |
| bladder | NUDT1 | 0.51065 | FALSE |
| ovarian | POLR2K | 0.56453 | FALSE |
| ovarian | POLR2K | 0.62187 | TRUE |
| brain | ILF2 | 0.70695 | FALSE |
| lung | HMGB2 | 0.82642 | TRUE |

This data set was applied in turn to several basic types of supervised classifier in Weka, as well as numerous configurable variations on these basic types. The basic types and their variations, shown with their short-form identifiers, are shown [*Table 2*].

## Table 2. The classifiers and their variations

| identifier | description |
|------------|-------------|
| OneR | Simple one-attribute classifier |
| Naïve Bayes | Naïve Bayes classifier |
| Bayes Net | Bayesian Network classifier |
| AODE | Averaged, one-dependence estimator Bayes classifier |
| NB Tree | Decision tree with Naïve Bayes classifier at leaves |
| ID3 Tree | Simple decision tree |
| J4.8 Tree | Weka implementation of classic C4.5 decision tree |
| AD Tree | Alternating decision tree |
| LMT Tree | Logistic model tree |
| | |
| | *variations* |
| CV | Use stratified 10-fold cross-validation. |
| HO | Train on 2/3 of the data, test on 1/3. |
| HO80 | Train on 4/5 of the data, test on 1/5. |
| mdl | Pre-discretize numeric attributes using Minimum Description Length rule. |
| tan | BayesNet: Use TAN (tree-augmented naïve Bayes) search algorithm. |
| ad | BayesNet: Use all-dimensions tree enhancement for speed. |
| bma | BayesNet: Use BMA estimator. |
| gen | BayesNet: Use genetic algorithm for search. |
| hil | BayesNet: Use hill-climber search algorithm. |
| lap | J4.8 Tree: Use Laplacian smoothing. |
| rep | J4.8 Tree: Use reduced error pruning. |
| sro | J4.8 Tree: Turn off subtree raising. |
| 128 | AD Tree: Use 128 boosting iterations. (Default is 10.) |

Gene expression level is the only numeric attribute among the four in the data table. Some of the classifier types perform an automatic conversion of numeric attributes to nominal; for the Bayes classifiers this conversion was known to be based on the assumption of a normal distribution. Two of the classifiers used – AODE and ID3 – require the input to contain only nominal attributes. The Weka package contains a discretizer based on MDL (minimum description length). This utility was used to pre-discretize expression levels [*Table 3*] for input to AODE and ID3, and for testing other classifiers using the pre-discretized expression level. In some cases comparison was possible between the performance of a classifier using its internal discretizer vs. the MDL preprocessor.

**Table 3. MDL discretization levels**

| category | count |
|---|---|
| '(-inf--1.653375]' | 468 |
| '(-1.653375--1.64537]' | 30 |
| '(-1.64537--1.590325]' | 147 |
| '(-1.590325--1.58407]' | 30 |
| '(-1.58407--1.4549]' | 169 |
| '(-1.4549-0.403225]' | 5772 |
| '(0.403225-0.977525]' | 3434 |
| '(0.977525-2.44719]' | 2166 |
| '(2.44719-inf)' | 63 |

Most classifier variations were validated using both 2/3 holdout and stratified 10-fold cross methods. Some were validated with only one or the other method based on a perceived performance trend among variations on a single classifier type. The purpose throughout was to identify one or more classifier methods that yield, in order of desirability:

(1) high prediction accuracy
(2) low ratio of false-negatives to all wrong predictions
(3) high kappa statistic

False negatives are especially undesirable from a clinical standpoint, representing cancers that were diagnosed as high-grade, but during validation were assigned as FALSE by the classifier. In a clinical setting this could translate into a patient having a potentially severe prognosis, but being told it was unlikely to be severe. The opposite error (not severe but being told it is) would also be undesirable, but less so.

The kappa statistic is derived from the confusion matrix [*Figure 1*], wherein the diagonal values represent predictions that agree with the actual values of validation data, and the off-diagonal values represent "false positives" and "false negatives". A positive value for the kappa statistic indicates agreement between the classifier and a perfect classifier, with the effect of chance subtracted out. [7]. Its use is controversial [8], and for this study it was meant only for possible use as a tie-breaker.

## Results

Sample output from a single classifier test is shown. [*Figure 1*]. The Appendix contains all results in 4 tables: grouped by *classifier* type, ranked by *accuracy*, ranked by *false-negative* rate, and ranked by *kappa* value. The classifier naming keys [*Table 2*] are used in those tables.

**Figure 1. Sample output**

```
Correctly Classified Instances         8650            70.4455 %
Incorrectly Classified Instances       3629            29.5545 %
Kappa statistic                           0.3853
Mean absolute error                       0.3759
Root mean squared error                   0.4554
Relative absolute error                  77.412  %
Root relative squared error              92.4269 %
Total Number of Instances             12279

=== Detailed Accuracy By Class ===

TP Rate    FP Rate    Precision    Recall    F-Measure    Class
 0.772      0.39        0.736      0.772      0.753       TRUE
 0.61       0.228       0.655      0.61       0.632       FALSE

=== Confusion Matrix ===

    a     b    <-- classified as
 5540 1640  |    a = TRUE
 1989 3110  |    b = FALSE
```

The highest 12 rankings for *accuracy* and for *kappa* value are the same, and their respective rankings below 12 do not differ markedly [*Appendix*].

Clear overall winners among the classifiers for this data set are AD Tree, NB Tree, and Logistic Model Tree [*Table 4*]. All 3, and only these 3, appear in the top ten for both *accuracy* and *false-negative*.

**Table 4. Top ten ranking**

| Rank | by Accuracy | | by False Neg | |
|------|-------------|---------|--------------|-------|
| 1 | LMT Tree, HO | 71.8802 | NB Tree, CV | 0.347 |
| 2 | AD Tree, 128, CV | 71.5775 | LMT Tree, HO | 0.353 |
| 3 | NB Tree, CV | 70.8853 | J4.8 Tree, mdl, CV | 0.370 |
| 4 | J4.8 Tree, CV | 70.4455 | J4.8 Tree, lap, mdl, CV | 0.370 |
| 5 | J4.8 Tree, lap, CV | 70.4210 | J4.8 Tree, sro, mdl, CV | 0.370 |
| 6 | AD Tree, 128, HO | 70.4192 | J4.8 Tree, sro+lap, mdl, CV | 0.370 |
| 7 | J4.8 Tree, sro, CV | 70.3800 | NB Tree, mdl, HO | 0.381 |
| 8 | J4.8 Tree, sro+lap, CV | 70.3559 | NB Tree, mdl, CV | 0.382 |
| 9 | J4.8 Tree, HO | 69.1737 | NB-Tree, HO | 0.383 |
| 10 | J4.8 Tree, sro, HO | 69.1737 | AD Tree, 128, CV | 0.386 |

Cross-validated runs are better represented than holdout runs in both top ten rankings. Also, there is a preponderance of holdout-validated runs in the lower third of the *accuracy* ranking [*Appendix*]. In general we expect cross-validation to correct for sample bias, which holdout validation can not do. Overall during testing on these data using these classifiers, the runs evaluated by "CV" tend to have higher performance indicators than those evaluated by "HO". Exceptions are found among the Bayes Net and ordinary Naïve Bayes classifiers, none of which ranked highly.

As of this writing (* see 'New Results Update' below), there is no listing in the tables for LMT using cross-validation. LMT has posed a problem in testing due to its running time for creating the model, which appears to be exponential. The vastly greater share of this running time is consumed in forming the model – not in evaluation. The holdout-validated run of LMT took about 3 hours. With complications from system issues, as of this writing I am trying for the third time to get an 8-hour run of 10-fold cross-validated LMT to go to completion. I believe its performance on the key indicators will outstrip all others (* not entirely), and the model, once working, can be saved for future use.

The AD Tree classifier is of special interest also. A parameter "boosting iterations" was increased in binary progression from 32 to 128. The results kept getting better, and the running time kept increasing. I stopped when its accuracy evaluation surpassed all the others I had gotten at that point, and its running time was about 15 minutes. Subsequently it was surpassed by LMT using holdout validation. It will be interesting to see whether increasing the AD Tree boosting iterations further will enable it to beat the accuracy of LMT (* not found to be the case) .

The AD Tree with its boosting parameter set to 128, ranked tenth in *false-negative* performance. Results were not kept for lower settings of the boosting parameter. It would be interesting to see whether the boosting parameter may affect *false-negative* performance.

LMT and NB Tree appear to be in their own subclass based on *false-negative* performance.

Interesting parameters and long runs aside, NB Tree may offer best "bang for the buck". I didn't have to fiddle with it; it ran to completion in moments; and its performance ranks in the top 3 of all key indicators.

The performance of non-tweaked J.48 is worthy of mention. It ranked in the top ten of *accuracy* by both validation methods, and ranked in the top 17 of *false-negative* by both validation methods.

The two known-valid uses of pre-discretizing (indicated by the 'mdl' tag) were for the AODE Bayes and ID3 classifiers. Both of these ranked low to medium. The effect of the 'mdl' option on performance for other classifiers was generally negative. Exceptions are among tunings of the Bayes Net classifier. Also, the 'mdl' option improved *false-negative* performance among a group of J4.8 classifier tunings, but had the opposite effect on accuracy.

## Summary

Three best-performance classifiers were identified for this data set:

- NB Tree, which uses Naïve Bayes classification at the leaves of a decision tree

- LMT, which uses a logistic model at the leaves of a decision tree

- AD Tree, which is a specialized option tree optimized for two-class problems

Of these, NB Tree was easiest to use. LMT requires a lengthy model-building phase. AD Tree is tunable for increasing accuracy, and merits further investigation, particularly with respect to its model-building time and its false-negative performance.

### *NEW RESULTS UPDATE*

| classifier | %Correct | %Wrong | kappa | False-neg |
|---|---|---|---|---|
| LMT Tree, CV | 72.7339 | 27.2661 | 0.4315 | 0.371 |
| AD Tree, 256, CV | 71.8381 | 28.1619 | 0.4144 | 0.373 |
| AD Tree, 512, CV | 71.6264 | 28.3736 | 0.4119 | 0.364 |

The highest *accuracy*, as expected, was registered using 10-fold CV on LMT, but its *false-negative* performance declined. For AD Tree, both *accuracy* and *false-negative* rate improved with increasing the boosting parameter from 128 to 256 (with corresponding increase in time required to build the model). With a further increase to 512, the model-building phase took about two hours, and the *accuracy* declined slightly, but the *false-negative* performance continued to improve.

# References

[1] Bajcsy P, Han J, Liu L, Yang J. "Survey of Biodata Analysis from a Data Mining Perspective". In *Data Mining in Bioinformatics*, ed. Wang JTL, Zaki M, Toivonen HTT, Shasha D. Springer (2005) pp.9-39.

[2] Walker PR, Smith B, Liu QY, Famili AF, Valdes JJ, Liu A, Lach B. Data mining of gene expression changes in Alzheimer brain. Artificial Intelligence in Medicine (2004) 31:137-154.

[3] Ramaswamy S, et.al. Multiclass cancer diagnosis using tumor gene expression signatures. (2001) PNAS (2001) 98,26:15149-15154.

[4] Van't Veer LJ, et.al. Gene expression profiling predicts clinical outcome of breast cancer. Nature (2002) 415:530-536.

[5] Rhodes DR, et.al. Large-scale meta-analysis of cancer microarray data identifies common transcriptional profiles of neoplastic transformation and progression. PNAS (2004) 101,25:9309-9314.

[6] Oncomine 'meta' study web site http://141.214.6.21:8080/Array/meta/index_html

[7] Witten IH, Frank E. *Data Mining: Practical Machine Learning Tools and Techniques*. Elsevier (2005)

[8] Kappa statistic web site http://ourworld.compuserve.com/homepages/jsuebersax/kappa.htm

# Appendix

## Classifier performance (grouped by classifier type)

| classifier | %correct | %wrong | kappa | false neg |
|---|---|---|---|---|
| OneR, CV | 57.4395 | 42.5605 | 0.1013 | 0.598 |
| OneR, HO | 56.5988 | 43.4012 | 0.0815 | 0.613 |
| Naïve Bayes, CV | 63.5964 | 34.4036 | 0.2220 | 0.564 |
| Naïve Bayes, HO | 64.3353 | 35.6647 | 0.2414 | 0.536 |
| Bayes Net, CV | 63.7837 | 36.2163 | 0.2487 | 0.462 |
| Bayes Net, HO | 64.4551 | 35.5449 | 0.2611 | 0.458 |
| Bayes Net, tan, CV | 66.2187 | 33.7813 | 0.2975 | 0.440 |
| Bayes Net, tan, HO | 65.9401 | 34.0599 | 0.2945 | 0.428 |
| Bayes Net, tan, mdl, CV | 67.0657 | 32.9343 | 0.3141 | 0.436 |
| Bayes Net, tan, mdl, HO | 66.7066 | 33.2934 | 0.3044 | 0.448 |
| Bayes Net, mdl, CV | 64.2723 | 35.7277 | 0.2582 | 0.459 |
| Bayes Net, mdl, HO | 64.7425 | 35.2575 | 0.2672 | 0.454 |
| Bayes Net, ad, mdl, HO | 64.7425 | 35.2575 | 0.2672 | 0.454 |
| Bayes Net, bma, mdl, HO | 66.2754 | 33.7246 | 0.2923 | 0.467 |
| Bayes Net, gen, mdl, HO | 65.9641 | 34.0359 | 0.2909 | 0.447 |
| Bayes Net, hil, mdl, HO | 66.2754 | 33.7246 | 0.2923 | 0.467 |
| AODE Bayes, mdl, CV | 67.2575 | 32.7425 | 0.3242 | 0.400 |
| AODE Bayes, mdl, HO | 67.8394 | 32.1606 | 0.3329 | 0.412 |
| NB Tree, CV | 70.8853 | 29.1147 | 0.4011 | 0.347 |
| NB Tree, HO | 68.6946 | 31.3054 | 0.3538 | 0.383 |
| NB Tree, mdl, CV | 68.2466 | 31.7534 | 0.3463 | 0.382 |
| NB Tree, mdl, HO | 68.2156 | 31.7844 | 0.3454 | 0.381 |
| ID3 Tree, mdl, CV | 66.8214 | 33.1786 | 0.3248 | 0.434 |
| ID3 Tree, mdl, HO | 65.4371 | 34.5629 | 0.3015 | 0.439 |
| J4.8 Tree, CV | 70.4455 | 29.5545 | 0.3853 | 0.390 |
| J4.8 Tree, HO | 69.1737 | 30.8263 | 0.3591 | 0.401 |
| J4.8 Tree, HO-80 | 68.6482 | 31.3518 | 0.3510 | 0.403 |
| J4.8 Tree, lap, CV | 70.4210 | 29.5790 | 0.3847 | 0.391 |
| J4.8 Tree, rep, CV | 68.5642 | 31.4358 | 0.3437 | 0.426 |
| J4.8 Tree, sro, CV | 70.3800 | 29.6197 | 0.3836 | 0.393 |
| J4.8 Tree, sro+lap, CV | 70.3559 | 29.6441 | 0.3830 | 0.394 |
| J4.8 Tree, lap, HO | 69.1257 | 30.8743 | 0.3580 | 0.402 |
| J4.8 Tree, rep, HO | 67.0419 | 32.9581 | 0.2976 | 0.510 |
| J4.8 Tree, sro, HO | 69.1737 | 30.8263 | 0.3591 | 0.401 |
| J4.8 Tree, mdl, HO | 67.4251 | 32.5749 | 0.3246 | 0.413 |
| J4.8 Tree, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| J4.8 Tree, lap, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| J4.8 Tree, sro, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| J4.8 Tree, sro+lap, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| AD Tree, 128, CV | 71.5775 | 28.4225 | 0.4072 | 0.386 |
| AD Tree, 128, HO | 70.4192 | 29.5808 | 0.3838 | 0.393 |
| AD Tree, 128, mdl, CV | 68.7515 | 31.2485 | 0.3492 | 0.415 |
| AD Tree, 128, mdl, HO | 67.8323 | 32.1677 | 0.3303 | 0.422 |
| LMT Tree, HO | 71.8802 | 28.1198 | 0.4181 | 0.353 |

## Classifier performance (ranked by accuracy)

| classifier | %correct | %wrong | kappa | false neg |
|---|---|---|---|---|
| LMT Tree, HO | 71.8802 | 28.1198 | 0.4181 | 0.353 |
| AD Tree, 128, CV | 71.5775 | 28.4225 | 0.4072 | 0.386 |
| NB-Tree, CV | 70.8853 | 29.1147 | 0.4011 | 0.347 |
| J4.8 Tree, CV | 70.4455 | 29.5545 | 0.3853 | 0.390 |
| J4.8 Tree, lap, CV | 70.4210 | 29.5790 | 0.3847 | 0.391 |
| AD Tree, 128, HO | 70.4192 | 29.5808 | 0.3838 | 0.393 |
| J4.8 Tree, sro, CV | 70.3800 | 29.6197 | 0.3836 | 0.393 |
| J4.8 Tree, sro+lap, CV | 70.3559 | 29.6441 | 0.3830 | 0.394 |
| J4.8 Tree, HO | 69.1737 | 30.8263 | 0.3591 | 0.401 |
| J4.8 Tree, sro, HO | 69.1737 | 30.8263 | 0.3591 | 0.401 |
| J4.8 Tree, lap, HO | 69.1257 | 30.8743 | 0.3580 | 0.402 |
| AD Tree, 128, mdl, CV | 68.7515 | 31.2485 | 0.3492 | 0.415 |
| NB-Tree, HO | 68.6946 | 31.3054 | 0.3538 | 0.383 |
| J4.8 Tree, HO-80 | 68.6482 | 31.3518 | 0.3510 | 0.403 |
| J4.8 Tree, rep, CV | 68.5642 | 31.4358 | 0.3437 | 0.426 |
| NB Tree, mdl, CV | 68.2466 | 31.7534 | 0.3463 | 0.382 |
| NB Tree, mdl, HO | 68.2156 | 31.7844 | 0.3454 | 0.381 |
| J4.8 Tree, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| J4.8 Tree, lap, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| J4.8 Tree, sro, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| J4.8 Tree, sro+lap, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| AODE Bayes, mdl, HO | 67.8394 | 32.1606 | 0.3329 | 0.412 |
| AD Tree, 128, mdl, HO | 67.8323 | 32.1677 | 0.3303 | 0.422 |
| J4.8 Tree, mdl, HO | 67.4251 | 32.5749 | 0.3246 | 0.413 |
| AODE Bayes, mdl, CV | 67.2575 | 32.7425 | 0.3242 | 0.400 |
| Bayes Net, tan, mdl, CV | 67.0657 | 32.9343 | 0.3141 | 0.436 |
| J4.8 Tree, rep, HO | 67.0419 | 32.9581 | 0.2976 | 0.510 |
| ID3 Tree, mdl, CV | 66.8214 | 33.1786 | 0.3248 | 0.434 |
| Bayes Net, tan, mdl, HO | 66.7066 | 33.2934 | 0.3044 | 0.448 |
| Bayes Net, bma, mdl, HO | 66.2754 | 33.7246 | 0.2923 | 0.467 |
| Bayes Net, hil, mdl, HO | 66.2754 | 33.7246 | 0.2923 | 0.467 |
| Bayes Net, tan, CV | 66.2187 | 33.7813 | 0.2975 | 0.440 |
| Bayes Net, gen, mdl, HO | 65.9641 | 34.0359 | 0.2909 | 0.447 |
| Bayes Net, tan, HO | 65.9401 | 34.0599 | 0.2945 | 0.428 |
| ID3 Tree, mdl, HO | 65.4371 | 34.5629 | 0.3015 | 0.439 |
| Bayes Net, mdl, HO | 64.7425 | 35.2575 | 0.2672 | 0.454 |
| Bayes Net, ad, mdl, HO | 64.7425 | 35.2575 | 0.2672 | 0.454 |
| Bayes Net, HO | 64.4551 | 35.5449 | 0.2611 | 0.458 |
| Naïve Bayes, HO | 64.3353 | 35.6647 | 0.2414 | 0.536 |
| Bayes Net, mdl, CV | 64.2723 | 35.7277 | 0.2582 | 0.459 |
| Bayes Net, CV | 63.7837 | 36.2163 | 0.2487 | 0.462 |
| Naïve Bayes, CV | 63.5964 | 34.4036 | 0.2220 | 0.564 |
| OneR, CV | 57.4395 | 42.5605 | 0.1013 | 0.598 |
| OneR, HO | 56.5988 | 43.4012 | 0.0815 | 0.613 |

## Classifier performance (ranked by false-negative ratio)

| classifier | %correct | %wrong | kappa | false neg |
|---|---|---|---|---|
| NB-Tree, CV | 70.8853 | 29.1147 | 0.4011 | 0.347 |
| LMT Tree, HO | 71.8802 | 28.1198 | 0.4181 | 0.353 |
| J4.8 Tree, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| J4.8 Tree, lap, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| J4.8 Tree, sro, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| J4.8 Tree, sro+lap, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| NB Tree, mdl, HO | 68.2156 | 31.7844 | 0.3454 | 0.381 |
| NB Tree, mdl, CV | 68.2466 | 31.7534 | 0.3463 | 0.382 |
| NB-Tree, HO | 68.6946 | 31.3054 | 0.3538 | 0.383 |
| AD Tree, 128, CV | 71.5775 | 28.4225 | 0.4072 | 0.386 |
| J4.8 Tree, CV | 70.4455 | 29.5545 | 0.3853 | 0.390 |
| J4.8 Tree, lap, CV | 70.4210 | 29.5790 | 0.3847 | 0.391 |
| J4.8 Tree, sro, CV | 70.3800 | 29.6197 | 0.3836 | 0.393 |
| AD Tree, 128, HO | 70.4192 | 29.5808 | 0.3838 | 0.393 |
| J4.8 Tree, sro+lap, CV | 70.3559 | 29.6441 | 0.3830 | 0.394 |
| AODE Bayes, mdl, CV | 67.2575 | 32.7425 | 0.3242 | 0.400 |
| J4.8 Tree, HO | 69.1737 | 30.8263 | 0.3591 | 0.401 |
| J4.8 Tree, sro, HO | 69.1737 | 30.8263 | 0.3591 | 0.401 |
| J4.8 Tree, lap, HO | 69.1257 | 30.8743 | 0.3580 | 0.402 |
| J4.8 Tree, HO-80 | 68.6482 | 31.3518 | 0.3510 | 0.403 |
| AODE Bayes, mdl, HO | 67.8394 | 32.1606 | 0.3329 | 0.412 |
| J4.8 Tree, mdl, HO | 67.4251 | 32.5749 | 0.3246 | 0.413 |
| AD Tree, 128, mdl, CV | 68.7515 | 31.2485 | 0.3492 | 0.415 |
| AD Tree, 128, mdl, HO | 67.8323 | 32.1677 | 0.3303 | 0.422 |
| J4.8 Tree, rep, CV | 68.5642 | 31.4358 | 0.3437 | 0.426 |
| Bayes Net, tan, HO | 65.9401 | 34.0599 | 0.2945 | 0.428 |
| ID3 Tree, mdl, CV | 66.8214 | 33.1786 | 0.3248 | 0.434 |
| Bayes Net, tan, mdl, CV | 67.0657 | 32.9343 | 0.3141 | 0.436 |
| ID3 Tree, mdl, HO | 65.4371 | 34.5629 | 0.3015 | 0.439 |
| Bayes Net, tan, CV | 66.2187 | 33.7813 | 0.2975 | 0.440 |
| Bayes Net, gen, mdl, HO | 65.9641 | 34.0359 | 0.2909 | 0.447 |
| Bayes Net, tan, mdl, HO | 66.7066 | 33.2934 | 0.3044 | 0.448 |
| Bayes Net, mdl, HO | 64.7425 | 35.2575 | 0.2672 | 0.454 |
| Bayes Net, ad, mdl, HO | 64.7425 | 35.2575 | 0.2672 | 0.454 |
| Bayes Net, HO | 64.4551 | 35.5449 | 0.2611 | 0.458 |
| Bayes Net, mdl, CV | 64.2723 | 35.7277 | 0.2582 | 0.459 |
| Bayes Net, CV | 63.7837 | 36.2163 | 0.2487 | 0.462 |
| Bayes Net, bma, mdl, HO | 66.2754 | 33.7246 | 0.2923 | 0.467 |
| Bayes Net, hil, mdl, HO | 66.2754 | 33.7246 | 0.2923 | 0.467 |
| J4.8 Tree, rep, HO | 67.0419 | 32.9581 | 0.2976 | 0.510 |
| Naïve Bayes, HO | 64.3353 | 35.6647 | 0.2414 | 0.536 |
| Naïve Bayes, CV | 63.5964 | 34.4036 | 0.2220 | 0.564 |
| OneR, CV | 57.4395 | 42.5605 | 0.1013 | 0.598 |
| OneR, HO | 56.5988 | 43.4012 | 0.0815 | 0.613 |

## Classifier performance (ranked by kappa value)

| classifier | %correct | %wrong | kappa | false neg |
|---|---|---|---|---|
| LMT Tree, HO | 71.8802 | 28.1198 | 0.4181 | 0.353 |
| AD Tree, 128, CV | 71.5775 | 28.4225 | 0.4072 | 0.386 |
| NB-Tree, CV | 70.8853 | 29.1147 | 0.4011 | 0.347 |
| J4.8 Tree, CV | 70.4455 | 29.5545 | 0.3853 | 0.390 |
| J4.8 Tree, lap, CV | 70.4210 | 29.5790 | 0.3847 | 0.391 |
| AD Tree, 128, HO | 70.4192 | 29.5808 | 0.3838 | 0.393 |
| J4.8 Tree, sro, CV | 70.3800 | 29.6197 | 0.3836 | 0.393 |
| J4.8 Tree, sro+lap, CV | 70.3559 | 29.6441 | 0.3830 | 0.394 |
| J4.8 Tree, HO | 69.1737 | 30.8263 | 0.3591 | 0.401 |
| J4.8 Tree, sro, HO | 69.1737 | 30.8263 | 0.3591 | 0.401 |
| J4.8 Tree, lap, HO | 69.1257 | 30.8743 | 0.3580 | 0.402 |
| NB-Tree, HO | 68.6946 | 31.3054 | 0.3538 | 0.383 |
| J4.8 Tree, HO-80 | 68.6482 | 31.3518 | 0.3510 | 0.403 |
| AD Tree, 128, mdl, CV | 68.7515 | 31.2485 | 0.3492 | 0.415 |
| NB Tree, mdl, CV | 68.2466 | 31.7534 | 0.3463 | 0.382 |
| NB Tree, mdl, HO | 68.2156 | 31.7844 | 0.3454 | 0.381 |
| J4.8 Tree, rep, CV | 68.5642 | 31.4358 | 0.3437 | 0.426 |
| J4.8 Tree, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| J4.8 Tree, lap, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| J4.8 Tree, sro, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| J4.8 Tree, sro+lap, mdl, CV | 67.9208 | 32.0792 | 0.3426 | 0.370 |
| AODE Bayes, mdl, HO | 67.8394 | 32.1606 | 0.3329 | 0.412 |
| AD Tree, 128, mdl, HO | 67.8323 | 32.1677 | 0.3303 | 0.422 |
| ID3 Tree, mdl, CV | 66.8214 | 33.1786 | 0.3248 | 0.434 |
| J4.8 Tree, mdl, HO | 67.4251 | 32.5749 | 0.3246 | 0.413 |
| AODE Bayes, mdl, CV | 67.2575 | 32.7425 | 0.3242 | 0.400 |
| Bayes Net, tan, mdl, CV | 67.0657 | 32.9343 | 0.3141 | 0.436 |
| Bayes Net, tan, mdl, HO | 66.7066 | 33.2934 | 0.3044 | 0.448 |
| ID3 Tree, mdl, HO | 65.4371 | 34.5629 | 0.3015 | 0.439 |
| J4.8 Tree, rep, HO | 67.0419 | 32.9581 | 0.2976 | 0.510 |
| Bayes Net, tan, CV | 66.2187 | 33.7813 | 0.2975 | 0.440 |
| Bayes Net, tan, HO | 65.9401 | 34.0599 | 0.2945 | 0.428 |
| Bayes Net, bma, mdl, HO | 66.2754 | 33.7246 | 0.2923 | 0.467 |
| Bayes Net, hil, mdl, HO | 66.2754 | 33.7246 | 0.2923 | 0.467 |
| Bayes Net, gen, mdl, HO | 65.9641 | 34.0359 | 0.2909 | 0.447 |
| Bayes Net, mdl, HO | 64.7425 | 35.2575 | 0.2672 | 0.454 |
| Bayes Net, ad, mdl, HO | 64.7425 | 35.2575 | 0.2672 | 0.454 |
| Bayes Net, HO | 64.4551 | 35.5449 | 0.2611 | 0.458 |
| Bayes Net, mdl, CV | 64.2723 | 35.7277 | 0.2582 | 0.459 |
| Bayes Net, CV | 63.7837 | 36.2163 | 0.2487 | 0.462 |
| Naïve Bayes, HO | 64.3353 | 35.6647 | 0.2414 | 0.536 |
| Naïve Bayes, CV | 63.5964 | 34.4036 | 0.2220 | 0.564 |
| OneR, CV | 57.4395 | 42.5605 | 0.1013 | 0.598 |
| OneR, HO | 56.5988 | 43.4012 | 0.0815 | 0.613 |

.