# A Comparison of Classification Methods on the Avila Bible Dataset

Robert Edwards
2416963E

## 1. INTRODUCTION

Palaeography is a field that studies analyses ancient and medieval texts to infer metadata such as when and where the text was written, determine the number of scribes, and the workload of each. Traditionally, Palaeographers manually analysed the texts to identify the unique writing patterns of a scribe or group of scribes, see [1] and references included. Modern techniques use computers to extract quantitative information from digitalized versions of ancient texts. The "Avila" Bible is a Latin copy of the Bible dating back to the 12th century that has been digitized via imaging. A dataset has been created from 800 images of the "Avila" Bible with 10 measurable features that are used to distinguish one of 12 copyists to each unique pattern, see [1] for further details. This project analyses a subset of 2000 observations created from "Avila" dataset to classify a pattern, or observation, to a copyist using classification methods covered in the course.

## 2. DATA

Before fitting models an attempt is made to understand the data through summary statistics and visualisation. The frequency of each class is shown in *Table 1*. There are three classes (A, B, C) that have very low frequencies and may need to be excluded from some models depending on the random sampling of the training, validation, and test subsets. Only class A and class F have a large number of observations. The low observations in the other classes may affect the variability of some methods.

### Class Frequency

| Class | A | B | C | D | E | F | G | H | I | W | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Frequency | 816 | 1 | 19 | 79 | 194 | 377 | 71 | 97 | 163 | 11 | 106 | 66 |

*Table 1: Frequencies for each class. Class B, class C, and class W are of note because they have low counts. These classes may need to be exempted from some models.*

The features extracted from the "Avila" Bible are visualised in *Figure 1*. At first glance, all the features seem distributed around 0 and have similar interquartile ranges, with the exception of feature F5. The first set (F1, F2, F3) has many outliers and are heavily right-skewed. The features in the first set do not appear normal and have high variability, a more flexible method may have higher accuracy classifying these data. The second set (F4, F5, F6, F7) has less variability than the first set but there are a number of outliers. Features F4 and F6 seem normally distributed whereas features F5 and F7 are skewed left. The third set of features (F8, F9, F10) all have a number of outliers but are symmetric and appear normally distributed. All the features are distributed around 0 and seem to be standardized. According to the source of

this dataset [2], the data have been normalized using Z-normalization. No further normalizing or standardizing of the data is done for the purposes of this report.
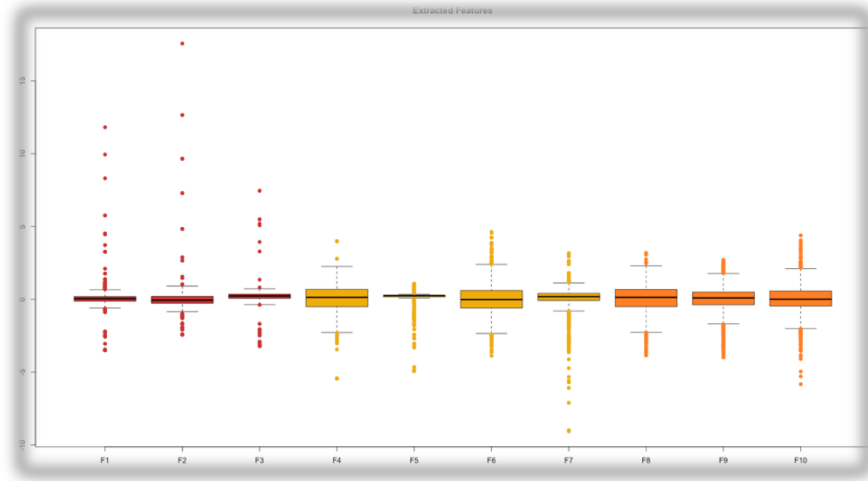


*Figure 1: Boxplot of extracted features from the Avila Bible.*

Another important way to consider the data is via a pairs plot showing the data as described by a feature against each other feature. Pairs plots are good for finding interactions between features; which are highly correlated and which are uncorrelated. The pairs plot in *Figure 2* shows a strong positive correlation between features F6 and F10 and also between features F8 and F9. The interactions between the other features look to be either uncorrelated or have a complex interactions. The diagonal of the pairs plot shows histograms for each feature. Features F4, F6, F8, F9, and F10 seem to have a normal or close to normal distribution. Methods that assume normality may have a higher accuracy with these features.
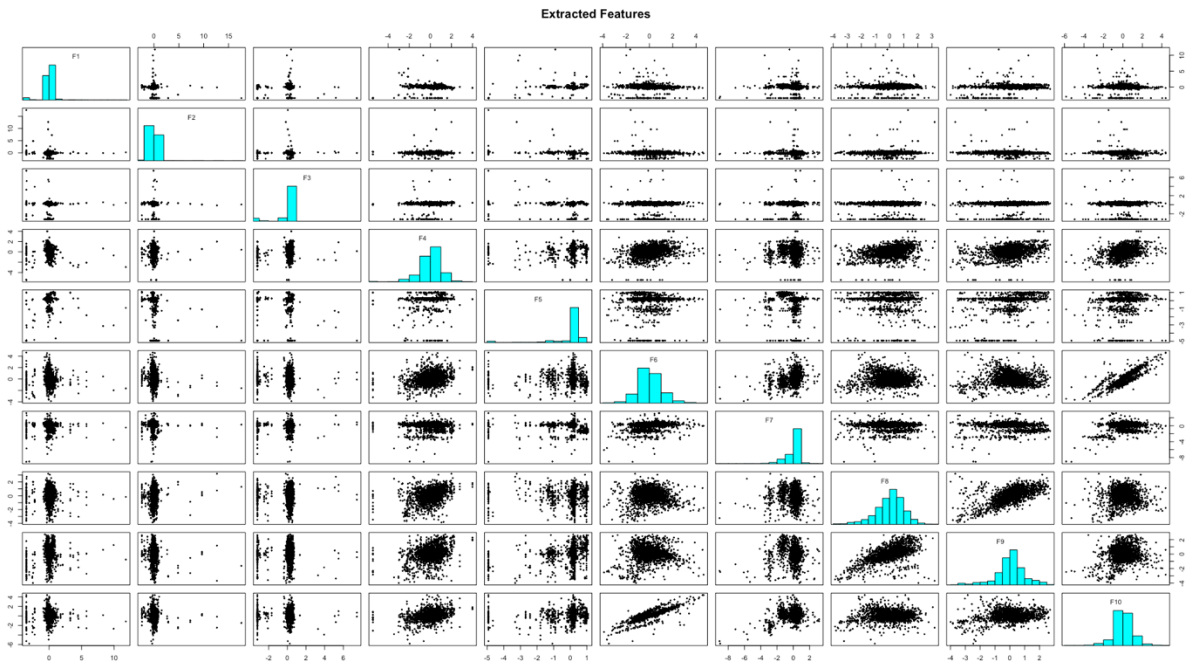


*Figure 2: Pairs plot of extracted features from "Avila" Bible*

2

## 3.  PROCEDURE

For each of the seven methods explored the data are randomly split into three subsets: training (50%), validation (25%), and testing (25%).  A classification model is created using the training data and a prediction of the validation data is created using the models.  A cross-classification table is created from the predicted classifications and the true classifications of the validation data and the error rate is calculated for each model.  Before model selection, the training and validation data are randomly sampled 100 times and the overall error rates for each model are averaged.  The best model is then chosen by the smallest overall average error rate of the seven methods.  The trained models are then tested on the test set and the expected future error rates are calculated.  Error rates for each class are also calculated in an attempt to identify the model that most accurately predicts each class.

## 4.  METHODS

Seven methods were explored for classifying the Avila Bible data: Linear Regression, K Nearest Neighbours, Canonical Variate Analysis, Linear Discriminant Analysis, Quadratic Discriminant Analysis, Classification Trees, and Random Forests.  A brief description of each method is given including the strengths, weaknesses, limitations, and implementation.

*Linear Regression*

Linear regression is a widely used approach for classification of a quantitative response because it is simple and results are interpretable.  Linear Regression is a parametric model that assumes the regression function parameters are linear, statistical independence of the error, homoscedasticity, and normality of the error distribution.  Many of the assumptions for linear regression are not valid for real world datasets but linear regression serves as a good baseline classification method.  Many of the other methods explored in this report and in the course text [3,4] are extensions of linear regression.  To use linear regression to classify the "Avila" Bible features to one of the 12 copyists, categorical response variables, each class is assigned to an integer; the copyists are assigned 1 for copyist A through 12 for copyist Y.  Prediction values are then rounded to the closest whole integer in the range of classifiers.

*KNN*

A more flexible technique for classifying data is K-Nearest Neighbors (KNN).  The method estimates the class of an observation by looking at the surrounding observations.  More formally, KNN is a nonparametric approach that identifies the K observations in the training set that are closest to the observation being classified.  An estimate for the conditional probability of each class is calculated and the observation is assigned to the class with the highest probability, see the class text [3,4].  KNN is a nonparametric method that works well for even small numbers of observations.  Since the distance between observations is calculated, the data should be scaled.  There are a number of tuning parameters for KNN.   An attempt is made to optimize KNN by training models with a range of K values, various types of distance, and various kernel types.  The kknn library is used for modelling because it has options for tuning the three parameters being considered.  With some classes having very low numbers of observations, KNN with high

K values are not expected to perform well. Additionally, knowing that KNN performs worse as the number of dimensions increase, see class text [3,4]. KNN models may not perform well if all 10 features are considered for modelling.

*Discriminant Analysis*

Discriminant analysis involves slightly more abstract technique that classifies observations by first transforming the data in a way that maximizes the ratio of group means to group variance to find the optimal separation between classes. The three versions explored in this paper are Canonical Variate Analysis (CVA), Linear Discriminant Analysis (LDA), and Quadratic Discriminant Analysis (QDA). CVA, also known as Fisher's Discriminant Analysis, does not assume a distribution for the data but it does assume equal prior probabilities. CVA works by finding ratios of variance between classes to variance within classes and assigning the observation to the class with the highest ratio. If a normal distribution is assumed then LDA and QDA can be used and the prior probabilities can be estimated as the proportion of observations belonging to each class, whereas with CVA the prior probabilities were assumed equal, see [1,2] for further details. Using Bayes Discriminant Rule the posterior probabilities for each class and the observation is classified to the class with the highest posterior probability. CVA and LDA both assume equal covariance matrices and Bayes Discriminant Rule simplifies to a linear function. Since QDA assumes unequal covariance matrix the Bayes Discriminant Rule simplifies to a quadratic function. QDA must estimate separate covariance matrices for each group and must have more observations than predictors for it to be used, see [5] for more details. LDA and QDA are both parametric models as they assume a multivariate normal distribution for the data in each class. There is a tradeoff between bias and variance in these models as more is assumed about the data. Discriminant analysis works well as separation between classes increases. LDA and QDA also work well with a small number of observations.

*Trees*

Classification Trees are a nonparametric technique that work well when the data is nonlinear and not normal. Classification Trees are grown via recursive binary splitting using a measure of the entropy as the split criterion. An observation is assigned to the class of the most abundant class of training observations in the same region. If the data has been overfitted in training then the tree can be pruned to reduce variability at the cost of increased bias. Trees are simple to explain and interpret but are not robust, in that a variability in data can grow a very different tree.

*Random Forests*

One method of reducing variability in a tree is to use a technique called Bagging. Many different training sets are sampled, unpruned trees are grown from the trainings sets, and the resulting predictions are averaged. Bagging tries to reduce variability by averaging many different trees but if there are particularly strong predictors then the trees can be highly correlated. The strong predictors tend to be chosen as the criterion split and the trees are grown similarly. The Random Forests method attempts to decorrelate trees by limiting the number of predictors considered at each split allowing less strong predictors to be considered. If there are a

number of correlated predictors, then only a small number of predictors should be considered at each split. The default setting of the square root of the number of parameters is used to build the Random Forests. There are many more methods that can increase the accuracy of classification trees, see [3,4].

## 5. RESULTS

Before comparing against the other methods some parameters of K Nearest Neighbours are tuned using the training data. Error rates were determined for a range of K values, types of distance, and kernel types. Odd K values from 1 to 19 are used to give a range of potential estimates that tradeoff between increased bias and decreased variability as the value of K increases. Four distances were evaluated using the Minkowski distance parameter in the kknn package: 1, 2, and 10 corresponding to Manhattan distance, Euclidean distance, and Chebyshev distance. Additionally, a distance parameter of 0.1 is considered. Lastly, four kernel types are considered: rectangular, triangular, gaussian, and optimal. The set of parameters with the smallest error rate were a K value of 11, a distance measure of 10, and a rectangular kernel with an error rate of 0.516. *Figure 3* shows the outputs for tuning the KNN model.
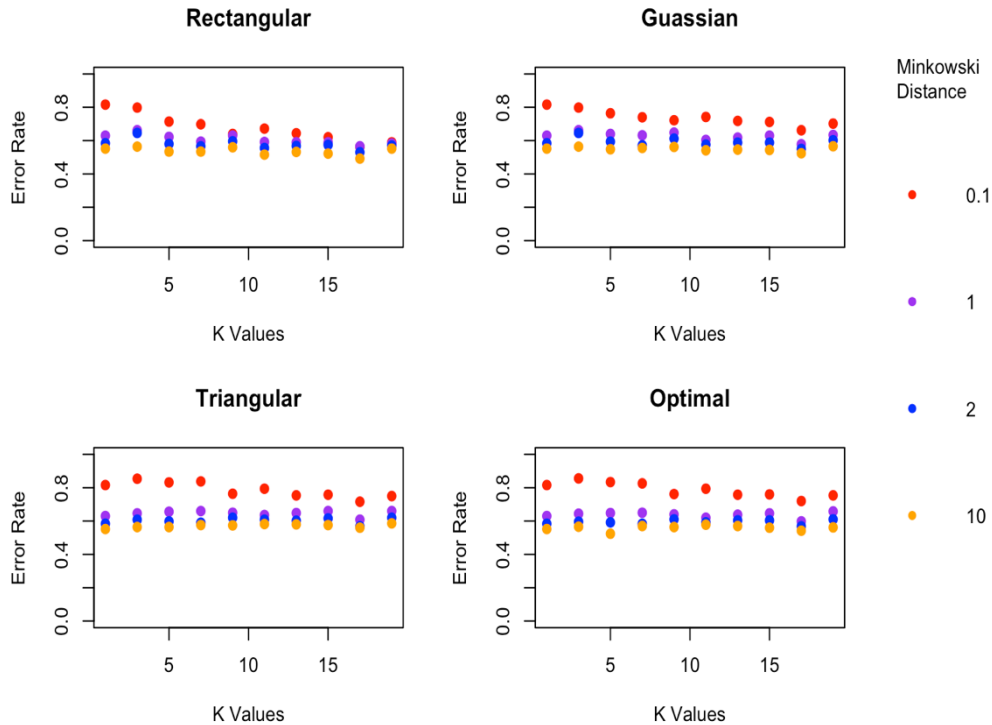


*Figure 3: Optimization of K-Nearest Neighbours by finding the smallest error rate for various parameters. K values are odd integers ranging from 1 to 19. Minkowski distances are 0.1, 1, 2, and 10. The four kernel types are rectangular, triangular, gaussian, and optimal.*

After tuning the KNN method, all the methods were trained and cross validated for 100 different random samplings of the "Avila" Bible dataset. The error rates for the 100 cross validations are shown in *Figure 4*. During this process classes were excluded from the training and validation sets if there were no observations or if the observations were less than 10 when training QDA.
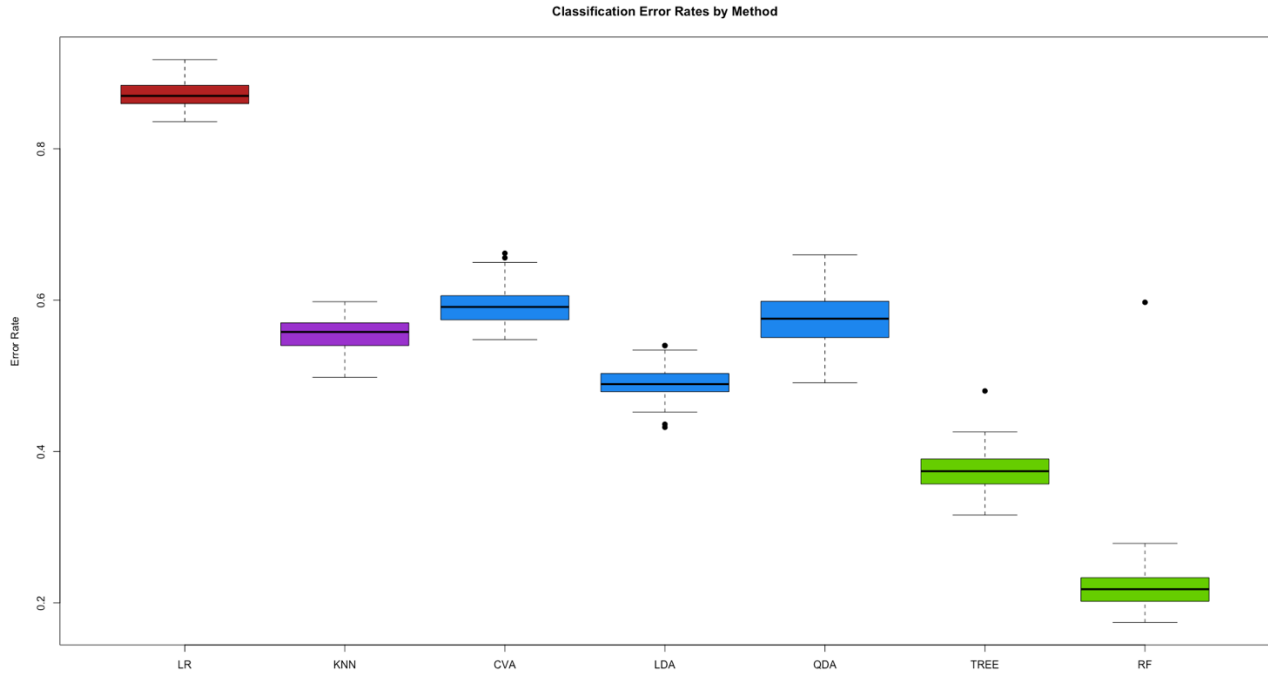


*Figure 4: Error rates from 100 random samplings of training and validation data.*

Methods such as Classification Trees and Random Forests trade off high variability for reduced bias. *Figure* 4 shows outliers in CVA, LDA, Classification Trees, and Random Forests. By averaging the error rates calculated from cross validating training data from 100 random samples, the variability of each method is reduced. Random Forests has the lowest overall average classification rate and is selected as the optimal model to use on the test data. All trained methods are tested on the test set to find the expected prediction error rate on future datasets. The overall classification rates are shown in *Table* 2 after the models were cross validated on the test set. Random Forests has the lowest error rate at 0.238 and should be used as the expected overall error rate for classification of future predictions.

| LR | KNN | CVA | LDA | QDA | TREE | RF |
|------|------|------|------|--------|------|------|
| 0.844 | 0.538 | 0.586 | 0.48 | 0.6369 | 0.392 | 0.238 |

*Table 2: Overall test error rates by class.*

If patterns are to be classified on a class-by-class basis then the method with the lowest error rate for each class can be used instead of the method with the overall lowest error rate. Some interesting results can be seen in *Table 3*, namely that no method is particularly reliable at predicting classes B, C, and W due to low observations of those classes. Class G and class H are

6

best predicted by QDA, suggesting the assumption of normality for those data is reasonable and can be modelled with a quadratic function. Linear Regression does not perform well on any of the classes suggesting complex nonlinear relationships between the predictors and the classes. This is not surprising from looking at boxplots and pairs plots in *Figure 1* and *Figure 2* of the data. Class I and Class Y are predicted very well by Random Forests with error rates under 10%. KNN was expected to perform better than more biased models but has about the same error rate as CVA, LDA, and QDA. The higher than expected error rate in KNN is likely due to the disproportionate number of classes in the observations and the relatively high number of predictors used in modelling.

|   | LR | KNN | CVA | LDA | QDA | TREE | RF |
|---|---|---|---|---|---|---|---|
| A | 0.9845 | 0.4811 | 0.5440 | 0.1088 | 0.7306 | 0.1865 | 0.1088 |
| B | 1.0000 | NaN | 1.0000 | NA | NA | 1.0000 | NA |
| C | 0.6667 | NaN | 1.0000 | 1.0000 | NA | 1.0000 | 0.6667 |
| D | 0.6897 | 1.0000 | 0.8276 | 1.0000 | 0.7931 | 1.0000 | 0.7586 |
| E | 0.6383 | 0.5833 | 0.6596 | 0.9149 | 0.7234 | 0.3404 | 0.2553 |
| F | 0.9775 | 0.6163 | 0.6854 | 0.9213 | 0.7416 | 0.6966 | 0.4045 |
| G | 0.9412 | 1.0000 | 0.4118 | 0.9412 | 0.2353 | 1.0000 | 0.6471 |
| H | 1.0000 | 0.6500 | 0.5714 | 0.9643 | 0.2143 | 0.5714 | 0.3929 |
| I | 0.7826 | 0.1905 | 0.1739 | 0.1304 | 0.2174 | 0.0652 | 0.0217 |
| W | 1.0000 | NaN | 0.8000 | 1.0000 | NA | 1.0000 | 0.8000 |
| X | 1.0000 | 0.1000 | 0.3913 | 0.5652 | 0.5652 | 0.4783 | 0.2174 |
| Y | 1.0000 | 0.2500 | 0.2500 | 0.5000 | 0.1875 | 0.2500 | 0.0625 |

*Table 3: Error rates by class for each method.*

## 6. CONCLUSION

Seven classification methods were evaluated in the classification of copyists for a subset of the Avila Bible dataset. Due to the highly complex and non-normal relationships between predictors, more flexible methods performed better than parametric models. Random Forests performed the best with an overall classification rate of 0.238. If classification of a specific copyist were desired, then the method with the lowest error rate in *Table 3* can be used for prediction rather than Random Forests for increased accuracy. Namely, Classes G and H are best predicted by QDA. Classes B, C, and W were often taken out of the training and validation sets due to too few observations. Interestingly, KNN performs about average compared to other classification methods.

There are a number of techniques that are beyond the scope of this report that could further improve the prediction accuracy of the methods considered. Only a subset of 2000 observations were used in modelling. By using the full dataset of 20,000 observations

available from the website [2] classes won't be removed as often when training the model and accuracy for those classes should increase. Another way to improve prediction accuracy is to perform dimension reduction on the predictors. Some classes could be best predicted by only a few predictors. By considering all the predictors, reduceable error is contained in the prediction estimates for the classes. Reducing the number of dimensions could improve accuracy of all models but particularly: Linear Regression, KNN, CVA, LDA, and QDA. The error rate of Classification Trees could be improved by boosting the trees. The methods considered in this report attempted to classify the observations to all classes. An attempt at reducing overfitting of more rigid models could be attained by predicting a singular class rather than all classes. In addition to finer tuning of the models, other methods such as Neural Networks, State Vector Machines, and Logistic Regression could be considered.

## 7. Bibliography

[1] De Stefano, C., Maniaci, M., Fontanella, F., Scotto di Freca, A. (2018) Reliable writer identification in medieval manuscripts through page layout features: The "Avila" Bible case. *Engineering Applications of Artificail Intelligence*, 72 pp. 99-110

[2] De Stefano, C., Maniaci, M., Fontanella, F., Scotto di Freca, A, "Avila Data Set", *UCI Machine Learning* Repository, archive.ics.uci.edu/ml/datasets/Avila

[3] James, G., Witten, D., Hastie, T., Tibshirani, R. (2013) *An Introduction to Statistical Learning* 1ed. Springer-Verlag New York.

[4] Hastie, T., Tibshirani, R., Friedman, J. (2009) *The Elements of Statistical Learning*. 2ed. Springer-Verlag New York.

[5] Hewson, P., (2009), *Multivariate Statistics with R*

```r
1.  ## Author: Robert Edwards
2.  ## University of Glasgow
3.  ## Multivariate Methods 2018
4.  ## Classification Project
5.
6.  setwd("~/OneDrive - University of Glasgow/University of Glasgow/Multivariate Methods/Pr
    oject")
7.  source("Methods.R")
8.  library(class)
9.  library(MASS)
10. library(e1071)
11. library(knitr)
12. library(kableExtra)
13.
14. # Load data
15. avila <- Avila.48
16.
17. ## Creating a data frame with the Class outcome and the 12 explanatory variables
18. auth.lab <- avila[,11]
19. data.avila <- as.data.frame(cbind(auth.lab, avila[,-11]))
20. colnames(data.avila) <- c("Class", "F1", "F2", "F3", "F4", "F5", "F6", "F7", "F8", "F9"
    , "F10")
21.
22.
23. ############################################
24. ## Inspecting the Data
25. ############################################
26. colors1 <- c("brown3", "brown3", "brown3", "darkgoldenrod2", "darkgoldenrod2", "darkgol
    denrod2", "darkgoldenrod2", "chocolate1", "chocolate1", "chocolate1")
27. boxplot(data.avila[-1], main="Extracted Features",
28.         col=colors1,
29.         pars=list(outcol=colors1))
30.
31. kable(round(cov(data.avila[-1]), 4), caption="Covariance Matrix") %>%
32.   kable_styling(bootstrap_options = "striped",
33.                 full_width = F,
34.                 position = "center")
35.
36. kable(round(cor(data.avila[-1]), 4), caption="Correlation Matrix") %>%
37.   kable_styling(bootstrap_options = "striped",
38.                 full_width = F,
39.                 position = "center")
40.
41. class.count <- as.data.frame(table(data.avila[1]))
42. colnames(class.count) <- c("Class", "Frequency")
43. class.count <- class.count[c("Frequency", "Class")]
44.
45. kable(t(class.count), caption="Class Frequency") %>%
46.   kable_styling(bootstrap_options = "striped",
47.                 full_width = F,
48.                 position = "center")
49.
50. dddpanel.hist <- function(x, ...)
51. {
52.   usr <- par("usr"); on.exit(par(usr))
53.   par(usr = c(usr[1:2], 0, 1.5) )
54.   h <- hist(x, plot = FALSE)
55.   breaks <- h$breaks; nB <- length(breaks)
56.   y <- h$counts; y <- y/max(y)
57.   rect(breaks[-nB], 0, breaks[-1], y, col = "cyan", ...)
58. }
```

```r
59.
60. pairs(data.avila[, -1], main="Extracted Features",
61.       diag.panel=panel.hist,
62.       pch=20)
63.
64.
65. ############################################
66. ## Splitting the data frame into 3 sets: training, validation and test.
67. ############################################
68. n <- nrow(avila)
69. ind1 <- sample(c(1:n), round(n / 2))          # Training set
70. ind2 <- sample(c(1:n)[-ind1], round(n / 4))   # Validation set
71. ind3 <- setdiff(c(1:n), c(ind1, ind2))        # Test set
72. train.avila <- data.avila[ind1, ]
73. valid.avila <- data.avila[ind2, ]
74. test.avila <- data.avila[ind3, ]
75.
76. table(train.avila$Class)
77. table(valid.avila$Class)
78. table(test.avila$Class)
79.
80.
81. ############################################
82. ## Optimizing KNN
83. ############################################
84. rect.err.rates <- as.data.frame(matrix(data=NA, nrow=itr/2, ncol=5))
85. names(rect.err.rates) <- c("K", "dist=0.1", "dist=1", "dist=2", "dist=100")
86. tri.err.rates <- as.data.frame(matrix(data=NA, nrow=itr/2, ncol=5))
87. names(tri.err.rates) <- c("K", "dist=0.1", "dist=1", "dist=2", "dist=100")
88. gaus.err.rates <- as.data.frame(matrix(data=NA, nrow=itr/2, ncol=5))
89. names(gaus.err.rates) <- c("K", "dist=0.1", "dist=1", "dist=2", "dist=100")
90. opt.err.rates <- as.data.frame(matrix(data=NA, nrow=itr/2, ncol=5))
91. names(opt.err.rates) <- c("K", "dist=0.1", "dist=1", "dist=2", "dist=100")
92.
93. itr <- 20
94. for (i in seq(1, itr, 2)) {
95.   set.seed(i)
96.   n <- nrow(avila)
97.   ind1 <- sample(c(1:n), round(n / 2))          # Training set
98.   ind2 <- sample(c(1:n)[-ind1], round(n / 4))   # Validation set
99.   ind3 <- setdiff(c(1:n), c(ind1, ind2))        # Test set
100.        train.avila <- data.avila[ind1, ]
101.        valid.avila <- data.avila[ind2, ]
102.        test.avila <- data.avila[ind3, ]
103.
104.        rect.err.rates[ceiling(i/2),1] <- i
105.        rect.err.rates[ceiling(i/2),2] <- classify.knn(train.avila, valid.avila, k=i,
     dist=0.1, ker="rectangular") #
106.        rect.err.rates[ceiling(i/2),3] <- classify.knn(train.avila, valid.avila, k=i,
     dist=1, ker="rectangular") # Manhattan
107.        rect.err.rates[ceiling(i/2),4] <- classify.knn(train.avila, valid.avila, k=i,
     dist=2, ker="rectangular") # Euclidean
108.        rect.err.rates[ceiling(i/2),5] <- classify.knn(train.avila, valid.avila, k=i,
     dist=10, ker="rectangular") # Chebyshev's
109.
110.        tri.err.rates[ceiling(i/2),1] <- i
111.        tri.err.rates[ceiling(i/2),2] <- classify.knn(train.avila, valid.avila, k=i, d
     ist=0.1, ker="triangular") #
112.        tri.err.rates[ceiling(i/2),3] <- classify.knn(train.avila, valid.avila, k=i, d
     ist=1, ker="triangular") # Manhattan
```

```r
113.          tri.err.rates[ceiling(i/2),4] <- classify.knn(train.avila, valid.avila, k=i, d
     ist=2, ker="triangular") # Euclidean
114.          tri.err.rates[ceiling(i/2),5] <- classify.knn(train.avila, valid.avila, k=i, d
     ist=10, ker="triangular") # Chebyshev's
115.
116.          gaus.err.rates[ceiling(i/2),1] <- i
117.          gaus.err.rates[ceiling(i/2),2] <- classify.knn(train.avila, valid.avila, k=i,
     dist=0.1, ker="gaussian") #
118.          gaus.err.rates[ceiling(i/2),3] <- classify.knn(train.avila, valid.avila, k=i,
     dist=1, ker="gaussian") # Manhattan
119.          gaus.err.rates[ceiling(i/2),4] <- classify.knn(train.avila, valid.avila, k=i,
     dist=2, ker="gaussian") # Euclidean
120.          gaus.err.rates[ceiling(i/2),5] <- classify.knn(train.avila, valid.avila, k=i,
     dist=10, ker="gaussian") # Chebyshev's
121.
122.          opt.err.rates[ceiling(i/2),1] <- i
123.          opt.err.rates[ceiling(i/2),2] <- classify.knn(train.avila, valid.avila, k=i, d
     ist=0.1, ker="optimal") #
124.          opt.err.rates[ceiling(i/2),3] <- classify.knn(train.avila, valid.avila, k=i, d
     ist=1, ker="optimal") # Manhattan
125.          opt.err.rates[ceiling(i/2),4] <- classify.knn(train.avila, valid.avila, k=i, d
     ist=2, ker="optimal") # Euclidean
126.          opt.err.rates[ceiling(i/2),5] <- classify.knn(train.avila, valid.avila, k=i, d
     ist=10, ker="optimal") # Chebyshev's
127.        }
128.
129.        par(mfcol=c(2,2), mar=c(4.1, 4.1, 4.1, 2), oma=c(0,0,0,6), xpd=NA)
130.        plot(rect.err.rates[,1], rect.err.rates[,2],
131.            main="Rectangular",
132.            xlab="K Values",
133.            ylab="Error Rate",
134.            ylim=c(0,1),
135.            col="red")
136.        points(rect.err.rates[,1], rect.err.rates[,3], col="purple")
137.        points(rect.err.rates[,1], rect.err.rates[,4], col="blue")
138.        points(rect.err.rates[,1], rect.err.rates[,5], col="orange")
139.
140.        plot(tri.err.rates[,1], tri.err.rates[,2],
141.            main ="Triangular",
142.            xlab="K Values",
143.            ylab="Error Rate",
144.            ylim=c(0,1),
145.            col="red")
146.        points(tri.err.rates[,1], tri.err.rates[,3],col="purple")
147.        points(tri.err.rates[,1], tri.err.rates[,4],col="blue")
148.        points(tri.err.rates[,1], tri.err.rates[,5],col="orange")
149.
150.        plot(gaus.err.rates[,1], gaus.err.rates[,2],
151.            main ="Guassian",
152.            xlab="K Values",
153.            ylab="Error Rate",
154.            ylim=c(0,1),
155.            col="red")
156.        points(gaus.err.rates[,1], gaus.err.rates[,3],col="purple")
157.        points(gaus.err.rates[,1], gaus.err.rates[,4],col="blue")
158.        points(gaus.err.rates[,1], gaus.err.rates[,5],col="orange")
159.
160.        plot(opt.err.rates[,1], opt.err.rates[,2],
161.            main ="Optimal",
162.            xlab="K Values",
163.            ylab="Error Rate",
```

```r
164.            ylim=c(0,1),
165.            col="red")
166.        points(opt.err.rates[,1], opt.err.rates[,3],col="purple")
167.        points(opt.err.rates[,1], opt.err.rates[,4],col="blue")
168.        points(opt.err.rates[,1], opt.err.rates[,5],col="orange")
169.        leg <- c("0.1", "1", "2", "10")
170.        legend("topright", legend = leg, col=c("red","purple","blue","orange"), pch=16,
    inset=c(-0.75, -2), bty="n")
171.        #legend("topright", legend = leg, col=c("red","purple","blue","orange"),  pch=16
    , inset=c(-1.3, -4), bty="n")
172.        legend("topright", legend="Minkowski\nDistance", inset=c(-1.15, -3), bty="n")
173.
174.
175.        ###########################################
176.        ## Model Selection
177.        ###########################################
178.        set.seed(2)
179.        itr=1
180.        err.rates <- as.data.frame(matrix(data=NA, nrow=itr, ncol=8))
181.        names(err.rates) <- c("LR", "KNN", "CVA", "LDA", "QDA", "TREE", "RF")
182.
183.
184.        for (i in 1:itr) {
185.          n <- nrow(avila)
186.          ind1 <- sample(c(1:n), round(n / 2))          # Training set
187.          ind2 <- sample(c(1:n)[-ind1], round(n / 4))   # Validation set
188.          ind3 <- setdiff(c(1:n), c(ind1, ind2))        # Test set
189.          train.avila <- data.avila[ind1, ]
190.          valid.avila <- data.avila[ind2, ]
191.          test.avila <- data.avila[ind3, ]
192.
193.          mcr <- classify.lr(train.avila, valid.avila)
194.          err.rates[i,1] <- mcr[1]
195.          err.rate.lr <- as.data.frame(unlist(mcr[-1]))
196.          mcr <- classify.knn(train.avila, valid.avila, k=11, dist=10, ker="rectangular"
    )
197.          err.rates[i,2] <- mcr[1]
198.          err.rate.knn <- as.data.frame(unlist(mcr[-1]))
199.          mcr <- classify.cva(train.avila, valid.avila)
200.          err.rates[i,3] <- mcr[1]
201.          err.rate.cva <- as.data.frame(unlist(mcr[-1]))
202.          mcr <- classify.lda(train.avila, valid.avila)
203.          err.rates[i,4] <- mcr[1]
204.          err.rate.lda <- as.data.frame(unlist(mcr[-1]))
205.          mcr <- classify.qda(train.avila, valid.avila)
206.          err.rates[i,5] <- mcr[1]
207.          err.rate.qda <- as.data.frame(unlist(mcr[-1]))
208.          mcr <- classify.tree(train.avila, valid.avila)
209.          err.rates[i,6] <- mcr[1]
210.          err.rate.tree <- as.data.frame(unlist(mcr[-1]))
211.          mcr <- classify.rf(train.avila, valid.avila)
212.          err.rates[i,7] <- mcr[1]
213.          err.rate.rf <- as.data.frame(unlist(mcr[-1]))
214.        }
215.
216.        err.rates
217.        boxplot(err.rates[-8], main="Classification Error Rates by Method",
218.                ylab="Error Rate",
219.                col=c("firebrick", "darkorchid3", "dodgerblue2", "dodgerblue2", "dodgerb
    lue2", "chartreuse3", "chartreuse3"))
220.
```

```r
221.        kable(round(err.rates, 4)) %>%
222.          kable_styling(bootstrap_options = "striped", full_width = F, position = "cente
    r")
223.          #row_spec(0, angle = -45)
224.
225.        ####################################
226.        ## Model Selection
227.        ####################################
228.        test.err.rates <- as.data.frame(matrix(data=NA, nrow=itr, ncol=7))
229.        names(test.err.rates) <- c("LR", "KNN", "CVA", "LDA", "QDA", "TREE", "RF")
230.        mcr <- classify.lr(train.avila, test.avila)
231.        test.err.rates[1] <- mcr[1]
232.        mcr <- classify.knn(train.avila, test.avila, k=17, dist=10, ker="rectangular")
233.        test.err.rates[2] <- mcr[1]
234.        mcr <- classify.cva(train.avila, test.avila)
235.        test.err.rates[3] <- mcr[1]
236.        mcr <- classify.lda(train.avila, test.avila)
237.        test.err.rates[4] <- mcr[1]
238.        mcr <- classify.qda(train.avila, test.avila)
239.        test.err.rates[5] <- mcr[1]
240.        mcr <- classify.tree(train.avila, test.avila)
241.        test.err.rates[6] <- mcr[1]
242.        mcr <- classify.rf(train.avila, test.avila)
243.        test.err.rates[7] <- mcr[1]
244.        test.err.rates
245.
246.        kable(round(test.err.rates, 4)) %>%
247.          kable_styling(bootstrap_options = "striped", full_width = F, position = "cente
    r")
248.
249.        ####################################
250.        ## Class Specific Error Rates
251.        ####################################
252.        err.rates.class <- cbind(err.rate.lr,
253.                                 err.rate.knn,
254.                                 err.rate.cva,
255.                                 err.rate.lda,
256.                                 err.rate.qda,
257.                                 err.rate.tree,
258.                                 err.rate.rf)
259.        rownames(err.rates.class) <- c("A","B","C","D","E","F","G","H","I","W","X","Y")

260.        colnames(err.rates.class) <- c("LR", "KNN", "CVA", "LDA", "QDA", "TREE", "RF")
261.        err.rates.class
262.
263.        kable(round(err.rates.class, 4)) %>%
264.          kable_styling(bootstrap_options = "striped", full_width = F, position = "cente
    r")
```

```r
1.  ## Author: Robert Edwards
2.  ## University of Glasgow
3.  ## Multivariate Methods 2018
4.  ## All Classifying Methods
5.
6.  library(kknn)
7.
8.  ##########################################
9.  ### Linear Regression
```

```
10. ############################################
11. classify.lr <- function(train.data, valid.data) {
12.
13.     # Change Class labels to numeric
14.     train.data <- as.data.frame(cbind(as.numeric(train.data[,1]), train.data[,-1]))
15.     valid.data <- as.data.frame(cbind(as.numeric(valid.data[,1]), valid.data[,-1]))
16.
17.     print("as numeric train data")
18.     print(table(as.numeric(train.data[,1])))
19.     print(table(as.numeric(train.data[,1]))[1])
20.     colnames(train.data) <- c("Class", "F1", "F2", "F3", "F4", "F5", "F6", "F7", "F8", "F
    9", "F10")
21.     colnames(valid.data) <- c("Class", "F1", "F2", "F3", "F4", "F5", "F6", "F7", "F8", "F
    9", "F10")
22.
23.
24.     ## Fitting the linear model to the training data.
25.     res <- lm(Class~.,data=train.data)
26.
27.     ## Using the predict function to predict for the validation data.
28.     pred.valid <- predict(res, valid.data[, -1])
29.
30.     ## Changing the predictions to predited labels.
31.     pred.valid.label <- round(pred.valid)
32.
33.     pred.valid.label[pred.valid.label < 1] <- 1
34.     pred.valid.label[pred.valid.label > 12] <- 12
35.
36.
37.     ## Cross-classification table.
38.     xtab <- table(valid.data[, 1], pred.valid.label)
39.
40.     ## Total Correctly Classiffied
41.     corr.class.total <- sum(pred.valid.label == valid.data[, 1])
42.     ## Correct classification rate (CCR).
43.     corr.class.rate <- sum(pred.valid.label == valid.data[, 1]) / length(pred.valid.label
    )
44.     ## Class-specific CCRs.
45.     ccr.class <- diag(xtab) / rowSums(xtab)
46.
47.     lr.err <- 1 - sum(diag(xtab)) / sum(xtab)
48.
49.     errors <- list(lr.err, 1-ccr.class)
50.     return(errors)
51. } # End classify.lr
52.
53.
54. ############################################
55. ### KNN
56. ############################################
57. classify.knn <- function(train.data, test.data, k=1, dist=2, ker="rectangular") {
58.
59.     n.test <- nrow(test.data)
60.
61.     k.nn.tr <- kknn(formula=formula(Class~.), train=train.data, test=train.data[-
    1], k=k, distance=dist, kernel=ker)
62.     k.nn.tst <- kknn(formula=formula(Class~.), train=train.data, test=test.data[-
    1], k=k, distance=dist, kernel=ker)
63.
64.     ## Class-specific CCRs.
65.     xtab <- table(k.nn.tst$fitted.values, test.data[, 1])
```

```r
66.    ccr.class <- diag(xtab) / rowSums(xtab)
67.    ## Calculate and store the test set correct classification rate
68.    knn.err <- sum(k.nn.tst$fitted.values == test.data[, 1]) / n.test
69.
70.    errors <- list(knn.err, 1-ccr.class)
71.    return(errors)
72. } # end KNN
73.
74.
75.
76. ############################################
77. ### CVA
78. ############################################
79. classify.cva <- function(train.data, valid.data) {
80.    ## Set Prior probabilites equal
81.    data.lda <- lda(train.data[, -
    1], train.data[, 1], prior=rep(1/nlevels(train.data[, 1]), 12) )
82.
83.
84.    ## predict the labels based on the model.
85.    data.valid.ld <- predict(data.lda, newdata=valid.data[, -1])
86.
87.    ## cross-classification table
88.    xtab <- table(valid.data[, 1], data.valid.ld$class)
89.
90.    ## Total Correctly Classiffied
91.    corr.class.total <- sum(diag(xtab))
92.    ## Correct Classification Rate
93.    corr.class.rate <- sum(diag(xtab)) / sum(xtab)
94.    ## Misclassification rate =  1 - correct classifcation rate
95.    mcr <- 1 - sum(diag(xtab)) / sum(xtab)
96.
97.    ## Class-specific CCRs.
98.    ccr.class <- diag(xtab) / rowSums(xtab)
99.
100.        ## Error Rate
101.        cva.err <- 1 - sum(diag(xtab)) / sum(xtab)
102.
103.        errors <- list(cva.err, 1-ccr.class)
104.        return(errors)
105.      } #end classify.cva
106.
107.
108.
109.      ############################################
110.      ## LDA
111.      ############################################
112.      classify.lda <- function(train.data, valid.data) {
113.
114.        ## Check if n > 0 and remove those classes
115.        if ( sum(table(train.data$Class) < 1 ) != 0) {
116.          too.few <- levels(train.data$Class)[table(train.data$Class) < 1]
117.          cat("Removing from train data: ", too.few, "\n")
118.
119.          for (i in 1:length(too.few)) {
120.            train.data <- subset(train.data, Class != too.few[i]) # if n < 0 for each
    class
121.            valid.data <- subset(valid.data, Class != too.few[i]) # if n < 0 for each
    class
122.          }
123.
```

```r
124.          train.data$Class <- factor(train.data$Class)
125.          valid.data$Class <- factor(valid.data$Class)
126.        }
127.        if ( sum(table(valid.data$Class) < 1 ) != 0) {
128.          too.few <- levels(valid.data$Class)[table(valid.data$Class) < 1]
129.          cat("Removing from valid data: ", too.few, "\n")
130.
131.          for (i in 1:length(too.few)) {
132.            train.data <- subset(train.data, Class != too.few[i]) # if n < 0 for each
      class
133.            valid.data <- subset(valid.data, Class != too.few[i]) # if n < 0 for each
      class
134.          }
135.          train.data$Class <- factor(train.data$Class)
136.          valid.data$Class <- factor(valid.data$Class)
137.        }
138.
139.
140.        ## Let lda set priors
141.        data.lda <- lda(Class~., data=train.data)
142.
143.        ## Predicting the classes on the data.
144.        data.valid.pred.lda<- predict(data.lda, valid.data)
145.
146.        ## Constructing a cross-
      classification table of the real versus predicted classifications.
147.        xtab<-table(valid.data$Class, data.valid.pred.lda$class)
148.
149.        ## Class-specific CCRs.
150.        ccr.class <- diag(xtab) / rowSums(xtab)
151.
152.        ## Calculating the misclassification rate.
153.        lda.err <- 1 - sum(diag(xtab)) / sum(xtab)
154.
155.        errors <- list(lda.err, 1-ccr.class)
156.        return(errors)
157.      } # end classify.lda
158.
159.
160.      ###########################################
161.      ## QDA
162.      ###########################################
163.      classify.qda <- function(train.data, valid.data) {
164.
165.        ## Check if n < p and remove those classes
166.        if ( sum(table(train.data$Class) <= ncol(train.data[-1]) ) != 0) {
167.          too.few.train <- levels(train.data$Class)[table(train.data$Class) <= ncol(tr
    ain.data[-1])]
168.          cat("Removing from train data: ", too.few.train, "\n")
169.        }
170.        if ( sum(table(valid.data$Class) <= ncol(valid.data[-1]) ) != 0) {
171.          too.few.valid <- levels(valid.data$Class)[table(valid.data$Class) <= ncol(va
    lid.data[-1])]
172.          cat("Removing from train data: ", too.few.valid, "\n")
173.        }
174.
175.        too.few <- c(too.few.train, too.few.valid)
176.
177.        if ( length(too.few) != 0) {
178.          for (i in 1:length(too.few)) {
```

```r
179.            train.data <- subset(train.data, Class != too.few[i]) # if n < p for each
     class
180.            valid.data <- subset(valid.data, Class != too.few[i]) # if n < p for each
     class
181.           }
182.         train.data$Class <- factor(train.data$Class)
183.         valid.data$Class <- factor(valid.data$Class)
184.       }
185.
186.       data.qda <- qda(Class~., data=train.data)
187.
188.       ## Prediction
189.       data.valid.pred.qda <- predict(data.qda, valid.data)
190.
191.       ## Constructing a cross-
     classification table of the real versus predicted classifications.
192.       xtab <- table(valid.data$Class, data.valid.pred.qda$class)
193.
194.       ## Class-specific CCRs.
195.       ccr.class <- diag(xtab) / rowSums(xtab)
196.
197.       ## Calculating the error rate.
198.       qda.err <- 1 - sum(diag(xtab)) / sum(xtab)
199.
200.       errors <- list(qda.err, 1-ccr.class)
201.       return(errors)
202.     } # end classify.qda
203.
204.
205.       #############################################
206.       ## TREES
207.       #############################################
208.     classify.tree <- function(train.data, valid.data) {
209.
210.       set.seed(126)
211.       library(rpart)
212.       ## Set"class" to ensure we get a classification
213.       data.rp <- rpart(Class~. , data=train.data, method="class")
214.
215.       ## We predict the classifications based on the fitted tree using the predict c
     ommand with the argument type set to "class", otherwise it returns a matrix of posterio
     r probabilities of each obseration belonging to each class.
216.       data.valid.rp.pred <- predict(data.rp, valid.data, type="class")
217.
218.       ## Constructing a cross-
     classification table of the real versus predicted classifications.
219.       xtab <- table(valid.data$Class, data.valid.rp.pred)
220.
221.       ## Taking a look at the posterior probabilities for the observations we got th
     e prediction wrong for.
222.       data.valid.rp.prob <- predict(data.rp, valid.data)
223.       data.valid.rp.prob[valid.data$Class != data.valid.rp.pred, ]
224.
225.       ## Class-specific CCRs.
226.       ccr.class <- diag(xtab) / rowSums(xtab)
227.       ## Calculating the misclassification rate.
228.       tree.err <- 1 - sum(diag(xtab)) / sum(xtab)
229.
230.       errors <- list(tree.err, 1-ccr.class)
231.       return(errors)
232.     } # end classify.tree
```

```r
233.
234.
235.
236.        #################################################
237.        ## RANDOM FORESTS
238.        #################################################
239.        classify.rf <- function(train.data, valid.data) {
240.          library(randomForest)
241.          set.seed(126)
242.
243.          too.few <- NA
244.          ## Remove empty classes
245.          if ( sum(table(train.data$Class) < 1 ) > 0) {
246.            too.few <- levels(train.data$Class)[table(train.data$Class) < 1]
247.            cat("Removing from train data: ", too.few, "\n")
248.
249.            for (i in 1:length(too.few)) {
250.              train.data <- subset(train.data, Class != too.few[i]) # if n < p for each
      class
251.              valid.data <- subset(valid.data, Class != too.few[i]) # if n < p for each
      class
252.            }
253.            train.data$Class <- factor(train.data$Class)
254.            valid.data$Class <- factor(valid.data$Class)
255.          }
256.          if ( sum(table(valid.data$Class) < 1 ) > 0) {
257.            too.few <- NA
258.            too.few <- levels(valid.data$Class)[table(valid.data$Class) < 1]
259.            cat("Removing from valid data: ", too.few, "\n")
260.
261.            for (i in 1:length(too.few)) {
262.              train.data <- subset(train.data, Class != too.few[i]) # if n < p for each
      class
263.              valid.data <- subset(valid.data, Class != too.few[i]) # if n < p for each
      class
264.            }
265.            train.data$Class <- factor(train.data$Class)
266.            valid.data$Class <- factor(valid.data$Class)
267.          }
268.
269.
270.          data.rf <- randomForest(Class~., data=train.data)
271.
272.          ## Prediction
273.          data.valid.rf.pred <- predict(data.rf, newdata=valid.data)
274.
275.          ## Cross classification
276.          xtab <- table(valid.data$Class, data.valid.rf.pred)
277.
278.          ## Class-specific CCRs.
279.          ccr.class <- diag(xtab) / rowSums(xtab)
280.          ## Calculating the misclassification rate.
281.          rf.err <- 1 - sum(diag(xtab)) / sum(xtab)
282.
283.          errors <- list(rf.err, 1-ccr.class)
284.          return(errors)
285.        } # end classify.rf
```