

TDDE01 - Computer lab 2

Group report

Emil Frid (emifr995)
Emir Alkazhami (emial637)
Nigel Cole (nigco547)

December 2019

Abstract

Each member contributed an assignment each to this report. Nigel contributed with assignment 1, Emir contributed with assignment 2 and Emil contributed with assignment 4. We were all active with solving the lab individually and shared ideas orally.

1 Assignment 1 - LDA and logistic regression

1.1 Task 1

Yes, the data for Male and Female are nicely split and the LDA boundary can be placed between them. The distributions have similar shapes but with different angles.

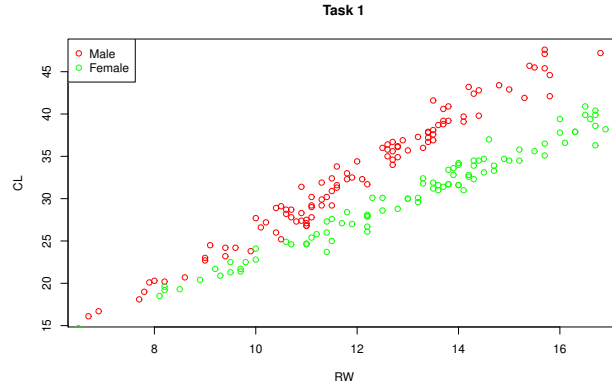


Figure 1: The male and female data groups for task 1.

1.2 Task 2

With LDA we see that a few of the data points shift sex. The points that shifted are located where the groups are connected. The fit is good. The misclassification error rate is 0.035. Confusion matrix:

$$\begin{array}{c} \text{True} \\ 0 \quad 1 \\ \text{Prediction} \begin{pmatrix} 97 & 4 \\ 3 & 96 \end{pmatrix} \end{array} \quad (1)$$

1.3 Task 3

With the priors $p(\text{Male}) = 0.9, p(\text{Female}) = 0.1$ we instead have a lot more points that shift sex from Female to Male with the misclassification error rate of 0.08. The prior probabilities are the probabilities of an observation coming from a particular group. Because the prior probability for Male is higher than the prior probability for Female, the function will make it more likely that an

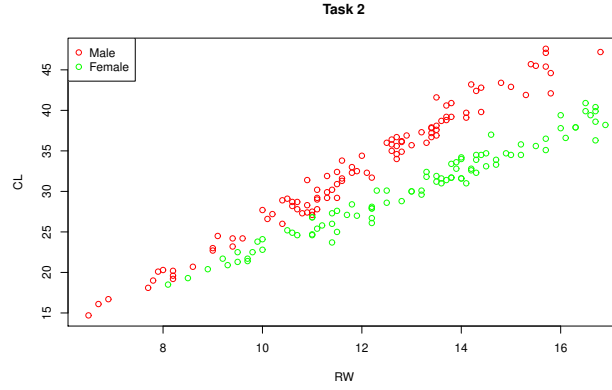


Figure 2: LDA implemented, a few points changed sex for task 2.

observation will be classified as a Male. The confusion matrix is:

$$\begin{array}{c} \text{True} \\ \begin{array}{cc} 0 & 1 \end{array} \\ \text{Prediction} \begin{array}{c} 0 \\ 1 \end{array} \begin{pmatrix} 84 & 0 \\ 16 & 100 \end{pmatrix} \end{array} \quad (2)$$

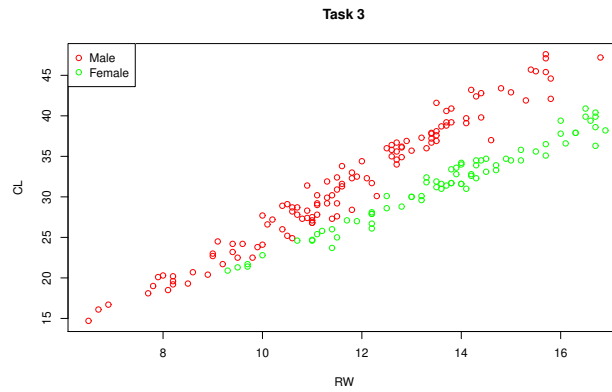


Figure 3: LDA with added priors, many points shift sex for task 3.

1.4 Task 4

With the *glm()* function we get a misclassification error rate of 0.035 which is the same as for task 2. Certain points have, compared to task 3 with LDA,

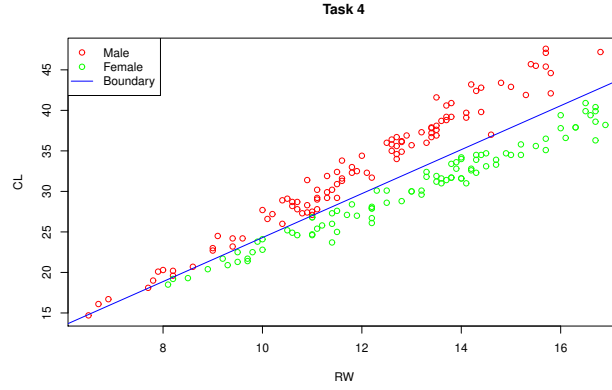


Figure 4: Using the $glm()$ function instead and also adding the decision boundary for task 4.

shifted from Male to Female. The actual points that are Male and Female are not the same as in task 2, even though the misclassification error rate is the same. The equation used for the decision boundary is

$$boundary = RW * x + CL * y + intercept \quad (3)$$

where y in R is

$$y = (boundary - m_coef[3] * x - m_coef[1]) / m_coef[2] \quad (4)$$

where the elements in m_coef are the coefficients of the $glm()$ model. The confusion matrix is:

$$\begin{array}{c} \text{True} \\ \begin{array}{cc} 0 & 1 \end{array} \\ \text{Prediction} \begin{array}{c} 0 \\ 1 \end{array} \begin{pmatrix} 97 & 4 \\ 3 & 96 \end{pmatrix} \end{array} \quad (5)$$

2 Assignment 2 - Analysis of credit scoring

2.1 Task 2

		Train	Test
Misclassification rates:	Gini	23.8%	37.2%
	Deviance	21.2%	26.8%

The deviance measure gives the best result.

2.2 Task 3

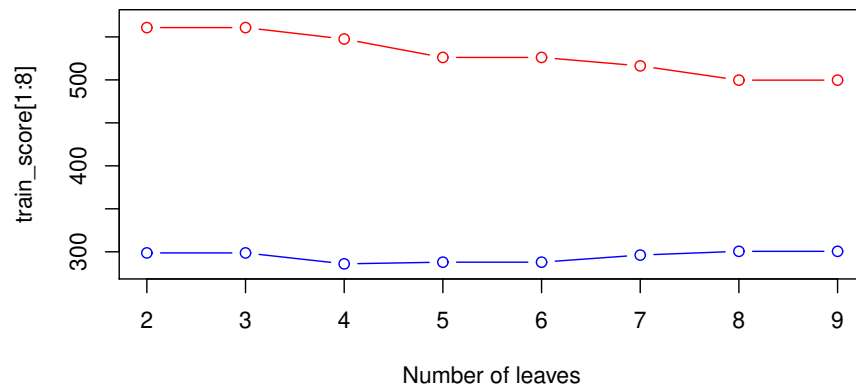


Figure 5: Red: train, Blue: test

Four leaves seem to give the best results.

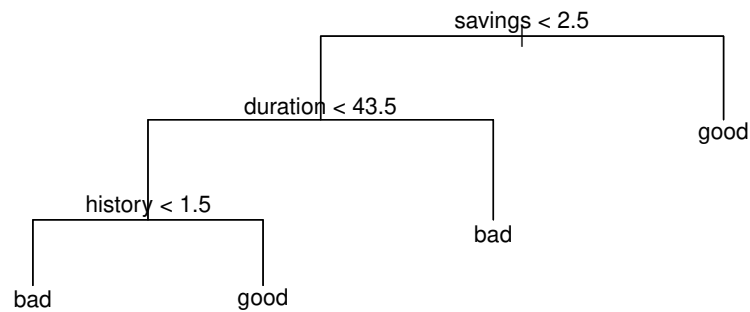


Figure 6: The resulting best pruned tree

The final tree has depth 3 and uses the variables: savings, duration and history. The variable names are not self-explanatory, but it seems like a person with large savings who is taking a loan over a short time span is more likely to pay the

debt back.

The misclassification rate on the test data is: 25.6%

Confusion matrix for test data:

$$\begin{array}{c} \text{Prediction} \end{array} \begin{array}{c} \text{True} \\ \text{Bad} \quad \text{Good} \end{array} \begin{pmatrix} 18 & 6 \\ 58 & 168 \end{pmatrix} \quad (6)$$

2.3 Task 4

Misclassification rate on training data: 30.0%

Misclassification rate on testing data: 31.6%

Confusion matrix for training data:

$$\begin{array}{c} \text{Prediction} \end{array} \begin{array}{c} \text{True} \\ \text{Bad} \quad \text{Good} \end{array} \begin{pmatrix} 95 & 98 \\ 52 & 225 \end{pmatrix} \quad (7)$$

Confusion matrix for test data:

$$\begin{array}{c} \text{Prediction} \end{array} \begin{array}{c} \text{True} \\ \text{Bad} \quad \text{Good} \end{array} \begin{pmatrix} 46 & 49 \\ 30 & 125 \end{pmatrix} \quad (8)$$

The decision tree gets a better misclassification rate than the Naive Bayes classifier. This is because the decision tree is better at classifying the 'Good' class than the 'Bad' class, and there is many more samples of class 'Good'.

This does not have to mean that the tree is a better classifier though. It seems biased towards the 'Good' class since it's the larger one, and it's possible to argue that the the Naive Bayes solution generalizes better for new data because the error rate is about the same for the two classes.

We try to explain this by the fact that the training data is unbalanced. There are 147 samples of class 'Bad' and 323 samples of class 'Good'. This will affect the tree classification since the trees use majority voting when assigning class labels to its leaves. Therefore the tree will be more prone the predict a sample as 'Good'. This is shown in confusion matrix (6) with many misclassifications of class 'Bad'.

The Naive Bayes classifier on the other hand, is based on Bayes theorem and takes the class-proportions in to account.

$$P(Y = y|X = x) \propto P(X = x|Y = y)P(Y = y)$$

We assume that the class distributions are Gaussian and that we can estimate $P(X = x|Y = y)$ from the data. $P(Y = y)$ are the class-proportions.

2.4 Task 5

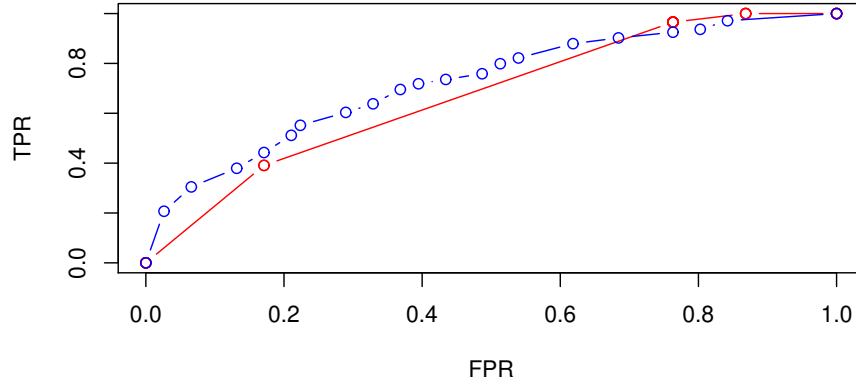


Figure 7: ROC curve. Red: Decision tree, Blue: Naive Bayes

The Naive Bayes classifier have more area under the curve and is therefore better for this data.

2.5 Task 6

Misclassification rate for the training data: 54.6%

Misclassification rate for the test data: 50.8%

Confusion matrix for training data:

$$\begin{array}{cc} & \begin{array}{cc} \text{True} \\ \text{Bad} & \text{Good} \end{array} \\ \begin{array}{c} \text{Prediction} \\ \text{Bad} \\ \text{Good} \end{array} & \begin{pmatrix} 137 & 263 \\ 10 & 90 \end{pmatrix} \end{array} \quad (9)$$

Confusion matrix for test data:

$$\begin{array}{cc} & \begin{array}{cc} \text{True} \\ \text{Bad} & \text{Good} \end{array} \\ \begin{array}{c} \text{Prediction} \\ \text{Bad} \\ \text{Good} \end{array} & \begin{pmatrix} 71 & 122 \\ 5 & 52 \end{pmatrix} \end{array} \quad (10)$$

In task 6 we have a worse misclassification rate than in task 4. This is because we are misclassifying a lot more samples of class 'Good' than in task 4. This happens because we are weighting the loss-values for classifying a sample of type 'Bad' by 10, basically shifting the our prediction boundary towards the 'Good' class.

In other words, we say that classifying a sample of class 'Bad' wrong is worse than misclassifying a sample of class 'Good'. Therefore we take a safe approach and say that we have to be really certain in our prediction to classify it as 'Good'.

3 Assignment 4 - Principal components

3.1 Task 1

A standard PCA was conducted using the feature space. The PCA showing how much variation is explained by each feature is presented in Figure 8.

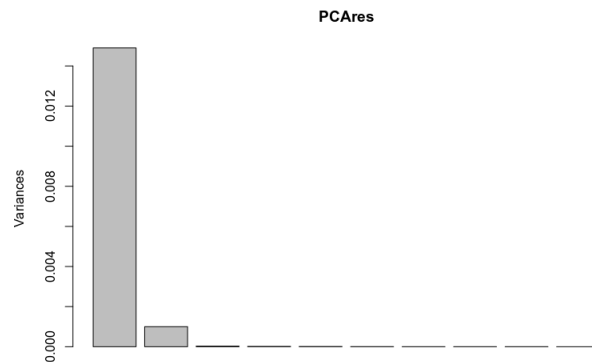


Figure 8: The PCA showing how much variation is explained by each feature.

In Figure 8 it is clear that the first two features stand for the majority of the variations, but not all variations and in some situations those small percentages are important. So it is not necessarily clear how many PC that should be extracted.

The first two components are enough to explain at least 99 % of the total variance.

Figure 9 presents the scores in the coordinates (PC1,PC2). In the figure it is seen that there are some outliers.

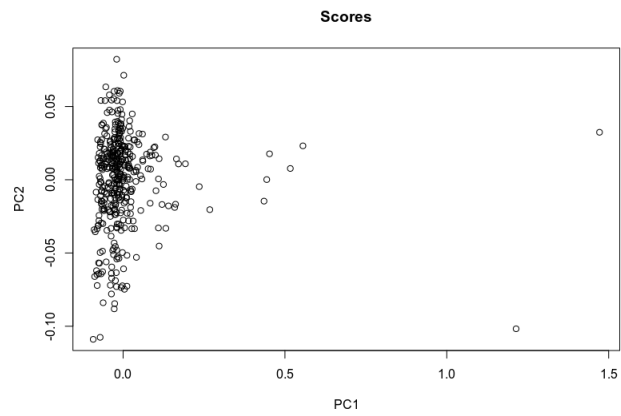


Figure 9: The first two scores.

3.2 Task 2

Figures 10 and 11 presents the trace plots of the loadings of the components selected in Task 1.

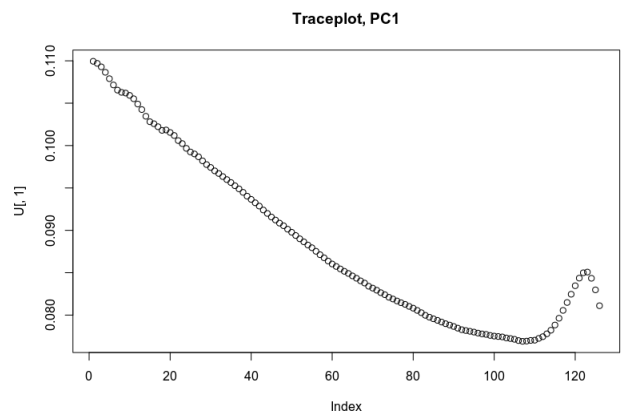


Figure 10: The trace plot for PC1.

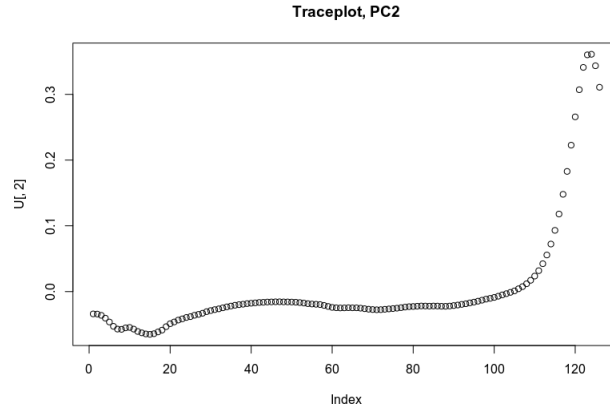


Figure 11: The trace plot for PC2.

The larger absolute value in the trace plots the more those features explain the PC. Therefore, it is seen in Figure 11 that PC2 is mostly explained by the last 14, or so, features.

3.3 Task 3

An Independent Component Analysis (ICA) was performed with the number of components selected in Task 1. The $W' = KW$ was computed and the columns of W' are presented in form of trace plots in Figures 12 and 13.

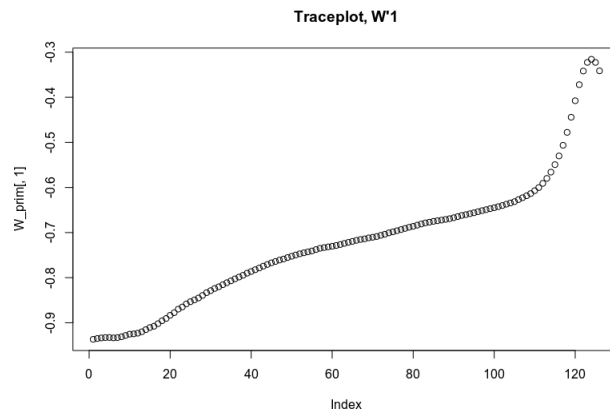


Figure 12: The trace plot for the first column in W' .

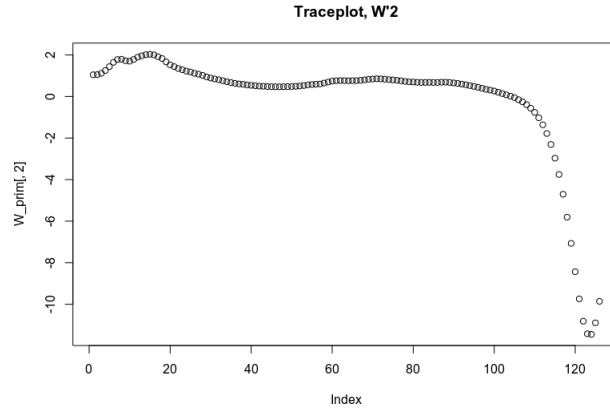


Figure 13: The trace plot for the first column in W' .

Compared to the trace plots in Task 2 these seem to be inverted and scaled larger. W' is the complete transformation and selection matrix for the ICA, the U_{ICA} .

Figure 14 presents the scores from ICA.

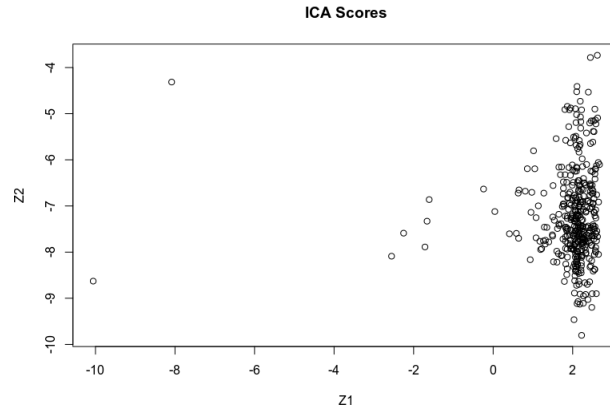


Figure 14: The ICA scores.

Just as for the trace plots the scores are inverted and scaled larger compared to the scores from PCA in Task 1.

A R code

Here is the code presented.

A.1 Assignment 1

```
1 RNGversion('3.5.1')
2 setwd("/home/nigel/Documents/Liu/TDDE01/Lab2")
3
4 # Task 1
5 data = read.csv("australian-crabs.csv")
6 CL = data$CL
7 RW = data$RW
8 sex = data$sex
9 sexMale = which(sex %in% "Male")
10 sexFemale = which(sex %in% "Female")
11
12 plot(RW[sexMale], CL[sexMale], col = "red", xlab = "RW", ylab = "CL",
13      , main = "Task 1")
14 points(RW[sexFemale], CL[sexFemale], col = "green")
15 legend("topleft", pch = 1, col = c("red", "green"), legend = c("Male", "Female"))
16
17 # Yes, data groups are nicely split and the LDA boundary can be placed nicely.
18
19 # Task 2
20 library(MASS)
21 model = lda(formula = sex ~ CL + RW, data = data)
22 pred_model = predict(object = model, newdata = data, type = "response")
23
24 pred_male = pred_model$class == "Male"
25 pred_female = pred_model$class == "Female"
26
27 plot(RW[pred_male], CL[pred_male], col = "red", xlab = "RW", ylab = "CL",
28      , main = "Task 2")
29 points(RW[pred_female], CL[pred_female], col = "green")
30 legend("topleft", pch = 1, col = c("red", "green"), legend = c("Male", "Female"))
31
32 # Misclassification
33 misclass=function(X,X1){n=length(X)
34 return(1-sum(diag(table(X,X1)))/n)
35 }
36
37 # Confusion matrix and misclass. error rate = 0.035
38 compare = 1:length(sex)
39 compare[pred_male] = "Male"
40 compare[pred_female] = "Female"
41 misclass(compare, sex)
42 table(compare, sex)
43
44 # Task 3
45 library(MASS)
46 model = lda(formula = sex ~ CL + RW, data = data, prior = c(0.1, 0.9))
```

```

44 | pred_model = predict(object = model, newdata = data, type = "
      response")
45
46 | pred_male = pred_model$class == "Male"
47 | pred_female = pred_model$class == "Female"
48
49 | plot(RW[pred_male], CL[pred_male], col = "red", xlab = "RW", ylab =
      "CL", main = "Task 3")
50 | points(RW[pred_female], CL[pred_female], col = "green")
51 | legend("topleft", pch = 1, col = c("red", "green"), legend = c("
      Male", "Female"))
52
53 | # Confusion matrix and misclass. error rate = 0.08
54 | compare = 1:length(sex)
55 | compare[pred_male] = "Male"
56 | compare[pred_female] = "Female"
57 | misclass(compare, sex)
58 | table(compare, sex)
59
60 | # Task 4
61 | model = glm(formula = sex ~ CL + RW, family = "binomial", data =
      data)
62 | pred_model = predict(object = model, newdata = data, type = "
      response")
63
64 | pred_male = pred_model > 0.5
65 | pred_female = pred_model <= 0.5
66
67 | plot(RW[pred_male], CL[pred_male], col = "red", xlab = "RW", ylab =
      "CL", main = "Task 4")
68 | points(RW[pred_female], CL[pred_female], col = "green")
69 | legend("topleft", pch = c(1, 1, NA), lty = c(NA, NA, 1), col = c("
      red", "green", "blue"), legend = c("Male", "Female", "Boundary"
      ))
70
71 | # boundary = RW*x + CL*y + intercept
72 | boundary = 0.5
73 | m_coef = model$coefficients
74 | x = c(6, 18)
75 | y = (boundary - m_coef[3] * x - m_coef[1]) / m_coef[2]
76 | lines(x, y, col = "blue")
77
78 | # Confusion matrix and misclass. error rate = 0.035
79 | compare = 1:length(sex)
80 | compare[pred_male] = "Male"
81 | compare[pred_female] = "Female"
82 | misclass(compare, sex)
83 | table(compare, sex)

```

Code/ass1.R

A.2 Assignment 2

```

1 | RNGversion('3.5.1')
2 | cat(rep("\n", 30))

```

```

3 rm(list = ls())
4
5 #####
6 #---Task 1---
7 #####
8 library(readxl)
9 source("lab1/missclass.R")
10
11 data = read_excel("lab2/creditscoring.xls")
12 data[["good_bad"]] = as.factor(data[["good_bad"]])
13
14 n=dim(data)[1]
15 set.seed(12345)
16 id=sample(1:n, floor(n*0.5))
17 train=data[id,]
18 id1=setdiff(1:n, id)
19 set.seed(12345)
20 id2=sample(id1, floor(n*0.25))
21 valid=data[id2,]
22 id3=setdiff(id1, id2)
23 test=data[id3,]
24
25 #####
26 #---Task 2---
27 #####
28 library(tree)
29 model = tree(good_bad~., data=train, split="deviance")
30 preds_train = predict(model, newdata = train, type="class")
31 preds_test = predict(model, newdata = test, type="class")
32 plot(model)
33 text(model, pretty=0)
34
35 deviance_train = missclass(train$good_bad, preds_train)
36 deviance_test = missclass(test$good_bad, preds_test)
37
38
39 model = tree(good_bad~., data=train, split="gini")
40 preds_train = predict(model, newdata = train, type="class")
41 preds_test = predict(model, newdata = test, type="class")
42 plot(model)
43 text(model, pretty=0)
44
45 gini_train = missclass(train$good_bad, preds_train)
46 gini_test = missclass(test$good_bad, preds_test)
47
48 #####
49 #---Task 3---
50 #####
51 model = tree(good_bad~., data=train, split="deviance")
52
53 train_score=rep(0,8)
54 test_score=rep(0,8)
55
56 for(i in 2:9) {
57   pruned_tree = prune.tree(model, best=i)
58   pred = predict(pruned_tree, newdata=valid, type="tree")
59   train_score[i-1] = deviance(pruned_tree)

```

```

60 | test_score[i-1] = deviance(pred)
61 | }
62 |
63 | plot(2:9, train_score[1:8], type="b", col="red", ylim=c(280,570),
64 |      xlab="Number of leaves")
65 | points(2:9, test_score[1:8], type="b", col="blue")
66 |
67 | BEST_INDEX = 4
68 | final_tree = prune.tree(model, best=BEST_INDEX)
69 | plot(final_tree)
70 | text(final_tree, pretty=0)
71 |
72 | pred_test = predict(final_tree, newdata=test, type="class")
73 | miss_test = missclass(test$good_bad, pred_test)
74 | test_table = table(pred_test, test$good_bad)
75 | #####
76 | #---Task 4---
77 | #####
78 | library(e1071)
79 | model = naiveBayes(good_bad~., data=train)
80 |
81 | preds_train = predict(model, newdata=train)
82 | preds_test = predict(model, newdata=test)
83 |
84 | train_table = table(preds_train, train$good_bad)
85 | train_rate = missclass(train$good_bad, preds_train)
86 |
87 | test_table = table(preds_test, test$good_bad)
88 | test_rate = missclass(test$good_bad, preds_test)
89 | print(test_table)
90 |
91 | #####
92 | #---Task 5---
93 | #####
94 |
95 | #Tree model
96 | model = tree(good_bad~., data=train, split="deviance")
97 | tree_model = prune.tree(model, best=BEST_INDEX)
98 |
99 | pred = predict(final_tree, newdata=test)
100 |
101 | PI = seq(0.00, 1.00, 0.05)
102 | tpr_vec_tree = rep(0, length(PI))
103 | fpr_vec_tree = rep(0, length(PI))
104 |
105 | for(ind in 1:length(PI)){
106 |   pi = PI[ind]
107 |   new_pred = rep("bad", length(pred[,2]))
108 |   new_pred[pred[,2] > rep(pi, length(pred[,2]))] = "good"
109 |
110 |   TP = sum(new_pred=="good" & test$good_bad=="good")
111 |   FP = sum(new_pred=="good" & test$good_bad=="bad")
112 |   TN = sum(new_pred=="bad" & test$good_bad=="bad")
113 |   FN = sum(new_pred=="bad" & test$good_bad=="good")
114 |
115 |   tpr_vec_tree[ind] = TP/(TP+FN)

```



```

116 |   fpr_vec_tree[ind] = FP/(FP+TN)
117 | }
118 |
119 | #Naive bayes model
120 | bayes_model = naiveBayes(good_bad~., data=train)
121 | pred = predict(bayes_model, newdata=test, type="raw")
122 |
123 | tpr_vec_bayes = rep(0, length(PI))
124 | fpr_vec_bayes = rep(0, length(PI))
125 |
126 | for(ind in 1:length(PI)){
127 |   pi = PI[ind]
128 |   new_pred = rep("bad", length(pred[,2]))
129 |   new_pred[pred[,2] > rep(pi, length(pred[,2]))] = "good"
130 |
131 |   TP = sum(new_pred=="good" & test$good_bad=="good")
132 |   FP = sum(new_pred=="good" & test$good_bad=="bad")
133 |   TN = sum(new_pred=="bad" & test$good_bad=="bad")
134 |   FN = sum(new_pred=="bad" & test$good_bad=="good")
135 |
136 |   tpr_vec_bayes[ind] = TP/(TP+FN)
137 |   fpr_vec_bayes[ind] = FP/(FP+TN)
138 | }
139 |
140 | plot(fpr_vec_tree, tpr_vec_tree, type="b", col="red", xlab="FPR",
141 |      ylab="TPR")
142 | points(fpr_vec_bayes, tpr_vec_bayes, type="b", col="blue")
143 | #####
144 | #---Task 6---
145 | #####
146 | model = naiveBayes(good_bad~., data=train)
147 |
148 | good_scale = 1
149 | bad_scale = 10
150 |
151 | preds_train = predict(model, newdata=train, type="raw")
152 | preds_train_temp = rep("bad", length(preds_train[,1]))
153 | preds_train_temp[bad_scale*preds_train[,1] < good_scale*preds_train
154 |   [,2]] = "good"
155 | preds_train = preds_train_temp
156 |
157 | preds_test = predict(model, newdata=test, type="raw")
158 | preds_test_temp = rep("bad", length(preds_test[,1]))
159 | preds_test_temp[bad_scale*preds_test[,1] < good_scale*preds_test
160 |   [,2]] = "good"
161 | preds_test = preds_test_temp
162 |
163 | train_table = table(preds_train, train$good_bad)
164 | train_rate = missclass(train$good_bad, preds_train)
165 |
166 | test_table = table(preds_test, test$good_bad)
167 | test_rate = missclass(test$good_bad, preds_test)
168 | print(test_table)

```

Code/lab2-2.R

A.3 Assignment 4

```

1 #setwd("./lab2")
2 library(readxl)
3 library(kknn)
4 library(MASS)
5 require(ggplot2)
6 require(scales)
7 require(gridExtra)
8 library(fastICA)
9 source("~/Documents/University/TDDE01/tdde01/lab1/missclass.R")
10
11 ## Task 1 #####
12 data = read.csv2("NIRSpectra.csv")
13
14 PCAres=prcomp(formula=~.-Viscosity, data = data)
15 # The eigenvalues of the covariance matrix
16 lambda=PCAsres$sdev^2
17 #Proportion of variation. First 2 is needed for >99 %
18 sprintf("%2.3f",lambda/sum(lambda)*100)
19
20 screeplot(PCAsres)
21
22 #Score plot
23 plot(PCAsres$x[,1], PCAsres$x[,2], main = "Scores", xlab = "PC1", ylab =
    "PC2")
24
25 ## Task 2 #####
26 U=PCAsres$rotation
27 #head(U)
28
29 #The larger abs value in the trace plots the more important those
30 #variables are.
31 plot(U[,1], main="Traceplot, PC1")
32 plot(U[,2], main="Traceplot, PC2")
33
34 ## Task 3 #####
35 # a)
36 set.seed(12345)
37
38 data2 = data
39 data2$Viscosity = c()
40
41 ICAsres$K
42
43 ICAsres = fastICA(data2,2)
44 # Down here is really a matrix multiplication
45 # but latex couldn't handle the strange syntax.
46 W_prim = ICAsres$K*ICAsres$W
47 plot(W_prim[,1], main="Traceplot, W1")
48 plot(W_prim[,2], main="Traceplot, W2")
49
50 # b) score
51 X = as.matrix(data2)
52
53 Z=ICAsres$S
54

```

```

55 plot(Z[,1], Z[,2], xlab = "Z1", ylab="Z2", main="ICA Scores")
56
57 # U_ICA = U*V <=> W= K*W
58 # where U=K from PCA and V=W from ICA.
59
60 # Z = S = X*U_ICA

```

Code/assignment4.R