# lab1

*Erik Tedhamre*

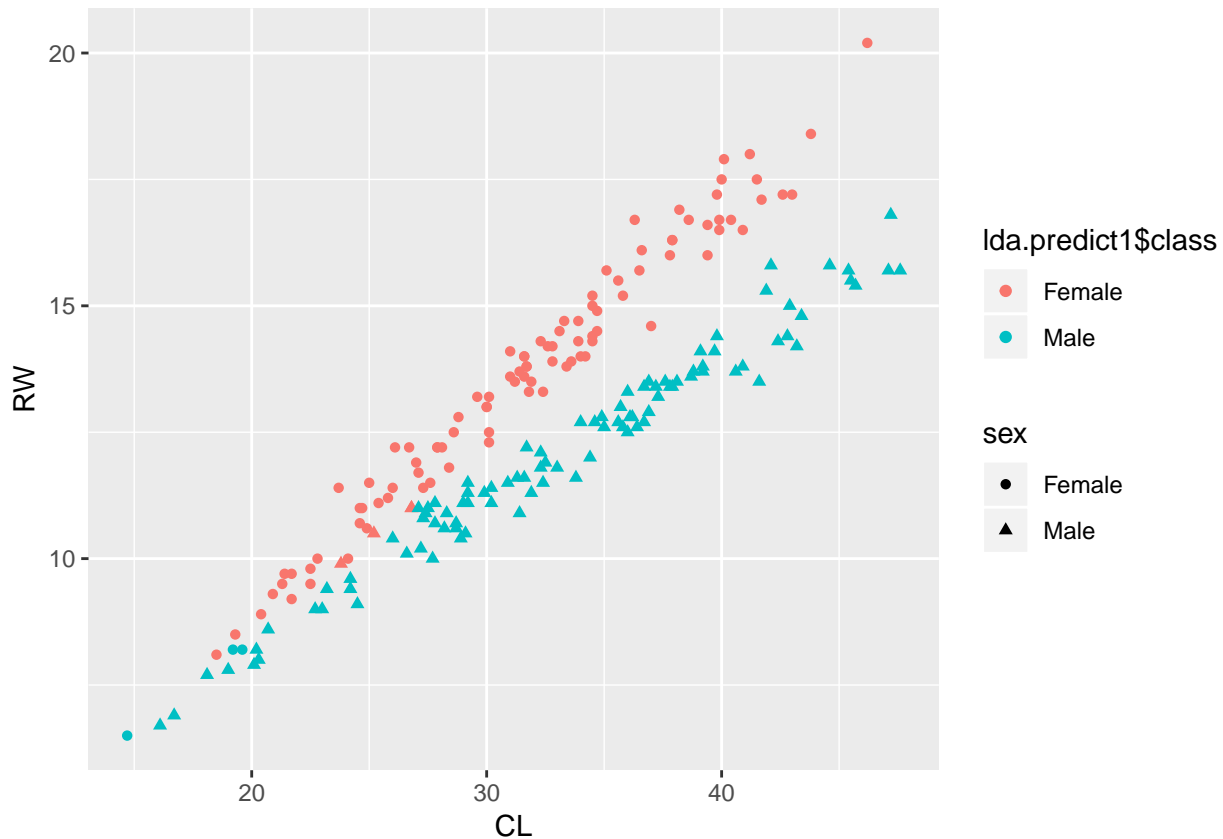*8 December 2018*

## Assignment 1.1

```
library("ggplot2")
library("MASS")
data <- data.frame(read.csv("australian-crabs.csv"))
ggplot(data = data, mapping = aes(CL, RW, color = sex)) + geom_point()
```



A plot of carapace length versus rear width where the observations are colored by sex. Looking at the graph the data seems reasonably easy to classify by linear discriminant analysis. Because there seems to be a line between the two sexes.

## Assignment 1.2

```
lda.model1 <- lda(sex ~ CL + RW, data = data)
lda.predict1 <- predict(lda.model1, data)
ggplot.0.5 <- ggplot(data = data, mapping = aes(CL, RW, color = lda.predict1$class, shape = sex )) + ge
ggplot.0.5
```
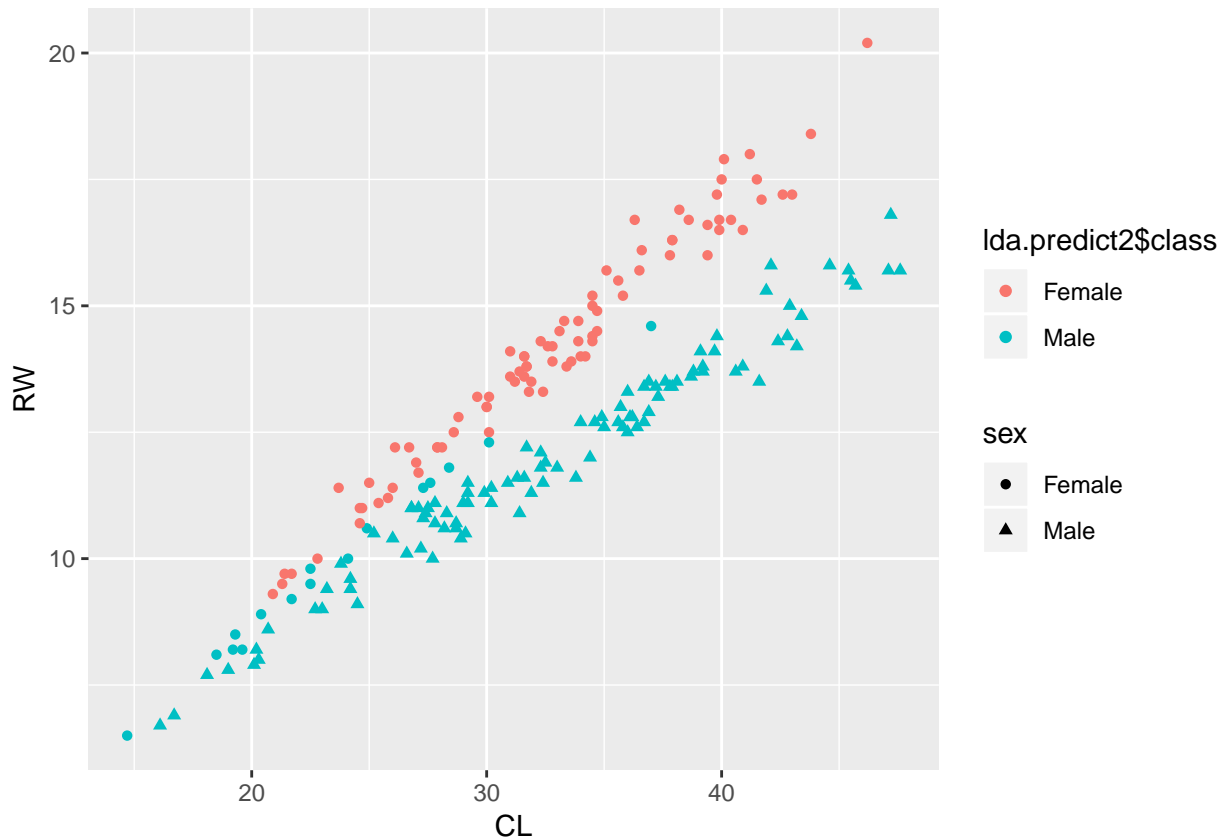
```
mcr.0.5 <- mean(lda.predict1$class != data$sex)
```

The missclassification rate for the linear discriminant analysis is **0.035**. This is pretty reasonable considering we saw on the original graph that there was one area with a bit of an overlap. If this is good enough for actual use is hard to say, it mostly depends on how much we would lose on an incorrect classification.

## Assignment 1.3

```
lda.model2 <- lda(sex ~ CL + RW, data = data, prior = c(Female = 0.1, Male = 0.9))
lda.predict2 <- predict(lda.model2, data)
ggplot(data = data, mapping = aes(CL, RW, color = lda.predict2$class, shape = sex )) + geom_point()
```

The number of males increased since we are assuming a wheighted distribution. Especially the areas containing both types of observations are now classified as only males instead of both.

```
mcr.0.9 <- mean(lda.predict2$class != data$sex)
```

The missclassification rate for the weighted linear discriminant analysis is **0.08**.

## Assignment 1.4

```
glm.model <- glm(as.factor(sex) ~ CL + RW, family = binomial, data = data)
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
glm.predict <- predict(glm.model, data, type = 'response')
mcr.glm <- mean(glm.predict != data$sex)
```
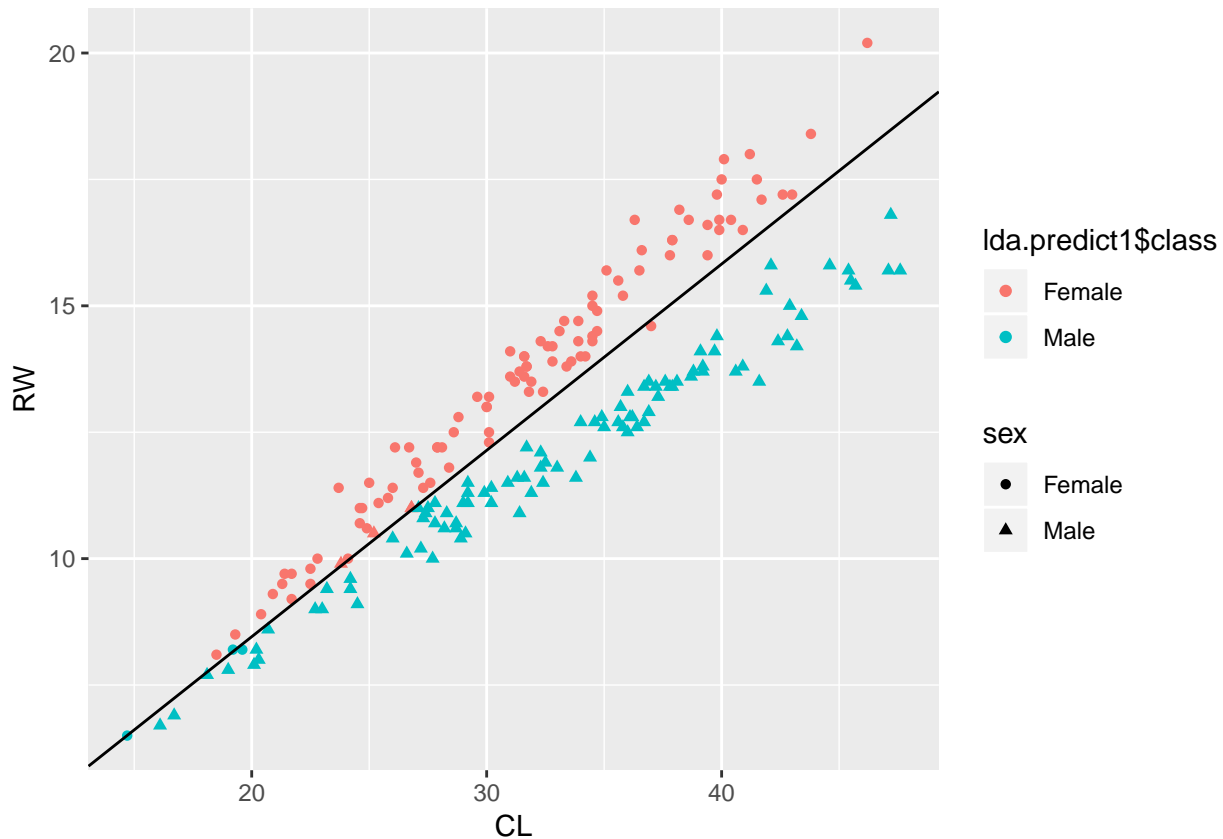
The missclassification rate is **0.035** which is the same as the original linear discriminant analysis.

```
glm.predict.0.5 <- ifelse(glm.predict > 0.5, "Male", "Female")
glm.slope <- coef(glm.model)[2]/(-coef(glm.model)[3])
glm.intercept <- coef(glm.model)[1]/(-coef(glm.model)[3])
ggplot.0.5 + geom_abline(slope = glm.slope, intercept = glm.intercept)
```

The decision line is drawn in the graph.

## Assignment 2

Splitting the data into partitions

```r
library("e1071")
library("MASS")
library("tree")
library("ggplot2")
setwd("~/TDDE01/lab2")
data.credit <- data.frame(read.csv("creditscoring.csv"))
n=dim(data.credit)[1]
set.seed(12345)
id=sample(1:n, floor(n*0.5))
train=data.credit[id,]
id1=setdiff(1:n, id)
set.seed(12345)
id2=sample(id1, floor(n*0.25))
valid=data.credit[id2,]
id3=setdiff(id1,id2)
test=data.credit[id3,]
```

The models used in calculating the following confusion matrices.

```r
tree.fit.dev <- tree(good_bad ~., data = train, split = "deviance")
tree.fit.gini <- tree(good_bad ~., data = train, split = "gini")
```

```
pred.dev.train <- predict(tree.fit.dev, newdata = train, type="class")
mean(train$good_bad != pred.dev.train)
```

## [1] 0.212

Missclassification rate for deviance on train data.

```
pred.dev.test <- predict(tree.fit.dev, newdata = test, type="class")
mean(test$good_bad != pred.dev.test)
```

## [1] 0.268

Missclassification rate for deviance on test data.

```
pred.gini.train <- predict(tree.fit.gini, newdata = train, type="class")
mean(train$good_bad != pred.gini.train)
```
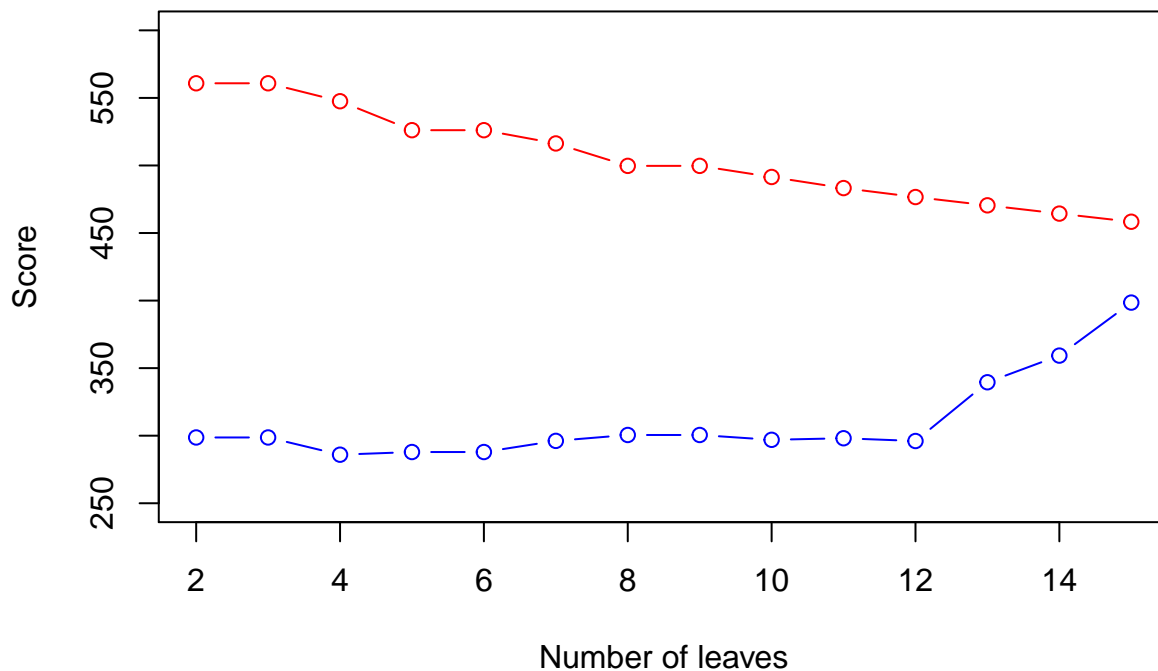
## [1] 0.238

Missclassification rate for gini on train data.

```
pred.gini.test <- predict(tree.fit.gini, newdata = test, type="class")
mean(test$good_bad != pred.gini.test)
```

## [1] 0.372

Missclassification rate for gini on test data.

The confusion matrix is the best for deviance compared to gini based on the number of correct predictions for the test data.
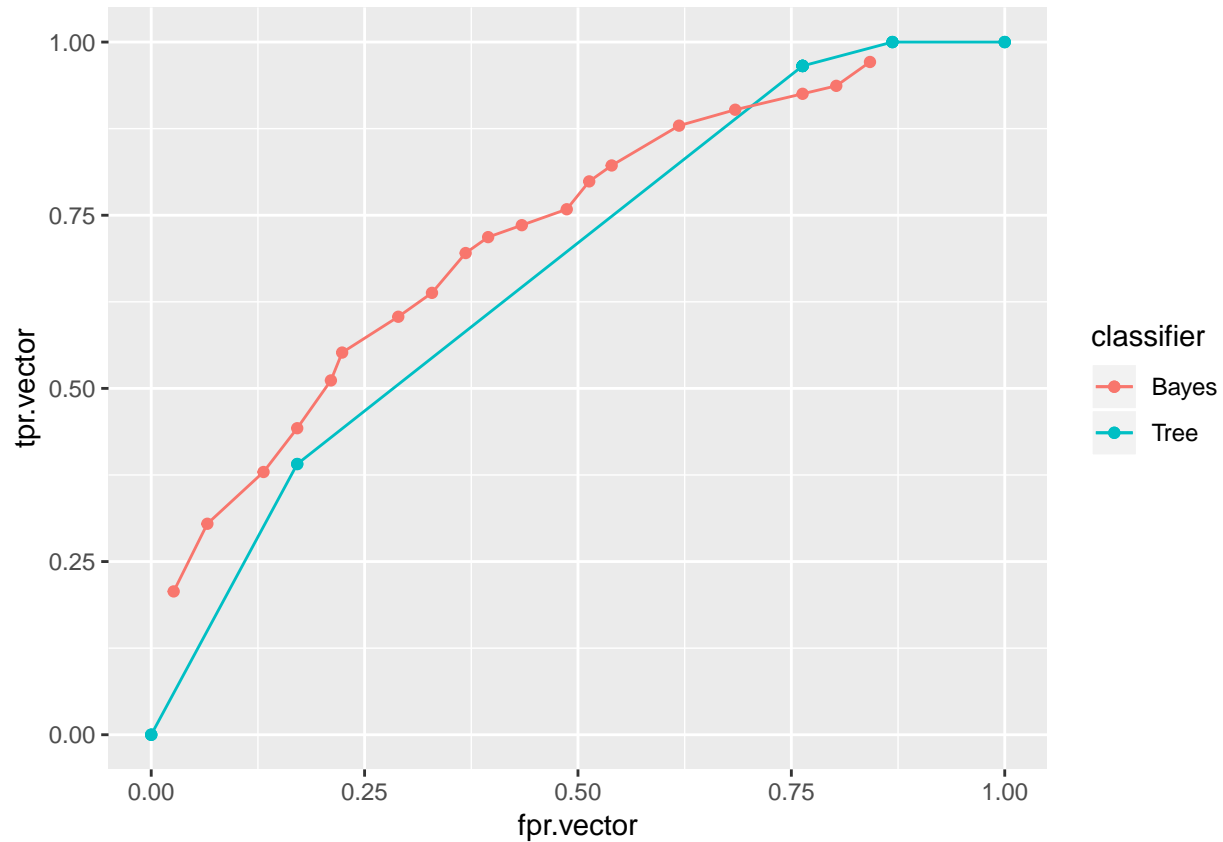
## Assignment 2.3



Looking at the graph we see a minimum value for 4.

The optimal depth of the tree is three which can be seen in the graph.

```
##        predicted
## actual bad good
##    bad   23   54
##    good  12  161
```

Confusion matrix for the validation data for the tree data

```
##
## Classification tree:
## snip.tree(tree = fit, nodes = c(5L, 3L, 9L))
## Variables actually used in tree construction:
## [1] "savings"  "duration" "history"
## Number of terminal nodes:  4
## Residual mean deviance:  1.117 = 547.5 / 490
## Misclassification error rate: 0.251 = 124 / 494
```

As seen above the variables used in the tree are "savings", "duration" and "history".

```
## [1] 0.264
```

Missclassification rate is **0.264**.

## Assignment 2.4

```
##        predicted
## actual bad good
##    bad   95   98
##    good  52  255
```

```
## [1] 0.3
```

Missclassification rate for bayes on train data.

```
##        predicted
## actual bad good
##    bad   46   49
##    good  30  125
```

```
## [1] 0.316
```

Missclassification rate for bayes on test data.

The tree prediction is a bit better than the bayesian prediction.



Graph of radius of convergence

```
##       predict
## actual bad good
##   bad   71    5
##   good 122   52
```

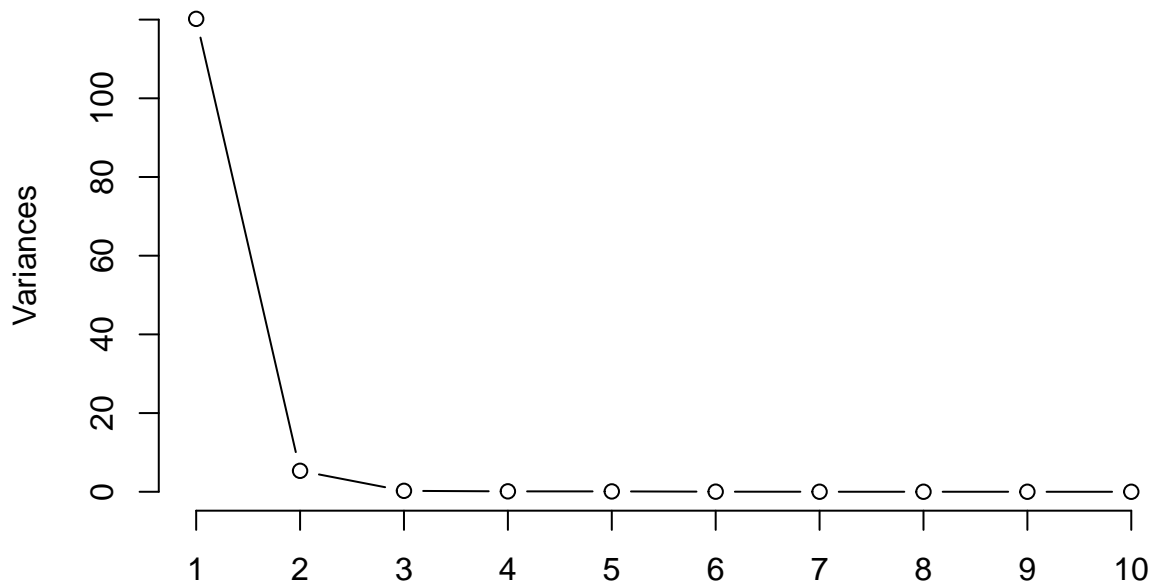Confusion matrix for bayes with loss matrix on test data.

```
##       predict
## actual bad good
##   bad  137   10
##   good 263   90
```

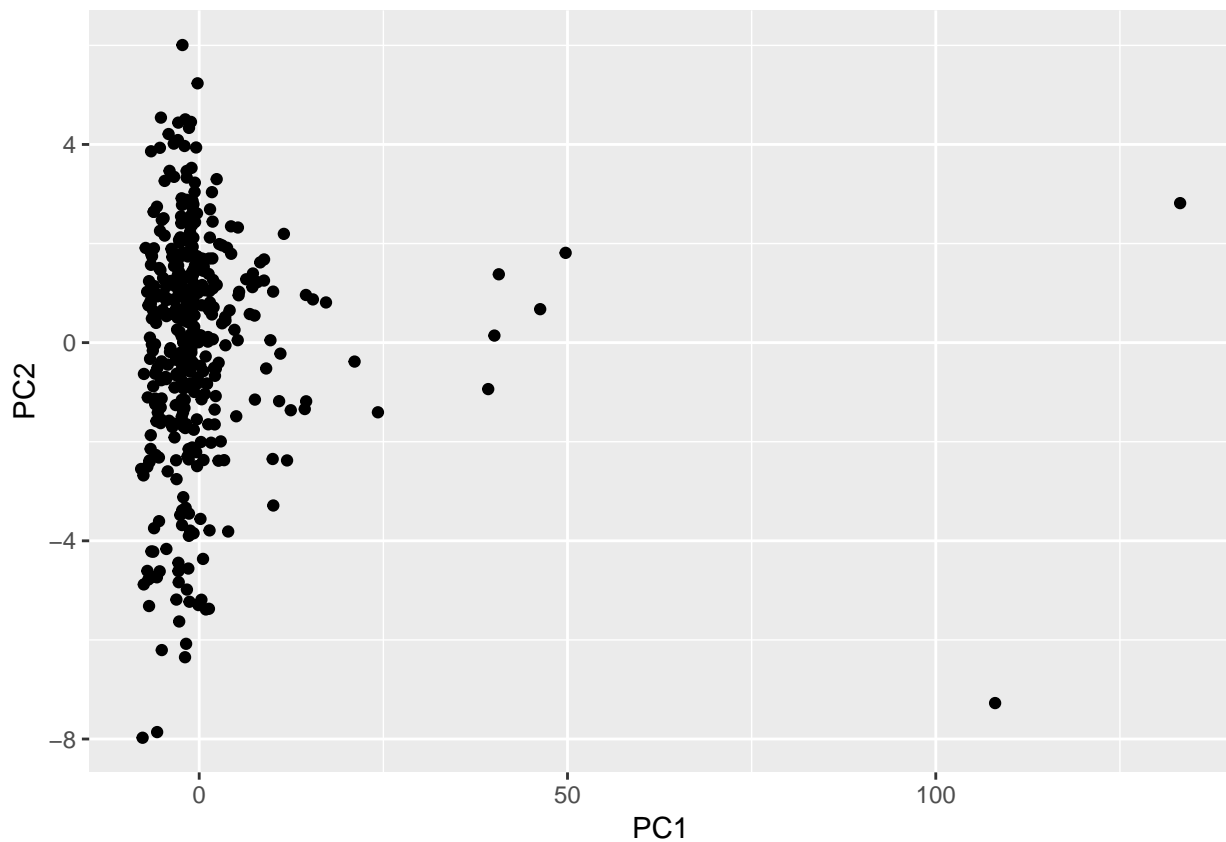Confusion matrix for bayes with loss matrix on train data.

It's a lot more expensive to have a false positive than it is to have a false negative according to a loss matrix. Meaning that we will be more cautious with the "good" classification. This means that we should have a much higher frequency of data being classified as "bad".

**Assignment 4.1**

## Principal component variance dependency



By looking at summary(pca.fit), can see that PC1 and PC2 cumulatively explains 99% of the variance. The summary(pca.fit) is not printed here since it's output is very large.
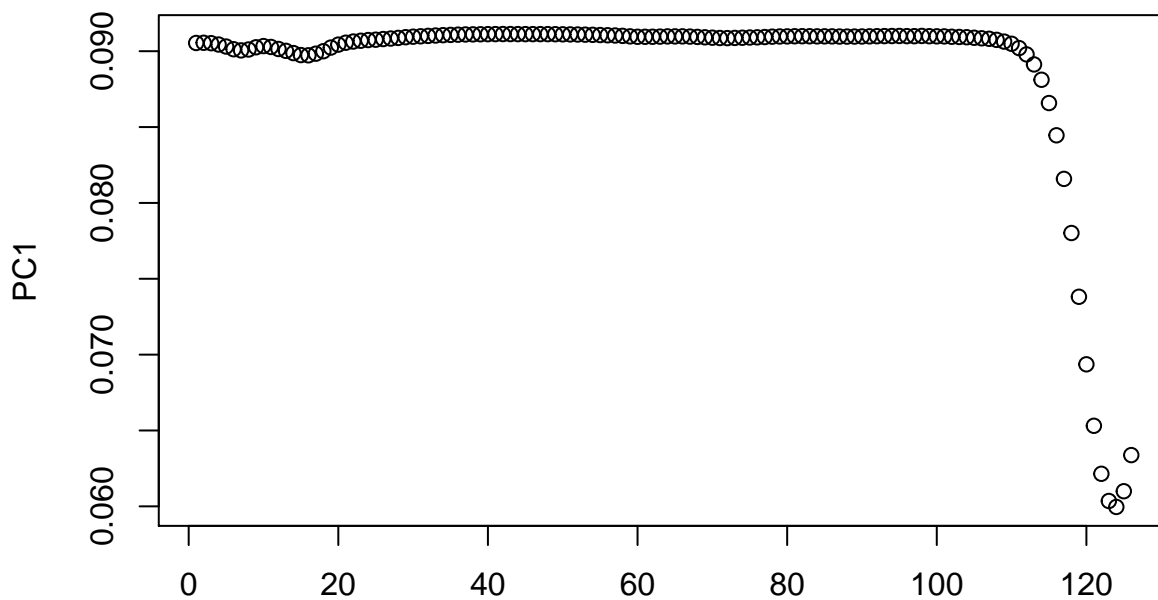


Plot of the scores in the in the (PC1, PC2) coordinates. There seems to be two fuels that differ greatly from
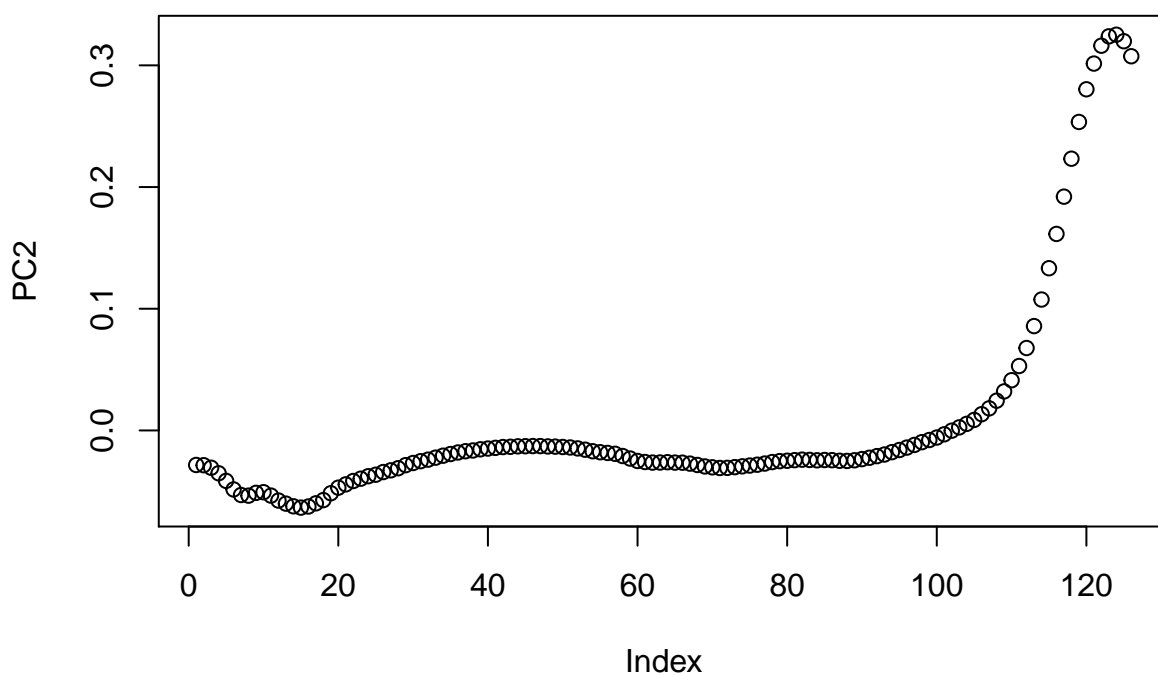
the other, in that they have a much higher coefficent for PC1, while a lot of the other observations seem to lie close to zero.

**Assignment 4.2**

## Traceplot



Index

## Traceplot



Index

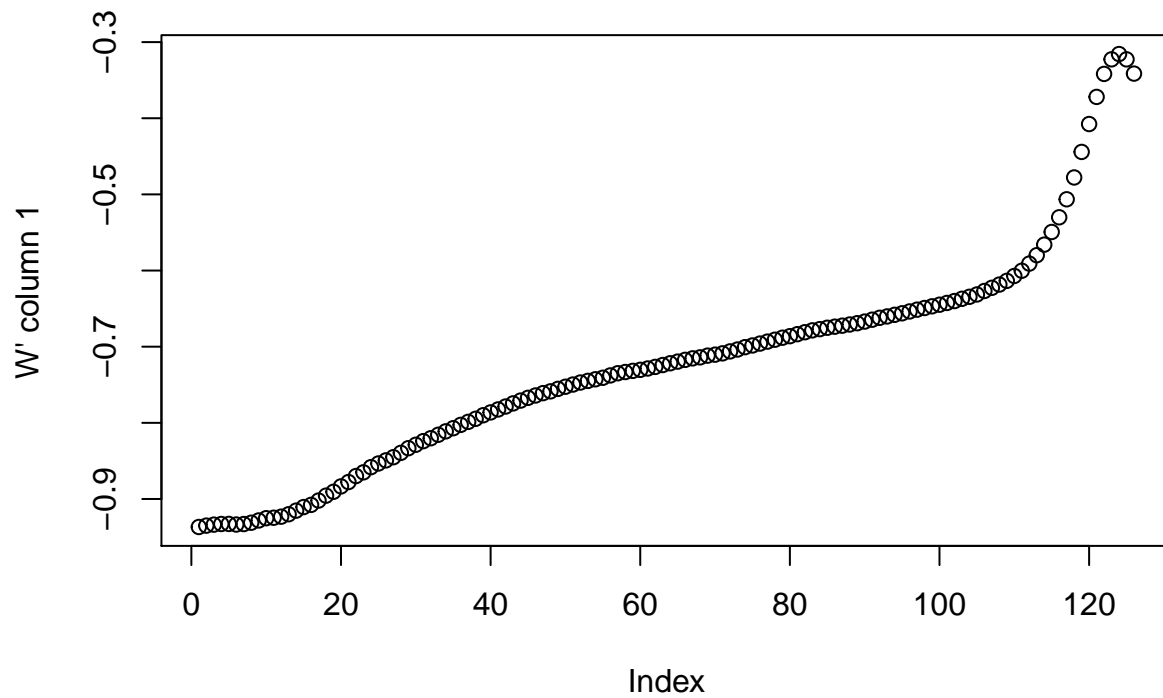**Assignment 4.3**

```
## Centering
## Whitening
## Symmetric FastICA using logcosh approx. to neg-entropy function
## Iteration 1 tol = 0.01930239
## Iteration 2 tol = 0.01303959
## Iteration 3 tol = 0.002393582
## Iteration 4 tol = 0.0006708454
## Iteration 5 tol = 0.0001661602
## Iteration 6 tol = 3.521604e-05
```

**Traceplot**

# Traceplot'