

# Compulsory exercise 2: Group 24

TMA4268 Statistical Learning V2019

Silje Anfindsen and Clara Panchaud

25 March, 2020

## Problem 1

a)

Let's find the ridge regression estimator. Remember that  $\hat{\beta}_{Ridge}$  minimizes  $RSS + \lambda \sum_{j=1}^p \beta_j^2$ . Let's rewrite this in matrix notation.

$$\begin{aligned} \min_{\beta} \{ (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta \} &= \text{develop the expression} \\ \min_{\beta} \{ y^T y - 2\beta^T X^T y + \beta^T X^T X \beta + \lambda \beta^T \beta \} & \quad \text{take the derivative with respect to beta and set equal to 0} \\ -2X^T y + 2X^T X \beta + 2\lambda \beta &= 0 \\ (X^T X + 2\lambda I) \beta &= X^T y \\ \beta &= (X^T X + \lambda I)^{-1} X^T y \end{aligned}$$

Therefore the estimator is  $\hat{\beta}_{Ridge} = (X^T X + \lambda I)^{-1} X^T y$ .

b)

To find the expected value and the variance-covariance matrix of  $\hat{\beta}_{Ridge}$  we need to remember the distribution of  $y$ ,  $y \sim N(X\beta, \sigma^2 I)$ . Therefore we get the expected value:

$$E(\hat{\beta}_{Ridge}) = E((X^T X + \lambda I)^{-1} X^T y) = (X^T X + \lambda I)^{-1} X^T E(y) = (X^T X + \lambda I)^{-1} X^T X \beta$$

and the variance-covariance matrix:

$$\begin{aligned} Var(\hat{\beta}_{Ridge}) &= Var((X^T X + \lambda I)^{-1} X^T y) = \text{by property of the variance} \\ (X^T X + \lambda I)^{-1} X^T Var(y) (X^T X + \lambda I)^{-1} X^T &= \text{develop the expression} \\ \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} & \end{aligned}$$

c)

TRUE, FALSE, FALSE, TRUE

d)

```
library(ISLR)
library(leaps)
library(glmnet)
```

We want to work with the *College* data. First we split it into a training and a testing set.

```
set.seed(1)

#make training and testing set
train.ind = sample(1:nrow(College), 0.5 * nrow(College))
college.train = College[train.ind, ]
college.test = College[-train.ind, ]

#the structure of the data
str(College)
```

```
## 'data.frame': 777 obs. of 18 variables:
## $ Private : Factor w/ 2 levels "No","Yes": 2 2 2 2 2 2 2 2 2 2 ...
## $ Apps : num 1660 2186 1428 417 193 ...
## $ Accept : num 1232 1924 1097 349 146 ...
## $ Enroll : num 721 512 336 137 55 158 103 489 227 172 ...
## $ Top10perc : num 23 16 22 60 16 38 17 37 30 21 ...
## $ Top25perc : num 52 29 50 89 44 62 45 68 63 44 ...
## $ F.Undergrad: num 2885 2683 1036 510 249 ...
## $ P.Undergrad: num 537 1227 99 63 869 ...
## $ Outstate : num 7440 12280 11250 12960 7560 ...
## $ Room.Board : num 3300 6450 3750 5450 4120 ...
## $ Books : num 450 750 400 450 800 500 500 450 300 660 ...
## $ Personal : num 2200 1500 1165 875 1500 ...
## $ PhD : num 70 29 53 92 76 67 90 89 79 40 ...
## $ Terminal : num 78 30 66 97 72 73 93 100 84 41 ...
## $ S.F.Ratio : num 18.1 12.2 12.9 7.7 11.9 9.4 11.5 13.7 11.3 11.5 ...
## $ perc.alumni: num 12 16 30 37 2 11 26 37 23 15 ...
## $ Expend : num 7041 10527 8735 19016 10922 ...
## $ Grad.Rate : num 60 56 54 59 15 55 63 73 80 52 ...
```

Now we will apply forward selection, using *Outstate* as a response. We have 18 variables including the response so we will obtain a model including up to 17 variables.

```
nb_predictors<-17
forward<-regsubsets(Outstate~.,college.train,nvmax=17,method="forward")
sum<-summary(forward)
```

In Figure 1 we can look at the RSS and the adjusted  $R^2$  in order to pick the number of variables that gives the optimal result. Remember that if the difference is not very significant we would rather pick the simplest model. It seems like 5 variables would be good here.

```
par(mfrow=c(1,2))
plot(sum$rss,xlab="Number of Variables",ylab="RSS",type="l")
plot(sum$adjr2,xlab="Number of Variables",ylab="Adjusted RSq",type="l")
```

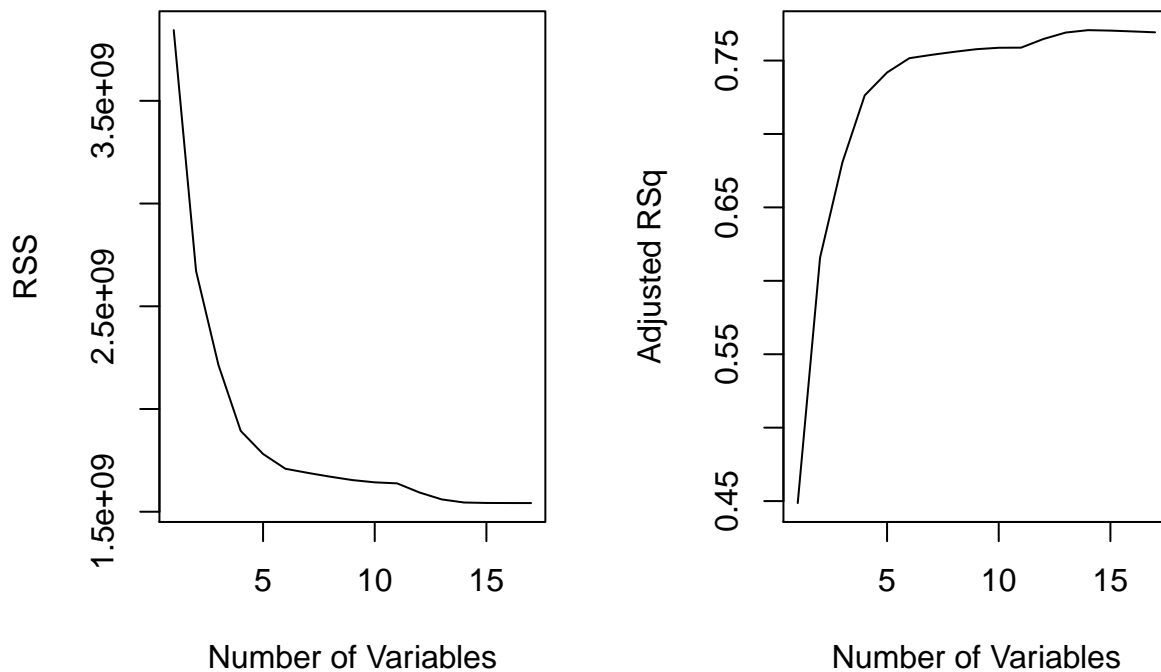


Figure 1: Comparison of models with different number of variables.

Below are the chosen variables when we decide to include 5 variables in the reduced model.

```
nb_selected_pred<-5
variables<-names( coef( forward,id=nb_selected_pred ) )
variables
```

```
## [1] "(Intercept)" "PrivateYes" "Room.Board" "perc.alumni" "Expend"
## [6] "Grad.Rate"
```

We will now find the reduced model as well as the MSE (mean squared error) on the test set.

```
#fit the reduced model
reduced.model<-lm(Outstate~Private+Room.Board+Grad.Rate+perc.alumni+Expend, data =college.train)
summary(reduced.model)
```

```
##
## Call:
```

```
## lm(formula = Outstate ~ Private + Room.Board + Grad.Rate + perc.alumni +
##      Expend, data = college.train)
##
## Residuals:
##      Min        1Q    Median        3Q        Max
## -7293.1 -1537.5  -159.9   1286.7   9254.8
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2.711e+03  5.155e+02  -5.259 2.41e-07 ***
## PrivateYes   2.250e+03  2.787e+02   8.075 8.92e-15 ***
## Room.Board   1.241e+00  1.205e-01  10.296 < 2e-16 ***
## Grad.Rate    3.855e+01  7.850e+00   4.910 1.35e-06 ***
## perc.alumni  6.446e+01  1.113e+01   5.792 1.45e-08 ***
## Expend       2.182e-01  2.317e-02   9.417 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2159 on 382 degrees of freedom
## Multiple R-squared:  0.7452, Adjusted R-squared:  0.7419
## F-statistic: 223.4 on 5 and 382 DF,  p-value: < 2.2e-16
```

The reduced model is

$$\begin{aligned} \text{Outstate} = & -2711.4329907 + 2250.1100562\text{Private} + 1.2410466\text{Room.Board} \\ & + 38.5491289\text{Grad.Rate} + 64.4580901\text{perc.alumni} + 0.218216\text{Expend}, \end{aligned}$$

```
#find test MSE
p<-predict(reduced.model,newdata=college.test)
mse_fwd <- mean(((college.test$Outstate)-p)^2)
mse_fwd
```

```
## [1] 4112680
```

The test MSE is

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 = 4.1126804 \times 10^6$$

e)

We will now select a model for the same dataset as in (d) but this time with the Lasso method. Again, we use both a training and testing set for the data.

```
#Make a x matrix and y vector for both the training and testing set
x_train<-model.matrix(Outstate~.,college.train)[-1]
y_train<-college.train$Outstate
x_test<-model.matrix(Outstate~.,college.test)[-1]
y_test<-college.test$Outstate
```

In order to select the best value for the tuning parameter  $\lambda$  we will use cross validation.

```
set.seed(5555)

#perform the Lasso method and choose the best model using CV
lasso.mod = glmnet(x_train,y_train,alpha=1) #lasso method on train set
cv.lasso = cv.glmnet(x_train,y_train,alpha=1) #CV on train set
lambda.best = cv.lasso$lambda.min #select best lambda
lambda.best
```

```
## [1] 5.093201
```

```
#find the test MSE
predictions<-predict(lasso.mod,s=lambda.best,newx=x_test)
mse_lasso <- mean((predictions-y_test)^2) #test MSE
mse_lasso
```

```
## [1] 3717020
```

%%Check if log scale!! From cross validation we can observe that the optimal tuning parameter is  $\lambda = 5.0932009$  as this is the parameter that minimizes the MSE for the training set.

The test MSE is now  $3.7170197 \times 10^6$ , which is lower than what we found for the reduced model using forward selection in d).

The lasso yields sparse models which involves only a subset of variables. Lasso performs variable selection by forcing some of the coefficient estimates to be exactly zero. The selected variables that was not put to zero are displayed below.

```
c<-coef(lasso.mod,s=lambda.best,exact=TRUE)
inds<-which(c!=0)
variables<-row.names(c)[inds]
variables
```

```
## [1] "(Intercept)" "PrivateYes" "Apps" "Accept" "Enroll"
## [6] "Top10perc" "Top25perc" "F.Undergrad" "P.Undergrad" "Room.Board"
## [11] "Books" "Personal" "PhD" "Terminal" "S.F.Ratio"
## [16] "perc.alumni" "Expend" "Grad.Rate"
```

## Problem 2

a)

FALSE, FALSE, TRUE, FALSE

b)

The basis functions for a cubic spline with knots at each quartile, of variable  $X$  are,

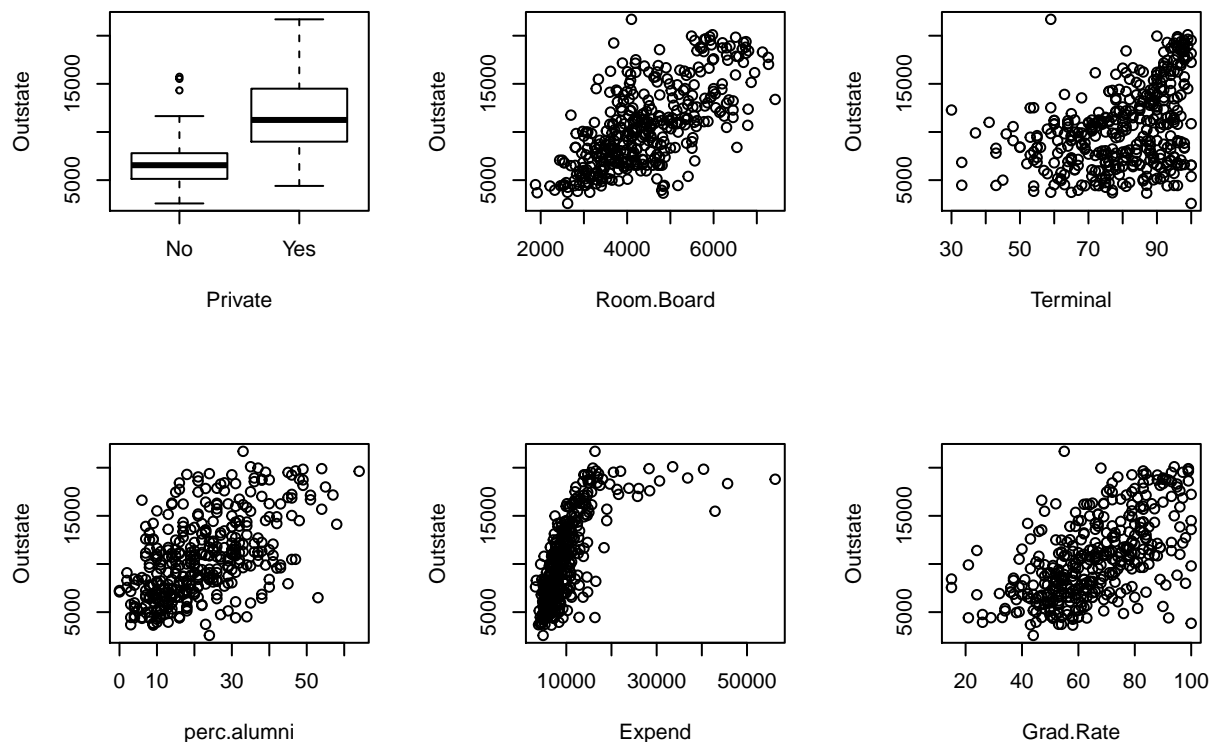
$$\begin{aligned} b_0(X) &= 1 & b_4(X) &= (X - q_1)_+^3 \\ b_1(X) &= x & b_5(x) &= (X - q_2)_+^3 \\ b_2(X) &= x^2 & b_6(X) &= (X - q_3)_+^3 \\ b_3(X) &= x^3 \end{aligned}$$

c)

We will now investigate the relationship between *Outstate* and the 6 of the predictors, *Private*, *Room.Board*, *Terminal*, *perc.alumni*, *Expend*, and *Grad.Rate*.

```
ds1 = college.train[c("Private", "Outstate")] #binary variable
ds2 = college.train[c("Room.Board", "Outstate")]
ds3 = college.train[c("Terminal", "Outstate")]
ds4 = college.train[c("perc.alumni", "Outstate")]
ds5 = college.train[c("Expend", "Outstate")]
ds6 = college.train[c("Grad.Rate", "Outstate")]

par(mfrow=c(2,3))
plot(ds1)
plot(ds2)
plot(ds3)
plot(ds4)
plot(ds5)
plot(ds6)
```



From each of the plots above we can conclude that at least *Terminal* and *Expend* seems to have a non-linear relationship with *Outstate*. These two variables therefore might benefit from a non-linear transformation. The others variables, *Room.Board*, *perc.alumni* and *Grad.Rate* seem to have a linear relationship with the response variable. The binary variable *Private* is presented through a boxplot. Generally the data seem to categorize quite well into these two classes of private and public universities where the trend is a higher out-of-state tuition for private universities, except for some outliers for public universities where the outcome

is very high and therefore can seem to belong to a private universities. Anyway, we cannot transform a binary variable.

d)

We will now fit several polynomial regression models for *Outstate* with *Terminal* as the only covariate. Each polynomial will have a degree from  $d = 1, \dots, 10$ .

```
library(ggplot2)

#make a dataframe
ds = College[c("Terminal", "Outstate")]
n = nrow(ds)

# chosen degrees
deg = 1:10

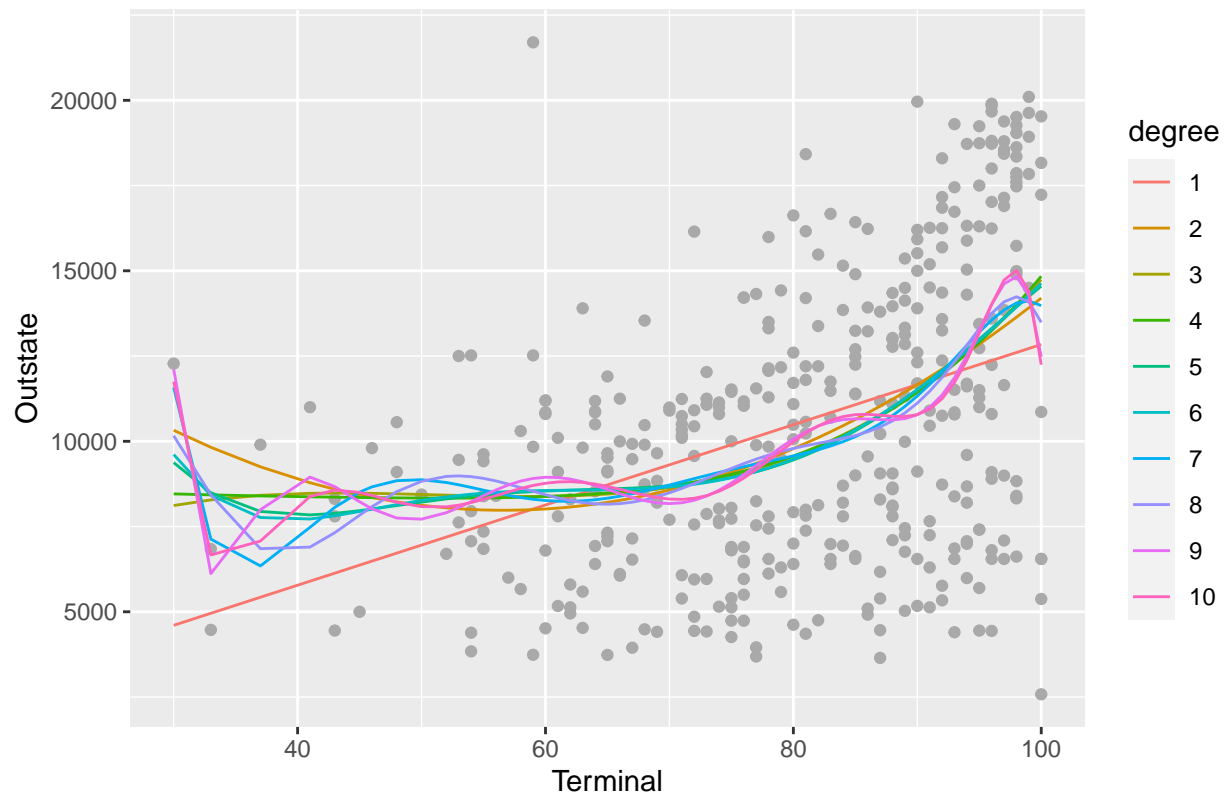
#now iterate over each degree d
dat = c() #make a empty variable to store predicted values for each degree
MSE_poly = c(rep(0,10)) #make a empty variable to store MSE for each degree
for (d in deg) {
  # fit model with this degree
  mod = lm(Outstate ~ poly(Terminal, d), ds[train.ind, ])

  #dataframe for Terminal and Outstate showing result for each degree over all samples
  dat = rbind(dat, data.frame(Terminal = ds[train.ind, 1], Outstate = mod$fit,
                              degree = as.factor(rep(d,length(mod$fit)))))

  # training MSE
  MSE_poly[d] = mean((predict(mod, ds[-train.ind, ]) - ds[-train.ind, 2])^2)
}

# plot fitted values for different degrees
ggplot(data = ds[train.ind, ], aes(x = Terminal, y = Outstate)) +
  geom_point(color = "darkgrey") + labs(title = "Polynomial regression")+
  geom_line(data = dat, aes(x = Terminal, y = Outstate, color = degree))
```

## Polynomial regression



We will now choose a suitable smoothing spline model to predict *Outstate* as a function of *Expend* and plot the fitted function.

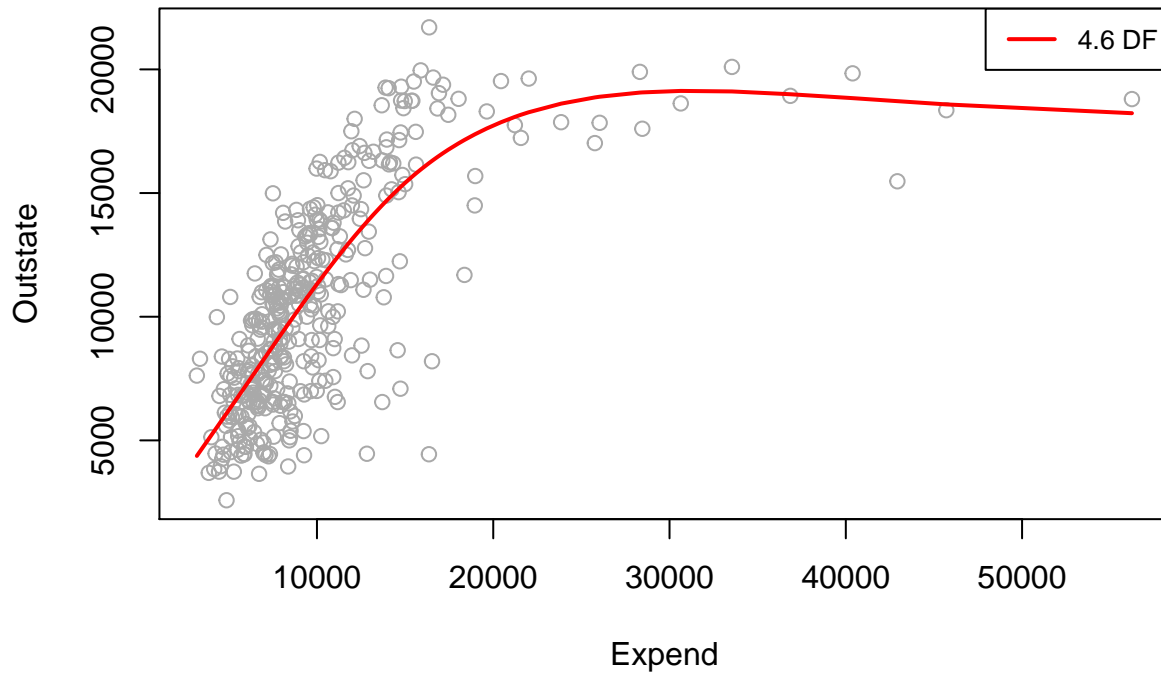
```
library(splines)

#plot training set for Expend as only covariate
plot(college.train$Expend, college.train$Outstate, col = "darkgrey", main="Smoothing spline", xlab="Expend")
#perform CV in order to find optimal number of df
fit = smooth.spline(college.train$Expend, college.train$Outstate, cv=TRUE)
df <- fit$df #choose df from CV
l <- fit$lambda

#add fitted function from smoothing spline
lines(fit, col="red", lwd=2)
legend("topright", legend=c("4.6 DF"), col="red", lty=1, lwd=2, cex=.8)
```



## Smoothing spline

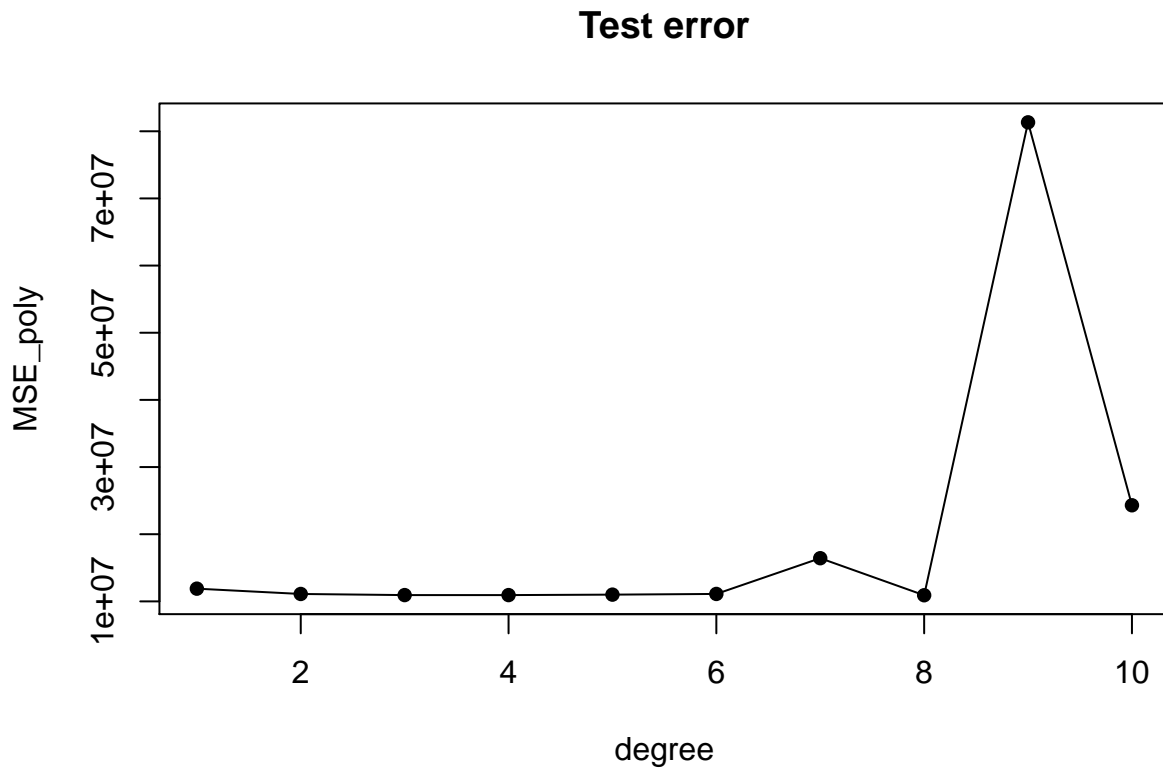


```
#training MSE
pred = predict(fit, newdata=college.train)
MSE_spline = mean( (college.train$Outstate - pred$y)^2 )
```

In order to choose a suitable model we did cross validation. The optimal number of degrees of freedom is  $df = 4.660711$ , which gives the smoothing parameter  $\lambda = 0.0075385$

We will now calculate the training MSE for the polynomial regression models and the smoothing spline model. For the polynomial regression models we find it easiest to look at the MSE by presenting a plot with MSE for each degree.

```
#MSE for polynomial regression models (1-10)
plot(1:10,MSE_poly, type = "o", pch = 16, xlab = "degree", main = "Test error")
```



```
min(MSE_poly)
```

```
## [1] 10914136
```

```
#MSE for smoothing spline  
MSE_spline
```

```
## [1] 31072818
```

The training MSE for the two methods are (Read about this!)  
(discuss is it is as expected!)

### Problem 3

a)

FALSE, FALSE, TRUE, FALSE

b)

To predict *Outstate* we will choose

PROS and CONS! One known disadvantage of regression trees are that they usually have a high variance, and are not very robust. The idea of bagging is to grow many trees to fix those two problems. Bagging is a special case of another approach called random forests. In random forests we allow only certain number of predictors  $m$  to be selected at each node, which decreases the correlation between each tree.

The tuning parameter in random forests is the number of predictors you are able to choose between at each split. It is recommended to use  $m = p/3$  variables when building a regression trees. Recall that we have 17 parameters in our dataset. We will therefore pick  $m = 5$ .

```
library(randomForest)

set.seed(1)
bag.college <- randomForest(Outstate ~ ., data=college.train, mtry=5, ntrees=500)
yhat.college <- predict(bag.college, newdata=college.test)
MSE_rf <- mean((yhat.college - college.test$Outstate)^2)
MSE_rf

## [1] 2607985
```

The test MSE for the random forests is  $MSE_{rf} = 2.6079851 \times 10^6$ .

c)

We will now compare the test MSEs among the methods used on the data set *College* so far. That is: the two linear model selection methods, forward selection and Lasso method, non-linear methods, polynomial regression and smoothing splines and at last the tree-based method Random forests of regression trees.

```
MSE <- c(mse_fwd, mse_lasso, min(MSE_poly), MSE_spline, MSE_rf)
Method <- c("Forward selection", "Lasso", "Polynomial regression", "Smoothing spline",
            "Random Forest")
df <- data.frame(Method, MSE)
kable(df)
```

Method	MSE
Forward selection	4112680
Lasso	3717020
Polynomial regression	10914136
Smoothing spline	31072818
Random Forest	2607985

The method performing best in terms of prediction error is the random forest of regression trees. But if the aim is to develop a interpretable model Lasso is the best choice.

## Problem 4

Start by loading the data of *diabetes* from a population of women.

```
id <- "1Fv6xwKLSZHldRAC1MrcK2mzd0Ynbgv0E" # google file ID
d.diabetes <- dget(sprintf("https://docs.google.com/uc?id=%s&export=download",
  id))
d.train = d.diabetes$ctrain
d.test = d.diabetes$ctest
```

a)

In order to answer on the Multiple choice we will have a look at the training data.

```
library(GGally)
library(gridExtra)

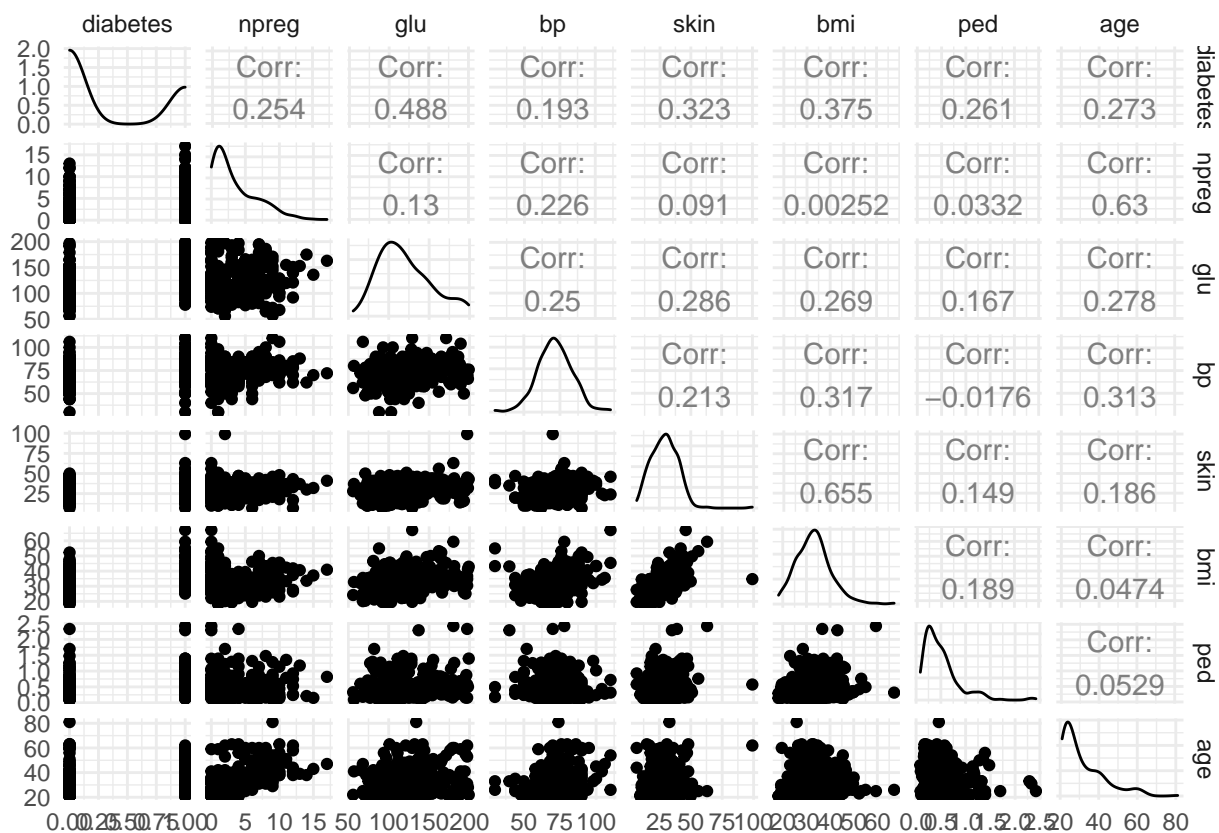
max(d.train$npreg) #max nr of pregnancies
```

```
## [1] 17
```

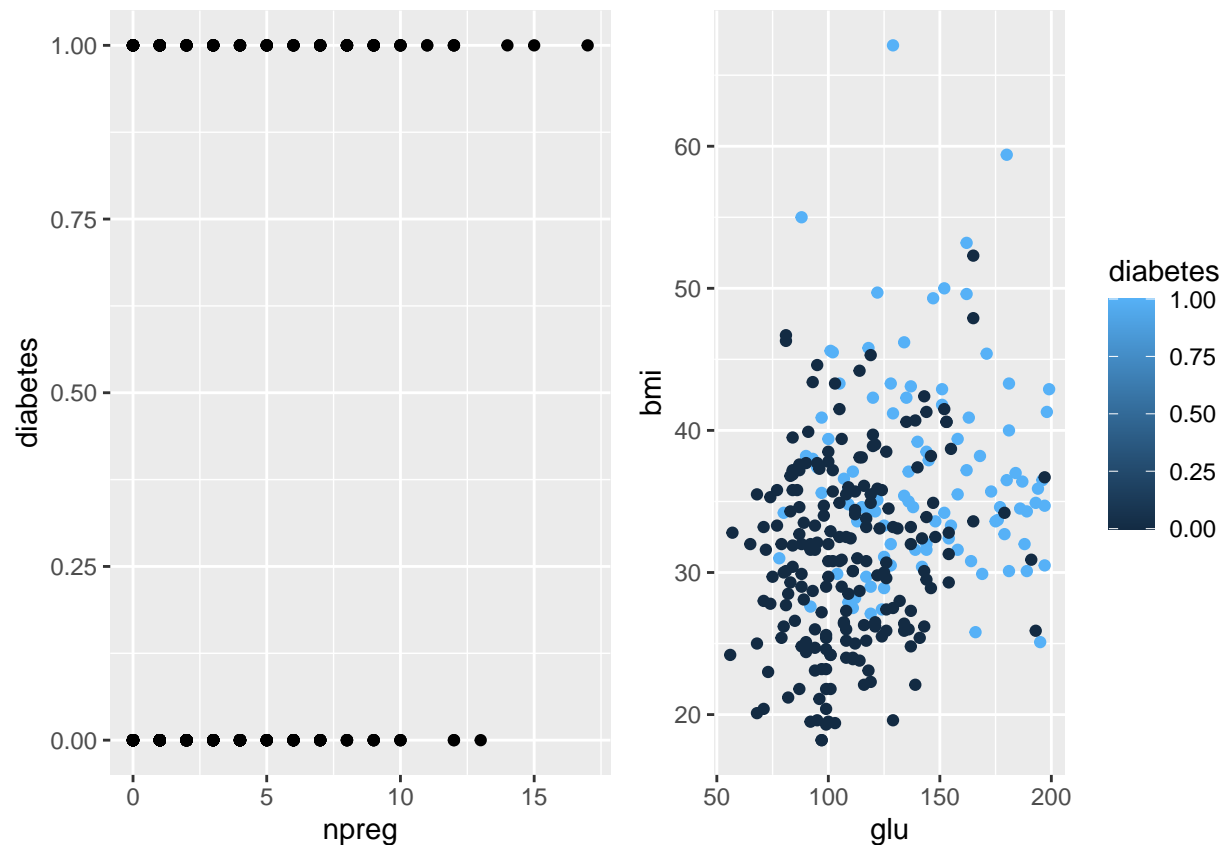
```
head(d.train) #overview of data
```

```
##      diabetes npreg glu bp skin  bmi   ped age
## 339         1     7 109 80   31 35.9 1.127 43
## 270         1     7 152 88   44 50.0 0.337 36
## 210         1     9 119 80   35 29.0 0.263 29
## 117         1     9 112 82   32 34.2 0.260 36
## 390         1     8 151 78   32 42.9 0.516 36
## 217         1     2  90 68   42 38.2 0.503 27
```

```
#plots
ggpairs(d.train) + theme_minimal() #look at correlation between variables
```



```
plot2 <- ggplot(d.train, aes(x=npreg, y=diabetes))+geom_point()
plot3 <- ggplot(d.train, aes(x=glu, y=bmi, color=diabetes))+geom_point()
grid.arrange(plot2, plot3, ncol=2, nrow = 1)
```



TRUE, TRUE, TRUE, FALSE (not sure about last)

b)

We will now fit a support vector classifier with linear boundary and a support vector machine with radial boundary to find good functions that predict the diabetes status of a patient.

```
library(e1071)

set.seed(10111)

d.train$diabetes <- as.factor(d.train$diabetes)
d.test$diabetes <- as.factor(d.test$diabetes)

#Fit a support vector classifier (linear boundary)
svm.linear = svm(diabetes~.,
                 data = d.train,
                 kernel = 'linear')

#fit a support vector machine (radial boundary)
svm.radial = svm(diabetes~.,
                 data = d.train,
                 kernel = 'radial')
```

```

#CV to find the best parameters for each model
cv.linear <- tune(svm, diabetes ~ ., data=d.train, kernel = "linear",
                 ranges = list(cost = c(.001, .01, .1, 1, 10, 100) ) )

cv.radial <- tune(svm, diabetes ~ ., data = d.train, kernel = "radial",
                 ranges = list(cost = c(0.1,1,10,100,1000), gamma = c(0.5,1,2,3,4) ))

#fit new models with the optimized parameters from CV
bestmod.linear = cv.linear$best.model
bestmod.radial = cv.radial$best.model

# Predict the response for the test set
pred.lin = predict(bestmod.linear, newdata = d.test)
pred.rad = predict(bestmod.radial, newdata = d.test)
# Confusion tables (0: no diabetes, 1: diabetes)
#for SVC (linear)
table(Prediction = pred.lin, Truth = d.test$diabetes)

```

```

##           Truth
## Prediction  0   1
##           0 137  35
##           1  18  42

```

```

#for SVM (radial)
table(Prediction = pred.rad, Truth = d.test$diabetes)

```

```

##           Truth
## Prediction  0   1
##           0 133  38
##           1  22  39

```

From the two confusion tables above we can find the missclassification rates. The left rate belongs to the support vector classifier (linear boundary), and the right rate to the support vector machine (radial boundary).

$$\frac{35 + 18}{137 + 35 + 18 + 42} = 0.2284 \quad \frac{38 + 22}{133 + 38 + 22 + 39} = 0.2586$$

We know that the missclassification error rate is a good performance measure for classification rules. So by looking at these rates the support vector classifier with a linear boundary seems to make the best classification rule for the prescence of diabetes for our data. This is because the number of missclassifications is remarkable lower for this classifier.

c)

We will now compare the performance of the support vector classifier and support vector machine with a new method: logistic regression. As we have observed above, our data seems to performe well when being classified with a linear decision boundary. Logistic regression models the probability that the response belongs to one of the two classes, producing a linear decision boundary. Therefore this method seems to be a good choice for our data.

```

#fit a logistic regression model using training data
glm.fit <- glm(diabetes ~ .,data=d.train, family="binomial")

#predict the response using testing data
glm.probs <- predict(glm.fit, newdata=d.test, type="response")

#sort the probabilities for whether the observations are > or > than 0.5
glm.pred = rep("0",length(d.test$diabetes)) #create vector of nr. of elements = dataset
glm.pred[glm.probs>0.5]="1"

#confusion table
table(Prediction = glm.pred, Truth = d.test$diabetes)

```

```

##           Truth
## Prediction    0    1
##           0 136  32
##           1   19  45

```

The missclassification rate for logistic regression is

$$\frac{32 + 19}{136 + 45 + 32 + 19} = 0.2198$$

We notice that the missclassification rate for logistic regression is even lower than for the two SVMs. This indicates that logistic regression could be a even better choice in order to predict the prescense of diabetes based on our data.

Both logistic regression and support vector machines solve classification problems. In logistic regression all observations contribute weighted to the decision boundary, while for SVMs, only the support vectors (the points closest to the decision boundary) contribute to the margin. A consequence of this is that LR is more sensitive to outliers than SVM. On the other hand, LR can be easily updated as new data comes in.

For classes that are well separated SVM tend to perform better than LR, while in more overlapping regimes we usually prefer LR. We also know that logistic regression produces probabilistic values, while SVM produces binary values. This can be an advantage if we want an estimation rather than just the resulting class for each observation.

d)

FALSE, FALSE, ??, TRUE? (really not sure about any of these answers!)

e)

## Problem 5

```

id <- "1VfVCQvWt121UN39NXZ4aR9Dmsbj-p90U" # google file ID
GeneData <- read.csv(sprintf("https://docs.google.com/uc?id=%s&export=download", id), header = F)
colnames(GeneData)[1:20] = paste(rep("H", 20), c(1:20), sep = "")
colnames(GeneData)[21:40] = paste(rep("D", 20), c(1:20), sep = "")
row.names(GeneData) = paste(rep("G", 1000), c(1:1000), sep = "")
GeneData<-t(GeneData)

```



a)

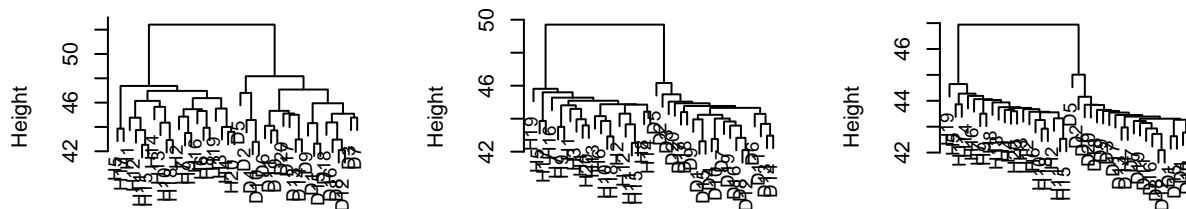
We start by performing hierarchical clustering on the dataset. We try the complete, single and average linkage for both the Euclidian and correlation-based distance. We transposed the data in order to cluster the tissues into the two patient groups.

```
hc.eucl.complete=hclust(dist(GeneData,method="euclidian"), method="complete")
hc.eucl.average=hclust(dist(GeneData,method="euclidian"), method="average")
hc.eucl.single=hclust(dist(GeneData,method="euclidian"), method="single")
```

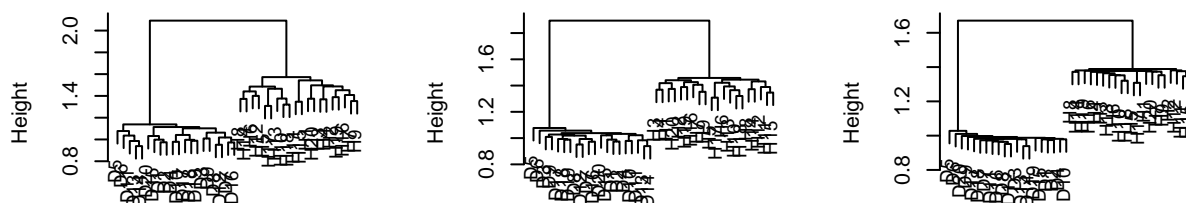
```
correlation<-dist(cor(t(GeneData)))
hc.corr.complete=hclust(correlation, method="complete")
hc.corr.average=hclust(correlation, method="average")
hc.corr.single=hclust(correlation, method="single")
```

```
par(mfrow=c(2,3))
plot(hc.eucl.complete,main="Complete Linkage, Euclidian distance", xlab="", sub="", cex=.9)
plot(hc.eucl.average, main="Average Linkage, Euclidian distance", xlab="", sub="", cex=.9)
plot(hc.eucl.single, main="Single Linkage, Euclidian distance", xlab="", sub="", cex=.9)
plot(hc.corr.complete,main="Complete Linkage, correlation-based distance", xlab="", sub="", cex=.9)
plot(hc.corr.average, main="Average Linkage, correlation-based distance", xlab="", sub="", cex=.9)
plot(hc.corr.single, main="Single Linkage, correlation-based distance", xlab="", sub="", cex=.9)
```

Complete Linkage, Euclidian dist:Average Linkage, Euclidian dista Single Linkage, Euclidian distar



Complete Linkage, correlation-based Average Linkage, correlation-based Single Linkage, correlation-based d



The dendrograms seem to recognize that there are two different groups. ## b)

We now use the dendrograms to cluster the tissues into two groups.

```
cutree(hc.eucl.complete, 2)
```

```
## H1 H2 H3 H4 H5 H6 H7 H8 H9 H10 H11 H12 H13 H14 H15 H16 H17 H18 H19 H20
## 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## D1 D2 D3 D4 D5 D6 D7 D8 D9 D10 D11 D12 D13 D14 D15 D16 D17 D18 D19 D20
## 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
```

```
cutree(hc.eucl.average, 2)
```

```
## H1 H2 H3 H4 H5 H6 H7 H8 H9 H10 H11 H12 H13 H14 H15 H16 H17 H18 H19 H20
## 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## D1 D2 D3 D4 D5 D6 D7 D8 D9 D10 D11 D12 D13 D14 D15 D16 D17 D18 D19 D20
## 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
```

```
cutree(hc.eucl.single, 2)
```

```
## H1 H2 H3 H4 H5 H6 H7 H8 H9 H10 H11 H12 H13 H14 H15 H16 H17 H18 H19 H20
## 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## D1 D2 D3 D4 D5 D6 D7 D8 D9 D10 D11 D12 D13 D14 D15 D16 D17 D18 D19 D20
## 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
```

```
cutree(hc.corr.complete, 2)
```

```
## H1 H2 H3 H4 H5 H6 H7 H8 H9 H10 H11 H12 H13 H14 H15 H16 H17 H18 H19 H20
## 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## D1 D2 D3 D4 D5 D6 D7 D8 D9 D10 D11 D12 D13 D14 D15 D16 D17 D18 D19 D20
## 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
```

```
cutree(hc.corr.average, 2)
```

```
## H1 H2 H3 H4 H5 H6 H7 H8 H9 H10 H11 H12 H13 H14 H15 H16 H17 H18 H19 H20
## 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## D1 D2 D3 D4 D5 D6 D7 D8 D9 D10 D11 D12 D13 D14 D15 D16 D17 D18 D19 D20
## 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
```

```
cutree(hc.corr.single, 2)
```

```
## H1 H2 H3 H4 H5 H6 H7 H8 H9 H10 H11 H12 H13 H14 H15 H16 H17 H18 H19 H20
## 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## D1 D2 D3 D4 D5 D6 D7 D8 D9 D10 D11 D12 D13 D14 D15 D16 D17 D18 D19 D20
## 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
```

We know that the first 20 tissues come from healthy patients and the last 20 from diseased ones. Therefore it seems like all linkage and distance measures perform perfectly.

**ASK PEOPLE???**

c)

The elements of the vector  $\phi$  are called loadings and define a direction in the feature space along which the data varies the most. The data is a  $n \times p$  matrix  $X$ .

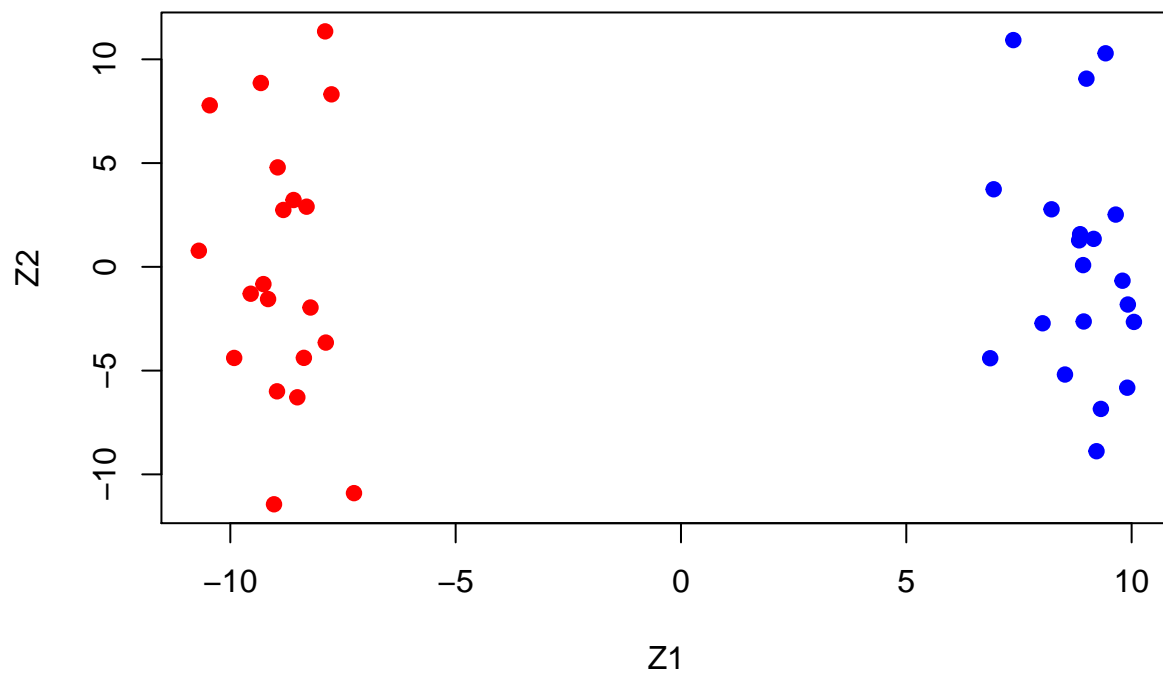
For the first principal component we want  $Z_1 = \phi_1^T X$  subject to  $\|\phi_1\|_2 = 1$ . We want  $Z_1$  to have the highest possible variance  $V(Z_1) = \phi_1^T \Sigma \phi_1$ , where  $\Sigma$  is the covariance matrix of  $X$ . The first principal component scores are then the column eigenvector corresponding to the largest eigenvalue of  $\Sigma$ .

d)

```
color<-c(rep(1,20),rep(2,20))
```

How to color based on tissues group of patients???

```
pca_gene=prcomp(GeneData, scale=TRUE)
plot(pca_gene$x[,1:2], col=c("red","blue")[color],pch=19,xlab="Z1",ylab="Z2")
```



Now we calculate the proportion of variance explained (PVE) by the 5 first components.

```
pve=100*pca_gene$sdev^2/sum(pca_gene$sdev^2)
cumsum(pve)[5]
```

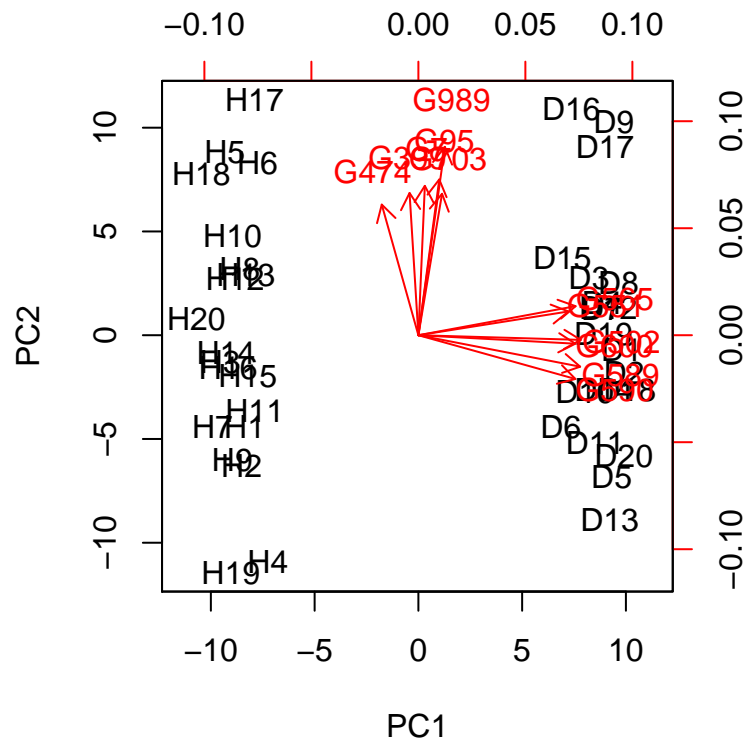
```
## [1] 21.09659
```

About 21 percent of the variance is explained by first 5 PCs.

e)

Use your results from PCA to find which genes that vary the most across the two groups. ??

```
nyt_loading = pca_gene$rotation[, 1:2]
informative_loadings = rbind(head(nyt_loading[order(nyt_loading[, 1], decreasing = TRUE),
]), head(nyt_loading[order(nyt_loading[, 2], decreasing =
biplot(x = pca_gene$x[, 1:2], y = informative_loadings, scale = 0)
```



## f)

```
km.out = kmeans(GeneData,2, nstart = 20)
km.out$cluster
```

```
## H1 H2 H3 H4 H5 H6 H7 H8 H9 H10 H11 H12 H13 H14 H15 H16 H17 H18 H19 H20
## 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## D1 D2 D3 D4 D5 D6 D7 D8 D9 D10 D11 D12 D13 D14 D15 D16 D17 D18 D19 D20
## 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
```

## References

James, G., D. Witten, T. Hastie, and R. Tibshirani. 2013. An Introduction to Statistical Learning with Applications in R. New York: Springer.