

Compulsory exercise 2: Group 24

TMA4268 Statistical Learning V2019

Silje Anfindsen and Clara Panchaud

02 April, 2020

Problem 1

a)

Let's find the ridge regression estimator. Remember that $\hat{\beta}_{Ridge}$ minimizes $RSS + \lambda \sum_{j=1}^p \beta_j^2$. Here we assume that all covariates and the response have been mean-centered, so that we do not include β_0 . 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)

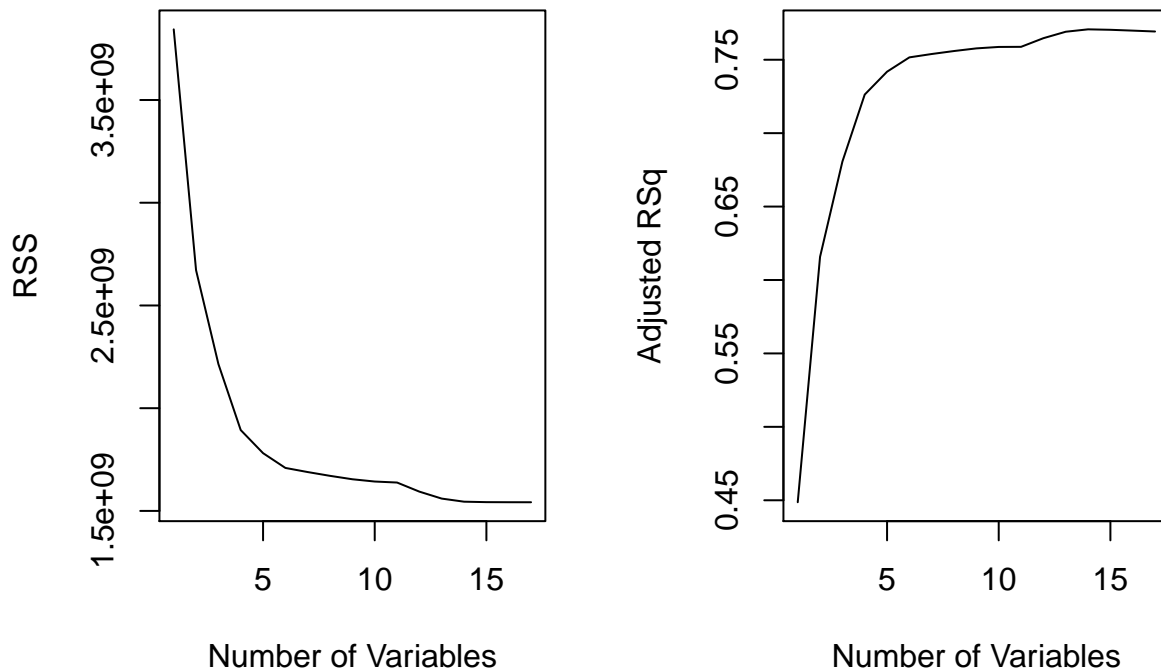
#make training and testing set
train.ind = sample(1:nrow(College), 0.5 * nrow(College))
college.train = College[train.ind, ]
college.test = College[-train.ind, ]
#str(College)
```

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 the figure below 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")
```



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)
```

The reduced model is

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

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

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 λ 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

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

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 3.7170197×10^6 , 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 q_1, q_2, q_3 , of variable X are,

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

where the power basis functions are defined as

$$b(X, q_j) = (x_i - q_j)_+^3 \begin{cases} (X - q_j)^3 & \text{if } X > q_j \\ 0 & \text{otherwise} \end{cases}$$

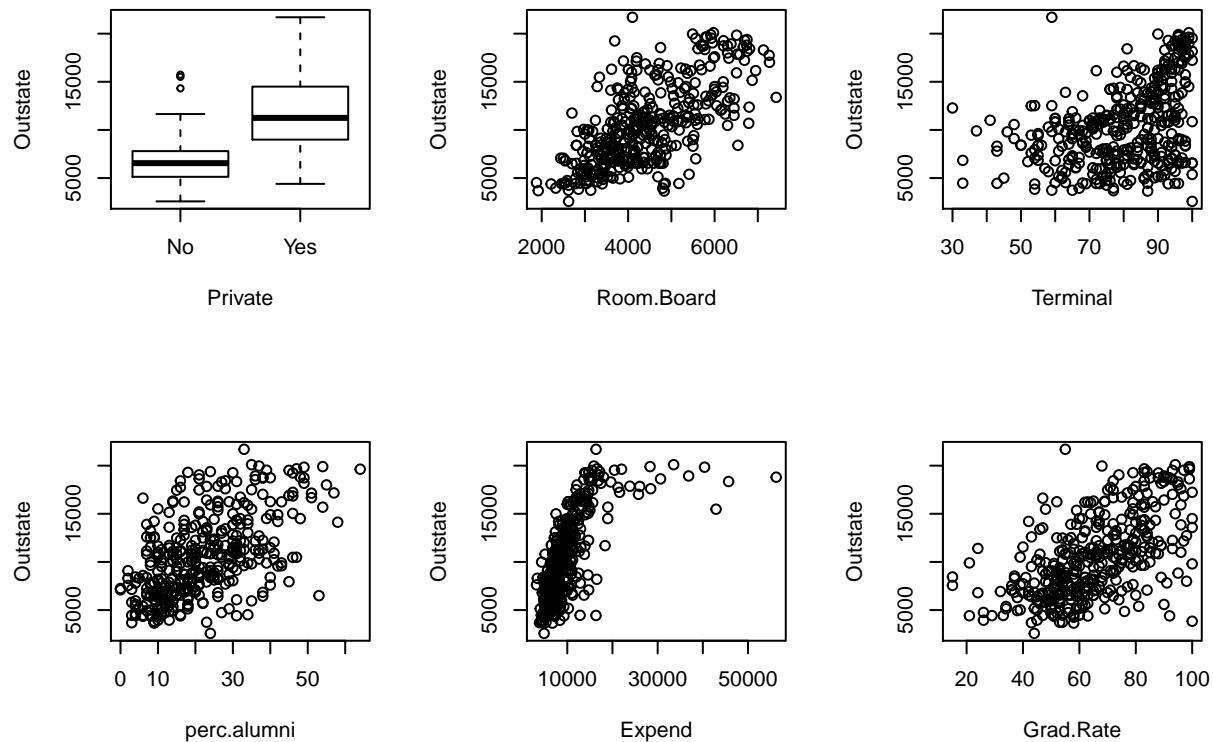
This amounts to estimate 7 regression coefficients β_k for $k = 0, \dots, 7$ for the intercept and predictors.

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. We cannot transform a binary variable.

d)

(i)

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_train_poly = c(rep(0,10)) #make a empty variable to store MSE for each degree
```

```

MSE_test_poly = c(rep(0,10))

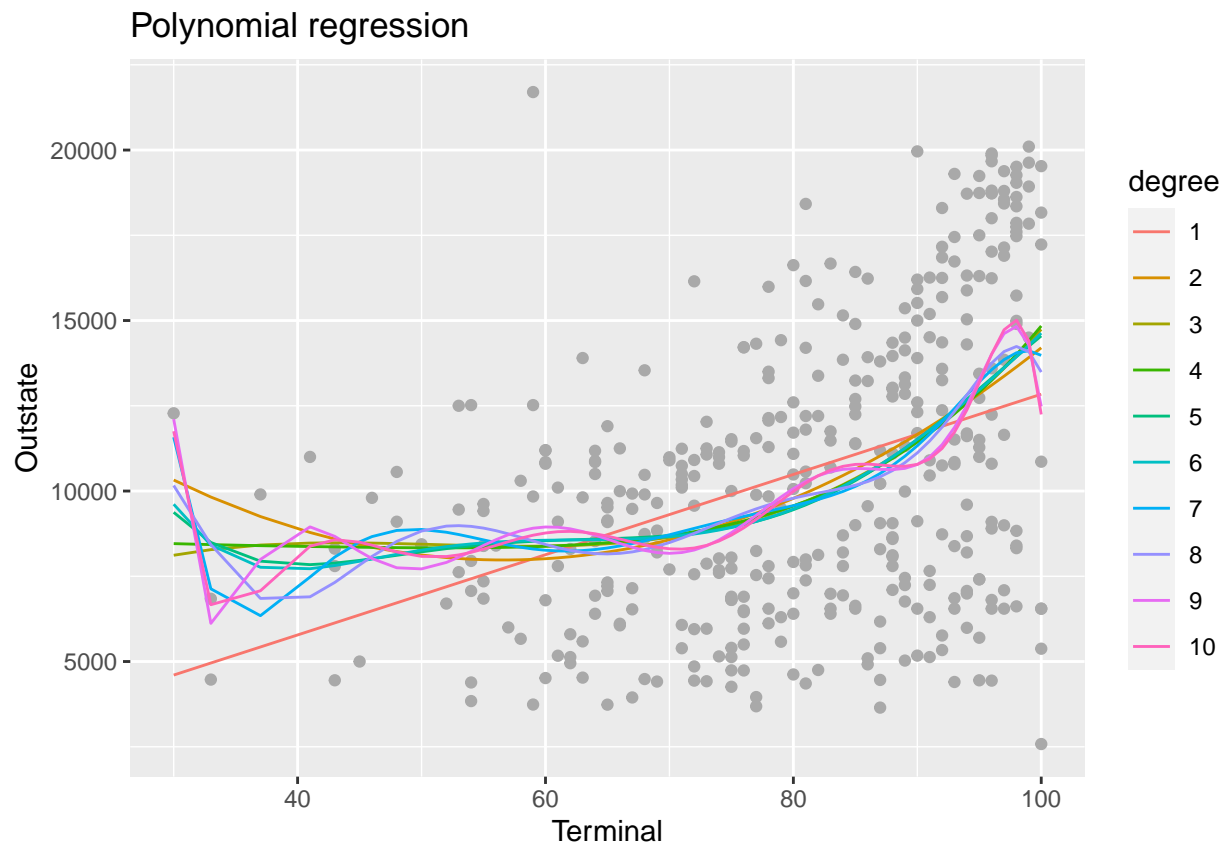
for (d in deg) {
  # fit model with this degree
  mod = lm(Outstate ~ poly(Terminal, d), data= college.train)

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

  # training MSE
  MSE_train_poly[d] = mean((predict(mod, newdata=college.train) - college.train$Outstate)^2)
  MSE_test_poly[d] = mean((predict(mod, newdata= college.test) - college.test$Outstate)^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))

```



(ii)

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

```

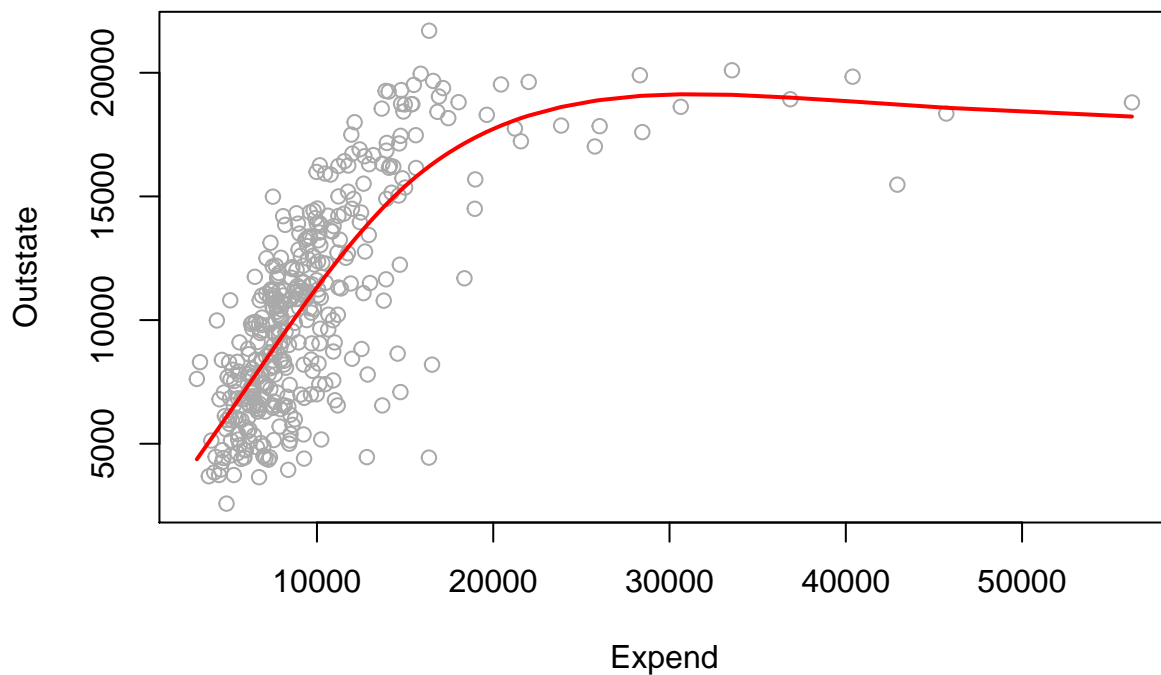
library(splines)
set.seed(1)
#perform CV in order to find optimal number of df
fit = smooth.spline(college.train$Expend, college.train$Outstate, cv=TRUE)

#plot training set for Expend as only covariate
plot(college.train$Expend, college.train$Outstate, col = "darkgrey",
     main=paste("Smoothing spline, df = ", round(fit$df,3)), xlab="Expend", ylab="Outstate")

#add fitted function from smoothing spline
lines(fit, col="red", lwd=2)

```

Smoothing spline, df = 4.661



```

#training MSE
pred = predict( fit, college.train$Expend)
MSE_spline_train <- mean((pred$y - college.train$Outstate )^2)

#test MSE
pred = predict(fit, college.test$Expend)
MSE_spline_test = mean( (college.test$Outstate - pred$y)^2 )

```

In order to choose the optimal number of degrees of freedom , which is $df = 4.660711$, we did cross validation.

(iii)

We will now report the training MSE for the polynomial regression models and the smoothing spline model. For the polynomial regression models we present the MSE for each degree of polynomial $d = 1, \dots, 10$,

```
#MSE for polynomial regression models (1-10)
Degree <- seq(1,10,by=1)
MSE = MSE_train_poly
df <- data.frame(Degree, MSE)
kable(df)
```

Degree	MSE
1	15075161
2	14330586
3	14249448
4	14247330
5	14231485
6	14230392
7	14153207
8	14097911
9	13841526
10	13822205

For the smoothing spline the training MSE is,

$$MSE = 6.8712814 \times 10^6$$

If we look back at the plots from c) we can see that when we plotted *Outstate* against *Expend* we had way less noise than with *Terminal*. Therefore we expected a smaller MSE for the model that uses *Expend*, that is the smoothing spline. We can confirm this as the measured training MSE for the smoothing spline model is smaller than each of the ten polynomial regression models. The plot from d) (ii) confirms this as it indicated a pretty good fit.

From the plot in d) (i) of the polynomial regression, it seems like all degrees of polynomials were somehow similar, at least we cannot see a huge difference from the plot. So even though we normally expect the MSE to decrease with higher degrees, this remark gives a reason to expect that the MSE not will decrease too much, which also seems to be the case when looking at the table that reports the training MSE for each degree of polynomial.

Problem 3

a)

FALSE, TRUE, TRUE, FALSE

b)

In order to predict *Outstate* we will choose the tree-based method: random forests. We know that regression trees usually suffer from high variance and non-robustness. By aggregating many decision trees and averaging

the resulting prediction, random forests solves these two problems. Another advantage with random forests is the fact that the method injects more randomness than for example bagging in order to avoid strong predictors dominating the decision trees. This is done by allowing only certain number of predictors m to be selected at each node, which decreases the correlation between each tree. Even though random forests is a flexible and accurate method, one big disadvantage is its high complexity as the method creates a lot of trees and therefore requires more computational resources than one simple regression 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)

# fit a model using random forests with a sufficiently large number of trees
rf.college <- randomForest(Outstate ~ ., data=college.train, mtry=5, ntree=500)
#predict the response using test data
yhat.college <- predict(rf.college, newdata=college.test)
#test MSE
MSE_rf <- mean((yhat.college - college.test$Outstate)^2)
```

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, and at the tree-based method random forests of regression trees.

```
MSE <- c(mse_fwd, mse_lasso, MSE_rf)
Method <- c("Forward selection", "Lasso",
            "Random Forest")
df <- data.frame(Method, MSE)
kable(df)
```

Method	MSE
Forward selection	4112680
Lasso	3717020
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 an 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)

```
#library(GGally)
#library(gridExtra)
#max(d.train$npreg)#max nr of pregnancies
#head(d.train) #overview of data
#ggpairs(d.train) + theme_minimal() #look at correlation between variables
#plot2 <- ggplot(d.train, aes(npreg)) + geom_histogram(binwidth = 2) + labs(title = "Histogram") + theme
#plot3 <- ggplot(d.train, aes(x=glu, y=bmi, color=diabetes))+geom_point()
#grid.arrange(plot2, plot3, ncol=2, nrow = 1)
```

TRUE, TRUE, TRUE, TRUE

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. We also observe that the number of false positive and false negative findings are both lower for this chosen model compared to the support vector machine with radial boundary.

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 p = 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 presense of diabetes based on our data.

In logistic regression all observations contribute 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. 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, TRUE

e)

We will now show the connection between hinge loss and logistic loss for the $y = -1, 1$ encoding. Let's find the loss function of the logistic regression model in this case. We can see that with this encoding we can write that $P(Y_i = y_i | \mathbf{x}_i) = g(y_i \beta^T \mathbf{x}_i)$ where $g(\mathbf{x}_i)$ is the logistic function. We will use the notation $f(\mathbf{x}_i) = \beta^T \mathbf{x}_i$. We verify that this corresponds to the model:

$$P(Y_i = 1 | \mathbf{x}_i) = \frac{1}{1 + e^{-f(\mathbf{x}_i)}}$$

$$P(Y_i = -1 | \mathbf{x}_i) = \frac{1}{1 + e^{f(\mathbf{x}_i)}} = 1 - \frac{1}{1 + e^{-f(\mathbf{x}_i)}} = 1 - P(Y_i = 1 | \mathbf{x}_i)$$

So indeed that is a correct way to write the model. We will now calculate the contribution off one observation to the deviance. Remember that the deviance is $-2 \log(\text{Likelihood})$, but we will ignore the 2 here.

$$-\log(g(y_i \beta^T \mathbf{x}_i)) =$$

$$-\log\left(\frac{1}{1 + e^{-y_i f(\mathbf{x}_i)}}\right) = \log(1 + e^{-y_i f(\mathbf{x}_i)})$$

We found the desired result.

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)

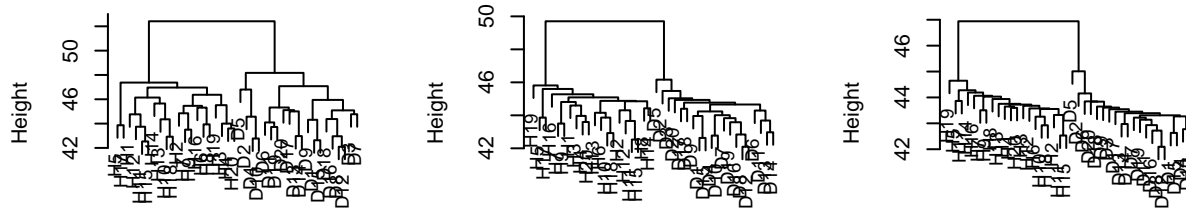
We start by performing hierarchical clustering on the dataset. We try the complete, single and average linkage for both the Euclidian and correlation-based distance. We transposed the data in order to cluster the tissues into the two patient groups.

```
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")
```

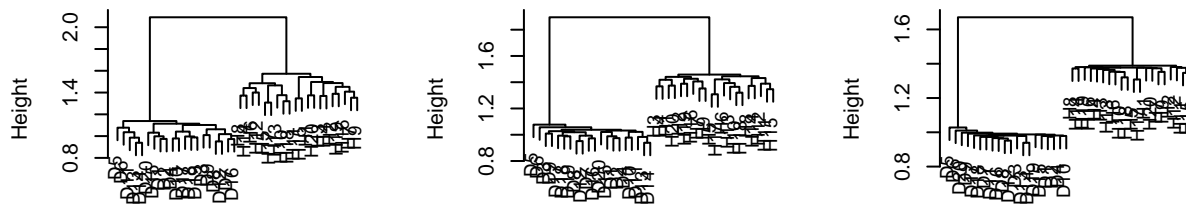
```
correlation<-dist(cor(t(GeneData)))
hc.corr.complete=hclust(correlation, method="complete")
hc.corr.average=hclust(correlation, method="average")
hc.corr.single=hclust(correlation, method="single")
```

```
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 dist: Single Linkage, Euclidian dist



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



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.

c)

The elements of the vector ϕ 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 .

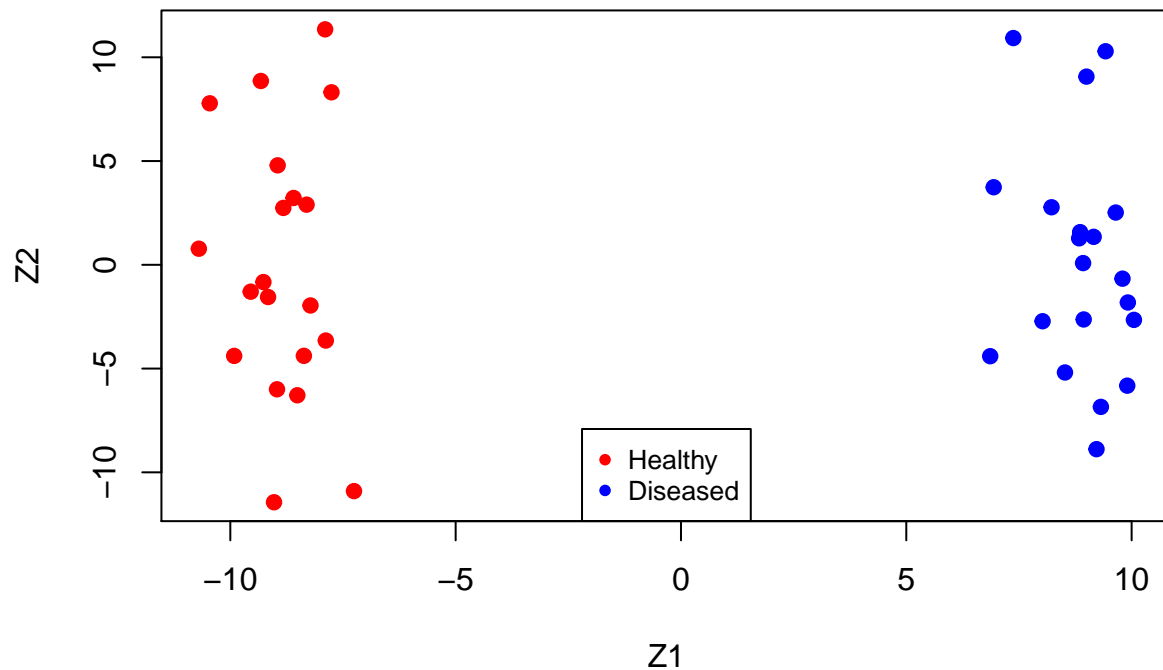
For the first principal component we want $Z_1 = \phi_1^T X$ subject to $\|\phi_1\|_2 = 1$. We want Z_1 to have the highest possible variance $V(Z_1) = \phi_1^T \Sigma \phi_1$, where Σ is the covariance matrix of X . The first principal component scores are then the column eigenvector corresponding to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

d)

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

```
pca_gene=prcomp(GeneData, scale=TRUE)
plot(pca_gene$x[,1:2], col=c("red","blue")[color],pch=19,xlab="Z1",ylab="Z2",main="Plot of tissue samples",
legend("bottom", legend=c("Healthy", "Diseased"),
col=c("red", "blue"), cex=0.8,pch=16)
```


Plot of tissue samples using PCA



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, which is very low.

e)

We now will use the results from PCA to find which genes vary the most accross the two groups. When looking at the plot of the two first principal components above, we see that the two groups have distinct values along the first principal components. Therefore we will define the genes that vary the most accross the two groups as the ones who have the highest value on the first principal component score vector. We show the 10 best ones below.

```
gene_loading = pca_gene$rotation[, 1]
sort(abs(gene_loading), decreasing = TRUE)[1:10]
```

```
##      G502      G589      G565      G590      G600      G551      G593
## 0.09485044 0.09449766 0.09183823 0.09173169 0.09167322 0.08768360 0.08758616
##      G538      G584      G509
## 0.08745400 0.08690858 0.08661015
```

f)

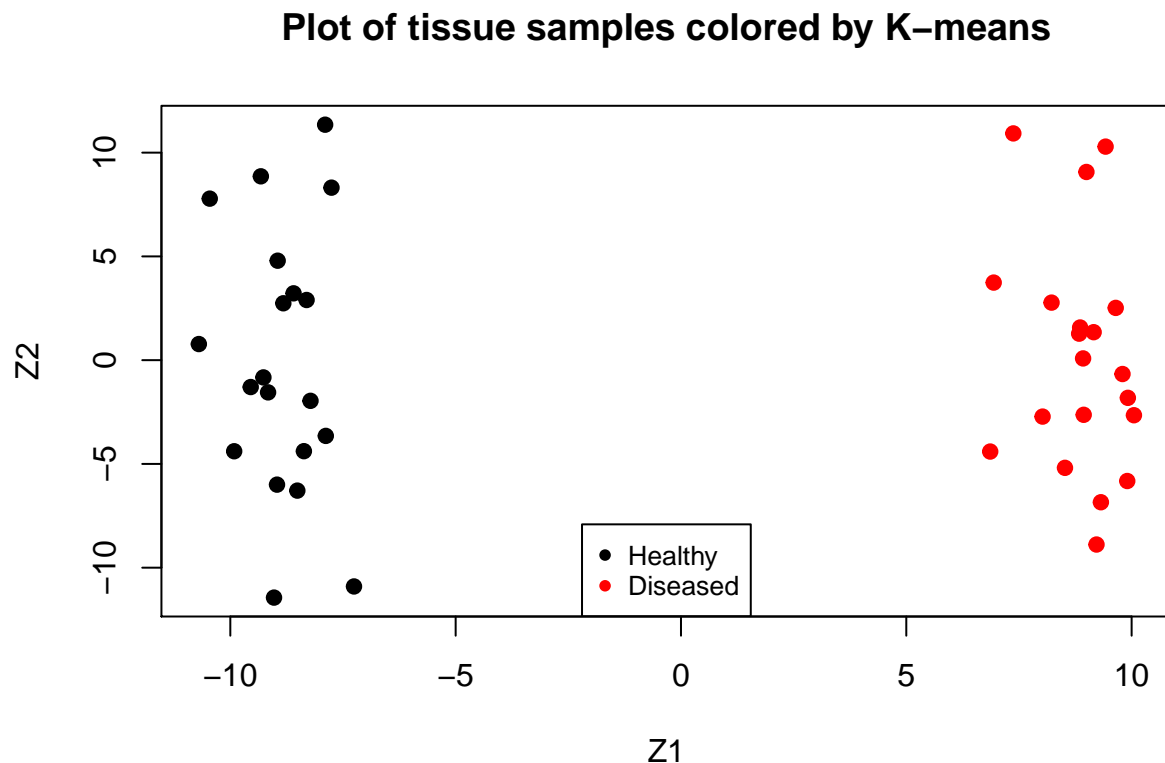
We now use K-means to separate the tissue samples into two groups.

```
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
```

This method separates the data perfectly and therefore yields an error rate of zero. Below is a plot of the data using PCA, where the observations are colored according to the results of K-means.

```
plot(pca_gene$x[,1:2], col=km.out$cluster,pch=19,xlab="Z1",ylab="Z2",main="Plot of tissue samples colored by K-means",
legend("bottom", legend=c("Healthy", "Diseased"),
col=c("black", "red"), cex=0.8,pch=16)
```



References

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