

# An example of ordinary least squares regression

```
mydata=read.csv2("Bilexempel.csv")
fit1=lm(Price~Year, data=mydata)
summary(fit1)
fit2=lm(Price~Year+Mileage+Equipment, data=mydata)
summary(fit2)
```

## Response variable:

Requested price of used Porsche cars (1000 SEK)

```
> summary(fit1)
```

```
Call:
lm(formula = Price ~ Year, data = mydata)
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-167683  -14683   20056   35933   72317
```

```
Coefficients:
```

```
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  78161027    8448038  -9.252 6.00e-13 ***
Year           39246         4226    9.288 5.25e-13 ***
```

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 57270 on 57 degrees of freedom
Multiple R-squared:  0.8821, Adjusted R-squared:  0.5952
F-statistic: 86.26 on 1 and 57 DF, p-value: 5.248e-13
```

## Inputs:

$X_1$  = Manufacturing year

$X_2$  = Milage (km)

$X_4$  = Equipment (0 or 1)

732A99/TDDE01

1

# An example of ordinary least squares regression

```
> summary(fit2)
```

```
Call:
lm(formula = Price ~ Year + Mileage + Equipment, data = mydata)
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-66223  -10525   -739   14128   65332
```

```
Coefficients:
```

```
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -2.083e+07  6.309e+06  -3.302  0.00169 **
Year          1.062e+04  3.154e+03   3.366  0.00139 **
Mileage       -2.077e+00  2.022e-01 -10.269 2.14e-14 ***
Equipment     5.790e+04  1.041e+04   5.563 8.08e-07 ***
```

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 29270 on 55 degrees of freedom
Multiple R-squared:  0.8997, Adjusted R-squared:  0.8942
F-statistic: 164.5 on 3 and 55 DF, p-value: < 2.2e-16
```

732A99/TDDE01

2

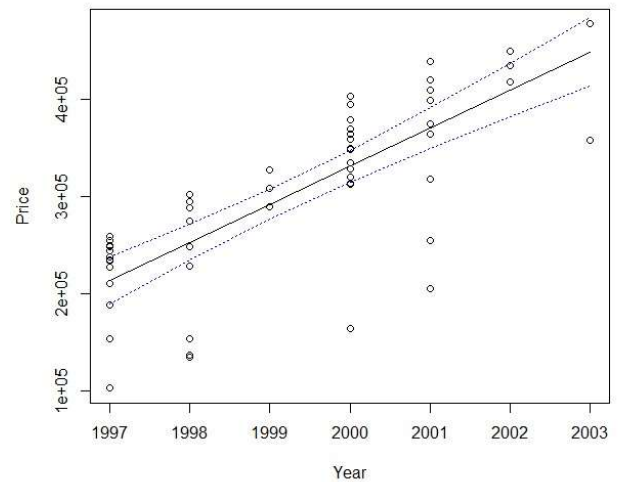
# An example of ordinary least squares regression

- Prediction

```
fitted <- predict(fit1, interval =  
"confidence")
```

```
# plot the data and the fitted line  
attach(mydata)  
plot(Year, Price)  
lines(Year, fitted[, "fit"])
```

```
# plot the confidence bands  
lines(Year, fitted[, "lwr"], lty = "dotted",  
col="blue")  
lines(Year, fitted[, "upr"], lty = "dotted",  
col="blue")  
detach(mydata)
```



732A99/TDDE01

3

## Data scaling

```
set.seed(12345)  
n=nrow(data)  
id=sample(1:n, floor(n*0.5))  
train=data[id,]  
test=data[-id,]
```

```
library(caret)  
scaler=preProcess(train)  
trainS=predict(scaler,train)  
testS=predict(scaler,test)
```

```
fit3=lm(Price~Year+Mileage+Equipment,  
data=trainS)  
summary(fit3)
```

```
> summary(fit3)
```

Call:

```
lm(formula = Price ~ Year + Mileage + Equipment, data = trainS)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.89677	-0.21339	-0.01089	0.18955	0.49303

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	1.160e-14	6.814e-02	0.000	1.00000
Year	3.362e-01	1.038e-01	3.239	0.00337 **
Mileage	-5.594e-01	8.603e-02	-6.502	8.26e-07 ***
Equipment	1.951e-01	9.742e-02	2.003	0.05616 .

732A99/TDDE01

4



# Optimization

- How to optimize a given log-likelihood?
  - In R, use `optim(start_point, function, method)`
- **Example:** minimizing  $(x_1 - 1)^2 + (x_2 - 2)^2$

```
df=data.frame(x1=1,x2=2)
mylikelihood<-function(x){
  x1=x[1]
  x2=x[2]
  return((x1-df$x1)^2+(x2-df$x2)^2)
}
```

```
optim(c(0,0), fn=mylikelihood, method="BFGS")
```

```
> optim(c(0,0), fn=mylikelihood, method="BFGS")
$par
[1] 1 2

$value
[1] 6.357135e-26

$counts
function gradient
          9         3

$convergence
[1] 0
```

732A99/TDDE01

5

# Decision trees in R

- **tree** package
  - Alternative: **rpart**

```
tree(formula, data, weights, control, split = c("deviance", "gini"), ...)
print(), summary(), plot(), text()
```

**Example:** breast cancer as a function of biological measurements

```
library(tree)
n=dim(biopsy)[1]
fit=tree(class~., data=biopsy)
plot(fit)
text(fit, pretty=0)
fit
summary(fit)
```

732A99/TDDE01

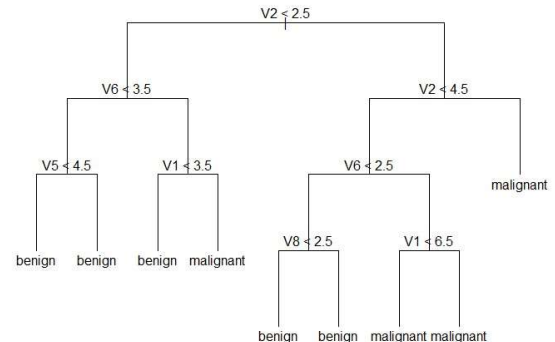
6

# Decision trees in R

- Adjust the splitting in the tree with *control* parameter (leaf size for ex)

```
> fit
node), split, n, deviance, yval, (yprob)
* denotes terminal node

1) root 683 884.400 benign ( 0.650073 0.349927 )
2) v2 < 2.5 418 108.900 benign ( 0.971292 0.028708 )
4) v6 < 3.5 395 25.130 benign ( 0.994937 0.005063 )
8) v5 < 4.5 389 0.000 benign ( 1.000000 0.000000 ) *
9) v5 > 4.5 6 7.638 benign ( 0.666667 0.333333 ) *
5) v6 > 3.5 23 31.490 benign ( 0.565217 0.434783 )
10) v1 < 3.5 11 0.000 benign ( 1.000000 0.000000 ) *
11) v1 > 3.5 12 10.810 malignant ( 0.166667 0.833333 ) *
3) v2 > 2.5 265 217.900 malignant ( 0.143396 0.856604 )
6) v2 < 4.5 90 120.300 malignant ( 0.388889 0.611111 )
12) v6 < 2.5 30 27.030 benign ( 0.833333 0.166667 )
24) v8 < 2.5 19 0.000 benign ( 1.000000 0.000000 ) *
25) v8 > 2.5 11 15.160 benign ( 0.545455 0.454545 ) *
13) v6 > 2.5 60 54.070 malignant ( 0.166667 0.833333 )
26) v1 < 6.5 28 35.160 malignant ( 0.321429 0.678571 ) *
27) v1 > 6.5 32 8.900 malignant ( 0.031250 0.968750 ) *
7) v2 > 4.5 175 30.350 malignant ( 0.017143 0.982857 ) *
```



```
> summary(fit)

Classification tree:
tree(formula = class ~ ., data = biopsy)
Variables actually used in tree construction:
[1] "V2" "V6" "V5" "V1" "V8"
Number of terminal nodes: 9
Residual mean deviance: 0.1603 = 108 / 674
Misclassification error rate: 0.03221 = 22 / 683
```

732A99/TDDE01

7

# Decision trees in R

- Misclassification results

```
Yfit=predict(fit, newdata=biopsy, type="class")
table(biopsy$class,Yfit)
```

```
> table(biopsy$class,Yfit)
      Yfit
      benign malignant
benign    440        18
malignant   7       234
```

732A99/TDDE01

8



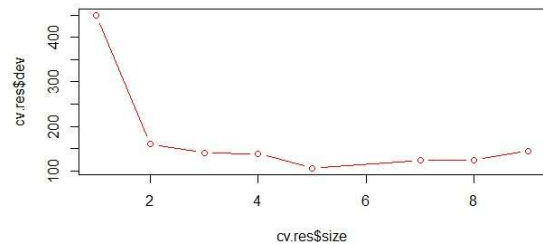
# Decision trees in R

- Selecting optimal tree by penalizing

## – Cv.tree()

```
set.seed(12345)
ind=sample(1:n, floor(0.5*n))
train=biopsy[ind,]
valid=biopsy[-ind,]

fit=tree(class~., data=train)
set.seed(12345)
cv.res=cv.tree(fit)
plot(cv.res$size, cv.res$dev, type="b",
     col="red")
```



What is optimal number of leaves?

732A99/TDDE01

9

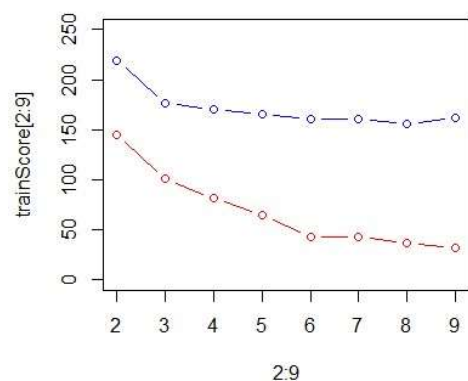
# Decision trees in R

- Selecting optimal tree by train/validation

```
fit=tree(class~., data=train)

trainScore=rep(0,9)
testScore=rep(0,9)

for(i in 2:9) {
  prunedTree=prune.tree(fit,best=i)
  pred=predict(prunedTree, newdata=valid,
               type="tree")
  trainScore[i]=deviance(prunedTree)
  testScore[i]=deviance(pred)
}
plot(2:9, trainScore[2:9], type="b", col="red",
     ylim=c(0,250))
points(2:9, testScore[2:9], type="b", col="blue")
```



What is optimal number of leaves?

732A99/TDDE01

10

# Decision trees in R

- Final tree: 5 leaves

```
finalTree=prune.tree(fit, best=5)
Yfit=predict(finalTree, newdata=valid,
type="class")
table(valid$class,Yfit)
```

```
> table(valid$class,Yfit)
      Yfit
      benign malignant
benign    222         8
malignant   6       114
```

732A99/TDDE01

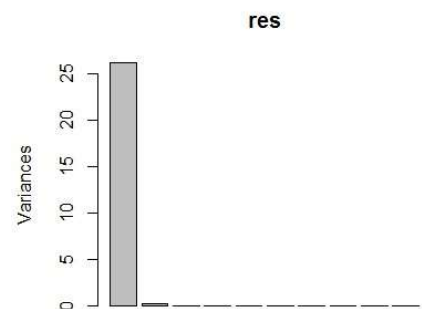
11

# PCA in R

- Prcomp(), biplot(), screeplot()

```
mydata=read.csv2("tecator.csv")
data1=mydata
data1$Fat=c()
res=prcomp(data1)
lambda=res$sdev^2
#eigenvalues
lambda
#proportion of variation
sprintf("%.2.3f",lambda/sum(lambda)*100)
screeplot(res)
```

```
> lambda
[1] 2.612713e+01 2.385369e-01 7.844883e-02 3.018501e-
[7] 2.052212e-04 1.084213e-04 2.077326e-05 1.150359e-
> sprintf("%.2.3f",lambda/sum(lambda)*100)
[1] "98.679" "0.901" "0.296" "0.114" "0.006"
[9] "0.000" "0.000" "0.000" "0.000" "0.000"
```



Only 1 component captures the 99% of variation!

732A99/TDDE01

12



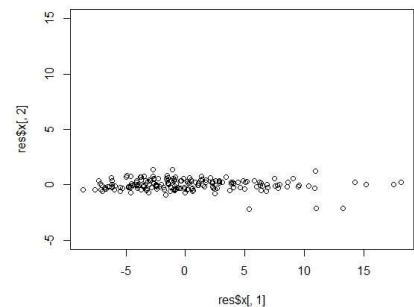
# PCA in R

- Principal component **loadings** (U)

```
U=res$rotation
head(U)
```

```
> head(U)
```

	PC1	PC2	PC3	
Channel1	0.07938192	0.1156228	0.08073156	-0.0927
Channel2	0.07987445	0.1170972	0.07887873	-0.0981
Channel3	0.08036498	0.1185571	0.07702127	-0.1031
Channel4	0.08085611	0.1200006	0.07515015	-0.1077
Channel5	0.08135022	0.1214075	0.07323819	-0.1119
Channel6	0.08184806	0.1227401	0.07125048	-0.1156



- Data in (PC1, PC2) = **scores** (Z)

Do we need the second dimension?

732A99/TDDE01

13

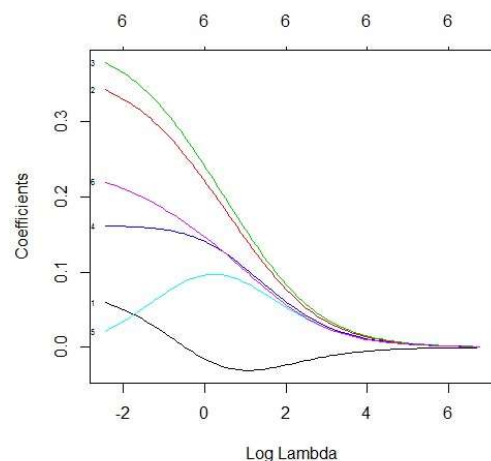
# Ridge regression

- R code: use package **glmnet** with alpha=0 (Ridge regression)
- Seeing how Ridge converges

```
data=read.csv("machine.csv", header=F)

library(caret)
library(glmnet)
scaler=preProcess(data)
data1=predict(scaler, data)
covariates=data1[,3:8]
response=data1[, 9]

model0=glmnet(as.matrix(covariates),
response, alpha=0,family="gaussian")
plot(model0, xvar="lambda", label=TRUE)
```



732A99/TDDE01

14

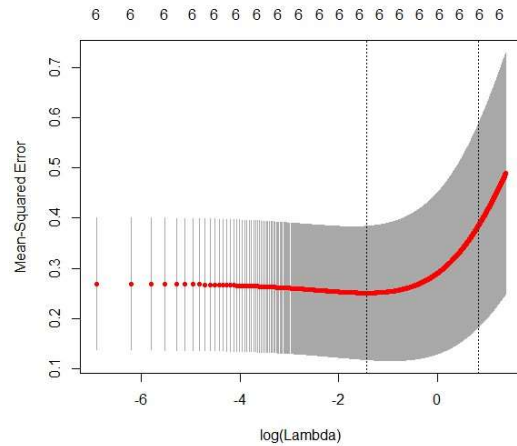
# Ridge regression

- Choosing the best model by cross-validation:

```
model=cv.glmnet(as.matrix(covariates),
response, alpha=0,family="gaussian")
model$lambda.min
plot(model)
coef(model, s="lambda.min")
```

```
> coef(model, s="lambda.min")
7 x 1 sparse Matrix of class "dgCMatrix"

(Intercept) 1
v3 3.420739e-02
v4 3.085696e-01
v5 3.403839e-01
v6 1.593470e-01
v7 5.489116e-02
v8 1.970982e-01
```



```
> model$lambda.min
[1] 0.046
```

732A99/TDDE01

15

# Ridge regression

- How good is this model in prediction?

```
covariates=train[,1:6]
response=train[, 7]
model=cv.glmnet(as.matrix(covariates), response, alpha=0,family="gaussian",
lambda=seq(0,1,0.001))
y=test[,7]
ynew=predict(model, newx=as.matrix(test[, 1:6]), type="response")
```

```
#Coefficient of determination
sum((ynew-mean(y))^2)/sum((y-mean(y))^2)

sum((ynew-y)^2)
```

Note that data are so small so numbers change much for other train/test

```
> sum((ynew-mean(y))^2)/sum((y-mean(y))^2)
[1] 0.5438148
> sum((ynew-y)^2)
[1] 18.04988
> |
```

732A99/TDDE01

16



# LASSO vs Ridge

- In R, use glmnet with **alpha=1**
- Only 5 variables selected by LASSO

```
> coef(model, s="lambda.min")
7 x 1 sparse Matrix of class "dgCMatrix"
      1
(Intercept) -5.091825e-17
v3          6.350488e-02
v4          3.578607e-01
v5          4.033670e-01
v6          1.541329e-01
v7          .
v8          2.287134e-01
> l
```

```
> sum((ynew-mean(y))^2)/sum((y-mean(y))^2)
[1] 0.5826904
> sum((ynew-y)^2)
[1] 16.63756
```

732A99/TDDE01

17

# Optimization methods in R

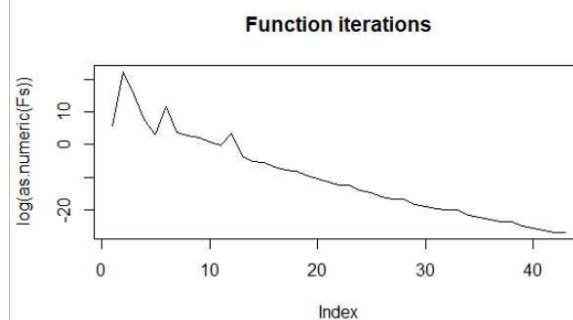
#Workaround: optim does not return iterations

```
Fs=list()
Params=list()
k=0

myf<- function(x){
  f=(x[1]-2)^4+(x[2]-4)^4
  .GlobalEnv$k= .GlobalEnv$k+1
  .GlobalEnv$Fs[[k]]=f
  .GlobalEnv$Params[[k]]=x
  return(f)
}
myGrad <-function(x) c(4*(x[1]-2)^3, 4*(x[2]-4)^3)

res<-optim(c(0,0), fn=myf, gr=myGrad, method="BFGS")

plot(log(as.numeric(Fs)), type="l", main="Function
iterations")
```



732A99/TDDE01

18

# LAB 1

The data file `optdigits.csv` contains information about normalized bitmaps of handwritten digits from a preprinted form from a total of 43 people. The data were first derived as 32x32 bitmaps which were then divided into nonoverlapping blocks of 4x4 and the number of on pixels are counted in each block. This has generated the resulting image of size 8x8 where each element is an integer in the range 0..16. Accordingly, each row in the data file is a sequence corresponding to 8x8 matrix, and the last element shows the actual digit from 0 to 9

1. Import the data into R and divide it into training, validation and test sets (50%/25%/25%)

```
library("kknn")

# Read in the optdigits.csv file and store the data in a variable named "data"
# Set the "header" argument to "F" to specify that the file does not contain a header row
data = read.csv("optdigits.csv", header = F)

# Get the number of rows in the data and store it in a variable named "n"
n = dim(data)[1]

# Set the seed of the random number generator to 12345
set.seed(12345)

# Randomly select 50% of the rows in the data and store them in a variable named "train"
id=sample(1:n, floor(n*0.5))
train=data[id,]

# Get the remaining rows that were not selected for the training set
id1=setdiff(1:n, id)
set.seed(12345)

# Randomly select 25% of the remaining rows and store them in a variable named "valid"
id2=sample(id1, floor(n*0.25))
valid=data[id2,]

# Get the remaining rows that were not selected for the training or validation sets
id3=setdiff(id1,id2)

# Store the remaining rows in a variable named "test"
test=data[id3,]

# Fit a k-NN classifier to the training data, with k=30 and a rectangular kernel
# Use the classifier to predict the class labels for the training data and store the predictions in "k_train"
k_train <- kknn(as.factor(V65)~., train = train, test = train, k = 30, kernel = "rectangular")

# Use the classifier to predict the class labels for the validation data and store the predictions in "k_valid"
k_valid <- kknn(as.factor(V65)~., train = train, test = valid, k = 30, kernel = "rectangular")

# Use the classifier to predict the class labels for the test data and store the predictions in "k_test"
k_test <- kknn(as.factor(V65)~., train = train, test = test, k = 30, kernel = "rectangular")

# Create a confusion matrix for the training data by comparing the predicted labels to the true labels
train_table = table(k_train$fitted.values, train$V65)

# Create a confusion matrix for the validation data by comparing the predicted labels to the true labels
valid_table = table(k_valid$fitted.values, valid$V65)

# Create a confusion matrix for the test data by comparing the predicted labels to the true labels
test_table = table(k_test$fitted.values, as.factor(test$V65))

# Define a function to calculate the misclassification rate for a given set of predictions and true labels
missclass = function(X, X1) {
  n = length(X)
  return (1 - sum(diag(table(X1, X)))/n)
}

missclass_test = missclass(k_test$fitted.values, as.factor(test$V65))
missclass_train = missclass(train$V65, k_train$fitted.values)
missclass_valid = missclass(k_valid$fitted.values, valid$V65)
```

Confusion matrices

```
test_table
```



```
##
##      0  1  2  3  4  5  6  7  8  9
## 0 77  0  0  0  0  0  0  0  0
## 1  0 81  0  0  0  1  0  0  7  1
## 2  0  2 98  0  0  1  0  0  0  1
## 3  0  0  0 107  0  0  0  1  1  1
## 4  1  0  0  0 94  0  0  0  0  0
## 5  0  0  0  2  0 93  0  0  0  0
## 6  0  0  0  0  2  2 90  0  0  0
## 7  0  0  0  0  6  1  0 111  0  1
## 8  0  0  3  1  2  0  0  0 70  0
## 9  0  3  0  1  5  5  0  0  0 85
```

```
train_table
```

```
##
##      0  1  2  3  4  5  6  7  8  9
## 0 202  0  0  0  1  0  0  0  0  1
## 1  0 179  1  0  3  0  2  3 10  3
## 2  0  11 190  0  0  0  0  0  0
## 3  0  0  0 185  0  1  0  0  2  5
## 4  0  0  0  0 159  0  0  0  0  2
## 5  0  0  0  1  0 171  0  0  0  0
## 6  0  0  0  0  0  0 190  0  2  0
## 7  0  1  1  1  7  1  0 178  0  3
## 8  0  1  0  0  1  0  0  1 188  3
## 9  0  3  0  1  4  8  0  0  2 183
```

Missclass Values

```
missclass_test
```

```
## [1] 0.05329154
```

```
missclass_train
```

```
## [1] 0.04500262
```

Comment on the quality of predictions for different digits and on the overall prediction quality.

Some numbers have worse predictions than others, for example number 4 was predicted to be a 7 a couple of times. This is most likely due to slopy writing that makes different numbers look more similar to others. The overall prediction quality is good.

#Exercise 3

- Find any 2 cases of digit "8" in the training data which were easiest to classify and 3 cases that were hardest to classify (i.e. having highest and lowest probabilities of the correct class). Reshape features for each of these cases as matrix 8x8 and visualize the corresponding digits (by using e.g. `heatmap()` function with parameters `Colv=NA` and `Rowv=NA`) and comment on whether these cases seem to be hard or easy to recognize visually

```

# Get the probability of each training sample being classified as the digit "8"
prob_eight <- k_train$prob[, 9]

# Order the probabilities in decreasing order and store the resulting indices in "ordered_high"
# Keep only the top two indices
ordered_high <- order(prob_eight, decreasing = T)
ordered_high <- ordered_high[1:2]

# Order the probabilities in increasing order and store the resulting indices in "ordered_low"
ordered_low <- order(prob_eight, decreasing = F)

# Get the indices of the "ordered_low" vector that correspond to samples with a non-zero probability of being classified as "8"
ordered_index <- which(prob_eight > 0)

# Create an empty vector to store the indices of the first three samples in "ordered_low" that are actually the digit "8"
ordered_low_eights <- c()

# Set an index variable to 1
index = 1

# Use a loop to append the first three indices from "ordered_low" that correspond to samples that are the digit "8" to "ordered_low_eights"
while(length(ordered_low_eights) < 3) {
  # Get the class label of the current sample
  value <- train[ordered_low[index], 65]

  # If the sample is the digit "8", append its index to "ordered_low_eights"
  if (value == 8) {
    ordered_low_eights <- append(ordered_low_eights, ordered_low[index])
  }

  # Increment the index variable
  index <- index + 1
}

# Print the "ordered_low_eights" vector
print(ordered_low_eights)

```

```
## [1] 520 431 1294
```

```

# Combine the "ordered_low_eights" and "ordered_high" vectors to create the "ordered_matrix" vector
ordered_matrix <- c(ordered_low_eights, ordered_high)

# Use a loop to create a heatmap for each of the samples corresponding to the indices in the "ordered_matrix" vector
for (i in ordered_matrix){
  # Reshape the sample's pixel values into a 8x8 matrix
  my_heatmap <- matrix(as.numeric(train[i,1:64]), nrow=8, ncol=8, byrow = T)

  # Create a heatmap of the matrix and set the main title to the true class label of the sample
  heatmap(my_heatmap, Colv = NA, Rowv = NA, main = train[i, 65])
}

```

The eights that were easy to classify are also very easy to see as eights. The eights that were most difficult to spot are very difficult to see without knowing that they're eights beforehand.

#### #Exercise 4

Fit a K-nearest neighbor classifiers to the training data for different values of  $K = 1, 2, \dots, 30$  and plot the dependence of the training and validation misclassification errors on the value of  $K$  (in the same plot). How does the model complexity change when  $K$  increases and how does it affect the training and validation errors? Report the optimal  $K$  according to this plot. Finally, estimate the test error for the model having the optimal  $K$ , compare it with the training and validation errors and make necessary conclusions about the model quality.



```

# Initialize two empty vectors to store the misclassification rates for the training and validation sets
# for different values of k
train_miss_error <- numeric(30)
val_miss_error <- numeric(30)

# Use a loop to fit a k-NN classifier to the training data for each value of k from 1 to 30
# Use the classifier to predict the class labels for the training and validation sets
# Calculate the misclassification rate for each set and store the results in the "train_miss_error" and "val_miss_error" vectors
for (i in 1:30){
  k_valid <- kknn(as.factor(V65)~., train = train, test = valid, k = i, kernel = "rectangular")
  k_train <- kknn(as.factor(V65)~., train = train, test = train, k = i, kernel = "rectangular")

  train_miss_error[i] <- missclass(k_train$fitted.values, train$V65)
  val_miss_error[i] <- missclass(k_valid$fitted.values, valid$V65)
}

# Create a plot of the "train_miss_error" and "val_miss_error" vectors
# Use k as the x-axis and the misclassification rate as the y-axis
plot = plot(c(1:30), train_miss_error, ylab = "train_miss_error", xlab = "k", col="pink")
points(c(1:30), val_miss_error, col="green")

```

```

# Fit a k-NN classifier to the training data with k=4
# Use the classifier to predict the class labels for the training, validation, and test sets
k_test <- kknn(as.factor(V65)~., train = train, test = test, k = 4, kernel = "rectangular")
k_train <- kknn(as.factor(V65)~., train = train, test = train, k = 4, kernel = "rectangular")
k_valid <- kknn(as.factor(V65)~., train = train, test = valid, k = 4, kernel = "rectangular")

missclass_test = missclass(k_test$fitted.values, as.factor(test$V65))
missclass_train = missclass(train$V65, k_train$fitted.values)
missclass_valid = missclass(k_valid$fitted.values, valid$V65)

```

How does the model complexity change when K increases and how does it affect the training and validation errors?

The model complexity remains the same as k increases, since it's not learning more parameters or features as k increases. One can argue that the overall complexity is increased since there are more data points added, but this does not mean that the complexity of the model is higher. The best K values are the ones with the lowest missclass error, in this case k=3 and k=4 are the best values.

We can see that train gives the lowest missclass error because it is the data we made the model from. Both test and validation is slightly higher which is to be expected. With the right k value computed with the validation data we get a test error which is just above 2.5% Which results in a pretty high quality model.

#### #Exercise 5

Fit K-nearest neighbor classifiers to the training data for different values of  $KK = 1, 2, \dots, 30$ , compute the error for the validation data as cross-entropy ( when computing log of probabilities add a small constant within log, e.g.  $1e-15$ , to avoid numerical problems) and plot the dependence of the validation error on the value of  $KK$ . What is the optimal  $KK$  value here? Assuming that response has multinomial distribution, why might the cross-entropy be a more suitable choice of the error function than the misclassification error for this problem?

```

cross.entropy <- function(p, phat){
  x <- 0
  for (i in 1:length(p)){
    x <- x + (p[i] * log(phat[i]))
  }
  return(-x)
}

train_miss_error <- as.vector(matrix(0,ncol = 30))
val_miss_error <- as.vector(matrix(0,ncol = 30))

entropy_error <- as.vector(matrix(0,ncol = 30))

for (i in 1:30){
  k_valid <- kknn(as.factor(V65)~., train = train, test = valid, k = i, kernel = "rectangular")
  k_train <- kknn(as.factor(V65)~., train = train, test = train, k = i, kernel = "rectangular")

  for (j in 0:9){
    cross_val <- valid$V65 == j
    cross_train <- train$V65 == j

    bool_val <- (which(cross_val, useNames = T))
    prob_val <- k_valid$prob[cross_val, as.character(j)] + 1e-15

    bool_train <- (which(cross_train, useNames = T))
    prob_train <- k_train$prob[cross_train, as.character(j)] + 1e-15

    val_miss_error[i] <- val_miss_error[i] + sum(-(bool_val*log(prob_val)))
    train_miss_error[i] <- train_miss_error[i] + sum(-(bool_train*log(prob_train)))

    #val_miss_error[i] <- cross.entropy(bool_val, prob_val)
    #train_miss_error[i] <- cross.entropy(bool_train, prob_train)

    entropy_error[i] <- abs(val_miss_error[i]-train_miss_error[i])
  }
}

plot(c(1:30), entropy_error, ylab = "cross_entropy_error", xlab = "k", col="blue")

```

```
which.min(entropy_error)
```

```
## [1] 5
```

What is the optimal  $KK$  value here? Assuming that response has multinomial distribution, why might the cross-entropy be a more suitable choice of the error function than the missclassification error for this problem?

This model might be more suitable because of lower complexity levels than missclass error. The optimal  $K$  value here is  $K = 5$ .

## #Assignment 2

The data file parkinson.csv is composed of a range of biomedical voice measurements from 42 people with early-stage Parkinson's disease recruited to a six-month trial of a telemonitoring device for remote symptom progression monitoring. The purpose is to predict Parkinson's disease symptom score (motor UPDRS) from the following voice characteristics: • Jitter(%),Jitter(Abs),Jitter:RAP,Jitter:PPQ5,Jitter:DDP - Several measures of variation in fundamental frequency • Shimmer,Shimmer(dB),Shimmer:APQ3,Shimmer:APQ5,Shimmer:APQ11,Shimmer:DDA - Several measures of variation in amplitude • NHR,HNR - Two measures of ratio of noise to tonal components in the voice • RPDE - A nonlinear dynamical complexity measure • DFA - Signal fractal scaling exponent • PPE - A nonlinear measure of fundamental frequency variation

Divide it into training and test data (60/40) and scale it appropriately. In the coming steps, assume that motor\_UPDRS is normally distributed and is a function of the voice characteristics, and since the data are scaled, no intercept is needed in the modelling.

```
library(caret)
```

```
## Warning: package 'caret' was built under R version 4.1.3
```

```
## Loading required package: ggplot2
```

```
## Warning: package 'ggplot2' was built under R version 4.1.3
```

```
## Loading required package: lattice
```



```
##  
## Attaching package: 'caret'
```

```
## The following object is masked from 'package:kkn':  
##  
##      contr.dummy
```

```
parkin = read.csv("parkinsons.csv", header = T)  
parkin_corr <- data.frame(parkin[c(5,7:22)]) #Remove unused voice characteristics  
parkin_scaled <- as.data.frame(scale(parkin_corr))  
  
n = dim(parkin_corr)[1]  
set.seed(12345)  
  
id=sample(1:n, floor(n*0.6))  
  
train_pre=parkin_corr[id,]  
test_pre=parkin_corr[-id,]  
  
data_scaler<-preProcess(train_pre)  
train<-predict(data_scaler,train_pre)  
test<-predict(data_scaler,test_pre)
```

## #Exercise 2

Compute a linear regression model from the training data, estimate training and test MSE and comment on which variables contribute significantly to the model.

```
fit = lm(motor_UPDRS ~ .-1, data = train)  
sum = summary(fit)  
MSE_train<-mean(sum$residuals^2)  
pred_test<-predict(fit,test)  
MSE_test<-mean((test$motor_UPDRS-pred_test)^2)  
print(MSE_train)
```

```
## [1] 0.8785431
```

```
print(MSE_test)
```

```
## [1] 0.9354477
```

```
print(sum)
```

```
##
## Call:
## lm(formula = motor_UPDRS ~ . - 1, data = train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.0255 -0.7363 -0.1087  0.7333  2.1960
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## Jitter...      0.186931   0.149561   1.250 0.211431
## Jitter.Abs.    -0.169609   0.040805  -4.157 3.31e-05 ***
## Jitter.RAP     -5.269544  18.834160  -0.280 0.779658
## Jitter.PPQ5    -0.074568   0.087766  -0.850 0.395592
## Jitter.DDP      5.249558  18.837525   0.279 0.780510
## Shimmer        0.592436   0.205981   2.876 0.004050 **
## Shimmer.dB.    -0.172655   0.139316  -1.239 0.215315
## Shimmer.APQ3   32.070932  77.159242   0.416 0.677694
## Shimmer.APQ5   -0.387507   0.113789  -3.405 0.000668 ***
## Shimmer.APQ11  0.305546   0.061236   4.990 6.34e-07 ***
## Shimmer.DDA   -32.387241  77.158814  -0.420 0.674695
## NHR            -0.185387   0.045567  -4.068 4.84e-05 ***
## HNR            -0.238543   0.036395  -6.554 6.41e-11 ***
## RPDE           0.004068   0.022664   0.179 0.857556
## DFA            -0.280318   0.020136 -13.921 < 2e-16 ***
## PPE            0.226467   0.032881   6.887 6.70e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.9394 on 3509 degrees of freedom
## Multiple R-squared:  0.1212, Adjusted R-squared:  0.1172
## F-statistic: 30.25 on 16 and 3509 DF, p-value: < 2.2e-16
```

The variables p-values that are higher than 0.05 does not contribute at all (insignificant P-values). The values that contribute significantly are the ones with low p-values. The variables that will contribute significantly to the model are the ones with the highest amount of stars in the table, which represents the lowest p-values.

MSE\_train = 0.8785431 MSE\_test = 0.9354477

### #Exercise3

Implement 4 following functions by using basic R commands only (no external packages): a. Loglikelihood function that for a given parameter vector  $\theta\theta$  and dispersion  $\sigma\sigma$  computes the log-likelihood function  $\log PP(TT|\theta\theta, \sigma\sigma)$  for the stated model and the training data

- Ridge function that for given vector  $\theta\theta$ , scalar  $\sigma\sigma$  and scalar  $\lambda\lambda$  uses function from 3a and adds up a Ridge penalty  $\lambda\lambda\|\theta\theta\|^2$  to the minus loglikelihood.
- RidgeOpt function that depends on scalar  $\lambda\lambda$ , uses function from 3b and function `optim()` with `method="BFGS"` to find the optimal  $\theta\theta$  and  $\sigma\sigma$  for the given  $\lambda\lambda$ .
- Df function that for a given scalar  $\lambda\lambda$  computes the degrees of freedom of the Ridge model based on the training data.

```

log_likelihood <- function(train, Y, theta, sigma){

  n <- dim(train)[1]
  train_theta = train%%theta

  sum1 <- n*log(sigma^2)/2
  sum2 <- n*log(2*pi)/2
  sum3 <- sum((train_theta-Y)^2)
  sum4 <- sum3/(2*sigma^2)

  return (-sum1-sum2-sum4)
}

ridge <- function(train, theta, lambda, Y){
  n<-dim(train)[2]
  sigma <- theta[n+1]
  theta <-as.matrix(theta[1:n])
  log_like <- log_likelihood(theta=theta,Y=Y,sigma=sigma,train=train)
  ridge <- -log_like + lambda*sum(theta^2)
  return(ridge)
}

ridgeOpt <- function(lambda, train, Y){

  train <- as.matrix(train)
  N = dim(train)[2]
  init_theta = integer(N)
  init_sigma = 1

  opt <- optim(par = c(init_theta,init_sigma), fn = ridge, lambda = lambda, train = train, Y = Y, method = "BFGS"
)
  return(opt)
}

dF <- function(X, lambda){
  #From the course formula
  X <- as.matrix(X)
  Xt <- t(X)
  n <- dim(X)[2]
  I <- diag(n)
  P <- X%%solve((Xt%%X + (lambda*I)))%%Xt
  return(sum(diag(P)))
}

```

```

AIC = function(train,Y,theta, sigma, lambda){
  log_like = log_likelihood(train = train, Y=Y, theta = theta, sigma = sigma)
  N = dim(train)[1] # No of data points
  df = dF(train,lambda)
  aic = (-2*log_like/N) + (2*df/N) #(-2*Log-likelihood/N) + 2*(df/N)
  return(aic)
}

```

#### #Exercise 4

By using function RidgeOpt, compute optimal  $\theta\theta$  parameters for  $\lambda\lambda = 1$ ,  $\lambda\lambda = 100$  and  $\lambda\lambda = 1000$ . Use the estimated parameters to predict the motor\_UPDRS values for training and test data and report the training and test MSE values. Which penalty parameter is most appropriate among the selected ones? Compute and compare the degrees of freedom of these models and make appropriate conclusions.



```
xtrain<-as.matrix(train[2:17])
ytrain<-as.matrix(train[1])
xtest=as.matrix(test[2:17])
ytest<-as.matrix(test[1])

for (lambda in c(1, 100, 1000)){
  opt = ridgeOpt(lambda, xtrain, ytrain)
  theta <- as.matrix(opt$par[1:16])
  sigma<- opt$par[17]
  MSE_train = mean((xtrain%*%theta - ytrain)^2)
  MSE_test = mean((xtest%*%theta - ytest)^2)
  aic = AIC(train=xtrain,Y= ytrain,theta = theta,sigma= sigma, lambda= lambda)
  print(paste("Lambda:",lambda))
  print(paste("TrainMSE:",MSE_train))
  print(paste("TestMSE:",MSE_test))
  print(paste("AIC:", aic)) #Unecessary
}
```

```
## [1] "Lambda: 1"
## [1] "TrainMSE: 0.878627086520567"
## [1] "TestMSE: 0.934996946852113"
## [1] "AIC: 2.71634659521451"
## [1] "Lambda: 100"
## [1] "TrainMSE: 0.884410410201325"
## [1] "TestMSE: 0.932331574068859"
## [1] "AIC: 2.72067414889541"
## [1] "Lambda: 1000"
## [1] "TrainMSE: 0.921121567431251"
## [1] "TestMSE: 0.953948178939008"
## [1] "AIC: 2.75891604085416"
```

Which penalty parameter is most appropriate among the selected ones? Compute and compare the degrees of freedom of these models and make appropriate conclusions.

We have read that to find the best aic value among mutiple aic values is the lowest. In our example it seems like  $\lambda = 1$  and  $\lambda = 100$  gives similar aic values which is hard to draw conclusions from.  $\lambda = 1$  gives the lowest AIC value. Since the the  $\lambda = 100$  has the lowest test\_MSE, this is the best value. The degrees of freedom reflects the complexity of the model, so the larger lambda, the smaller degrees of freedom we would get. We get more complex models, basically penalizing the complexity.

##Assignment 3.

The data file pima-indians-diabetes.csv contains information about the onset of diabetes within 5 years in Pima Indians given medical details. The variables are (in the same order as in the dataset):

Make a scatterplot showing a Plasma glucose concentration on Age where observations are colored by Diabetes levels. Do you think that Diabetes is easy to classify by a standard logistic regression model that uses these two variables as features? Motivate your answer.

```
prime = read.csv("pima-indians-diabetes.csv", header = F)

set.seed(12345)
```

```
colour <- function(x){
  if (x==1){
    c = "red"
  } else{
    c="green"
  }
  return(c)
}

colours = sapply(prime$V9, colour)
plot( prime$V2, prime$V8, xlab = "Plasma", ylab = "Age", main = "doabets", col = colours)
```

seems to be different, no obvious correlation between the variables. Making a linear regression would not provide a meaningful line.

#Exercise 2

Train a logistic regression model with  $y = \text{Diabetes}$  as target  $x_1 = \text{Plasma glucose concentration}$  and  $x_2 = \text{Age}$  as features and make a prediction for all observations by using  $rr = 0.5$  as the classification threshold. Report the probabilistic equation of the estimated model (i.e., how the target depends on the features and the estimated model parameters probabilistically). Compute also the training misclassification error and make a scatter plot of the same kind as in step 1 but showing the predicted values of Diabetes as a color instead. Comment on the quality of the classification by using these results

```
glm.fits = glm(V9~ V2 + V8, prime, family = "binomial" )
prob=predict(glm.fits, type="response")
pred=ifelse(prob>0.5, 'red','green')

table = table(pred, prime$V9)

miss <- missclass(pred, prime$V9)

plot(prime$V2, prime$V8,col=pred, xlab = "Plasma glucose levels", ylab = "Age", main = paste("Missclass error", t
oString(miss), sep=" = ") )
```

```
table
```

```
##
## pred      0    1
## green 436 138
## red    64 130
```

```
#summary(table)
#summary(glm.fits)
glm.fits$coefficients
```

```
## (Intercept)          V2          V8
## -5.91244906  0.03564404  0.02477835
```

Probabilistic equation:

$$p(y = 1 | x_*) = g(x_*, \hat{\theta}) = \frac{1}{1 + e^{-\hat{\theta}^T x_*}} = \frac{1}{1 + e^{-(-5.91244906 + 0.03564404x_1 + 0.02477835x_2)}}$$

It seems to be an acceptable classification with about 1/4 error but not perfect There seems to be a division when you reach above 150 in plasma-glucose levels. But this threshold seems to lower as you age. An older person is more likely to have diabetes even with lower plasma-glucose levels.

### #Exercise 3

Use the model estimated in step 2 to a) report the equation of the decision boundary between the two classes b) add a curve showing this boundary to the scatter plot in step 2. Comment whether the decision boundary seems to catch the data distribution well.

```
r=.5
glm.fits = glm(V9~ V2 + V8, prime, family = "binomial" )
cf = glm.fits$coefficients
prob=predict(glm.fits, type="response")
pred=ifelse(prob>0.5, 'red','green')

plot(prime$V2, prime$V8,col=pred, ylab = "Age", xlab= "Plasma", main = paste("Missclass Error", toString(miss), s
ep=" = "))

slope = coef(glm.fits)[2]/(-coef(glm.fits)[3])
intercept = (coef(glm.fits)[1] + log(-r/(r-1)))/(-coef(glm.fits)[3])

#clip(min(x1),max(x1), min(x2), max(x2))
abline(intercept , slope, lwd = 2, lty = 2)
```

```
he = log(-.2/(.2-1))
he
```

```
## [1] -1.386294
```

```
slope
```

```
##          V2
## -1.438516
```

```
intercept
```

```
## (Intercept)
##      238.6135
```

y = -1.438516X + 238.6135

We can see how the distribution is hard to define and capture in the plot from the first assignment. This attempt is a pretty good compromise, where most non-diabetics are captured in the green to the left, and most diabetes cases are captured to the right. But with a missclass error of 26% it still misses a lot, which we can also see when compared to the plot from the first assignment.

#### #Exercise 4

Make the same kind of plots as in step 2 but use thresholds  $r = 0.2$  and  $r = 0.8$ . By using these plots, comment on what happens with the prediction when  $r$  value changes.

```
for(r in c(.2,.5, .8)) {
  pred=ifelse(prob>r, 'red','green')
  table(pred, prime$V9)

  miss <- missclass(pred, prime$V9)

  plot(prime$V2, prime$V8,col=pred, ylab = "Age", xlab= "Plasma", main = paste("Missclass_Error = ", toString(miss), "\n r = ", r, sep= ""))

  slope = coef(glm.fits)[2]/(-coef(glm.fits)[3])
  intercept = (coef(glm.fits)[1] - log(-r/(r-1)))/(-coef(glm.fits)[3])

  #clip(min(x1),max(x1), min(x2), max(x2))
  abline(intercept , slope, lwd = 2, lty = 2)
}
```

With  $r=0.2$  and  $r=0.8$ , the missclassification error increased.  $r=0.2$  was clearly worse (37% error).

#### #Exercise 5

Perform a basis function expansion trick by computing new features  $z_1 = x_1^4$ ,  $z_2 = x_1^3 x_2$ ,  $z_3 = x_1^2 x_2^2$ ,  $z_4 = x_1 x_2^3$ ,  $z_5 = x_2^4$ , adding them to the data set and then computing a logistic regression model with  $y$  as target and  $x_1$ ,  $x_2$ ,  $z_1$ , ...,  $z_5$  as features. Create a scatterplot of the same kind as in step 2 for this model and compute the training misclassification rate. What can you say about the quality of this model compared to the previous logistic regression model? How have the basis expansion trick affected the shape of the decision boundary and the prediction accuracy?

```
expanded <- prime

expanded$z1 <- expanded$V2 ** 4
expanded$z2 <- expanded$V2 ** 3 * expanded$V8
expanded$z3 <- expanded$V2 ** 2 * expanded$V8 ** 2
expanded$z4 <- expanded$V2 * expanded$V8 ** 3
expanded$z5 <- expanded$V8 ** 4

glm.fits = glm(V9~ V2 + V8 + z1 + z2 + z3 + z4 + z5, expanded, family = "binomial" )

for(r in c(.2,.5, .8)) {
  prob=predict(glm.fits, type="response")
  pred=ifelse(prob>r, 'red','green')
  table(pred, expanded$V9)

  miss <- missclass(pred, prime$V9)

  plot(expanded$V2, expanded$V8,col=pred, ylab = "Age", xlab= "Plasma", main = paste("Missclass_Error = ", toString(miss), "\n r = ", r, sep= ""))
}
```

This model has a lower missclass error, so the quality of this model is slightly better than the previous one. The linearity of the model is however gone, since we've added non-linear variables. This causes the plotted colors to look more like a "slice"



# Untitled

```
library(glmnet)
```

```
## Loading required package: Matrix
```

```
## Loaded glmnet 4.1-4
```

```
library(tree)  
library(caret)
```

```
## Loading required package: ggplot2
```

```
## Loading required package: lattice
```

```
library(dplyr)
```

```
##  
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':  
##  
## filter, lag
```

```
## The following objects are masked from 'package:base':  
##  
## intersect, setdiff, setequal, union
```

```

# Assume that Fat can be modeled as a linear regression in which absorbance characteristics (Channels) are used as features. Report the underlying probabilistic model, fit the linear regression to the training data and estimate the training and test errors. Comment on the quality of fit and prediction and therefore on the quality of model.

# Read in data from "tecator.csv" and store in a data frame called "tecator"
tecator = read.csv("tecator.csv", header = T)

# Get the number of rows in the data frame and store in a variable called "n"
n = dim(tecator)[1]

# Set the seed for the random number generator
set.seed(12345)

# Create a new data frame called "df" that contains all the columns from the "tecator" data frame except for the first column
df = data.frame(tecator[c(2:102)])

# Generate a random sample of row indices from the "tecator" data frame
id = sample(1:n, floor(n*0.5))

# Use the "id" indices to select a subset of rows from the "df" data frame and store them in a new data frame called "train"
train = df[id,]

# Use the "id" indices to select all rows except those in the "id" subset from the "df" data frame and store them in a new data frame called "test"
test = df[-id,]

# Fit a linear regression model to the "train" data frame, with the "Fat" column as the response variable and all the other columns as predictor variables
fit = lm(Fat ~ ., data = train)

# Use the "fit" model to make predictions on the training data and store the predictions in a vector called "train_preds"
train_preds = predict(fit, train)

# Use the "fit" model to make predictions on the test data and store the predictions in a vector called "test_preds"
test_preds = predict(fit, test)

# Generate a summary of the "fit" model and store it in an object called "sum"
sum = summary(fit)

# Calculate the mean squared error (MSE) of the model's predictions on the training data and store in a variable called "MSE_train"
MSE_train = mean((train_preds - train$Fat)^2)

# Calculate the MSE of the model's predictions on the test data and store in a variable called "MSE_test"
MSE_test = mean((test_preds - test$Fat)^2)

# Print "Test error" and the MSE of the model's predictions on the test data
print("Test error")

```

```
## [1] "Test error"
```

```
MSE_test
```

```
## [1] 722.4294
```

```

# Print "Train error" and the MSE of the model's predictions on the training data
print("Train error")

```

```
## [1] "Train error"
```

```
MSE_train
```

```
## [1] 0.005709117
```

*#Fit the LASSO regression model to the training data. Present a plot illustrating how the regression coefficients depend on the log of penalty factor ( $\log \lambda$ ) and interpret this plot. What value of the penalty factor can be chosen if we want to select a model with only three features?"*

```
# Get the "Fat" column from the "train" data frame and store in a variable called "y"
y = train$Fat
```

```
# Get the first 100 columns from the "train" data frame and store in a variable called "x"
x = train[1:100]
```

```
# Fit a lasso regression model to the "x" and "y" data using the "glmnet" function, with a Gaussian error distribution
model_lasso= glmnet(as.matrix(x), as.matrix(y), alpha=1,family="gaussian")
```

```
# Plot the model's coefficients as a function of the regularization parameter "lambda"
plot(model_lasso, xvar = "lambda")
```

```
# Use the "model_lasso" model to make predictions on the "x" data and store the predictions in a variable called "ynew"
ynew=predict(model_lasso, newx=as.matrix(x), type="response")
```

*#Repeat step 3 but fit Ridge instead of the LASSO regression and compare the plots from steps 3 and 4. Conclusion s?"*

```
# Get the "Fat" column from the "train" data frame and store in a variable called "y"
y = train$Fat
```

```
# Get the first 100 columns from the "train" data frame and store in a variable called "x"
x = train[1:100]
```

```
# Fit a ridge regression model to the "x" and "y" data using the "glmnet" function, with a Gaussian error distribution
model_lasso= glmnet(as.matrix(x), as.matrix(y), alpha=0,family="gaussian")
```

```
# Plot the model's coefficients as a function of the regularization parameter "lambda"
plot(model_lasso, xvar = "lambda")
```

```
# Use the "model_lasso" model to make predictions on the "x" data and store the predictions in a variable called "ynew"
ynew=predict(model_lasso, newx=as.matrix(x), type="response")
```

*#Use cross-validation with default number of folds to compute the optimal LASSO model. Present a plot showing the dependence of the CV score on  $\log \lambda$  and comment how the CV score changes with  $\log \lambda$ . Report the optimal  $\lambda$  and how many variables were chosen in this model. Does the information displayed in the plot suggests that the optimal  $\lambda$  value results in a statistically significantly better prediction than  $\log \lambda = -4$ ? Finally, create a scatter plot of the original test versus predicted test values for the model corresponding to optimal  $\lambda$  and comment whether the model predictions are good."*

```
# Fit a lasso regression model to the "x" and "y" data using cross-validation and the "cv.glmnet" function, with a Gaussian error distribution
model_lasso= cv.glmnet(as.matrix(x), as.matrix(y), alpha=1,family="gaussian")
```

```
# Get the value of the regularization parameter "lambda" that resulted in the lowest cross-validation error
lambda_min = model_lasso$lambda.min
```

```
# Plot the model's coefficients as a function of the regularization parameter "lambda"
plot(model_lasso, xvar = "lambda")
```

```
# Fit a lasso regression model to the "x" and "y" data using the "glmnet" function and the "lambda_min" value, with a Gaussian error distribution
better_model = glmnet(as.matrix(x), as.matrix(y), lambda = lambda_min, alpha = 1, family = "gaussian")
```

```
# Use the "better_model" model to make predictions on the "x" data and store the predictions in a variable called "ynew"
ynew=predict(better_model, newx=as.matrix(x), s = lambda_min , type="response")
```

```
# Plot the original "y" values against the predicted "ynew" values, with a red color and a labeled x-axis and y-axis
plot(y, ynew, xlab = "Original", ylab = "Predicted",col = "red", main = "Scatter plot")
```

```
# Add a line with slope 1 and intercept 0 to the plot
abline(0,1)
```



*#Import the data to R, remove variable "duration" and divide into training/validation/test as 40/30/30: use data partitioning code specified in Lecture 2a."*

```
# Read in data from "bank-full.csv" and store in a data frame called "d"
d = read.csv("bank-full.csv", sep = ";", stringsAsFactors = TRUE)

# Create a copy of the "d" data frame called "data"
data = d

# Remove the "duration" column from the "data" data frame
data$duration = c()

# Get the "y" column from the "d" data frame and store in a variable called "output"
output = d['y']

# Get the number of rows in the "data" data frame and store in a variable called "n"
n = dim(data)[1]

# Set the seed for the random number generator
set.seed(12345)

# Generate a random sample of row indices from the "data" data frame
id=sample(1:n, floor(n*0.4))

# Use the "id" indices to select a subset of rows from the "data" data frame and store them in a new data frame called "train"
train=data[id,]

# Get the indices of the rows in the "data" data frame that are not in the "id" subset
id1=setdiff(1:n, id)

# Set the seed for the random number generator
set.seed(12345)

# Generate a random sample of row indices from the "id1" subset
id2=sample(id1, floor(n*0.3))

# Use the "id2" indices to select a subset of rows from the "data" data frame and store them in a new data frame called "valid"
valid=data[id2,]

# Get the indices of the rows in the "id1" subset that are not in the "id2" subset
id3=setdiff(id1,id2)

# Use the "id3" indices to select a subset of rows from the "data" data frame and store them in a new data frame called "test"
test=data[id3,]

#Fit decision trees to the training data so that you change the default settings one by one (i.e. not simultaneously): a. Decision Tree with default settings. b. Decision Tree with smallest allowed node size equal to 7000. c. Decision trees minimum deviance to 0.0005. and report the misclassification rates for the training and validation data. Which model is the best one among these three? Report how changing the deviance and node size affected the size of the trees and explain why."

# Fit a decision tree model to the "train" data using the "tree" function, using all of the columns except the "y" column as predictors
fit=tree(as.factor(y)~., data=train)

# Plot the decision tree model
plot(fit)

# Add text labels to the plot to show the splits at each node of the tree
text(fit, pretty=0)
```

```
# Fit a decision tree model to the "train" data using the "tree" function, using all of the columns except the "y" column as predictors, with a minimum size of 7000 for each leaf
fit2=tree(as.factor(y)~., data=train, minsize=7000)

# Plot the decision tree model
plot(fit2)

# Add text labels to the plot to show the splits at each node of the tree
text(fit2, pretty=0)
```

```
# Fit a decision tree model to the "train" data using the "tree" function, using all of the columns except the "y"
# column as predictors, with a minimum relative decrease in impurity of 0.0005 required to split a node
fit3=tree(as.factor(y)~., data=train, mindev=0.0005)
```

```
# Plot the decision tree model
plot(fit3)
```

```
# Add text labels to the plot to show the splits at each node of the tree
text(fit3, pretty=0)
```

```

#Use training and validation sets to choose the optimal tree depth in the model 2c: study the trees up to 50 leaves. Present a graph of the dependence of deviances for the training and the validation data on the number of leaves and interpret this graph in terms of bias-variance tradeoff. Report the optimal amount of leaves and which variables seem to be most important for decision making in this tree. Interpret the information provided by the tree structure (not everything but most important findings).

# Make predictions on the "train" data using the first decision tree model and store the predictions in a variable called "Yfit_t"
Yfit_t=predict(fit, newdata=train, type="class")

# Create a confusion matrix for the predictions on the "train" data using the first decision tree model
t1<-table(train$y,Yfit_t)

# Calculate the misclassification rate for the predictions on the "train" data using the first decision tree model
mis_t1 <- 1-sum(diag(t1))/sum(t1)

# Make predictions on the "train" data using the second decision tree model and store the predictions in a variable called "Yfit_t2"
Yfit_t2=predict(fit2, newdata=train, type="class")

# Create a confusion matrix for the predictions on the "train" data using the second decision tree model
t2<-table(train$y,Yfit_t2)

# Calculate the misclassification rate for the predictions on the "train" data using the second decision tree model
mis_t2 <- 1-sum(diag(t2))/sum(t2)

# Make predictions on the "train" data using the third decision tree model and store the predictions in a variable called "Yfit_t3"
Yfit_t3=predict(fit3, newdata=train, type="class")

# Create a confusion matrix for the predictions on the "train" data using the third decision tree model
t3<-table(train$y,Yfit_t3)

# Calculate the misclassification rate for the predictions on the "train" data using the third decision tree model
mis_t3 <- 1-sum(diag(t3))/sum(t3)

# Make predictions on the "valid" data using the first decision tree model and store the predictions in a variable called "Yfit_v"
Yfit_v=predict(fit, newdata=valid, type="class")

# Create a confusion matrix for the predictions on the "valid" data using the first decision tree model
v1<-table(valid$y,Yfit_v)

# Calculate the misclassification rate for the predictions on the "valid" data using the first decision tree model
mis_v1<-1-sum(diag(v1))/sum(v1)

# Make predictions on the "valid" data using the second decision tree model and store the predictions in a variable called "Yfit_v2"
Yfit_v2=predict(fit2, newdata=valid, type="class")

# Create a confusion matrix for the predictions on the "valid" data using the second decision tree model
v2<-table(valid$y,Yfit_v2)

# Calculate the misclassification rate for the predictions on the "valid" data using the second decision tree model
mis_v2<-1-sum(diag(v2))/sum(v2)

# Make predictions on the "valid" data using the third decision tree model and store the predictions in a variable called "Yfit_v3"
Yfit_v3=predict(fit3, newdata=valid, type="class")

# Create a confusion matrix for the predictions on the "valid" data using the third decision tree model
v3<-table(valid$y,Yfit_v3)

# Calculate the misclassification rate for the predictions on the "valid" data using the third decision tree model
mis_v3<-1-sum(diag(v3))/sum(v3)

# Print the misclassification rates for the predictions on the "train" and "valid" data using the first decision tree model
print("1) Training and validation")

## [1] "1) Training and validation"

```



```
print(mis_t1)
```

```
## [1] 0.1048441
```

```
print(mis_v1)
```

```
## [1] 0.1092679
```

```
# Print the misclassification rates for the predictions on the "train" and "valid" data using the second decision tree model  
print("2) Training and validation")
```

```
## [1] "2) Training and validation"
```

```
print(mis_t2)
```

```
## [1] 0.1048441
```

```
print(mis_v2)
```

```
## [1] 0.1092679
```

```
# Print the misclassification rates for the predictions on the "train" and "valid" data using the third decision tree model  
print("3) Training and validation")
```

```
## [1] "3) Training and validation"
```

```
print(mis_t3)
```

```
## [1] 0.09400575
```

```
print(mis_v3)
```

```
## [1] 0.1119221
```

*#Estimate the confusion matrix, accuracy and F1 score for the test data by using the optimal model from step 3. Comment whether the model has a good predictive power and which of the measures (accuracy or F1-score) should be preferred here."*

```
# Create two vectors called "trainScore" and "testScore" with length 50 and filled with 0s
trainScore=rep(0,50)
testScore=rep(0,50)

# Loop through the numbers 2 to 50
for(i in 2:50) {
  # Prune the third decision tree model to have "i" leaves
  prunedTree=prune.tree(fit3,best=i)

  # Make predictions on the "valid" data using the pruned decision tree model and store the predictions in a variable called "pred"
  pred=predict(prunedTree, newdata=valid, type="tree")

  # Calculate the deviance for the pruned decision tree model on the "train" data and store it in the "trainScore" vector
  trainScore[i]=deviance(prunedTree)

  # Calculate the deviance for the predictions on the "valid" data using the pruned decision tree model and store it in the "testScore" vector
  testScore[i]=deviance(pred)
}

# Plot the deviances for the pruned decision tree models on the "train" data (in red) and the deviances for the p
# redictions on the "valid" data (in blue)
plot(2:50, trainScore[2:50], type="b", col="red", ylim=c(min(testScore[-1]), max(trainScore[-1])))
points(2:50, testScore[2:50], type="b", col="blue")
```

```
# Find the number of leaves that gives the minimum deviance on the "valid" data
optimal_leaves = which.min(testScore[2:50])

# Prune the third decision tree model to have the optimal number of leaves
finalTree=prune.tree(fit3, best=optimal_leaves)

# Make predictions on the "valid" data using the final pruned decision tree model and store the predictions in a variable called "finalfit"
finalfit=predict(finalTree, newdata=valid, type="class")

# Create a confusion matrix for the predictions on the "valid" data using the final pruned decision tree model
tab = table(valid$y,finalfit)

# Plot the final pruned decision tree model
plot(finalTree)
```

```
# Print the summary of the third decision tree model
summary(fit3)
```

```
##
## Classification tree:
## tree(formula = as.factor(y) ~ ., data = train, mindev = 5e-04)
## Variables actually used in tree construction:
## [1] "poutcome" "month" "contact" "marital" "day" "campaign"
## [7] "job" "pdays" "age" "balance" "housing" "education"
## [13] "previous"
## Number of terminal nodes: 122
## Residual mean deviance: 0.5213 = 9363 / 17960
## Misclassification error rate: 0.09362 = 1693 / 18084
```

```
# Print the summary of the final pruned decision tree model
summary(finalTree)
```

```
##
## Classification tree:
## snip.tree(tree = fit3, nodes = c(581L, 17L, 577L, 79L, 37L, 77L,
## 576L, 153L, 580L, 6L, 1157L, 16L, 5L, 1156L, 156L, 152L, 579L,
## 7L))
## Variables actually used in tree construction:
## [1] "poutcome" "month" "contact" "pdays" "age" "day" "balance"
## [8] "housing"
## Number of terminal nodes: 21
## Residual mean deviance: 0.5706 = 10310 / 18060
## Misclassification error rate: 0.1041 = 1882 / 18084
```

```
#text(fit3, pretty=0)
```

*#Estimate the confusion matrix, accuracy and F1 score for the test data by using the optimal model from step 3. Comment whether the model has a good predictive power and which of the measures (accuracy or F1-score) should be preferred here."*

```
# Prune the decision tree model fit3 using the optimal number of leaves
opt_tree <- prune.tree(fit3, best=optimal_leaves)
```

```
# Make predictions on the test set using the pruned tree
ffitTest <- predict(opt_tree, newdata=test, type="class")
```

```
# Calculate the confusion matrix
c_m <- table(test$y, ffitTest)
```

```
# Print the confusion matrix
c_m
```

```
##      ffitTest
##      no  yes
## no 11812 167
## yes 1294 291
```

```
# Calculate the accuracy
acc <- sum(c_m[1], c_m[4])/sum(c_m[1:4])
```

```
# Calculate the precision
prec <- c_m[4] / sum(c_m[4], c_m[2])
```

```
# Calculate the recall
recall <- c_m[4] / sum(c_m[4], c_m[3])
```

```
# Calculate the F1 score
f1_score <- (2*(recall*prec)) / (recall + prec)
```

*#Perform a decision tree classification of the test data with the following loss matrix, 732A99/732A68/ TDDE01 Machine Learning Division of Statistics and Machine Learning Department of Computer and Information Science and report the confusion matrix for the test data. Compare the results with the results from step 4 and discuss how the rates has changed and why."*

```
# Make predictions on the test set using the pruned tree model
predtree5 <- predict(opt_tree, newdata=test, type="vector")
```

```
# Extract the probability of the positive class from the predictions
probY <- predtree5[,2]
```

```
# Extract the probability of the negative class from the predictions
probN <- predtree5[,1]
```

```
# Create a new vector of predictions based on the ratio of the probabilities
pred5 <- ifelse((probY/probN)>1/5, "yes", "no")
```

```
# Calculate the confusion matrix for the new predictions
c_m <- table(test$y, pred5)
```

```
# Print the confusion matrix
c_m
```

```
##      pred5
##      no  yes
## no 11030 949
## yes 771 814
```

```

# Calculate the accuracy
acc <- sum(c_m[1], c_m[4])/sum(c_m[1:4])

# Calculate the precision
prec <- c_m[4] / sum(c_m[4], c_m[2])

# Calculate the recall
recall <- c_m[4] / sum(c_m[4], c_m[3])

# Calculate the F1 score
f1_score <- (2*(recall*prec)) / (recall + prec)

#Use the optimal tree and a logistic regression model to classify the test data by using the following principle:
where  $\pi = 0.05, 0.1, 0.15, \dots, 0.9, 0.95$ . Compute the TPR and FPR values for the two models and plot the corresponding ROC curves. Conclusion? Why precisionrecall curve could be a better option here?"

# Train a decision tree model on the training data
optimalTree <- tree(as.factor(y)~., data=train, mindev=0.0005)

# Prune the tree using the optimal number of leaves
optimalTree <- prune.tree(optimalTree, best=21)

# Create a sequence of probability thresholds
pi <- seq(0.05, 0.95, 0.05)

# Train a logistic regression model on the training data
logic_model <- glm(as.factor(y)~., data = train, family="binomial")

# Make predictions on the test set using the logistic regression model
pred6_probY = predict(logic_model, newdata = test, type = "response")

# Calculate the probability of the negative class
pred6_probN = 1 -pred6_probY

# Make predictions on the test set using the pruned decision tree model
tree_pred = predict(optimalTree, newdata = test, type = "vector")

# Initialize vectors for false positive rate and true positive rate for both models
fpr_1 <- c(1:length(pi))
tpr_1 <- c(1:length(pi))
fpr_2 <- c(1:length(pi))
tpr_2 <- c(1:length(pi))

# Initialize the predictions for the decision tree model
pred6 <- tree_pred[, 1]

# Loop through each probability threshold
for (i in 1:length(pi)){

  # Decision tree model
  tpr_1[i] = 0
  fpr_1[i] = 0

  # Create predictions for the decision tree model based on the current threshold
  pred6 <- ifelse(tree_pred[,2]>pi[i], "yes", "no")

  # Calculate the confusion matrix for the decision tree model
  pred6_matrix <- table(test$y, pred6)

  # If the confusion matrix has more than one column, calculate TPR and FPR
  if(ncol(pred6_matrix) > 1){
    tpr_1[i] <- pred6_matrix[2,2] / (pred6_matrix[2,1]+pred6_matrix[2,2])
    fpr_1[i] <- pred6_matrix[1,2] / (pred6_matrix[1,1]+pred6_matrix[1,2])
  }

  # Logistic regression model
  tpr_2[i] = 0
  fpr_2[i] = 0

  # Create predictions for the logistic regression model based on the current threshold
  pred6_logic <- ifelse(pred6_probY > pi[i], "yes", "no")

  # Calculate the confusion matrix for the logistic regression model
  pred6_logic_matrix <- table(test$y, pred6_logic)

  # Calculate TPR and FPR for the logistic regression model
  tpr_2[i] <- pred6_logic_matrix[2,2] / (pred6_logic_matrix[2,1] +pred6_logic_matrix[2,2])

```

```

    fpr_2[i] <- (pred6_logic_matrix[1,2] / (pred6_logic_matrix[1,1]
                                           +pred6_logic_matrix[1,2]))
  }

# Plot the TPR and FPR for the decision tree model
plot(fpr_1, tpr_1, type='l', xlim = c(0,1), ylim = c(0,1),
     xlab='FPR', ylab='TPR', col='red')

# Plot the TPR and FPR for the logistic regression model
lines(fpr_2, tpr_2, type='l', xlim = c(0,1),
     xlab='FPR', ylab='TPR', col='blue')

# Add a reference line to the plot
abline(0,1, lty = 5)

# Add a legend to the plot
legend(x = "bottomright", col = c("red", "blue", "black"),
      legend = c("Tree", "Logistic", "Reference"), lwd = 2, title = "Lines",
      lty = c(1,1,5))

```

*#Scale all variables except of ViolentCrimesPerPop and implement PCA by using function eigen(). Report how many components are needed to obtain at least 95% of variance in the data. What is the proportion of variation explained by each of the first two principal components? "*

```

rm(list = ls(all = TRUE))
graphics.off()
shell("cls")
# Read in the communities data from a CSV file
data = read.csv(file = "communities.csv", header = TRUE)

# Create a logical vector indicating which column is "ViolentCrimesPerPop"
index <- names(data) %in% "ViolentCrimesPerPop"

# Scale the data, excluding "ViolentCrimesPerPop"
data.scaled <- scale(x = data[, !index], center = TRUE, scale = TRUE)

# Calculate the eigenvalues and eigenvectors of the covariance matrix of the original data
e = eigen(cov(data[, -1]))

# Calculate the eigenvalues and eigenvectors of the covariance matrix of the scaled data
e.scaled = eigen(cov(data.scaled))

# Calculate the cumulative sum of the scaled eigenvalues
cum_var = cumsum(e.scaled$values/sum(e.scaled$values))

# Calculate the number of eigenvalues needed to explain 95% of the variance
sum(cum_var<0.95)

```

```
## [1] 34
```

```

# Calculate the proportion of the variance explained by the first two eigenvalues
e.scaled$values[1:2]/sum(e.scaled$values)

```

```
## [1] 0.2501699 0.1693597
```



*#Repeat PCA analysis by using princomp() function and make the trace plot of the first principle component. Do many features have a notable contribution to this component? Report which 5 features contribute mostly (by the absolute value) to the first principle component. Comment whether these features have anything in common and whether they may have a logical relationship to the crime level. Also provide a plot of the PC scores in the coordinates (PC1, PC2) in which the color of the points is given by ViolentCrimesPerPop. Analyse this plot (hint: use ggplot2 package )."*

```
# Read in the communities data from a CSV file
data = read.csv(file = "communities.csv", header = TRUE)

# Create a logical vector indicating which column is "ViolentCrimesPerPop"
index <- names(data) %in% "ViolentCrimesPerPop"

# Scale the data, excluding "ViolentCrimesPerPop"
data.scaled <- scale(x = data[, !index], center = TRUE, scale = TRUE)

# Perform principal component analysis on the scaled data
pr=princomp(data.scaled)

# Calculate the eigenvalues
lambda=pr$sdev^2

# Calculate the proportion of variance explained by each eigenvalue
var = sprintf("%2.3f",lambda/sum(lambda)*100)

# Extract the loadings for the first principal component
ev1 = pr$loadings[,1]

# Find the top five absolute loadings for the first principal component
ev1[order(abs(ev1),decreasing = TRUE)[1:5]]
```

##	medFamInc	medIncome	PctKids2Par	pctWInvInc	PctPopUnderPov
##	-0.1833080	-0.1819830	-0.1755423	-0.1748683	0.1737978

*#Load the ggfortify library for plot visualization*  
**library(ggfortify)**

*# Create a scatterplot matrix of the first two principal components, colored by "ViolentCrimesPerPop"*  
 autoplot(pr, data = data, colour = "ViolentCrimesPerPop")

*#Split the original data into training and test (50/50) and scale both features and response appropriately, and estimate a linear regression model from training data in which ViolentCrimesPerPop is target and all other data columns are features. Compute training and test errors for these data and comment on the quality of model"*

```
# Read in the communities data from a CSV file
df = read.csv("communities.csv")

# Scale the data
df = scale(df, TRUE, TRUE)

# Set a seed for reproducibility
set.seed(12345)

# Split the data into a training set and a test set
n <- dim(df)[1]
id <- sample(1:n,floor(n*0.5))
df_train <- data.frame(df[id,])
df_test <- data.frame(df[-id,])

# Fit a linear regression model to the training data
lr = lm(ViolentCrimesPerPop ~ .,df_train)

# Make predictions on the training and test sets
train.pred = predict(lr, df_train)
test.pred = predict(lr, df_test)

# Calculate the mean squared error (MSE) for the training and test sets
train_MSE = mean((train.pred - df_train$ViolentCrimesPerPop) ^ 2)
test_MSE = mean((test.pred - df_test$ViolentCrimesPerPop) ^ 2)

# Print the MSE for the training and test sets
print("Train error")
```

```
## [1] "Train error"
```

```
train_MSE
```

```
## [1] 0.2591772
```

```
print("Test error")
```

```
## [1] "Test error"
```

```
test_MSE
```

```
## [1] 0.4000579
```

*#Implement a function that depends on parameter vector  $\theta$  and represents the cost function for linear regression without intercept on the training data set. Afterwards, use BFGS method (optim() function without gradient specified) to optimize this cost with starting point  $\theta_0 = 0$  and compute training and test errors for every iteration number. Present a plot showing dependence of both errors on the iteration number and comment which iteration number is optimal according to the early stopping criterion. Compute the training and test error in the optimal model, compare them with results in step 3 and make conclusions. a. Hint 1: don't store parameters from each iteration (otherwise it will take a lot of memory), instead compute and store test errors directly. b. Hint 2: discard some amount of initial iterations, like 500, in your plot to make the dependences visible."*

```
# Initialize empty vectors for storing training and test errors
train_error <- numeric(0)
test_error <- numeric(0)

# Set the seed for reproducibility
set.seed(12345)

# Define the cost function
cost <- function(theta, train, acc_train, test, acc_test){
  # Calculate predictions on the training and test data using theta
  pred_train = train %*% theta
  pred_test = test %*% theta

  # Calculate MSE of predictions on training and test data
  mse_train = mean((acc_train-pred_train)^2)
  mse_test = mean((acc_test - pred_test)^2)

  # Append MSEs to the train_error and test_error vectors
  train_error <- append(train_error,mse_train)
  test_error <- append(test_error,mse_test)

  # Return MSE on training data as cost
  return(mse_train)
}

# Create matrices of predictors and responses for the training and test data
trainy = as.matrix(df_train[,1:(dim(df_train)[2]-1)])
acc_train = as.matrix(df_train['ViolentCrimesPerPop'])
testy = as.matrix(df_test[,1:(dim(df_test)[2]-1)])
acc_test = as.matrix(df_test['ViolentCrimesPerPop'])

# Initialize theta as a matrix of all 0s with the same number of columns as trainy
theta = numeric(dim(trainy)[2])
theta = as.matrix(theta)

# Use the optim function to find the optimal value of theta
opt = optim(par=theta, fn=cost, train = trainy, acc_train = acc_train,
           test=testy, acc_test=acc_test, method = "BFGS")

# Extract the optimal value of theta from the optim function output
opt_theta = opt$par

# Calculate the MSE of the predictions made using opt_theta on the training and test data
train_opt_error = opt$value
test_opt_error = mean((acc_test - (testy %*% opt_theta))^2)

# Print the MSEs
print("calculated optimal train")
```

```
## [1] "calculated optimal train"
```

```
train_opt_error
```

```
## [1] 0.2592247
```

```
print("Lm train error")
```

```
## [1] "Lm train error"
```

```
train_MSE
```

```
## [1] 0.2591772
```

```
print("calculated optimal test")
```

```
## [1] "calculated optimal test"
```

```
test_opt_error
```

```
## [1] 0.3997238
```

```
print("Lm test error")
```

```
## [1] "Lm test error"
```

```
test_MSE
```

```
## [1] 0.4000579
```

```
# Create a logical vector indicating which elements of train_error and test_error should be plotted
excluded = c(TRUE,rep(FALSE,500))

# Extract the relevant elements of train_error and test_error
rest_train = train_error[excluded]
rest_test = test_error[excluded]

# Find the index of the minimum test error
test_min_ind = which(test_error==min(test_error))

print("Early stopping index and MSE")
```

```
## [1] "Early stopping index and MSE"
```

```
test_min_ind
```

```
## [1] 2183
```

```
min(test_error)
```

```
## [1] 0.3769468
```

```
plot(rest_train, xlim=c(0,length(rest_train)), ylim=c(0,1.5), col = "blue")
points(rest_test, col="red")

# Add horizontal lines at the final MSE values for the training and test data
lines(c(0,1000), rep(train_MSE, 2), col="blue")
lines(c(0,1000), rep(test_MSE, 2), col="red")

# Print the early stopping index and minimum test error
print("Early stopping index and MSE")
```

```
## [1] "Early stopping index and MSE"
```

```
test_min_ind
```

```
## [1] 2183
```

```
min(test_error)
```

```
## [1] 0.3769468
```

# lab3

2022-12-08

## ##Task 1

Implement a kernel method to predict the hourly temperatures for a date and place in Sweden. The forecast should consist of the predicted temperatures from 4 am to 24 pm in an interval of 2 hours.

## #Setup

```
set.seed(1234567890)
library(geosphere)
stations <- read.csv("stations.csv")
temps <- read.csv("temps50k.csv")

date <- as.Date("2013-11-04") # The date to predict (up to the students)
filtered_temps <- temps[temps$date < as.Date(date),] #Filter the tempratures to discard irrelevant dates.No time
traverller.

st <- merge(stations,filtered_temps,by="station_number")
# These three values are up to the students
h_distance <- 40000 #Unsure how much we want to consider stations further away.
h_date <- 9
h_time <- 3
a <- 58.4274 # The point to predict (up to the students) #latitud
b <- 14.826 #longitud
pos_vec <- cbind(b,a)

times <- c("04:00:00", "06:00:00", "08:00:00",
           "10:00:00", "12:00:00", "14:00:00",
           "16:00:00", "18:00:00", "20:00:00",
           "22:00:00", "24:00:00")
times_numbers <-c(4,6,8,10,12,14,16,18,20,22,24) #for the plot
temp <- vector(length=length(times)) # length for times: used in forloop.
```

- The first to account for the physical distance from a station to the point of interest. For this purpose, use the function distHaversine from the R package geosphere.
- The second to account for the distance between the day a temperature measurement was made and the day of interest.
- The third to account for the distance between the hour of the day a temperature measurement was made and the hour of interest.



```

#Calculations for relative distances in location, date (day) and time (hours).
# Students' code here

#physical distance with distHaversine
st_loc<-cbind(st$longitude,st$latitude)
dist_hav<-distHaversine(p1=pos_vec,p2=st_loc) #Compare our pos with all the other pos
m<-cbind(dist_hav)

#Gaussian Kernel
#(x_* - x_i)/h from lectures
#diff represents (x_* - x_i)
gaussian_kernel<-function(diff, h_val) {

  u <- diff/h_val
  return(exp(-u*u))
}

#Relative day distance
relative_day_dist <- function(d1,d2){
  diff<- as.Date(d1) - as.Date(d2) #Difference between our date and the date we're comparing (in days).
  return (as.numeric(diff))
}

#relative hour distance, compare our time with another given time. Returns amount of hours.
relative_hour_dist <- function(time1, time2) {
  time_obj1 <-strptime(time1,format="%H:%M:%S") #Create time object so that we can extract hour
  time_obj2 <-strptime(time2,format="%H:%M:%S")
  h1<-as.integer(format(time_obj1,"%H")) #take the hour value as an integer
  h2<-as.integer(format(time_obj2,"%H"))
  # Convert hours to minutes
  minutel <- h1 * 60
  minute2 <- h2 * 60

  # Compute the absolute difference in minutes
  minute_diff <- abs(minutel - minute2)

  # Compute the relative distance in hours
  hour_diff <- minute_diff / 60

  return(hour_diff)
}

```

Use a kernel that is the sum of three Gaussian kernels:

Choose an appropriate smoothing coefficient or width for each of the three kernels above. No cross-validation should be used. Instead, choose manually a width that gives large kernel values to closer points and small values to distant points. Show this with a plot of the kernel value as a function of distance.

#Calculations with sum of kernels.

```

#Calculations
predictions = rep(0,11)
k_distance <- gaussian_kernel(dist_hav,h_distance) #kernel distance
k_days <- gaussian_kernel(relative_day_dist(date,filtered_temps$date),h_date) #kernel days

#Loop over the different times in a day.
for (i in 1:length(temp)) {
  rel_h<-relative_hour_dist(times[i],filtered_temps$time)

  #Relative time difference for each time specified in the assignment
  k_time <- gaussian_kernel(relative_hour_dist(times[i],filtered_temps$time),h_time)

  #Summation of kernels
  k_sum <- cbind(k_distance + k_days + k_time)

  # Normalize the values in the k_sum matrix by dividing each element by the sum of the elements in each row
  k_sum <- (k_sum/sum(k_sum))

  #get weighted temperatures
  weighted_temps <- k_sum * filtered_temps$air_temperature

  predictions[i] <- sum(weighted_temps)
}

#Used to look at h values
plot(rel_h,k_time,main = "h value for time")

```

```

plot(dist_hav,k_distance,xlim=c(0,100000),main = "h value for distance")

```

```

plot(relative_day_dist(date,filtered_temps$date),k_days,xlim=c(0,100),main = "h value for days")

```

The different plots for h-values has been given the following h values: h\_distance <- 40000 h\_date <- 9 h\_time <- 3 This has given us plots that look reasonably good, since closer points are considered more, while points further away will be less and less considered.

```

#sum of kernel plot
plot(times_numbers,predictions,type="o",xlab = "Time of day (hours)",ylab = "Predicted temp",main = "Prediction o
f temperature using summation")

```

These values look very reasonable. Looking at the curve going from colder at the earlier and later hours, while having higher temperatures at the middle of the day. November of 2013 was a very cold month in general, so having a prediction of +6 C is not unreasonable.

Finally, repeat the exercise above by combining the three kernels into one by multiplying them, instead of summing them up. Compare the results obtained in both cases and elaborate on why they may differ.

```

predictions2 = rep(0,11)
k_distance <- gaussian_kernel(dist_hav,h_distance) #kernel distance
k_days <- gaussian_kernel(relative_day_dist(date,filtered_temps$date),h_date) #kernel days

for (i in 1:length(temp)) {
  rel_h<-relative_hour_dist(times[i],filtered_temps$time)

  #Relative time difference for each time specified in the assignment
  k_time <- gaussian_kernel(relative_hour_dist(times[i],filtered_temps$time),h_time)

  #Multiply kernels
  k_sum <- cbind(k_distance * k_days * k_time)

  # Normalize the values in the k_sum matrix by dividing each element by the sum of the elements in each row
  k_sum <- (k_sum/sum(k_sum))

  #get weighted temperatures
  weighted_temps <- k_sum * filtered_temps$air_temperature

  predictions2[i] <- sum(weighted_temps)
}

plot(times_numbers,predictions2,type="o",xlab = "Time of day (hours)",ylab = "Predicted temp",main = "Prediction
of temperature using multiplication")

```

This graph looks a bit more unclear. It's difficult to interpret if there's something wrong with how the data has been divided when multiplying, since the curve is not distributed like a bell curve.

The reason for this graph being less "well distributed" is because the multiplied kernels will give more weight to the higher elements in the kernels, while suppressing the lower ones. This means that values that are not very close to the points we're looking at will not be considered nearly as much (even if they're already not being considered much).

The summed kernel graph gives equal weights to all values which will give our predictions more "real" data to use, since values a bit further away from our initial ones will still be considered. In this case these values are very relevant, since weather predictions (more specifically, temperature predictions) require a lot of data, which should be distributed over time (and not only looking at the exact same day/time/location as what we're trying to predict).

##Task 2

#Error estimates:

```
# Load the kernlab library
library(kernlab)

# Set seed for reproducibility
set.seed(1234567890)

# Load the spam data set
data(spam)

# Shuffle the rows of the data set
foo <- sample(nrow(spam))
spam <- spam[foo,]

# Scale the features (except for the target column)
spam[, -58] <- scale(spam[, -58])

# Split the data into training, validation, and test sets
tr <- spam[1:3000, ]
va <- spam[3001:3800, ]
trva <- spam[1:3800, ]
te <- spam[3801:4601, ]

# Set the step size for the loop
by <- 0.3

# Initialize a vector to store the validation errors
err_va <- NULL

# Loop over different values of the regularization parameter C
for(i in seq(by,5,by)){
  # Train a SVM classifier with the RBF kernel and the current value of C
  filter <- ksvm(type~., data=tr, kernel="rbfdot", kpar=list(sigma=0.05), C=i, scaled=FALSE)

  # Make predictions on the validation set
  mailtype <- predict(filter, va[, -58])

  # Calculate the error on the validation set
  t <- table(mailtype, va[, 58])
  err_va <- c(err_va, (t[1,2]+t[2,1])/sum(t))
}

# Find the value of C that achieved the lowest validation error
best_C <- which.min(err_va) * by

# Retrain the classifier with the optimal value of C
filter0 <- ksvm(type~., data=tr, kernel="rbfdot", kpar=list(sigma=0.05), C=best_C, scaled=FALSE)

# Evaluate the error on the validation set
mailtype <- predict(filter0, va[, -58])
t <- table(mailtype, va[, 58])
err0 <- (t[1,2]+t[2,1])/sum(t)
print("Error 0")
```

```
## [1] "Error 0"
```

```
err0
```

```
## [1] 0.0675
```

```
# Retrain the classifier on the training set and evaluate the error on the test set
filter1 <- ksvm(type~., data=tr, kernel="rbfdot", kpar=list(sigma=0.05), C=best_C, scaled=FALSE)
mailtype <- predict(filter1, te[,58])
t <- table(mailtype, te[,58])
err1 <- (t[1,2]+t[2,1])/sum(t)
print("Error 1")
```

```
## [1] "Error 1"
```

```
err1
```

```
## [1] 0.08489388
```

```
# Retrain the classifier on the combined training and validation sets and evaluate the error on the test set
filter2 <- ksvm(type~., data=trva, kernel="rbfdot", kpar=list(sigma=0.05), C=best_C, scaled=FALSE)
mailtype <- predict(filter2, te[,58])
t <- table(mailtype, te[,58])
err2 <- (t[1,2]+t[2,1])/sum(t)
print("Error 2")
```

```
## [1] "Error 2"
```

```
err2
```

```
## [1] 0.082397
```

```
# Retrain the classifier on the entire data set and evaluate the error on the test set
filter3 <- ksvm(type~., data=spam, kernel="rbfdot", kpar=list(sigma=0.05), C=best_C, scaled=FALSE)
mailtype <- predict(filter3, te[,58])
t <- table(mailtype, te[,58])
err3 <- (t[1,2]+t[2,1])/sum(t)

print("Error 3")
```

```
## [1] "Error 3"
```

```
err3
```

```
## [1] 0.02122347
```

## Questions

1. Which filter do we return to the user ? filter0, filter1, filter2 or filter3? Why?

The filter we return is filter1. Because the SVM model is trained on training data (tr) and validated on validation data (va). Then we should use the optimal  $err\_va$  to make prediction on the test data. Using ksvm on training data and not any other. In general you do your training on your training set, evaluate its performance on the validation set, and then use the best model to predict on the test dataset.

New: The filter we should return is the filter trained on the biggest dataset which is filter3.

2. What is the estimate of the generalization error of the filter returned to the user?  $err_0$ ,  $err_1$ ,  $err_2$  or  $err_3$ ? Why?

it is  $err_1 = 0.08489388$ . Filter0 gives low error estimate due to the fact that it is filtered on validation, the same way we choose the optimal  $err\_va$ . Filter3 gives the lowest error estimate because of the the model is trained on the same data it is predicted on. As for filter1 and filter2, they give similar error estimates. Filter2 error rate is slightly lower due to the fact that it is trained on a larger data set (training and validation), where filter1 is only trained on training data.

New: filter 2 gives an error of 0.082397, this is the nearest estimate we can get of the filter (filter3) returned to the user. Filter3 is expected to be a bit better because of more training data. If we knew that the user wanted to be able to get an error estimate, then it could be worth to instead return filter2.

3. Implementation of SVM predictions. Once a SVM has been fitted to the training data, a new point is essentially classified according to the sign of a linear combination of the kernel function values between the support vectors and the new point. You are asked to implement this linear combination for filter3. You should make use of the functions `alphaindex`, `coef` and `b` that return the indexes of the support vectors, the linear coefficients for the support vectors, and the negative intercept of the linear combination. See the help file of the kernlab package for more information. You can check if your results are correct by comparing them with the output of the function `predict` where you set `type = "decision"`. Do so for the first 10 points in the spam dataset. Feel free to use the template provided in the Lab3Block1 2021 SVMs St.R file.

```

# Get the indices of the support vectors in the SVM classifier
sv <- alphasindex(filter3)[[1]]

# Get the coefficients of the support vectors in the SVM classifier
co <- coef(filter3)[[1]]

# Get the bias term of the SVM classifier and negate it
inte <- -b(filter3)

# Create an RBF kernel with a standard deviation of 0.05
rbfkernel <- rbfdot(sigma = 0.05)

# Initialize an empty vector to store the predicted values
k <- c()

# Loop over the first 10 rows of the spam dataset
for(i in 1:10) {

  # Initialize a variable to store the predicted value for the current row
  k2 <- 0

  # Loop over each support vector
  for(j in 1:length(sv)) {

    # Apply the RBF kernel to the jth support vector and the ith row of the spam dataset. We unlist the matrixes
    # to perform the calculations.
    f <- rbfkernel(unlist(spam[sv[j], -58]), unlist(spam[i, -58]))

    # Update the predicted value for the current row by adding the product of the jth coefficient and the value o
    # f the kernel. F is a 1x1 matrix.
    k2 <- k2 + co[j] * f[1]
  }

  # Append the predicted value for the current row to the k vector
  k <- c(k, k2 + inte)
}

# Print the k vector
k

```

```

## [1] -1.998999  1.560584  1.000278 -1.756815 -2.669577  1.291312 -1.068444
## [8] -1.312493  1.000184 -2.208639

```

```

# Use the predict function to make predictions on the first 10 rows of the spam dataset using the trained SVM cla
ssifier
prediction = predict(filter3, spam[1:10, -58], type = "decision")

plot(k, col = "red")
lines(prediction)

```

The red dotted points is the predicted values that generated from the linear combination for filter3. The lines is the predictions from the predict function.

## Task 3

1:



```

library(neuralnet)
set.seed(1234567890)

# Generate 500 random values between 0 and 10
Var <- runif(500, 0, 10)
mydata <- data.frame(Var, Sin=sin(Var))
tr <- mydata[1:25,] # Training
te <- mydata[26:500,] # Test
winit = runif(31, -1,1) #weights: one for each of the 10 hidden nodes, plus one for the bias term for each of the
10 hidden nodes,
                        #plus one for the output node, plus one for the bias term for the output node

# Fit the neural network using the training data
# The neural network has 10 hidden nodes and uses the initial weights stored in winit
nn <- neuralnet(Sin ~ Var, data = tr,
                hidden = 10, startweights = winit)

# Plot of the training data (black), test data (blue), and predictions (red)
plot(tr, cex=2)
points(te, col = "blue", cex=1)
points(te[,1],predict(nn,te), col="red", cex=1)

```

Question: Train a neural network to learn the trigonometric sine function. To do so, sample 500 points uniformly at random in the interval [0,10]. Apply the sine function to each point. The resulting value pairs are the data points available to you. Use 25 of the 500 points for training and the rest for test. Use one hidden layer with 10 hidden units. You do not need to apply early stopping. Plot the training and test data, and the predictions of the learned NN on the test data. You should get good results. Comment your results.

Answer: The predictions seems pretty accurate. We can see that there is less training data around the points where it kind of diverges from the testdata. this is reasonable and it then finds its way back.

2:

```

library(neuralnet)

# Set seed for reproducibility
set.seed(1234567890)

# Generate data
Var <- runif(500, 0, 10)
mydata <- data.frame(Var, Sin=sin(Var))
tr <- mydata[1:25,] # Training
te <- mydata[26:500,] # Test
winit = runif(31, -1,1)

# Define custom activation functions
h1 <- function(x) x
h2 <- function(x) ifelse(x>0,x,0)
h3 <- function(x) log(1 + exp(x))

# Train neural network with custom activation functions (h1)
nn1 <- neuralnet(Sin ~ Var, data = tr,
                hidden = 10, act.fct = h1, startweights = winit)

# Train neural network with custom activation functions (h2)
nn2 <- neuralnet(Sin ~ Var, data = tr,
                hidden = 10, act.fct = h2, startweights = winit)

# Train neural network with custom activation functions (h3)
nn3 <- neuralnet(Sin ~ Var, data = tr,
                hidden = 10, act.fct = h3, startweights = winit)

# Plot results
plot(tr, cex=3, ylim=c(-1.5, 1.5))
points(te, col = "blue", cex=2)
points(te[,1],predict(nn1,te), col="red", cex=1)
points(te[,1],predict(nn2,te), col="green", cex=1)
points(te[,1],predict(nn3,te), col="pink", cex=1)
legend(1, -0.5, legend=c("train", "test","linear", "ReLU", "Softplus"),
      col=c("black","blue","red", "green", "pink"), lty=1:2, cex=0.8)

```

Question: In question (1), you used the default logistic (a.k.a. sigmoid) activation function, i.e. `act.fct = "logistic"`. Repeat question (1) with the following custom activation functions:  $h_1(x) = x$ ,  $h_2(x) = \max\{0, x\}$  and  $h_3(x) = \ln(1 + \exp x)$  (a.k.a. linear, ReLU and softplus). See the help file of the neuralnet package to learn how to use custom activation functions. Plot and comment your results.

Answer: The third custom activation softplus seems to fit the best. while the first, linear is very bad and the second, ReLU manages predict some values.

3:

```
library(neuralnet)
set.seed(1234567890)
Var <- runif(500, 0, 50)

# Create a data frame with two columns: Var and Sin
# Sin contains the sine of the values in Var
mydata <- data.frame(Var, Sin=sin(Var))

# Plot of the training data (black), test data (blue), and predictions (red)

# Use the nn neural network to make predictions on mydata
prediction = predict(nn,mydata)

plot(mydata, col = "blue", cex=1, ylim=c(-15, 2), xlim=c(0, 55))
points(mydata[,1],prediction, col="red", cex=1, )
```

```
# Find the index of the minimum value in the prediction vector
smallestIndex = which.min(prediction)

# Retrieve the minimum value of the prediction vector using the index
smallestvalue = prediction[smallestIndex]

weights = nn$weights

weights
```

```
## [[1]]
## [[1]][[1]]
##           [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]
## [1,] -11.870831 -0.9153178 -3.0316259  1.5313105  6.555437 -11.768525  1.397325
## [2,]  4.042494 -0.5368677  0.2603355 -0.5487656 -2.233839  1.787862 -1.259834
##           [,8]      [,9]      [,10]
## [1,]  0.2032402  0.2948804 -0.21326487
## [2,] -2.1974036 -0.5886005 -0.03177278
##
## [[1]][[2]]
##           [,1]
## [1,] -0.1088032
## [2,]  0.7716196
## [3,] -1.0436754
## [4,] -16.0529801
## [5,] -1.2718986
## [6,]  3.2805756
## [7,]  4.3076898
## [8,]  0.1220703
## [9,] -1.2234185
## [10,] -2.7618205
## [11,]  2.3696931
```

```
smallestvalue
```

```
## [1] -10.74471
```

Question: Sample 500 points uniformly at random in the interval [0,50], and apply the sine function to each point. Use the NN learned in question (1) to predict the sine function value for these new 500 points. You should get mixed results. Plot and comment your results.

Answer: The predictions seem to still be accurate to the data up to Var 10, and after that it converges to -10.74471 up to Var=50.

4:

Question: In question (3), the predictions seem to converge to some value. Explain why this happens. To answer this question, you may need to get access to the weights of the NN learned. You can do it by running nn or nn\$weights where nn is the NN learned.

Answer: The predictions seem to still be accurate to the data up to Var 10, and after that it converges to -10.74471 at Var=50. The nn is only trained of var values between 0-10 but attempts to predict 0-50. Same for the hidden units. Not sure why it converges to exactly -10.74.

5:

```

# Sample 500 points uniformly at random in the interval [0,10]
Var <- runif(500, 0, 10)

# Apply the sine function to each point
mydata <- data.frame(Var, Sin=sin(Var))
otherWayData <- data.frame(Sin=sin(Var), Var)

# Generate 31 random numbers between 0 and 10
winit = runif(31, 0,10)

# Use all the points in mydata as training points for a neural network
# The target variable is set to be Var instead of Sin
nn <- neuralnet(Var ~ Sin, data = mydata, hidden = 10, threshold = 0.1, startweights = winit )

# Plot the predictions of the neural network
plot(te, col = "blue", cex=1,ylim=c(-10, 10), xlim=c(0, 10))
points(tr[,1],predict(nn,tr), col="red", cex=1)

plot(otherWayData,col = "blue", cex=1) #kanske borde plotta enligt gamla training data istället?
points(otherWayData[,1],predict(nn,otherWayData), col="red", cex=1)

```

Question: Sample 500 points uniformly at random in the interval [0,10], and apply the sine function to each point. Use all these points as training points for learning a NN that tries to predict  $x$  from  $\sin(x)$ , i.e. unlike before when the goal was to predict  $\sin(x)$  from  $x$ . Use the learned NN to predict the training data. You should get bad results. Plot and comment your results. Help: Some people get a convergence error in this question. It can be solved by stopping the training before reaching convergence by setting threshold = 0.1.

Answer: Looks very bad, it is reasonable because we are trying to map from multiple values and not a clear function. One sin value corresponds to multiple var values which makes it impossible to predict. It kinda finds an average.

# Examination

Linköping University, Department of Computer and Information Science, Statistics

---

Course code and name	TDDE01 Machine Learning
Date and time	2022-01-14, 14.00-19.00
Assisting teacher	Oleg Sysoev
Allowed aids	PDF of the course book + your help file (if submitted to LISAM in due time)
Grades:	
	5=18-20 points
	4=14-17 points
	3=10-13 points
	U=0-9 points

---

Provide a detailed report that includes plots, conclusions and interpretations. Give motivated answers to the questions. If an answer is not motivated, the points are reduced. Provide all necessary codes in the appendix.

**Note: seed 12345 should be used in all codes that assumes randomness unless stated otherwise!**

**To start work in RStudio, type this in the Terminal application:**

```
module add courses/TDDE01
rstudio
```

**To submit your report:**

1. Create one file (DOC, DOCX, ODT, PDF)
2. Use Exam Client to submit, and choose Assignment 1 in the drop box
3. Attach your report
4. Submit.
5. "Request Received" status implies that your report is successfully submitted.

## Assignment 1 (10p)

File **adult.csv** shown data collected in a population census in 1994. The following metrics are available:

1. age: continuous.
2. workclass: Private, Self-emp-not-inc, etc.
3. fnlwgt: a population index.
4. education: Bachelors, Some-college, etc.
5. education-num: ordered Education variable.
6. marital-status: Married-civ-spouse, Divorced, etc.
7. occupation: Tech-support, Craft-repair, etc.
8. relationship: Wife, Own-child, etc.
9. capital-gain: continuous. (when you sell an asset and get more for it than you paid originally)
10. capital-loss: continuous.
11. hours-per-week: continuous.
12. native-country: United-States, Cambodia etc.
13. Income level

1. Read data into R with option `stringsAsFactors=T` and divide data randomly (60/20/20) into training, validation and test sets. Use the hold-out principle to first grow a decision tree with default settings and with Income level as target and all other variables except Native Country as inputs. Then use this tree to compute decision trees with various number of leaves and select the optimal tree according to the deviance criterion. Provide the plot of the dependence of the training and validation errors on the number of leaves and interpret this plot from the perspective of bias-variance tradeoff. Report the optimal tree and make some interpretations. Why training a logistic regression model with these features and target is not possible directly? **(4p)**
2. Use the optimal tree computed in step 1 to predict the target probabilities for the test data and estimate predictions for the decision thresholds  $\pi = 0.1, 0.2, \dots, 0.9$  in the rule  $\hat{Y} = '>50K'$  if  $p(Y = '> 50K'|X) > \pi$ , otherwise  $\hat{Y} = '\leq 50K'$ . Use these predictions to report a table showing test accuracies and test F1 scores for various  $\pi$ . Which threshold is the best to use for these data and why? **(3p)**
3. Assume now that hours-per-week is a target variable and  $x_1$ =age,  $x_2$ =Capital gain,  $x_3$ =Capital loss are features. Use cross-validation to estimate a LASSO model with this target and these features from the training data. Report the cross-validation error plot, the optimal penalty factor value and how many features are selected by the model. Report also an equation of the estimated LASSO model. Does the cross-validation plot suggest that having 0 features is significantly better than having 3 features? **(3p)**

## Assignment 2 (10p)

### KERNEL MODELS – 6 POINTS

In the course, you have learned about kernel models for classification and regression. Kernel models can also be used for density estimation, i.e. to model a probability distribution or density function  $p(x_*)$ . In particular,

$$p(x_*) = \frac{1}{n} \sum_{i=1}^n k\left(\frac{x_* - x_i}{h}\right)$$

where the kernel function  $k()$  must integrate to 1. To ensure this, you will hereinafter consider  $k()$  to be the density function of a Gaussian distribution with mean equal to 0 and standard deviation equal to 1. You can get it by using the command `dnorm` in R.

Run the code below to produce some learning data, which consist of 1000 samples from class 1 and 1000 samples from class 2. These points are stored in the variables `data_class1` and `data_class2`. Implement the kernel model presented above to estimate the density function of the data sampled from class 1. Do the same for class 2. Use only 800 samples from class 1 and 800 samples from class 2.

```
set.seed(123456789)

N_class1 <- 1000
N_class2 <- 1000

data_class1 <- NULL
for(i in 1:N_class1){
  a <- rbinom(n = 1, size = 1, prob = 0.3)
  b <- rnorm(n = 1, mean = 15, sd = 3) * a + (1-a) * rnorm(n = 1, mean = 4, sd = 2)
  data_class1 <- c(data_class1,b)
}

data_class2 <- NULL
for(i in 1:N_class2){
  a <- rbinom(n = 1, size = 1, prob = 0.4)
  b <- rnorm(n = 1, mean = 10, sd = 5) * a + (1-a) * rnorm(n = 1, mean = 15, sd = 2)
  data_class2 <- c(data_class2,b)
}
```

Once you have kernel models for the class conditional density functions  $p(x_* | \text{class}=1)$  and  $p(x_* | \text{class}=2)$ , you can use them to produce posterior class probabilities  $p(\text{class} | x_*)$  via Bayes theorem. Specifically,

$$p(\text{class}=1 | x_*) = p(x_* | \text{class}=1) p(\text{class}=1) / [ p(x_* | \text{class}=1) p(\text{class}=1) + p(x_* | \text{class}=2) p(\text{class}=2) ]$$

Use these probabilities to compute the correct classification rate on 200 samples that you did not use before, 100 from class 1 and 100 from class 2. Use this classification rate to select the kernel width  $h$  from among the values 0.1, 0.2, ..., 4.9, 5. Finally, use the 200 samples that you have not used so far to estimate the generalization error of the kernel model selected.



In summary, you should use 1600 samples to build kernel models of the class conditional density functions that you should convert into a probabilistic classifier via Bayes theorem. To select the kernel width, you should use 200 samples as validation set. Finally, you should use 200 samples to estimate the generalization error of the model selected.

#### NEURAL NETWORKS – 4 POINTS

Train a neural network to learn the trigonometric sine function. To do so, sample 50 points uniformly at random in the interval  $[0, 10]$ . Apply the sine function to each point. The resulting pairs are the data available to you. Use 25 of the 50 points for training and the rest for validation. The validation set is used for early stop of the gradient descent. That is, you should use the validation set to detect when to stop the gradient descent and so avoid overfitting (you can read more about early stopping in page 95 of the course textbook). Stop the gradient descent when the partial derivatives of the error function are below a given threshold value. Check the argument `threshold` in the documentation of the R package `neuralnet`. Consider threshold values  $i/1000$  with  $i = 1, \dots, 10$ . Initialize the weights of the neural network to random values in the interval  $[-1, 1]$ . Use a neural network with a single hidden layer of 10 units. **Use the default values for the arguments not mentioned here.** Choose the most appropriate value for the threshold. **Motivate your choice.** Provide the final neural network learned with the chosen threshold. Feel free to use the following template.

```
library(neuralnet)
set.seed(1234567890)

Var <- runif(50, 0, 10)
trva <- data.frame(Var, Sin=sin(Var))
tr <- trva[1:25,] # Training
va <- trva[26:50,] # Validation

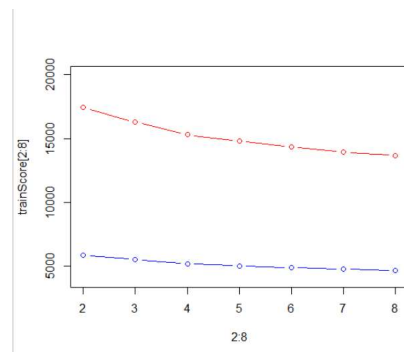
# Random initialization of the weights in the interval [-1, 1]
# Your code here

for(i in 1:10) {
  nn <- neuralnet(# Your code here)

  # Your code here
}

# Your code here
```

## 1.1



Optimal tree with 8 leaves.

Bayes variance: the more leaves, the higher variance but smaller bias.

```
> print(IT)
```

node), split, n, deviance, yval, (yprob)

\* denotes terminal node

```
1) root 19536 21610.00 <=50K ( 0.758344 0.241656 )
2) C8: Not-in-family, Other-relative, Own-child, Unmarried 10666 5224.00 <=50K ( 0.933340
0.066660 )
4) C9 < 7073.5 10479 4190.00 <=50K ( 0.949518 0.050482 )
8) C4: 10th, 11th, 12th, 1st-4th, 5th-6th, 7th-8th, 9th, Assoc-acdm, Assoc-voc, HS-grad,
Preschool, Some-college 8431 2031.00 <=50K ( 0.974024 0.025976 ) *
9) C4: Bachelors, Doctorate, Masters, Prof-school 2048 1741.00 <=50K ( 0.848633 0.151367 ) *
5) C9 > 7073.5 187 46.08 >50K ( 0.026738 0.973262 ) *
3) C8: Husband, Wife 8870 12210.00 <=50K ( 0.547914 0.452086 )
6) C7: ?, Adm-clerical, Craft-repair, Farming-fishing, Handlers-cleaners, Machine-op-inspct, Other-
service, Priv-house-serv, Protective-serv, Transport-moving 4947 5988.00 <=50K ( 0.706489
0.293511 )
12) C9 < 5095.5 4736 5458.00 <=50K ( 0.736909 0.263091 ) *
13) C9 > 5095.5 211 47.30 >50K ( 0.023697 0.976303 ) *
7) C7: Armed-Forces, Exec-managerial, Prof-specialty, Sales, Tech-support 3923 5070.00 >50K (
0.347948 0.652052 )
14) C9 < 5095.5 3412 4592.00 >50K ( 0.399766 0.600234 )
28) C4: 10th, 11th, 12th, 1st-4th, 5th-6th, 7th-8th, 9th, Assoc-acdm, Assoc-voc, HS-grad, Some-
college 1614 2224.00 <=50K ( 0.545229 0.454771 ) *
29) C4: Bachelors, Doctorate, Masters, Prof-school 1798 2094.00 >50K ( 0.269188 0.730812 ) *
15) C9 > 5095.5 511 14.47 >50K ( 0.001957 0.998043 ) *
```

Different interpretations can be done.

- 1) People who luckily sell their assets have high income

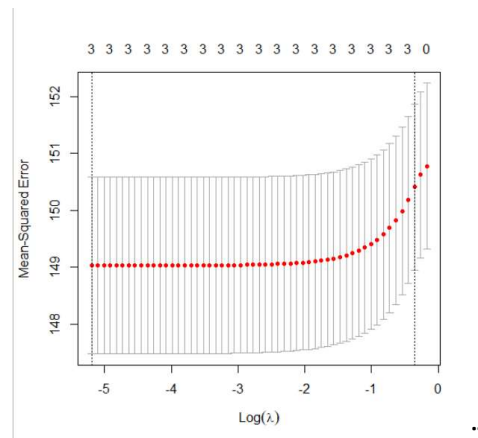
- 2) Married people with good education have high income
- 3) Married with Armed-Forces, Exec-managerial, Prof-specialty, Sales, Tech-support high income

## 1.2

```
> rbind(pis,acc,F1)
      [,1] [,2] [,3] [,4] [,5] [,6] [,7]
pis 0.1000000 0.2000000 0.3000000 0.4000000 0.5000000 0.6000000 0.7000000
acc 0.6549977 0.7299248 0.8410871 0.8410871 0.8481499 0.8481499 0.8481499
F1 0.5712650 0.6126404 0.6580773 0.6580773 0.6013704 0.6013704 0.6013704
      [,8] [,9]
pis 0.8000000 0.9000000
acc 0.8063872 0.8063872
F1 0.3288984 0.3288984
```

Optimal threshold is 0.3 or 0.4, F1 should be prioritized since classes are quite imbalanced.

## 1.3



```
model$lambda.min
[1] 0.005552479
> coef(model, s="lambda.min")
4 x 1 sparse Matrix of class "dgCMatrix"
      1
(Intercept) 38.000806508
C1          0.052593047
C9          0.000107178
C10         0.001783425
3 variables selected, equation y=28+0.052C1+0.0001C9+0.0018C10+epsilon
```

The plot does not suggest that 0 features better than 3 (intervals overlap)

```

# Read in the csv file and store the resulting data frame in `df`
df = read.csv("adult.csv", stringsAsFactors = T)

# Determine the number of rows in `df` and store the result in `n`
n = dim(df)[1]

# Set the seed for the random number generator to 12345
set.seed(12345)

# Take a random sample of `id` indices from 1 to `n`, with the size of the
sample being the floor of `n` multiplied by 0.6
id = sample(1:n, floor(n*0.6))

# Create a new data frame `train` consisting of the rows of `df` at the
indices specified by `id`
train = df[id,]

# Create a new variable `id1` as the set difference between the set of all
indices from 1 to `n` and `id`
id1 = setdiff(1:n, id)

# Take another random sample of indices from `id1`, with the size of the
sample being the floor of `n` multiplied by 0.2
set.seed(12345)
id2 = sample(id1, floor(n*0.2))

# Create a new data frame `valid` consisting of the rows of `df` at the
indices specified by `id2`
valid = df[id2,]

# Create a new variable `id3` as the set difference between `id1` and `id2`
id3 = setdiff(id1, id2)

# Create a new data frame `test` consisting of the rows of `df` at the indices
specified by `id3`
test = df[id3,]

# Load the `tree` library
library(tree)

# Fit a decision tree model with `C13` as the response variable and all other
variables except `C12` as predictors, using the `train` data
fit = tree(C13~.-C12, data=train)

# Print the summary of the fit
summary(fit)

# Store the number of leaves in the tree in `nleaves`
nleaves = 8

# Initialize two vectors of length `nleaves` with all values set to 0
trainScore = rep(0, nleaves)
testScore = rep(0, nleaves)

```

```

# Start a loop that will iterate from 2 to `nleaves`
for(i in 2:nleaves) {
  # Print the current iteration number
  print(i)

  # Prune the tree using `prune.tree` and store the result in `prunedTree`
  prunedTree = prune.tree(fit, best=i)

  # Make predictions on the `valid` data using the pruned tree, and store the
  result in `pred`
  pred = predict(prunedTree, newdata=valid, type="tree")

  # Calculate the deviance of the pruned tree and store the result in
  `trainScore`
  trainScore[i] = deviance(prunedTree)

  # Calculate the deviance of the predictions and store the result in
  `testScore`
  testScore[i] = deviance(pred)
}

# Create a plot showing the deviances in `trainScore` and `testScore`
plot(2:nleaves, trainScore[2:nleaves], type="b", col="red",
ylim=c(4000,20000))
points(2:nleaves, testScore[2:nleaves], type="b", col="blue")

# Identify the tree with the lowest deviance on the test set and prune it
which.min(testScore[-1])+1
lT = prune.tree(fit, best=which.min(testScore[-1])+1))

# Plot the pruned tree
plot(lT)

# Print the pruned tree as text
text(lT)

#2

# Set a range of probability thresholds from 0.1 to 0.9, with a step size of
0.1
pis = seq(0.1, 0.9, 0.1)

# Determine the number of probability thresholds
ps = length(pis)

# Initialize two vectors of length `ps` with all values set to 0
F1 = numeric(ps)
acc = numeric(ps)

# Make predictions on the `test` data using the pruned tree, and store the
result in `Pred`
Pred = predict(lT, newdata=test, type="vector")

# Start a loop that will iterate through the probability thresholds

```

```

for (i in 1:ps){
  # Create a vector `Pr` consisting of the class labels based on the
  probability thresholds
  Pr = ifelse(Pred[,2] > pis[i], ">50K", "<=50K")

  # Create a contingency table comparing the actual and predicted class labels
  tab = table(test$C13, Pr)

  # Calculate the F1 score and store the result in `F1`
  F1[i] = tab[2,2]/(tab[2,2] + 0.5*(tab[1,2] + tab[2,1]))

  # Calculate the accuracy and store the result in `acc`
  acc[i] = (tab[1,1] + tab[2,2]) / (sum(tab))
}

# Print the probability thresholds, accuracy, and F1 score in a table
rbind(pis, acc, F1)

#3

# Extract the first three columns of the `train` data as a matrix and store it
in `covariates`
covariates = as.matrix(train[, c(1,9,10)])

# Extract the 11th column of the `train` data as a vector and store it in
`response`
response = train[[11]]

# Load the `glmnet` library
library(glmnet)

# Fit a Lasso model using `cv.glmnet` on the `covariates` and `response` data,
with an alpha of 1 and a Gaussian distribution
model = cv.glmnet(as.matrix(covariates), response, alpha=1,family="gaussian")

# Print the minimum lambda value
model$lambda.min

# Plot the model results
plot(model)

# Print the coefficients of the model at the minimum
coef(model, s="lambda.min")

```

```
#####
# Kernel methods - TDDE01 - 6p
#####

# Set the seed for the random number generator to 123456789
set.seed(123456789)

# Set the number of data points in each class
N_class1 <- 1000
N_class2 <- 1000

# Initialize an empty vector to store the data for class 1
data_class1 <- NULL

# Start a loop that will iterate `N_class1` times
for(i in 1:N_class1){
  # Generate a random binary value with a probability of 0.3
  a <- rbinom(n = 1, size = 1, prob = 0.3)

  # Generate a random normal value with mean 15 and standard deviation 3 if
  `a` is 1, otherwise generate a random normal value with mean 4 and standard
  deviation 2
  b <- rnorm(n = 1, mean = 15, sd = 3) * a + (1-a) * rnorm(n = 1, mean = 4, sd
  = 2)

  # Append the value to the `data_class1` vector
  data_class1 <- c(data_class1,b)
}

# Initialize an empty vector to store the data for class 2
data_class2 <- NULL

# Start a loop that will iterate `N_class2` times
for(i in 1:N_class2){
  # Generate a random binary value with a probability of 0.4
  a <- rbinom(n = 1, size = 1, prob = 0.4)

  # Generate a random normal value with mean 10 and standard deviation 5 if
  `a` is 1, otherwise generate a random normal value with mean 15 and standard
  deviation 2
  b <- rnorm(n = 1, mean = 10, sd = 5) * a + (1-a) * rnorm(n = 1, mean = 15,
  sd = 2)

  # Append the value to the `data_class2` vector
  data_class2 <- c(data_class2,b)
}

# Estimate the class conditional density functions: 2p.

# Define a function that calculates the conditional probability of a given
value `t` belonging to class 1, given a bandwidth `h`
conditional_class1 <- function(t, h){
  # Initialize a variable to store the sum of the probabilities
```

```

d <- 0

# Start a loop that will iterate 800 times
for(i in 1:800)
  # Add the probability of `t` being within a distance `h` of the `i`-th
  data point in `data_class1`, calculated using a normal distribution
  d <- d + dnorm((t-data_class1[i])/h)

# Return the average probability
return (d/800)
}

# Define a function that calculates the conditional probability of a given
value `t` belonging to class 2, given a bandwidth `h`
conditional_class2 <- function(t, h){
  # Initialize a variable to store the sum of the probabilities
  d <- 0

  # Start a loop that will iterate 800 times
  for(i in 1:800)
    # Add the probability of `t` being within a distance `h` of the `i`-th
    data point in `data_class2`, calculated using a normal distribution
    d <- d + dnorm((t-data_class2[i])/h)

  # Return the average probability
  return (d/800)
}

# Estimate the class posterior probability distribution: 1p.

# Define a function that calculates the posterior probability of a given value
`t` belonging to class 1, given a bandwidth `h`
prob_class1 <- function(t, h){
  # Calculate the conditional probabilities of `t` belonging to each class,
  given the bandwidth `h`
  prob_class1 <- conditional_class1(t,h)*800/1600
  prob_class2 <- conditional_class2(t,h)*800/1600

  # Return the ratio of the probability of `t` belonging to class 1 to the
  total probability of `t` belonging to either class
  return (prob_class1/(prob_class1 + prob_class2))
}

# Select h value via validation: 1p.

# Initialize an empty vector to store the results
foo <- NULL

# Start a loop that will iterate over a sequence of values for the bandwidth
`h`
for(h in seq(0.1,5,0.1)){

```



```

# Calculate the number of values in the second set of 200 values for class 1
for which the posterior probability of belonging to class 1 is greater than
0.5, and the number of values in the second set of 200 values for class 2 for
which the posterior probability of belonging to class 1 is less than 0.5
foo <- c(foo, (sum(prob_class1(data_class1[801:900], h)>0.5)
+sum(prob_class1(data_class2[801:900], h)<0.5))/200)
}

# Plot the results
plot(seq(0.1,5,0.1), foo)

# Print the maximum value of `foo`
max(foo)

# Print the value of `h` at which `foo` is maximum
which(foo==max(foo))*0.1

# Estimate the generalization error: 2p.

# To estimate the generalization error, we use the best h value found
previously.
# Note that the training data is now the old training data union the
validation data.
# Using just the old training data results results in an estimate that is a
bit
# too pessimistic.

# Define a function that calculates the conditional probability of a given
value `t` belonging to class 1, given a bandwidth `h`
conditional_class1 <- function(t, h){
  # Initialize a variable to store the sum of the probabilities
  d <- 0

  # Start a loop that will iterate 900 times
  for(i in 1:900)
    # Add the probability of `t` being within a distance `h` of the `i`-th
data point in `data_class1`, calculated using a normal distribution
    d <- d + dnorm((t-data_class1[i])/h)

  # Return the average probability
  return (d/900)
}

# Define a function that calculates the conditional probability of a given
value `t` belonging to class 2, given a bandwidth `h`
conditional_class2 <- function(t, h){
  # Initialize a variable to store the sum of the probabilities
  d <- 0

  # Start a loop that will iterate 900 times
  for(i in 1:900)
    # Add the probability of `t` being within a distance `h` of the `i`-th
data point in `data_class2`, calculated using a normal distribution
    d <- d + dnorm((t-data_class2[i])/h)

```

```

    # Return the average probability
    return (d/900)
}

# Define a function that calculates the posterior probability of a given value
`t` belonging to class 1, given a bandwidth `h`
prob_class1 <- function(t, h){
  # Calculate the conditional probabilities of `t` belonging to each class,
  given the bandwidth `h`
  prob_class1 <- conditional_class1(t,h)*900/1800
  prob_class2 <- conditional_class2(t,h)*900/1800

  # Return the ratio of the probability of `t` belonging to class 1 to the
  total probability of `t` belonging to either class
  return (prob_class1/(prob_class1 + prob_class2))
}

# Find the optimal value of `h` by selecting the value that maximizes the
value of `foo`
h <- which(foo==max(foo))*0.1

# Calculate the number of values in the third set of 200 values for each class
(901 to 1000) for which the posterior probability of belonging to class 1,
given `h`, is above or below the threshold (0.5)
(sum(prob_class1(data_class1[901:1000], h)>0.5)
+sum(prob_class1(data_class2[901:1000], h)<0.5))/200

#####
# Neural networks - TDDE01 - 4p
#####

library(neuralnet)
set.seed(1234567890)

Var <- runif(50, 0, 10)
trva <- data.frame(Var, Sin=sin(Var))
tr <- trva[1:25,] # Training
va <- trva[26:50,] # Validation

restr <- vector(length = 10)
resva <- vector(length = 10)
winit <- runif(31, -1, 1) # Random initializaiton of the weights in the
interval [-1, 1]
for(i in 1:10) {
  nn <- neuralnet(formula = Sin ~ Var, data = tr, hidden = 10, startweights =
winit,
                  threshold = i/1000, lifesign = "full")

  aux <- predict(nn, tr) # Compute predictions for the trainig set and their
squared error
  restr[i] <- sum((tr[,2] - aux)**2)/2

  aux <- predict(nn, va) # The same for the validation set

```

```

    resva[i] <- sum((va[,2] - aux)**2)/2
  }
plot(restr, type = "o")
plot(resva, type = "o")
restr
resva

# The graphs show an example of overfitting, i.e. the threshold that achieves
the lowest squared error
# in the training set is not the one that achieves the lowest error in the
validation set. Therefore,
# early stopping is necessary, i.e. running gradient descent until convergence
is not the best option,
# as the lowest threshold gives the best error in the training set but not in
the validation set.
# Specifically, the validation set indicates that gradient descent should be
stoped when
# threshold = 4/1000. So, the output should be a NN learnt with all (!) the
data available and the
# threshold = 4/1000.

winit <- runif(31, -1, 1)
plot(nn <- neuralnet(formula = Sin ~ Var, data = trva, hidden = 10,
startweights = winit,
                    threshold = 4/1000, lifesign = "full"))

# Plot of the predictions (blue dots) and the data available (red dots)

plot(trva[,1],predict(nn,trva), col="blue", cex=3)
points(trva, col = "red", cex=3)

```