

# Lab 1 Report

## Group B33

2023/11/19

Group Members: Priyansh Gupta, Afshan Hashemi Hosseinabad, Ragini Gunda

Statement of contribution: Assignment 1 is mainly done by Ragini Gunda, Priyansh Gupta was mainly responsible for assignment 2 and, assignment 3 is mostly done by Afshan Hashemi Hosseinabad. Results from all assignments are discussed afterwards between all of us and the group report is formulated based on the discussion.

### Assignment 1: Handwritten digit recognition with K-nearest neighbours.

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

```
1. data <- read.csv("optdigits.csv",header = FALSE)
2. data$V65 <- as.factor(data$V65)
3. num <- dim(data)[1]
4. set.seed(12345)
5. id = sample(1:num, floor(num*0.5))
6. trainData <- data[id,]
7. id_1 <- setdiff(1:num,id)
8. set.seed(12345)
9. id_2 <- sample(id_1,floor(num*0.25))
10. validationData <- data[id_2,]
11. id_3 <- setdiff(id_1,id_2)
12. testData <- data[id_3,]
```

2. Use training data to fit 30-nearest neighbor classifier with function *kknn()* and kernel="rectangular" from package *kknn* and estimate

#### a) Confusion Matrices:

```
1. cm_1 <- table(trainData$V65,confmat_train)
2. print("Confusion matrix for the training data")
3. cm_1
4.
```

Confusion matrix for the training data:

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

```
1. cm_2 <- table(testData$V65,confmat_test)
2. print("Confusion matrix for the test data")
3. cm_2
```

Confusion matrix for the testing data:

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

### Comment on the quality of predictions for different digits:

It can be seen that model succeeded best on predicting digits 0 to fit the input image to the correct output in the training set, the quality of the model is as follows:

```
1. diagonal_values <- diag(cm_1)
2. print(diagonal_values)
3. sorted_indices <- order(diagonal_values,decreasing = TRUE)
4. print(sorted_indices)
5. sorted_values <- diagonal_values[sorted_indices]
6. print(sorted_values)
```

```
> print(diagonal_values)
 0    1    2    3    4    5    6    7    8    9
202 179 190 185 159 171 190 178 188 183
> sorted_indices <- order(diagonal_values,decreasing = TRUE)
> print(sorted_indices)
[1] 1 3 7 9 4 10 2 8 6 5
> sorted_values <- diagonal_values[sorted_indices]
>
> print(sorted_values)
 0    2    6    8    3    9    1    7    5    4
202 190 190 188 185 183 179 178 171 159
```

The model is succeeded on predicting digit 7 in the test set. The quality of predictions for different digits is:

```
1. diagonal_values_test <- diag(cm_2)
2. print(diagonal_values_test)
3. sorted_indices_test <- order(diagonal_values_test,decreasing = TRUE)
4. print(sorted_indices_test)
5. sorted_values_test <- diagonal_values_test[sorted_indices_test]
6. print(sorted_values_test)
```

```
> diagonal_values_test <- diag(cm_2)
> print(diagonal_values_test)
 0    1    2    3    4    5    6    7    8    9
77  81  98 107  94  93  90 111  70  85
> sorted_indices_test <- order(diagonal_values_test,decreasing = TRUE)
> print(sorted_indices_test)
[1] 8 4 3 5 6 7 10 2 1 9
> sorted_values_test <- diagonal_values_test[sorted_indices_test]
> print(sorted_values_test)
 7    3    2    4    5    6    9    1    0    8
111 107  98  94  93  90  85  81  77  70
```

b). Miss-classification errors for the training and test data:

```
1. miss_class = function(y1,y2){
2.   n=length(y1)
3.   return(1-sum(diag(table(y1,y2)))/n)
4. }
5. print("Miss-classification of training data")
6. miss_class(trainData$V65,confmat_train)
```

```
[1] "Miss-classification of training data"
> miss_class(trainData$V65,confmat_train)
[1] 0.04500262
```

```
1. print("Miss-classification of test data")
2. miss_class(testData$V65,confmat_test)
```

```
[1] "Miss-classification of test data"
> miss_class(testData$V65,confmat_test)
[1] 0.05329154
```

```
1. accuracy = function(y1,y2)
2. {
3.   n=length(y1)
4.   return(sum(diag(table(y1,y2)))/n)
5. }
6. }
7. print("Accuracy of the training data")
8. accuracy(trainData$V65,confmat_train)
9.
```

```
[1] "Accuracy of the training data"
> accuracy(trainData$V65,confmat_train)
[1] 0.9549974
```

```
1. print("Accuracy of the test data")
2. accuracy(testData$V65,confmat_test)
```

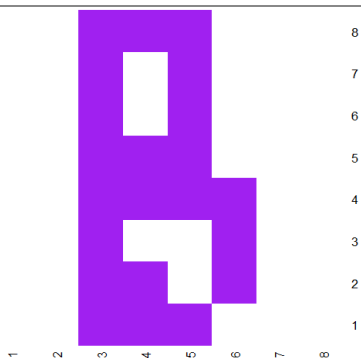
```
[1] "Accuracy of the test data"
> accuracy(testData$V65,confmat_test)
[1] 0.9467085
```

3. 3 Hardest and 2 easiest cases to classify:

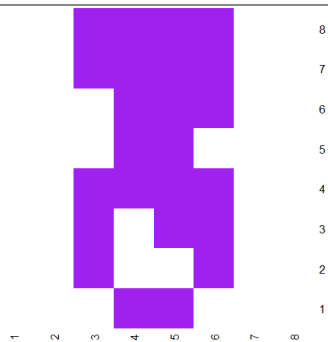
```
1. trainData_new <- trainData
2. trainData_new$prob <- train_kknn$prob[, "8"]
3. df <- data.frame(trainData_new)
4. hardest <- df[order(df$prob), ] %>% head(3) # 3 hardest cases to classify
5. easiest <- df[rev(order(df$prob)), ] %>% head(2) #2 easiest cases to classify
```

The easiest case:

```
1. heatmap(matrix(as.numeric(easiest[1,1:64]),byrow = TRUE,nrow = 8,ncol = 8),Colv = NA,Rowv =
NA,col=c("white","purple"))
```

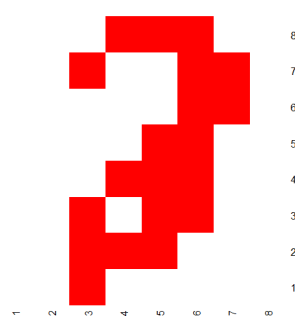
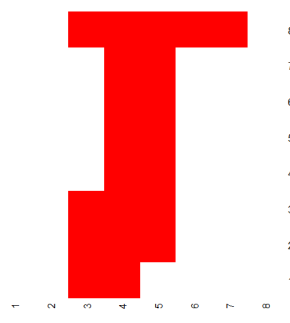
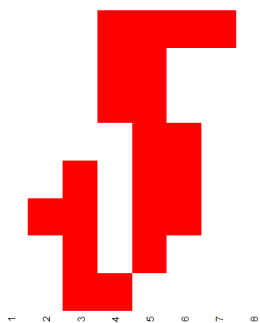


```
1. heatmap(matrix(as.numeric(easiest[2,1:64]),byrow = TRUE,nrow = 8,ncol = 8),Colv = NA,Rowv = NA,col=c("white","purple"))
2.
```



The Hardest case:

```
1. heatmap(matrix(as.numeric(hardest[1,1:64]),byrow = TRUE,nrow = 8,ncol = 8),Colv = NA,Rowv = NA,col=c("white","red"))
2.
3. heatmap(matrix(as.numeric(hardest[2,1:64]),byrow = TRUE,nrow = 8,ncol = 8),Colv = NA,Rowv = NA,col=c("white","red"))
4.
5. heatmap(matrix(as.numeric(hardest[3,1:64]),byrow = TRUE,nrow = 8,ncol = 8),Colv = NA,Rowv = NA,col=c("white","red"))
6.
```



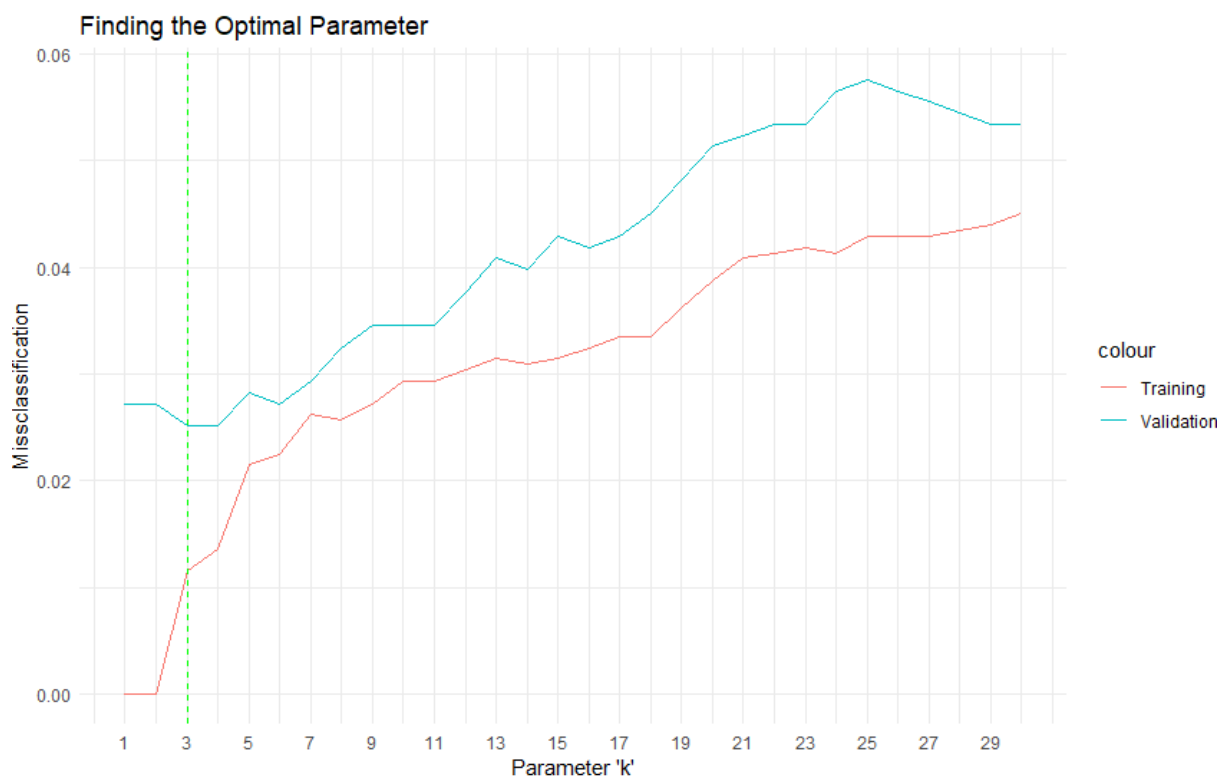
The heatmap for the easiest-to-classify cases shows clear and recognizable patterns of the digit '8.' The features are well-aligned with the typical structure of the digit, contributing to the model's accurate and confident predictions where as the heatmaps for the hardest-to-classify cases exhibit more ambiguous patterns. It appears challenging to visually discern the digit '8' in these instances, which aligns with the model's lower confidence in its predictions."

#### 4. Fitting a K-nearest neighbor classifiers to the training data for different values of $k=1,2,\dots,30$ .

```
1. missclass_training <- c()
2. missclass_validation <- c()
3. for(k in 1:30){
4.   knn_training <- knn(formula = V65~.,kernel = "rectangular",train =trainData,test = trainData,k=k
5. )
6.   knn_validation <- knn(formula = V65~.,kernel = "rectangular",train =trainData,test =
7. validationData,k=k)
8.   missclass_training <- c(missclass_training,miss_class(trainData$V65,knn_training$fitted.values))
9.   missclass_validation <-
10. c(missclass_validation,miss_class(validationData$V65,knn_validation$fitted.values))
11. }
12. }
```

Plotting the dependence of the training and validation misclassification errors on the value of K:

```
1. k<-1:30
2. df1<-data.frame(missclass_training,k1 =k)
3. df2<-data.frame(missclass_validation,k2 =k)
4. df3 <- data.frame(df1,df2)
5. ggplot( ) +
6.   geom_line(aes(x=df3$k1,y=df3$missclass_training,
7.                 colour="Training")) +
8.   geom_line(aes(x=df3$k2,y=df3$missclass_validation,
9.                 colour="Validation")) +
10.  theme_minimal()+ggtitle("Finding the Optimal Parameter") +
11.  scale_x_continuous(breaks = seq(1,30,2)) +
12.  xlab("Parameter 'k'") + ylab("Missclassification") +
13.  geom_vline(xintercept = 3
14.             , linetype = "dashed",color = "green")
```



According to the plot, it can be seen that the optimal parameter for our model is  $k = 3$  in which has the lowest Missclassification Error. In addition, when the model is tested with  $k = 3$ , we get the result as:

```
1. knn_test <- knn(formula = V65~.,kernel = "rectangular",train =trainData,test = testData,k=3 )
2. cat("for k=3 :\n")
3. cat("MissClassification error on Training =",missclass_training[3],"\n")
```

```

4. cat("MissClassification error on Validation = ",missclass_validation[3],"\\n")
5. cat("MissClassification error on Test = ",
6.     miss_class(testData$V65,kknn_test$fitted.values),"\\n")

```

```

> cat("for k=3 :\\n")
for k=3 :
> cat("MissClassification error on Training =",missclass_training[3],"\\n")
MissClassification error on Training = 0.0115123
> cat("MissClassification error on Validation = ",missclass_validation[3],"\\n")
MissClassification error on Validation = 0.02513089
> cat("MissClassification error on Test = ",
+     miss_class(testData$V65,kknn_test$fitted.values),"\\n")
MissClassification error on Test = 0.02403344

```

So, the miss-classification Error rate has the lowest value in the training set as expected. The miss-classification error is nearly equal on the Validation set, where the Test set shows that the model does not overfit. According to the plot, it can be perceived that the complexity of the model increases as the parameter  $k$  increases. Finally, we find that the optimal  $k$  value is 3.

In general, the complexity of the model increases as well as the variance in which means that when the complexity is high the observed variance is also high. When there is model complexity, we deal with a model with high bias. Our aim is to minimize these 2 variables so that there is always a trade-off between bias and variance. For  $k = 3$ ; the bias is very small since the miss-classification error has the lowest and there is a small difference between the miss-classification Error on the training set and validation. That is why choosing  $k = 3$  is the most suitable option.

## 5. Fit K-nearest neighbor classifiers to training data for different values of $k = 1, 2, \dots, 30$ and compute the empirical risk for the validation data as cross-entropy.

```

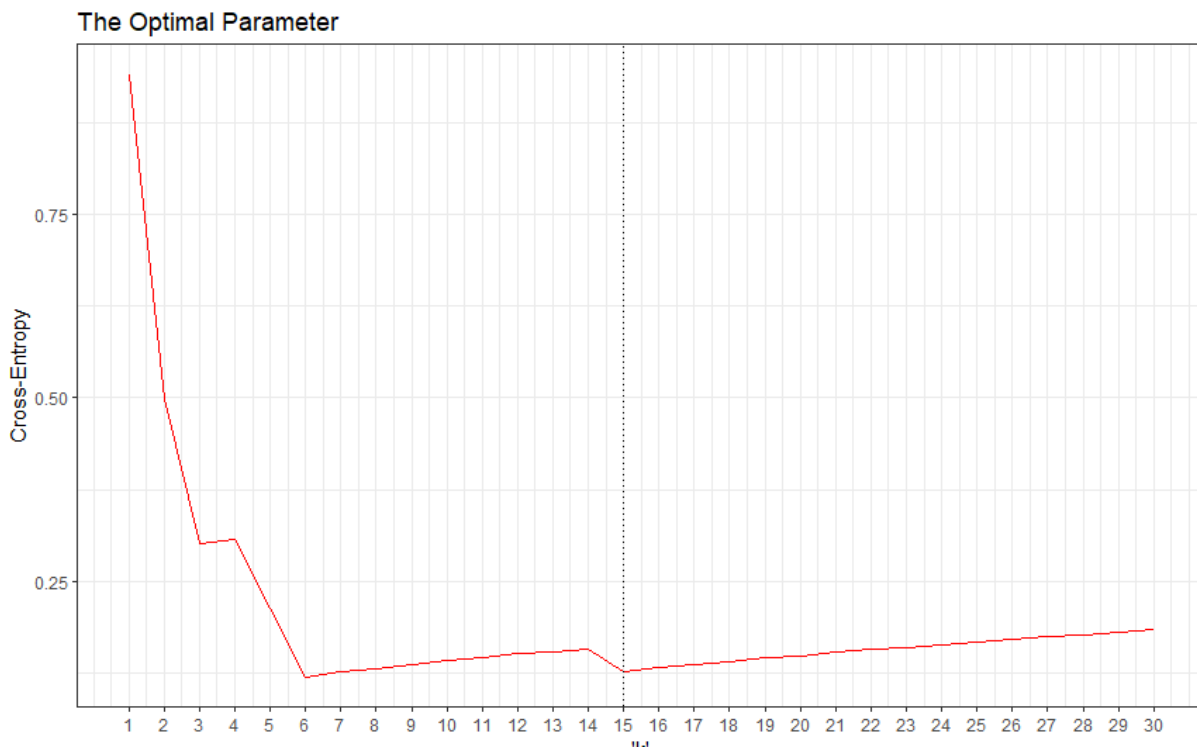
1. library(data.table)
2. library(kknn)
3. crossEntropy_training <- c()
4. crossEntropy_validation <- c()
5. noOfClasses <- length(unique(trainData$V65))
6. noOfRows_training <- nrow(trainData)
7. noOfRows_validation <- nrow(validationData)
8. training_onehot <- t(sapply(as.numeric(trainData$V65),function(x)
9.   {c(rep(0,x-1), 1 ,rep(0 , noOfClasses-x ))}))
10.
11. validation_onehot <- t(sapply(as.numeric(validationData$V65),function(x)
12.   {c(rep(0,x-1), 1 ,rep(0 , noOfClasses-x ))}))
13. for(i in 1:30)
14. {
15.   fitModel_training <- kknn(V65~.,kernel = "rectangular", train = trainData, test = trainData, k=i)
16.   fitModel_validation <- kknn(V65~.,kernel = "rectangular", train = trainData, test =
validationData, k=i)
17.
18.   #cross-entropy
19.   crossEntropy_training <- c(crossEntropy_training,sum(training_onehot * -log(fitModel_training$prob
+ 10^- 15))/noOfRows_training)
20.
21.   crossEntropy_validation <- c(crossEntropy_validation,sum(validation_onehot * -
log(fitModel_validation$prob + 10^- 15))/noOfRows_validation)
22.
23. }
24. crossEntropy_validation

```

```

25.
26. crossEntropy <- melt(data.table(k = 1:30, Validation = crossEntropy_validation), "k", variable.name =
  " ")
27. ggplot(crossEntropy) + geom_line(aes(k, value), col="red") + theme_bw() +
28.   geom_vline(xintercept = 15, linetype = "dotted") + scale_x_continuous(breaks = seq(1,30)) +
29.   xlab("'k'") + ylab("Cross-Entropy") +
30.   ggtitle("The Optimal Parameter")
31.

```



The plot represents that  $k = 6$  gives the lowest empirical risk error, so it can be defined as the optimal parameter for our model. From a probabilistic overview, miss-classification error is derived through Linear regression, while Cross-Entropy is from Logistic Regression. We want to classify the discrete return in which choosing Cross-Entropy is the suitable option. Moreover, unlikely to Misclassification Error, Cross Entropy penalizes the model more for wrong classification since it's derivation is higher.

## Assignment 2: Linear regression and Ridge regression

1. 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.

```
1. parkinsons = read.csv("parkinsons.csv")
2. parkinson_data = as.data.frame(parkinsons[,c(5,7:22)])
3.
4. n = dim(parkinson_data)[1]
5. set.seed(12345)
6. id = sample(1:n, floor(n*0.6))
7. traindata = parkinson_data[id,]
8. testdata = parkinson_data[-id,]
9. scaler <- preProcess(as.data.frame(traindata))
10. train_s <- predict(scaler,traindata)
11. test_s <- predict(scaler,testdata)
```

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.

```
1. lrm <- lm(motor_UPDRS ~ .-1, data = train_s)
2. train_predicted <- predict(lrm, newdata = train_s )
3. test_predicted <- predict(lrm, newdata = test_s)
4. MSETrain <- mean((train_predicted - train_s$motor_UPDRS)^2)
5. MSETest <- mean((test_predicted - test_s$motor_UPDRS)^2)
6. (c(MSE_of_train_data = MSETrain, MSE_of_test_data = MSETest ))
7. summary(lrm)
```

Output of the above code is:

```
(c(MSE_of_train_data = MSETrain, MSE_of_test_data = MSETest ))
MSE_of_train_data MSE_of_test_data
0.8785431         0.9354477
> summary(lrm)
Call:
lm(formula = motor_UPDRS ~ . - 1, data = train_s)
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
```

As we can see from summary variables Jitter.Abs, Shimmer.APQ5, Shimmer.APQ11, NHR, HNR, DFA, PPE contribute significantly to model with 0.001 significance level.



3. Implement 4 following basic functions by using basic R commands only (not using any external packages):

**a. loglikelihood function**

```
1. loglikelihood <- function(parameters) {
2.   theta <- parameters[1:16]
3.   sigma <- parameters[17]
4.   y <- as.matrix(train_s[, -1]) %*% theta
5.   n <- nrow(train_s)
6.   return (-n/2*log(2*pi) -n/2*log(sigma^2)-sum((train_s[,1]-y)^2)/(2*sigma^2))
7. }
```

**b. Ridge function**

```
1. ridgePenalty <- function(parameters, lambda){
2.   return(-loglikelihood(parameters)+ lambda*norm(as.matrix(parameters[1:16]),type="2")^2)
3. }
```

**c. RidgeOpt function**

```
1. ridgeOpt <- function(lambda){
2.   return (optim( par = c(rep(0,length=16),0.01),fn = ridgePenalty, lambda = lambda, method = "BFGS"))
3. }
```

**d. DF function**

```
1. df <- function(lambda){
2.   X <- as.matrix(train_s[, -1])
3.   P <- X %*% solve(t(X)%*%X + lambda*diag(16))%*%t(X)
4.   return(sum(diag(P)))
5. }
```

4. By using the RidgeOpt() method compute the optimal parameter  $\theta$  for  $\lambda = 1, 10, 1000$ .

```
1. MSE_Matrix<-matrix(c(MSE_Train_1,MSE_Test_1,MSE_Train_100,MSE_Test_100,
MSE_Train_1000,MSE_Test_1000), nrow =2, ncol = 3)
2. rownames(MSE_Matrix) <-c("Train_MSE", 'Test_MSE')
3. colnames(MSE_Matrix) <-c("λ=1", "λ=100", "λ=1000")
4. print(MSE_Matrix)
```

Output is as follow:

```
print(MSE_Matrix)
           λ=1      λ=100      λ=1000
Train_MSE 0.8786272 0.8844091 0.9211176
Test_MSE  0.9349961 0.9323369 0.9539460
```

- Which penalty parameter is most appropriate among the selected one?

As MSE value for testing data is lowest among for  $\lambda = 100$  among all the selected values. It is the most appropriate value for penalty parameter.

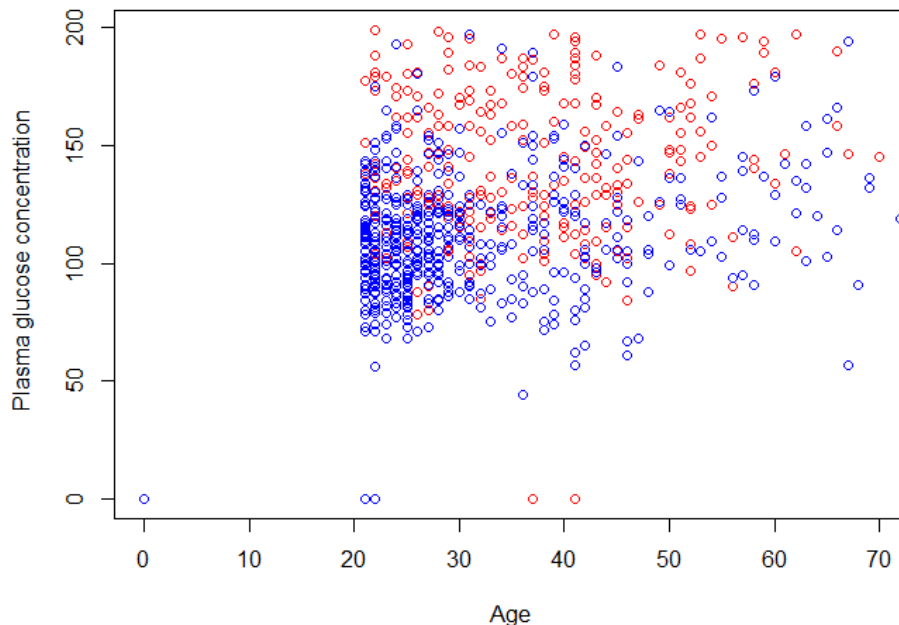
**compute and compare the degree of freedom of these models and make appropriate conclusion.**

```
1. c(DF_1 = df(1),DF_100 = df(100),DF_1000 = df(1000))
2.      DF_1      DF_100      DF_1000
3. 13.860736  9.924887  5.643925
```

As it can be clearly seen from the result DF is inversely propositional to penalty parameter  $\lambda$ , i.e., df decreases as penalty parameter increases.

## Assignment 3: Logistic regression and basic function

1. It is not easy to classify diabetes only using Plasma glucose concentration and Age since as we see in the scatterplot the points are not well separated and it is hard to draw a line or a curve between two groups that divide them nicely.



2. In this assignment we have divided the data 50/50 and used the code below to train our logistic models on the "train" data set and get predictions for the "test" data set

```
1) logistic=glm(Diabetes ~ Age + Plasma,data = train,family = "binomial")
2) LRpred=predict(logistic,newdata = test,type="response")
3) cutoff=0.5
4) LRpred[LRpred<cutoff]=0
5) LRpred[LRpred>=cutoff]=1
6) summary(logistic)
```

As we run the above code we will get the following results:

```
Call:
glm(formula = Diabetes ~ Age + Plasma, family = "binomial", data = train)

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept) -6.047850   0.664918  -9.096 < 2e-16 ***
Age           0.047579   0.010939   4.349 1.37e-05 ***
Plasma        0.031033   0.004524   6.860 6.88e-12 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 497.98  on 383  degrees of freedom
Residual deviance: 399.86  on 381  degrees of freedom
```

AIC: 405.86

Number of Fisher Scoring iterations: 4

The coefficients values for both age and plasma are fairly small so we have chosen good parameters. Based on the coefficient values we get the following equation:

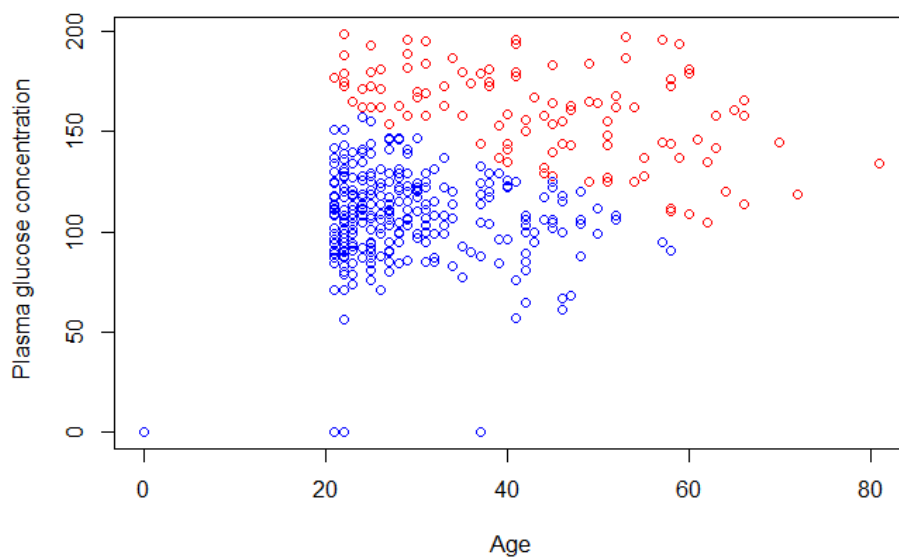
$$P(\text{Diabetes}) = \text{sigm}(-6.047 + 0.047 \cdot \text{Age} + 0.031 \cdot \text{Plasma})$$

To calculate misclassification rate we should:

```
1) confusionMX=table(LRpred,test$Diabetes)
2) misclassificationRate=(1 - sum(diag(confusionMX))/sum(confusionMX))
```

Which outputs: 0.255

The scatterplot for the test data:



Unlike step one in this plot the blue and red pints can be divided easily so the quality of the classification is good.

3.

A) The equation of the decision boundary can be calculated using coefficients in order to that we need to calculate the slope and intercept of the line:

```
1. coeff=logistic$coefficients
2. slope=coeff[2]/(-1*coeff[3])
3. intercept=coeff[1]/(-1*coeff[3])
```

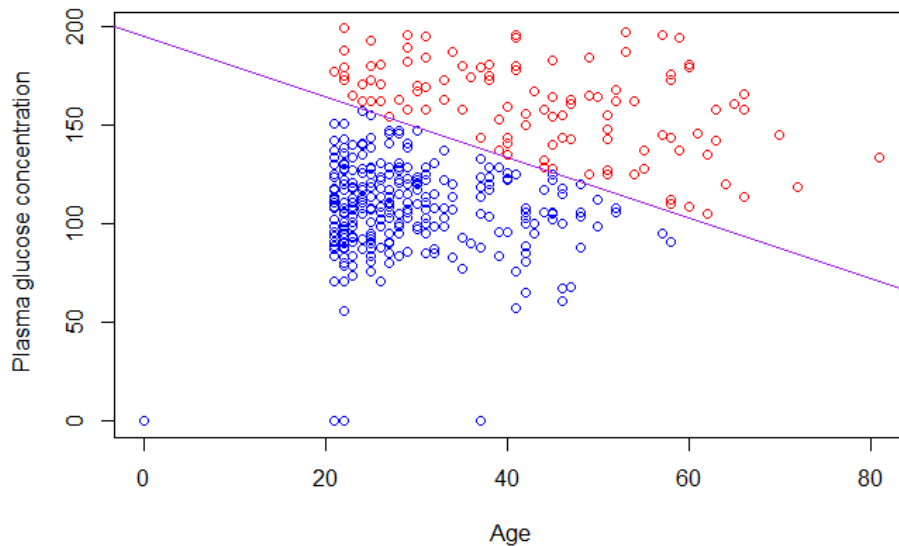
Using this calculation *slope* = -1.53 and *intercept* = 195 so we have:

$$y = 195 - x(1.53)$$

B) We can add a line in the data using the following code:

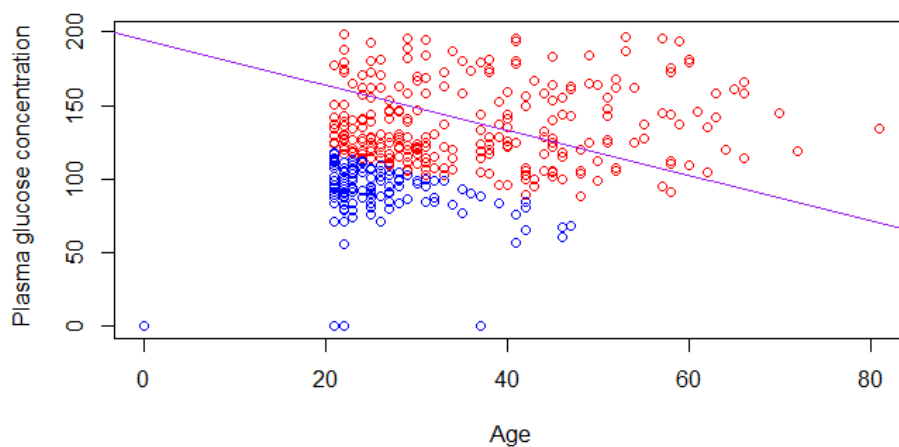
```
1) abline(intercept,slope,col="purple")
```

We will get the following plot:

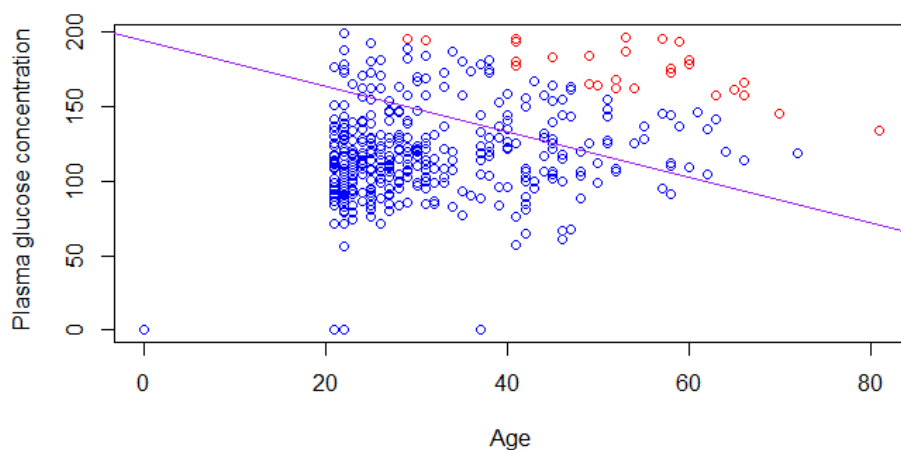


As we can see the two sides of the purple line has different colours, so the decision boundary has cached the data distribution well.

4. If we set  $r = 0.2$  we will get the plot below and as we can see the model is not very accurate since there are some patients without diabetes that are classified as having it.



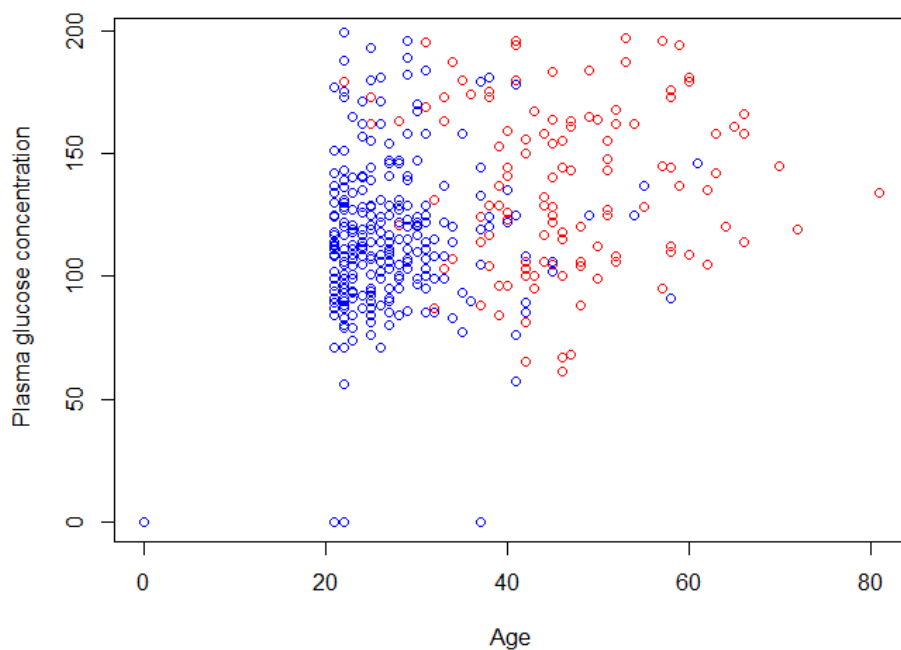
If we set  $r = 0.8$ , this is even worse than the previous case in medical terms since there are people with diabetes that are classified as non-diabetes, and this can endanger their life.



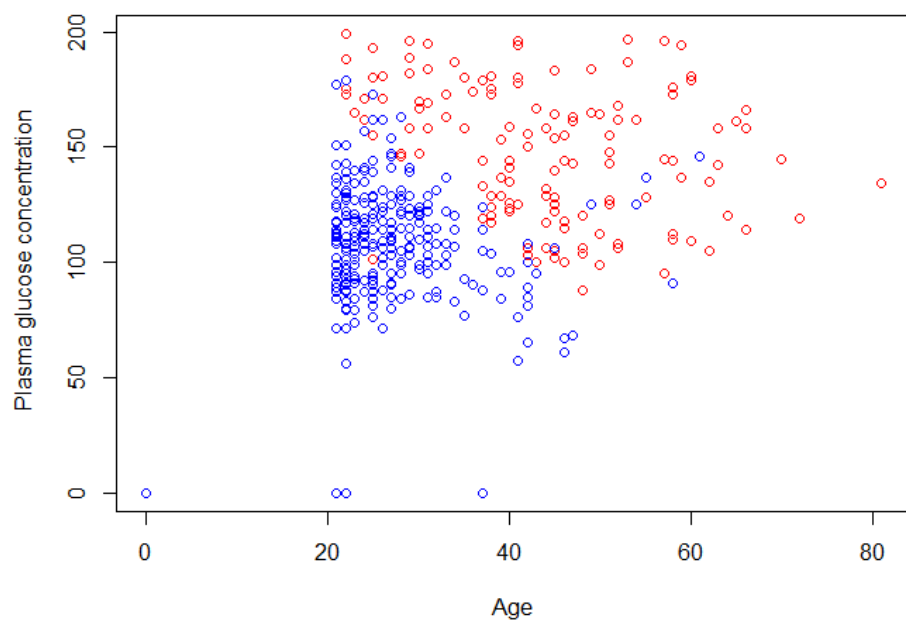
5. We can use the following code to add new binomial features to train our model with:

```
1) train$phi1=phi1=train$Plasma^4
2) train$phi2=phi2=train$Plasma^3*train$Age
3) train$phi3=phi3=train$Plasma^2*train$Age^1
4) train$phi4=phi4=train$Plasma^1*train$Age^3
5) train$phi5=phi5=train$Age^4
```

As we train the model with the new features, we will get the following plot:



Misclassification rate in this model is 0.372 so, it does not work well and does not have a good accuracy compared to the model in part 2 but it has more of a curve compared to the previous decision boundary. But we can improve it by using more variables. For instance, we can use 7 variables, and we get the plot below with 0.299 misclassification rate.



## Appendix:

### Code for assignment 1:

```
1. #install.packages("tidyverse")
2. library(tidyverse)
3. library(kknn)
4. library(dplyr)
5.
6. data <- read.csv("optdigits.csv",header = FALSE)
7. data$V65 <- as.factor(data$V65)
8. num <- dim(data)[1]
9. set.seed(12345)
10. id = sample(1:num, floor(num*0.5))
11. trainData <- data[id,]
12. id_1 <- setdiff(1:num,id)
13. set.seed(12345)
14. id_2 <- sample(id_1,floor(num*0.25))
15. validationData <- data[id_2,]
16. id_3 <- setdiff(id_1,id_2)
17. testData <- data[id_3,]
18.
19. #install.packages("kknn")
20.
21. train_kknn <- kknn(V65~., train = trainData, test = trainData, kernel = "rectangular" , k=30)
22. test_kknn <- kknn(V65~., train = trainData, test = testData, kernel="rectangular",k=30)
23. confmat_train <- train_kknn$fitted.values
24. confmat_test <- test_kknn$fitted.values
25. cm_1 <- table(trainData$V65,confmat_train)
26. print("Confusion matrix for the training data")
27. cm_1
28. cm_2 <- table(testData$V65,confmat_test)
29. print("Confusion matrix for the test data")
30. cm_2
31. diagonal_values <- diag(cm_1)
32. print(diagonal_values)
33. sorted_indices <- order(diagonal_values,decreasing = TRUE)
34. print(sorted_indices)
35. sorted_values <- diagonal_values[sorted_indices]
36.
37. print(sorted_values)
38.
39. diagonal_values_test <- diag(cm_2)
40. print(diagonal_values_test)
41. sorted_indices_test <- order(diagonal_values_test,decreasing = TRUE)
42. print(sorted_indices_test)
43. sorted_values_test <- diagonal_values_test[sorted_indices_test]
44.
45. print(sorted_values_test)
46.
47. miss_class = function(y1,y2){
48.   n=length(y1)
49.   return(1-sum(diag(table(y1,y2)))/n)
50. }
51. print("Miss-classification of training data")
52. miss_class(trainData$V65,confmat_train)
53.
54. print("Miss-classification of test data")
55. miss_class(testData$V65,confmat_test)
56.
57. accuracy = function(y1,y2)
58. {
59.   n=length(y1)
60.   return(sum(diag(table(y1,y2)))/n)
61. }
62.
63. print("Accuracy of the training data")
64. accuracy(trainData$V65,confmat_train)
```

```

65. print("Accuracy of the test data")
66. accuracy(testData$V65,confmat_test)
67. #install.packages("dbplyr")
68. trainData_new <- trainData
69. trainData_new$prob <- train_kknn$prob[, "8"]
70. df <- data.frame(trainData_new)
71. hardest <- df[order(df$prob), ] %>% head(3) # 3 hardest cases to classify
72. easiest <- df[rev(order(df$prob)), ] %>% head(2) #2 easiest cases to classify
73. heatmap(matrix(as.numeric(easiest[1,1:64]),byrow = TRUE,nrow = 8,ncol = 8),Colv = NA,Rowv =
NA,col=c("white","purple"))
74.
75. heatmap(matrix(as.numeric(easiest[2,1:64]),byrow = TRUE,nrow = 8,ncol = 8),Colv = NA,Rowv =
NA,col=c("white","purple"))
76.
77. heatmap(matrix(as.numeric(hardest[1,1:64]),byrow = TRUE,nrow = 8,ncol = 8),Colv = NA,Rowv =
NA,col=c("white","red"))
78.
79. heatmap(matrix(as.numeric(hardest[2,1:64]),byrow = TRUE,nrow = 8,ncol = 8),Colv = NA,Rowv =
NA,col=c("white","red"))
80.
81. heatmap(matrix(as.numeric(hardest[3,1:64]),byrow = TRUE,nrow = 8,ncol = 8),Colv = NA,Rowv =
NA,col=c("white","red"))
82.
83. missclass_training <- c()
84. missclass_validation <- c()
85. for(k in 1:30){
86.   kknn_training <- kknn(formula = V65~.,kernel = "rectangular",train =trainData,test =
trainData,k=k )
87.   kknn_validation <- kknn(formula = V65~.,kernel = "rectangular",train =trainData,test =
validationData,k=k)
88.   missclass_training <- c(missclass_training,miss_class(trainData$V65,kknn_training$fitted.values))
89.   missclass_validation <-
c(missclass_validation,miss_class(validationData$V65,kknn_validation$fitted.values))
90. }
91.
92. library(ggplot2)
93. k<-1:30
94. df1<-data.frame(missclass_training,k1 =k)
95. df2<-data.frame(missclass_validation,k2 =k)
96. df3 <- data.frame(df1,df2)
97. ggplot( ) +
98.   geom_line(aes(x=df3$k1,y=df3$missclass_training,
99.                 colour="Training")) +
100.  geom_line(aes(x=df3$k2,y=df3$missclass_validation,
101.                colour="Validation")) +
102.  theme_minimal()+ggtitle("Finding the Optimal Parameter") +
103.  scale_x_continuous(breaks = seq(1,30,2)) +
104.  xlab("Parameter 'k'") + ylab("Missclassification") +
105.  geom_vline(xintercept = 3
106.             , linetype = "dashed",color = "green")
107.
108. kknn_test <- kknn(formula = V65~.,kernel = "rectangular",train =trainData,test = testData,k=3 )
109. cat("for k=3 :\n")
110. cat("MissClassification error on Training =",missclass_training[3],"\n")
111. cat("MissClassification error on Validation = ",missclass_validation[3],"\n")
112. cat("MissClassification error on Test = ",
113.      miss_class(testData$V65,kknn_test$fitted.values),"\n")
114. library(data.table)
115. library(kknn)
116. crossEntropy_training <- c()
117. crossEntropy_validation <- c()
118. noOfClasses <- length(unique(trainData$V65))
119. noOfRows_training <- nrow(trainData)
120. noOfRows_validation <- nrow(validationData)
121. training_onehot <- t(sapply(as.numeric(trainData$V65),function(x)
122.   {c(rep(0,x-1), 1 ,rep(0 , noOfClasses-x ))}))
123.
124. validation_onehot <- t(sapply(as.numeric(validationData$V65),function(x)
125.   {c(rep(0,x-1), 1 ,rep(0 , noOfClasses-x ))}))
126. for(i in 1:30)
127. {

```



```

128. fitModel_training <- kknn(V65~.,kernel = "rectangular", train = trainData, test = trainData, k=i)
129. fitModel_validation <- kknn(V65~.,kernel = "rectangular", train = trainData, test =
validationData, k=i)
130.
131. #cross-entropy
132. crossEntropy_training <- c(crossEntropy_training,sum(training_onehot * -
log(fitModel_training$prob + 10^-15))/noOfRows_training)
133.
134. crossEntropy_validation <- c(crossEntropy_validation,sum(validation_onehot * -
log(fitModel_validation$prob + 10^-15))/noOfRows_validation)
135.
136. }
137. crossEntropy_validation
138.
139. crossEntropy <- melt(data.table(k = 1:30, Validation = crossEntropy_validation), "k",variable.name
= " ")
140. ggplot(crossEntropy) + geom_line(aes(k, value),col="red") + theme_bw() +
141.   geom_vline(xintercept = 15, linetype = "dotted") + scale_x_continuous(breaks = seq(1,30)) +
142.   xlab("'k'") + ylab("Cross-Entropy") +
143.   ggtitle("The Optimal Parameter")
144.

```

## Code for assignment2:

```

1. require(ggplot2)
2. require(lattice)
3. library(caret)
4. #1. Divide it into training and test data (60/40) and scale it appropriately.
5. #In the coming steps, assume that motor_UPDRS is normally distributed and is a
6. #function of the voice characteristics, and since the data are scaled, no
7. #intercept is needed in the modelling.
8.
9. parkinsons = read.csv("parkinsons.csv")
10. parkinson_data = as.data.frame(parkinsons[,c(5,7:22)])
11.
12. n = dim(parkinson_data)[1]
13. set.seed(12345)
14. id = sample(1:n, floor(n*0.6))
15. traindata = parkinson_data[id,]
16. testdata = parkinson_data[-id,]
17. scaler <- preProcess(as.data.frame(traindata))
18. train_s <- predict(scaler,traindata)
19. test_s <- predict(scaler,testdata)
20.
21. #2.compute a linear regression model from the training data, estimate training
22. #and test MSE and comment on which variables contribute significantly to the
23. #model.
24.
25. lrm <- lm(motor_UPDRS ~ .-1, data = train_s)
26. train_predicted <- predict(lrm, newdata = train_s )
27. test_predicted <- predict(lrm, newdata = test_s)
28. MSETrain <- mean((train_predicted - train_s$motor_UPDRS)^2)
29. MSETest <- mean((test_predicted - test_s$motor_UPDRS)^2)
30. (c(MSE_of_train_data = MSETrain, MSE_of_test_data = MSETest ))
31. summary(lrm)
32.
33. #3.Implement 4 following basic functions by using basic R commands only ( not using any external
packages):
34. #a. loglikelihood function that for a given vector $\theta$  and dispersion  $\sigma$  compute the log-likelihood
function
35. #logP(T| $\theta,\sigma$ ) for the states model and training data.
36.
37. loglikelihood <- function(parameters) {
38.   theta <- parameters[1:16]
39.   sigma <- parameters[17]
40.   y <- as.matrix(train_s[,1]) %*% theta
41.   n <- nrow(train_s)
42.   return (-n/2*log(2*pi) -n/2*log(sigma^2)-sum((train_s[,1]-y)^2)/(2*sigma^2))
43. }
44.

```

```

45. #b.Ridge function that for a given vector  $\theta$ , scalar  $\sigma$  and scalar  $\lambda$  uses function from 3a
46. # and adds up a Ridge penalty  $\lambda\|\theta\|^2$  to the minus log likelihood.
47.
48. ridgePenalty <- function(parameters, lambda){
49.   return(-loglikelihood(parameters)+ lambda*norm(as.matrix(parameters[1:16]),type="2")^2)
50. }
51.
52. #c. RidgeOpt function that depends on scalar  $\lambda$ , uses function from 3b and function optim()
53. # with method="BFGS" to find the optimal  $\theta$  and  $\sigma$  for given value of  $\lambda$ .
54.
55. ridgeOpt <- function(lambda){
56.   return (optim( par = c(rep(0,length=16),0.01), fn = ridgePenalty, lambda = lambda, method =
"BFSGS"))
57. }
58.
59. #d. DF function that for a given scalar  $\lambda$  computes the degree of freedom of
60. #the Ridgemodel based on training data.
61. df <- function(lambda){
62.   X <- as.matrix(train_s[,-1])
63.   P <- X %>% solve(t(X)%%X + lambda*diag(16))%>%t(X)
64.   return(sum(diag(P)))
65. }
66.
67. #4. by using RidgeOpt(), compute optimal  $\theta$  for  $\lambda=1$ ,  $\lambda=100$ ,  $\lambda=1000$ , use the estimated parameter to
68. #predict the motor_UPDRS value for training and test data and report the training and Test MSE
values.
69. # which penalty parameter is most appropriate among the selected one? compute and compare the degree
of freedom
70. # of these model and make appropriate conclusion.
71.
72. ridgeOpt(lambda = 1)
73. train_predict <- as.matrix(train_s[,-1]) %>% ridgeOpt(lambda = 1)$par[1:16]
74. test_predict <- as.matrix(test_s[,-1]) %>% ridgeOpt(lambda = 1)$par[1:16]
75. MSE_Train_1 <- sum((train_predict - train_s$motor_UPDRS)^2) / nrow(train_s)
76. MSE_Test_1 <- sum((test_predict - test_s$motor_UPDRS)^2) / nrow(test_s)
77.
78.
79. ridgeOpt(lambda = 100)
80. train_predict <- as.matrix(train_s[,-1]) %>% ridgeOpt(lambda = 100)$par[1:16]
81. test_predict <- as.matrix(test_s[,-1]) %>% ridgeOpt(lambda = 100)$par[1:16]
82. MSE_Train_100 <- sum((train_predict - train_s$motor_UPDRS)^2) / nrow(train_s)
83. MSE_Test_100 <- sum((test_predict - test_s$motor_UPDRS)^2) / nrow(test_s)
84.
85.
86. ridgeOpt(lambda = 1000)
87. train_predict <- as.matrix(train_s[,-1]) %>% ridgeOpt(lambda = 1000)$par[1:16]
88. test_predict <- as.matrix(test_s[,-1]) %>% ridgeOpt(lambda = 1000)$par[1:16]
89. MSE_Train_1000 <- sum((train_predict - train_s$motor_UPDRS)^2) / nrow(train_s)
90. MSE_Test_1000 <- sum((test_predict - test_s$motor_UPDRS)^2) / nrow(test_s)
91.
92. MSE_Matrix<-
matrix(c(MSE_Train_1,MSE_Test_1,MSE_Train_100,MSE_Test_100,MSE_Train_1000,MSE_Test_1000), nrow =2, ncol
= 3)
93. rownames(MSE_Matrix) <-c("Train_MSE", "Test_MSE")
94. colnames(MSE_Matrix) <-c(" $\lambda=1$ ", " $\lambda=100$ ", " $\lambda=1000$ ")
95. print(MSE_Matrix)
96.
97. c(DF_1 = df(1),DF_100 = df(100),DF_1000 = df(1000))
98.

```

## Code for assignment 3

### Part1

```

1. data=read.csv("pima-indians-diabetes.csv",header = FALSE)
2.
3. colnames(data)=c("Number of times pregnant",
4.   "Plasma",
5.   " Diastolic blood pressure",
6.   "Triceps skinfold thickness",

```

```

7.  "2-Hour serum insulin",
8.  "Body mass index",
9.  "Diabetes pedigree function",
10. "Age",
11. "Diabetes")
12.
13. haveDiabetes= (as.integer(unlist((data$Diabetes))))
14. dontHaveDiabetes= -1*(as.integer(unlist((data$Diabetes)))-1)
15. plasma= as.integer(unlist((data$`Plasma`)))
16. age= as.integer(unlist((data$Age)))
17. plot(haveDiabetes*age,haveDiabetes*plasma,ylab = " Plasma glucose concentration",
18.       xlab = "Age",col="red")
19. points(dontHaveDiabetes*age,dontHaveDiabetes*plasma,col="blue")
20.

```

## Part2,3

```

1. data=read.csv("pima-indians-diabetes.csv",header = FALSE)
2.
3. colnames(data)=c("Number of times pregnant",
4.                  "Plasma",
5.                  "Diastolic blood pressure",
6.                  "Triceps skinfold thickness",
7.                  "2-Hour serum insulin",
8.                  "Body mass index",
9.                  "Diabetes pedigree function",
10.                 "Age",
11.                 "Diabetes")
12. # Devide date
13. n=dim(data)[1]
14. set.seed(12345)
15. id=sample(1:n, floor(n*0.5))
16. train=data[id,]
17. test=data[-id,]
18. # Train model
19.
20. logistic=glm(Diabetes ~ Age + Plasma,data = train,family = "binomial")
21. LRpred=predict(logistic,newdata = test,type="response")
22. cutoff=0.5
23. LRpred[LRpred<cutoff]=0
24. LRpred[LRpred>=cutoff]=1
25. summary(logistic)
26.
27. #Misclassification rate
28.
29. confusionMX=table(LRpred,test$Diabetes)
30. misclassificationRate=(1 - sum(diag(confusionMX))/sum(confusionMX))
31.
32. #Scatterplot
33. haveDiabetes=as.integer(unlist(LRpred))
34. dontHaveDiabetes= -1*(as.integer(unlist((LRpred)))-1)
35. plasma= as.integer(unlist((test$`Plasma`)))
36. age= as.integer(unlist((test$Age)))
37. plot(haveDiabetes*age,haveDiabetes*plasma,ylab = " Plasma glucose concentration",
38.       xlab = "Age",col="red")
39. points(dontHaveDiabetes*age,dontHaveDiabetes*plasma,col="blue")
40.
41. ##### Part3
42.
43. #Calculate equation
44. coeff=logistic$coefficients
45. slope=coeff[2]/(-1*coeff[3])
46. intercept=coeff[1]/(-1*coeff[3])
47. #Plot decision boundary
48. abline(intercept,slope,col="purple")

```

## Part4

```

1. ###Seting r=0.2
2. data=read.csv("pima-indians-diabetes.csv",header = FALSE)

```

```

3.
4. colnames(data)=c("Number of times pregnant",
5.                  "Plasma",
6.                  " Diastolic blood pressure",
7.                  "Triceps skinfold thickness",
8.                  "2-Hour serum insulin",
9.                  "Body mass index",
10.                 "Diabetes pedigree function",
11.                 "Age",
12.                 "Diabetes")
13. # Devide date
14. n=dim(data)[1]
15. set.seed(12345)
16. id=sample(1:n, floor(n*0.5))
17. train=data[id,]
18. test=data[-id,]
19. # Train model
20.
21. logistic=glm(Diabetes ~ Age + Plasma,data = train,family = "binomial")
22. LRpred=predict(logistic,newdata = test,type="response")
23. cutoff=0.2
24. LRpred[LRpred<cutoff]=0
25. LRpred[LRpred>=cutoff]=1
26.
27. haveDiabetes=as.integer(unlist(LRpred))
28. dontHaveDiabetes= -1*(as.integer(unlist((LRpred)))-1)
29. plasma= as.integer(unlist((test$`Plasma`)))
30. age= as.integer(unlist((test$Age)))
31. plot(haveDiabetes*age,haveDiabetes*plasma,ylab = " Plasma glucose concentration",
32.       xlab = "Age",col="red")
33. points(dontHaveDiabetes*age,dontHaveDiabetes*plasma,col="blue")
34. coeff=logistic$coefficients
35. slope=coeff[2]/(-1*coeff[3])
36. intercept=coeff[1]/(-1*coeff[3])
37.
38. #Plot decision boundary
39. abline(intercept,slope,col="purple")
40.
41.

```

## Part5

```

1. data=read.csv("pima-indians-diabetes.csv",header = FALSE)
2.
3. colnames(data)=c("Number of times pregnant",
4.                  "Plasma",
5.                  " Diastolic blood pressure",
6.                  "Triceps skinfold thickness",
7.                  "2-Hour serum insulin",
8.                  "Body mass index",
9.                  "Diabetes pedigree function",
10.                 "Age",
11.                 "Diabetes")
12. # Devide date
13. n=dim(data)[1]
14. set.seed(12345)
15. id=sample(1:n, floor(n*0.5))
16. train=data[id,]
17. test=data[-id,]
18.
19. #new features
20. train$phi1=phi1=train$Plasma^4
21. train$phi2=phi2=train$Plasma^3*train$Age
22. train$phi3=phi3=train$Plasma^2*train$Age^1
23. train$phi4=phi4=train$Plasma^1*train$Age^3
24. train$phi5=phi5=train$Age^4
25.
26. logistic=glm(Diabetes ~ Plasma +Age+phi1+phi2+phi3+phi4+phi5
27.              ,data = train,family = "binomial")
28. LRpred=predict(logistic,newdata = test,type="response")

```

```
29. cutoff=0.5
30. LRpred[LRpred<cutoff]=0
31. LRpred[LRpred>=cutoff]=1
32. summary(logistic)
33.
34. #Misclassification rate
35.
36. confusionMX=table(LRpred,test$Diabetes)
37. misclassificationRate=(1 - sum(diag(confusionMX))/sum(confusionMX))
38.
39. #Scatterplot
40. haveDiabetes=as.integer(unlist(LRpred))
41. dontHaveDiabetes= -1*(as.integer(unlist((LRpred)))-1)
42. plasma= as.integer(unlist((test$`Plasma`)))
43. age= as.integer(unlist((test$Age)))
44. plot(haveDiabetes*age,haveDiabetes*plasma,ylab = " Plasma glucose concentration",
45.       xlab = "Age",col="red")
46. points(dontHaveDiabetes*age,dontHaveDiabetes*plasma,col="blue")
47.
```