

Lab1

liume102@student.liu.se
hanxi898@student.liu.se
xiali125@student.liu.se

12 Nov 2024

1 Statement of Contribution

In Assignment 1, Xiaochen Liu was mainly responsible for code writing while Liuxi Mei was responsible for the analyses. Assignment 2 was mainly contributed by Han Xia. In assignment 3, Liuxi Mei was responsible for code writing while Han Xia was responsible for the analysis. Assignment 4 was mainly contributed by Xiaochen Liu and Liuxi Mei. Results from all assignments have been discussed afterwards between Liuxi Mei, Xiaochen Liu and Han Xia and the group report was created based on this discussion.

2 Introduction

This is the first lab in the Machine Learning In this lab, contains the following tasks:1. Handwritten digit recognition with K-nearest neighbors.2. Linear regression and ridge regression.3. Logistic regression and basis function expansion.4. Theory

3 Assignment 1: Handwritten digit recognition with K-nearest neighbors

3.1 Load and check data

```
# Load packages  
library('ggplot2') # visualization  
library('ggthemes') # visualization
```

```
## Warning:  'ggthemes' R 4.4.2
```

```
library('scales') # visualization  
library('dplyr') # data manipulation  
library('randomForest') # classification algorithm
```

```
## Warning:  'randomForest' R 4.4.2
```

```
library('caret')
```

Now that our packages are loaded and we divide it into training, validation and test sets (50%/25%/25%)

```
# do not use StringAsFact = FALSE
digitals <- read.csv('../data/optdigits.csv',header = FALSE)
# change all the columns to factor
#digitals <- digitals %>% mutate_all(as.factor)
digitals$V65 <- as.factor(digitals$V65)
train_index <- createDataPartition(digitals$V65, p = 0.5, list = F)
train_digitals <- digitals[train_index,]
remainingData <- digitals[-train_index, ]
validationIndex <- createDataPartition(remainingData$V65, p = 0.5, list = FALSE)
valid_digitals <- remainingData[validationIndex, ]
test_digitals <- remainingData[-validationIndex, ]
cat("train length:", nrow(train_digitals),'\n')
```

```
## train length: 1914
```

```
cat("test length:", nrow(valid_digitals),'\n')
```

```
## test length: 956
```

```
cat("valid length:", nrow(test_digitals),'\n')
```

```
## valid length: 953
```

3.2 KNN to fit classification model using train data

```
library(kknn)
formula <- V65~.
# if kernel = 'rectangular' , so every point in the neighborhood is weighted equally
# both of the parameters of train and test use train_digital data
# if your predict columns is continuous, kknn will recognized as a regression task
# under this situation, you can not get a probability of the prediction
knn_train_model <- kknn(formula, train_digitals, train_digitals,
                        kernel = 'rectangular',distance = 1,)
train_predictions <- fitted(knn_train_model)
print(length(train_predictions))
```

```
## [1] 1914
```

3.3 Predict on the test data

```
knn_test_model <- kknn(formula, train = train_digitals, test = test_digitals,
                      k = 30, kernel = 'rectangular')
print(length(knn_test_model$fitted.values))
```

```
## [1] 953
```

3.4 Confusion matrices and Misclassification errors for train data and test data

```
train_confusion <- table(train_digitals$V65, train_predictions)
test_confusion <- table(test_digitals$V65, knn_test_model$fitted.values)
test_error_rate <- 1 - sum(diag(test_confusion)) / sum(test_confusion)
train_error_rate <- 1 - sum(diag(train_confusion)) / sum(train_confusion)
# only observer the top 10 rows
cat("Misclassification errors on train data:", train_error_rate, '\n')
```

```
## Misclassification errors on train data: 0.02194357
```

```
cat("train_confusion:")
```

```
## train_confusion:
```

```
table(train_digitals$V65[1:10], train_predictions[1:10])
```

```
##
##      0 1 2 3 4 5 6 7 8 9
## 0 2 0 0 0 0 0 0 0 0
## 1 0 1 0 0 0 0 0 0 0
## 2 0 0 1 0 0 0 0 0 0
## 3 0 0 0 1 0 0 0 0 0
## 4 0 0 0 0 0 0 0 0 1
## 5 0 0 0 0 0 1 0 0 0
## 6 0 0 0 0 0 0 1 0 0
## 7 0 0 0 0 0 0 0 1 0
## 8 0 0 0 0 0 0 0 0 1
## 9 0 0 0 0 0 0 0 0 0
```

```
cat("Misclassification errors on test data:", test_error_rate, '\n')
```

```
## Misclassification errors on test data: 0.04826863
```

```
cat("test confusion:")
```

```
## test confusion:
```

```
table(test_digitals$V65[1:10], knn_test_model$fitted.values[1:10])
```

```
##
##      0 1 2 3 4 5 6 7 8 9
## 0 3 0 0 0 0 0 0 0 0
## 1 0 0 0 0 0 0 0 0 0
## 2 0 0 0 0 0 0 0 0 0
## 3 0 0 0 0 0 0 0 0 0
## 4 0 0 0 0 0 0 0 0 1
## 5 0 0 0 0 0 1 0 0 0
## 6 0 0 0 0 0 0 0 0 0
## 7 0 0 0 0 0 0 0 4 0
## 8 0 0 0 0 0 0 0 0 0
## 9 0 0 0 0 0 0 0 0 1
```

3.5 Filter 2 cases of digit “8” in the training data which were easiest to classify and 3 cases that were hardest to classify

```
# filter the digital '8'
library(dplyr)

train_predict <- data.frame(train_digitals$V65, train_predictions,knn_train_model$prob)
train_predict$max_prob <- apply(train_predict[,3:12], 1, max)

train_predict_8 <- train_predict[train_predict$train_digitals.V65 == 8,]
# do not change the index while sorting
train_predict_8 <- train_predict_8[order(train_predict_8$X8), , drop = FALSE]

# get the 3 cases that were hardest to classify
hardest_cases_for_8 <- train_predict_8 %>% head(3)
easy_cases_for_8 <- train_predict_8 %>% tail(2)
```

3.6 Analysis the difference of the hardest case and easiest cases

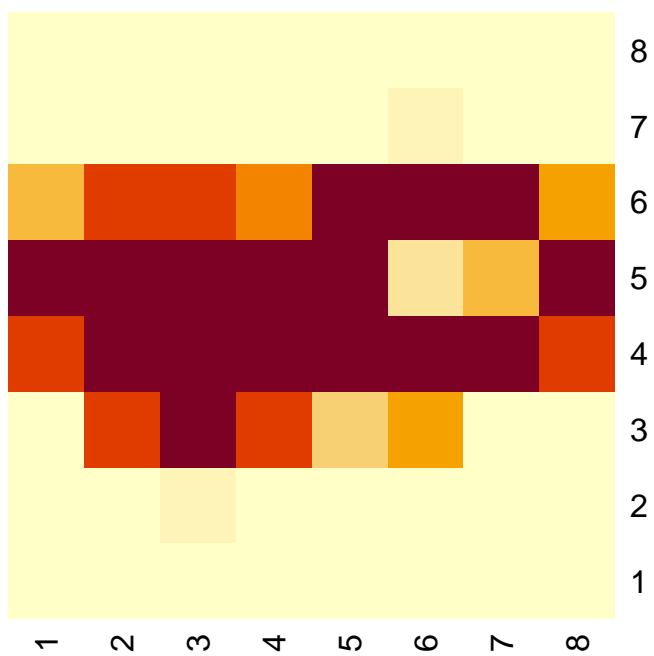
we can see on the heatmap that the hardest cases are more complex than the easiest cases. Dark-colored squares concentrated in the middle of the matrix while the easiest cases are more concentrated on the edges which looking more like the number 8.

```
hardest_cases_index <- rownames(hardest_cases_for_8)
est_cases_index <- rownames(easy_cases_for_8)
# reindex the row index
row.names(train_digitals) <- NULL
full_hardest_cases <- train_digitals[hardest_cases_index,1:64]
full_est_cases <- train_digitals[est_cases_index,1:64]

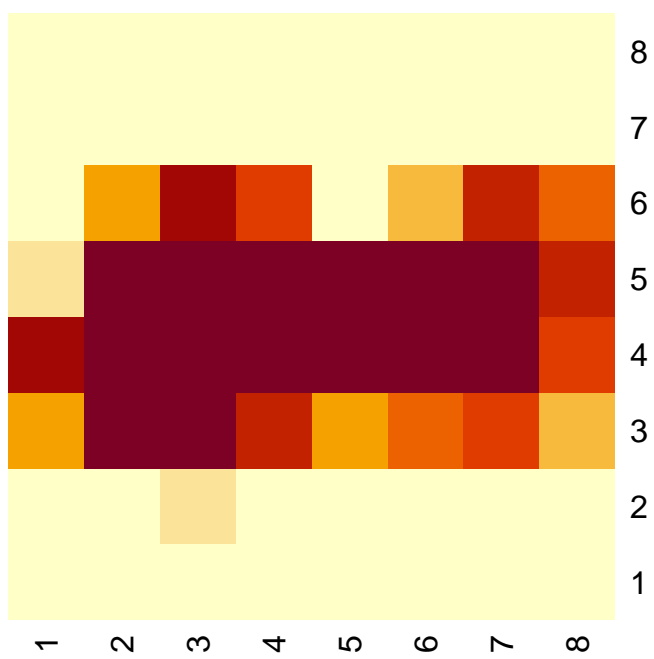
hardest_matrixs <- lapply(1:nrow(full_hardest_cases),
  function(i) matrix(as.numeric(full_hardest_cases[i, ,
    drop = FALSE]),nrow = 8,ncol = 8))
est_matrixs <- lapply(1:nrow(full_est_cases),
  function(i) matrix(as.numeric(full_est_cases[i, ,
    drop = FALSE]),nrow = 8,ncol = 8))

for (i in 1:length(hardest_matrixs)) {
  mat <- hardest_matrixs[[i]]
  heatmap(mat, Colv = NA, Rowv = NA, scale = "none", main = paste("Hard Case", i))
}
```

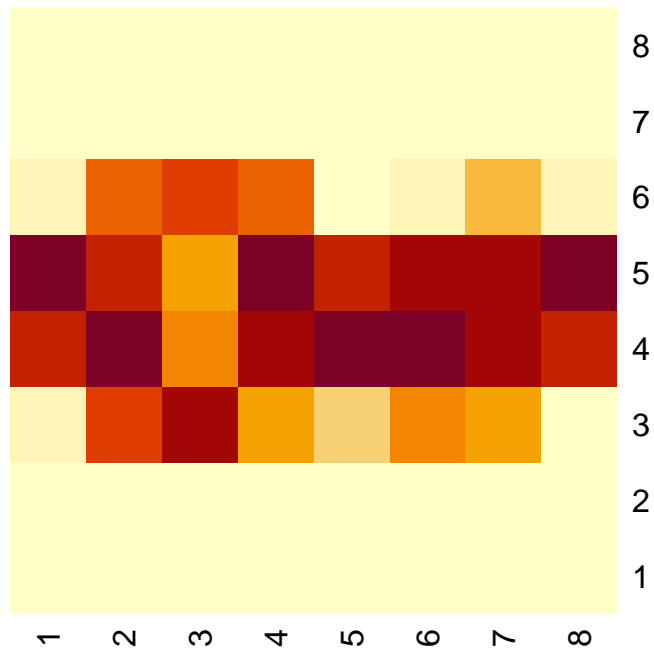
Hard Case 1



Hard Case 2

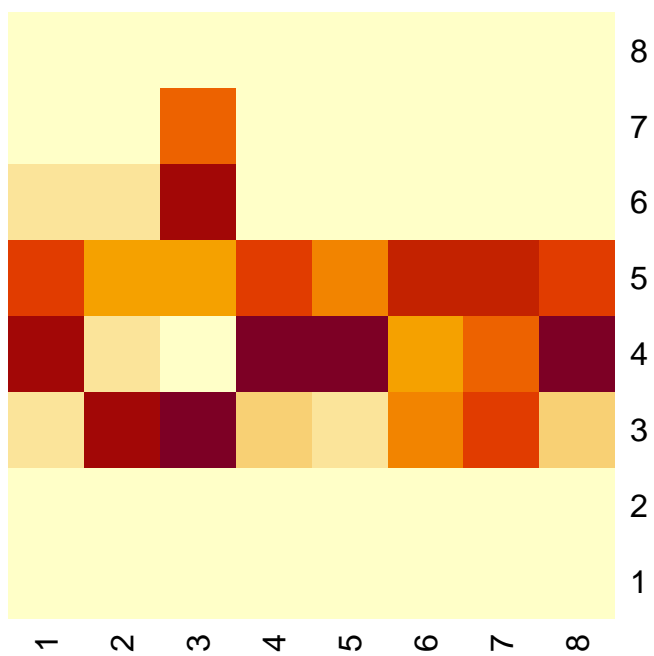


Hard Case 3

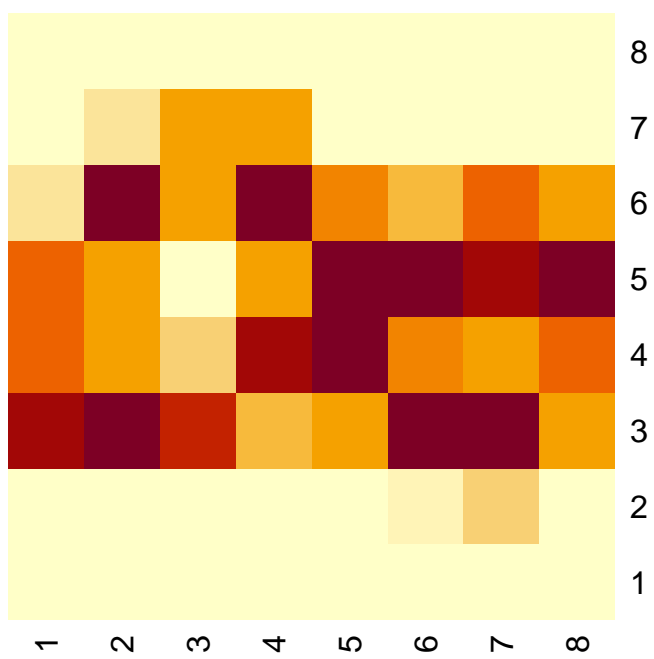


```
for (i in 1:length(est_matrixs)) {  
  mat <- est_matrixs[[i]]  
  heatmap(mat, Colv = NA, Rowv = NA, scale = "none", main = paste("Hard Case", i))  
}
```

Hard Case 1



Hard Case 2



3.7 Training via different k in on the training and validation data

according the plot , $k = 3$ is best value on training data and validation data, though the performance of $k = 1$ is better than $k = 3$ on training data, it is not the best value on validation data due to the weak generalization ability, but when we apply it on test data, its performance is not as good as predicted

```
library(ggplot2)
train_error_rates <- list()
valid_error_rates <- list()
test_error_rates <- list()

for (ki in 1:30) {
  # cat(paste("current k:",ki,"\n",sep=""))
  train_ki_model <- kknnc(formula, train = train_digitals, test = train_digitals,
    k = ki, kernel = 'rectangular')
  valid_ki_model <- kknnc(formula, train = train_digitals, test = valid_digitals,
    k = ki, kernel = 'rectangular')

  test_ki_model <- kknnc(formula, train = train_digitals, test = test_digitals,
    k = ki, kernel = 'rectangular')

  train_confusion <- table(train_digitals$V65, train_ki_model$fitted.values)
  valid_confusion <- table(valid_digitals$V65, valid_ki_model$fitted.values)
  test_confusion <- table(test_digitals$V65, test_ki_model$fitted.values)

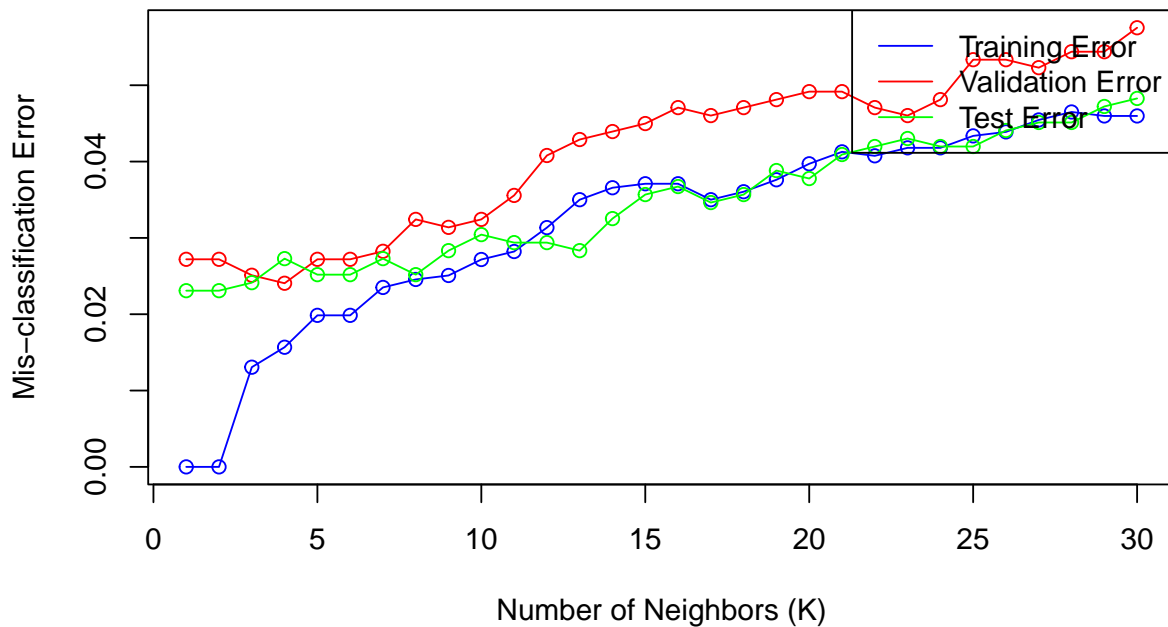
  train_error_rate <- sum(diag(train_confusion)) / sum(train_confusion)
  valid_error_rate <- sum(diag(valid_confusion)) / sum(valid_confusion)

  test_error_rate <- sum(diag(test_confusion)) / sum(test_confusion)

  # print(train_error_rate)
  # print(valid_error_rate)
  train_error_rates[[ki]] <- 1 - train_error_rate
  valid_error_rates[[ki]] <- 1 - valid_error_rate
  test_error_rates[[ki]] <- 1 - test_error_rate
}

plot(1:30, train_error_rates, type = "o", col = "blue",
  ylim = range(c(train_error_rates, valid_error_rates)),
  xlab = "Number of Neighbors (K)", ylab = "Mis-classification Error",
  main = "Training and Validation Errors")
lines(1:30, valid_error_rates, type = "o", col = "red")
lines(1:30, test_error_rates, type = "o", col = "green")
legend("topright", legend = c("Training Error", "Validation Error", "Test Error"),
  col = c("blue", "red", "green"), lty = 1)
```


Training and Validation Errors



#

3.8 Change mis-classification error to cross-entropy

```
valid_cross_entropy_errors <- list()
train_cross_entropy_errors <- list()
test_cross_entropy_errors <- list()

for (ki in 1:30) {

  print(ki)
  valid_ki_model <- kknn(formula, train = train_digitals, test = valid_digitals,
                          k = ki, kernel = 'rectangular')
  train_ki_model <- kknn(formula, train = train_digitals, test = train_digitals,
                          k = ki, kernel = 'rectangular')
  test_ki_model <- kknn(formula, train = train_digitals, test = test_digitals,
                         k = ki, kernel = 'rectangular')
  valid_probs <- valid_ki_model$prob
  train_probs <- train_ki_model$prob
  test_probs <- test_ki_model$prob

  valid_log_probs <- log(valid_probs + 1e-15) # Add small constant to avoid log(0)
  train_log_probs <- log(train_probs + 1e-15) # Add small constant to avoid log(0)
  test_log_probs <- log(test_probs + 1e-15) # Add small constant to avoid log(0)

  # -1 means do not contain intercept
```

```

# One-hot encoding
#This type of matrix is typically used in machine learning and statistical modeling for feature
#engineering, particularly when converting categorical variables into dummy variables.
valid_correct_class <- model.matrix(~V65 - 1, data = valid_digitals) # One-hot encoding
train_correct_class <- model.matrix(~V65 - 1, data = train_digitals) # One-hot encoding
test_correct_class <- model.matrix(~V65 - 1, data = test_digitals) # One-hot encoding

valid_cross_entropy_errors[[ki]] <- -sum(valid_correct_class
                                         * valid_log_probs) / nrow(valid_digitals)
train_cross_entropy_errors[[ki]] <- -sum(train_correct_class
                                         * train_log_probs) / nrow(train_digitals)
test_cross_entropy_errors[[ki]] <- -sum(test_correct_class
                                         * test_log_probs) / nrow(test_digitals)

print(-sum(valid_correct_class * valid_log_probs) )
print(-sum(train_correct_class * train_log_probs))
print(-sum(test_correct_class * test_log_probs) )
}

```

```

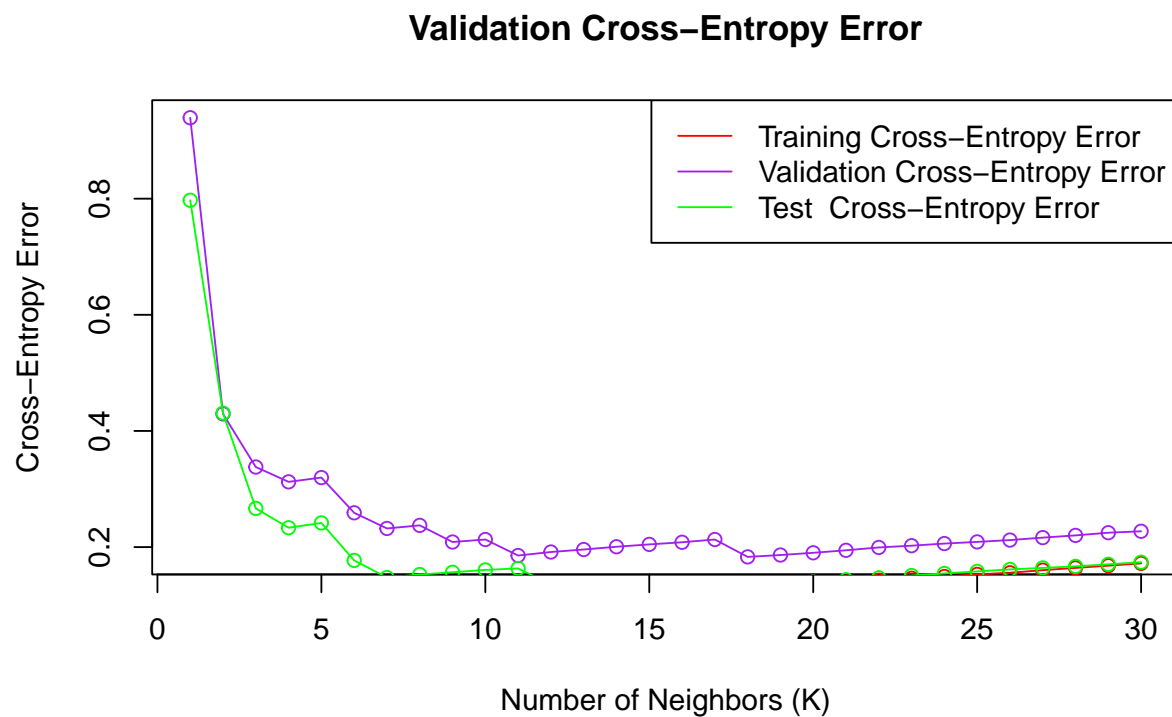
## [1] 1
## [1] 898.0082
## [1] -2.124967e-12
## [1] 759.8531
## [1] 2
## [1] 410.425
## [1] 36.7368
## [1] 410.425
## [1] 3
## [1] 323.0158
## [1] 65.97094
## [1] 254.1082
## [1] 4
## [1] 298.6595
## [1] 85.77618
## [1] 222.4149
## [1] 5
## [1] 305.6495
## [1] 105.5458
## [1] 230.2473
## [1] 6
## [1] 247.6512
## [1] 120.7046
## [1] 168.7506
## [1] 7
## [1] 221.8669
## [1] 135.2842
## [1] 140.6422
## [1] 8
## [1] 226.9749
## [1] 147.1415
## [1] 145.5915
## [1] 9
## [1] 199.5943
## [1] 158.2855

```

```
## [1] 149.3081
## [1] 10
## [1] 203.6969
## [1] 169.4683
## [1] 153.1582
## [1] 11
## [1] 177.3317
## [1] 180.9608
## [1] 155.6134
## [1] 12
## [1] 183.0979
## [1] 193.9443
## [1] 129.5844
## [1] 13
## [1] 187.2329
## [1] 203.0935
## [1] 101.5549
## [1] 14
## [1] 191.7639
## [1] 212.6724
## [1] 105.561
## [1] 15
## [1] 195.6415
## [1] 220.7923
## [1] 111.4436
## [1] 16
## [1] 199.1468
## [1] 228.887
## [1] 116.8026
## [1] 17
## [1] 203.7113
## [1] 236.3774
## [1] 120.183
## [1] 18
## [1] 175.0995
## [1] 244.4901
## [1] 124.9017
## [1] 19
## [1] 178.1482
## [1] 252.7682
## [1] 129.0026
## [1] 20
## [1] 181.8277
## [1] 260.1861
## [1] 133.1465
## [1] 21
## [1] 185.9388
## [1] 267.3539
## [1] 137.0905
## [1] 22
## [1] 190.6113
## [1] 273.7041
## [1] 140.2434
## [1] 23
```

```
## [1] 193.5203
## [1] 280.1422
## [1] 143.9528
## [1] 24
## [1] 197.0174
## [1] 285.7942
## [1] 147.3564
## [1] 25
## [1] 199.7635
## [1] 292.8862
## [1] 150.6803
## [1] 26
## [1] 202.7811
## [1] 298.2519
## [1] 153.982
## [1] 27
## [1] 206.8199
## [1] 306.7626
## [1] 156.437
## [1] 28
## [1] 210.453
## [1] 314.2168
## [1] 159.0041
## [1] 29
## [1] 214.8819
## [1] 321.5783
## [1] 162.3622
## [1] 30
## [1] 217.2431
## [1] 329.1383
## [1] 165.8107
```

```
# plot(1:30, train_error_rates, type = "o", col = "blue",
# ylim = range(c(train_error_rates, valid_error_rates)),
# xlab = "Number of Neighbors (K)", ylab = "Mis-classification Error",
# main = "Training and Validation Errors")
# lines(1:30, valid_error_rates, type = "o", col = "red")
# lines(1:30, test_error_rates, type = "o", col = "green")
# legend("topright", legend = c("Training Error", "Validation Error", "Test Error"),
# col = c("blue", "red", "green"), lty = 1)
plot(1:30, valid_cross_entropy_errors, type = "o", col = "purple",
xlab = "Number of Neighbors (K)", ylab = "Cross-Entropy Error",
main = "Validation Cross-Entropy Error")
lines(1:30, train_cross_entropy_errors, type = "o", col = "red")
lines(1:30, test_cross_entropy_errors, type = "o", col = "green")
legend("topright", legend = c("Training Cross-Entropy Error",
"Validation Cross-Entropy Error", "Test Cross-Entropy Error"),
col = c("red", "purple", "green"), lty = 1)
```



4 Assignment 2: Linear regression and ridge regression

4.1 set up

We need to download some useful packages before the start.

```
install.packages("caret")
```

```
## Warning: 'caret'
```

```
library(caret)
```

4.2 Prepare the dataset

Firstly, we read the file and divided the data into training and test data (60/40).

```
data <- read.csv("../data/parkinsons.csv") #
set.seed(42)
ini_sample<- sample(1:nrow(data),0.6*nrow(data))
train_data<- data[ini_sample,]
test_data<- data[-ini_sample,]
```

And then we scaled the dataset appropriately.

```
sacale_data<- train_data[,names(train_data)!="motor_UPDRS"]
scale_para<- preProcess(sacale_data)
train_data_scaled<- predict(scale_para,train_data)
test_data_scaled<- predict(scale_para,test_data)
train_data_scaled$motor_UPDRS <- train_data$motor_UPDRS
test_data_scaled$motor_UPDRS <- test_data$motor_UPDRS
```

4.3 Build models

Next, we computed a linear regression model , estimate training and test MSE

```
model<- lm(motor_UPDRS ~ .,train_data_scaled)
train_prediction<- predict(model,train_data_scaled)
train_mse<- mean((train_prediction - train_data_scaled$motor_UPDRS)^2)
test_prediction<- predict(model,test_data_scaled)
test_mse<- mean((test_prediction - test_data_scaled$motor_UPDRS)^2)
summary(model)
```

```
##
## Call:
## lm(formula = motor_UPDRS ~ ., data = train_data_scaled)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
	-8.4962	-1.3230	0.1978	1.6722	6.7627

```
##
## Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	21.27888	0.04240	501.883	< 2e-16 ***
subject.	-0.12900	0.04908	-2.628	0.008618 **
age	-0.23517	0.04615	-5.095	3.66e-07 ***
sex	0.45948	0.05173	8.883	< 2e-16 ***
test_time	-0.05169	0.04293	-1.204	0.228664
total_UPDRS	7.81829	0.04892	159.827	< 2e-16 ***
Jitter...	1.48149	0.41073	3.607	0.000314 ***
Jitter.Abs.	-0.55937	0.11441	-4.889	1.06e-06 ***
Jitter.RAP	-49.61179	52.25757	-0.949	0.342498
Jitter.PPQ5	-0.31807	0.23637	-1.346	0.178500
Jitter.DDP	48.67488	52.25602	0.931	0.351675
Shimmer	1.00327	0.54793	1.831	0.067183 .
Shimmer.dB.	-0.06932	0.39046	-0.178	0.859095
Shimmer.APQ3	71.26529	209.14155	0.341	0.733311
Shimmer.APQ5	-1.36309	0.30045	-4.537	5.90e-06 ***
Shimmer.APQ11	0.57158	0.15977	3.577	0.000352 ***
Shimmer.DDA	-71.24366	209.14168	-0.341	0.733389
NHR	0.09299	0.12189	0.763	0.445542
HNR	0.16756	0.09855	1.700	0.089157 .
RPDE	-0.23042	0.06213	-3.709	0.000212 ***
DFA	-0.03420	0.05639	-0.606	0.544251
PPE	0.46026	0.09176	5.016	5.54e-07 ***

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

```
##
## Residual standard error: 2.517 on 3503 degrees of freedom
## Multiple R-squared: 0.905, Adjusted R-squared: 0.9045
## F-statistic: 1590 on 21 and 3503 DF, p-value: < 2.2e-16
```

Implement 4 following functions:

loglikelihood function that for a given parameter vector theta and dispersion sigma.

```
logLikelihood <- function(theta, sigma, x, y) {
  n <- length(y)
  predictions <- x %*% theta
  residuals <- y - predictions
  log_likelihood <- -0.5 * n * log(2 * pi * sigma^2) - (t(residuals) %*% residuals) / (2 * sigma^2)
  return(as.numeric(log_likelihood))
}
```

Ridge function that for given vector theta, scalar sigma and scalar lambda and adds up a Ridge penalty to the minus loglikelihood.

```
ridge <- function(theta, sigma, lambda, x, y) {
  log_likelihood <- logLikelihood(theta, sigma, x, y)
  ridge_penalty <- lambda * sum(theta^2)
  return(-log_likelihood + ridge_penalty)
}
```

Use function optim() with method="BFGS" to find the optimal theta and sigma for the given lambda.

```
ridgeopt <- function(lambda, x, y) {
  n <- ncol(x)
  init_params <- c(rep(0, n), 1)
  ridge_obj <- function(params) {
    theta <- params[1:n]
    sigma <- params[n + 1]
    return(ridge(theta, sigma, lambda, x, y))
  }
  opt <- optim(init_params, ridge_obj, method = "BFGS")
  theta_opt <- opt$par[1:n]
  sigma_opt <- opt$par[n + 1]
  return(list(theta = theta_opt, sigma = sigma_opt))
}
```

computes the degrees of freedom of the Ridge model based on the training data.

```
freedom_degree <- function(lambda, x) {
  xT <- t(x) %*% x
  heat <- solve(xT + lambda * diag(ncol(x))) %*% t(x)
  df <- sum(diag(heat)) #trace
  return(df)
}
```

4.4 predict the values

Finally, we can compute optimal theta parameters for different lambda values by using function RidgeOpt.

```
train_data2 <- as.matrix(train_data[,names(train_data)!="motor_UPDRS"])
test_data2<- as.matrix(test_data[,names(test_data)!="motor_UPDRS"])
train_value <- train_data$motor_UPDRS
test_value <- test_data$motor_UPDRS

lambda_values <- c(1, 100, 1000)

train_mse2<- c()
test_mse2<- c()
df<- c()
theta_value<- list()
for (i in seq_along(lambda_values)){
  lambda<- lambda_values[i]
  ridgemodel<- ridgeopt(lambda,train_data2,train_value)
  thetavalue<- ridgemodel$theta

  theta_value[[i]]<- thetavalue

  train_predictions<- train_data2 %*% thetavalue
  train_mse2[i]<- mean((train_value - train_predictions)^2)

  test_predictions<- test_data2 %*% thetavalue
  test_mse2[i]<- mean((test_value - test_predictions)^2)

  df[i] <- freedom_degree(lambda,train_data2)

  result <- list(
    train_mse2 = train_mse2,
    test_mse2 = test_mse2,
    df = df,
    theta_value = theta_value
  )
}
print(result)

## $train_mse2
## [1] 6.465518 6.653863 6.846509
##
## $test_mse2
## [1] 6.387589 6.609421 6.787192
##
## $df
## [1] 0.009998311 0.001516241 0.000442187
##
## $theta_value
## $theta_value[[1]]
## [1] -0.0052580550 -0.0202829652 1.0515160348 -0.0008380227 0.7336342412
## [6] 0.0077671729 -0.0008418776 -0.0209940314 -0.0001519663 -0.0628711878
## [11] 0.0367333247 0.8557988044 -0.1154339829 -0.0815733362 0.3909625855
## [16] -0.3463559230 -0.2125478171 0.0539468078 -1.0994645885 -0.1441446539
```



```
## [21] 2.0815333869
##
## $theta_value[[2]]
## [1] 6.193772e-03 -8.526021e-03 3.650537e-01 -3.961347e-04 7.249791e-01
## [6] 2.022569e-03 7.405812e-05 1.225060e-03 1.374508e-03 2.701054e-03
## [11] 9.127526e-03 8.931548e-02 3.340734e-03 4.754327e-03 9.767899e-03
## [16] 9.350434e-03 1.924272e-02 2.616445e-02 -1.872273e-02 6.979124e-03
## [21] 4.345821e-02
##
## $theta_value[[3]]
## [1] 1.590115e-02 4.205629e-03 4.653460e-02 1.342374e-04 6.968982e-01
## [6] 2.437015e-04 1.712245e-05 1.013762e-04 1.472975e-04 3.740594e-04
## [11] 9.610700e-04 9.501445e-03 3.243690e-04 4.837515e-04 1.003988e-03
## [16] 1.009101e-03 2.221216e-03 1.972672e-02 -1.510650e-03 1.318646e-04
## [21] 4.337676e-03
```

In general, a lower test MSE indicates that the model generalizes better. Higher degrees of freedom mean that models are more flexible and tend to fit details in the data, but can lead to overfitting; Lower degrees of freedom mean that the model is smoother, limiting the fit to the training data.

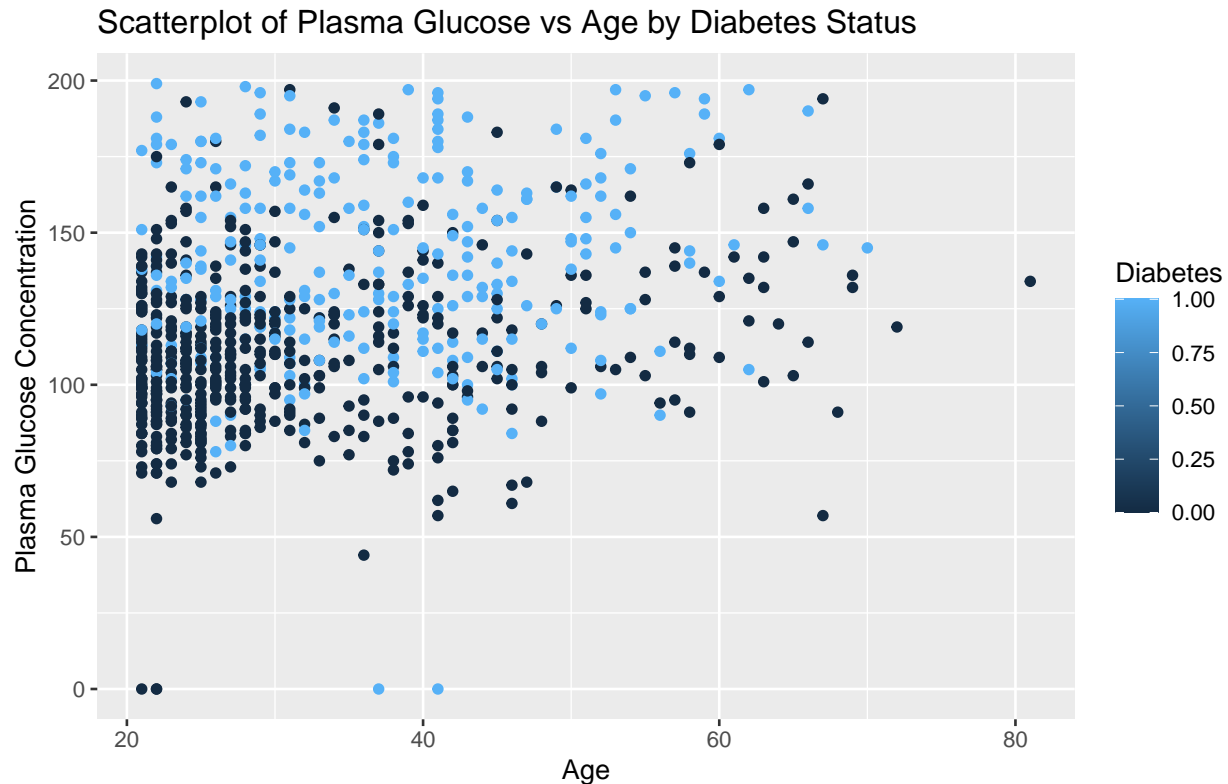
In this example, under penalty parameter equals to 1, the model's train_mse2 and test_mse2 are the lowest, and the degree of freedom is small but not too low. So it is the most appropriate parameter choice.

5 Assignment 3. Logistic regression and basis function expansion

5.1 Read Data and show scatter plot

read data and give a scatter plot showing a Plasma glucose concentration on Age where observations are colored by Diabetes levels

```
diabetes <- read.csv('../data/pima-indians-diabetes.csv', header = FALSE)
colnames(diabetes) <- c('Pregnancies', 'Plasma_glucose', 'blood_pressure', 'TricepsSkinFoldThickness', 'Serum_insulin', 'Diabetes_status')
#
ggplot(diabetes, aes( x = diabetes$Age, y = diabetes$Plasma_glucose, color = diabetes$Diabetes)) +
  geom_point()+labs(x = "Age", y = "Plasma Glucose Concentration", color = "Diabetes") +
  ggtitle("Scatterplot of Plasma Glucose vs Age by Diabetes Status")
```



5.2 Train a logistic regression model when the threshold $r = 0.5$

```
formula <- Diabetes ~ Age + Plasma_glucose
diabetes$Diabetes <- as.factor(diabetes$Diabetes)
gml_model <- caret::train(formula, data = diabetes, method = "glm", family = "binomial")
#type = "prob" predict probability
#type = "raw" predict the raw value/ class
#diabetes_pred <- predict(gml_model, type = "prob")

classify_pred_res <- function(r,gml_model) {

  diabetes_pred <- predict(gml_model, type = "prob")
  diabetes_pred$predict <- lapply(1:nrow(diabetes_pred),
    function(x) ifelse(diabetes_pred[x,2] > r, 1, 0))
  diabetes_pred$predict <- unlist(diabetes_pred$predict)
  diabetes_pred$raw <- diabetes$Diabetes
  diabetes_pred[, 3:4] <- lapply(diabetes_pred[, 3:4], as.factor)

  trainingData <- gml_model$trainingData %>% select(-.outcome)
  diabetes_pred <- cbind(diabetes_pred, trainingData)

  diabetes_pred$Age <- gml_model$trainingData$Age
  diabetes_pred$Plasma_glucose <- gml_model$trainingData$Plasma_glucose
  return(diabetes_pred)
}
```

```
diabetes_pred <- classify_pred_res(0.5,gml_model)

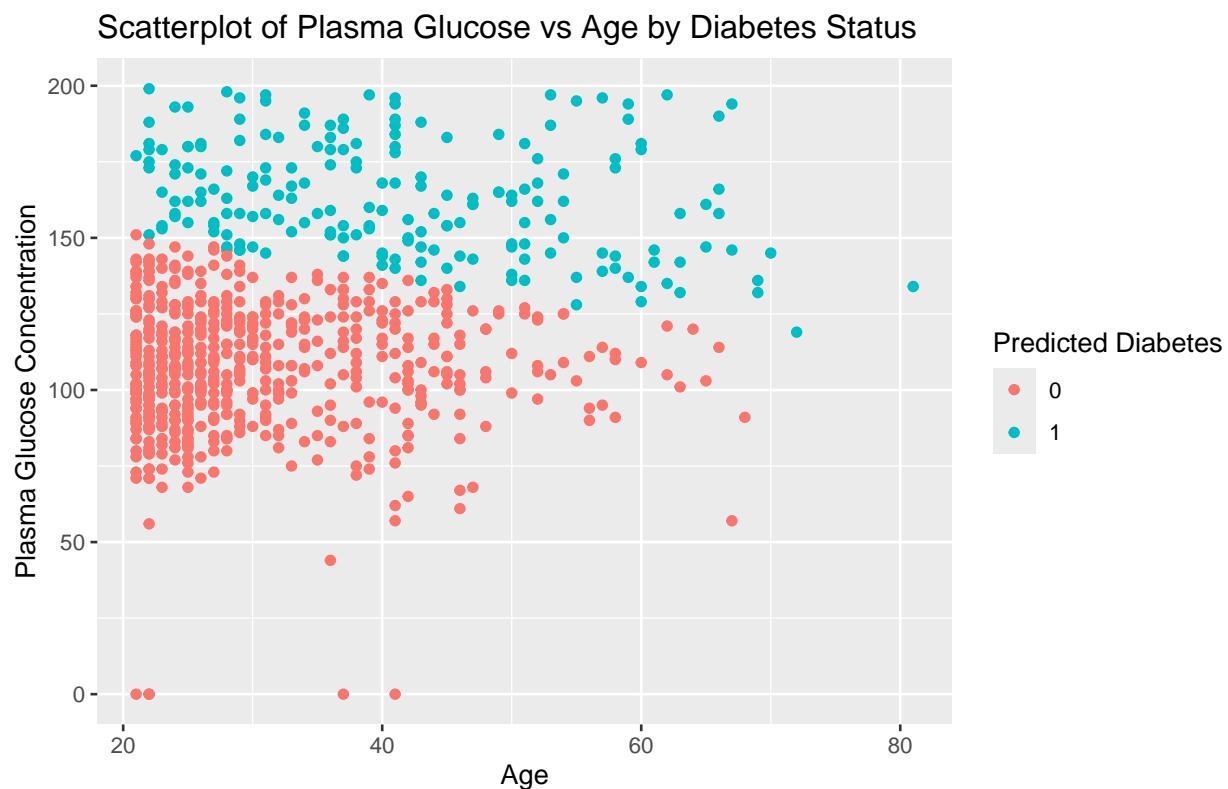
diabetes_confusion <- table(diabetes_pred$raw, diabetes_pred$predict)
error_rate <- 1 - (sum(diag(diabetes_confusion)) / sum(diabetes_confusion))
cat(" training misclassification error:",error_rate)
```

```
## training misclassification error: 0.2630208
```

5.3 Draw a scatter plot showing the predicted diabetes status

we can see that the logistic regression visually separates the two classes of diabetes status well, but the mis classification error is high due to the overlap of the two classes, maybe change can improve the performance.we will try later

```
ggplot(diabetes,aes( x = diabetes_pred$Age, y = diabetes_pred$Plasma_glucose,
                    color = diabetes_pred$predict)) + geom_point()+
labs(x = "Age", y = "Plasma Glucose Concentration", color = "Predicted Diabetes") +
ggtitle("Scatterplot of Plasma Glucose vs Age by Diabetes Status")
```



5.4 Draw a decision boundary between the two predicted classes

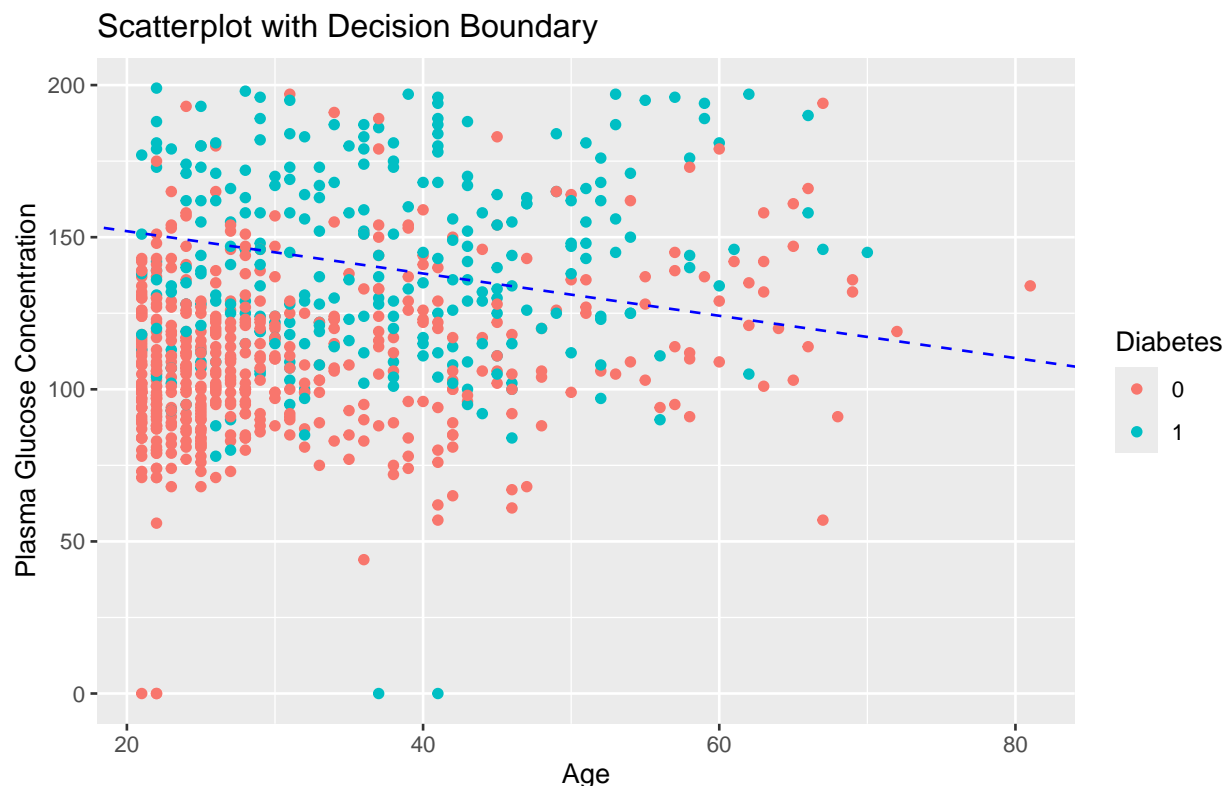
we can see that the boundary line try to split the dots into two classes and put the most the red dots below the line and the most blue dots above the line, but when the age exceed the 50, the performance of the model is not good.it seems that the number of red dot below the boundary line is same as the number above the line, it results the high misclassification error

```

get_boundary_line <- function(gml_model,r, y_name) {
  coefficients <- gml_model$finalModel$coefficients
  boundary_parameter <- list()
  coef_names <- names(coefficients)
  y_value <- coefficients[[y_name]]
  boundary_parameter$Intercept <- -(coefficients[['(Intercept)']] / y_value) - (log((1/r) - 1)/y_value)
  # boundary_parameter$intercept <- intercept
  for (name in coef_names){
    if (name != '(Intercept)'){
      boundary_parameter[[name]] <- -coefficients[[name]] / y_value
    }
  }
  return(boundary_parameter)
}

boundary_parameter <- get_boundary_line(gml_model,0.5,'Plasma_glucose')
ggplot(diabetes,aes( x = diabetes$Age, y = diabetes$Plasma_glucose, color = diabetes$Diabetes)) +
  geom_point()+
  geom_abline(slope = boundary_parameter$Age, intercept = boundary_parameter$Intercept,color = "blue", lty = 2) +
  labs(x = "Age", y = "Plasma Glucose Concentration", color = "Diabetes") +
  ggtitle("Scatterplot with Decision Boundary")

```



5.5 Change the thresholds r to 0.2 , 0.8 to see the what happened

we can see that when $r = 0.2$, for the red dots , its TP is relatively high , but the Recall is low ,for the blue dots,its TP is lower than red dots, but the Recall is higher than red dots, it means that the model is

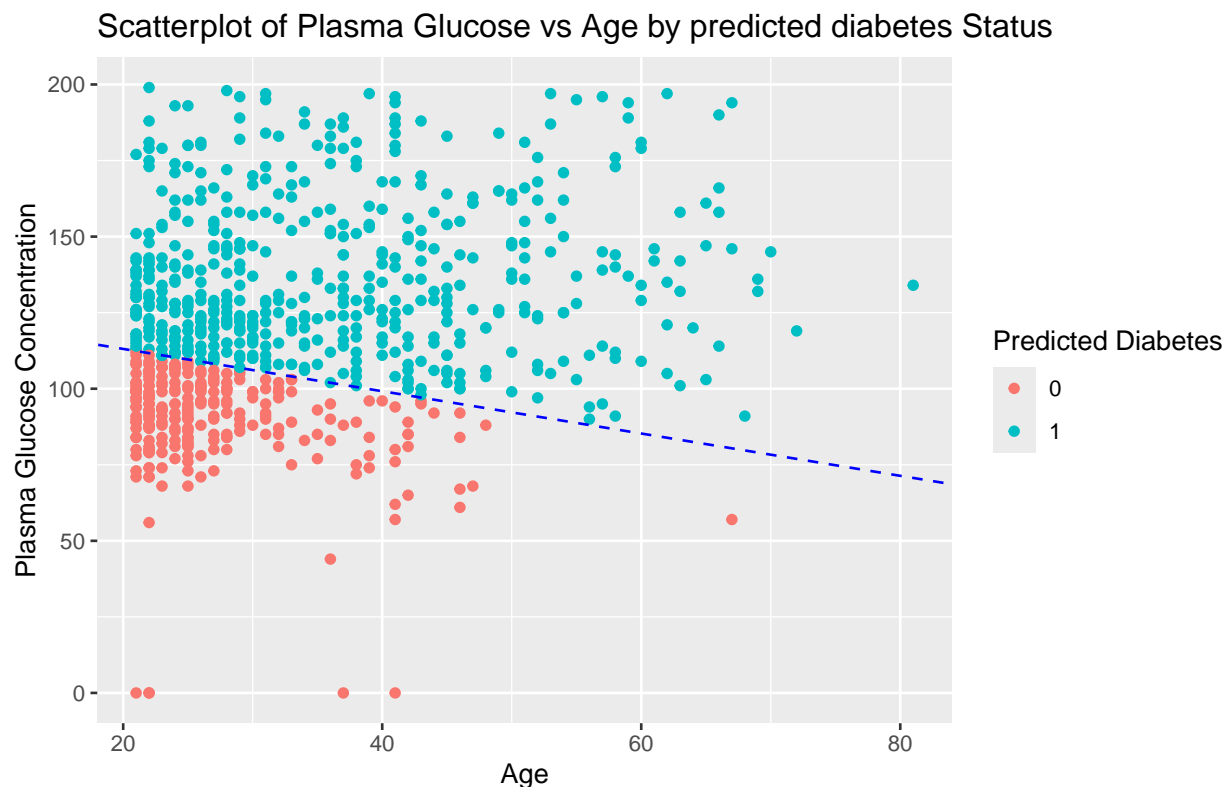
more likely to predict the blue dots as the positive class, but the blue dots are more likely to be the negative class, it results in the high misclassification error, when $r = 0.8$, the model is more likely to predict the red dots as the positive class, but the red dots are more likely to be the negative class, it results in the high misclassification error

```
pred_res_0.2 <- classify_pred_res(0.2,gml_model)
pred_res_0.8 <- classify_pred_res(0.8,gml_model)
```

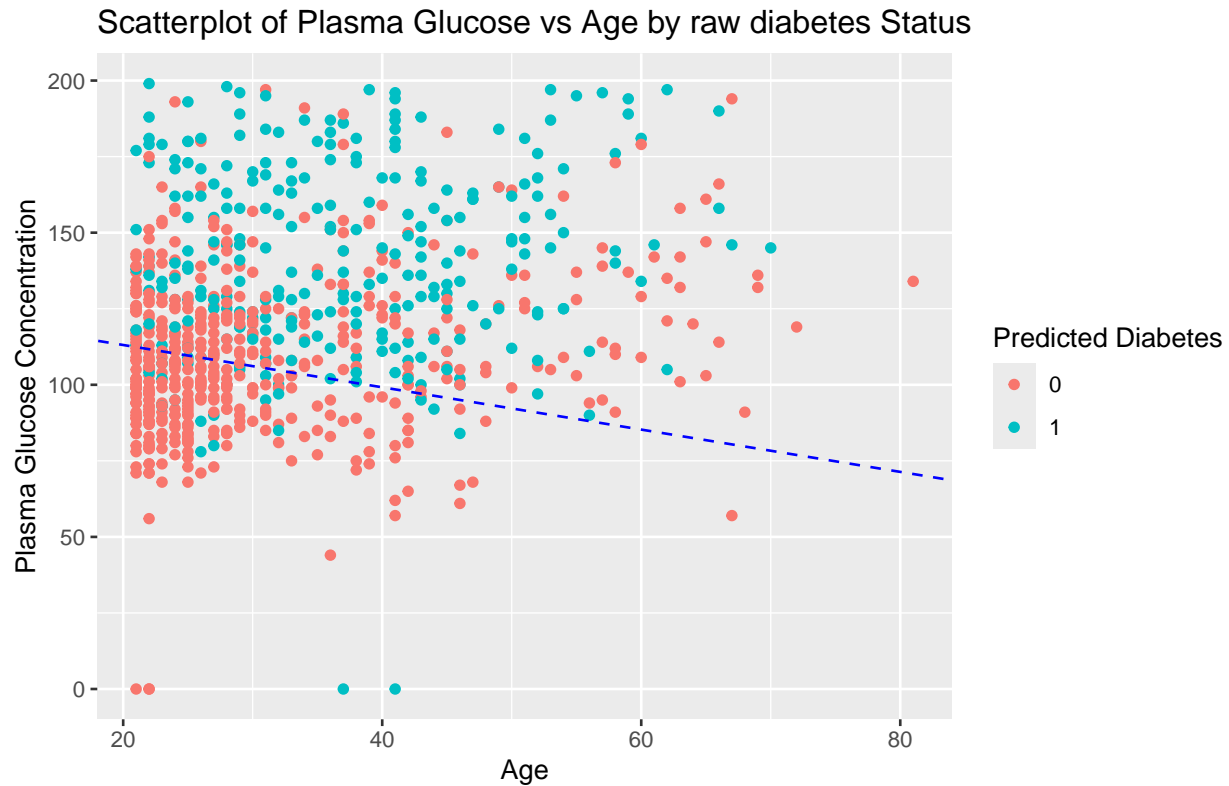
5.5.1 plot the scatter when $r = 0.2$

```
boundary_parameter_0.2 <- get_boundary_line(gml_model,0.2,'Plasma_glucose')

ggplot(diabetes,aes( x = pred_res_0.2$Age, y = pred_res_0.2$Plasma_glucose,
  color = pred_res_0.2$predict)) + geom_point()+
geom_abline(slope = boundary_parameter_0.2$Age,
  intercept = boundary_parameter_0.2$Intercept, color = "blue", linetype = "dashed") +
labs(x = "Age", y = "Plasma Glucose Concentration", color = "Predicted Diabetes") +
ggtitle("Scatterplot of Plasma Glucose vs Age by predicted diabetes Status")
```



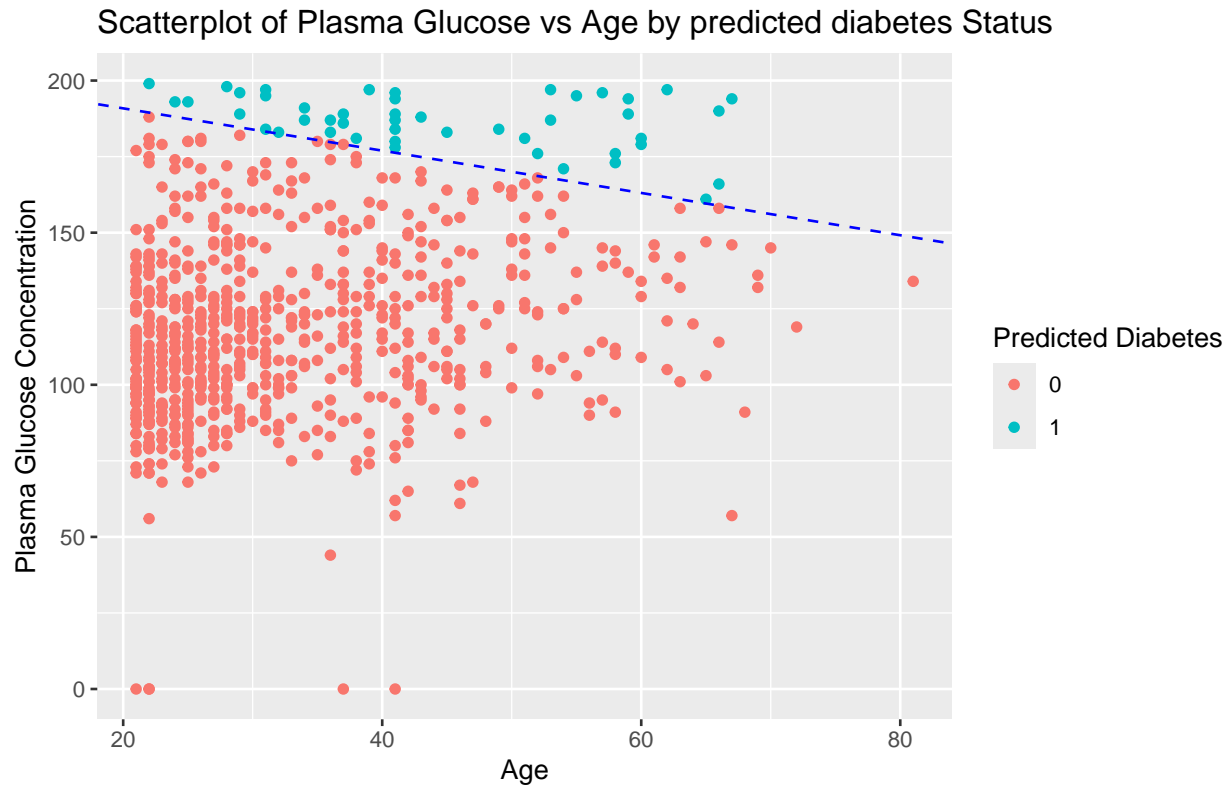
```
ggplot(diabetes,aes( x = pred_res_0.2$Age, y = pred_res_0.2$Plasma_glucose,
  color = pred_res_0.2$raw)) + geom_point()+
geom_abline(slope = boundary_parameter_0.2$Age,
  intercept = boundary_parameter_0.2$Intercept,
  color = "blue", linetype = "dashed") +
labs(x = "Age", y = "Plasma Glucose Concentration", color = "Predicted Diabetes") +
ggtitle("Scatterplot of Plasma Glucose vs Age by raw diabetes Status")
```



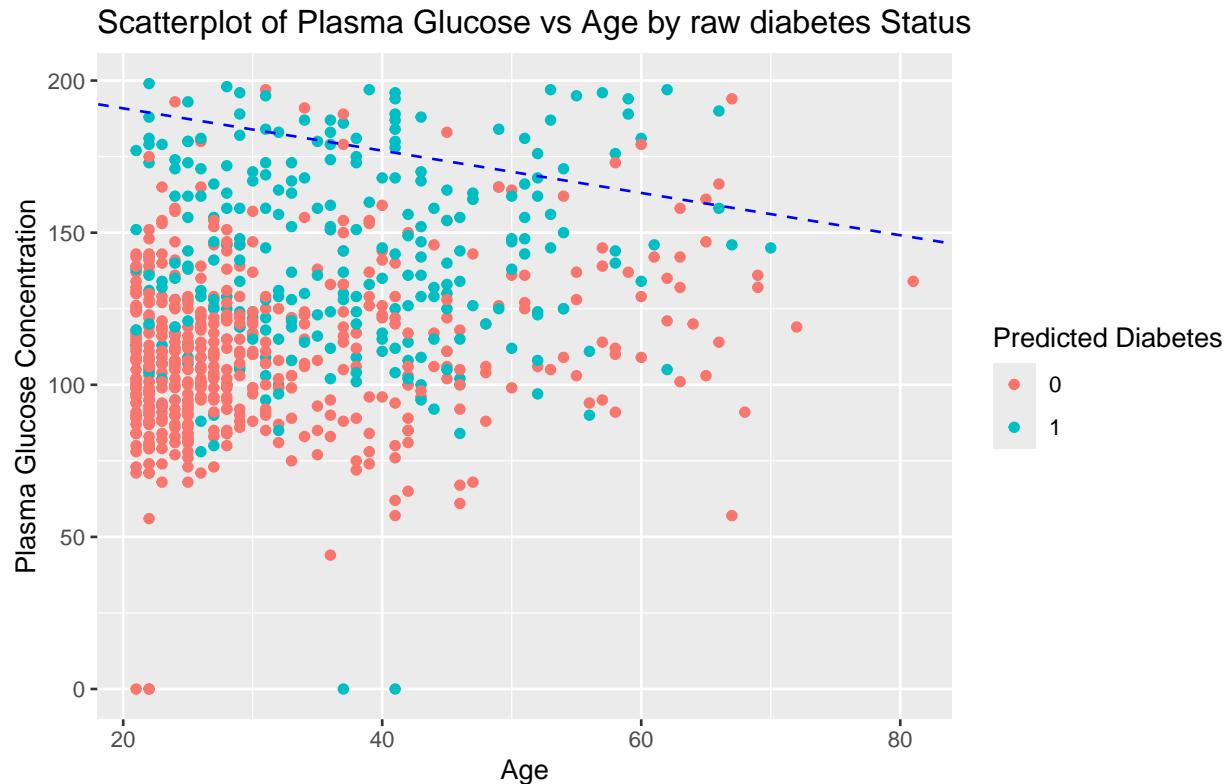
5.5.2 plot the scatter when $r = 0.8$

```
boundary_parameter_0.8 <- get_boundary_line(gml_model,0.8,'Plasma_glucose')

ggplot(diabetes,aes( x = pred_res_0.8$Age, y = pred_res_0.8$Plasma_glucose,
                    color = pred_res_0.8$predict)) + geom_point()+
geom_abline(slope = boundary_parameter_0.8$Age,
            intercept = boundary_parameter_0.8$Intercept,
            color = "blue", linetype = "dashed") +
labs(x = "Age", y = "Plasma Glucose Concentration", color = "Predicted Diabetes") +
ggtitle("Scatterplot of Plasma Glucose vs Age by predicted diabetes Status")
```



```
ggplot(diabetes,aes( x = pred_res_0.8$Age, y = pred_res_0.8$Plasma_glucose,
                    color = pred_res_0.8$raw)) + geom_point()+
geom_abline(slope = boundary_parameter_0.8$Age,
            intercept = boundary_parameter_0.8$Intercept,
            color = "blue", linetype = "dashed") +
labs(x = "Age", y = "Plasma Glucose Concentration", color = "Predicted Diabetes") +
ggtitle("Scatterplot of Plasma Glucose vs Age by raw diabetes Status")
```



5.5.3 Perform a basis function expansion trick

we can see that after add the basis function expansion, the misclassification error is lower than the previous model, it means that the basis function expansion can improve the performance of the model, look at the coefficients furtherly, the new added variables slightly affect the prediction, it means that the new added variables affect the prediction positively and the decision boundary become from a line to a multidimensional graphics

```
diabetes$z1 <- diabetes$Plasma_glucose^4
diabetes$z2 <- diabetes$Plasma_glucose^3 * diabetes$Age
diabetes$z3 <- diabetes$Plasma_glucose^2 * diabetes$Age^2
diabetes$z4 <- diabetes$Plasma_glucose * diabetes$Age^3
diabetes$z5 <- diabetes$Age^4
formula <- Diabetes ~ Age + Plasma_glucose + z1 + z2 + z3 + z4 + z5
new_gml_model <- caret::train(formula, data = diabetes, method = "glm", family = "binomial")
new_pred_res <- classify_pred_res(0.5,new_gml_model)
new_diabetes_confusion <- table(new_pred_res$raw, new_pred_res$predict)
error_rate <- 1 - (sum(diag(new_diabetes_confusion)) / sum(new_diabetes_confusion))
cat(" training misclassification error:",error_rate)
```

```
## training misclassification error: 0.2447917
```

```
new_boundary_parameter <- get_boundary_line(new_gml_model,0.5,'Plasma_glucose')
cat(new_gml_model$finalModel$coefficients)
```

```
## -9.309821 0.1456805 0.03793014 1.278015e-08 -1.7796e-07 8.51515e-07 -1.698011e-06 8.126623e-07
```


6 Assignment 4: Handwritten digit recognition with K-nearest neighbors

- Why can it be important to consider various probability thresholds in the classification problems, according to the book?

Probability thresholds serve as the reference point for evaluating performance of the model. Usually, a baseline is defined to indicate the model's worst performance level. And the achievable performance is defined by the maximum performance level. (Page 290 Baseline and Achievable Performance Level)

- What ways of collecting correct values of the target variable for the supervised learning problems are mentioned in the book?

In supervised learning problems, the target variables can be manually labelled by a domain expert. Target variables can also be labelled from predictive models based. Or the output is labelled naturally during the collection of training data. (Page 6, paragraph 2)

- How can one express the cost function of the linear regression in the matrix form, according to the book? The cost function for the linear regression model can be written with matrix notations as (Page 41):

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (\hat{y}(x_i; \boldsymbol{\theta}) - y_i)^2 = \frac{1}{n} \|\hat{y} - y\|_2^2 = \frac{1}{n} \|X\boldsymbol{\theta} - y\|_2^2 = \frac{1}{n} \|\boldsymbol{\epsilon}\|_2^2$$