

# Machine Learning Computer Lab 1 (Group A7)

Qinyuan Qi(qinqi464)      Satya Sai Naga Jaya Koushik Pilla (satpi345)  
Daniele Bozzoli(danbo826)

2023-11-19

## Assignment 1: Handwritten digit recognition with K- nearest neighbors (Solved by Qinyuan Qi - qinqi464)

### Answer:

(1) The following code will input the data and divide the data into training, validation and test sets.

```
##### Assignment 1.1 #####  
# read data  
data <- read.csv("optdigits.csv")  
row_num <- nrow(data)  
cols_num <- ncol(data)  
  
# set the last column name as "label_value" since we need to create a formula later  
names(data)[cols_num] <- "label_value"  
  
# set data split ratio to 0.5, 0.25 and 0.25  
ratio <- c(train = .5, test = .25, validate = .25)  
  
# data pre-processing  
# columns 1-64 are number based and do not need to be normalized, last column  
# is integer represent number from 0-9 which don't need to process again  
  
# set random seed  
set.seed(12345)  
  
# split data to training, test and validation set  
train_id <- sample(1:row_num, floor(row_num * ratio[1]))  
train_set <- data[train_id, ]  
  
set.seed(12345)  
test_val_id <- setdiff(1:row_num, train_id)  
test_id <- sample(test_val_id, floor(row_num * ratio[2]))  
test_set <- data[test_id, ]  
  
valid_id <- setdiff(test_val_id, test_id)  
valid_set <- data[valid_id, ]
```

(2)

The following code contain a function call to knn with k=30 and kernel = "rectangular". We calculate the confusion matrices and the misclassification errors. The code is attached in the appendix.

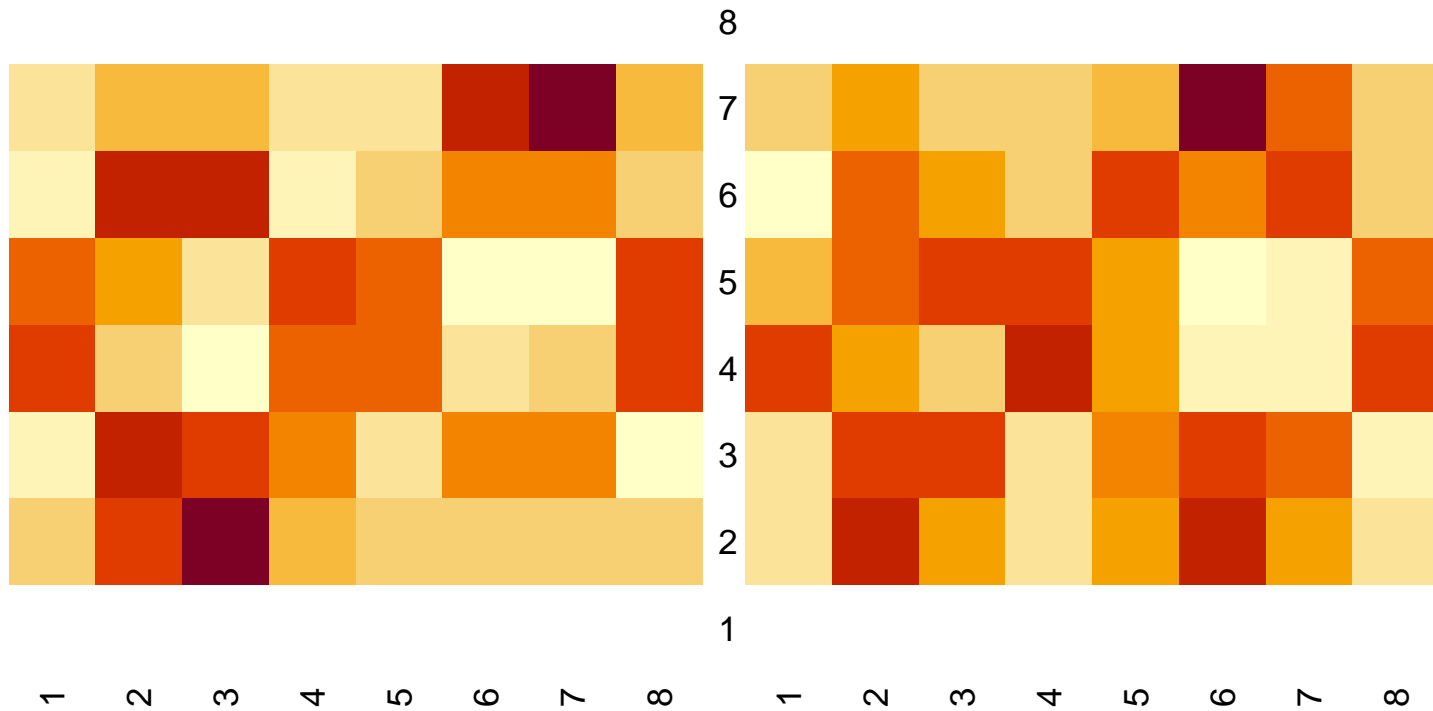
The confusion\_matrices is as follows:

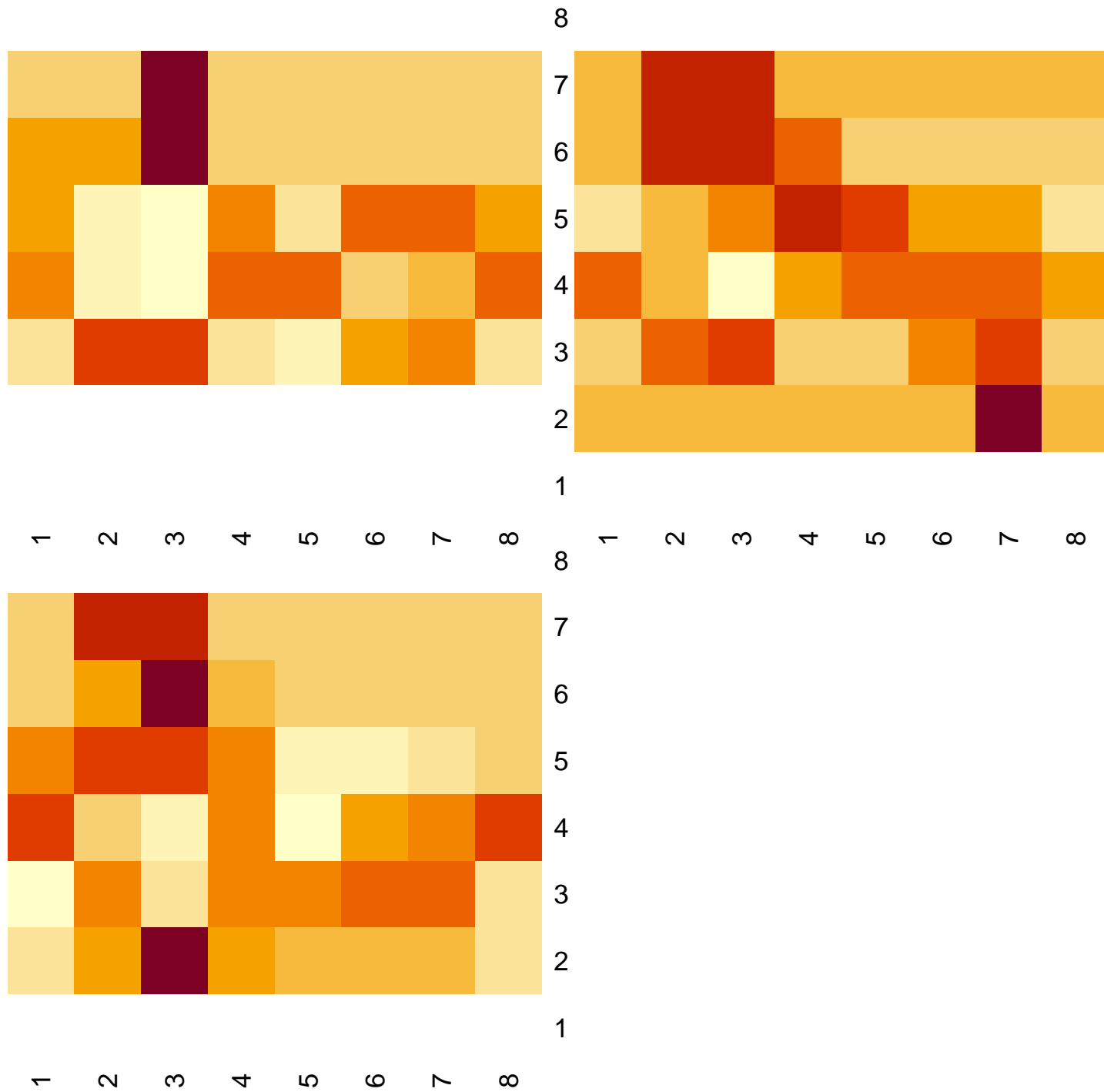
	0	1	2	3	4	5	6	7	8	9
0	98	0	0	0	0	0	0	0	0	0
1	0	101	0	0	0	0	1	0	2	0
2	1	0	108	1	0	1	0	0	0	0
3	0	0	0	80	0	0	0	0	0	0
4	0	1	1	1	99	0	1	0	2	0
5	0	0	2	1	0	83	0	0	1	1
6	0	1	1	6	7	1	81	0	2	0
7	0	0	0	0	1	5	0	96	0	3
8	0	1	0	1	0	1	0	1	77	4
9	0	0	0	1	1	1	0	0	0	77

As what the data show in the confusion matrix, we can find that the accuracy of number 0,3 are super high. For most of the numbers, the accuracy are still very. And we got a total accuracy rate of 0.9424084, and total error rate is 0.05759162.

**(3)** We use the following code to generate heatmaps. By show 205 heat images one by one, we found first 2 images(index 18,50) looks like 8, but next 3 images(index 1, 2,55) not like 8 at all.

The image is as follows. The code attached as appendix.

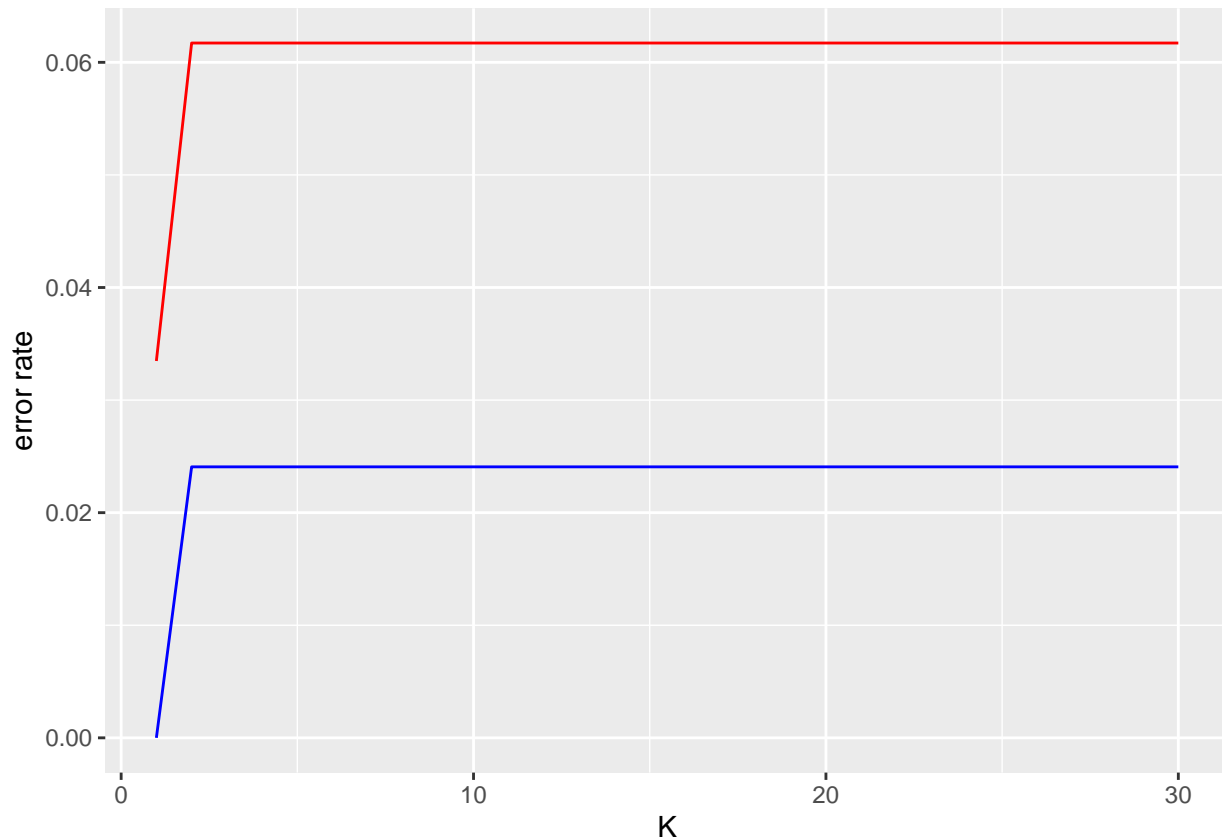




(4)

We fit the KNN with  $k = 1$  to 30, we plot the error rates as follows, the code is attached in appendix.

When  $K$  increase, the model not become very complex, it seems it does not change a lot when  $k \geq 2$ . And according to the graph, we can know that the optimal  $K$  is 1.



(5)

The cross-entropy is defined as  $-\sum_{i=1}^N \sum_{m=1}^M I(Y_i = C_m) * \log \hat{p}(y_i = C_m)$ .

So we will calculate it using the following code(Not implemented).

The reason why we use cross-entropy is more suitable is: 1) It is sensitive to the predicted probabilities. 2) In multinomial distribution, the predicted value will more likely to be explained as a probability, and cross-entropy loss function is designed to compare the probability based prediction.

## Assignment 2: Linear regression and ridge regressions (Solved by Satya Sai Naga Jaya Koushik Pilla)

**Answer:**

(1)

We read the data and divide data into training and test data (60/40) and normalize them.

```
##### Assignment 2.1 #####
# Load the data
data <- read.csv("parkinsons.csv")
row_num <- nrow(data)
cols_num <- ncol(data)

# Divide the data into training and test data (60/40)
# set data split ratio
ratio <- c(train = .6, test = .4)

# set random seed
```

```

set.seed(12345)

# split data
train_id <- sample(1:row_num, floor(row_num * ratio[1]))
train_set <- data[train_id, ]
test_id <- setdiff(1:row_num, train_id)
test_set <- data[test_id, ]

# normalize data.
train_set <- as.data.frame(lapply(train_set, normalize))
test_set <- as.data.frame(lapply(test_set, normalize))

```

(2) We create a linear model using following code, we get the parameters as show below.

```

##### Assignment 2.2 #####
# apply linear regression
model <- lm(motor_UPDRS ~ ., data = train_set)

# predict the test value using the linear regression just created
test_pred <- predict(model, test_set)

# calculate test MSE
test_mse <- mean((test_pred - test_set$motor_UPDRS)^2)
model

##
## Call:
## lm(formula = motor_UPDRS ~ ., data = train_set)
##
## Coefficients:
## (Intercept)      subject.          age          sex      test_time
##    0.041153    -0.020612    -0.045811     0.029830    -0.002456
## total_UPDRS  Jitter...  Jitter.Abs.  Jitter.RAP  Jitter.PPQ5
##    1.012833     0.679831    -0.231111     3.062194    -0.149610
## Jitter.DDP      Shimmer  Shimmer.dB.  Shimmer.APQ3  Shimmer.APQ5
##   -3.477175     0.404012    -0.059469    24.111950    -0.309738
## Shimmer.APQ11  Shimmer.DDA          NHR          HNR          RPDE
##    0.253228    -24.256718     0.004315    -0.008594    -0.054478
##          DFA          PPE
##   -0.007923     0.100989

```

And we get test MSE = 0.005370424

we know that Shimmer.APQ3 and Shimmer.DDA contribute significantly to the model.

(3)

```

##### Assignment 2.3 #####
# Loglikelihood function
loglikelihood <- function(theta, sigma) {
  n <- length(Y)
  log_likelihood_values <- -0.5 * (log(2 * pi * sigma^2) + ((Y - theta %*% X) / sigma)^2)
  return(total_log_likelihood <- sum(log_likelihood_values))
}

# ridge
ridge <- function(param) {

```

```

param_length <- length(param)
theta <- param[1:param_length-2]
sigma <- param[param_length-1]
lambda <- param[param_length]
loglikelihood_result <- loglikelihood(theta, sigma)
ridge_penalty <- lambda * sum(theta^2)
return(loglikelihood_result - ridge_penalty)
}

# ridgeopt
ridgeopt <- function(theta, sigma, lambda) {
  size_theta <- length(theta)
  param <- rep(0, size_theta+2)
  for(i in 1:size_theta){
    param[i] <- theta[i]
  }
  param[size_theta + 1] <- sigma
  param[size_theta + 2] <- lambda
  ridge_result <- optim(par = param, fn = ridge, method="BFGS")
}

# df
df <- function(X, lambda) {
  n <- nrow(X)
  # calculate hat_matrix
  hat_matrix <- X %*% solve(t(X) %*% X + lambda * diag(ncol(X)))
  # get degree of the hat_matrix
  df_ridge <- sum(diag(hat_matrix))
  return(df_ridge)
}

```

(4)

```

##### Assignment 2.4 #####

#sigma <- 0.01
#theta <- rep(1,21)
#X <- train_set[]
#lambda <- 1
#ridgeopt(theta, sigma, lambda)

#lambda <- 100
#ridgeopt(theta, sigma, lambda)

#lambda <- 1000
#ridgeopt(theta, sigma, lambda)

```

### Assignment 3. Logistic regression and basis function expansion (Solved by Daniele Bozzoli)

**Answer:**

(1) The plot was generated as following.

As can be seen on the plot, we can know that diabetes is easy to classify by a standard logistic regression model by just using 2 features.

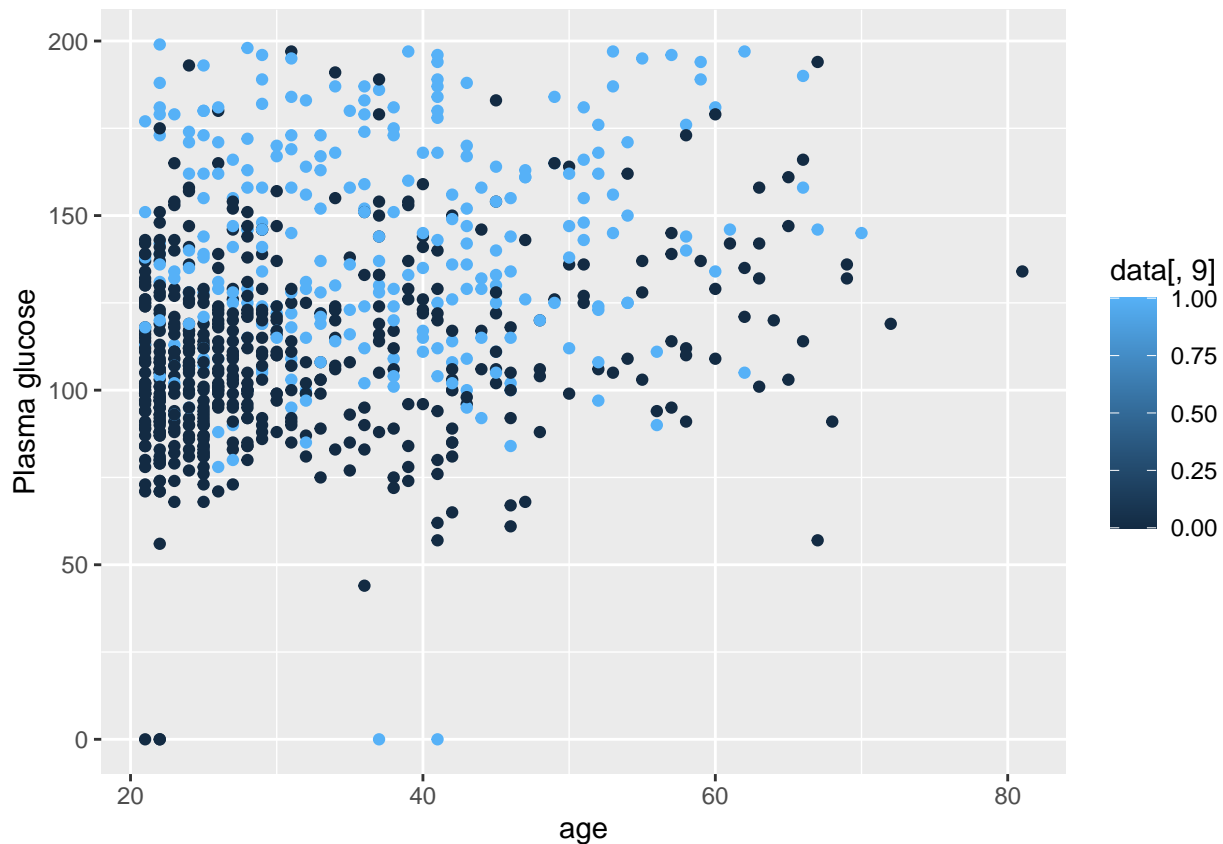
```
data <- read.csv("pima-indians-diabetes.csv")

age <- data[, 8]
plasma <- data[, 2]

names(data)[9] <- "diabetes"
names(data)[2] <- "plasma"
names(data)[8] <- "age"

x <- age
y <- plasma
df <- data.frame(x, y)

ggplot2::ggplot(df, ggplot2::aes(age, plasma)) +
  ggplot2::geom_point(ggplot2::aes(colour = data[, 9]))+
  ggplot2::labs(x = "age", y="Plasma glucose")
```



(2) The logistic regression model was trained using glm function as following.

```
model <- glm(diabetes ~ plasma + age, family = gaussian, data = data)

summary(model)

##
## Call:
## glm(formula = diabetes ~ plasma + age, family = gaussian, data = data)
```

```
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.5984068  0.0657429  -9.102  < 2e-16 ***
## plasma      0.0064646  0.0004901  13.191  < 2e-16 ***
## age         0.0049734  0.0013335   3.729  0.000206 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 0.175189)
##
##    Null deviance: 174.05  on 766  degrees of freedom
## Residual deviance: 133.84  on 764  degrees of freedom
## AIC: 845.62
##
## Number of Fisher Scoring iterations: 2
```

(3)

(4)

(5)



## Appendix: All code for this report

```
##### Init code #####
rm(list = ls())
knitr::opts_chunk$set(echo = TRUE)
library(kknn)
##### Assignment 1.1 #####
# read data
data <- read.csv("optdigits.csv")
row_num <- nrow(data)
cols_num <- ncol(data)

# set the last column name as "label_value" since we need to create a formula later
names(data)[cols_num] <- "label_value"

# set data split ratio to 0.5, 0.25 and 0.25
ratio <- c(train = .5, test = .25, validate = .25)

# data pre-processing
# columns 1-64 are number based and do not need to be normalized, last column
# is integer represent number from 0-9 which don't need to process again

# set random seed
set.seed(12345)

# split data to training, test and validation set
train_id <- sample(1:row_num, floor(row_num * ratio[1]))
train_set <- data[train_id, ]

set.seed(12345)
test_val_id <- setdiff(1:row_num, train_id)
test_id <- sample(test_val_id, floor(row_num * ratio[2]))
test_set <- data[test_id, ]

valid_id <- setdiff(test_val_id, test_id)
valid_set <- data[valid_id, ]
##### Assignment 1.2 #####
# calculate error rate for k = 30
k_value <- 30
train_kknn <- train.kknn(label_value ~ .,
                        data = train_set,
                        kmax = k_value,
                        kernel = "rectangular")

predict_data <- predict(train_kknn, newdata=test_set)

# generate confusion matrix
confusion_matrices <- table(round(predict_data), test_set$label_value)

# calculate accuracy
accuracy <- sum(diag(confusion_matrices)) / sum(confusion_matrices)

# calculate error rate
error_rate <- 1 - accuracy
```

```
##### Assignment 1.3 #####
# Find all the 8 images
eight_images <- train_set[train_set$label_value == 8, ]

# get the number of 8 images in training data, we got 205 images
eight_images_number <- nrow(eight_images)

image_number <- 18
eight_images_matrix <- matrix(c(as.numeric(eight_images[image_number, 1:64])), ncol = 8)
heatmap(x = eight_images_matrix, Rowv = NA, Colv = NA)

image_number <- 50
eight_images_matrix <- matrix(c(as.numeric(eight_images[image_number, 1:64])), ncol = 8)
heatmap(x = eight_images_matrix, Rowv = NA, Colv = NA)

image_number <- 1
eight_images_matrix <- matrix(c(as.numeric(eight_images[image_number, 1:64])), ncol = 8)
heatmap(x = eight_images_matrix, Rowv = NA, Colv = NA)

image_number <- 2
eight_images_matrix <- matrix(c(as.numeric(eight_images[image_number, 1:64])), ncol = 8)
heatmap(x = eight_images_matrix, Rowv = NA, Colv = NA)

image_number <- 55
eight_images_matrix <- matrix(c(as.numeric(eight_images[image_number, 1:64])), ncol = 8)
heatmap(x = eight_images_matrix, Rowv = NA, Colv = NA)

##### Assignment 1.4 #####
# calculate error rate for different k values

error_rates_train <- rep(30, 0)
error_rates_validate <- rep(30, 0)

for(k_value in 1:30){
  # apply kkn function
  train_kknn <- train.kknn(label_value ~ .,
                           data = train_set,
                           kmax = k_value,
                           kernel = "rectangular")

  predict_data_train <- predict(train_kknn, newdata=train_set)
  predict_data_validation <- predict(train_kknn, newdata=valid_set)

  # generate confusion matrix
  confusion_matrices_train <- table(round(predict_data_train), train_set$label_value)
  confusion_matrices_validation <- table(round(predict_data_validation), valid_set$label_value)

  # calculate accuracy
  accuracy_train <- sum(diag(confusion_matrices_train)) / sum(confusion_matrices_train)
  accuracy_validation <- sum(diag(confusion_matrices_validation)) / sum(confusion_matrices_validation)

  # calculate error rate
  error_rates_train[k_value] <- 1 - accuracy_train
}
```

```

    error_rates_validate[k_value] <- 1 - accuracy_validation
  }

  # plot the error rate graph
  k <- 1:30

  error_rate_data <- data.frame(k, error_rates_train,error_rates_validate)
  ggplot2::ggplot() +
    ggplot2::geom_line(error_rate_data,mapping=ggplot2::aes(x=k,y=error_rates_train),colour="blue") +
    ggplot2::geom_line(error_rate_data,mapping=ggplot2::aes(x=k,y=error_rates_validate),colour="red") +
    ggplot2::labs(x = "K",y="error rate")
  #error_rates_cross <- rep(30, 0)

  #for(k_value in 1:30){
  #  # apply kknn function
  #  train_kknn <- train.kknn(label_value ~ .,
  #                           data = train_set,
  #                           kmax = k_value,
  #                           kernel = "rectangular")

  #  predict_data_test <- predict(train_kknn,newdata=test_set)

  #  confusion_matrices_test <- table(round(predict_data_test), test_set$label_value)

  #  calculate cross-entropy
  #ce <- 0
  #for(i in 1:10){
  #  total_count <- sum(confusion_matrices_test[i])
  #  for(j in 1:10){
  #    p <- confusion_matrices_test[i,j] / total_count
  #    print(p)
  #    ce <- ce + log(p) * confusion_matrices_test[i,j]
  #  }
  #}
  #error_rates_cross[k_value] <- ce * (-1)
#}

#k <- 1:30

#error_rate_data <- data.frame(k,error_rates_cross)
#ggplot2::ggplot() +
#  ggplot2::geom_line(error_rate_data,mapping=ggplot2::aes(x=k,y=error_rates_cross),colour="blue") +
#  ggplot2::labs(x = "K",y="error rate")

##### Init code #####
rm(list = ls())
knitr::opts_chunk$set(echo = TRUE)

##### Common Functions #####
# normalize data
normalize <- function(x) {
  return((x - min(x)) / (max(x) - min(x)))
}

```

```

}

##### Assignment 2.1 #####
# Load the data
data <- read.csv("parkinsons.csv")
row_num <- nrow(data)
cols_num <- ncol(data)

# Divide the data into training and test data (60/40)
# set data split ratio
ratio <- c(train = .6, test = .4)

# set random seed
set.seed(12345)

# split data
train_id <- sample(1:row_num, floor(row_num * ratio[1]))
train_set <- data[train_id, ]
test_id <- setdiff(1:row_num, train_id)
test_set <- data[test_id, ]

# normalize data.
train_set <- as.data.frame(lapply(train_set, normalize))
test_set <- as.data.frame(lapply(test_set, normalize))
##### Assignment 2.2 #####
# apply linear regression
model <- lm(motor_UPDRS ~ ., data = train_set)

# predict the test value using the linear regression just created
test_pred <- predict(model, test_set)

# calculate test MSE
test_mse <- mean((test_pred - test_set$motor_UPDRS)^2)
model

##### Assignment 2.3 #####
# Loglikelihood function
loglikelihood <- function(theta, sigma) {
  n <- length(Y)
  log_likelihood_values <- -0.5 * (log(2 * pi * sigma^2) + ((Y - theta %*% X) / sigma)^2)
  return(total_log_likelihood <- sum(log_likelihood_values))
}

# ridge
ridge <- function(param) {
  param_length <- length(param)
  theta <- param[1:param_length-2]
  sigma <- param[param_length-1]
  lambda <- param[param_length]
  loglikelihood_result <- loglikelihood(theta, sigma)
  ridge_penalty <- lambda * sum(theta^2)
  return(loglikelihood_result - ridge_penalty)
}

```

```

# ridgeopt
ridgeopt <- function(theta, sigma, lambda) {
  size_theta <- length(theta)
  param <- rep(0, size_theta+2)
  for(i in 1:size_theta){
    param[i] <- theta[i]
  }
  param[size_theta + 1] <- sigma
  param[size_theta + 2] <- lambda
  ridge_result <- optim(par = param, fn = ridge, method="BFGS")
}

# df
df <- function(X, lambda) {
  n <- nrow(X)
  # calculate hat_matrix
  hat_matrix <- X %*% solve(t(X) %*% X + lambda * diag(ncol(X)))
  # get degree of the hat_matrix
  df_ridge <- sum(diag(hat_matrix))
  return(df_ridge)
}

##### Assignment 2.4 #####

#sigma <- 0.01
#theta <- rep(1,21)
#X <- train_set[]
#lambda <- 1
#ridgeopt(theta, sigma, lambda)

#lambda <- 100
#ridgeopt(theta, sigma, lambda)

#lambda <- 1000
#ridgeopt(theta, sigma, lambda)

##### Init code #####
rm(list = ls())
knitr::opts_chunk$set(echo = TRUE)
data <- read.csv("pima-indians-diabetes.csv")

age <- data[, 8]
plasma <- data[, 2]

names(data)[9] <- "diabetes"
names(data)[2] <- "plasma"
names(data)[8] <- "age"

x <- age
y <- plasma
df <- data.frame(x, y)

ggplot2::ggplot(df, ggplot2::aes(age, plasma)) +

```

```
ggplot2::geom_point(ggplot2::aes(colour = data[, 9]))+  
ggplot2::labs(x = "age",y="Plasma glucose")  
  
model <- glm(diabetes ~ plasma + age, family = gaussian, data = data)  
  
summary(model)
```