

Machine Learning Lab1

Lepeng Zhang, Xuan Wang, Priyarani Patil

2023-11-08

Assignment 1. Handwritten digit recognition with Knearest neighbors.

1.

Answer:

2.

Answer:

Confusion Matrix for Training Data:

```
##
##      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
```

Confusion Matrix for Test Data:

```
##
##      0   1   2   3   4   5   6   7   8   9
## 0 82   0   0   0   1   0   1   0   0   0
## 1   0 90   2   0   0   0   0   0   0   3
## 2   0   1 92   0   0   0   0   1   1   1
## 3   0   0   0 85   0   2   0   3   1   1
## 4   0   1   0   0 89   0   1   6   0   5
## 5   0   1   0   1   0 97   1   1   0   7
## 6   0   0   0   0   0   0 97   0   0   0
## 7   0   1   0   1   0   0   0 99   0   0
## 8   0   7   0   0   0   0   0   0 84   0
## 9   0   2   0   0   0   0   0   2   1 86
```

Misclassification errors for the training data are: 0.04500262

Misclassification errors for the test data are: 0.0585162

The accuracy of prediction of for digit 0 is: 0.9761905

```

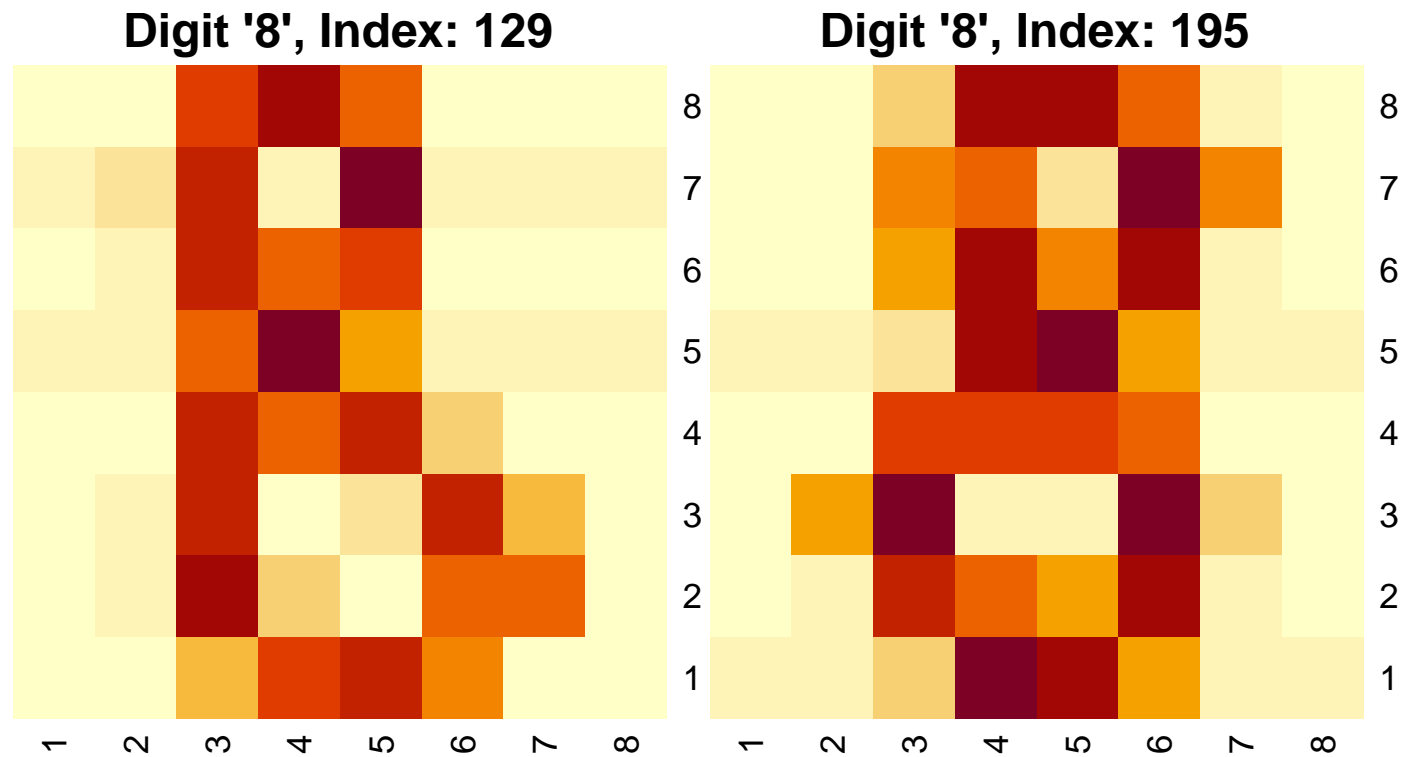
## The accuracy of prediction of for digit 1 is: 0.9473684
## The accuracy of prediction of for digit 2 is: 0.9583333
## The accuracy of prediction of for digit 3 is: 0.923913
## The accuracy of prediction of for digit 4 is: 0.872549
## The accuracy of prediction of for digit 5 is: 0.8981481
## The accuracy of prediction of for digit 6 is: 1
## The accuracy of prediction of for digit 7 is: 0.980198
## The accuracy of prediction of for digit 8 is: 0.9230769
## The accuracy of prediction of for digit 9 is: 0.9450549

## The overall accuracy of prediction is: 0.9414838

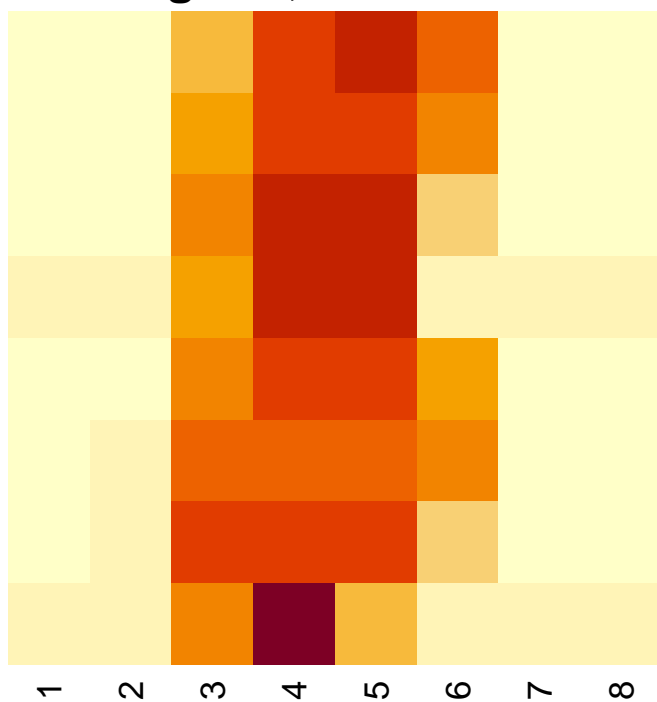
```

3.

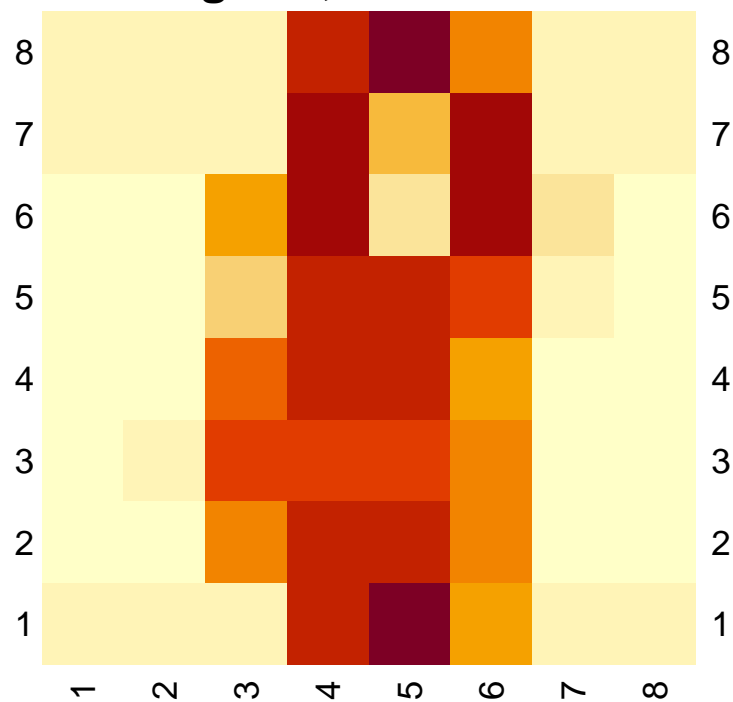
Answer:



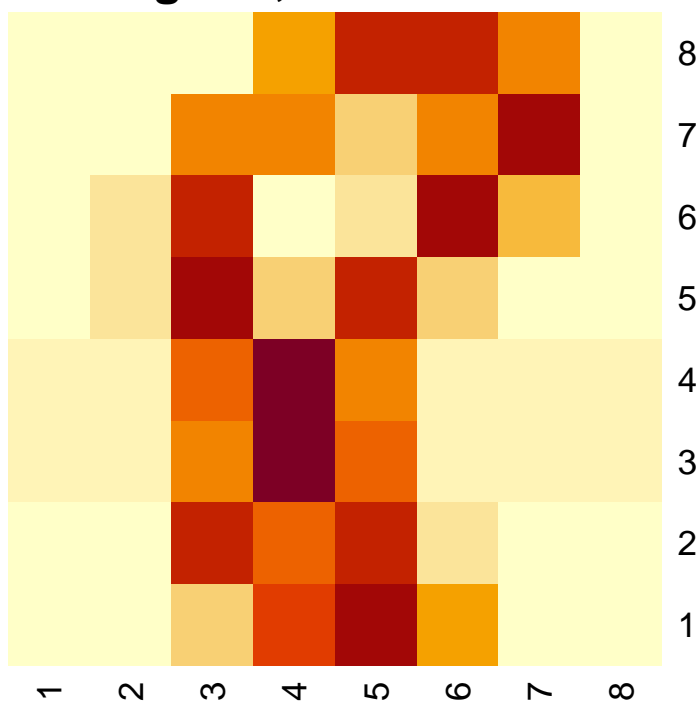
Digit '8', Index: 520



Digit '8', Index: 431



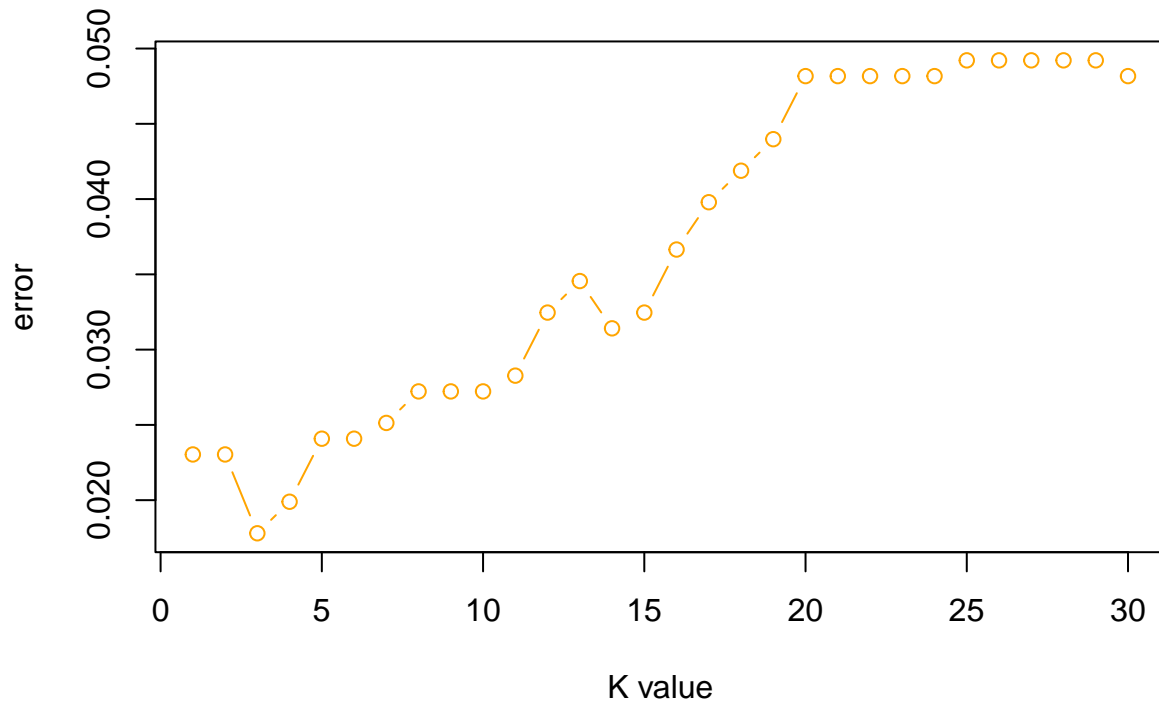
Digit '8', Index: 1294



4.

Answer:

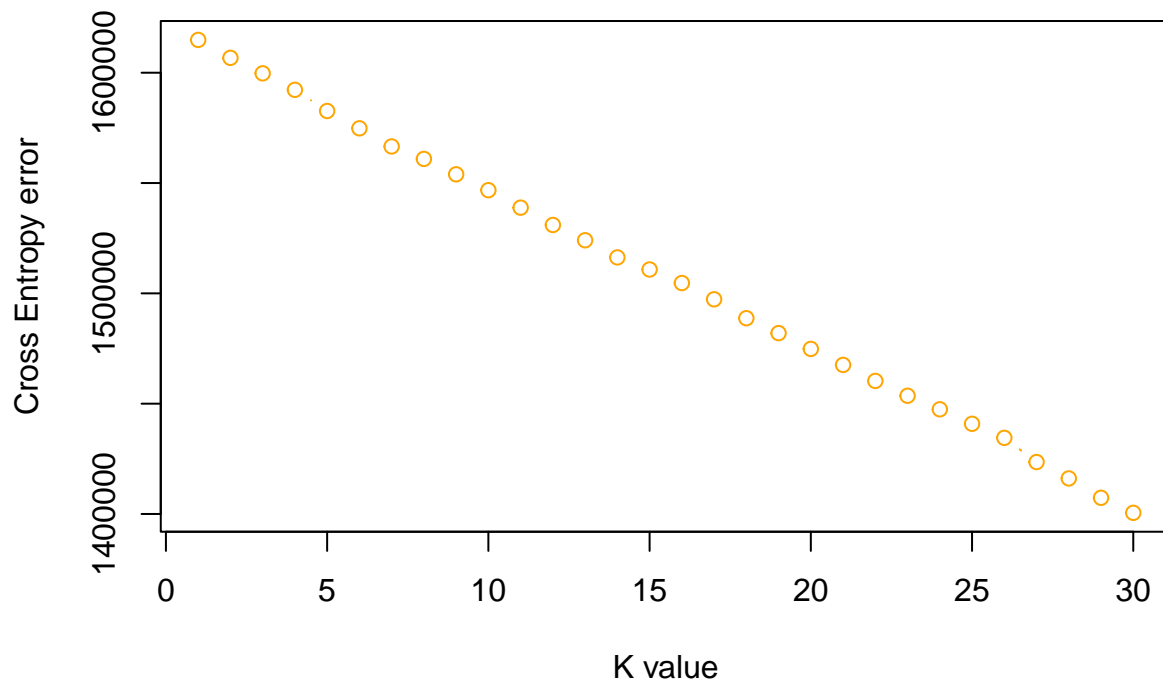
Misclassification errors on the value of K



5.

Answer:

Cross Entropy errors on the value of K



The optimal k is: 30

Assignment 2. Linear regression and ridge regression

1.

Answer:

2.

Answer:

```
## The training data's MSE is: 0.8785431
## The training data's MSE is: 0.9191113
##
## Call:
## lm(formula = motor_UPDRS ~ ., data = scaled_train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.0255 -0.7363 -0.1087  0.7333  2.1960
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  6.563e-15  1.583e-02   0.000 1.000000
## Jitter...    1.869e-01  1.496e-01   1.250 0.211496
## Jitter.Abs. -1.696e-01  4.081e-02 -4.156 3.32e-05 ***
## Jitter.RAP   -5.270e+00  1.884e+01  -0.280 0.779688
## Jitter.PPQ5  -7.457e-02  8.778e-02  -0.850 0.395659
## Jitter.DDP    5.250e+00  1.884e+01   0.279 0.780541
## Shimmer      5.924e-01  2.060e-01   2.876 0.004055 **
## Shimmer.dB.  -1.727e-01  1.393e-01  -1.239 0.215380
## Shimmer.APQ3  3.207e+01  7.717e+01   0.416 0.677738
## Shimmer.APQ5 -3.875e-01  1.138e-01  -3.405 0.000669 ***
## Shimmer.APQ11 3.055e-01  6.124e-02   4.989 6.37e-07 ***
## Shimmer.DDA  -3.239e+01  7.717e+01  -0.420 0.674739
## NHR          -1.854e-01  4.557e-02  -4.068 4.85e-05 ***
## HNR          -2.385e-01  3.640e-02  -6.553 6.45e-11 ***
## RPDE         4.068e-03  2.267e-02   0.179 0.857576
## DFA          -2.803e-01  2.014e-02 -13.919 < 2e-16 ***
## PPE          2.265e-01  3.289e-02   6.886 6.75e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.9396 on 3508 degrees of freedom
## Multiple R-squared:  0.1212, Adjusted R-squared:  0.1172
## F-statistic: 30.24 on 16 and 3508 DF,  p-value: < 2.2e-16
```

3.

Answer:

4.

Answer:

```
## Ridge Training MSE: 7.323041e+28 1.357493e+29 4.687817e+28
## Ridge Test MSE: 7.299556e+28 1.356537e+29 4.692501e+28
## Degrees of Freedom: 13.86074 9.924887 5.643925
```

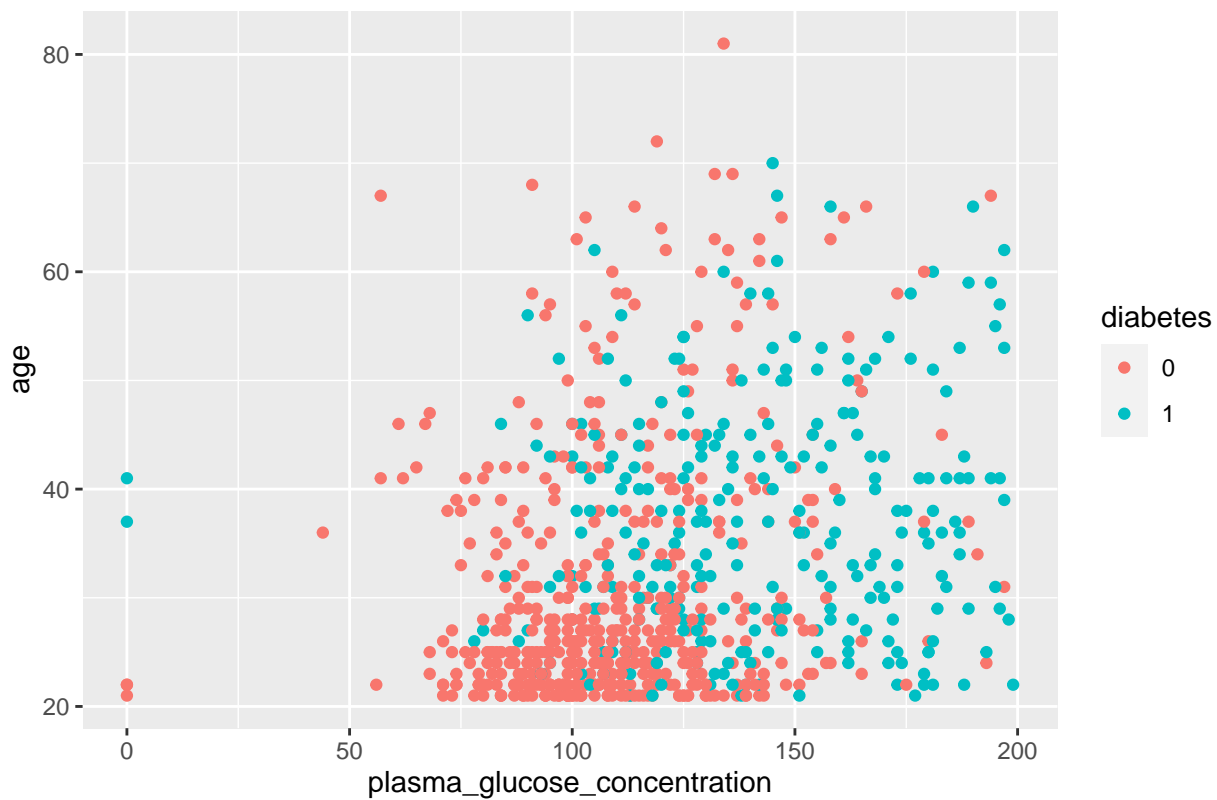
Assignment 3. Logistic regression and basis function expansion

1.

Answer:

```
## Loading required package: lattice
##
## Attaching package: 'caret'
## The following object is masked from 'package:kkn':
##
##   contr.dummy
```

Scatter Plot of Plasma Glucose Concentration on Age



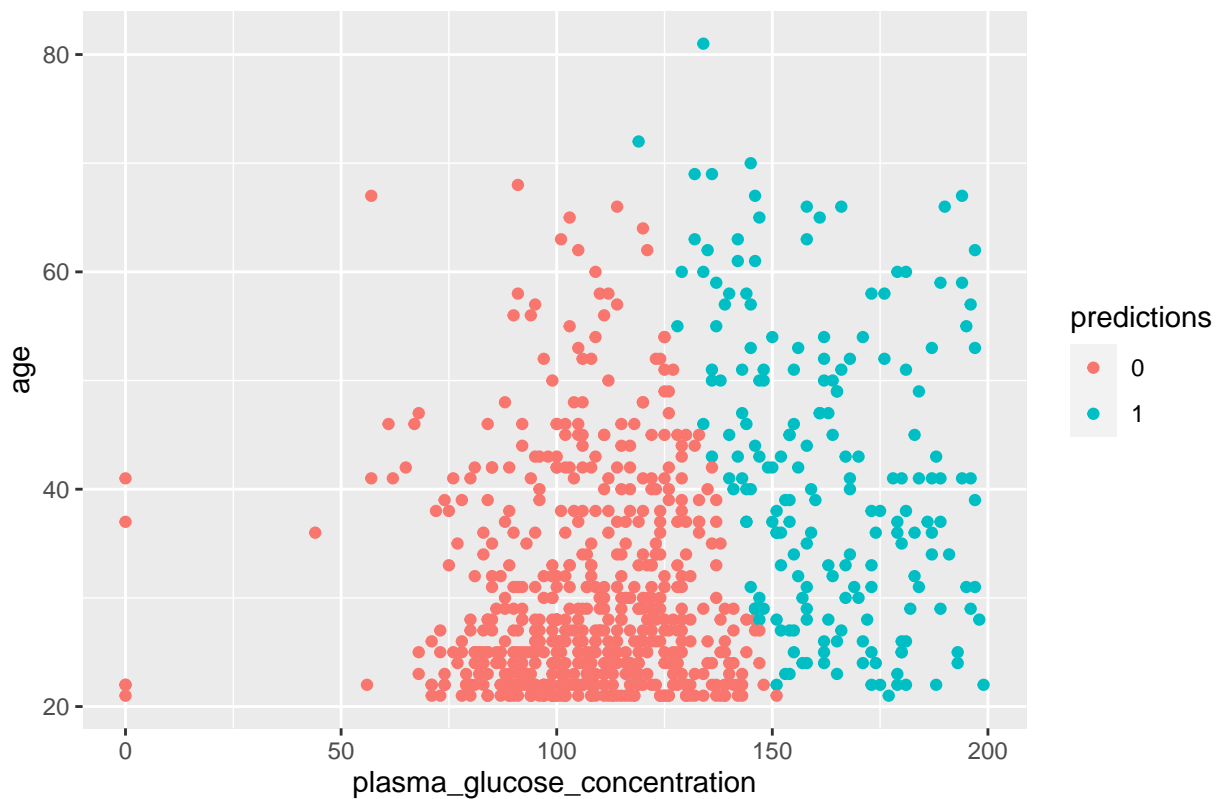
2.

Answer:

```
##
## Call:
## glm(formula = diabetes ~ plasma_glucose_concentration + age,
##     family = binomial, data = pima_data)
##
## Deviance Residuals:
##   Min       1Q   Median       3Q      Max
## -2.3367  -0.7775  -0.5087   0.8367   3.1630
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
```

```
## (Intercept)                -5.912449    0.462620   -12.78   < 2e-16 ***
## plasma_glucose_concentration 0.035644    0.003290    10.83   < 2e-16 ***
## age                        0.024778    0.007374     3.36  0.000778 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 993.48  on 767  degrees of freedom
## Residual deviance: 797.36  on 765  degrees of freedom
## AIC: 803.36
##
## Number of Fisher Scoring iterations: 4
## [1] "Probability(Diabetes=1) = 1 / (1 + exp(-(-5.91244906318139 + 0.0356440425816302*x1 + 0.02477835
## Misclassification Error: 0.2630208
```

Scatter Plot with Predicted Diabetes Values

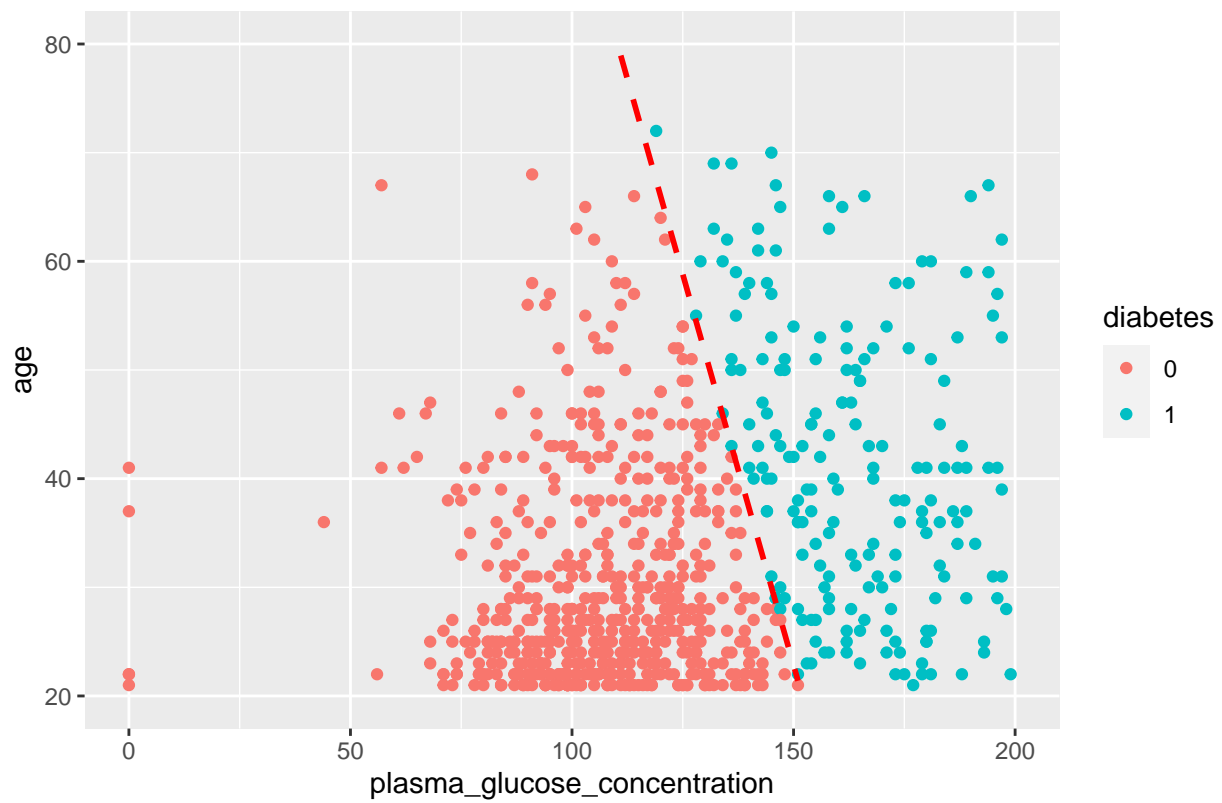


3.

Answer:

```
## Warning: Removed 1 rows containing missing values (`geom_point()`).
## Warning: Removed 454 rows containing missing values (`geom_line()`).
```

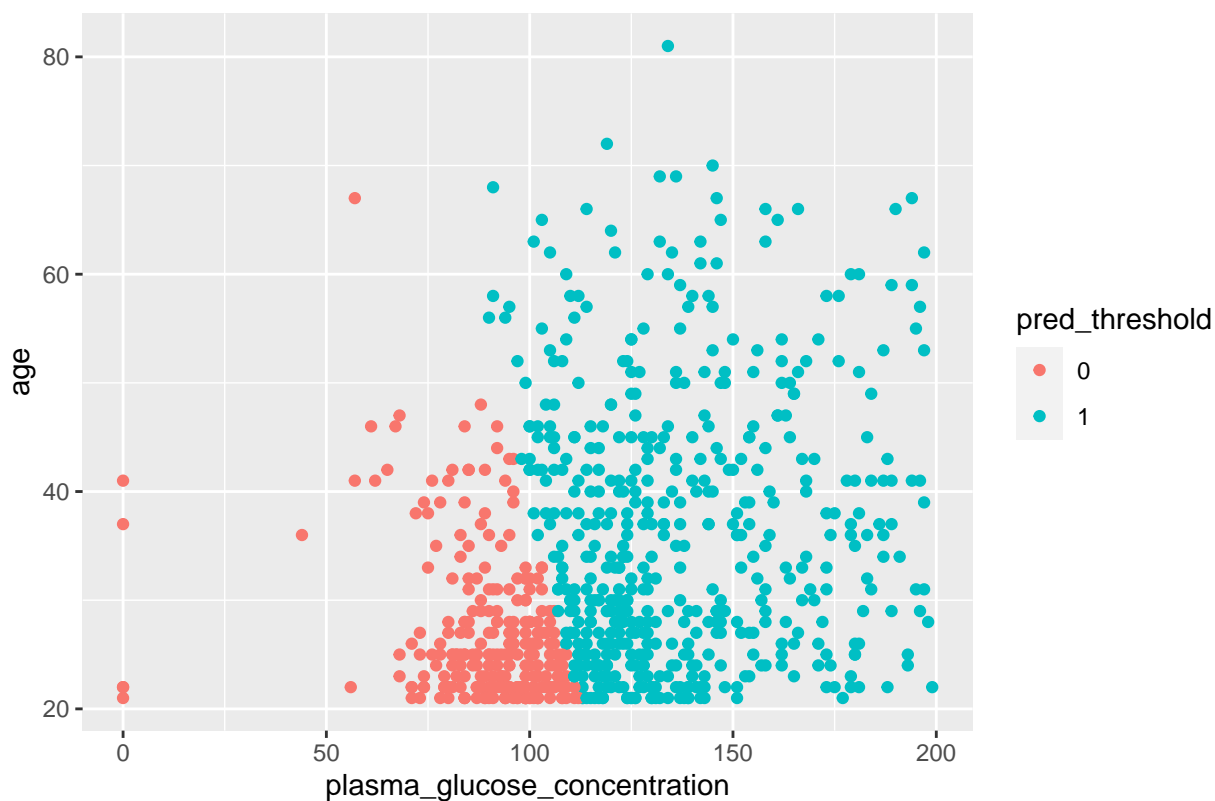
Scatter Plot of Plasma Glucose Concentration on Age with Decision Bound



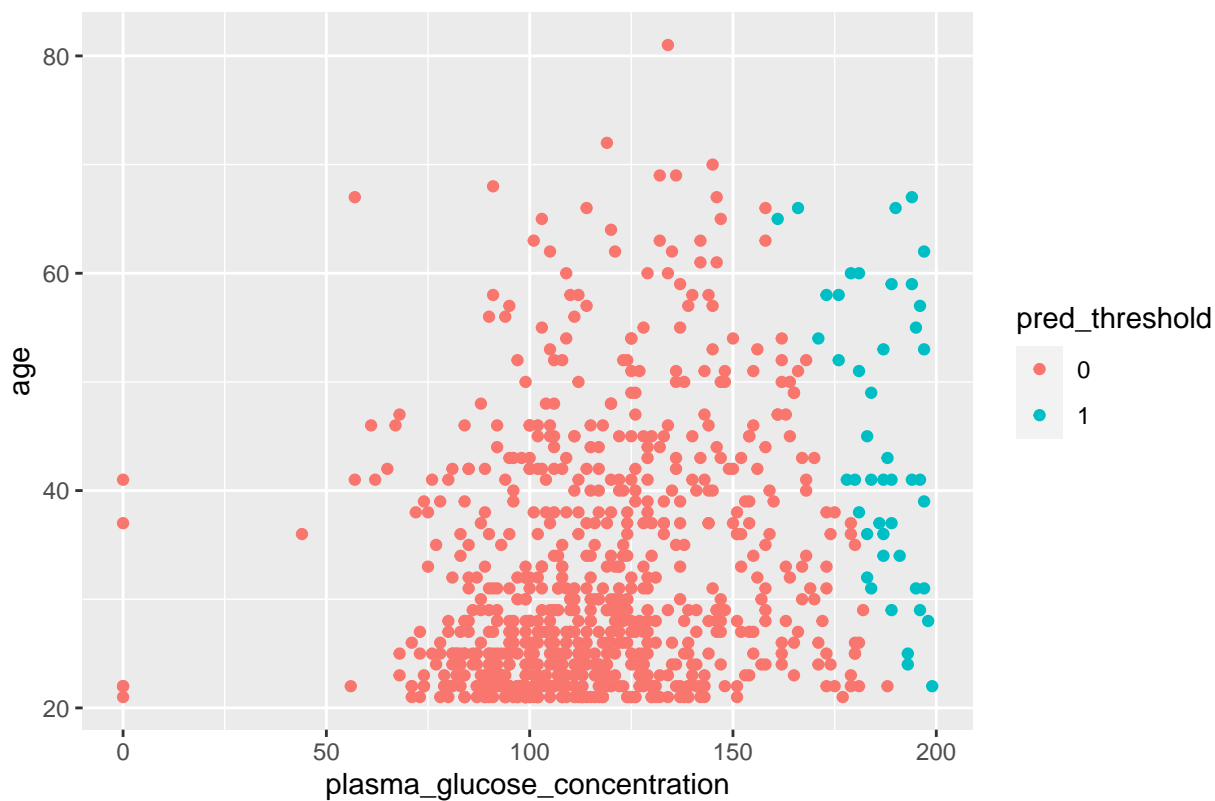
4.

Answer:

Scatter Plot with Predicted Diabetes Values (Threshold = 0.2)



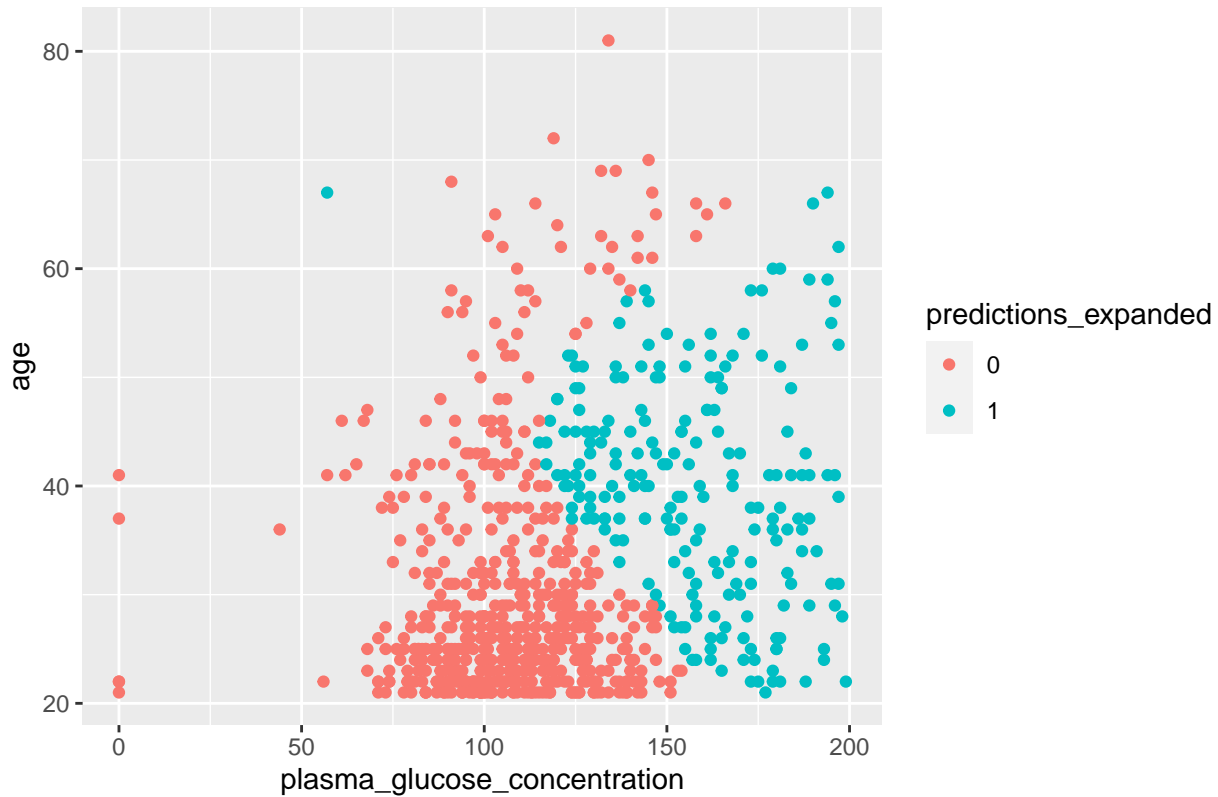
Scatter Plot with Predicted Diabetes Values (Threshold = 0.8)



5.

Answer:

Scatter Plot with Predicted Diabetes Values (Expanded Features)



Misclassification Error (Expanded Features): 0.2447917

Appendix:

knearest.R

```
# load necessary libraries
library(ggplot2)
library(kknn)

# import data set
optdigits_data <- read.csv('optdigits.csv', header = FALSE)
colnames(optdigits_data) <- c(paste0("a",1:64),"digit")
optdigits_data$digit <- as.factor(optdigits_data$digit)
#head(optdigits_data, 5)

n=dim(optdigits_data)[1]
set.seed(12345)
id=sample(1:n, floor(n*0.5))
train_data=optdigits_data[id,]
id1=setdiff(1:n, id)
id2=sample(id1, floor(n*0.25))
valid_data=optdigits_data[id2,]
id3=setdiff(id1,id2)
test_data=optdigits_data[id3,]
```

```

# 30-nearest neighbor classification
k_fit_train <- kknns(formula = digit ~ ., train_data, train_data, k = 30, kernel = "rectangular")
k_fit_test <- kknns(formula = digit ~ ., train_data, test_data, k = 30, kernel = "rectangular")

# Confusion matrices for the training and test data
train_confusion <- table(train_data$digit, fitted(k_fit_train))
cat("Confusion Matrix for Training Data:\n")
print(train_confusion)

test_confusion <- table(test_data$digit, fitted(k_fit_test))
cat("Confusion Matrix for Test Data:\n")
print(test_confusion)

# Misclassification errors for the training and test data
train_error <- 1 - sum(diag(train_confusion)) / sum(train_confusion)
cat("Misclassification errors for the training data are: ", train_error, "\n" )

test_error <- 1 - sum(diag(test_confusion)) / sum(test_confusion)
cat("Misclassification errors for the test data are: ", test_error, "\n" )

# the quality of predictions for different digits
for ( i in 1:nrow(test_confusion)) {
  digit_accuracy <- test_confusion[i,i] / sum(test_confusion[i,])
  cat("The accuracy of prediction of for digit ", i-1, " is: ", digit_accuracy, "\n")
}

overall_accuracy <- sum(diag(test_confusion)) / sum(test_confusion)
cat("The overall accuracy of prediction is:", overall_accuracy, "\n")

# Get probabilities of class "8"
probabilities <- k_fit_train$prob[, "8"]

# Get indices of training data for class "8"
indices_8 <- which(train_data$digit == "8")

# Get probabilities for class "8"
probabilities_8 <- probabilities[indices_8]

# Find 2 easiest (highest probability) and 3 hardest (lowest probability) to classify cases
easiest_indices <- indices_8[order(probabilities_8, decreasing = TRUE)[1:2]]
hardest_indices <- indices_8[order(probabilities_8)[1:3]]

# Reshape features as 8x8 matrix and visualize
for (index in c(easiest_indices, hardest_indices)) {
  digit_8 <- matrix(as.numeric(train_data[index, 1:64]), nrow = 8, byrow = TRUE)
  heatmap(digit_8, Colv = NA, Rowv = NA, main = paste("Digit '8', Index:", index))
}

# Fit KNN for different K values and plot errors
errors <- data.frame()
for (k in 1:30) {
  fit <- kknns(digit ~ ., train_data, valid_data, k = k, kernel = "rectangular")

```

```

pred <- fit$fitted.values
confusion_matrix <- table(valid_data$digit, pred)
error <- 1 - sum(diag(confusion_matrix)) / sum(confusion_matrix)
errors <- rbind(errors, data.frame(K = k, Error = error))
}
# Plot the misclassification errors on the value of K
plot(errors$K, errors$Error, type = "b", col='orange', main="Misclassification errors on the value of K")

# Initialize a data frame to store the results
ce_errors <- data.frame(k_value = integer(), cross_entropy_error = numeric())

for (k in 1:30) {
  # Fit K-nearest neighbor classifier
  fit <- kknnc(digit ~ ., train_data, valid_data, k = k, kernel = "rectangular")

  # Get predicted probabilities
  predicted_probabilities <- fit$prob

  # Compute cross-entropy error
  #actual_probabilities <- ifelse(valid_data$digit == "8", 1, 0)
  #cross_entropy <- -sum(actual_probabilities * log(predicted_probabilities + 1e-15))
  cross_entropy <- -sum(as.numeric(valid_data$digit) * log(predicted_probabilities + 1e-15))

  ce_errors <- rbind(ce_errors, data.frame(k_value = k, cross_entropy_error = cross_entropy))
}

# Plot the results
plot(ce_errors$k_value, ce_errors$cross_entropy_error, type = "b", col='orange', main="Cross Entropy error vs K")

# Find the optimal K
optimal_k <- ce_errors$k_value[which.min(ce_errors$cross_entropy_error)]
cat("The optimal k is:", optimal_k, "\n")

```

linear_ridge.R

```

# import data set
parkinson_data <- read.csv('parkinsons.csv', header = TRUE)
#head(parkinson_data, 5)

# Split the data into training and test sets
parkinson_data <- subset(parkinson_data, select = -c(1:4,6))
n=dim(parkinson_data)[1]
set.seed(12345)
id=sample(1:n, floor(n*0.6))
train_data=parkinson_data[id,]
test_data=parkinson_data[-id,]

# Scale the data
#scaled_train_data <- as.data.frame(cbind(scale(subset(train_data, select = -c(motor_UPDRS))), motor_UPDRS))
#scaled_test_data <- as.data.frame(cbind(scale(subset(test_data, select = -c(motor_UPDRS))), motor_UPDRS))
scaled_train_data <- as.data.frame(scale(train_data))
scaled_test_data <- as.data.frame(scale(test_data))

# linear regression model

```

```

lm_model <- lm(motor_UPDRS ~ ., data = scaled_train_data)

# Training and test MSE
train_pred <- predict(lm_model, newdata = scaled_train_data)
train_mse <- mean((scaled_train_data$motor_UPDRS - train_pred)^2)
cat("The training data's MSE is: ", train_mse)

test_pred <- predict(lm_model, newdata = scaled_test_data)
test_mse <- mean((scaled_test_data$motor_UPDRS - test_pred)^2)
cat("The training data's MSE is: ", test_mse)

# Significant variables
summary(lm_model)

# Log-likelihood function
log_likelihood_fun <- function(theta, sigma, y, X) {
  n <- length(y)
  log_likelihood <- -n/2 * log(2*pi*sigma^2) - 1/(2*sigma^2) * sum((y - X%%theta)^2)
  return(log_likelihood)
}

# Ridge log-likelihood function
ridge_log_Likelihood_fun <- function(theta, sigma, lambda, y, X) {
  log_like <- log_likelihood_fun(theta, sigma, y, X)
  ridge_penalty <- lambda * sum(theta^2)
  return(-log_like + ridge_penalty)
}

# Ridge log-likelihood optimization function
ridge_log_likelihood_opt <- function(lambda, y, X) {
  start <- c(rep(0, ncol(X)), sd(y)) # combine theta and sigma into one numeric vector
  theta_length <- ncol(X)

  # Define a function for the optim() call
  fn_to_optim <- function(params) {
    theta <- params[1:theta_length]
    sigma <- params[theta_length + 1]
    return(-ridge_log_Likelihood_fun(theta, sigma, lambda, y, X))
  }

  # Call optim()
  opt_res <- optim(start, fn = fn_to_optim, method = "BFGS")

  # Return theta and sigma separately
  return(list(theta = opt_res$par[1:theta_length], sigma = opt_res$par[theta_length + 1]))
}

# Degrees of freedom function
df_fun <- function(lambda, X) {
  H <- solve(t(X) %*% X + lambda * diag(ncol(X))) %*% t(X) %*% X
  return(sum(diag(H)))
}

```

```

y <- scaled_train_data$motor_UPDRS
X <- as.matrix(subset(scaled_train_data, select = -c(motor_UPDRS)))
# Ridge optimization for different lambdas
lambdas <- c(1, 100, 1000)
opt_params <- lapply(lambdas, function(lambda) ridge_log_likelihood_opt(lambda, y, X))

# Predictions and MSE for different lambdas
ridge_train_preds <- lapply(1:length(opt_params), function(i) {
  theta <- opt_params[[i]]$theta
  as.matrix(subset(scaled_train_data, select = -c(motor_UPDRS))) %*% theta
})

#train_mses <- colMeans((scaled_train_data$motor_UPDRS - ridge_train_preds)^2)
ridge_train_mse <- sapply(ridge_train_preds, function(pred) mean((pred - scaled_train_data$motor_UPDRS)^2))
cat("Ridge Training MSE:", ridge_train_mse, "\n")

ridge_test_preds <- lapply(1:length(opt_params), function(i) {
  theta <- opt_params[[i]]$theta
  as.matrix(subset(scaled_test_data, select = -c(motor_UPDRS))) %*% theta
})

ridge_test_mse <- sapply(ridge_test_preds, function(pred) mean((pred - scaled_test_data$motor_UPDRS)^2))
cat("Ridge Test MSE:", ridge_test_mse, "\n")

# Degrees of freedom for different lambdas
degrees_of_freedom <- sapply(lambdas, function(lambda) df_fun(lambda, X))
cat("Degrees of Freedom:", degrees_of_freedom, "\n")

```

logistic.R

```

library(ggplot2)
library(caret)

# import data set
pima_data <- read.csv('pima-indians-diabetes.csv', header = FALSE)
colnames(pima_data) <- c("num_of_pregnant", "plasma_glucose_concentration", "blood_pressure",
                        "skinfold_thickness", "serum_insulin", "bmi",
                        "diabetes_predigree", "age", "diabetes")

#head(pima_data, 5)

# Scatterplot
ggplot(pima_data, aes(x=plasma_glucose_concentration, y=age, color=as.factor(diabetes))) +
  geom_point() +
  labs(color="diabetes", title='Scatter Plot of Plasma Glucose Concentration on Age')

# Logistic Regression
logistic_model <- glm(diabetes ~ plasma_glucose_concentration + age, data=pima_data, family=binomial)
summary(logistic_model)

# Predict probabilities
probabilities <- predict(logistic_model, type="response")
print(paste0("Probability(Diabetes=1) = 1 / (1 + exp(-(",coef(logistic_model)[1]," + ", coef(logistic_m

```

```

# Classify observations
predictions <- ifelse(probabilities >= 0.5, 1, 0)

# Compute misclassification error
mis_error <- mean(predictions != pima_data$diabetes)
cat("Misclassification Error:", mis_error, "\n")

# Scatter plot with predicted values
ggplot(data.frame(pima_data, predictions), aes(x=plasma_glucose_concentration, y=age, color=as.factor(p
  geom_point() +
  labs(color="predictions",title='Scatter Plot with Predicted Diabetes Values')

# Decision boundary equation
# Decision boundary is where the logistic function equals 0.5
#  $0 = \text{intercept} + \text{coef1} \cdot x_1 + \text{coef2} \cdot x_2$ 
#  $x_2 = -(\text{intercept} + \text{coef1} \cdot x_1) / \text{coef2}$ 

#decision_boundary <- -(coef(logsitic_model)[1] + coef(logsitic_model)[2]*pima_data$age) / coef(logsiti
decision_boundary_x1 <- pima_data$plasma_glucose_concentration
decision_boundary_x2 <- -(coef(logsitic_model)[1] + coef(logsitic_model)[2]*decision_boundary_x1) / coe

# Add decision boundary to scatter plot
ggplot(data.frame(pima_data, predictions), aes(x=plasma_glucose_concentration, y=age, color=as.factor(p
  geom_point() +
  geom_line(aes(x=decision_boundary_x1, y=decision_boundary_x2), color='red', linetype='dashed', linewid
  labs(color="diabetes", title='Scatter Plot of Plasma Glucose Concentration on Age with Decision Bound

# Predictions with different thresholds
thresholds <- c(0.2, 0.8)

for ( threshold in thresholds) {
  pred_threshold <- ifelse(predict(logsitic_model, newdata=pima_data, type='response') >= threshold, 1,

  # Scatter plot with predicted values and threshold
  scatter_plot <- ggplot(data.frame(pima_data, pred_threshold), aes(x=plasma_glucose_concentration, y=a
    geom_point() +
    labs(color="pred_threshold", title=paste('Scatter Plot with Predicted Diabetes Values (Threshold = '
    print(scatter_plot)
}

# Create new features
pima_data$z1 <- pima_data$plasma_glucose_concentration^4
pima_data$z2 <- pima_data$plasma_glucose_concentration^3 * pima_data$age
pima_data$z3 <- pima_data$plasma_glucose_concentration^2 * pima_data$age^2
pima_data$z4 <- pima_data$plasma_glucose_concentration * pima_data$age^3
pima_data$z5 <- pima_data$age^4

# Logistic Regression with new features
model_expanded <- glm(diabetes ~ plasma_glucose_concentration + age + z1 + z2 + z3 + z4 + z5, data=pima,

```

```

# Predict probabilities and classify observations
probabilities_expanded <- predict(model_expanded, type="response")
predictions_expanded <- ifelse(probabilities_expanded >= 0.5, 1, 0)

# Scatter plot for the model with expanded features
ggplot(data.frame(pima_data, predictions_expanded), aes(x=plasma_glucose_concentration, y=age, color=predictions_expanded)) +
  geom_point() +
  labs(color="predictions_expanded", title='Scatter Plot with Predicted Diabetes Values (Expanded Features)')

# Compute misclassification error
mis_error_expanded <- mean(predictions_expanded != pima_data$diabetes)
cat("Misclassification Error (Expanded Features):", mis_error_expanded, "\n")

```