

ASSIGNMENT 1

DSCI 6607 – Programmatic Data Analysis Using Python and R

Name: Sahil Khan

Student ID: 202482066

Submission Date: 15-Oct-2024

Question 1

Simulate tossing an unfair coin (probability of heads = 0.63) for 301 trials and compute the proportion of heads.

$$\hat{p} = \frac{\text{Number of Heads}}{\text{Total number of trials}}$$

```
import random

probHeads = 0.63
tossTrials = 301

headsCount = 0

for i in range(tossTrials):
    if random.random() < probHeads:
        headsCount += 1

estProbHeads = headsCount / tossTrials

print(f"Number of heads: {headsCount}")

## Number of heads: 187

print(f"Estimated probability of Heads: {estProbHeads}")

## Estimated probability of Heads: 0.6212624584717608
```

Question 2

- a) Implementation of the Bisection Method

```

import math

def bisection_root(a, b, f):
    """Finds the root of a function using the bisection method.

    Args:
        a: Left endpoint of the interval.
        b: Right endpoint of the interval.
        f: The function for which to find the root.

    Returns:
        The root of the function.
    """

    while abs(f((a + b) / 2)) > 1e-5: # Changed the stop condition to be based on f(m)
        m = (a + b) / 2
        if f(m) > 0:
            a = m
        else:
            b = m
    return (a + b) / 2

```

The bisection method requires an interval $[a, b]$ where $f(a)$ and $f(b)$ have opposite signs (one positive, one negative). This guarantees that the function crosses zero at least once within the interval.

For given interval $[-2.3, -0.23]$:

```

def f(x):
    return x**3 + 4 * x**2 - 3

a = -2.3
b = -0.23
print(f"The sign of f({a}) = x^3 + 4x^2 - 3 for a={a} is: {f(a)}")

## The sign of f(-2.3) = x^3 + 4x^2 - 3 for a=-2.3 is: 5.9929999999999986

print(f"The sign of f({b}) = x^3 + 4x^2 - 3 for a={b} is: {f(b)}")

```

```

## The sign of f(-0.23) = x^3 + 4x^2 - 3 for a=-0.23 is: -2.800567

```

Since $f(-2.3)$ and $f(-0.23)$ have opposite signs, the bisection method is guaranteed to find a root in the interval $[-2.3, -0.23]$.

b) Finding the root of $f(x) = x^3 + 4x^2 - 3$ in interval $[-2.3, -0.23]$.

```

root = bisection_root(a, b, f)
print(f"The root of f(x) = x^3 + 4x^2 - 3 in the interval [-2.3,-0.23] is: {root}")

## The root of f(x) = x^3 + 4x^2 - 3 in the interval [-2.3,-0.23] is: -0.9999999809265135

```

Question 3

```
# library(datasets)
head(rivers)
```

```
## [1] 735 320 325 392 524 450
```

a) R Function to Compute Cumulative Sample Means

```
cumulative_means <- function(x) {
  n <- length(x)
  cumsum(x) / (1:n)
}
```

b) $n = 1000$ observations with replacement from 'rivers' dataset.

```
n <- 1000
rivers_sample <- sample(rivers, n, replace = TRUE)
head(rivers_sample)
```

```
## [1] 360 420 1000 710 620 250
```

c) Applying function (a) to the data of part (b).

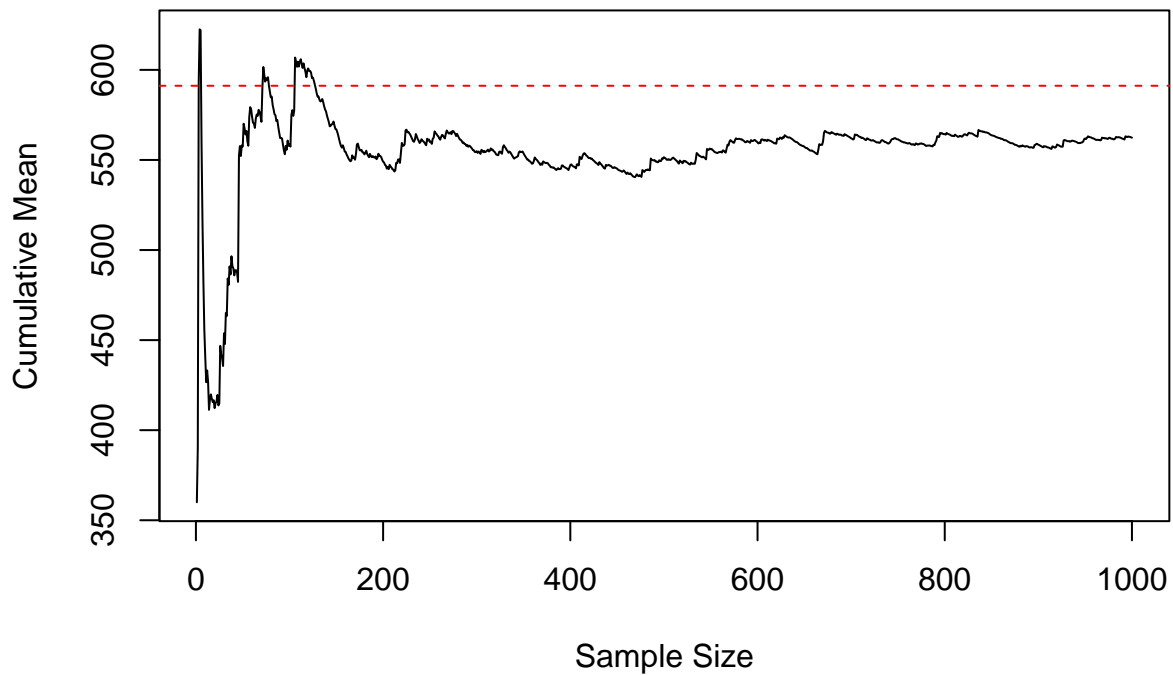
```
cumulative_sample_means <- cumulative_means(rivers_sample)
head(cumulative_sample_means)
```

```
## [1] 360.0000 390.0000 593.3333 622.5000 622.0000 560.0000
```

d) Plotting the population mean.

```
plot(cumulative_sample_means, type = "l", xlab = "Sample Size", ylab = "Cumulative Mean",
     main = "Cumulative Sample Means of River Lengths with 1000 Obs")
abline(h=mean(rivers), col="red", lty=2)
```

Cumulative Sample Means of River Lengths with 1000 Obs



Explanation:

The plot shows how the cumulative sample mean converges towards the population mean (red dashed line) as the sample size increases.

The Central Limit Theorem suggests that as 'n' grows larger, the distribution of sample means approaches a normal distribution centered around the population mean. We can observe this convergence visually in the plot.

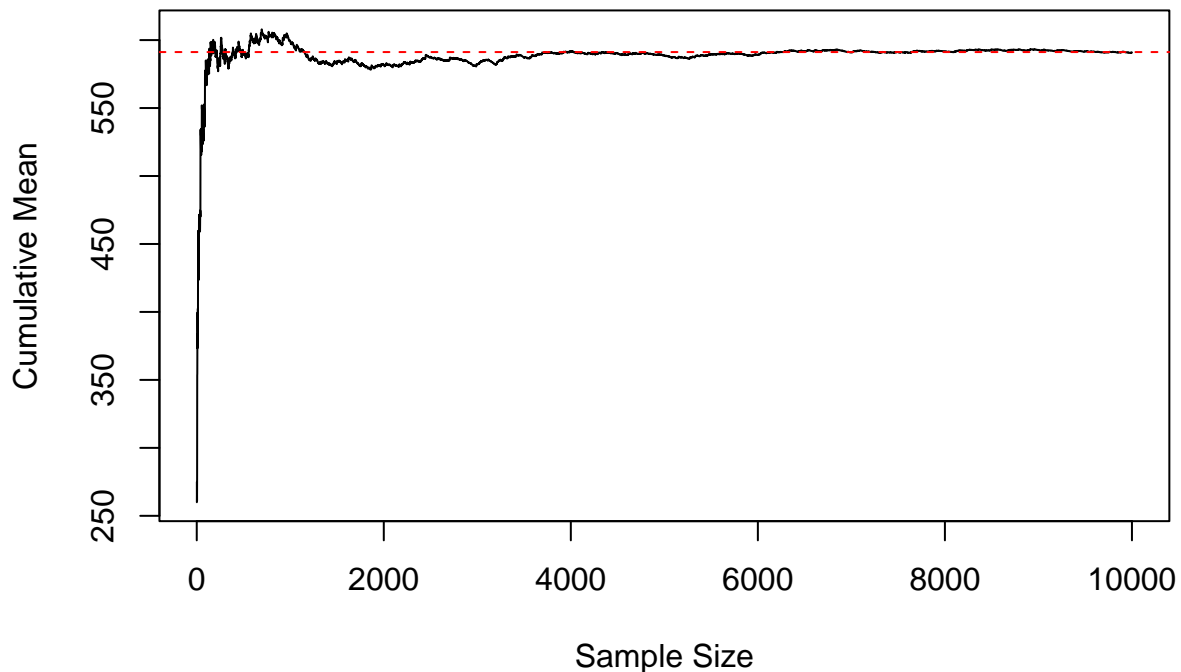
The more samples we have, the closer the cumulative sample mean gets to the true population mean.

Plotting with 10000 Observations:

```
n <- 10000
rivers_sample <- sample(rivers, n, replace = TRUE)
cumulative_sample_means <- cumulative_means(rivers_sample)

plot(cumulative_sample_means, type = "l", xlab = "Sample Size", ylab = "Cumulative Mean",
     main = "Cumulative Sample Means of River Lengths with 10000 Obs")
abline(h = mean(rivers), col = "red", lty = 2)
```

Cumulative Sample Means of River Lengths with 10000 Obs



The plot clearly shows that as the sample size increases, the cumulative sample mean converges towards the true population mean (represented by the red dashed line). The fluctuations around the population mean decrease with increasing sample size, illustrating the principle of the Central Limit Theorem: with larger samples, the sample mean becomes a more precise estimate of the population mean.

Question 4

The `airquality` data set reports the daily air quality measurements in New York, May to September 1973. The data set includes 153 observations and 6 variables.

```
# library(datasets)
head(airquality)
```

```
##      Ozone Solar.R Wind Temp Month Day
## 1      41     190   7.4   67     5    1
## 2      36     118   8.0   72     5    2
## 3      12     149  12.6   74     5    3
## 4      18     313  11.5   62     5    4
## 5      NA      NA  14.3   56     5    5
## 6      28      NA  14.9   66     5    6
```

a) Ridge Regression Model Coefficients:

```
ridge_regression <- function(X, y, lambda) {
  n <- nrow(X)
  p <- ncol(X)
  I <- diag(p)
  beta_ridge <- solve(t(X) %*% X + lambda * I) %*% t(X) %*% y
  return(beta_ridge)
}
```

b) Normalize Variables:

```
ozone <- airquality$Ozone
solar_r <- airquality$Solar.R
wind <- airquality$Wind
temp <- airquality$Temp

data <- data.frame(Ozone = ozone, Solar.R = solar_r, Wind = wind, Temp = temp)
data <- na.omit(data) # Remove rows with missing values

normalize <- function(x) {
  (x - mean(x)) / sd(x)
}

data_norm <- as.data.frame(lapply(data, normalize))

X_norm <- as.matrix(data_norm[, c("Solar.R", "Wind", "Temp")])
y_norm <- data_norm$Ozone

print("Normalized Explanatory Variables (X_norm):")
```

```
## [1] "Normalized Explanatory Variables (X_norm):"
```

```
print(head(X_norm, 5))
```

```
##           Solar.R           Wind           Temp
## [1,]  0.05702761 -0.7138405 -1.1325108
## [2,] -0.73285918 -0.5451928 -0.6078501
## [3,] -0.39276904  0.7477726 -0.3979858
## [4,]  1.40641756  0.4385852 -1.6571715
## [5,]  1.25282846 -0.3765451 -1.3423751
```

```
print("Normalized Response Variable (y_norm):")
```

```
## [1] "Normalized Response Variable (y_norm):"
```

```
print(head(y_norm, 10))
```

```
## [1] -0.03302982 -0.18328840 -0.90452961 -0.72421931 -0.57396073 -0.69416759
## [7] -1.02473648 -0.78432275 -0.93458133 -0.84442618
```

c) computes $\hat{\beta}_R$ of model for $\lambda = 0.1$:

```
lambda <- 0.1
beta_r <- ridge_regression(X_norm, y_norm, lambda)

print(beta_r)
```

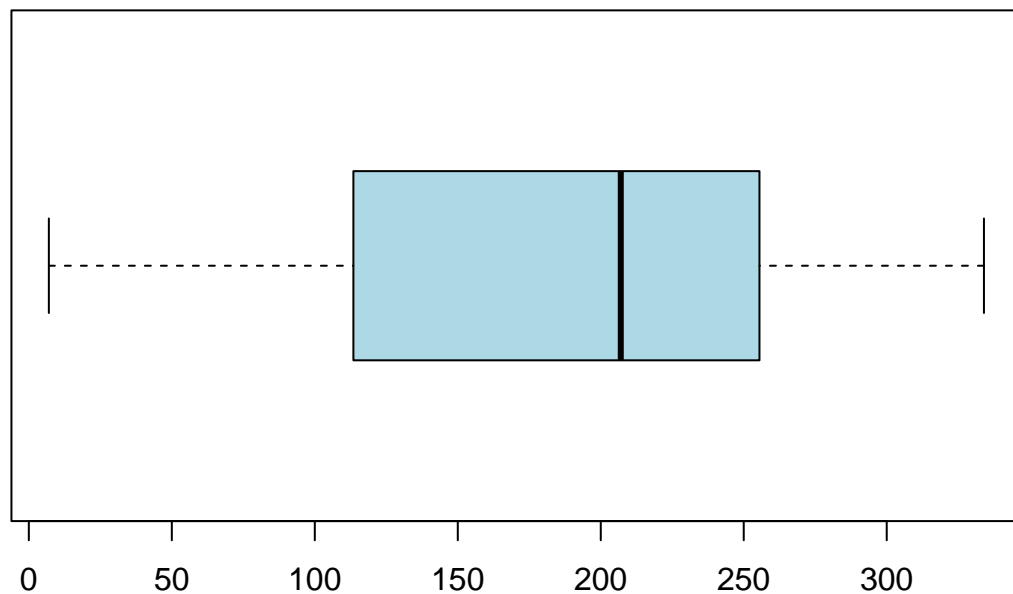
```
##           [,1]
## Solar.R  0.1638378
## Wind    -0.3562652
## Temp     0.4727975
```

d) Box Plot of 'OZONE', 'SOLAR', 'WIND' AND 'TEMP' :

```
# Load necessary library
library(ggplot2)

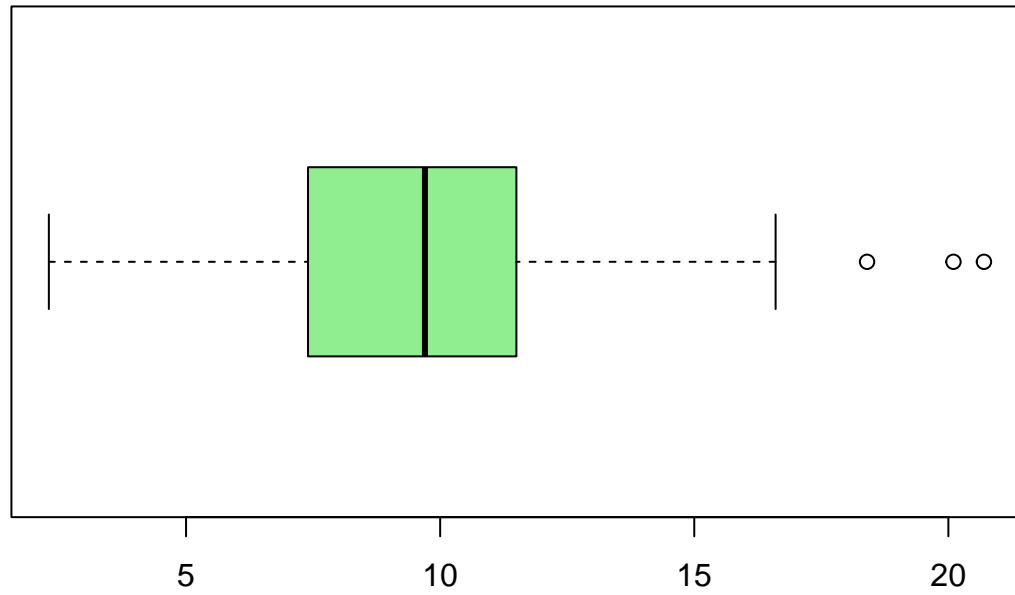
# Display boxplots of the variables
boxplot(data$Solar.R, main="Boxplot of Solar.R", col="lightblue", horizontal=TRUE)
```

Boxplot of Solar.R



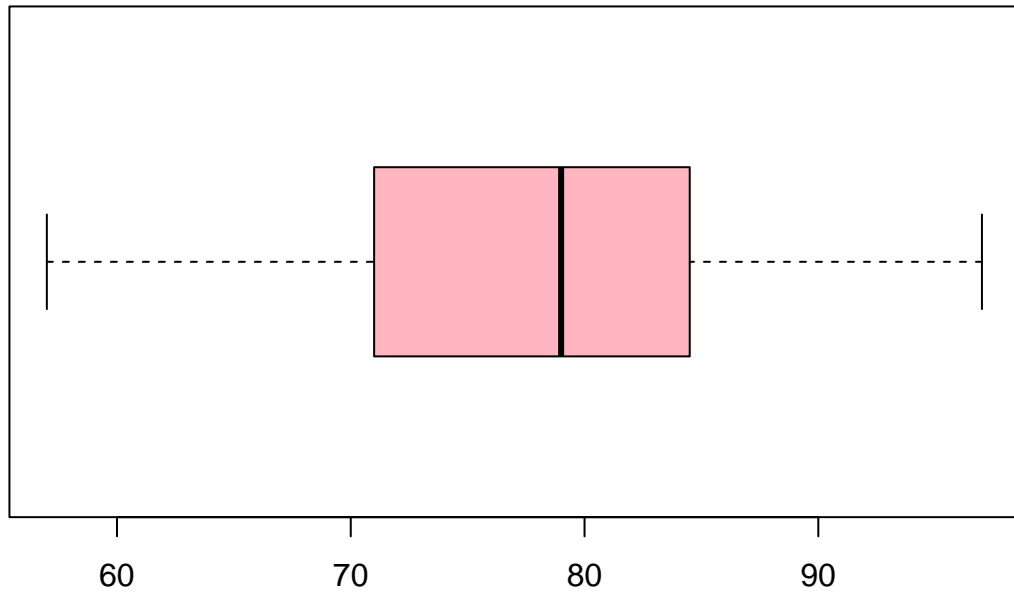
```
boxplot(data$Wind, main="Boxplot of Wind", col="lightgreen", horizontal=TRUE)
```

Boxplot of Wind



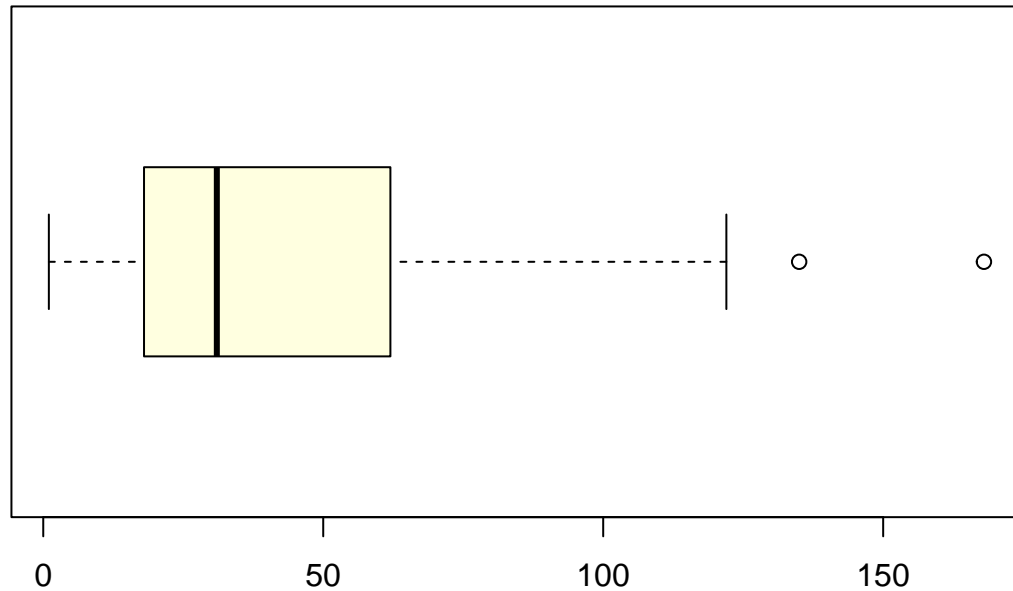
```
boxplot(data$Temp, main="Boxplot of Temp", col="lightpink", horizontal=TRUE)
```


Boxplot of Temp



```
boxplot(data$Ozone, main="Boxplot of Ozone", col="lightyellow", horizontal=TRUE)
```

Boxplot of Ozone



Analysis of BoxPlot:

- a) Ozone: The Ozone boxplot shows some outliers to the right (higher values). This indicates a right-skewed distribution. When a distribution is skewed, the mean is pulled in the direction of the skew (in this case, towards the higher values), making it less representative of the typical value. Therefore, the median is a better measure of the center for Ozone.
- b) Solar.R: The Solar.R boxplot appears relatively symmetrical, with the median roughly in the center of the box and the whiskers fairly even. There's a possible minor outlier or two, but they don't seem to drastically skew the distribution. In such cases, the mean is a reasonable measure of the center.
- c) Wind: Similar to Ozone, the Wind boxplot indicates a right-skewed distribution due to the presence of a few outliers on the right. The skew, while not as pronounced as Ozone, is enough to suggest that the median is a more appropriate measure of central tendency.
- d) Temp: The Temp boxplot looks quite symmetrical, similar to Solar.R. The median is in the middle of the box, and the whiskers are relatively even. Therefore, the mean is a suitable measure of central tendency for Temp.

Question 5

- a) Divide the AirQuality Dataset [n=100] into training and testing:

```

set.seed(42)
sample_data <- data_norm[sample(nrow(data_norm), 100), ]
data <- na.omit(sample_data)

# Split the data into training (70%) and testing (30%)
n <- nrow(data)
train_indices <- sample(1:n, 70) #70% for training
test_indices <- setdiff(1:n, train_indices)
print(n)

```

```
## [1] 100
```

b) Function to Calculate Coefficients of Regression Model and RMSE:

```

ridge_regression <- function(X, y, lambda) {
  n <- nrow(X)
  p <- ncol(X)
  I <- diag(p)
  beta_ridge <- solve(t(X) %*% X + lambda * I) %*% t(X) %*% y
  return(beta_ridge)
}

calculate_rmse <- function(X_train, y_train, X_test, y_test, lambda){
  beta_r <- ridge_regression(X_train, y_train, lambda)
  y_pred <- X_test %*% beta_r
  rmse <- sqrt(mean((y_test - y_pred)^2))
  return(rmse)
}

```

c) Find optimal λ using RMSE:

```

lambda_values <- seq(-2, 2, length.out = 100)
rmse_values <- numeric(length(lambda_values))

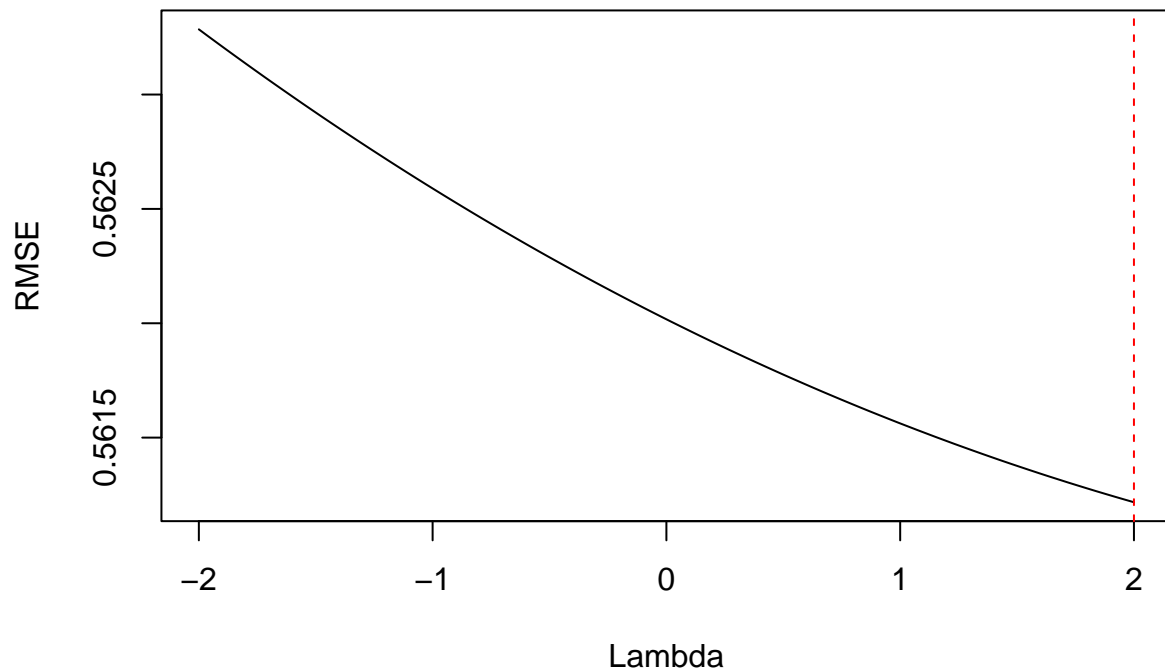
for (i in 1:length(lambda_values)) {
  X_train <- X_norm[train_indices, ]
  y_train <- y_norm[train_indices]
  X_test <- X_norm[test_indices, ]
  y_test <- y_norm[test_indices]

  rmse_values[i] <- calculate_rmse(X_train, y_train, X_test, y_test, lambda_values[i])
}

optimal_lambda <- lambda_values[which.min(rmse_values)]

plot(lambda_values, rmse_values, type = "l", xlab = "Lambda", ylab = "RMSE")
abline(v = optimal_lambda, col = "red", lty = 2) # Show the optimal lambda

```



```
print(paste("Optimal Lambda:", optimal_lambda))
```

```
## [1] "Optimal Lambda: 2"
```

The plot of RMSE vs. lambda shows the impact of the regularization parameter on model performance. The optimal lambda value '2' (indicated by the vertical red dashed line) is where the RMSE is minimized, indicating that, for this dataset split, this degree of regularization achieves the optimal balance between overfitting and underfitting on unseen data.

Question 6

a) Function to simulate from a standard normal distribution:

```
import math
import random
import matplotlib.pyplot as plt
import numpy as np

def box_muller(u1, u2):
    """Generates standard normal random variables using the Box-Muller transform.
```

```

Args:
    u1: A random number between 0 and 1.
    u2: A random number between 0 and 1.

Returns:
    A tuple containing two standard normal random variables.
    """
    x = np.sqrt(-2 * np.log(u1)) * np.cos(2 * np.pi * u2)
    y = np.sqrt(-2 * np.log(u1)) * np.sin(2 * np.pi * u2)

    return x, y

```

b) Generate 1000 observations:

```

n_samples = 1000
box_muller_samples = []
for _ in range(n_samples):
    u1 = random.random()
    u2 = random.random()
    x, y = box_muller(u1, u2)
    box_muller_samples.extend([x,y])

```

c) Plot histogram and compare:

```
plt.figure(figsize=(12, 6))
```

```
## <Figure size 1200x600 with 0 Axes>
```

```
plt.subplot(1, 2, 1)
```

```
## <AxesSubplot:>
```

```
plt.hist(box_muller_samples, bins=50, color='b', alpha=0.7, label="Box-Muller Samples")
```

```

## (array([ 2.,  1.,  3.,  4.,  3.,  1.,  8., 10., 13., 12., 14.,
##         26., 24., 36., 42., 54., 67., 59., 80., 68., 103., 98.,
##         101., 89., 88., 95., 107., 102., 108., 82., 79., 77., 62.,
##         54., 44., 42., 34., 22., 17., 30., 13., 11.,  5.,  1.,
##         3.,  1.,  1.,  2.,  1.,  1.]), array([-3.29092569, -3.1597504 , -3.02857511, -2.897399
##        -2.63504925, -2.50387397, -2.37269868, -2.24152339, -2.11034811,
##        -1.97917282, -1.84799753, -1.71682224, -1.58564696, -1.45447167,
##        -1.32329638, -1.1921211 , -1.06094581, -0.92977052, -0.79859523,
##        -0.66741995, -0.53624466, -0.40506937, -0.27389409, -0.1427188 ,
##        -0.01154351,  0.11963178,  0.25080706,  0.38198235,  0.51315764,
##        0.64433292,  0.77550821,  0.9066835 ,  1.03785878,  1.16903407,
##        1.30020936,  1.43138465,  1.56255993,  1.69373522,  1.82491051,
##        1.95608579,  2.08726108,  2.21843637,  2.34961166,  2.48078694,
##        2.61196223,  2.74313752,  2.8743128 ,  3.00548809,  3.13666338,
##        3.26783867])), <BarContainer object of 50 artists>)

```

```
plt.title('Histogram of Box-Muller Observations')
```

```
## Text(0.5, 1.0, 'Histogram of Box-Muller Observations')
```

```
plt.xlabel('Value')
```

```
## Text(0.5, 0, 'Value')
```

```
plt.ylabel('Density')
```

```
## Text(0, 0.5, 'Density')
```

```
plt.subplot(1, 2, 2)
```

```
## <AxesSubplot:>
```

```
plt.hist(box_muller_samples, bins=50, alpha=0.5, label="Box-Muller Samples")
```

```
## (array([ 2.,  1.,  3.,  4.,  3.,  1.,  8., 10., 13., 12., 14.,
##         26., 24., 36., 42., 54., 67., 59., 80., 68., 103., 98.,
##        101., 89., 88., 95., 107., 102., 108., 82., 79., 77., 62.,
##        54., 44., 42., 34., 22., 17., 30., 13., 11.,  5.,  1.,
##        3.,  1.,  1.,  2.,  1.,  1.]), array([-3.29092569, -3.1597504 , -3.02857511, -2.897399
##       -2.63504925, -2.50387397, -2.37269868, -2.24152339, -2.11034811,
##       -1.97917282, -1.84799753, -1.71682224, -1.58564696, -1.45447167,
##       -1.32329638, -1.1921211 , -1.06094581, -0.92977052, -0.79859523,
##       -0.66741995, -0.53624466, -0.40506937, -0.27389409, -0.1427188 ,
##       -0.01154351,  0.11963178,  0.25080706,  0.38198235,  0.51315764,
##        0.64433292,  0.77550821,  0.9066835 ,  1.03785878,  1.16903407,
##        1.30020936,  1.43138465,  1.56255993,  1.69373522,  1.82491051,
##        1.95608579,  2.08726108,  2.21843637,  2.34961166,  2.48078694,
##        2.61196223,  2.74313752,  2.8743128 ,  3.00548809,  3.13666338,
##        3.26783867])), <BarContainer object of 50 artists>)
```

```
plt.hist(np.random.randn(1000), bins=50, alpha=0.5, label='True Normal')
```

```
## (array([ 1.,  0.,  0.,  1.,  1.,  3.,  6.,  6.,  9.,  6., 11., 15., 16.,
##        18., 19., 16., 18., 35., 24., 33., 27., 51., 38., 43., 53., 50.,
##        44., 41., 47., 55., 41., 33., 35., 32., 22., 26., 16., 30., 14.,
##        15., 10., 13.,  3.,  7.,  3.,  6.,  0.,  4.,  2.,  1.]), array([-3.01412208e+00, -2.89821941e
##       -2.55051138e+00, -2.43460870e+00, -2.31870602e+00, -2.20280335e+00,
##       -2.08690067e+00, -1.97099799e+00, -1.85509532e+00, -1.73919264e+00,
##       -1.62328996e+00, -1.50738729e+00, -1.39148461e+00, -1.27558193e+00,
##       -1.15967925e+00, -1.04377658e+00, -9.27873901e-01, -8.11971224e-01,
##       -6.96068547e-01, -5.80165870e-01, -4.64263193e-01, -3.48360517e-01,
##       -2.32457840e-01, -1.16555163e-01, -6.52485900e-04,  1.15250191e-01,
##        2.31152868e-01,  3.47055545e-01,  4.62958222e-01,  5.78860899e-01,
##        6.94763575e-01,  8.10666252e-01,  9.26568929e-01,  1.04247161e+00,
##        1.15837428e+00,  1.27427696e+00,  1.39017964e+00,  1.50608231e+00,
##        1.62198499e+00,  1.73788767e+00,  1.85379034e+00,  1.96969302e+00,
##        2.08559570e+00,  2.20149837e+00,  2.31740105e+00,  2.43330373e+00,
##        2.54920641e+00,  2.66510908e+00,  2.78101176e+00])), <BarContainer object of 50 artists>)
```

```
plt.xlabel("Value")

## Text(0.5, 0, 'Value')

plt.ylabel("Density")

## Text(0, 0.5, 'Density')

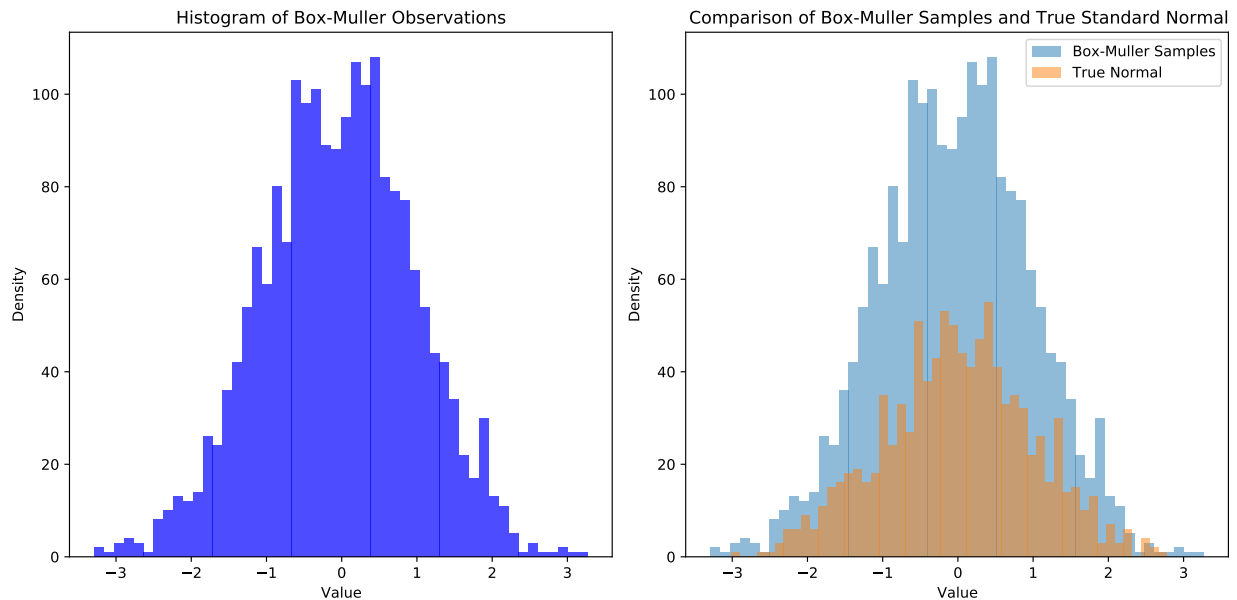
plt.title("Comparison of Box-Muller Samples and True Standard Normal")

## Text(0.5, 1.0, 'Comparison of Box-Muller Samples and True Standard Normal')

plt.legend()

## <matplotlib.legend.Legend object at 0x000002465E9C2970>

plt.tight_layout()
plt.show()
```



The histogram of the Box-Muller generated samples closely resembles the histogram of the true standard normal distribution. This visually confirms that the Box-Muller transform effectively generates random numbers that follow a standard normal distribution.

The overlapping histograms indicate that the generated samples mimic the bell-shaped curve, centered around zero, characteristic of a standard normal distribution.