# IE 5390 – Assignment - 09

Name: Simran Abhay Sinha

## Lab 17

1. Use the cars dataset that we worked with in class. But this time, instead of using the lm() function for linear regression, create the linear regression model from scratch. You have the math preliminaries to help you with the model components. Some hints on what you will need to do:

i) Establish what x and y are. In this case, x is speed and y is the stopping distance.
ii) Calculate Sxx, Sxy and Syy
iii) Then calculate Beta0 and Beta1

Loading the dataset
Cars is built-in dataset

```r
5 ```{r}
6  # Lab 17_1
7  data(cars)  # Load built-in dataset
8  x <- cars$speed  # Independent variable (speed)
9  y <- cars$dist   # Dependent variable (stopping distance)
10 ```
11
```

Answer the following questions:

a) What are the values for Sxx, Sxy and Syy that you obtained?

Code:

```r
12 - ```{r}
13   #Lab 17_A
14   x_mean <- mean(x)
15   y_mean <- mean(y)
16
17   Sxx <- sum((x - x_mean)^2)
18   Sxy <- sum((x - x_mean) * (y - y_mean))
19   Syy <- sum((y - y_mean)^2)
20
21   Sxx
22   Sxy
23   Syy
24 - ```
```
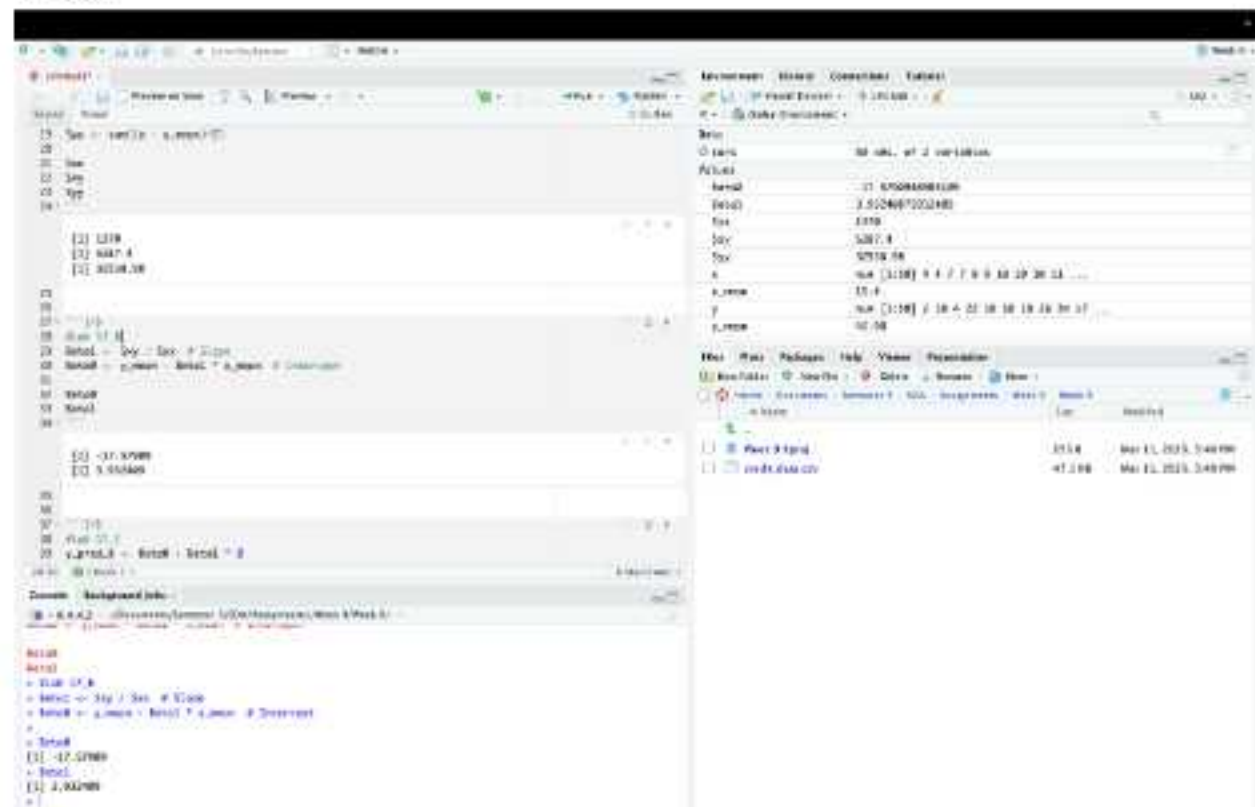
Output:

b) Compare these coefficients you obtained via manual calculate to the ones obtained using the lm() function used in class. Are they same/similar?

Code:

```r
26
27 ```{r}
28 #Lab 17_B
29 Beta1 <- Sxy / Sxx  # Slope
30 Beta0 <- y_mean - Beta1 * x_mean  # Intercept
31
32 Beta0
33 Beta1
34 ```
```

```
[1] -17.57909
```

Output:



Ans: - Intercept (Beta0): -17.57909 and Slope (Beta1): 3.932409

- They are different

c) Predict the stopping distance for a speed of 8 mph. You MUST use the Beta0 and Beta1 you found in this calculation.

Code:

```
36
37 * ```{r}
38  #Lab 17_C
39  y_pred_8 <- Beta0 + Beta1 * 8
40  y_pred_8
41 * ```
```

Output:

d) Calculate the residual for a car travelling at 8 mph. What you need to do is obtain the stopping distance for a car travelling at 8 mph using the regression equation manually with Beta0 and Beta1, then compare this with data from row 5 from the cars dataset. What is the residual?

Code:

```r
43
44  ```{r}
45  #Lab 17_D
46  actual_y_8 <- cars[5, "dist"]
47  residual_8 <- actual_y_8 - y_pred_8
48  residual_8
49  ```
```

Output:

2. Load the faithful dataset by typing data(faithful). Do the following:

a) Plot the waiting times against the duration of eruption.

Code:

```r
63
64 - ```{r}
65  # Lab 17_2A
66  plot(faithful$eruptions, faithful$waiting,
67       xlab = "Eruption Duration (mins)",
68       ylab = "Waiting Time (mins)",
69       main = "Eruption Duration vs. Waiting Time",
70       col = "blue", pch = 16)
71
72 - ```
```

Output:

b) Compute the Pearson correlation coefficient between waiting times and duration of eruption. Is there evidence of a linear relationship? Is it: weak, strong, positive, etc.

Code:

```
74
75 ▾ ```{r}
76    # Lab 17_2B
77    correlation <- cor(faithful$eruptions, faithful$waiting)
78    correlation
79 ▾ ```
```

Output:



Ans:

c) Use lm to fit waiting times as a linear function of eruption durations and save the result of the regression function to the variable 'lab_model'.

Code:

```r
# Lab 17_2C
# Fit a linear regression model
lab_model <- lm(waiting ~ eruptions, data = faithful)
```

Output:

d) Print the summary of the model. What does the summary tell you about model performance from the metrics shown?

Code:

```r
88
89 ```{r}
90 # Lab 17_2D
91 summary(lab_model)
92 ```
```

Output:



Ans:- The model shows a strong relationship between eruption duration and waiting time, with 81.15% of waiting time variation explained by eruption length ($R^2$ = 0.8115).

- The relationship is highly significant (F = 1162, p < 2.2e-16), with predictions expected to be within ±5.9 minutes of actual values.

e) Using the 'lab_model' to predict the waiting time for an eruption duration of 3.667 mins.

Code:

```r
93
94  ```{r}
95  # Lab 17_2E
96  new_data <- data.frame(eruptions = 3.667)
97  predicted_waiting <- predict(lab_model, new_data)
98  predicted_waiting
99  ```
```

Output:

## 3. kNN Preliminaries

a) Write a function that returns the Euclidean distance between two Vectors of the same length. Function must return an error if vectors are of different lengths.

Code:

```r
# Lab 17_3A
euclidean_distance <- function(v1, v2) {
  # Check if vectors are of the same length
  if (length(v1) != length(v2)) {
    stop("Error: Vectors must be of the same length.")
  }

  # Compute Euclidean distance
  distance <- sqrt(sum((v1 - v2)^2))
  return(distance)
}
```

Output:

b) Find the Euclidean distance between the following two vectors:

v1 = (10.09, 2.33, 9.71, 101.46) and v2 = (12.21, 9.41, 7.65, 163.12)

Code:

```r
L15
L16
L17 ▾  ```{r}
L18    # Lab 17_38
L19    # Define the vectors
L20    v1 <- c(10.09, 2.33, 9.71, 101.46)
L21    v2 <- c(12.21, 9.41, 7.65, 163.12)
L22
L23    # Compute the Euclidean distance
L24    distance_result <- euclidean_distance(v1, v2)
L25    distance_result
L26 ▾  ```
```

Output:

4. Use the vehicle_info.csv dataset provided to you and load it into R.

Importing the vehicle_info.csv

Code:

```r
128 ▾  ```{r}
129    # Lab 17_4
130    # Load the dataset
131    vehicle_data <- read.csv("vehicle_info.csv")
132
133    # Browse the data
134    head(vehicle_data)
135 ▾  ```
```

Output:

a) Browse the data, view its structure and its dimensionality. How many variables does it have? How many rows does it have?

Code:

```
137
138 ```{r}
139  # Lab 17_4A
140  # View structure and dimensionality
141  str(vehicle_data)
142  dim(vehicle_data)
143 ```
```

Output:



Ans: It has 74 obs of 10 variables

b) Is there a "label" variable suited for classification purposes? Which variable is that?

Code:

```
44
45 ▾ ```{r}
46   #lab 17_4B
47   # Identify the label variable for classification
48   names(vehicle_data)
49   cat("\nUnique vales in the 'origin'column:\n")
50   table(vehicle_data$origin)
51 ▾ ```
```

Output:



That is the reason we can use "Origin" column for classification

c) Normalize the data (excluding the "label" variable). The resulting values of all variables should be in the range from 0 to 1.

Code:

```r
145
146  ```{r}
147  #lab 17_4C
148  # Function to normalize values between 0 and 1
149  normalize <- function(x) {
150      return((x - min(x)) / (max(x) - min(x)))
151  }
152
153  # Create a copy of the dataframe for normalization
154  vehicle_data_norm <- vehicle_data
155
156  # Normalize all numeric columns (exclude the "origin" column)
157  vehicle_data_norm[, 2:ncol(vehicle_data)] <- lapply(vehicle_data[, 2:ncol(vehicle_data)], normalize)
158
159  # Check the normalized data
160  head(vehicle_data_norm)
161  ```
```

Output:

d) Create a training and test set. Use 80% of the records for the training set and 20% for test set. Make sure to sample records so that the class labels in the training and testing sets are not skewed.

Code:

```r
163
164 - ```{r}
165  #Lab 17_4D
166  # Load required libraries
167  library(caret)
168  library(ggplot2)
169  |
170  # Set seed for reproducibility
171  set.seed(123)
172
173  # Stratified sampling to maintain class distribution
174  train_index <- createDataPartition(vehicle_data_norm$origin,
175                                      p = 0.8,
176                                      list = FALSE,
177                                      times = 1)
178
179  # Create training and test sets
180  train_data <- vehicle_data_norm[train_index, ]
181  test_data <- vehicle_data_norm[-train_index, ]
182
183  # Check distribution of classes in both sets
184  table(train_data$origin)
185  table(test_data$origin)
186 - ```
```

Output:

e) Build a knn() classifier using the "class" package. Set k=9.

Code:

```r
188
189 ```{r}
190  # Lab 17_4E
191  # Load required package
192  library(class)
193
194  # Prepare training and test sets
195  train_labels <- train_data$origin
196  test_labels <- test_data$origin
197
198  # Remove the label column from the feature sets
199  train_features <- train_data[, -1]  # Excluding the "origin" column
200  test_features <- test_data[, -1]    # Excluding the "origin" column
201
202  # Build KNN model with k=9
203  knn_pred <- knn(train = train_features,
204                  test = test_features,
205                  cl = train_labels,
206                  k = 9)
207 ```
208
```

Output:

f) Get the confusionMatrix using the "caret" package. Are there any misclassified labels?

Code:

```r
09
10    ```{r}
11    # Lab 17_4F
12    # Load required library
13    library(caret)
14
15    # Create confusion matrix
16    conf_matrix <- confusionMatrix(data = as.factor(knn_pred),
17                                   reference = as.factor(test_labels))
18
19    # Display confusion matrix
20    conf_matrix
21    ```
```

Output:

g) What is the accuracy? Is it a good performance?

```
         usa         1    8

                 Accuracy : 0.7857
                   95% CI : (0.492, 0.9534)
     No Information Rate : 0.7143
     P-Value [Acc > NIR] : 0.4001

                    Kappa : 0.5116

 Mcnemar's Test P-Value : 1.0000

              Sensitivity : 0.7500
              Specificity : 0.8000
           Pos Pred Value : 0.6000
           Neg Pred Value : 0.8889
               Prevalence : 0.2857
           Detection Rate : 0.2143
     Detection Prevalence : 0.3571
        Balanced Accuracy : 0.7750

         'Positive' Class : other
```

Ans: - The accuracy of the model is 0.7857 (78.57%)

- The model performs moderately well but may not be significantly better than random guessing.

- If the application requires high precision or recall, further model tuning or feature engineering might be needed.

## 5. Gen AI Prompt Engineering:

A telecommunications company wants to predict which customers are likely to churn (cancel their service) based on their usage patterns and demographic information. The company has provided historical data on customers who have either stayed with the service or churned. Your task is to develop a K-Nearest Neighbors (KNN) classifier in Excel to predict customer churn and analyze the factors that contribute most to customer attrition.

**Bot Used: ChatGPT**

Prompt:



customer_data.csv
Spreadsheet

A telecommunications company wants to predict which customers are likely to churn (cancel their service) based on their usage patterns and demographic information. The company has provided historical data on customers who have either stayed with the service or churned. Your task is to develop a K-Nearest Neighbors (KNN) classifier in Excel to predict customer churn and analyze the factors that contribute most to customer attrition.

Give me a guide of steps to run the vba code

Bot:

### Steps Covered in the VBA Code

- Load data from the Excel sheet
- Normalize numerical features using min-max Scaling
- Calculate Euclidean distance for each test instance
- Select k-nearest neighbors and classify based on majority voting
- Output predictions in a new column

### VBA Code for KNN in Excel

1. Open Excel and press ALT + F11 to open the VBA Editor
2. Click Insert > Module and paste the following code:

```
Option Explicit

Function MinMaxNormalize(ByVal value As Double, ByVal minVal As Double, ByVal maxVal ...
    ' Perform Min-Max normalization
    If maxVal = minVal Then
        MinMaxNormalize = 0
    Else
```

```
Option Explicit

Function MinMaxNormalize(ByVal value As Double, ByVal minVal As Double, ByVal maxVal ...
    ' Perform Min-Max normalization
    If maxVal = minVal Then
        MinMaxNormalize = 0
    Else
        MinMaxNormalize = (value - minVal) / (maxVal - minVal)
    End If
End Function

Function EuclideanDistance(ByRef row1 As Range, ByRef row2 As Range) As Double
    ' Compute Euclidean distance between two rows
    Dim i As Integer
    Dim sumSquares As Double
    sumSquares = 0

    For i = 1 To row1.Cells.Count
        sumSquares = sumSquares + (row1.Cells(1, i).Value - row2.Cells(1, i).Value) ^ ...
    Next i

    EuclideanDistance = Sqr(sumSquares)
End Function

Sub RunKNN()
```

Code:

```vba
Sub IdentifyOutliers()
    Dim ws As Worksheet
    Dim lastRow As Long, lastCol As Long
    Dim col As Integer, i As Long
    Dim Q1 As Double, Q3 As Double, IQR As Double
    Dim lowerBound As Double, upperBound As Double

    Set ws = ThisWorkbook.Sheets("customer_data") ' Change sheet name if needed
    lastRow = ws.Cells(ws.Rows.Count, 1).End(xlUp).Row
    lastCol = ws.Cells(1, ws.Columns.Count).End(xlToLeft).Column

    ' Loop through numerical columns (skip categorical ones)
    For col = 2 To lastCol - 1 ' Adjust as needed (exclude ID and target variable)
        ' Sort data in ascending order
        ws.Range(ws.Cells(1, col), ws.Cells(lastRow, col)).Sort Key1:=ws.Cells(1, col), Order1:=xlAscending, Header:=xlYes

        ' Compute Q1 and Q3
        Q1 = Application.WorksheetFunction.Quartile(ws.Range(ws.Cells(2, col), ws.Cells(lastRow, col)), 1)
        Q3 = Application.WorksheetFunction.Quartile(ws.Range(ws.Cells(2, col), ws.Cells(lastRow, col)), 3)
        IQR = Q3 - Q1

        ' Define outlier bounds
        lowerBound = Q1 - 1.5 * IQR
        upperBound = Q3 + 1.5 * IQR

        ' Highlight outliers in red
        For i = 2 To lastRow
            If ws.Cells(i, col).Value < lowerBound Or ws.Cells(i, col).Value > upperBound Then
                ws.Cells(i, col).Interior.Color = RGB(255, 0, 0) ' Red highlight
            End If
        Next i
    Next col

    MsgBox "Outliers identified and highlighted in red", vbInformation, "Outlier Detection"
End Sub


Sub NormalizeFeatures()
    Dim ws As Worksheet
    Dim lastRow As Long, lastCol As Long
    Dim minVal As Double, maxVal As Double
    Dim col As Integer, i As Long

    Set ws = ThisWorkbook.Sheets("customer_data") ' Change sheet name if needed
    lastRow = ws.Cells(ws.Rows.Count, 1).End(xlUp).Row
    lastCol = ws.Cells(1, ws.Columns.Count).End(xlToLeft).Column
```

```vba
Sub NormalizeFeatures()
    Dim ws As Worksheet
    Dim lastRow As Long, lastCol As Long
    Dim minVal As Double, maxVal As Double
    Dim col As Integer, i As Long

    Set ws = ThisWorkbook.Sheets("customer_data") ' Change sheet name if needed
    lastRow = ws.Cells(ws.Rows.Count, 1).End(xlUp).Row
    lastCol = ws.Cells(1, ws.Columns.Count).End(xlToLeft).Column

    ' Loop through numerical columns (skip categorical ones)
    For col = 2 To lastCol - 1
        minVal = Application.WorksheetFunction.Min(ws.Range(ws.Cells(2, col), ws.Cells(lastRow, col)))
        maxVal = Application.WorksheetFunction.Max(ws.Range(ws.Cells(2, col), ws.Cells(lastRow, col)))

        ' Normalize values
        For i = 2 To lastRow
            ws.Cells(i, col).Value = ws.Cells(i, col).Value - minVal / (maxVal - minVal)
        Next i
    Next col

    MsgBox "Features normalized using Min-Max Scaling", vbInformation, "Normalization Complete"
End Sub

Sub PrepareData()
    Dim ws As Worksheet
    Dim lastRow As Long, lastCol As Integer
    Dim trainRow As Long, testRow As Long
    Dim i As Integer, j As Integer
    Dim minVal As Double, maxVal As Double
    Dim rng As Range

    ' Set the worksheet
    Set ws = ThisWorkbook.Sheets("customer_data") ' Adjust sheet name

    ' Find last row and column
    lastRow = ws.Cells(ws.Rows.Count, 1).End(xlUp).Row
    lastCol = ws.Cells(1, ws.Columns.Count).End(xlToLeft).Column

    ' Define train/test split
    trainSize = Application.RoundDown(0.7 * lastRow, 0)
    testSize = lastRow - trainSize

    ' Ensure no empty dataset
    If lastRow < 2 Then
        MsgBox "Not enough data to split", vbExclamation, "Error"
```

Msgbox:

**6. Gen AI Prompt Engineering**: use logistic regression to predict the churn probability (between 0 and 1)

**Bot Used: Claude.ai**

Prompt:



A telecommunications company wants to predict whichcustomers are likely to churn (cancel their service) based on their usage patterns anddemographic information. The company has provided historical data on customers whohave either stayed with the service or churned.Your task is to develop a K-Nearest Neighbors (KNN) classifier in Excel to predictcustomer churn and analyze the factors that contribute most to customer attrition

But you have to Use logistic regression to predict the churn probability (between 0 and 1)

Bot:

Code:

```vba
Sub LogisticRegression()

Dim ws As Worksheet, wsEval As Worksheet
Dim lastRow As Long, featureCount As Integer
Dim i As Long, j As Integer, iter As Integer
Dim learningRate As Double, iterations As Integer
Dim prediction As Double, actual As Double, errorVal As Double
Dim numError As Double
Dim weights() As Double, features() As Double
Dim prob As Double
Dim threshold As Double
Dim TP As Integer, FP As Integer, TN As Integer, FN As Integer
Dim accuracy As Double, precision As Double, recall As Double, F1Score As Double

' Set worksheet
Set ws = ThisWorkbook.Sheets("customer_data")

' Find last row
lastRow = ws.Cells(ws.Rows.count, 1).End(xlUp).Row

' Define features count (Assuming last column is the target variable "Churn")
featureCount = ws.Cells(1, ws.Columns.count).End(xlToLeft).Column - 1

' Hyperparameters
learningRate = 0.01
iterations = 500
threshold = 0.5 ' Probability threshold for classification

' Initialize weights
ReDim weights(1 To featureCount)
For j = 1 To featureCount
    weights(j) = Rnd() ' Assign random initial values
Next j

' Train using Gradient Descent
For iter = 1 To iterations
    numError = 0
    For i = 2 To lastRow ' 0.2 ' Training data (first 70%)
        ReDim features(1 To featureCount)

        ' Extract feature values
        For j = 1 To featureCount
            If IsNumeric(ws.Cells(i, j).value) And Not IsEmpty(ws.Cells(i, j).value) Then
                features(j) = CDbl(ws.Cells(i, j).value)
            Else
                features(j) = 0 ' Handle missing data
            End If
        Next j
```

```vba
        ' Extract feature values
        For j = 1 To featureCount
            If IsNumeric(ws.Cells(i, j).value) And Not IsEmpty(ws.Cells(i, j).value) Then
                features(j) = CDbl(ws.Cells(i, j).value)
            Else
                features(j) = 0 ' Handle missing data
            End If
        Next j

        ' Compute prediction
        prediction = 0
        For j = 1 To featureCount
            prediction = prediction + weights(j) * features(j)
        Next j
        prob = 1 / (1 + Exp(-prediction)) ' Sigmoid function

        ' Assign class based on threshold
        If prob >= threshold Then
            predictedClass = 1
        Else
            predictedClass = 0
        End If

        actual = ws.Cells(i, featureCount + 1).value

        ' Compute confusion matrix
        If predictedClass = 1 And actual = 1 Then TP = TP + 1
        If predictedClass = 1 And actual = 0 Then FP = FP + 1
        If predictedClass = 0 And actual = 0 Then TN = TN + 1
        If predictedClass = 0 And actual = 1 Then FN = FN + 1
    Next i

    ' Compute evaluation metrics
    accuracy = (TP + TN) / (TP + TN + FP + FN)
    If TP + FP > 0 Then precision = TP / (TP + FP) Else precision = 0
    If TP + FN > 0 Then recall = TP / (TP + FN) Else recall = 0
    If precision + recall > 0 Then
        F1Score = 2 * (precision * recall) / (precision + recall)
    Else
        F1Score = 0
    End If

    ' Create Evaluation Sheet
    On Error Resume Next
    Set wsEval = ThisWorkbook.Sheets("Logistic Regression Model")
    If wsEval Is Nothing Then
```

# Lab 18

## Task 1:

Use the 'countries_data.xlsx' dataset.

Code:

```r
5
6  ```{r}
7  # Load required libraries
8  library(readxl)
9  library(dplyr)
10 ```
11
12
13 ```{r}
14 # Lab 18_Task 1
15 file_path <- "countries_data.xlsx"  # Update with the correct path if needed
16 df <- read_excel(file_path)
17 ```
18
```

Output:

## 1. Filter countries where income level is 'High income'.

Code:

```r
18
19 - '''{r}
20   # Lab 18_Task 1-1
21   # 1. Filter countries where income level is 'High income'
22   high_income_countries <- df %>%
23     filter(income_level == "High income")
24
25   # Display result for #1
26   print("High Income Countries:")
27   print(high_income_countries)
28   '''
```

Output:

## 2. Filter countries where region is 'Europe' or 'Americas' and life expectancy is greater than 80 years.

Code:

```r
29
30 - ```{r}
31  # Lab 18_Task 1-2
32  # 2. Filter countries where region is 'Europe' or 'Americas' and life expectancy > 80
33  europe_americas_high_life_exp <- df %>%
34    filter(region %in% c("Europe", "Americas") & life_expect > 80)
35
36  # Display result for #2
37  print("Countries in Europe or Americas with Life Expectancy > 80:")
38  print(europe_americas_high_life_exp)
39  ```
```

Output:

## 3. Filter the top 5 countries based on GDP growth.

Code:

```
41
42 ▾ ```{r}
43   # Lab 18_Task 1-3
44   # 3. Filter the top 5 countries based on GDP growth
45   top_5_gdp_growth <- df %>%
46     arrange(desc(gdp_growth)) %>%
47     head(5)
48
49   # Display result for #3
50   print("Top 5 Countries Based on GDP Growth:")
51   print(top_5_gdp_growth)
52 ▾ ```
```

Output:

4. Create a new variable named NED with a value of 1 when country is equal to Netherlands and 0 when it is not. Select the country name and NED columns side by side and arrange in descending order of NED.

Code:

```r
54
55 - ```{r}
56   # Lab 18_Task 1-4
57   # 4. Create a new variable NED (1 if country is Netherlands, otherwise 0)
58   ned_sorted <- df %>%
59     mutate(NED = ifelse(country == "Netherlands", 1, 0)) %>%
60     select(country, NED) %>%
61     arrange(desc(NED))
62
63   # Display result for #4
64   print("Netherlands Flag Variable (NED), Sorted:")
65   print(ned_sorted)
66 - ```
```

Output:

Task 2:

Use the 'countries_data.xlsx' dataset. Write code that does the following:

1. Filter for income level equal to Lower middle income equal to Education

Code:

```r
67
68
69 - ```{r}
70  # Lab 18_Task 2-1
71  # 1. Filter for income level equal to 'Lower middle income'
72  lower_middle_income <- df %>%
73    filter(income_level == "Lower middle income")
74
75  # Display result for #1
76  print("Filtered for Lower Middle Income:")
77  print(lower_middle_income)
78 - ```
```

Output:

## 2. Group by region

Code:

```r
80
81 - ```{r}
82   # Lab 18_Task 2-2
83   # 2. Group by region
84   grouped_by_region <- lower_middle_income %>%
85     group_by(region)
86
87   # Display result for #2
88   print("Grouped by Region:")
89   print(grouped_by_region)
90 - ```
```

Output:

## 3. Find (1) the median value for unemployment and (2) the number of countries in each region.

Code:

```r
# Lab 18_Task 2-3
# 3. Find median unemployment and count of countries per region
summary_stats <- grouped_by_region %>%
  summarise(
    median_unemployment = median(unemployment, na.rm = TRUE),  # Median unemployment
    num_countries = n()  # Count of countries in each region
  )

# Display result for #3
print("Median Unemployment & Number of Countries by Region:")
print(summary_stats)
```

Output:

4. Sort the results from the highest median optimal level downward.

Code:

```r
108
109   ```{r}
110   # Lab 18_Task 2-4
111   # 4. Sort from highest to lowest median unemployment
112   sorted_result <- summary_stats %>%
113     arrange(desc(median_unemployment))
114
115   # Display result for #4
116   print("Sorted by Highest Median Unemployment:")
117   print(sorted_result)
118   ```
```

Output:

Task 3:

Suppose we work for a housing developer like Toll Brothers (NYSE: TOL) and want to allocate resources to marketing and financing the building of luxury homes in major US metropolitan areas. We have data for a test market in the file 'hprice.csv'.

Do the following:

Code:

```r
# Load required libraries
library(tidyverse)
library(kableExtra)
```

```r
# Lab 18_Task 3
# Load the dataset
hprice <- read.csv("hprice.csv")
```

1. Create a Pivot table to filter the most valuable (higher price) neighborhoods? (Greater than 9999)

Code:

```r
134
135 ```{r}
136 # Lab 18_Task 3-1
137 # Filter neighborhoods where price > 9999
138 valuable_neighborhoods <- hprice %>%
139   filter(Price > 9999)
140
141 # Count the number of valuable neighborhoods
142 num_valuable_neighborhoods <- nrow(valuable_neighborhoods)
143 print(num_valuable_neighborhoods)
144 ```
```

Output:

## 2. Group the Pivot table data by Brick and Neighborhood.

Code:

```
145
146 ˅ ```{r}
147   # Lab 18_Task 3-2
148   # Group the data by Brick and Neighborhood
149   grouped_data <- high_value_properties %>%
150     group_by(Brick, Neighborhood)
151 ˅ ```
152
```

Output:

## 3. Fetch the Average Price, Average Square Foot, and Price Per Square Foot

Code:

```r
```{r}
# Lab 18_Task 3-3
# Calculate Average Price, Average Square Foot, and Price Per Square Foot
summary_data <- grouped_data %>%
  summarise(
    Avg_Price = mean(Price),
    Avg_Sqft = mean(Sqft),
    Price_Per_SqFt = mean(Price / SqFt)
  )
```
```

Output:

## 4. Visualize 1-3 in a table (Use the kableExtra package and knitr:: command)
Code:

```r
165
166
167 ```{r}
168 # Lab 18_Task 3-4
169 # Visualize the data in a table
170 kable(summary_data, format = "html", caption = "Summary of High-Value Neighborhoods") %>%
171   kable_styling(full_width = FALSE, bootstrap_options = c("striped", "hover", "condensed"))
172 ```
```

Output:

Task 4:

## Gen AI Prompt Engineering:

You work as a data analyst for a retail company. The marketing department has provided you with customer transaction data in an Excel file, customer_sales.csv, and wants you to perform customer segmentation analysis based on how much customers spend and calculate potential impact of targeted marketing campaigns. Use spending-based customer classification system that segments customers into three categories:

"High Spender" - Customers who spent $300 or more

"Medium Spender" - Customers who spent between $150 and $299.99

"Low Spender" - Customers who spent less than $150

## Bot Used: Claude.AI

Prompt:

Bot:

## Code:

This the excel we got after running R script.

- As we can see the bot's query resulted in a clear response and we have also achieved the required result.