UNIVERSITI TUN HUSSEIN ONN MALAYSIA

FACULTY OF COMPUTER SCIENCE AND INFORMATION TECHNOLOGY

(FSKTM)

SEMESTER II 2024/2025

DATA MINING

BIT 33603

SECTION 03

LAB ASSIGNMENT 04

**TITLE**

DATA PREPROCESSING USING R

**LECTURER'S NAME**

DR. ROZITA BINTI ABDUL JALIL

| NAME | TUAN KHALIDAH SYAZWANA BINTI TUAN MOHD KASMAWI |
|---|---|
| MATRIC NUMBER | AI220118 |
| DATE SUBMISSION | April 09, 2025 |

# LAB ACTIVITY 4

**Topic:** Data Preprocessing Using R

**Objectives:**

1. To understand basic data preprocessing techniques in R including data cleaning, encoding, and splitting datasets.
2. To apply missing value imputation and data normalization methods to prepare data for analysis.
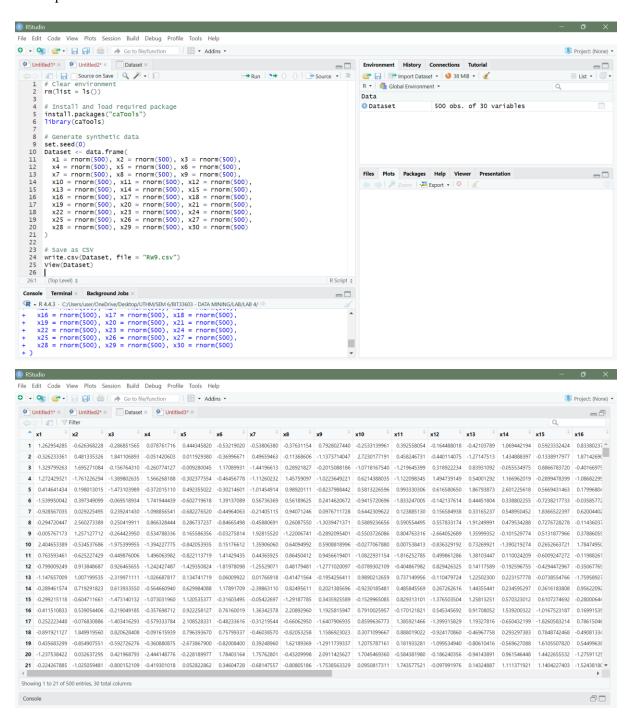
**Duration:** 2 hours

---

**Assessment Question:**

1. Run the provided code in R (Activity 1-4) and understanding the data preprocessing.
2. Submit the visualizations as image/data snapshots for each activity (before and after) along with a brief explanation of the insights gained.

---

**Submission Guidelines:**

1. Submit your solution/answer as a report or document in a single file (.pdf or .docx format).
2. Include a cover page that contains your name, matrix number, and lab name.
3. On the following page, insert screenshots of each activity.
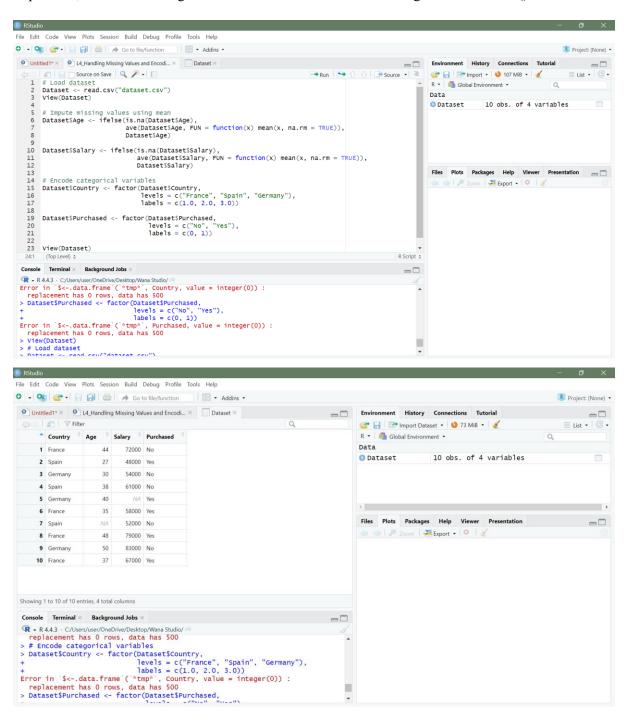4. Submit your lab exercise through AUTHOR.

# Activity 1: Creating a Synthetic Dataset

In this activity, you will generate a synthetic dataset using random value rnorm() and save it as a CSV file. This exercise helps you understand how data can be simulated for practice in data mining techniques.
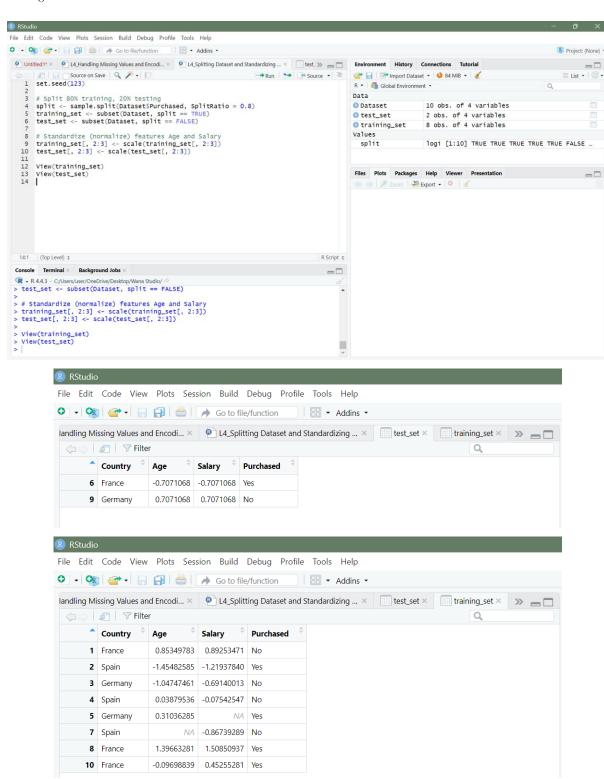
```r
# Clear environment
rm(list = ls())

# Install and load required package
install.packages("caTools")
library(caTools)

# Generate synthetic data
set.seed(0)
Dataset <- data.frame(
  x1 = rnorm(500), x2 = rnorm(500), x3 = rnorm(500),
  x4 = rnorm(500), x5 = rnorm(500), x6 = rnorm(500),
  x7 = rnorm(500), x8 = rnorm(500), x9 = rnorm(500),
  x10 = rnorm(500), x11 = rnorm(500), x12 = rnorm(500),
  x13 = rnorm(500), x14 = rnorm(500), x15 = rnorm(500),
  x16 = rnorm(500), x17 = rnorm(500), x18 = rnorm(500),
  x19 = rnorm(500), x20 = rnorm(500), x21 = rnorm(500),
  x22 = rnorm(500), x23 = rnorm(500), x24 = rnorm(500),
  x25 = rnorm(500), x26 = rnorm(500), x27 = rnorm(500),
  x28 = rnorm(500), x29 = rnorm(500), x30 = rnorm(500)
)

# Save as CSV
write.csv(Dataset, file = "RW9.csv")
View(Dataset)
```

**Activity 2: Handling Missing Values and Encoding Categorical Data**

In this activity, you will learn to load an existing dataset, detect and handle missing values using mean imputation, and encode categorical variables into numeric form using factor the factor() function.

## Activity 3: Splitting Dataset and Standardizing Features

You will split the dataset into training and testing subsets (80/20 split) and apply feature scaling using scale().

## Activity 4: Handling outlier and Normalization

Detect and remove outliers using the IQR method. Then, apply min-max normalization.