# K.R. MANGALAM UNIVERSITY, GURUGRAM-122103

## SCHOOL OF ENGENIERRING AND TECHNOLOGY

## ASSIGNMENT 1

## Data Analysis with Power BI & KNIME

## ETMMML174



| Department:  SOET | Session:  2025-27 |
|---|---|
| Program:  MCA (AI & ML) | Semester:  1 |
| Course Code:  ETMMML174 | College Roll no:  2501940069 |
| Course Name: Data Analysis with Power BI & KNIME | |
| Submitted by: Tanya Bhatia | Faculty: Mr. Mohammad Ajaz |

1) Read the adult.csv file available in the data folder on the KNIME Hub. The data are provided by the UCI Machine Learning Repository.

2) Calculate the count and average age of women with income >50K

3) Calculate the averages of all numerical columns for each one of the 4 groups defined by sex and income values

4) Calculate

- the number of missing values in the occupation column
- the number of non-missing rows in the occupation column
- the number of rows in the occupation column
- the number of rows in the marital-status column

Notice that the last two aggregations should provide the same numbers!

1) Read the adult.csv file



2) A) Filter Female and Income >50k using Row Filter

2) B) Calculate the Count and Average age of women with income >50k



3) Calculate the averages of all numerical columns for each one of the 4 groups defined by sex and income value

4) Calculate:
- the number of **missing values** in the *occupation* column
- the number of **non-missing rows** in the *occupation* column
- the **number of rows** in the *occupation* column
- the **number of rows** in the *marital-status* column