# Assignment 1

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!

**Step 1:** Read CSV File "adult.csv"

# Power BI and KNIME Assignment

**Step 2:** Filter Row for Women with income >50K



**Step 3:** Use GroupBy node to calculate the count and average age of women with income >50K

# Power BI and KNIME Assignment

**Step 4:** Use GroupBy node to calculate the average of all numerical column for each of the 4-group defined by sex and income value



**Step 5:** Use GroupBy node to calculate Missing value count for occupation, non-missing value count for occupation, no of rows in occupation column, no of rows in martial-status

Aman Kumar| MCA (AI & ML) | 2501940051

# Assignment 2

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 average age and count for each one of the 4 groups defined by sex and income values

3) Join the two aggregated values to the original table

**Step 1:** Read the adult.csv file

# Power BI and KNIME Assignment

**Step 2:** Calculate the average age and count for each one of the 4 groups defined by sex and income values



**Step 3:** Join the two aggregated values to the original value

# Assignment 3

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) Extract people with age between 20 and 40 (both included) and working in a workclass starting with "S"

3) Extract people with age between 40 and 60 (both included) and working in a workclass starting with "P"

4) Concatenate both subsets into a single data table

**Step 1:** Read the adult.csv file

# Power BI and KNIME Assignment

**Step 2:** Extract people with age between 20 and 40 (both included) and working in a work class starting with "S"



**Step 3:** Extract People with age between 40 and 60 (both included) and working in a work class starting with "P"



Aman Kumar| MCA (AI & ML) | 2501940051

**Step 4:** Concatenate both subsets into a single data