# Knime - 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
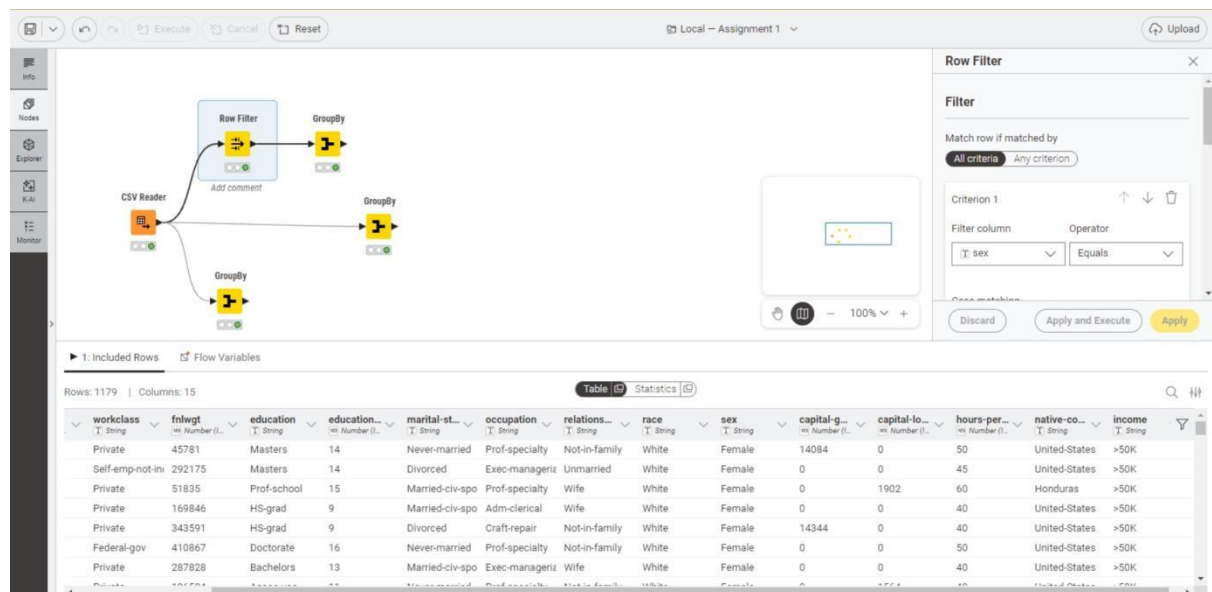- the number of rows in the marital-status column

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

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

**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

**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



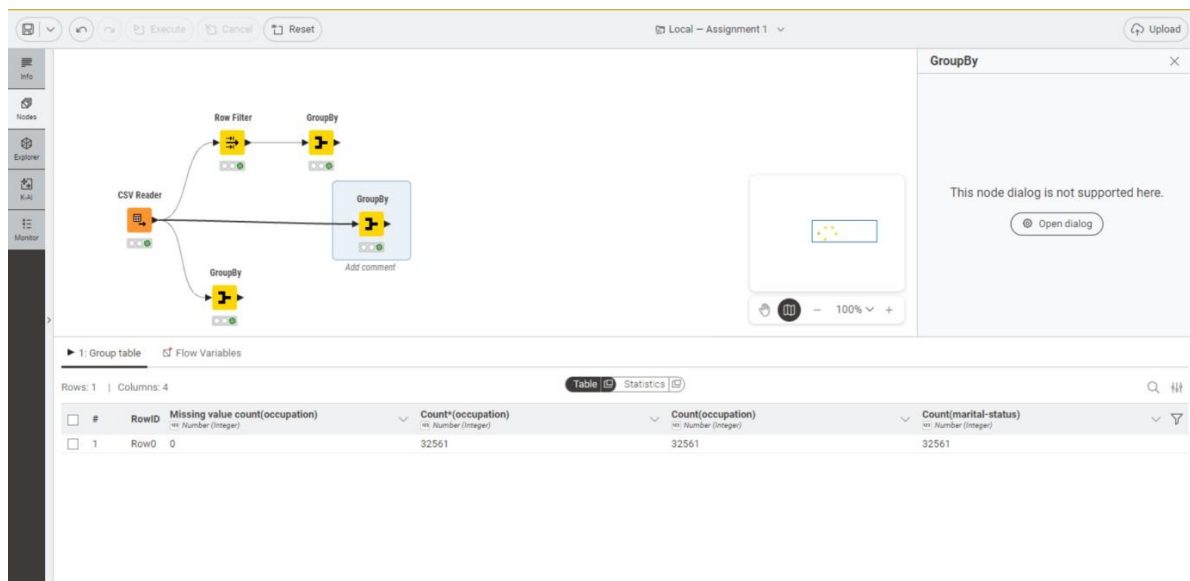| # | RowID | sex (String) | income (String) | Mean(age) (Number (Float)) | Mean(capital-gain) (Number (Float)) | Mean(capital-loss) (Number (Float)) | Mean(education-num) (Number (Float)) | Mean(hours-per-week) (Number (Float)) |
|---|-------|--------|---------|-----------|----------------|----------------|------------------|--------------------|
| 1 | Row0 | Female | <=50K | 36.211 | 121.986 | 47.364 | 9.82 | 35.917 |
| 2 | Row1 | Female | >50K | 42.126 | 4,200.389 | 173.649 | 11.787 | 40.427 |
| 3 | Row2 | Male | <=50K | 37.147 | 165.724 | 56.807 | 9.452 | 40.694 |
| 4 | Row3 | Male | >50K | 44.626 | 3,971.766 | 198.78 | 11.581 | 46.366 |

**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



| # | RowID | Missing value count(occupation) (Number (Integer)) | Count*(occupation) (Number (Integer)) | Count(occupation) (Number (Integer)) | Count(marital-status) (Number (Integer)) |
|---|-------|-----------------------------------|------------------|-----------------|----------------------|
| 1 | Row0 | 0 | 32561 | 32561 | 32561 |