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

**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