```
Widgets
                                                                                                                      Python 3 (ipykernel) O
                                                                                                            Trusted
                                                    ► Run ■ C → Markdown
         Occupation
         Introduction:
         Special thanks to: <a href="https://github.com/justmarkham">https://github.com/justmarkham</a> for sharing the dataset and materials.
         Step 1. Import the necessary libraries
 In [1]: import pandas as pd
         import numpy as np
         Step 2. Import the dataset from this address.
         Step 3. Assign it to a variable called users.
 In [5]: df = pd.read_csv("https://raw.githubusercontent.com/justmarkham/DAT8/master/data/u.user", sep="||")
 Out[5]:
              user_id age gender occupation zip_code
                                            85711
                  2 53
                                            94043
                                            32067
                                            43537
                                            15213
                  5 33
                                     other
                                            33319
                 939 26
                                    student
                                            02215
                                            97229
                 943 22
                                   student
         943 rows x 5 columns
         Step 4. Discover what is the mean age per occupation
 In [6]: df.groupby("occupation").age.mean()
 Out[6]: occupation
         administrator
                         38.746835
                         31.392857
         artist
                         43.571429
         doctor
                         42.010526
         educator
         engineer
                         36.388060
                         29.222222
         entertainment
                         38.718750
         executive
         healthcare
                         41.562500
                         32.571429
         homemaker
         lawyer
                         36.750000
                         40.000000
         librarian
         marketing
                         37.615385
                         26.555556
         none
                         34.523810
         other
                         33.121212
         programmer
                         63.071429
         retired
                         35.666667
         salesman
                         35.548387
         scientist
                         22.081633
         student
         technician
                         33.148148
         writer
                         36.311111
         Name: age, dtype: float64
         Step 5. Discover the Male ratio per occupation and sort it from the most to the least
 In [7]: df['is_male'] = df.gender.apply(lambda x: True if x == 'M' else False)
         df['is_male']
 Out[7]: 0
                 True
                False
                 True
                 True
                False
                ...
         938
                False
         939
                True
         940
                True
         941
                False
         942
                 True
         Name: is_male, Length: 943, dtype: bool
         Step 6. For each occupation, calculate the minimum and maximum ages
 In [9]: df.groupby('occupation').age.agg(['min', 'max'])
 Out[9]:
                     min max
            occupation
          administrator 21 70
                     19
                artist
               doctor 28 64
             educator 23 63
                     22 70
              engineer
                      15 50
          entertainment
                      22 69
             executive
                      22 62
                     20
            homemaker
               lawyer 21 53
              librarian 23 69
                      24 55
             marketing
                      11 55
                none
                      13 64
           programmer 20 63
               retired 51 73
                     18 66
             salesman
              scientist 23 55
                       7 42
               student
                     21 55
             technician
                writer 18 60
         Step 7. For each combination of occupation and gender, calculate the mean age
In [10]: df.groupby(['occupation', 'gender']).age.mean()
Out[10]: occupation
                       gender
         administrator F
                                 40.638889
                                 37.162791
         artist
                                 30.307692
                                 32.333333
                                 43.571429
         doctor
                                 39.115385
         educator
                                 43.101449
                                 29.500000
         engineer
                                 36.600000
         entertainment F
                                 31.000000
                                 29.000000
         executive
                                 44.000000
                                 38.172414
         healthcare
                                 39.818182
                                 45.400000
         homemaker
                                 34.166667
                                 23.000000
         lawyer
                                 39.500000
                                 36.200000
         librarian
                                 40.000000
                                 40.000000
         marketing
                                 37.200000
                                 37.875000
                                 36.500000
         none
                                 18.600000
         other
                                 35.472222
                                 34.028986
                                 32.166667
         programmer
                                 33.216667
         retired
                                 70.000000
                                 62.538462
         salesman
                                 27.000000
                                 38.555556
         scientist
                                 28.333333
                                 36.321429
         student
                                 20.750000
                                 22.669118
         technician
                                 38.000000
                                 32.961538
         writer
                                 37.631579
                                 35.346154
         Name: age, dtype: float64
         Step 8. For each occupation present the percentage of women and men
In [12]: occup_count = df.groupby(['occupation']).count()
         occup_count
Out[12]:
                     user_id age gender zip_code is_male
            occupation
                                                   79
                         79 79
                                    79
                                            79
          administrator
                         28 28
                                            28
                                                   28
                                    28
                artist
               doctor
                         95 95
                                            95
                                                   95
              educator
                                            67
                                                   67
                         67 67
              engineer
                         18 18
                                    18
                                            18
                                                   18
          entertainment
                         32 32
                                    32
                                            32
                                                   32
             executive
                         16 16
                                    16
                                            16
                                                   16
             healthcare
                          7 7
            homemaker
                         12 12
                                            12
               lawyer
                                                   12
              librarian
                         51 51
                                            51
                                                   51
                         26 26
                                            26
                                                   26
             marketing
                none
                          9 9
                        105 105
                                           105
                                   105
                                                   105
                other
                         66 66
                                            66
           programmer
                         14 14
                                    14
                                            14
               retired
                         12 12
                                            12
             salesman
                         31 31
                                   31
                                            31
                                                   31
              scientist
                         196 196
                                   196
                                           196
                                                   196
               student
                         27 27
                                                   27
             technician
                writer
                         45 45
In [13]: df.describe(include="all")
Out[13]:
                   user_id
                               age gender occupation zip_code is_male
                                                943
                                                       943
                                                               943
           count 943.000000 943.000000
                                      943
                                                       795
                     NaN
                               NaN
                                                21
          unique
                                                      55414
                     NaN
                               NaN
                                             student
                                                              True
                     NaN
                               NaN
                                      670
                                                196
                                                               670
            freq
           mean 472.000000
                          34.051962
                                                              NaN
                                      NaN
                                               NaN
                                                       NaN
             std 272.364951
                           12.192740
                                      NaN
                                               NaN
                                                       NaN
                                                              NaN
                  1.000000
                           7.000000
                                      NaN
                                               NaN
                                                       NaN
                                                              NaN
            25% 236.500000
                          25.000000
                                      NaN
                                               NaN
                                                       NaN
                                                              NaN
            50% 472.000000
                          31.000000
                                      NaN
                                               NaN
                                                       NaN
                                                              NaN
            75% 707.500000
                          43.000000
                                      NaN
                                               NaN
                                                       NaN
                                                              NaN
```

max 943.000000 73.000000

NaN

NaN

NaN

NaN

Jupyter Project 12 Occupation (autosaved)

Logout