```
Widgets
                                                                                                                        Python 3 (ipykernel) O
                                                                                                             Trusted
               Insert
                                                     ✓ ☐ Autosave interval (min): off
         Ex3 - Getting and Knowing your Data
         This time we are going to pull data directly from the internet. 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 and use the 'user_id' as index
In [25]: df = pd.read_table("https://raw.githubusercontent.com/justmarkham/DAT8/master/data/u.user", sep='|', index_col='user_id')
Out[25]:
                 age gender occupation zip_code
          user_id
                                         85711
              1 24
              2 53
                                         94043
              3 23
                                         32067
              4 24
                                         43537
               5 33
                                         15213
                                         33319
                                         02215
                         M administrator
                                         97229
                                student
                                         78209
                                librarian
                                         77841
                                student
         943 rows × 4 columns
         Step 4. See the first 25 entries
In [26]: df.head(25)
                                         29206
                                educator
                                         55106
                                scientist
                                         97301
              15 49
                                educator
              16 21
                                         10309
                         M entertainment
              17 30
                                         06355
                         M programmer
              18 35
                                         37212
                                  other
              19 40
                                         02138
                                librarian
                                         95660
                             homemaker
              21 26
                                         30068
              22 25
                                         40206
                         М
                                  writer
              23 30
                                         48197
                                  artist
              24 21
                                         94533
                                  artist
              25 39
                                         55107
                                engineer
         Step 5. See the last 10 entries
In [27]: df.tail(10)
Out[27]:
                 age gender occupation zip_code
          user_id
             934 61
                                         22902
                               engineer
                                         66221
             935 42
                                doctor
             936 24
                                        32789
                                 other
             937 48
                                         98072
                               educator
                                         55038
                              technician
             939 26
                                         33319
                                student
             940 32
                                        02215
                         M administrator
                                        97229
                 20
                                student
                                         78209
                                librarian
             943 22
                                        77841
                                student
         Step 6. What is the number of observations in the dataset?
In [28]: df.shape[0]
Out[28]: 943
         Step 7. What is the number of columns in the dataset?
In [29]: df.shape[1]
Out[29]: 4
         Step 8. Print the name of all the columns.
In [30]: df.columns
Out[30]: Index(['age', 'gender', 'occupation', 'zip_code'], dtype='object')
         Step 9. How is the dataset indexed?
In [31]: df.index
Out[31]: Int64Index([ 1, 2, 3, 4, 5, 6, 7, 8, 9, 10,
                     934, 935, 936, 937, 938, 939, 940, 941, 942, 943],
                    dtype='int64', name='user_id', length=943)
         Step 10. What is the data type of each column?
In [32]: df.dtypes
Out[32]: age
                        int64
         gender
                       object
         occupation
                      object
         zip_code
                       object
         dtype: object
         Step 11. Print only the occupation column
In [35]: df.occupation
Out[35]: user_id
                        other
                       writer
                   technician
                        other
         939
                      student
         940
                administrator
         941
                      student
         942
                    librarian
         943
                      student
         Name: occupation, Length: 943, dtype: object
         Step 12. How many different occupations are in this dataset?
In [38]: df.occupation.nunique()
Out[38]: 21
         Step 13. What is the most frequent occupation?
In [45]: df.occupation.value_counts().head(3)
Out[45]: student
                     196
                     105
         other
                     95
         educator
         Name: occupation, dtype: int64
         Step 14. Summarize the DataFrame.
In [48]: df.describe()
Out[48]:
                      age
          count 943.000000
                34.051962
            std 12.192740
                 7.000000
           25% 25.000000
           50% 31.000000
           75% 43.000000
           max 73.000000
         Step 15. Summarize all the columns
In [50]: df.describe(include="all")
Out[50]:
                      age gender occupation zip_code
                                              943
           count 943.000000
                                       943
                                              795
                      NaN
                                       21
                      NaN
                                    student
                                             55414
            freq
                      NaN
                            NaN
                 34.051962
                                      NaN
                 12.192740
                            NaN
                                              NaN
                                      NaN
                            NaN
                                              NaN
                  7.000000
                                      NaN
                            NaN
                 25.000000
                                      NaN
                                              NaN
                 31.000000
                            NaN
                                              NaN
                                      NaN
                 43.000000
                            NaN
                                      NaN
            max 73.000000
                            NaN
                                      NaN
                                              NaN
         Step 16. Summarize only the occupation column
In [53]: df.occupation.describe()
Out[53]: count
         unique
         top
                   student
         freq
         Name: occupation, dtype: object
In [54]: df.age.describe()
Out[54]: count
                  943.000000
                   34.051962
         std
                   12.192740
                   7.000000
         min
         25%
                   25.000000
         50%
                   31.000000
         75%
                   43.000000
                   73.000000
         Name: age, dtype: float64
         Step 17. What is the mean age of users?
In [57]: round(df.age.mean())
Out[57]: 34
         Step 18. What is the age with least occurrence?
In [58]: df.age.value_counts().tail()
```

Out[58]: 7 1

Name: age, dtype: int64

Jupyter project 7 Users Last Checkpoint: an hour ago (autosaved)

Logout