

```
In [1]: pip install numpy

Requirement already satisfied: numpy in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (1.21.2)
Note: you may need to restart the kernel to use updated packages.
WARNING: You are using pip version 21.1.3; however, version 21.2.4 is available.
You should consider upgrading via the 'c:\users\yogesh\appdata\local\programs\python\python39\python.exe -m pip install --upgrade pip' command.

In [2]: pip install pandas

Requirement already satisfied: pandas in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (1.3.2)
Requirement already satisfied: pytz>=2017.3 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from pandas) (2021.1)
Requirement already satisfied: numpy>=1.17.3 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from pandas) (1.21.2)
Requirement already satisfied: python-dateutil>=2.7.3 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from pandas) (2.8.1)
Requirement already satisfied: six>=1.5 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from python-dateutil>=2.7.3->pandas) (1.16.0)
Note: you may need to restart the kernel to use updated packages.
WARNING: You are using pip version 21.1.3; however, version 21.2.4 is available.
You should consider upgrading via the 'c:\users\yogesh\appdata\local\programs\python\python39\python.exe -m pip install --upgrade pip' command.

In [3]: pip install matplotlib

Requirement already satisfied: matplotlib in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (3.4.3)
Requirement already satisfied: cycler>=0.10 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from matplotlib) (0.10.0)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from matplotlib) (1.3.1)
Requirement already satisfied: pillow>=6.2.0 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from matplotlib) (8.3.1)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from matplotlib) (2.8.1)
Requirement already satisfied: pyparsing>=2.2.1 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from matplotlib) (2.4.7)
Requirement already satisfied: numpy>=1.16 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from matplotlib) (1.21.2)
Requirement already satisfied: six in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from cycler>=0.10->matplotlib) (1.16.0)
Note: you may need to restart the kernel to use updated packages.
WARNING: You are using pip version 21.1.3; however, version 21.2.4 is available.
You should consider upgrading via the 'c:\users\yogesh\appdata\local\programs\python\python39\python.exe -m pip install --upgrade pip' command.

In [4]: pip install seaborn

Collecting seaborn
Note: you may need to restart the kernel to use updated packages.
WARNING: You are using pip version 21.1.3; however, version 21.2.4 is available.
You should consider upgrading via the 'c:\users\yogesh\appdata\local\programs\python\python39\python.exe -m pip install --upgrade pip' command.
Downloading seaborn-0.11.2-py3-none-any.whl (282 kB)
Collecting scipy>=1.0
  Downloading scipy-1.7.1-cp39-cp39-win_amd64.whl (33.8 MB)
Requirement already satisfied: matplotlib>=2.2 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from seaborn) (3.4.3)
Requirement already satisfied: numpy>=1.15 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from seaborn) (1.21.2)
Requirement already satisfied: pandas>=0.23 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from seaborn) (1.3.2)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from matplotlib>=2.2->seaborn) (1.3.1)
Requirement already satisfied: pyparsing>=2.2.1 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from matplotlib>=2.2->seaborn) (2.4.7)
Requirement already satisfied: pillow>=6.2.0 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from matplotlib>=2.2->seaborn) (8.3.1)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from matplotlib>=2.2->seaborn) (2.8.1)
Requirement already satisfied: cycler>=0.10 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from matplotlib>=2.2->seaborn) (0.10.0)
Requirement already satisfied: six in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from cycler>=0.10->matplotlib>=2.2->seaborn) (1.16.0)
Requirement already satisfied: pytz>=2017.3 in c:\users\yogesh\appdata\local\programs\python\python39\lib\site-packages (from pandas>=0.23->seaborn) (2021.1)
Installing collected packages: scipy, seaborn
Successfully installed scipy-1.7.1 seaborn-0.11.2

In [5]: import pandas as pd
import numpy as np

import matplotlib
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

Import the dataset from this(<https://raw.githubusercontent.com/justmarkham/DAT8/master/data/u.user>).

Use sep="|" while reading the data

```
In [6]: url = 'https://raw.githubusercontent.com/justmarkham/DAT8/master/data/u.user'

In [7]: df = pd.read_csv(url, sep="|")

In [8]: df

Out[8]:
```

	user_id	age	gender	occupation	zip_code
0	1	24	M	technician	85711
1	2	53	F	other	94043
2	3	23	M	writer	32067
3	4	24	M	technician	43537
4	5	33	F	other	15213
...
938	939	26	F	student	33319
939	940	32	M	administrator	02215
940	941	20	M	student	97229
941	942	48	F	librarian	78209
942	943	22	M	student	77841

943 rows × 5 columns

Assign it to a variable called users and use the 'user_id' as index

```
In [9]: users=df.set_index("user_id")

In [10]: users

Out[10]:
```

	age	gender	occupation	zip_code
user_id				
1	24	M	technician	85711
2	53	F	other	94043
3	23	M	writer	32067
4	24	M	technician	43537
5	33	F	other	15213
...
939	26	F	student	33319
940	32	M	administrator	02215
941	20	M	student	97229
942	48	F	librarian	78209
943	22	M	student	77841

943 rows × 4 columns

See the first 10 and last 10 entries

```
In [11]: print("-----First 10 entries -----")
print(users.head(10))
print("-----Last 10 entries -----")
print(users.tail(10))

-----First 10 entries -----
   age gender  occupation zip_code
user_id
1    24      M   technician   85711
2    53      F      other    94043
3    23      M      writer    32067
4    24      M   technician   43537
5    33      F      other    15213
6    42      M   executive   98101
7    57      M administrator   91344
8    36      M administrator   05201
9    29      M      student   01602
10   53      M      lawyer    90703
-----Last 10 entries -----
   age gender  occupation zip_code
user_id
934    61      M   engineer   22902
935    42      M     doctor   66221
936    24      M      other   32789
937    48      M   educator   98072
938    38      F   technician   55038
939    26      F      student   33319
940    32      M administrator   02215
941    20      M      student   97229
942    48      F   librarian   78209
943    22      M      student   77841
```

What is the number of observations in the dataset?

```
In [12]: print("Number of Observations : ",users.shape[0])

Number of Observations : 943
```

What is the number of columns in the dataset?

```
In [13]: print("Number of Columns : ",users.shape[1])

Number of Columns : 4
```

Print the name of all the columns.

```
In [14]: users.columns

Out[14]: Index(['age', 'gender', 'occupation', 'zip_code'], dtype='object')
```

How is the dataset indexed?

```
In [15]: users.index

Out[15]: 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)
```

What is the data type of each column?

```
In [17]: users.dtypes

Out[17]: age          int64
gender       object
occupation   object
zip_code     object
dtype: object
```

Print only the occupation column

```
In [19]: users['occupation']

Out[19]: user_id
1      technician
2      other
3      writer
4      technician
5      other
...
939     student
940  administrator
941     student
942     librarian
943     student
Name: occupation, Length: 943, dtype: object
```

How many different occupations are in this dataset?

```
In [20]: users['occupation'].nunique()

Out[20]: 21
```

What is the most frequent occupation?

```
In [21]: users['occupation'].mode()

Out[21]: 0      student
dtype: object
```

DataFrame Info.

```
In [22]: users.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 943 entries, 1 to 943
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
--  --
 0   age        943 non-null    int64
 1   gender     943 non-null    object
 2   occupation 943 non-null    object
 3   zip_code   943 non-null    object
dtypes: int64(1), object(3)
memory usage: 36.8+ KB
```

Describe all the columns

```
In [23]: users.describe(include="all")

Out[23]:
```

	age	gender	occupation	zip_code
count	943.000000	943	943	943
unique	NaN	2	21	795
top	NaN	M	student	55414
freq	NaN	670	196	9
mean	34.051962	NaN	NaN	NaN
std	12.192740	NaN	NaN	NaN
min	7.000000	NaN	NaN	NaN
25%	25.000000	NaN	NaN	NaN
50%	31.000000	NaN	NaN	NaN
75%	43.000000	NaN	NaN	NaN
max	73.000000	NaN	NaN	NaN

Summarize only the occupation column

```
In [24]: users['occupation'].value_counts()

Out[24]: student      196
other      105
educator     95
administrator 79
engineer     67
programmer   66
librarian    51
writer       45
executive    32
scientist    31
artist       28
technician   27
marketing    26
entertainment 18
healthcare   16
retired      14
lawyer       12
salesman     12
none         9
homemaker    7
doctor        7
Name: occupation, dtype: int64
```

What is the mean age of users?¶

```
In [25]: users["age"].mean()

Out[25]: 34.05196182396607
```

What is the age with least occurrence?

```
In [26]: users['age'].min()

Out[26]: 7
```