Data Mining - Lab - 2

Numpy & Perform Data Exploration with Pandas

Mohil Parmar

23010101192

Numpy

- 1. NumPy (Numerical Python) is a powerful open-source library in Python used for numerical and scientific computing.
- 2. It provides support for large, multi-dimensional arrays and matrices, along with a collection of mathematical functions to operate on them efficiently.
- 3. NumPy is highly optimized and written in C, making it much faster than using regular Python lists for numerical operations.
- 4. It serves as the foundation for many other Python libraries in data science and machine learning, like pandas, TensorFlow, and scikit-learn.
- 5. With features like broadcasting, vectorization, and integration with C/C++ code, NumPy allows for cleaner and faster code in numerical computations.

Step 1. Import the Numpy library

In [5]: import numpy as np

Step 2. Create a 1D array of numbers

```
In [11]: arr = np.arange(10)
         arr
Out[11]: array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
In [13]: arr = np.arange(10,21)
         arr
Out[13]: array([10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20])
In [15]: arr1 =np.array([22,33,44,55])
         arr1
Out[15]: array([22, 33, 44, 55])
         Step 3. Reshape 1D to 2D Array
In [19]: arr2 = np.arange(20).reshape(4,5)
```

```
arr2
Out[19]: array([[ 0, 1, 2, 3, 4],
               [5, 6, 7, 8, 9],
               [10, 11, 12, 13, 14],
               [15, 16, 17, 18, 19]])
```

Step 4. Create a Linspace array

```
In [21]: np.linspace(14,18)
         # by default 50 decimal
```

```
Out[21]: array([14.
                           , 14.08163265, 14.16326531, 14.24489796, 14.32653061,
                14.40816327, 14.48979592, 14.57142857, 14.65306122, 14.73469388,
                14.81632653, 14.89795918, 14.97959184, 15.06122449, 15.14285714,
                15.2244898 , 15.30612245 , 15.3877551 , 15.46938776 , 15.55102041 ,
                15.63265306, 15.71428571, 15.79591837, 15.87755102, 15.95918367,
                16.04081633, 16.12244898, 16.20408163, 16.28571429, 16.36734694,
                16.44897959, 16.53061224, 16.6122449 , 16.69387755, 16.7755102 ,
                16.85714286, 16.93877551, 17.02040816, 17.10204082, 17.18367347,
                17.26530612, 17.34693878, 17.42857143, 17.51020408, 17.59183673,
                17.67346939, 17.75510204, 17.83673469, 17.91836735, 18.
In [23]: np.linspace(14,18,20)
         # in third argument use specify length of number
Out[23]: array([14.
                           , 14.21052632, 14.42105263, 14.63157895, 14.84210526,
                15.05263158, 15.26315789, 15.47368421, 15.68421053, 15.89473684,
                16.10526316, 16.31578947, 16.52631579, 16.73684211, 16.94736842,
                17.15789474, 17.36842105, 17.57894737, 17.78947368, 18.
         Step 5. Create a Random Numbered Array
In [25]: arr4 = np.random.rand(6)
         arr4
Out[25]: array([0.76539575, 0.54073109, 0.61656565, 0.50856795, 0.9464426,
                0.5109816 ])
In [27]: arr5=np.random.rand(6,5)
         arr5
         # for 2D array use give second argument
Out[27]: array([[0.46938984, 0.75824529, 0.79401908, 0.667392 , 0.86403358],
                 [0.06399315, 0.35781957, 0.39032216, 0.40778384, 0.93098397],
                [0.59229555, 0.3266491, 0.20554348, 0.171566, 0.30838983],
                 [0.8795907, 0.31034884, 0.55286617, 0.58000734, 0.36457935],
                 [0.50805045, 0.40600399, 0.88392335, 0.17804145, 0.09062223],
                 [0.83127196, 0.62684204, 0.8553937 , 0.47073148, 0.75429648]])
```

Step 6. Create a Random Integer Array

```
In [29]: np.random.randint(30,60)
Out[29]: 55
In [31]: np.random.randint(20,40,size=5)
Out[31]: array([34, 25, 22, 36, 34])
In [33]: np.random.randint(30,50,size=(3,4))
Out[33]: array([[38, 37, 38, 44],
                [31, 47, 44, 45],
                [36, 32, 45, 47]])
         Step 7. Create a 1D Array and get Max, Min, ArgMax, ArgMin
In [53]: arr6 = np.random.randint(10,40,size=10)
         print(arr6)
         # max(arr6)
         arr6.max()
        [21 27 24 14 37 25 21 10 22 22]
Out[53]: 37
In [55]: arr6.min()
Out[55]: 10
In [57]: arr6.argmax()
Out[57]: 4
In [59]: arr6.argmin()
Out[59]: 7
```

Step 8. Indexing in 1D Array

Out[91]: array([11, 38, 32, 26, 24])

In [95]: arr8[0:3]

```
In [63]: arr7 = np.random.randint(10,50,size=7)
         arr7
Out[63]: array([28, 33, 47, 35, 48, 13, 41])
In [67]: arr7[5]
Out[67]: 13
In [69]: arr7[2:6]
Out[69]: array([47, 35, 48, 13])
In [77]: arr7.dtype
         print(type(arr7))
        <class 'numpy.ndarray'>
         Step 9. Indexing in 2D Array
In [89]: arr8 = np.random.randint(10,40,size=(5,5))
         arr8
Out[89]: array([[23, 17, 23, 10, 37],
                [15, 14, 14, 14, 39],
                [33, 34, 29, 34, 15],
                [11, 38, 32, 26, 24],
                [17, 31, 39, 18, 14]])
In [91]: arr8[3]
```

Step 10. Conditional Selection

```
In [121... arr10 = np.random.randint(20,40,size=10) a1=[i for i in arr10 if i>25] a1

Out[121... [35, 29, 26, 37, 27, 26]

In [123... arr10[arr10>25]

Out[123... array([35, 29, 26, 37, 27, 26])

In [127... arr10[(arr10>25) & (arr10<30)]

Out[127... array([29, 26, 27, 26])
```

You did it! 10 exercises down — you're on fire!

Pandas

Step 1. Import the necessary libraries

```
In [131... import pandas as pd
```

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 [148...

users = pd.read_csv("https://raw.githubusercontent.com/justmarkham/DAT8/master/data/u.user" , sep="|" , index_col="user_id") users

Out[148...

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

943 rows × 4 columns

Step 4. See the first 25 entries

In [150... users.head(25)

Out[150...

	-9-	9	occupation.	p
user_id				
1	24	М	technician	85711
2	53	F	other	94043
3	23	М	writer	32067
4	24	М	technician	43537
5	33	F	other	15213
6	42	М	executive	98101
7	57	М	administrator	91344
8	36	М	administrator	05201
9	29	М	student	01002
10	53	М	lawyer	90703
11	39	F	other	30329
12	28	F	other	06405
13	47	М	educator	29206
14	45	М	scientist	55106
15	49	F	educator	97301
16	21	М	entertainment	10309
17	30	М	programmer	06355
18	35	F	other	37212
19	40	М	librarian	02138
20	42	F	homemaker	95660
21	26	М	writer	30068

age gender occupation zip_code

	age	gender	occupation	zip_code
user_id				
22	25	М	writer	40206
23	30	F	artist	48197
24	21	F	artist	94533
25	39	М	engineer	55107

Step 5. See the last 10 entries

In [152... users.tail(10)

Out[152...

	age	gender	occupation	zip_code
user_id				
934	61	М	engineer	22902
935	42	М	doctor	66221
936	24	М	other	32789
937	48	М	educator	98072
938	38	F	technician	55038
939	26	F	student	33319
940	32	М	administrator	02215
941	20	М	student	97229
942	48	F	librarian	78209
943	22	М	student	77841

Step 6. What is the number of observations in the dataset?

```
In [162... users.shape[0]
```

Step 7. What is the number of columns in the dataset?

```
In [170... users.shape[1]
```

Out[170... 4

Out[162... 943

Step 8. Print the name of all the columns.

```
In [174... users.columns
Out[174... Index(['age', 'gender', 'occupation', 'zip_code'], dtype='object')
```

Step 9. How is the dataset indexed?

Step 10. What is the data type of each column?

```
In [178... users.dtypes
```

```
Out[178... age int64 gender object occupation object zip_code object dtype: object
```

Step 11. Print only the occupation column

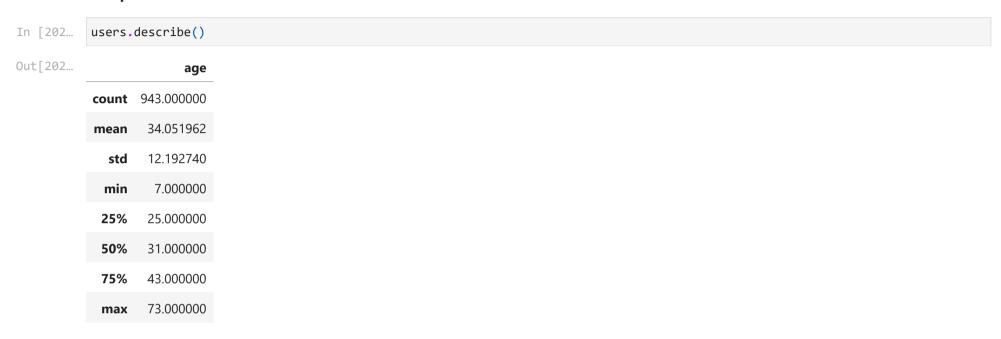
```
users["occupation"]
In [184...
Out[184...
           user id
           1
                      technician
           2
                           other
           3
                          writer
                      technician
                           other
           939
                         student
           940
                  administrator
           941
                         student
           942
                       librarian
                         student
           943
           Name: occupation, Length: 943, dtype: object
```

Step 12. How many different occupations are in this dataset?

Step 13. What is the most frequent occupation?

```
users["occupation"].value_counts()
In [196...
          occupation
Out[196...
           student
                            196
           other
                            105
                             95
           educator
           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
           Name: count, dtype: int64
          users["occupation"].value_counts().head(1)
In [198...
Out[198...
           occupation
           student
                      196
           Name: count, dtype: int64
          users["occupation"].value_counts().idxmax()
In [200...
Out[200...
           'student'
```

Step 14. Summarize the DataFrame.



Step 15. Summarize all the columns

```
In [214... users.describe(include="all")
```

\circ		г	1	1	4
U	uт	١.	7	Τ	4

	age	gender	occupation	zip_code
count	943.000000	943	943	943
unique	NaN	2	21	795
top	NaN	М	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

Step 16. Summarize only the occupation column

```
In [222... users["occupation"].describe()

Out[222... count 943
    unique 21
    top student
    freq 196
    Name: occupation, dtype: object
```

Step 17. What is the mean age of users?

```
In [224... users["age"].mean()
Out[224... 34.05196182396607
```

Step 18. What is the age with least occurrence?

You're not just learning, you're mastering it. Keep aiming higher! 💉