

Data Mining - Lab - 2

Numpy & Perform Data Exploration with Pandas

NAME: AYUSH J. MARADIA

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 [2]: import numpy as np
```

Step 2. Create a 1D array of numbers

```
In [7]: arr = np.arange(10)
arr
```

```
Out[7]: array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
```

```
In [6]: print(arr)
```

```
[0 1 2 3 4 5 6 7 8 9]
```

```
In [8]: arr = np.arange(15,30)
arr
```

```
Out[8]: array([15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29])
```

```
In [10]: arr1 = np.array([20,250,63,1,5,94])
arr1
```

```
Out[10]: array([ 20, 250,  63,   1,   5,  94])
```

```
In [12]: print(type(arr1))
```

```
<class 'numpy.ndarray'>
```

```
In [13]: arr1.dtype
```

```
Out[13]: dtype('int32')
```

Step 3. Reshape 1D to 2D Array

```
In [15]: arr3 = np.arange(20).reshape(4,5)
arr3
```

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

```
In [16]: arr4 = np.arange(30).reshape(6,5)
arr4
```

```
Out[16]: array([[ 0,  1,  2,  3,  4],
               [ 5,  6,  7,  8,  9],
               [10, 11, 12, 13, 14],
               [15, 16, 17, 18, 19],
               [20, 21, 22, 23, 24],
               [25, 26, 27, 28, 29]])
```

```
In [18]: arr5 = np.arange(30).reshape(7,6)
        arr4 #because it always in multiplication
```

```
-----
ValueError                                Traceback (most recent call last)
Cell In[18], line 1
----> 1 arr5 = np.arange(30).reshape(7,6)
      2 arr4

ValueError: cannot reshape array of size 30 into shape (7,6)
```

```
In [22]: arr = np.array([2,5,6,1,2,6]).reshape(2,3)
        print(arr)
```

```
[[2 5 6]
 [1 2 6]]
```

Step 4. Create a Linspace array

```
In [23]: np.linspace(14,18)
```

```
Out[23]: 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 [24]: np.linspace(14,18,25)
```

```
Out[24]: array([14.         , 14.16666667, 14.33333333, 14.5         , 14.66666667,
               14.83333333, 15.         , 15.16666667, 15.33333333, 15.5         ,
               15.66666667, 15.83333333, 16.         , 16.16666667, 16.33333333,
               16.5         , 16.66666667, 16.83333333, 17.         , 17.16666667,
               17.33333333, 17.5         , 17.66666667, 17.83333333, 18.         ])
```

Step 5. Create a Random Numbered Array

```
In [25]: np.random.rand(8)
```

```
Out[25]: array([0.41291694, 0.37105915, 0.21048365, 0.36919309, 0.93240785,
               0.13067258, 0.55556962, 0.349459  ])
```

```
In [27]: np.random.rand(5,7)
```

```
Out[27]: array([[0.78966493, 0.12643212, 0.04609329, 0.69075626, 0.96486792,
               0.70586393, 0.6852699 ],
               [0.09034164, 0.55006824, 0.48374078, 0.79382106, 0.09037752,
               0.80109187, 0.8785915 ],
               [0.14155285, 0.4608573 , 0.12223327, 0.96604649, 0.86228412,
               0.50605056, 0.58928582],
               [0.7779563 , 0.4560097 , 0.78406501, 0.03932901, 0.29236888,
               0.3336317 , 0.63670783],
               [0.45075189, 0.47405161, 0.21609842, 0.39692943, 0.09535095,
               0.89046195, 0.24911291]])
```

Step 6. Create a Random Integer Array

```
In [30]: np.random.randint(25,56)
```

```
Out[30]: 43
```

```
In [32]: np.random.randint(25,56, size = 6)
```

```
Out[32]: array([37, 37, 41, 35, 52, 27])
```

```
In [33]: np.random.randint(25,56, size = (3,3))
```

```
Out[33]: array([[34, 50, 31],
               [43, 47, 44],
               [34, 42, 52]])
```

Step 7. Create a 1D Array and get Max,Min,ArgMax,ArgMin

```
In [35]: arr1 = np.random.randint(10,20,10)
arr1
```

```
Out[35]: array([13, 12, 13, 14, 13, 13, 16, 13, 17, 10])
```

```
In [36]: arr1.max()
```

```
Out[36]: 17
```

```
In [37]: arr1.min()
```

```
Out[37]: 10
```

```
In [39]: arr1.argmax() #it give you index of max element (if repeated then first index will be returned.)
```

```
Out[39]: 8
```

```
In [41]: arr1.argmin() #it give you index of min element
```

```
Out[41]: 9
```

Step 8. Indexing in 1D Array

```
In [43]: arr4 = np.random.randint(20,40,size = 10)
arr4
```

```
Out[43]: array([37, 26, 38, 39, 36, 22, 38, 22, 33, 39])
```

```
In [44]: print(arr4[5])
```

```
22
```

```
In [45]: print(arr4[3:5])
```

```
[39 36]
```

Step 9. Indexing in 2D Array

```
In [48]: arr = np.random.randint(1,10, (5,5))
arr
```

```
Out[48]: array([[2, 3, 5, 6, 6],
               [5, 6, 5, 8, 5],
               [7, 7, 4, 7, 4],
               [9, 5, 6, 4, 6],
               [2, 2, 2, 4, 3]])
```

```
In [49]: arr[2]
```

```
Out[49]: array([7, 7, 4, 7, 4])
```

```
In [50]: arr[1:4]
```

```
Out[50]: array([[5, 6, 5, 8, 5],
               [7, 7, 4, 7, 4],
               [9, 5, 6, 4, 6]])
```

```
In [51]: arr[3][3]
```

```
Out[51]: 4
```

Step 10. Conditional Selection

```
In [52]: arr = np.random.randint(20,40,10)
print()
```

```
Out[52]: array([23, 24, 22, 24, 39, 24, 23, 30, 25, 32])
```

```
In [53]: arr > 25
```

```
Out[53]: array([False, False, False, False,  True, False, False,  True, False,
               True])
```

```
In [55]: arr[(arr > 25) & (arr < 35)]
```

```
Out[55]: array([30, 32])
```

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

Pandas

Step 1. Import the necessary libraries

```
In [56]: import pandas as pd
```

Step 2. Import the dataset from this [address](https://raw.githubusercontent.com/justmarkham/DAT8/master/data/u.user).

Step 3. Assign it to a variable called users and use the 'user_id' as index

```
In [58]: users = pd.read_csv("https://raw.githubusercontent.com/justmarkham/DAT8/master/data/u.user")
users
```

```
Out[58]:
```

	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

```
In [59]: users = pd.read_csv("https://raw.githubusercontent.com/justmarkham/DAT8/master/data/u.user", sep = "|")
users
```

```
Out[59]:
```

	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

```
In [60]: users = pd.read_csv("https://raw.githubusercontent.com/justmarkham/DAT8/master/data/u.user", sep = "|", index_col=0)
users
```

Out[60]:

	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

Step 4. See the first 25 entries

In [61]: `users.head(25)`

Out[61]:

	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	01002
10	53	M	lawyer	90703
11	39	F	other	30329
12	28	F	other	06405
13	47	M	educator	29206
14	45	M	scientist	55106
15	49	F	educator	97301
16	21	M	entertainment	10309
17	30	M	programmer	06355
18	35	F	other	37212
19	40	M	librarian	02138
20	42	F	homemaker	95660
21	26	M	writer	30068
22	25	M	writer	40206
23	30	F	artist	48197
24	21	F	artist	94533
25	39	M	engineer	55107

Step 5. See the last 10 entries

In [63]: `users.tail(10)`

```
Out[63]:
```

	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

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

```
In [64]: users.shape[0]
```

```
Out[64]: 943
```

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

```
In [65]: users.shape[1]
```

```
Out[65]: 4
```

Step 8. Print the name of all the columns.

```
In [67]: users.columns
```

```
Out[67]: Index(['age', 'gender', 'occupation', 'zip_code'], dtype='object')
```

Step 9. How is the dataset indexed?

```
In [68]: users.index
```

```
Out[68]: Index([ 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 [73]: users.dtypes
```

```
Out[73]: age          int64
gender          object
occupation      object
zip_code        object
dtype: object
```

Step 11. Print only the occupation column

```
In [74]: users["occupation"]
```

```
Out[74]: 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
```

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

```
In [75]: users["occupation"].nunique() #give you count of unique
```

```
Out[75]: 21
```

```
In [76]: users["occupation"].unique()
```

```
Out[76]: array(['technician', 'other', 'writer', 'executive', 'administrator',  
              'student', 'lawyer', 'educator', 'scientist', 'entertainment',  
              'programmer', 'librarian', 'homemaker', 'artist', 'engineer',  
              'marketing', 'none', 'healthcare', 'retired', 'salesman', 'doctor'],  
          dtype=object)
```

Step 13. What is the most frequent occupation?

```
In [78]: users["occupation"].value_counts().head(1)
```

```
Out[78]: occupation  
student    196  
Name: count, dtype: int64
```

```
In [83]: users["occupation"].value_counts().idxmax()
```

```
Out[83]: 'student'
```

Step 14. Summarize the DataFrame.

```
In [84]: users.describe()
```

```
Out[84]:
```

	age
count	943.000000
mean	34.051962
std	12.192740
min	7.000000
25%	25.000000
50%	31.000000
75%	43.000000
max	73.000000

Step 15. Summarize all the columns

```
In [81]: users.describe(include = "all")
```

```
Out[81]:
```

	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

Step 16. Summarize only the occupation column

```
In [82]: users["occupation"].describe()
```

```
Out[82]: count          943  
         unique         21  
         top      student  
         freq         196  
         Name: occupation, dtype: object
```

Step 17. What is the mean age of users?

```
In [85]: users['age'].mean()
```

```
Out[85]: 34.05196182396607
```

Step 18. What is the age with least occurrence?

```
In [86]: users['age'].value_counts().idxmin()
```

```
Out[86]: 7
```

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

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js