

Data Mining - Lab - 2

Numpy & Perform Data Exploration with Pandas

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 [1]: import numpy as np
```

Step 2. Create a 1D array of numbers

```
In [7]: arr=np.arange(15,25)
arr
```

```
Out[7]: array([15, 17, 19, 21, 23])
```

```
In [3]: arr =np.arange(10)
arr
```

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

```
In [10]: arr = np.array([10,25,34,55])
arr
type(arr)
```

```
Out[10]: numpy.ndarray
```

```
In [ ]:
```

Step 3. Reshape 1D to 2D Array

```
In [14]: arr = np.arange(20).reshape(4,5)
arr
```

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

```
In [15]: arr = np.array([10,20,40,60,80,70,60,50]).reshape(2,4)
arr
```

```
Out[15]: array([[10, 20, 40, 60],
 [80, 70, 60, 50]])
```

Step 4. Create a Linspace array

```
In [16]: np.linspace(14,18)
```

```
Out[16]: 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 [17]: np.linspace(14,18,25)
```

```
Out[17]: 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 [19]: np.random.rand()
```

```
Out[19]: 0.8903614358539341
```

```
In [20]: np.random.rand(8)
```

```
Out[20]: array([0.12765522, 0.80368603, 0.7605211 , 0.06785392, 0.75272059,
                0.46491016, 0.8904168 , 0.27680038])
```

```
In [21]: np.random.rand(5,6)
```

```
Out[21]: array([[0.33029975, 0.91221604, 0.76017104, 0.82116787, 0.68761145,
                0.06012767],
                [0.28742647, 0.98794472, 0.341425 , 0.12655574, 0.52085417,
                0.0014361 ],
                [0.91373875, 0.79935431, 0.58547435, 0.37281441, 0.44300592,
                0.33343973],
                [0.31825284, 0.06860888, 0.56751552, 0.66597111, 0.04407012,
                0.4128374 ],
                [0.58766781, 0.79826255, 0.41017247, 0.50405426, 0.14588308,
                0.78405554]])
```

Step 6. Create a Random Integer Array

```
In [22]: np.random.randint(20,40)
```

```
Out[22]: 37
```

```
In [27]: np.random.randint(20,40,10)
```

```
Out[27]: array([27, 22, 34, 25, 21, 21, 25, 36, 35, 36])
```

```
In [30]: np.random.randint(20,40,size=(5,5))
```

```
Out[30]: array([[20, 22, 25, 36, 37],
                [39, 37, 22, 35, 29],
                [28, 26, 36, 30, 36],
                [30, 20, 24, 29, 39],
                [22, 23, 27, 37, 28]])
```

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

```
In [37]: arr1=np.random.randint(10,50,size=10)
arr1
```

```
Out[37]: array([42, 35, 14, 10, 38, 31, 30, 43, 19, 46])
```

```
In [38]: arr1.max()
```

```
Out[38]: 46
```

```
In [39]: arr1.min()
```

```
Out[39]: 10
```

```
In [41]: arr1.argmax()
```

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

```
In [42]: arr1.argmin()
```

```
Out[42]: 3
```

Step 8. Indexing in 1D Array

```
In [45]: arr1=np.random.randint(10,50,size=10)
arr1
```

```
Out[45]: array([13, 19, 49, 17, 24, 23, 24, 43, 47, 39])
```

```
In [46]: arr1[5]
```

```
Out[46]: 23
```

```
In [47]: arr1[2:6]
```

```
Out[47]: array([49, 17, 24, 23])
```

Step 9. Indexing in 2D Array

```
In [49]: arr1=np.random.randint(10,50,size=(5,5))
arr1
```

```
Out[49]: array([[43, 30, 13, 17, 43],
               [40, 19, 24, 38, 12],
               [24, 12, 38, 24, 33],
               [48, 33, 41, 31, 37],
               [42, 10, 17, 20, 10]])
```

```
In [50]: arr1[2]
```

```
Out[50]: array([24, 12, 38, 24, 33])
```

```
In [51]: arr1[1:4]
```

```
Out[51]: array([[40, 19, 24, 38, 12],
               [24, 12, 38, 24, 33],
               [48, 33, 41, 31, 37]])
```

```
In [54]: arr1[3,3]
```

```
Out[54]: 31
```

```
In [55]: arr1[3][3]
```

```
Out[55]: 31
```

Step 10. Conditional Selection

```
In [59]: arr1=np.random.randint(20,40,size=10)
arr1
```

```
Out[59]: array([36, 22, 30, 35, 33, 26, 34, 29, 33, 20])
```

```
In [60]: arr1[arr1>25]
```

```
Out[60]: array([36, 30, 35, 33, 26, 34, 29, 33])
```

```
In [67]: arr1[(arr1>25)&(arr1<35)]
```

```
Out[67]: array([30, 33, 26, 34, 29, 33])
```

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

Pandas

Step 1. Import the necessary libraries

```
In [68]: import pandas as pd
```

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

(<https://raw.githubusercontent.com/justmarkham/DAT8/master/data/u.user>).

```
In [79]: users = pd.read_csv("https://raw.githubusercontent.com/justmarkham/DAT8/master/data/u.user", sep="|", index_col="user_id")
```

```
In [80]: users
```

Out[80]:

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

```
In [83]: users
```

Out[83]:

	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 entriesIn [84]: `users.head(25)`

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 entriesIn [85]: `users.tail(10)`

Out[85]:

	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 [90]: `users.shape[0]`

Out[90]: 943

Step 7. What is the number of columns in the dataset?In [91]: `users.shape[1]`

Out[91]: 4

Step 8. Print the name of all the columns.In [94]: `users.columns`

Out[94]: Index(['age', 'gender', 'occupation', 'zip_code'], dtype='object')

Step 9. How is the dataset indexed?In [95]: `users.index`

Out[95]: 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 [96]: users.dtypes
```

```
Out[96]: age            int64
gender          object
occupation      object
zip_code        object
dtype: object
```

Step 11. Print only the occupation column

```
In [97]: users['occupation']
```

```
Out[97]: 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 [98]: users['occupation'].nunique()
```

```
Out[98]: 21
```

```
In [99]: users['occupation'].unique()
```

```
Out[99]: 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 [102]: users['occupation'].value_counts()
```

```
Out[102]: 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
```

```
In [103]: users['occupation'].value_counts().head(1)
```

```
Out[103]: student      196
Name: occupation, dtype: int64
```

```
In [106]: users['occupation'].value_counts().idxmax()
```

```
Out[106]: 'student'
```

Step 14. Summarize the DataFrame.

In [109]: `users.describe()`

Out[109]:

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

Out[111]:

	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 [108]: `users['occupation'].describe()`

Out[108]:

```
count      943
unique      21
top      student
freq       196
Name: occupation, dtype: object
```

Step 17. What is the mean age of users?

In [113]: `users["age"].mean()`

Out[113]: 34.05196182396607

Step 18. What is the age with least occurrence?

In [114]: `users['occupation'].value_counts().tail(1)`

Out[114]:

```
doctor      7
Name: occupation, dtype: int64
```

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

In []:

