

pandas library

useful for data procssing & analysis

pandas dataframe

pandas: pandas dataframe is two dimensional tabular data structure with labeled axes(rows and columns)

```
In [1]: #importing the pandas library  
import pandas as pd
```

creating a pandas dataframe

```
In [2]: #importing the boston house price data  
from sklearn.datasets import load_boston
```

```
In [11]: boston_dataset=load_boston()
```

```
In [12]: import pandas as pd  
import numpy as np  
data_url = "http://lib.stat.cmu.edu/datasets/boston"  
raw_df = pd.read_csv(data_url, sep="\s+", skiprows=22, header=None)  
data = np.hstack([raw_df.values[::2, :], raw_df.values[1::2, :2]])  
target = raw_df.values[1::2, 2]
```

In [5]: target

```
Out[5]: array([24. , 21.6, 34.7, 33.4, 36.2, 28.7, 22.9, 27.1, 16.5, 18.9, 15. ,
18.9, 21.7, 20.4, 18.2, 19.9, 23.1, 17.5, 20.2, 18.2, 13.6, 19.6,
15.2, 14.5, 15.6, 13.9, 16.6, 14.8, 18.4, 21. , 12.7, 14.5, 13.2,
13.1, 13.5, 18.9, 20. , 21. , 24.7, 30.8, 34.9, 26.6, 25.3, 24.7,
21.2, 19.3, 20. , 16.6, 14.4, 19.4, 19.7, 20.5, 25. , 23.4, 18.9,
35.4, 24.7, 31.6, 23.3, 19.6, 18.7, 16. , 22.2, 25. , 33. , 23.5,
19.4, 22. , 17.4, 20.9, 24.2, 21.7, 22.8, 23.4, 24.1, 21.4, 20. ,
20.8, 21.2, 20.3, 28. , 23.9, 24.8, 22.9, 23.9, 26.6, 22.5, 22.2,
23.6, 28.7, 22.6, 22. , 22.9, 25. , 20.6, 28.4, 21.4, 38.7, 43.8,
33.2, 27.5, 26.5, 18.6, 19.3, 20.1, 19.5, 19.5, 20.4, 19.8, 19.4,
21.7, 22.8, 18.8, 18.7, 18.5, 18.3, 21.2, 19.2, 20.4, 19.3, 22. ,
20.3, 20.5, 17.3, 18.8, 21.4, 15.7, 16.2, 18. , 14.3, 19.2, 19.6,
23. , 18.4, 15.6, 18.1, 17.4, 17.1, 13.3, 17.8, 14. , 14.4, 13.4,
15.6, 11.8, 13.8, 15.6, 14.6, 17.8, 15.4, 21.5, 19.6, 15.3, 19.4,
17. , 15.6, 13.1, 41.3, 24.3, 23.3, 27. , 50. , 50. , 50. , 22.7,
25. , 50. , 23.8, 23.8, 22.3, 17.4, 19.1, 23.1, 23.6, 22.6, 29.4,
23.2, 24.6, 29.9, 37.2, 39.8, 36.2, 37.9, 32.5, 26.4, 29.6, 50. ,
32. , 29.8, 34.9, 37. , 30.5, 36.4, 31.1, 29.1, 50. , 33.3, 30.3,
34.6, 34.9, 32.9, 24.1, 42.3, 48.5, 50. , 22.6, 24.4, 22.5, 24.4,
20. , 21.7, 19.3, 22.4, 28.1, 23.7, 25. , 23.3, 28.7, 21.5, 23. ,
26.7, 21.7, 27.5, 30.1, 44.8, 50. , 37.6, 31.6, 46.7, 31.5, 24.3,
31.7, 41.7, 48.3, 29. , 24. , 25.1, 31.5, 23.7, 23.3, 22. , 20.1,
22.2, 23.7, 17.6, 18.5, 24.3, 20.5, 24.5, 26.2, 24.4, 24.8, 29.6,
42.8, 21.9, 20.9, 44. , 50. , 36. , 30.1, 33.8, 43.1, 48.8, 31. ,
36.5, 22.8, 30.7, 50. , 43.5, 20.7, 21.1, 25.2, 24.4, 35.2, 32.4,
32. , 33.2, 33.1, 29.1, 35.1, 45.4, 35.4, 46. , 50. , 32.2, 22. ,
20.1, 23.2, 22.3, 24.8, 28.5, 37.3, 27.9, 23.9, 21.7, 28.6, 27.1,
20.3, 22.5, 29. , 24.8, 22. , 26.4, 33.1, 36.1, 28.4, 33.4, 28.2,
22.8, 20.3, 16.1, 22.1, 19.4, 21.6, 23.8, 16.2, 17.8, 19.8, 23.1,
21. , 23.8, 23.1, 20.4, 18.5, 25. , 24.6, 23. , 22.2, 19.3, 22.6,
19.8, 17.1, 19.4, 22.2, 20.7, 21.1, 19.5, 18.5, 20.6, 19. , 18.7,
32.7, 16.5, 23.9, 31.2, 17.5, 17.2, 23.1, 24.5, 26.6, 22.9, 24.1,
18.6, 30.1, 18.2, 20.6, 17.8, 21.7, 22.7, 22.6, 25. , 19.9, 20.8,
16.8, 21.9, 27.5, 21.9, 23.1, 50. , 50. , 50. , 50. , 50. , 13.8,
13.8, 15. , 13.9, 13.3, 13.1, 10.2, 10.4, 10.9, 11.3, 12.3, 8.8,
7.2, 10.5, 7.4, 10.2, 11.5, 15.1, 23.2, 9.7, 13.8, 12.7, 13.1,
12.5, 8.5, 5. , 6.3, 5.6, 7.2, 12.1, 8.3, 8.5, 5. , 11.9,
27.9, 17.2, 27.5, 15. , 17.2, 17.9, 16.3, 7. , 7.2, 7.5, 10.4,
8.8, 8.4, 16.7, 14.2, 20.8, 13.4, 11.7, 8.3, 10.2, 10.9, 11. ,
9.5, 14.5, 14.1, 16.1, 14.3, 11.7, 13.4, 9.6, 8.7, 8.4, 12.8,
10.5, 17.1, 18.4, 15.4, 10.8, 11.8, 14.9, 12.6, 14.1, 13. , 13.4,
15.2, 16.1, 17.8, 14.9, 14.1, 12.7, 13.5, 14.9, 20. , 16.4, 17.7,
19.5, 20.2, 21.4, 19.9, 19. , 19.1, 19.1, 20.1, 19.9, 19.6, 23.2,
29.8, 13.8, 13.3, 16.7, 12. , 14.6, 21.4, 23. , 23.7, 25. , 21.8,
20.6, 21.2, 19.1, 20.6, 15.2, 7. , 8.1, 13.6, 20.1, 21.8, 24.5,
23.1, 19.7, 18.3, 21.2, 17.5, 16.8, 22.4, 20.6, 23.9, 22. , 11.9])
```

In [6]: data_url

```
Out[6]: 'http://lib.stat.cmu.edu/datasets/boston'
```

In [7]: data

```
Out[7]: array([[6.3200e-03, 1.8000e+01, 2.3100e+00, ..., 1.5300e+01, 3.9690e+02,
               4.9800e+00],
               [2.7310e-02, 0.0000e+00, 7.0700e+00, ..., 1.7800e+01, 3.9690e+02,
               9.1400e+00],
               [2.7290e-02, 0.0000e+00, 7.0700e+00, ..., 1.7800e+01, 3.9283e+02,
               4.0300e+00],
               ...,
               [6.0760e-02, 0.0000e+00, 1.1930e+01, ..., 2.1000e+01, 3.9690e+02,
               5.6400e+00],
               [1.0959e-01, 0.0000e+00, 1.1930e+01, ..., 2.1000e+01, 3.9345e+02,
               6.4800e+00],
               [4.7410e-02, 0.0000e+00, 1.1930e+01, ..., 2.1000e+01, 3.9690e+02,
               7.8800e+00]])
```

In [13]: *#pandas dataframe*

```
target=pd.DataFrame(data,columns=boston_dataset.feature_names)
```

In [14]: target

Out[14]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90
...
501	0.06263	0.0	11.93	0.0	0.573	6.593	69.1	2.4786	1.0	273.0	21.0	391.99
502	0.04527	0.0	11.93	0.0	0.573	6.120	76.7	2.2875	1.0	273.0	21.0	396.90
503	0.06076	0.0	11.93	0.0	0.573	6.976	91.0	2.1675	1.0	273.0	21.0	396.90
504	0.10959	0.0	11.93	0.0	0.573	6.794	89.3	2.3889	1.0	273.0	21.0	393.45
505	0.04741	0.0	11.93	0.0	0.573	6.030	80.8	2.5050	1.0	273.0	21.0	396.90

506 rows × 13 columns



In [16]: target.shape

Out[16]: (506, 13)

importing the data from a csv file to a pandas dataframe

```
In [17]: diabetes_df=pd.read_csv("diabetes.csv")
diabetes_df
```

Out[17]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFun
0	6	148	72	35	0	33.6	
1	1	85	66	29	0	26.6	
2	8	183	64	0	0	23.3	
3	1	89	66	23	94	28.1	
4	0	137	40	35	168	43.1	
...	
763	10	101	76	48	180	32.9	
764	2	122	70	27	0	36.8	
765	5	121	72	23	112	26.2	
766	1	126	60	0	0	30.1	
767	1	93	70	31	0	30.4	

768 rows × 9 columns



```
In [21]: diabetes_df.head()
```

Out[21]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunc
0	6	148	72	35	0	33.6	0.6
1	1	85	66	29	0	26.6	0.3
2	8	183	64	0	0	23.3	0.6
3	1	89	66	23	94	28.1	0.1
4	0	137	40	35	168	43.1	2.2



```
In [22]: diabetes_df.shape
```

Out[22]: (768, 9)

loading the data from a excel file to a pandas dataframe:

```
pa.read_excel("file path")
```

exporting a dataframe to a csv file

```
In [29]: target.to_csv("http://lib.stat.cmu.edu/datasets/boston")
```

In [30]: target

Out[30]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90
...
501	0.06263	0.0	11.93	0.0	0.573	6.593	69.1	2.4786	1.0	273.0	21.0	391.99
502	0.04527	0.0	11.93	0.0	0.573	6.120	76.7	2.2875	1.0	273.0	21.0	396.90
503	0.06076	0.0	11.93	0.0	0.573	6.976	91.0	2.1675	1.0	273.0	21.0	396.90
504	0.10959	0.0	11.93	0.0	0.573	6.794	89.3	2.3889	1.0	273.0	21.0	393.45
505	0.04741	0.0	11.93	0.0	0.573	6.030	80.8	2.5050	1.0	273.0	21.0	396.90

506 rows × 13 columns



In [31]: *#creating a dataframe with random values*
 random_df=pd.DataFrame(np.random.rand(20,10))

In [32]: random_df.head()

Out[32]:

	0	1	2	3	4	5	6	7	8
0	0.686625	0.636520	0.093915	0.488454	0.735997	0.416543	0.038713	0.793722	0.728184
1	0.105768	0.053672	0.150409	0.234073	0.210695	0.637186	0.356588	0.802642	0.064209
2	0.594303	0.284989	0.061488	0.024186	0.259927	0.111923	0.965723	0.057942	0.122898
3	0.514427	0.416036	0.104626	0.198550	0.035863	0.736287	0.933567	0.607442	0.575330
4	0.858736	0.635930	0.054070	0.794145	0.788379	0.441367	0.229441	0.950394	0.955095



In [33]: target.info

```
Out[33]: <bound method DataFrame.info of
RM      AGE      DIS  RAD      TAX  \
0      0.00632  18.0   2.31   0.0   0.538  6.575  65.2  4.0900  1.0  296.0
1      0.02731   0.0   7.07   0.0   0.469  6.421  78.9  4.9671  2.0  242.0
2      0.02729   0.0   7.07   0.0   0.469  7.185  61.1  4.9671  2.0  242.0
3      0.03237   0.0   2.18   0.0   0.458  6.998  45.8  6.0622  3.0  222.0
4      0.06905   0.0   2.18   0.0   0.458  7.147  54.2  6.0622  3.0  222.0
..      ...      ...      ...      ...      ...      ...      ...      ...      ...
501     0.06263   0.0  11.93   0.0   0.573  6.593  69.1  2.4786  1.0  273.0
502     0.04527   0.0  11.93   0.0   0.573  6.120  76.7  2.2875  1.0  273.0
503     0.06076   0.0  11.93   0.0   0.573  6.976  91.0  2.1675  1.0  273.0
504     0.10959   0.0  11.93   0.0   0.573  6.794  89.3  2.3889  1.0  273.0
505     0.04741   0.0  11.93   0.0   0.573  6.030  80.8  2.5050  1.0  273.0

      PTRATIO      B  LSTAT
0          15.3  396.90   4.98
1          17.8  396.90   9.14
2          17.8  392.83   4.03
3          18.7  394.63   2.94
4          18.7  396.90   5.33
..      ...      ...      ...
501         21.0  391.99   9.67
502         21.0  396.90   9.08
503         21.0  396.90   5.64
504         21.0  393.45   6.48
505         21.0  396.90   7.88

[506 rows x 13 columns]>
```

In [38]: target.isnull().sum()

```
Out[38]: CRIM      0
ZN          0
INDUS       0
CHAS        0
NOX         0
RM          0
AGE         0
DIS         0
RAD         0
TAX         0
PTRATIO     0
B           0
LSTAT       0
dtype: int64
```

In [40]: *#computing the values based on the labels*
diabetes_df.value_counts("Outcome")

```
Out[40]: Outcome
0      500
1      268
dtype: int64
```

```
In [42]: #group the value based on the mean
diabetes_df.groupby("Outcome").mean()
```

```
Out[42]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Diabetes
Outcome							
0	3.298000	109.980000	68.184000	19.664000	68.792000	30.304200	
1	4.865672	141.257463	70.824627	22.164179	100.335821	35.142537	

statistical measures

```
In [43]: #count of number of values
diabetes_df.count()
```

```
Out[43]: Pregnancies      768
Glucose      768
BloodPressure 768
SkinThickness 768
Insulin      768
BMI          768
DiabetesPedigreeFunction 768
Age          768
Outcome      768
dtype: int64
```

```
In [44]: #mean value-column wise
diabetes_df.mean()
```

```
Out[44]: Pregnancies      3.845052
Glucose      120.894531
BloodPressure 69.105469
SkinThickness 20.536458
Insulin      79.799479
BMI          31.992578
DiabetesPedigreeFunction 0.471876
Age          33.240885
Outcome      0.348958
dtype: float64
```

```
In [45]: #standard deviation-- column wise
diabetes_df.std()
```

```
Out[45]: Pregnancies      3.369578
Glucose      31.972618
BloodPressure 19.355807
SkinThickness 15.952218
Insulin      115.244002
BMI          7.884160
DiabetesPedigreeFunction 0.331329
Age          11.760232
Outcome      0.476951
dtype: float64
```

```
In [46]: #minimum_value  
diabetes_df.mean()
```

```
Out[46]: Pregnancies      3.845052  
Glucose      120.894531  
BloodPressure  69.105469  
SkinThickness  20.536458  
Insulin      79.799479  
BMI          31.992578  
DiabetesPedigreeFunction  0.471876  
Age          33.240885  
Outcome      0.348958  
dtype: float64
```

```
In [47]: #maximum value  
diabetes_df.max()
```

```
Out[47]: Pregnancies      17.00  
Glucose      199.00  
BloodPressure  122.00  
SkinThickness   99.00  
Insulin      846.00  
BMI           67.10  
DiabetesPedigreeFunction  2.42  
Age           81.00  
Outcome        1.00  
dtype: float64
```

```
In [48]: #all the statistical measure about the dataframe  
diabetes_df.describe()
```

```
Out[48]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Diabi
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	