

Pandas

Date:
Page:

Usage:

1. `dict1 = { "name": ["Ra", "Sha", "Ga", "Na"], "marks": [10, 20, 30, 40], "City": ["Rampur", "Cizer", "Norway", "Hetzuda"] }`

- we created a dictionary. Jupyter Notebook adds index itself.

2. `df = pd.DataFrame(dict1)`

`df`

- Turns into excel sheet

	Name	marks	city
0	Ra	10	Rampur
1	sha	20	Cizer
2	Ga	30	Norway
3	Na	40	Hetzuda

- `pd.DataFrame`: pd is typically the alias used for importing pandas library. DataFrame class is one of the fundamental data structures. used to store & manipulate 2-D labeled data.

- You can create DataFrames from various data structures like lists, dictionaries, or Numpy arrays.

3. `df.to_csv('std-record')`

It creates CSV file named 'std-record'

4. `df.to_csv('std-record', index=False)`

If you don't want index.

5. `df = pd.read_csv('...')`

6. `df.head()` || first five

7. `df.tail()` || last five

8. `df.describe()`

- It computes the numerical data from the data frame.

	marks
Count	4
mean	25
std	
min	10 min score
25%	25% quartile
50%	50% quartile
75%	75% quartile
max	40 max score

9. `df['marks']` \Rightarrow [0, 10]

- Selecting a whole column

0	10
1	20
2	30
3	40

10. `df['marks'][0]` \Rightarrow 10

- Select data from marks column index 0

11. `df['marks'][1] = 30`

Edit value of index 1 to 30.

11. If you have to change index:

```
df.index = ['first', 'Second', 'third', 'fourth']
```

Pandas has 2 types of data structures.

- (i) Series - one dimensional array with indexes, stores single column or row of data in a DataFrame.
- (ii) DataFrame - It is a tabular spreadsheet like structure representing rows each of which contains one or more columns.

Series - capable of holding any type of data.

Dataframe - potentially different data types

12. Ser = pd.DataFrame Series(np.random.rand(5))

0 0.12359

• type(Ser)

1 0.23503

↳ pandas.core.series.Series

2 0.21571

3 0.53124

• Ser.describe()

4 0.72596

↳ gives count, mean, std, min, max

13. new-df = pd.DataFrame(np.random.rand(334,5),
index = np.arange(334))

type(new-df) ⇒ pandas.core.frame.DataFrame

new_df.describe()

14. new_df.head()

	0	1	2	3	4
0	0.253	0.135	0.786	0.882	0.439
1	0.123	0.953	0.329	0.281	0.605
2	0.301	0.182	0.300	0.791	0.147
3	0.312	0.001	0.132	0.879	0.010
4	0.111	0.117	0.471	0.924	0.495

15. new_df.index

[0, 1, 2, ..., 300, ..., 334]

16. new_df.columns

RangeIndex (start=0, stop=5, step=1)

To convert floating point into numpy array.

17. new_df.to_numpy()

array([0.253, 0.135, 0.786, 0.882, 0.439,

[

334

])

18. `new_df.T` → rows & columns interchange

0 1 2 3 4 5 .. 334 columns

1

2

3

4

Rows

19. `new_df.sort_index(axis=0, ascending=False)`

Sorting according to index.

`axis=0`; it represents rows

by default; `ascending=True`.

0 1 2 3 4

334 ..

333 ..

:

:

20. If `axis=1`, it sorts columns.

4 3 2 1 0

0

1

2

21. `type(new_df[0])` `pandas.core.series.Series`

comb'n of series forms or data frame.

>> new_df2 = new_df

>> new_df2[0][0] = 0.8783

>> new_df - it gets updated, but we had edited new_df2

so, new_df2 is a view of new_df. Like a pointer, when you change view, original data gets updated.

>> new_df3 = new_df.copy()

>> new_df3[0][0] = 0.329
colm index

>> new_df

newdf remains the same.

But, we should not use chained indexing like df[0][0] = 2 in pandas because pandas might return a copy of the data instead of view of the data.

To ensure, modifications are made to the original Dataframe and avoid potential issues, it is recommended to use the '`.loc`' or '`.iloc`' accessors for setting values.

axis = 0 \rightarrow row

axis = 1 \rightarrow column

Page:

using loc / iloc instead of Chained indexing

[df.loc[0,0] = 2]

[df.iloc[0,0] = 2]

new_df.loc[0,0] = 654
↑ ↑
row column

To drop a column:

new_df = new_df.drop(0, axis=1) \rightarrow columnⁿ name

new_df.loc[[1,2], :]

row 1 & row 2 and all columnⁿ

» new_df.iloc[0,4] || '0' row & 4th column value.

» new_df.drop(['A', 'C'], axis=1) || drop whole column of A & C.

» new_df.drop(['B'], axis=1, inplace=True) || change in original form