

25th July

Python- tutorial 13

- from array to dataframe :

from numpy.random import randn as rn
np.random.seed(101) (for random data)
random but fixed every time we execute
matrix_data = rn(5,4) random array 5x4

matrix_data

array([[0.2, -0.4, 0.6, 1.0]
[2, 9, 8, 4]
[11, 2, 6, 4]
[2, 9, 5, 2]
[1, 2, 9, 5]])

pd.DataFrame(matrix_data)

	0	1	2	3
0	0.2	-0.4	0.6	1
1	2	9	8	4
2	11	2	6	4
3	2	9	5	2
4	1	2	9	5

(default heading & index)

if we wish to change heading, indexes :-

row-labels = ['A', 'B', 'C', 'D', 'E']

column_headings = ['w', 'x', 'y', 'z']

df = pd.DataFrame(matrix_data,
row-labels, column_headings)

df → same above Dataframe with changed head & index

$\text{iloc} \rightarrow$ default index location.
 $\text{loc} \rightarrow$ named by us indexes.
 $\text{ix} \rightarrow$ combination of loc & iloc

• extracting some section of dataframe:

① Some columns:

$\text{df}['x']$, $\text{df}[['x', 'z']]$
 ↑ ↑
 Single column multiple columns

② Some rows:

Same $\text{df.loc}['00']$ & type = series (single row or column)
 $\text{df.iloc}[2] \leftarrow$ by using default index

if I change my index with integer data
 ex, $['A', 'B', 'C', 'D', 'D'] \rightarrow [56, 23, 42, 8, 80]$

$\text{df.loc}[23] \rightarrow$ 2nd row
 $\text{df.iloc}[23] \rightarrow$ error (no data found)

Note: $\text{iloc} \rightarrow$ understands only default integer location & not the integer indexes given by us.

$\text{df.loc}[8] \rightarrow$ 4th & 5th rows

$\text{df.iloc}[3] \rightarrow$ 4th row only.

③ intersection of some rows / columns:

$\text{df.iloc}[[2, 3], [0, 1]]$

		W	X
42	C	11	9
8	D	9	2
8	D	1	2

$\text{df.loc}[[42, 8], ['W', 'X']]$

③ we can drop row/column from df, not a
can't drop data (yes we can replace)

④ using range function extracting section:
`df.loc[56:42, ['w': 'x']]`

• Dropping some rows & columns:

① drop column:

`df.drop('X')`
error (∵ by default
axis=0)

axis=0 → x-axis
(row)

axis=1 → y-axis
(column)

(by-default axis=0)

`df.drop('X', axis=1)` (∵ doesn't drop
permanently)

same previous dataframe with 'X' series

• if you want to make changes permanently (∵ `inplace=True`
by default, `inplace=False`)

`df.drop('X', axis=1, inplace=True)`
drops 'X' column permanently.

② drop row:

`df.drop(56)` named indexes not
able to understand
default indexes.

row-1 will be dropped

`df = df.drop(56)`

another way of permanent deletion

• filter data based on some conditions:

`print(df.loc[['A', 'B', 'C']] > 0)`

	W	X	Y	Z
A	True	False	True	True
B	True	True	True	True
C	True	True	True	True

if you want values, then put condition in ^{df}
`df[df.loc[['A', 'B', 'C']] > 0]`

	W	X	Y	Z
A	0.2	NaN	0.6	1
B	2	9	8	4
C	11	2	6	4
D	NaN	NaN	NaN	NaN
D	NaN	NaN	NaN	NaN

- multiple conditions ("and") :

`df[(df['Age'] > 30) & (df['Height'] > 65) & (df['Weight'] < 70)]`

- Resetting index to default index:

no
permanent
changes

`df.reset_index()`

default index + old index

`df.reset_index(drop=True)`

only default index

`df.reset_index(drop=True, inplace=True)`

only ~~df~~ permanently changed to default index.

[zip] → creates list of zip based on indexes.

• "My name is MADHU".split()
→ ['My', 'name', 'is', 'MADHU'] a list

adding new column inside df:

df["xyz"] = "My name is MADHU".split()
df

→ data frame with new 'xyz' column.

use this column as index-column:

df.set_index('xyz')

• Multiple index (hierarchy)

outside = ['G1', 'G1', 'G1', 'G2', 'G2', 'G2']

inside = [1, 2, 3, 1, 2, 3]

hier_index = list(zip(outside, inside))

hier_index = pd.MultiIndex.from_tuples(hier_index)

df1 = pd.DataFrame(data=np.random
(6, 3), index=hier_index,
columns=['A', 'B', 'C'])

print(df1)

		A	B	C
G1	1	1.0	1.0	-1.0
	2	-1.0	-0.0	2.0
	3	1.0	0.0	-1.0
G2	1	1.0	-0.0	-1.0
	2	-2.0	-2.0	1.0
	3	-1.0	2.0	-2.0

now check whether 'G1' is working as a row:

df1.loc['G1']

	A	B	C
1	1.0	1.0	-1.0
2	-1.0	-0.0	2.0
3	1.0	0.0	-1.0

df1.loc['G1'].loc[3, ['B', 'C']]

B	0.0
C	-1.0

(we may also use combination of loc & iloc)

- Drop row/column having NaN-values:

```
df = pd.DataFrame({ 'A': [1, 2, np.nan],
                    'B': [5, np.nan, np.nan], 'C': [1, 2, 3] })
df['States'] = 'MA DH UP'.split()
df.set_index('States', inplace=True)
```

df

	A	B	C
States	.	.	.
MA	1.0	5.0	1
DH	2.0	NaN	2
UP	NaN	NaN	3

print(df.dropna()) ← dropped rows with NaN because by default axis=0 (rows)

	A	B	C
States	.	.	.
MA	1.0	5.0	1

print(df.dropna(axis=1))

States	C
MA	1
DH	2
UP	3

- with conditions dropping: (read carefully) threshold=2

df.dropna(thresh=2)

	A	B	C
states			
MA	1.0	5.0	1
DH	2.0	NaN	2

Note: 'thresh' is applied on non-NaN.

Note: (for permanent changes use,)
inplace=True