

PANDAS

## use case of Pandas

- ① Data cleaning
- ② Data Analysis
- ③ Data transformation
- ④ Data visualization (Basic level)
- ⑤ Data Aggregation
- ⑥ File handling
- ⑦ Data filtering & Selection
- ⑧ Time series analysis.

[ 'S' , 'M' , 'L' ] =

[ 0.8 , 0.6 , 0.1 ] =

( [ 0.8 , 0.6 , 0.1 ] ) =

{ 'loc': 'S' , 'occ': 'd' , 'ctr': 'S' }

loc was fed from file name 11- (dict1) etc

## (1) Series

- ① Stacking a data into vertical way.
- ② 1D labeled array capable of holding any data type.
- ③ Axis labels are collectively called the index.

→ Creating Series :-

import numpy as np

import pandas as pd

labels = ['a', 'b', 'c']

my\_lists = [10, 20, 30]

arr = np.array([10, 20, 30])

d = {'a': 10, 'b': 20, 'c': 30}

Pd.Series (labels) # Same with my-list, arr & d.

→ 0 a

1 b

2 c

## ② Data frames

→ ① Data frame using Dict

data = {

'Name': ['John', 'Anna', 'Peter', 'Linda'],

'Age': [28, 34, 29, 42],

'City': ['Newyork', 'Paris', 'Berlin', 'London'],

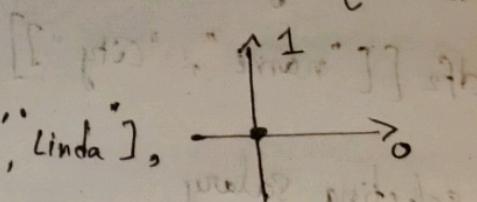
'Salary': [65000, 70000, 62000, 85000]

}

df = pd.DataFrame(data)

axis = 1 [col]

axis = 0 [row]



→ ② Data frame using list :-

data\_list = [

['John', 28, 'Newyork', 65000],

['Anna', 34, 'Paris', 70000],

['Peter', 29, 'Berlin', 62000],

['Linda', 42, 'London', 85000]

]

df2 = pd.DataFrame(data\_list)

→ giving col names to list :-

cols = ["name", "age", "city", "salary"]

df2 = pd.DataFrame(data\_list, columns=cols)

df2

## ② Selecting & indexing of cols

# Selecting Name & city

```
df2[['name', 'city']]
```

# selecting salary.

```
df2['salary']
```

# creating a new col

```
df2['Designation'] = ['Doctor', 'Scientist', 'Lawyer',  
                      'Designer']
```

# Removing col :-

```
df2.drop(['designation'], axis=1, inplace=True)
```

# Removing rows :-

```
df2.drop(3, axis=0)
```

# Selecting row using loc function

```
df2.loc[0, 1]
```

```
df2.iloc[2, 3]
```

M < > ?

Enter

Shift

# Selecting subsets of Rows & cols

① df2.loc[[0, 1]][["city", "salary"]]

② df2.loc[[2, 3]][["name", "age"]]

### ③ conditional Selection

# Selecting people whose age is greater than 30

df2[df2["age"] > 30]

# Selecting people with age > 30 & live in paris

df2[(df2["age"] > 30) & (df2["city"] == "Paris")]

## Finding Missing Data

import numpy as np

import pandas as pd

data = {

'A': [1, 2, np.nan, 4, 5],

'B': [1, 2, 3, 4, 5],

'C': [1, 2, 3, np.nan, np.nan],

'D': [1, np.nan, np.nan, np.nan, 5]

}

df = pd.DataFrame(data)

① df.isna() # returns Bool. True means NaN

② df.isna().sum()

③ df.isna().any() # whether the column contains  
at least one ~~is~~ missing value.

N M &lt; &gt; ?

(2) Removing a missing data

① df.dropna()

② df.dropna(thresh=1) # It will keep all rows in a data frame that have atleast one non-missing value.

(3) Filling the missing data

① df.fillna(0)

② values = {"A": 0,  
'B': 100,  
'C': 300,  
'D': 400}

{}

df.fillna(value=values)

③ df.fillna(df.mean())

## Merging joining & concatenation ..

```
employees = pd.DataFrame({  
    'employee_id': [1, 2, 3, 4, 5],  
    'name': ['John', 'Anna', 'Peter', 'Linda', 'Bob'],  
    'department': ['HR', 'IT', 'Finance', 'IT', 'HR']  
})
```

```
Salaries = pd.DataFrame({  
    'employee_id': [1, 2, 3, 6, 7],  
    'Salary': [60000, 80000, 65000, 70000, 90000],  
    'bonus': [5000, 10000, 7000, 8000, 12000]  
})
```

### ① Merging two dataframes :-

```
pd.merge(employees, Salaries, on = "employee_id",  
        how = 'inner')
```

# outer

# left

# right

B

N

M &lt; &gt; ?

② Concatenation of two Dataframes:  
[1]:

① pd.concat([df1, df2]) # col.

② pd.concat([df1, df2], axis=1) # row

Bob]

HR]

③ Joining two Dataframes:

① df1.join(df2, how='outer')

② df2.join(df1)

[0.02, 0.03, 0.81, 0.91, 0.82, 0.21, 0.05, 0.01] : [2]

{2.0, 2.1, 0.6, 3.81, 3.91, 2.1, 0.1} : [3]

(3-rows, 7cols) (10-10-2cols) [0.02, 0.03, 0.81, 0.91, 0.82, 0.21, 0.05, 0.01] : [2]

(7cols) (inner join) [7]

id,

(inner, [cols2]) ('propator') pd merge A

two propator pd merge & does slight diff

(inner, [cols2]) ('cols', 'propator')) pd merge

## Group by & aggregation

```
import numpy as np
```

```
import pandas as pd
```

```
data = {
```

```
    'category': ['A', 'B', 'A', 'B', 'A', 'B', 'A', 'B']
```

```
    'store': ['S1', 'S1', 'S2', 'S2', 'S1', 'S2', 'S2', 'S1']
```

```
    'sales': [100, 200, 150, 280, 120, 180, 200, 300],
```

```
    'quantity': [10, 15, 12, 18, 8, 20, 15, 25],
```

```
    'Date': pd.date_range('2023-01-01', periods=8)
```

```
}  
df = pd.DataFrame(data)
```

---

① Group by Category & calculate the sum of Sales

```
cat = df.groupby(['category'])['Sales'].sum()
```

```
cat.
```

② Group by multiple cols & Groupby Category and Sales: Sales:

```
cat = df.groupby(['category', 'Sales'])['Sales'].sum()
```

J K L : ; Enter  
N M < > ,  
→ df['sales'].mean()  
→ df['sales'].agg(['sum', 'mean', 'min', 'max',  
'count', 'std', 'median'])

→ X

Source code on  
Github

8)

8.

and