

PandasData Structures in pandas :

- ↳ Series (similar to 1 dimensional numpy array)
- ↳ Data Frame (similar to 2 dim. numpy array)

Series

$S = \text{pd.Series}(\text{data}, \text{index} = \text{idx})$

list → $\text{data} = [1, 2, 3, 4]$

dictionary → $\text{data} = \{'a': 1, 'b': 2\}$

numpy array → $\text{data} = \text{np.array}([1, 2, 3])$

scalars → $\text{data} = 5, \text{index} = ['a', 'b', 'c']$

properties / Attributes

- data types : $S.dtype$
- shape : $S.shape$
- values : $S.values$
- Array : $S.array$

Methods

- Head : $S.head()$
- Tail : $S.tail()$
- Info : $S.info()$
- To numpy : $S.to_numpy()$

IndexingSingle Index

$S[2]$ ← using Index value

$S['a']$ ← using Row label

$S[\text{start} : \text{stop} : \text{step}]$

slicingMultiple Indexes

$S[[1, 4]]$ @ $S[['a', 'c']]$

list of indexes

Data Frame

`df = pd.DataFrame (data, index = idxs, columns = cols)`

create data frame :

using
dictionary

`{data = {'Name' : ['Ravi', 'Shaheer'], 'Age' : [26, 23]}`

using
tuple

`{data = (('Ravi', 26), ('Shaheer', 23)),
columns = ['Name', 'age']}`

using
list

`{data = (['Ravi', 26], ['Shaheer', 23]), columns = ['Name', 'Age']}`

using
numpy array

`{data = np.array([1, 2, 3], [4, 5, 6]), columns = [col-1, col-2]}`

Attributes / properties

- shape of df : `df.shape`
- columns : `df.columns`
- data types : `df.dtypes`
- Axes : `df.axes`
- values : `df.values`

Methods

- Head : `df.head()`
- Tail : `df.tail()`
- Info : `df.info()`

- Load data to pandas :

`- pd.read_* (* = csv, excel, sql, json, etc.....)`

Eg: `pd.read_csv('path')`

- Export / save from pandas :

`- df.to_* (* = csv, excel, sql, json, etc.....)`

Ex: `df.to_excel('path', sheet_name = 'name', index = False)`

Filtering Rows :-

Way-1: `df[starting-row-index : ending-row-index : step]`

Way-2: `df[condition]`

`df[df['date'].isin(['10-5-2016', '10-4-2016'])]`

Filtering Specific Rows & columns from df :

`df.loc[row-label, column-label]`

`df.iloc[row-index, column-index]`

Accessing Rows based on condition :

`df.loc[condition, column-labels]`

Accessing Rows based on multiple condition :

`df.loc[(cond-1) & (cond-2) | (cond-3), column-labels]`

Renaming columns :

`df.rename(index=None, columns={'old-name': 'new-name'})`

Modify columns datatype :

• `astype()`

`df['quantity'] = df['quantity'].astype(int)`

`df = df.astype({'quantity': int})`

• `apply()`

`df['quantity'] = df['quantity'].apply(pd.to_numeric)`
`pd.to_datetime`
`pd.to_timedelta`

- creating a Derived column :

`df['amount'] = df['quantity'] * df['unit-price']`

- creating columns using apply() function :

`df.apply (function, axis =)`

`s.apply (function, axis =)`

`df['amount'] = df.apply (lambda row : row['quantity'] *
row['unit-price'], axis=1)`

date time

properties

- day `df['date'].dt.day`
- month `.dt.month`
- years `.dt.years`
- day-of-week `.dt.day_of_week`
- day-of-years `.dt.day_of_years`
- is-leap-years `.dt.is-leap-years`
- week `.dt.week`

methods

- day_name `.dt.day_name()`
- month_name `.dt.month_name()`

time delta

properties

- days `.dt.days`
- hours `.dt.hours`
- minutes `.dt.minutes`
- seconds `.dt.seconds`
- components `.dt.components`

methods

- total seconds `.dt.total_seconds()`

pd.to_datetime()

```
df['date'] = pd.to_datetime(df['date'], dayfirst=True)  
format = '%d/%m/%Y'
```

```
df = pd.read_csv('path', parse_dates=['cols'])
```

sort_index, sort_values, Reset-index :

```
df.sort_index(ascending=True)
```

```
df.sort_values(by='col-name')
```

```
df.reset_index()
```

split :

```
grouped_df = df.groupby('category')
```

```
grouped_df = df.groupby(['category', 'sub-category'])
```

Groups : grouped_df.groups

```
grouped_df.groups.keys()
```

filter : grouped_df.get_group('Group-key-name')

```
grouped_df.get_group(('Group-category', 'Group-sub-cat'))
```

first, last & nth row of each group :

```
grouped_df.first()
```

```
grouped_df.last()
```

```
grouped_df.nth()
```

Apply :

Aggregation :

```
cov()
```

```
idxmax()
```

```
grouped_df['sales'].agg(['min', 'max'])
```

} Built-in

} user defined

Filteration:

head ()
tail ()
nth ()

} Built-in

uses defined. { grouped-df.filter (lambda group: group['sales'].mean > 20)

Transformation:

bfill ()
cumsum ()

} Built-in

grouped-df.transform (lambda x: x+1) } uses defined.

Pivot-table:

df.pivot-table (values = 'sales', index = ['region'],
aggfunc = ['sum'], columns = ['category'])

pd.merge (, , ' ') :

pd.merge (left, right, how = 'inner', on = None, left_on = None,
right_on = None, left_index = None, right_index = ,
sort =)

how = 'inner'
= 'outer'
= 'left'
= 'right'
= 'cross'

- plotting :

```
df.plot(kind='kind')
```

```
kind = 'bar'
```

```
= 'barh'
```

```
= 'hist'
```

```
= 'kde'
```

- Delete columns :

Syntax-1

```
df.drop(['col1', 'col3'], axis=1, inplace=True)
```

Syntax-2

```
df.drop(df.loc[:, 'col2':'col4'], inplace=True)
```

Syntax-3

```
df.drop(df.columns[[0, 4, 2]], axis=1, inplace=True)
```

Syntax-4

```
df.drop(df.iloc[:, 1:3], inplace=True)
```

Syntax-5

```
df.pop('col4', inplace=True)
```

- Inserting Rows :

```
df = pd.concat([df, new-df], ignore_index=True, axis=0)
```

- Inserting a row using df.loc[]

```
df.loc[len(df)] = ['1/12/2017', 28, 2, 'Rain']
```

- Inserting a row at specific index

```
df.loc[8.5] = ['1/11/2017', 30, 3, 'Rain']
```

↘ Sort & Reset index

```
df = df.sort_index().reset_index(drop=True)
```