

## \*Pandas:

- used for working with data sets
- analyzing, cleaning, exploring and manipulating data
- Pandas - Panel data
- Python data analysis

import Pandas as Pd

```
mydata = { 'cars' : [ 'BMW', 'Volvo', 'Ford' ]  
          'Passings' : [ 3, 7, 2 ]  
        }
```

```
myVar = Pd.DataFrame(mydataset)
```

```
Print(myVar)
```

## Series:

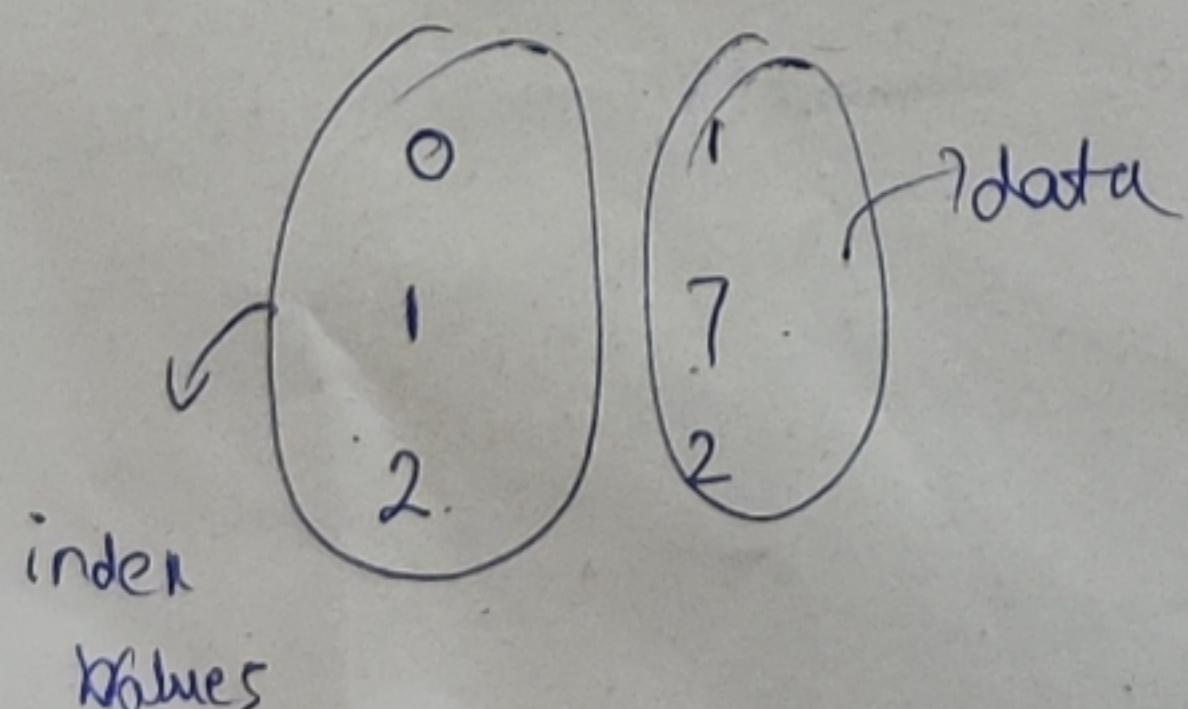
→ A Pandas series is like a column in a table.

→ It is a 1-D array holding data of any type.

```
a = [ 1, 7, 2 ]
```

```
myVar = Pd.Series(a)
```

if nothing specified the values are labelled with index number



index argument is used to specify the labels separated by ,

```
myVar = pd.Series(a, index = ['x', 'y', 'z'])
```

We can access item by referring label

```
Print(myVar["y"]) ↴
```

```
myVar = pd.Series(
```

```
a = [{"day1": 420, "day2": 320, "day3": 300}]
```

↳ [No need of using index in myVar as they are already labelled.]

```
a = {"Fruits": ['Mango', 'Apple', 'Grapes'],  
      "Vegetables": ['Potato', 'Tomato', 'Pumpkin']}
```

}

```
myVar = pd.DataFrame(a)
```

↳ series like a column, and it's a whole table

```
Print(df.loc[0]) → Prints the 1st row
```

↳ used to specify location of row (can be one/more)

| - - - - - |  
| we can also give index | ↴  
| - - - - - |

```
myVar = pd.DataFrame(a, index = ["day1", "day2"])
```

```
pd.read_csv('data.csv')
```

↳ loading a csv file into a Dataframe

```
df = pd.read_csv('data.csv')
```

```
Print(df.to_string())
```

↳ to Print entire Dataframe

if we don't use `to_string()` it will only Print first 5 rows and last 5 rows.

```
[Print(pd.options.display.max_rows)]
```

↳ display max no of rows that can be shown]

to change the no of rows that can be displayed

```
↳ [Import Pandas as pd]
```

```
| pd.options.display.max_rows = 9999
```

```
| df = pd.read_csv('data.csv')
```

```
| Print df
```

CSV - comma separated values

## JSON :

Big data sets are often stored/extracted in JSON.

```
df = pd.read_json('data.json')
```

↳ loading a json file

## Viewing the data :

The `head()` method returns the headers and specific no of rows, starting from top.

```
Print(df.head(10))
```

↳ 10 rows

```
Print(df.head(5))
```

↳ 5 rows

the tail() method returns the headers and a specified no of rows

Print(df.tail())

the info() give more information about data set.

Print(df.info())

### \* Cleaning Data :

Bad data

- Empty cells
- Data in wrong format
- Wrong data
- duplicates

df = pd.read\_csv('data.csv')

new\_df = df.dropna() → [dropna() method return new Dataframe, and will not change the original]

If we want to change original Dataframe use inplace = True

| new\_df = df.dropna(inplace = True) |

→ it will remove all rows containing NULL values from original DF.

To fill empty cell with a Value use

| new\_df = df.fillna(130, inplace = True) |

To replace Null Values from a column specify column name

| df.fillna({'calories': 130}, inplace = True) |

df['calories'].median() → median

df['calories'].mode()[0] → mode

`x = df[["calories"]].mean()`

`df.fillna({ "calories": x }, inplace=True)`

} replace null values in  
calories column with mean

`df['Date'] = pd.to_datetime(df['Date'], format='mixed')`

↳ to convert the messy date into date format

`df.dropna(subset=['Date'], inplace=True)`

### \*Replacing Values :

`df.loc[7, 'Duration'] = 45`

### to replace a large no of Values :

for X in df.index :

if df.loc[X, "Duration"] > 120 :

df.loc[X, "Duration"] = 120

↳ If value is greater than 120 replace with

120

for X in df.index :

if df.loc[X, "Duration"] > 120 :

df.drop(X, inplace=True)

↳ remove rows with no value

`df.drop_duplicates()`

↳ remove row duplicate values

`df.duplicated()`

↳ if any of the row has duplicates gives True/ False

## \* Working with missing data:

Pandas use different sentinel values to represent a missing value depending on data type.

`numpy.nan` for Numpy datatypes. The disadvantage is that the original datatype will be coerced to `np.float64` (or) `object`

Pd. Series([1, 2], dtype=np.int64).reindex([0, 1, 2])

\*output : 0 1.0      ↳ "The np.int64 is coerced  
              1 2.0  
              2 NaN  
              dtype : float64  
              to float64".

Pd. Series([True, False], dtype=np.bool\_).reindex([0, 1, 2])

\*output : 0 True      ↳ "bool coerced to object".  
              1 False  
              2 Nan  
              dtype : object

NAT for Numpy `np.datetime64`, `np.timedelta64` and `PeriodDtype`.

Pd. Series([1, 2], dtype=np.dtype("timedelta64[ns]")).reindex

([0, 1, 2])      ↳ "There will be no

0 0 days 00:00:00.000000001 conversion to  
1 0 days 00:11:11:000000002 another type  
2 NAT

And similar for above mentioned data types.

\* NA for StringDType, Int64DType and ArrowDType

Pd.Series([1, 2], dtype = "Int64").reindex([0, 1, 2])

O/P:

0	1
1	2
2	<NA>

Pd.Series([True, False], dtype = "boolean[Pyarrow]").reindex([0, 1, 2])

O/P:

0	True
1	False
2	<NA>

dtype = bool[Pyarrow]

\* To detect these missing values, use the isna() (or) notna() methods.

Sen = Pd.Series([Pd.Timestamp("2020-01-01"), Pd.NaT])

Pd.isna(Sen)

O/P:

0	False
1	True

Not null Value  
Null Value

dtype = bool

\* isna() / notna() will also consider "None" a missing value.

\* Equality comparisons between np.nan, NaT and NA don't act like none.

None == None

Pd.NaT == Pd.NaT

OP: True

OP: False

np.nan == np.nan

Pd.NA == Pd.NA

OP: False

OP: <NA>

\* An experimental NA value (singleton) is available to represent scalar missing values.

For ex, when having missing values in a series with nullable integer dtype, it will use NA:

s = Pd.Series([1, 2], None, dtype = "Int64")

s[2]

OP: <NA>

s[2] is Pd.NA

OP: True

\* Logical operations:

for logical operations NA follows the rules of the three-valued logic.

for ex, for the logical "or" operation (1), if one of the operands is True, we already know the result is true, regardless of the other value.

True | False

Pd. NA | True

oP : True

oP : True

True | Pd. NA

oP : True

\* If the operands is false, the result depends on the value of the other operand. Hence here NA Propagates.

False | True

oP : True

False | False

oP : False

False | Pd. NA  
oP : <NA>

\* Similar thing happens in the "&" operator.

False & True

False & Pd. NA

oP : False

oP : False

False & False

True & Pd. NA

oP : False

oP : <NA>