

PANDAS :

BOYA RATESH...

1. pandas is a python library used for working with data sets.
2. It has functions for analyzing, cleaning, exploring and manipulating data.
3. The name "Pandas" has reference to both "PANEL DATA" & "PYTHON DATA ANALYSIS".
4. It was created by "WES McKinney" in 2008.

FEATURE'S :

1. cleaning up data
2. Handling missing data
3. Alignment and indexing
4. Great handling of data
5. Handling missing data
6. python support
7. Lot of time series
8. Optimized performance
9. Grouping
10. Multiple file formats are supported.

→ Series + Series → Data frame.

| | a |
|---|---|
| 0 | 5 |
| 1 | 6 |
| 2 | 8 |

+

| | b |
|---|---|
| 0 | 7 |
| 1 | 9 |
| 2 | 4 |

→

| | a | b |
|---|---|---|
| 0 | 5 | 7 |
| 1 | 6 | 9 |
| 2 | 8 | 4 |

PANDAS OPERATIONS :

• Loading packages and Initialization

```
[ ] import numpy as np  
import pandas as pd.
```

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

```
my-data = [10, 20, 30]
```

```
arr = np.array(my-data)
```

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

```
Print(" Dictionary: ", d)
```

```
[ ] output:
```

Dictionary: { 'a': 10, 'b': 20, 'c': 30 }

→ • CREATING A SERIES (Pandas class)

→ From numerical data only

→ From numerical data and corresponding index (row labels)

→ From numpy array as the source of numerical data.

→ just using a pre-defined dictionary.

```
[ ] pd.Series(data=my-data)
```

```
[ ] output:
```

```
0    10  
1    20  
2    30
```

```
dtype: int64
```

```
[ ] Pd.Series(data = my-data, index = labels)
```

```
[ ] output: a 10  
          b 20  
          c 30  
          dtype: int64
```

```
[ ] Pd.Series(arr, labels)
```

```
[ ] output: a 10  
          b 20  
          c 30  
          dtype: int32
```

→ #using a pre-defined Direct Dictionary object

```
[ ] Pd.Series(d)
```

```
[ ] output: a 10  
          b 20  
          c 30  
          dtype: int64.
```

→ What type of values can a Pandas Series hold?

```
[ ] print("\nHolding numerical data\n", '\t25',  
        sep = '\t')  
print(Pd.Series(arr))
```

```
[ ] output:  
          Holding numerical data
```

```
0 10  
1 20  
2 30
```

```
dtype: int32.
```


[] Print("\n Holding text labels\n", '-' * 20, sep='')
Print(Pd.Series(labels))

[] output: Holding text labels

0 a
1 b
2 c

dtype: Object.

[] Print("\n Holding functions\n", '-' * 20, sep='')
Print(Pd.Series(data = [sum, print, len]))

[] output:

Holding functions:

0 <built-in function sum>

1 <built-in function print>

2 <built-in function len>

dtype: object.

[] Print("\n Holding Objects from a dictionary\n",
print ('-' * 40, sep=''))
Print(Pd.Series(data = [d.keys, d.items, d.values]))

[] output:

Holding objects from a dictionary

0 <built-in method keys of dict object
at 0X09A----

1 <built-in method items of dict
Object at 0X09A-----

2 <built-in method values of dict
Object at 0X09A-----

dtype: Object.

INDEXING AND SLICING:

```
[ ] ser1 = pd.Series([1,2,3,4], ['CA', 'OR', 'CO', 'AZ'])  
ser2 = pd.Series([1,2,5,4], ['CA', 'OR', 'NV', 'AZ'])
```

Print("In Indexing by name of the item/object
(String identifier) \n", ' ' * 50, sep=" ")

print("Value of CA in ser1:", ser1['CA'])

print("Value of AZ in ser1:", ser1['AZ'])

print("Value of NV in ser2:", ser2['NV'])

Print("In Indexing by number (positional value in the
series) \n", ' ' * 50, sep=" ")

print("Value for CA in ser1:", ser1[0])

print("Value for AZ in ser1:", ser1[3])

print("Value of NV in ser2:", ser2[2])

Print("In Indexing by a range \n", ' ' * 25, sep=" ")

print("Value for OR, CO, & AZ in ser1: \n", ser1[1:4],

[] Output: sep=" ")

Indexing by name of the item/object (string identifier):

Value of CA in ser1: 1

Value of AZ in ser1: 4

Value of NV in ser2: 5

Indexing by number (positional value in series):

Value for CA in ser1: 1

Value for AZ in ser1: 4

Value for NV in ser2: 5

Indexing by a range:

Value for OR, CO, and AZ in ser1

OR 2

CO 3

AZ 4

dtype: int64.

ADDING/MERGING TWO SERIES WITH COMMON INDEX:

```
[C] Ser1 = Pd.Series([1,2,3,4], ["CA", 'OR', 'CO', 'AZ'])  
Ser2 = Pd.Series([1,2,5,4], ['CA', 'OR', 'NV', 'AZ'])  
Ser3 = Ser1 + Ser2
```

```
Print("\n After adding the two series, the result  
looks like this...\n", '- '*59, sep=' ')  
Print(Ser3)
```

[C] Output:

After adding the two series, the result looks like this....

| | |
|----|-----|
| AZ | 8.0 |
| CA | 2.0 |
| CO | NaN |
| NV | NaN |
| OR | 4.0 |

dtype: float64

□ print('In python tries to add values where it finds common index name, and puts NaN where indices are missing in')

→ print("\n The idea works even for multiplication")
print(Ser1 * Ser2)

[C] Output:

| | |
|----|------|
| AZ | 16.0 |
| CA | 1.0 |
| CO | NaN |
| NV | NaN |
| OR | 4.0 |

dtype: float64


```
[1] print("In Or even for combination of mathematical operations!\n", i, '*5'), sep='')
print(np.exp(ser1) + np.log10(ser2)).
```

[2] Or even for combination of mathematical operations!

AZ 55.200210

CA 2.710282

CO nan

NV nan

OR 7.60006

dtype: float64.

⇒ DATAFRAME (THE REAL MEAT!)

[1] from numpy.random import randn as rn.

CREATING AND ACCESSING DATAFRAME

- Indexing
- Adding and deleting rows & columns.
- Subsetting Dataframe.

[1] np.random.seed(101)

matrix_data = rn(5, 4)

row_labels = ['A', 'B', 'C', 'D', 'E']

column_headings = ['W', 'X', 'Y', 'Z']

df = pd.DataFrame(data=matrix_data, index=row_labels, columns=column_headings)

print("\nTHE data frame look like\n", i, '*45,

print(df) sep='')

[2] Output: The data frame look like:

| | W | X | Y | Z |
|---|-----|------|-----|------|
| A | 2.7 | 5.8 | 9.9 | 8.0 |
| B | 2.4 | 10.4 | 2.2 | 9.7 |
| C | 5.6 | 3.6 | 3.3 | 2.8 |
| D | 2.6 | 4.3 | 9.4 | 9.7 |
| E | 8.7 | 7.1 | 5.6 | 10.2 |

INDEXING AND SLICING (columns)

- By bracket method . . .
- By Dot method (Not recommended).

[1] `Print("In The 'x' column\n", '- ' * 25, sep='')`
`Print(df['x'])`

[2] Output:

The 'x' column :

A 0.628133

B -0.319318

C 0.740122

D -0.758872

E 1.978757

Name: x, dtype: float64

[3] `print("In Type of the column : ",`
`type(df['x']), sep='')`

[4] Output:

Type of the column :

<class 'Pandas.Core.Series.Series'>

[5] `print("In The 'x' and 'z' column`
`indexed by passing a list\n", '- ' * 55)`
`print(df[['x', 'z']])`

[6] The 'x' and 'z' column indexed by
passing a list :

| | x | z |
|---|-----------|-----------|
| A | -0.628133 | 0.503826 |
| B | -0.319318 | 0.605965 |
| C | 0.740122 | -0.955057 |
| D | -0.758872 | 0.590001 |
| E | 1.978757 | 0.683500 |


```
[ ] print("In type of the pair of columns:",  
        type(df[['x', 'z']]), sep='')
```

[c] Type of the pair of columns: <class 'pandas.
core.frame.DataFrame'>

=> For more than one column, the object
turns into a dataframe.

```
[ ] .print("\n The 'x' column accessed by Dot  
method (NOT recommended) \n", '- ' * 55, sep='')  
print(df.x)
```

[c] The 'x' column accessed by Dot method

A 0.620133

B -0.319318

C 0.740122

D -0.758872

E 1.978757

Name: x,

dtype: float64

CREATING and Deleting a column or row

[] Print("In A column is created by assigning it in relation to an existing column in ", ' - ' * 75, sep = " ")

[] print
df['New'] = df['x'] + df['z']
Print(
df['New (sum of x and z)'] = df['x'] + df['z']
Print(df)

[] A column is created by assigning it in relation to an existing column:

| | x | y | z | New (sum of x and z) | |
|---|------|------|------|----------------------|------|
| A | 2.7 | 0.6 | 0.1 | 0.5 | 1.13 |
| B | 0.6 | -0.3 | -0.8 | 0.6 | 0.28 |
| C | -2.0 | 0.7 | 0.5 | -0.5 | 0.15 |
| D | 0.1 | -0.7 | -0.3 | 0.7 | 0.19 |
| E | 0.1 | 1.9 | 2.6 | 0.6 | 2.66 |

[] print("In A column is dropped by using df.drop() method in ", ' - ' * 75, sep = " ")
df = df.drop('New', axis=1)
Notice the axis=1 option, axis=0 is default, so one has to change it to 1
print(df)

[] A column is dropped by assigning drop():

| | x | y | z | New (sum of x & z) | |
|---|------|------|------|--------------------|------|
| A | 2.70 | 0.62 | 0.90 | 0.00 | 1.13 |

| | | | | | |
|---|-------|-------|-------|-------|------|
| B | 0.65 | -0.31 | -0.04 | 0.60 | 0.28 |
| C | -2.01 | 0.740 | 0.90 | -0.50 | 0.15 |
| D | 0.118 | -0.75 | -0.84 | 0.95 | 0.10 |
| E | 0.190 | 1.978 | 2.60 | 0.68 | 2.66 |

[] df1 = df.drop('A')

Print("An in place change can be done by making inplace")

Print("In A row(index) is dropped by using df.drop() method and axis=0 in", " *0.5")

Print(df1)

[ii] Output:

A row(index) is dropped by using drop=1 df.drop() method and axis=0:

| | W | X | Y | Z | New(sum of X and Z) |
|---|-------|-------|-------|-------|---------------------|
| B | 0.65 | -0.31 | -0.04 | 0.60 | 0.28 |
| C | -2.01 | 0.740 | 0.90 | -0.50 | 0.15 |
| D | 0.118 | -0.75 | -0.84 | 0.95 | 0.10 |
| E | 0.190 | 1.978 | 2.60 | 0.68 | 2.66 |

Print("In An in place change can be done by making inplace-TRUE u the drop method in")

[] df.drop('New(sum of X and Z)', axis=1, inplace=True)

Print(df)

[ii] Output: An in place change can be done by making New(sum of X and Z) in place-TRUE u drop method:

| | W | X | Y | Z |
|---|-------|-------|-------|-------|
| A | 2.70 | 0.62 | 0.90 | 0.50 |
| B | 0.65 | -0.31 | 0.90 | 0.60 |
| C | -2.01 | 0.74 | -0.84 | -0.50 |
| D | 0.118 | -0.75 | 0.90 | 0.95 |
| E | 0.190 | 1.97 | 2.60 | 0.68 |

Selecting/ Indexing Rows...

- Label-based 'loc' method.
- Index(enumerate) 'iloc' method.

```
[ ] Print("\n Label-based 'loc' method can be  
used for selecting row(s) in", '-' * 60)  
Print("\n Single row in")  
Print(df.loc['c'])
```

```
[c] Output:  
Label-based 'loc' method can be used for  
Selecting row(s):  
Single row:
```

| | |
|---|--------|
| W | -2.010 |
| X | 0.740 |
| Y | 0.52 |
| Z | -0.58 |

Name: c, dtype: float64

```
[ ] Print("\n Multiple rows in")  
Print(df.loc[['B', 'c']])
```

```
[c] Multiple rows:
```

| | | | | |
|---|--------|-------|------|-------|
| | W | X | Y | Z |
| B | 0.651 | -0.31 | -0.8 | 0.60 |
| C | -2.018 | 0.74 | 0.52 | -0.50 |

```
[ ] Print("\n Index position based 'iloc'  
method can be used for selecting row(s)", '-' * 60)  
Print("\n Single row in")  
Print(df.iloc[2]).
```

[5] Output:

Index position based 'iloc' method can be used for selecting row(s):

Single row:

W -2.01

X 0.74

Y 0.52

Z -0.58

Name: C, dtype: float64.

[] print("\n multiple rows\n").

print(df.iloc[[1,2]])

[6] output:

Multiple rows:

| | W | X | Y | Z |
|---|---------|--------|-------|-------|
| B | 0.6511 | -0.319 | -0.84 | 0.60 |
| C | -2.0101 | 0.74 | 0.52 | -0.50 |

SUBSETTING DataFrame...

[] `Print("\n The DataFrame\n ", ' - ' * 45, sep=' ')`
`Print(df)`

[c] Output:

The DataFrame

| | W | X | Y | Z |
|---|--------|--------|--------|--------|
| A | 2.706 | 0.6281 | 0.907 | 0.503 |
| B | 0.651 | -0.319 | -0.848 | 0.605 |
| C | -2.010 | 0.740 | 0.528 | -0.589 |
| D | 0.181 | -0.758 | -0.933 | 0.955 |
| E | 0.190 | 1.978 | 2.605 | 0.632 |

[] `Print("\n Element at row 'B' and Column 'Y' is\n")`
`Print(df.loc['B', 'Y'])`

[c] Output:

Element at row 'B' and Column 'Y' is -0.848076

[] `Print("\n Subset comparing of rows B & D, and column W and Y, is\n")`
`Print(df.loc[['B', 'D'], ['W', 'Y']])`

[c] Output:

Subset comparing of rows B and D, and column W and Y, is:

| | W | Y |
|---|-------|--------|
| B | 0.651 | -0.840 |
| D | 0.188 | -0.933 |

CONDITIONAL SELECTION, INDEX

RESETTING, MULTI-INDEX:

- Basic idea of conditional check and Boolean DataFrame

```
[ ] print("In Boolean DataFrame where we are  
checking if the values are greater than 0",  
1-1*75, sep=' - ')  
print(df > 0).
```

[c] Output:

Boolean DataFrame where we are checking if the
values are greater than 0:

| | W | X | Y | Z |
|---|-------|-------|-------|-------|
| A | True | True | True | True |
| B | True | False | False | True |
| C | False | True | True | False |
| D | True | False | False | True |
| E | True | True | True | True |

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

[c] Output:

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

```
[ ] booldf = df > 0
```

```
print("In DataFrame indexed by boolean dataframe",  
1-1*45, sep=' - ')  
print(df[booldf]).
```

[c] Output: DataFrame indexed by boolean dataframe:

| | W | X | Y | Z |
|---|---------|-------|--------|-------|
| A | 2.7068 | 0.629 | 0.9079 | 0.501 |
| B | 0.65111 | NaN | NaN | 0.66 |
| C | NaN | 0.740 | 0.5288 | NaN |
| D | 0.18869 | NaN | NaN | 0.955 |
| E | 0.19078 | 1.978 | 2.605 | 0.613 |

PASSING BOOLEAN SERIES TO CONDITIONALLY SUBSET THE DATAFRAME;

[]

```
matrix-data = np.matrix([22, 60, 140;  
42, 70, 148;  
30, 62, 125;  
15, 61, 160;  
25, 62, 152])
```

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

```
Column-headings = ['Age', 'Height', 'Weight']
```

```
pd.  
df = DataFrame(data=matrix-data, index=row-  
labels, columns=column-headings)
```

```
Print("In The data-frame n", "\n", sep=' ')  
print(df)
```

[]

Output:

A new Dataframe

| | Age | Height | Weight |
|---|-----|--------|--------|
| A | 22 | 60 | 40 |
| B | 42 | 70 | 48 |
| C | 30 | 62 | 125 |
| D | 15 | 68 | 160 |
| E | 25 | 62 | 152 |

[]

```
print("In Rows with height > 65 inch n",  
      "\n", sep=' ')  
print(df[df['Height'] > 65])
```

[]

Output: Rows with height > 65 inch

| | Age | Height | Weight |
|---|-----|--------|--------|
| A | 22 | 60 | 140 |
| B | 42 | 70 | 148 |
| D | 35 | 67 | 160 |

```
[1] bool df1 = df['Height'] > 65
    bool df2 = df['Weight'] > 145
    Print("In Rows with height > 65 inch &
    weight > 145 in", 1 - 1 * 50, sep = ' ')
    Print(df[(bool df1) & (bool df2)])
```

[2] Output:

Rows with Height > 65 inch and weight
> 145 :

| | Age | Height | Weight |
|--|-----|--------|--------|
|--|-----|--------|--------|

| | | | |
|---|----|----|-----|
| B | 42 | 70 | 148 |
|---|----|----|-----|

| | | | |
|---|----|----|-----|
| D | 35 | 68 | 160 |
|---|----|----|-----|

```
[3] Print("In Dataframe with only Age and
    Weight column whose Height > 65 inch in",
    1 - 1 * 60, sep = ' ')
    print(df[bool df1][['Age', 'Weight']])
```

[4] Output:

DataFrame with only Age and weight
column whose Height > 65 inch :

| | Age | Weight |
|--|-----|--------|
|--|-----|--------|

| | | |
|---|----|-----|
| A | 22 | 160 |
|---|----|-----|

| | | |
|---|----|-----|
| B | 42 | 140 |
|---|----|-----|

| | | |
|---|----|-----|
| D | 35 | 160 |
|---|----|-----|

RE-SETTING AND SETTING INDEX :

[]

```
matrix_data = np.matrix ('22, 66, 140;  
42, 76, 148;  
30, 62, 125;  
35, 64, 160;  
25, 62, 152')
```

```
row_labels = ['A', 'B', 'C', 'D', 'E']
```

```
Column_labels = ['Age', 'Height', 'Weight']
```

```
df = pd.DataFrame(data = matrix_data,  
index = row_labels, columns = Column_labels)
```

```
print ("In After resetting index in", '- '*35)
```

```
print (df.reset_index())
```

[c]

output:

After resetting index

| | index | Age | Height | Weight |
|---|-------|-----|--------|--------|
| 0 | A | 22 | 66 | 140 |
| 1 | B | 42 | 76 | 148 |
| 2 | C | 30 | 62 | 125 |
| 3 | D | 35 | 64 | 160 |
| 4 | E | 25 | 62 | 152 |

[]

```
print ("In After resetting index with  
'drop' option TRUE in", '- '*45, sep = '\n')
```

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

[c]

output:

After resetting index with 'drop' option true:

| | Age | Height | Weight |
|---|-----|--------|--------|
| 0 | 22 | 66 | 140 |
| 1 | 42 | 76 | 148 |

| | | | |
|---|----|----|-----|
| 2 | 30 | 62 | 125 |
| 3 | 35 | 68 | 160 |
| 4 | 25 | 62 | 152 |

[] print("In Adding a new column 'profession'
in", ' - ' * 45, sep=' ')

df['profession'] = 'Student Teacher Engineer
Doctor Nurse ".split()' .

print(df)

[] output:

Adding a new column 'profession'

| | Age | Height | Weight | profession |
|---|-----|--------|--------|------------|
| A | 22 | 66 | 140 | Student |
| B | 42 | 70 | 148 | Teacher |
| C | 30 | 62 | 125 | engineer |
| D | 35 | 68 | 160 | Doctor |
| E | 25 | 62 | 152 | Nurse |

[] print("In setting 'profession' column at
index in", ' - ' * 45, sep=' ')

print(df.set_index('profession'))

[] output:

setting 'profession' column at index:

| Profession | Age | Height | Weight |
|------------|-----|--------|--------|
| Student | 22 | 66 | 140 |
| Teacher | 42 | 70 | 148 |
| engineer | 30 | 62 | 125 |
| Doctor | 35 | 68 | 160 |
| Nurse | 25 | 62 | 152 |

MULTI-INDEXING :

```
[ ] #Index Level.  
outside = ['G1', 'G1', 'G1', 'G2', 'G2', 'G2']  
inside = [1, 2, 3, 1, 2, 3]  
hier_index = list(zip(outside, inside))  
Print("In tuple pairs after the zip and list  
Command In", '\n', sep = '\n')  
Print(hier_index).
```

```
[c] Output :  
Tuple pairs after the zip and list command :  
[('G1', 1), ('G1', 2), ('G1', 3), ('G2', 1), ('G2', 2),  
 ('G2', 3)]
```

```
[ ] hier_index = pd.MultiIndex.from_tuples(hier_index)  
Print("In Index hierarchy In", '\n', sep = '\n')  
print(hier_index)
```

```
[c] Index Hierarchy :
```

```
MultiIndex ( [ ('G1', 1),  
                ('G1', 2),  
                ('G1', 3),  
                ('G2', 1),  
                ('G2', 2),  
                ('G2', 3)], )
```

```
[ ] print("In Index hierarchy type In", '\n', sep = '\n')  
print(type(hier_index))
```

```
[c] Output :  
Index hierarchy type :
```

```
<class 'pandas.core.indexes.multi.  
MultiIndex'>
```



```
[1] Print("In Creating Dataframe with multi-
Index\n", ' - ' * 37, sep = '\n')
[2] np.random.seed(101)
df1 = pd.DataFrame(data = np.round(rn(6,3),2),
index = hier_index, columns = ['A', 'B', 'C'])
print(df1)
```

[2] Output: creating Dataframe with multi-index

| | | A | B | C |
|----|---|-------|-------|-------|
| G1 | 1 | 2.71 | 0.63 | 0.91 |
| | 2 | 0.50 | 0.65 | -0.32 |
| | 3 | -0.05 | 0.61 | -2.02 |
| G2 | 1 | 0.79 | 0.53 | -0.55 |
| | 2 | 0.19 | -0.76 | -0.93 |
| | 3 | 0.96 | 0.19 | 1.98 |

```
[3] Print("In Subsetting multi-Index Dataframe using
two 'loc' methods\n", ' - ' * 60, sep = '\n')
print(df1.loc['G2'].loc[['1','3']][['B', 'C']])
```

[3] Output: subsetting multi-Index Dataframe using two 'loc' methods:

| | B | C |
|---|------|-------|
| 1 | 0.53 | -0.55 |
| 3 | 0.19 | 1.98 |

```
[4] print("In Naming the index by 'index.names'
method\n", ' - ' * 45, sep = '\n')
df1.index.names = ['outer', 'inner']
print(df1)
```

[4] Output: Naming the index by 'index.names' method:

| | | A | B | C |
|-------|-------|-------|-------|-------|
| outer | inner | | | |
| | 1 | 2.71 | 0.63 | 0.91 |
| | 2 | 0.50 | 0.65 | -0.32 |
| G1 | 3 | -0.05 | 0.61 | -2.02 |
| | 1 | 0.79 | 0.53 | -0.55 |
| | 2 | 0.19 | -0.76 | -0.93 |
| G2 | 3 | 0.96 | 0.19 | 1.98 |

CROSS-SECTION ('XS') COMMAND:

- [] Print ("in creating a cross-section from outer
- section from Outer level in", '...45')
Print(df1.XS('G1'))

[] Output:

Creating a cross-section from outer section level:

| | A | B | C |
|-------|-------|------|-------|
| Inner | | | |
| 1 | 2.71 | 0.63 | 0.91 |
| 2 | 0.56 | 0.65 | -0.32 |
| 3 | -0.85 | 0.61 | -2.02 |

- [] Print("in Grabbing a cross-section from inner
level (for all outer level) in")
Print(df1.XS(2, level = "Inner"))

[] Output:

Grabbing a cross-section from inner
level (for all outer level) 2:

| | A | B | C |
|-------|------|-------|-------|
| outer | | | |
| G1 | 0.56 | 0.65 | -0.32 |
| G2 | 0.19 | -0.76 | -0.93 |

MISSING VALUES

```
[1] df = pd.DataFrame({'A': [1, 2, np.nan],  
                    'B': [5, np.nan, np.nan],  
                    'C': [1, 2, 3]})  
df['states'] = "CA NV AZ".split()  
df.set_index('states', inplace=True)  
print(df)
```

[1] Output:

| | A | B | C |
|----|-----|-----|---|
| CA | 1.0 | 5.0 | 1 |
| NV | 2.0 | NaN | 2 |
| AZ | NaN | NaN | 3 |

[] Pandas 'dropna' method:
 Print("In Dropping any rows with a Nan Value")
 print(df.dropna(axis=0))

[c] output:

Dropping any rows with a Nan Value:

| | A | B | C |
|-------|-----|-----|---|
| stats | | | |
| CA | 1.0 | 5.0 | 1 |

[] Print("In Dropping any column with a Nan Value")
 print(df.dropna(axis=1))

[c] output:

Dropping any column with a Nan Value:

| | C |
|-------|---|
| stats | |
| CA | 1 |
| NV | 2 |
| AZ | 3 |

[] Print("In Dropping a row with a minimum 2
 Nan value using 'thresh' parameter")
 print(df.dropna(axis=0, thresh=2))

[c] Dropping a row with a minimum 2 Nan value
 using 'thresh' parameter:

| | A | B | C |
|-------|-----|-----|---|
| stats | | | |
| CA | 1.0 | 5.0 | 1 |
| NV | 2.0 | NAN | 2 |

PANDAS "FILLNA" METHOD:

[1] print("In Filling Values with a default
Value in", '-', 'ss, sep='')

```
print(df.fillna(value='FILL VALUE'))
```

[2] output:

Filling Values with a default value:

| | A | B | C |
|-------|------------|------------|---|
| State | | | |
| CA | 1 | 5 | 1 |
| NV | 2 | FILL VALUE | 2 |
| AZ | FILL VALUE | FILL VALUE | 3 |

[3] print("In Filling Values with a computed value
(mean of column A here) in", '-', '60, sep='')

```
print(df.fillna(value=df['A'].mean()))
```

[2] output:

Filling Values with a computed value (mean of
column A here):

| | A | B | C |
|-------|-----|-----|---|
| State | | | |
| CA | 1.0 | 5.0 | 1 |
| NV | 2.0 | 1.5 | 2 |
| AZ | 1.5 | 1.5 | 3 |

Group By Method :

[] Creating a dataframe

```
data = { 'company': ['Google', 'Google', 'PS', 'MSFT', 'FB', 'FB'],  
        'person': ['Sam', 'charlie', 'Amy', 'vanessa', 'Carl', 'Sarah'],  
        'sales': [200, 120, 340, 124, 243, 350]}
```

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

```
print(df)
```

[]

| | company | person | sales |
|--|---------|--------|-------|
|--|---------|--------|-------|

| | | | |
|---|--------|-----|-----|
| 0 | Google | Sam | 200 |
|---|--------|-----|-----|

| | | | |
|---|--------|---------|-----|
| 1 | Google | charlie | 120 |
|---|--------|---------|-----|

| | | | |
|---|------|-----|-----|
| 2 | MSFT | Amy | 340 |
|---|------|-----|-----|

| | | | |
|---|------|---------|-----|
| 3 | MSFT | vanessa | 124 |
|---|------|---------|-----|

| | | | |
|---|----|------|-----|
| 4 | FB | Carl | 243 |
|---|----|------|-----|

| | | | |
|---|----|-------|-----|
| 5 | FB | Sarah | 350 |
|---|----|-------|-----|

[]

```
bycomp = df.groupby('company')
```

```
print("In grouping by 'company' column and  
listing mean sales in", '-1*55, sep='\n')
```

```
print(bycomp.mean())
```

[]

output:

Grouping by 'company' column & list by mean sales:

| | sales |
|---------|-------|
| company | |
| FB | 296.5 |
| Google | 160.0 |
| MSFT | 232.0 |


```
print("\n Grouping by 'company' column and  
listing sum of sales in", '\n', '\t * $5, sep='')  
print(bycomp.sum())
```

output:

Grouping by 'company' column and listing sum of sales:

| company | Sales |
|---------|-------|
| FB | 593 |
| GOOG | 320 |
| MSFT | 464 |

```
print("\n All in one line of command (stats for 'FB')  
in", '\n', '\t * $6, sep='') # transpose  
print(pd.DataFrame(df.groupby('company').describe().  
loc['FB']).transpose())
```

output: All in one line of command (stats for 'FB'):

| | | Sales | | | | | | | |
|----|-------|-------|-------|-------|--------|-------|--------|-------|--|
| | count | mean | std | min | 25% | 50% | 75% | max | |
| FB | 2.0 | 296.5 | 75.66 | 293.0 | 269.75 | 296.5 | 323.25 | 358.0 | |

Print("\n same type of extraction with little different
command in", '\n', '\t * \$68, sep='')

```
print(df.groupby('company').describe().loc[['GOOG',  
MSFT]])
```

output:

Same type of extraction with little different command:

| | | Sales | | | | | | | |
|------|-------|-------|-------|-------|-------|-------|-------|-------|--|
| | count | mean | std | min | 25% | 50% | 75% | max | |
| GOOG | 2.0 | 160.0 | 8.5 | 120.0 | 100.0 | 160.0 | 160.0 | 200.0 | |
| MSFT | 2.0 | 232.0 | 17.23 | 129.0 | 175.0 | 232.0 | 285.0 | 340.0 | |

MERGING, JOINING, CONCATENATING

Concatenation:

[] #creating data frames

```
df1 = pd.DataFrame({'A': ['A0', 'A1', 'A2', 'A3'],  
                    'B': ['B0', 'B1', 'B2', 'B3'],  
                    'C': ['C0', 'C1', 'C2', 'C3'],  
                    'D': ['D0', 'D1', 'D2', 'D3']},  
                    index=[0, 1, 2, 3])
```

```
df2 = pd.DataFrame({'A': ['A4', 'A5', 'A6', 'A7'],  
                    'B': ['B4', 'B5', 'B6', 'B7'],  
                    'C': ['C4', 'C5', 'C6', 'C7'],  
                    'D': ['D4', 'D5', 'D6', 'D7']},  
                    index=[4, 5, 6, 7])
```

```
df3 = pd.DataFrame({'A': ['A8', 'A9', 'A10', 'A11'],  
                    'B': ['B8', 'B9', 'B10', 'B11'],  
                    'C': ['C8', 'C9', 'C10', 'C11'],  
                    'D': ['D8', 'D9', 'D10', 'D11']},  
                    index=[8, 9, 10, 11])
```

[] df_cat1 = pd.concat([df1, df2, df3], axis=0,
 print("In. after concatenation along row axis,
 i.e. * 30, sep=' '")
 print(df_cat1).

(c) After concatenation along row:

| | A | B | C | D |
|----|-----|-----|-----|-----|
| 0 | A0 | B0 | C0 | D0 |
| 1 | A1 | B1 | C1 | D1 |
| 2 | A2 | B2 | C2 | D2 |
| 3 | A3 | B3 | C3 | D3 |
| 4 | A4 | B4 | C4 | D4 |
| 5 | A5 | B5 | C5 | D5 |
| 6 | A6 | B6 | C6 | D6 |
| 7 | A7 | B7 | C7 | D7 |
| 8 | A8 | B8 | C8 | D8 |
| 9 | A9 | B9 | C9 | D9 |
| 10 | A10 | B10 | C10 | D10 |
| 11 | A11 | B11 | C11 | D11 |

[] `df_cat2 = pd.concat([df1, df2, df3], axis=1)`
`print("After concatenation along column is", 'is')`
`print(df_cat2)`

(c) Output: After concatenation along column:

| | A | B | C | D | A | B | C | D | A | B | C | D |
|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0 | A0 | B0 | C0 | D0 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 1 | A1 | B1 | C1 | D1 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 2 | A2 | B2 | C2 | D2 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 3 | A3 | B3 | C3 | D3 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 4 | NaN | NaN | NaN | NaN | A4 | B4 | C4 | D4 | NaN | NaN | NaN | NaN |
| 5 | NaN | NaN | NaN | NaN | A5 | B5 | C5 | D5 | NaN | NaN | NaN | NaN |
| 6 | NaN | NaN | NaN | NaN | A6 | B6 | C6 | D6 | NaN | NaN | NaN | NaN |
| 7 | NaN | NaN | NaN | NaN | A7 | B7 | C7 | D7 | NaN | NaN | NaN | NaN |
| 8 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | A8 | B8 | C8 | D8 |
| 9 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | A9 | B9 | C9 | D9 |
| 10 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | A10 | B10 | C10 | D10 |
| 11 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | A11 | B11 | C11 | D11 |

[] Print ("In After filling missing values with
Zero\n", '-', 'Go', sep='')

Print(df_cat2)

[C] df_cat.fillna(value=0, inplace=True)

Output:

After filling missing values with zero:

| | A | B | C | D | A | B | C | D | A | B | C | D |
|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|
| 0 | A0 | B0 | C0 | D0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | A1 | B1 | C1 | D1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | A2 | B2 | C2 | D2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | A3 | B3 | C3 | D3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | A4 | B4 | C4 | D4 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | A5 | B5 | C5 | D5 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | A6 | B6 | C6 | D6 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0 | A7 | B7 | C7 | D7 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | A8 | B8 | C8 | D8 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | A9 | B9 | C9 | D9 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | A10 | B10 | C10 | D10 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | A11 | B11 | C11 | D11 |

Merging by a common 'key':

→ The merge function allows you to merge Data frames together using a similar logic as merging SQL tables together.

```
[1] left = pd.DataFrame({'key': ['K0', 'K1', 'K2', 'K3'],  
                      'A': ['A0', 'A1', 'A2', 'A3'],  
                      'B': ['B0', 'B1', 'B2', 'B3']})
```

```
right = pd.DataFrame({'key': ['K0', 'K1', 'K2', 'K3'],  
                      'C': ['C0', 'C1', 'C2', 'C3'],  
                      'D': ['D0', 'D1', 'D2', 'D3']})
```

```
[2] print("In The DataFrame 'left' in", '- ' + 30, sep='')  
print(left)
```

[3] Output:

The DataFrame 'left'

| | key | A | B |
|---|-----|----|----|
| 0 | K0 | A0 | B0 |
| 1 | K1 | A1 | B1 |
| 2 | K2 | A2 | B2 |
| 3 | K3 | A3 | B3 |

```
[4] print("In The DataFrame 'right' in", '- ' + 30, sep='')  
print(right)
```

[5] The DataFrame 'right'

| | key | C | D |
|---|-----|----|----|
| 0 | K0 | C0 | D0 |
| 1 | K1 | C1 | D1 |
| 2 | K2 | C2 | D2 |
| 3 | K3 | C3 | D3 |

```
[ ] merge1 = pd.merge(left, right, how='inner', on='key')
print('In After simple merging with 'inner'
      method\n', '\n', sep='\n')
print(merge1).
```

[C] Output:

After Simple merging with 'inner' method

| | key | A | B | C | D |
|---|-----|----|----|----|----|
| 0 | k0 | A0 | B0 | C0 | D0 |
| 1 | k1 | A1 | B1 | C1 | D1 |
| 2 | k2 | A2 | B2 | C2 | D2 |
| 3 | k3 | A3 | B3 | C3 | D3 |

⇒ JOINING:

Joining is a convenient method for combining the columns of two potentially differently-indexed DataFrames into a single DataFrame based on 'index keys'.

```
[ ] left = pd.DataFrame({'A': ['A0', 'A1', 'A2'],
                        'B': ['B0', 'B1', 'B2'],
                        index = ['k0', 'k1', 'k2']})
right = pd.DataFrame({'C': ['C0', 'C2', 'C3'],
                      'D': ['D0', 'D2', 'D3'],
                      index = ['k0', 'k2', 'k3']})
print(left)
```

[C]

| | A | B |
|----|----|----|
| k0 | A0 | B0 |
| k1 | A1 | B1 |
| k2 | A2 | B2 |

[C] print(right)

| | C | D |
|----|----|----|
| k0 | C0 | D0 |
| k2 | C2 | D2 |
| k3 | C3 | D3 |

[] left.join(right)

[c]

| | A | B | C | D |
|----|----|----|-----|-----|
| K0 | A0 | B0 | C0 | D0 |
| K1 | A1 | B1 | NaN | NaN |
| K2 | A2 | B2 | C2 | D2 |

[] left.join(right, how = 'outer')

[c]

| | A | B | C | D |
|----|-----|-----|-----|-----|
| K0 | A0 | B0 | C0 | D0 |
| K1 | A1 | B1 | NaN | NaN |
| K2 | A2 | B2 | C2 | D2 |
| K3 | NaN | NaN | C3 | D3 |

⇒

USEFUL OPERATIONS:

head() and Unique Value

- head()
- Unique()
- nunique()
- Value-Count()

[] import pandas as pd

```
df = pd.DataFrame({'Col1': [1, 2, 3, 4, 5],  
                   'Col2': [444, 555, 666, 444, 333],  
                   'Col3': 'aaa bb c dd eee'.
```

[c]

| | Col1 | Col2 | Col3 |
|---|------|------|-------|
| 0 | 1 | 444 | aaa |
| 1 | 2 | 555 | bb |
| 2 | 3 | 666 | c |
| 3 | 4 | 444 | dd |
| 4 | 5 | 333 | eeee. |

[] print("\n Method head() is for showing
first few entries\n", 1, '\n 50')
df.head()

[c] Method head() is for showing first few
entries:

| | col1 | col2 | col3 |
|---|------|------|------|
| 0 | 1 | 444 | aaa |
| 1 | 2 | 555 | bb |
| 2 | 3 | 666 | c |
| 3 | 4 | 444 | dd |
| 4 | 5 | 333 | cccc |

[] print("\n Finding Unique Values in 'col2'\n",
1, '\n 40').

print(df['col2'].unique())

[c] Finding unique values in col2:

[444 555 666 333 222 777]

[] print("\n Finding number of unique values
in 'col2'\n", 1, '\n 40, sep='\n')

t1 = df['col2'].value_counts()

[c] Finding no. of unique values in 'col2':

6

[] print("\n Table of unique values in

'col2'\n", 1, '\n 40, sep='\n')

t1 = df['col2'].value_counts()

print(t1)

[1]

Table of unique values in 'col2':

| | |
|-----|---|
| 666 | 3 |
| 444 | 2 |
| 555 | 2 |
| 022 | 1 |
| 333 | 1 |
| 777 | 1 |

Name: col2, dtype: int64

APPLYING FUNCTIONS:

→ Pandas work with 'apply' method in accept any user-defined function.

[1]

Define a function.

```
def testfunc(x):
    if x > 500:
        return (10 * np.log10(x))
    else:
        return (x/10)
```

[1]

df['funcApplied'] = df['col2'].apply(testfunc)

print(df).

| | col1 | col2 | col3 | funcApplied |
|----|------|------|------|-------------|
| 0 | | | | 44.40 |
| 1 | 1 | 444 | aaa | 27.99 |
| 2 | 2 | 555 | bb | 28.32 |
| 3 | 3 | 666 | dd | 44.40 |
| 4 | 4 | 444 | eecc | 33.30 |
| 5 | 5 | 333 | ffcf | 22.20 |
| 6 | 6 | 202 | gg | 27.70 |
| 7 | 7 | 666 | h | 28.23 |
| 8 | 8 | 777 | iii | 27.64 |
| 9 | 9 | 666 | j | |
| 10 | 10 | 555 | | |

Rajesh Boya