# Data Handling using Pandas

## Pandas:

- It is a package useful for data analysis and manipulation.
- Pandas provide an easy way to create, manipulate and wrangle the data.
- Pandas provide powerful and easy-to-use data structures, as well
  as the means to quickly perform operations on these structures.

Data scientists use Pandas for its following advantages:

- Easily handles missing data.
- It uses Series for one-dimensional data structure and DataFrame for multi-dimensional data structure.
- It provides an efficient way to slice the data.
- It provides a flexible way to merge, concatenate or reshape the data.

#### **DATA STRUCTURE IN PANDAS**

A data structure is a way to arrange the data in such a way that so it can be accessed quickly and we can perform various operation on this data like- retrieval, deletion, modification etc.

Pandas deals with 3 data structure-

- 1. Series
- 2. Data Frame
- 3. Panel

We are having only series and data frame in our syllabus.

#### Series

Series-Series is a one-dimensional array like structure with homogeneous data, which can be used to handle and manipulate data. What makes it special is its index attribute, which has incredible functionality and is heavily mutable.

#### It has two parts-

- 1. Data part (An array of actual data)
- 2. Associated index with data (associated array of indexes or data labels)

#### e.g.-

Index	Data
0	10
1	15
2	18
3	22

- ✓ We can say that Series is a labeled one-dimensional array
  which can hold any type of data.
- ✓ Data of Series is always mutable, means it can be changed.
- ✓ But the size of Data of Series is always immutable, means it cannot be changed.
- ✓ Series may be considered as a Data Structure with two arrays out which one array works as Index (Labels) and the second array works as original Data.
- ✓ Row Labels in Series are called Index.

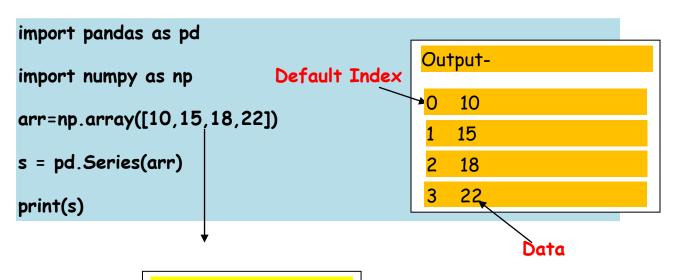
# Syntax to create a Series:

<Series Object>=pandas.Series (data, index=idx (optional))

√ Where data may be python sequence (Lists), ndarray, scalar value or a python dictionary.

#### **How to create Series with nd array**

#### Program-



Here we create an array of 4 values.

#### **How to create Series with Mutable index**

#### Program-

```
import pandas as pd

import numpy as np

arr=np.array(['a','b','c','d'])

s=pd.Series(arr,
    index=['first','second','third','fourth'])

print(s)
Output-

first a

second b

third c

fourth d
```

# Creating a series from Scalar value

To create a series from scalar value, an index must be provided. The scalar value will be repeated as per the length of index.

# Creating a series from a Dictionary

```
# import the pandas lib as pd
import pandas as pd

# create a dictionary

d = {'Name' : 'Hardik', 'Iplteam' : 'MI', 'Runs' : 1500}

# create a series

s = pd.Series(d)

print(s)
```

```
Name Hardik
Iplteam MI
Runs 1500
dtype: object
```

# Mathematical Operations in Series

```
import pandas as pd
s=pd.Series([1,2,3,4,5])
print('To Multiply all values in a series by 2')
print('-----
print(s*2)
print('To Find the Square of all the values in a series ')
print('-----
print(s**2)
print('To print all the values in a series that are greater than 2')
print(s[s>2])
To Multiply all values in a series by 2
                     Print all the values of the Series by multiplying them by 2.
    10
dtype: int64
To Find the Square of all the values in a series
     1
                     Print Square of all the values of the series.
    16
    25
dtype: int64
To print all the values in a series that are greater than 2
                    Print all the values of the Series that are greater than 2.
dtype: int64
```

# Example-2

```
import pandas as pd
s1=pd.Series([1,2,3,4,5],index=['a','b','c','d','e'])
s2=pd.Series([10,20,30,40,50],index=['a','b','c','d','e'])
s3=pd.Series([5,14,23,32],index=['a','b','c','d'])
print('To Add Series1 & series2')
print('----')
print(s1+s2)
print('To Add Series2 & Series3')
print('----')
print('To Add Series2 & series3 and Filled Non Matching Index with 0')
print(s2.add(s3,fill_value=0))
To Add Series1 & series2
    11
b 22
  33
   44
   55
dtype: int64
To Add Series2 & Series3
a 15.0
b
  34.0
               While adding two series, if Non-Matching Index is found in either of the
c 53.0
  72.0
               Series, Then NaN will be printed corresponds to Non-Matching Index.
    NaN-
dtype: float64
To Add Series2 & series3 and Filled Non Matching Index with 0
  15.0
b 34.0
  53.0
               If Non-Matching Index is found in either of the series, then this Non-
d 72.0
               Matching Index corresponding value of that series will be filled as 0.
    50.0-
dtype: float64
```

# Head and Tail Functions in Series

head (): It is used to access the first 5 rows of a series.

Note: To access first 3 rows we can call series\_name.head(3)

```
1 import pandas as pd
 2 import numpy as np
 3 arr=np.array([10,15,18,22,55,77,42,48,97])
 4 # create a series from array
 5 s = pd.Series(arr)
 6 # to print fiest 5 rows
 7 print (s.head())
 8 # To print first 3 rows
   print(s.head(3))
0
    10
                     Result of s.head()
1
    15
2
    18
3
    22
    55
dtype: int32
                     Result of s.head(3)
    10
1
    15
2
    18
dtype: int32
```

tail(): It is used to access the last 5 rows of a series.

Note: To access last 4 rows we can call series\_name.tail (4)

```
import pandas as pd
import numpy as np
arr=np.array([10,15,18,22,55,77,42,48,97])

# create a series from array
s = pd.Series(arr)
# to print last 5 rows
print (s.tail())
# To print last 4 rows
print(s.tail(4))

# 55
5 77
6 42
```

#### Selection in Series

Series provides index label loc and ilocand [] to access rows and columns.

loc index label : Syntax:-series\_name.loc[StartRange: StopRange]
 Example-

```
import pandas as pd
 2 import numpy as np
 3 arr=np.array([10,15,18,22,55,77])
 4 s = pd.Series(arr)
                               To Print Values from Index 0 to 2
 5 print(s)
   print(s.loc[:2])
    print(s.loc[3:4])
                                 To Print Values from Index 3 to 4
   s.loc[2:3]
0
    10
1
    15
2
    18
3
    22
    55
5
    77
dtype: int32
    10
1
    15
2
    18
dtype: int32
    22
    55
dtype: int32
2
    18
    22
dtype: int32
```

#### 2. Selection Using iloc index label:-

Syntax:-series\_name.iloc[StartRange : StopRange]

Example-

2

18 dtype: int32

```
import pandas as pd
 2 import numpy as np
 3 arr=np.array([10,15,18,22,55,77])
 4 s = pd.Series(arr)
 5 print(s)
 6 print(s.iloc[:2])-
                                 To Print Values from Index 0 to 1.
 7 print(s.iloc[3:4])
 8 s.iloc[2:3]
0
    10
    15
1
2
    18
3
    22
4
    55
    77
5
dtype: int32
    10
    15
1
dtype: int32
    22
3
dtype: int32
```

# 3. Selection Using []:

```
Syntax:-series_name[StartRange>: StopRange] or
       series_name[index]
```

#### Example-

1

2

15

18 dtype: int32

```
import pandas as pd
 2 import numpy as np
   arr=np.array([10,15,18,22,55,77])
 4 s = pd.Series(arr)
   print(s)
 6 print(s[1])
   print('\n')
                                   To Print Values at Index 3.
   print(s[3:4]) -
   s[:3]
0
     10
1
     15
2
     18
3
     22
     55
4
5
     77
dtype: int32
15
3
     22
dtype: int32
0
     10
```

# Indexing in Series

Pandas provide index attribute to get or set the index of entries or values in series.

#### Example-

```
import pandas as pd
import numpy as np
arr=np.array(['a','b','c','d'],)
s=pd.Series(arr,index=['first','second','third','fourth'])
print(s)
# To print only indexes in series
print('\n indexes in Series are:::')
print(s.index)
```

```
first    a
second    b
third    c
fourth    d
dtype: object

indexes in Series are:::
Index(['first', 'second', 'third', 'fourth'], dtype='object')
```

# Slicing in Series

Slicing is a way to retrieve subsets of data from a pandas object. A slice object syntax is -

#### SERIES\_NAME [start:end: step]

The segments start representing the first item, end representing the last item, and step representing the increment between each item that you would like.

#### Example:-

```
import pandas as pd
 2 import numpy as np
    arr=np.array([10,15,18,22,55,77])
 4 s = pd.Series(arr,index=['A','B','C','D','E','F'])
    print(s)
    print(s[1:5:2])
    print(s[0:6:2])
 8
     10
Α
В
     15
C
     18
D
     22
     55
     77
dtype: int32
     15
     22
dtype: int32
     10
     18
     55
dtype: int32
```

#### DATAFRAME

<u>DATAFRAME-</u>It is a two-dimensional object that is useful in representing data in the form of rows and columns. It is similar to a spreadsheet or an SQL table. This is the most commonly used pandas object. Once we store the data into the Dataframe, we can perform various operations that are useful in analyzing and understanding the data.

#### **DATAFRAME STRUCTURE**

columns_		PLAYERNAME	IPLTEAM	BASEPRICEINCR
0	Ī	ROHIT	MI	13
1	Ī	VIRAT	RCB	17
2	Ī	HARDIK	MI	14
$\uparrow$	_		<u> </u>	
INDEX			DATA	
				╗

- 1. A Dataframe has axes (indices)-
  - > Row index (axis=0)

**PROPERTIES OF DATAFRAME** 

- Column index (axes=1)
- 2. It is similar to a spreadsheet, whose row index is called index and column index is called column name.
- 3. A Dataframe contains Heterogeneous data.
- 4. A Dataframe Size is Mutable.
- 5. A Dataframe Data is Mutable.

#### A data frame can be created using any of the following-

- 1. Series
- 2. Lists
- 3. Dictionary
- 4. A numpy 2D array

#### How to create Empty Dataframe

: import pandas as pd
df=pd.DataFrame()
print(df)

Empty DataFrame Columns: [] Index: []

#### **How to create Dataframe From Series**

# Program import pandas as pd s = pd.Series(['a','b','c','d']) df=pd.DataFrame(s) print(df) Output 1 a 2 b Default Column Name As 0 3 c 4 d

# DataFrame from Dictionary of Series

#### Example-

```
import pandas as pd
name=pd.Series(['Hardik','Virat'])
team=pd.Series(['MI','RCB'])
dic={'Name':name,'Team':team}
df=pd.DataFrame(dic)
print(df)

Name Team
0 Hardik MI
1 Virat RCB
```

## DataFrame from List of Dictionaries

#### Example-

```
Name SirName
O Sachin Bhardwaj
1 Vinod Verma
2 Rajesh Mishra
```

# Iteration on Rows and Columns

If we want to access record or data from a data frame row wise or column wise then iteration is used. Pandas provide 2 functions to perform iterations-

- 1. iterrows ()
- 2. iteritems ()

# iterrows()

It is used to access the data row wise. Example-

```
Name
           SirName
0 Sachin Bhardwaj
   Vinod
             Verma
Row index is :: 0
Row Value is::
Name
            Sachin
         Bhardwaj
SirName
Name: 0, dtype: object
Row index is :: 1
Row Value is::
Name
          Vinod
SirName
          Verma
Name: 1, dtype: object
```

# iteritems()

It is used to access the data column wise.

Example-

```
0 Sachin Bhardwaj
1 Vinod Verma

Column Name is :: Name
Column Values are::
0 Sachin
1 Vinod
Name: Name, dtype: object

Column Name is :: SirName
Column Values are::
0 Bhardwaj
1 Verma
Name: SirName, dtype: object
```

Name SirName

# Select operation in data frame

To access the column data ,we can mention the column name as subscript.

```
e.g. - df[empid] This can also be done by using df.empid.

To access multiple columns we can write as df[ [col1, col2,---] ]
```

#### Example -

102

104

1

Vinod 15-01-2012

Anil 17-01- 2012

103 Lakhbir 05-09-2007

4 105 Devinder 05-09-2007 5 106 UmaSelvi 16-01-2012

```
>>df.empid or df['empid']
    101
0
1
    102
2
    103
    104
3
4
    105
5
    106
Name: empid, dtype: int64
>>df[['empid', 'ename']]
  empid
                 ename
    101
                 Sachin
0
                  Vinod
1
    102
2
    103
                 Lakhbir
3
    104
                    Anil
4
    105
               Devinder
```

UmaSelvi

# To Add & Rename a column in data frame

```
import pandas as pd
s = pd.Series([10, 15, 18, 22])
df=pd.DataFrame(s)
df.columns=['List1'] To Rename the default column of Data
                         Frame as List1
df['List2']=20 To create a new column List2 with all values
                  as 20
                                        Output-
df['List3']=df['List1']+df['List2']
                                          List1 List2 List3
Add Column1 and Column2 and store in
                                            10
                                                20
                                                     30
                                            15
                                                20
                                                     35
New column List3
                                            18
                                                20
                                                     38
                                            22
                                                     42
                                                20
print(df)
```

# To Delete a Column in data frame

We can delete the column from a data frame by using any of the the following -

```
    del
    pop()
    drop()
```

```
Output-

List1 List2

0 10 20

1 15 20

2 18 20

3 22 20
```

>>df.pop('List2') we can simply delete a column by passing column name in pop method.

# To Delete a Column Using drop()

```
import pandas as pd
s= pd.Series([10,20,30,40])
df=pd.DataFrame(s)
df.columns=['List1']
df['List2']=40
df1=df.drop('List2',axis=1) ____ (axis=1) means to delete Data
                                 column wise
df2=df.drop(index=[2,3],axis=0) (axis=0) means to delete
                              data row wise with given index
print(df)
print(" After deletion::")
print(df1)
print (" After row deletion::")
print(df2)
Output-
 List1 List2
0 10 40
1 20 40
2 30 40
3 40 40
After deletion::
  List1
0 10
1 20
2 30
3 40
After row deletion::
 List1
   10
   20
```

	_	

# Accessing the data frame through loc() and iloc() method or indexing using Labels

Pandas provide loc() and iloc() methods to access the subset from a data frame using row/column.

# Accessing the data frame through loc()

It is used to access a group of rows and columns.

Syntax-

#### Df.loc[StartRow: EndRow, StartColumn: EndColumn]

Note -If we pass: in row or column part then pandas provide the entire rows or columns respectively.

```
import pandas as pd
    Runs={ 'TCS': { 'Qtr1':2500, 'Qtr2':2000, 'Qtr3':3000, 'Qtr4':2000},
 3
             'WIPRO': {'Qtr1':2800,'Qtr2':2400,'Qtr3':3600,'Qtr4':2400},
 4
 5
             'L&T': { 'Qtr1':2100,'Qtr2':5700,'Qtr3':35000,'Qtr4':2100}}
 7
    df=pd.DataFrame(Runs)
    print(df)
                                     To access a single row
    print(df.loc['Qtr3', : ]) -
print(df.loc['Qtr1':'Qtr3'
10
11
      TCS WIPRO
                    L&T
Qtr1 2500
           2800
                   2100
Qtr2 2000
             2400
                    5700
                  35000
Otr3
     3000
             3600
                               To access multiple Rows Qtr1 to Qtr3
Qtr4 2000
             2400
                    2100
TCS
         3000
         3600
WIPRO
       35000
L&T
Name: Qtr3, dtype: int64
      TCS WIPRO
                    L&T
Otr1 2500 2800
                   2100
Qtr2 2000 2400
                  5700
Qtr3 3000 3600 35000
```

#### Example 2:-

Qtr1 2500

Otr2 2000

Otr3 3000

Qtr4 2000

2800

2400

3600

2400

```
import pandas as pd
    Runs={ 'TCS': { 'Qtr1':2500, 'Qtr2':2000, 'Qtr3':3000, 'Qtr4':2000},
 2
 3
            'WIPRO': {'Qtr1':2800,'Qtr2':2400,'Qtr3':3600,'Qtr4':2400},
 4
 5
            'L&T': { 'Qtr1':2100,'Qtr2':5700,'Qtr3':35000,'Qtr4':2100}}
 6
    df=pd.DataFrame(Runs)
                                     To access single column
    print(df)
 8
    print(df.loc[ : ,'TCS' ])
    print(df.loc[ : , 'TCS':'WIPRO'])
11
      TCS WIPRO
                    L&T
Qtr1 2500
            2800
                   2100
                            To access Multiple Column namely TCS and WIPRC
Otr2 2000
            2400
                   5700
Otr3 3000
            3600
                 35000
Otr4 2000
            2400
                   2100
0tr1
       2500
Otr2
       2000
Qtr3
      3000
Otr4
       2000
Name: TCS, dtype: int64
      TCS WIPRO
```

#### Example-3

0

2

102

101 Sachin 12-01-2012

103 Lakhbir 05-09-2007

Vinod 15-01-2012

```
import pandas as pd
 2
    empdata={ 'empid':[101,102,103,104,105,106],
              'ename':['Sachin','Vinod','Lakhbir','Anil','Devinder','UmaSelvi'],
 3
 4
              'Doj':['12-01-2012','15-01-2012','05-09-2007','17-01-2012','05-09-2007','16-01-2012']}
 5
    df=pd.DataFrame(empdata)
                                   To access first row
    print(df)
 6
    print(df.loc[0])
    df.loc[0:2]
                                              To access first 3 Rows
                           Doj
  empid
            ename
    101
           Sachin
                    12-01-2012
    102
            Vinod
                    15-01-2012
2
    103
          Lakhbir
                    05-09-2007
    104
             Anil 17-01- 2012
    105 Devinder
                    05-09-2007
    106 UmaSelvi
                    16-01-2012
               101
empid
            Sachin
ename
Doj
        12-01-2012
Name: 0, dtype: object
   empid ename
                     Doj
```

# Accessing the data frame through iloc()

It is used to access a group of rows and columns based on numeric index value.

Syntax-

Df.loc[StartRowindexs: EndRowindex, StartColumnindex: EndColumnindex]

Note -If we pass: in row or column part then pandas provide the entire rows or columns respectively.

```
import pandas as pd
    Runs={ 'TCS': { 'Qtr1':2500,'Qtr2':2000,'Qtr3':3000,'Qtr4':2000},
 3
 4
            'WIPRO': {'Otr1':2800,'Otr2':2400,'Otr3':3600,'Otr4':2400},
 5
            'L&T': { 'Qtr1':2100,'Qtr2':5700,'Qtr3':35000,'Qtr4':2100}}
 6
 7
    df=pd.DataFrame(Runs)
    print(df)
                                     To access First two Rows
    print(df.iloc[0 :2 ,1:2 ])_
    print(df.iloc[ : , 0:2])
10
                                     and Second column
11
      TCS WIPRO
                   L&T
                                   To access all Rows and First
Qtr1 2500
            2800
                   2100
Otr2 2000
            2400
                   5700
                                   Two columns Record
Otr3 3000
            3600 35000
Otr4 2000
            2400
                  2100
     WIPRO
0tr1
      2800
Otr2
      2400
      TCS WIPRO
Qtr1 2500
            2800
Otr2 2000
            2400
Qtr3 3000
            3600
Otr4 2000
            2400
```

# head() and tail() Method

The method head() gives the first 5 rows and the method tail() returns the last 5 rows.

#### Output-

	•			
	Doj	empid	ename	
0	12-01-2012	101	Sachin	
1	15-01-2012	102	Vinod	
2	05-09-2007	103	Lakhbir	Data Frame
3	17-01-2012	104	Anil	
4	05-09-2007	105	Devinder	
5	16-01-2012	106	UmaSelvi	
	Doj	empid	ename	
0	12-01-2012	101	Sachin	
1	15-01-2012	102	Vinod	head() displays first 5 rows
2	05-09-2007	103	Lakhbir	
3	17-01-2012	104	Anil	
4	05-09-2007	105	Devinder	
	Doj	empid	ename	
1	15-01-2012	102	Vinod	
2	05-09-2007	103	Lakhbir	
3	17-01-2012	104	Anil	tail() display last 5 rows
4	05-09-2007	105	Devinder	
5	16-01-2012	106	UmaSelvi	

To display first 2 rows we can use head(2) and to returns last2 rows we can use tail(2) and to return  $3_{rd}$  to  $4_{th}$  row we can write df[2:5].

```
import pandas as pd
empdata={ 'Doj':['12-01-2012','15-01-2012','05-09-2007',
              '17-01-2012', '05-09-2007', '16-01-2012'],
          'empid':[101,102,103,104,105,106],
          'ename':['Sachin','Vinod','Lakhbir','Anil','Devinder','UmaSelvi']
df=pd.DataFrame(empdata)
print(df)
print(df.head(2))
print(df.tail(2))
print(df[2:5])
           Doj empid
                        ename
Output-
0 12-01-2012
                101
                        Sachin
1 15-01-2012 102
                        Vinod
2 05-09-2007 103 Lakhbir
3 17-01- 2012 104
                         Anil
4 05-09-2007 105 Devinder
5 16-01-2012 106 UmaSelvi
          Doi
               empid ename
0 12-01-2012
                 101 Sachin
                                         head(2) displays first 2 rows
1 15-01-2012 102 Vinod
           Doj empid
                        ename
4 05-09-2007
                105 Devinder
                                       tail(2) displays last 2 rows
5 16-01-2012
                106 UmaSelvi
           Doj empid
                       ename
2 05-09-2007
                103 Lakhbir
3 17-01- 2012 104
                         Anil
                                        df[2:5] display 2<sub>nd</sub> to 4<sub>th</sub> row
4 05-09-2007 105 Devinder
```

# Boolean Indexing in Data Frame

Boolean indexing helps us to select the data from the DataFrames using a boolean vector. We create a DataFrame with a boolean index to use the boolean indexing.

```
import pandas as pd
   dic= {
           'Name': ['Sachin Bhardwaj', 'Vinod Verma', 'Rajesh Mishra'],
 3
 4
           'Age': [32, 35, 40]
 5
 6 # creating a DataFrame with boolean index vector
    df = pd.DataFrame(dic, index = [True, False, True])
 7
   print(df)
    print(df.loc[True]).
                                      To Return Data frame where index is True
10 print()
11 print('Result of iloc method')
12 print(df.iloc[1])
                                 We can pass only integer value in iloc
                  Name
                        Age
True
       Sachin Bhardwaj
                         32
False
          Vinod Verma
                         35
True
         Rajesh Mishra
                         40
                 Name Age
True Sachin Bhardwaj
                        32
True
        Rajesh Mishra
```

Result of iloc method Name Vinod Verma Age 35 dtype: object

# Concat operation in data frame

Pandas provides various facilities for easily combining together **Series**, **DataFrame**.

pd.concat(objs, axis=0, join='outer', join\_axes=None,ignore\_index=False)

- objs This is a sequence or mapping of Series, DataFrame, or Panel objects.
- axis {0, 1, ...}, default 0. This is the axis to concatenate along.
- join {'inner', 'outer'}, default 'outer'. How to handle indexes on other axis(es). Outer for union and inner for intersection.
- ignore\_index boolean, default False. If True, do not use the index values on the concatenation axis. The resulting axis will be labeled 0, ..., n - 1.
- join\_axes This is the list of Index objects. Specific indexes to use for the other (n-1) axes instead of performing inner/outer set logic.

The Concat() performs concatenation operations along an axis.

#### Example-1

```
id Value1 Value2
0 1
        A
        C
              D
1 2
2 3
        E
 4
3
        G
4 5
       I
0 2
       K
1 3
       M
2 6
        0
3 7
        0
4 8
        S
              Т
```

#### Example-2

```
1 import pandas as pd
 2 dic1= { 'id': ['1', '2', '3', '4', '5'], 'Value1': ['A', 'C', 'E', 'G', 'I'],
            'Value2': ['B', 'D', 'F', 'H', 'J']}
 3
 4 dic2= {'id': ['2', '3', '6', '7', '8'], 'Value1': ['K', 'M', '0', 'Q', 'S'],
           'Value2': ['L', 'N', 'P', 'R', 'T']}
 6 df1=pd.DataFrame(dic1)
 7 df2=pd.DataFrame(dic2)
 8 df3=pd.concat([df1,df2],ignore index=True)
 9 print(df3)
10
  id Value1 Value2
         A
0 1
1 2
         C
                              If you want the row labels to adjust automatically
2 3
         E
3 4
         G
               Н
                              according to the join, you will have to set the
5 2
               L
         K
                              argument ignore_index as True while
                                                                                 calling
6 3
               N
7 6
         0
                             the concat() function:
         Q
 7
                T
```

#### Example-3

```
1 import pandas as pd
   dic1= { 'id': ['1', '2', '3', '4', '5'], 'Value1': ['A', 'C', 'E', 'G', 'I'],
 6 df1=pd.DataFrame(dic1)
 7 df2=pd.DataFrame(dic2)
 8 merge={'Data1':df1,'Data2':df2}
 9 df3=pd.concat(merge)
10 print(df3)
11
      id Value1 Value2
Data1 0 1
             A
             C
                  D
    1 2
                  F
     2 3
             E
                         pandas also provides you with an option to label
             G
                  Н
     3 4
            Ι
     4 5
                  J
                         the DataFrames, after the concatenation, with
Data2 0 2
             K
                  L
                         a key so that you may know which data came
     1
       3
             M
                  N
     2 6
             0
     3 7
                         from which DataFrame.
     4 8
                  Т
```

#### Example-4

```
import pandas as pd
 2 dic1= { 'id': ['1', '2', '3', '4', '5'], 'Value1': ['A', 'C', 'E', 'G', 'I'],
           'Value2': ['B', 'D', 'F', 'H', 'J']}
 3
   dic2= {'id': ['2', '3', '6', '7', '8'], 'Value1': ['K', 'M', '0', 'Q', 'S'],
 4
           'Value2': ['L', 'N', 'P', 'R', 'T']}
 5
 6 df1=pd.DataFrame(dic1)
 7 df2=pd.DataFrame(dic2)
   df3=pd.concat([df1,df2],axis=1)
 9
    print(df3)
10
  id Value1 Value2 id Value1 Value2
0 1
                B 2
                                            To concatenate DataFrames
               D 3
1 2
         C
                               N
                                            along column, you can specify
2 3
         E
              F 6
3 4
         G
              H 7
                                R
                                            the axis parameter as 1.
               J 8
4 5
         Ι
                               T
```

# Merge operation in data frame

Two DataFrames might hold different kinds of information about the same entity and linked by some common feature/column. To join these DataFrames, pandas provides multiple functions like merge(), join() etc.

#### Example-1

```
import pandas as pd
    dic1= { 'id': ['1', '2', '3', '4', '5'], 'Value1': ['A', 'C', 'E', 'G', 'I'],
            'Value2': ['B', 'D', 'F', 'H', 'J']}
 3
    dic2= {'id': ['2', '3', '6', '7', '8'], 'Value1': ['K', 'M', '0', 'Q', 'S'],
 5
            'Value2': ['L', 'N', 'P', 'R', 'T']}
    dic3 = {'id': ['1', '2', '3', '4', '5', '7', '8', '9', '10', '11'],
            'Value3': [12, 13, 14, 15, 16, 17, 15, 12, 13, 23]}
 7
    df1=pd.DataFrame(dic1)
    df2=pd.DataFrame(dic2)
10 df3=pd.concat([df1,df2])
11 df4=pd.DataFrame(dic3)
12 df5=pd.merge(df3,df4,on='id')
13 print(df5)
 id Value1 Value2 Value3
  1
                       12
                            This will give the common rows between the
 2
         C
                       13
                D
                           two data frames for the corresponding column
2 2
         K
                L
                      13
```

3 F 14 3 М Ν 14 5 4 G Н 15 Ι J 16 7 7 Q R 17 S Τ 15

values ('id').

### Example-2

```
import pandas as pd
   dic1= { 'id': ['1', '2', '3', '4', '5'], 'Value1': ['A', 'C', 'E', 'G', 'I'],
            'Value2': ['B', 'D', 'F', 'H', 'J']}
 3
   dic2= {'id': ['2', '3', '6', '7', '8'], 'Value1': ['K', 'M', '0', 'Q', 'S'],
            'Value2': ['L', 'N', 'P', 'R', 'T']}
 5
   dic3 = {'id': ['1', '2', '3', '4', '5', '7', '8', '9', '10', '11'],
            'Value3': [12, 13, 14, 15, 16, 17, 15, 12, 13, 23]}
 7
   df1=pd.DataFrame(dic1)
 8
 9 df2=pd.DataFrame(dic2)
10 df3=pd.concat([df1,df2])
11 df4=pd.DataFrame(dic3)
12 df5=pd.merge(df3,df4,left on='id', right on='id')
13 print(df5)
```

id Value1 Value2 Value3 0 1 Α В 12 1 2 С D 13 K 13 2 2 L 3 Ε F 3 14 3 4 М 14 Ν 4 G Н 5 15 5 J 6 Ι 16 7 7 R 17 5 Τ 8 8 15

It might happen that the column on which you want to merge the Data Frames have different names (unlike in this case). For such merges, you will have to specify the arguments left\_on as the left DataFrame name and right\_on as the right DataFrame name.

# Join operation in data frame

It is used to merge data frames based on some common column/key.

1. Full Outer Join:- The full outer join combines the results of both the left and the right outer joins. The joined data frame will contain all records from both the data frames and fill in NaNs for missing matches on either side. You can perform a full outer join by specifying the how argument as outer in merge() function.

#### Example-

```
import pandas as pd
  dic1= { 'id': ['1', '2', '3', '4', '5'], 'Value1': ['A', 'C', 'E', 'G', 'I'],
         'Value2': ['B', 'D', 'F', 'H', 'J']}
  4
5
  df1=pd.DataFrame(dic1)
  df2=pd.DataFrame(dic2)
  df3=pd.merge(df1,df2,on='id',how='outer')-
                                             resulting DataFrame
                                         The
  print(df3)
id Value1_x Value2_x Value1_y Value2_y
                                               NaN
                                                     values
                                                              for
                                        with
1
        Α
               В
                     NaN
                            NaN
```

1 2 C D K L 2 3 Ε F М Ν 3 4 G Н NaN NaN Ι J NaN NaN 5 6 NaN NaN 0 Р 7 NaN NaN Q R 7 8 NaN NaN 5 Τ the entries from both the tables with NaN values for missing matches on either side. However, one more thing to notice is the suffix which got appended to the column names to show which column came from which DataFrame. The default suffixes are x and y, however, you can modify them by specifying the suffixes argument in the merge() function.

```
id Value1_left Value2_left Value1_right Value2_right
```

NaN	NaN	В	Α	1	0
L	K	D	C	2	1
N	М	F	Е	3	2
NaN	NaN	Н	G	4	3
NaN	NaN	J	I	5	4
Р	0	NaN	NaN	6	5
R	Q	NaN	NaN	7	6
T	S	NaN	NaN	8	7

2. Inner Join: The inner join produce only those records that match in both the data frame. You have to pass inner in how argument inside merge() function.

3. RightJoin: The right join produce a complete set of records from data frame B(Right side Data Frame) with the matching records (where available) in data frame A(Left side data frame). If there is no match right side will contain null. You have to pass right in how argument inside merge() function.

```
id Value1 x Value2 x Value1 y Value2 y
0 2
            С
                     D
                              K
            Ε
                     F
1
  3
                              М
                                       Ν
                              0
          NaN
                   NaN
3
  7
                                       R
          NaN
                   NaN
                              S
          NaN
                   NaN
                                       Τ
```

4. Left Join: The left join produce a complete set of records from data frame A(Left side Data Frame) with the matching records (where available) in data frame B(Right side data frame). If there is no match left side will contain null. You have to pass left in how argument inside merge() function.

```
id Value1_x Value2_x Value1_y Value2_y
0 1
           Α
                   В
                          NaN
                                  NaN
1 2
           C
                   D
                           Κ
                                   L
          E
2 3
                   F
                          М
                                   N
          G
3 4
                   Н
                         NaN
                                  NaN
4 5
          Ι
                   J
                         NaN
                                  NaN
```

**5. Joining on Index**:-Sometimes you have to perform the join on the indexes or the row labels. For that you have to specify right\_index( for the indexes of the right data frame ) and left\_index( for the indexes of left data frame) as True.

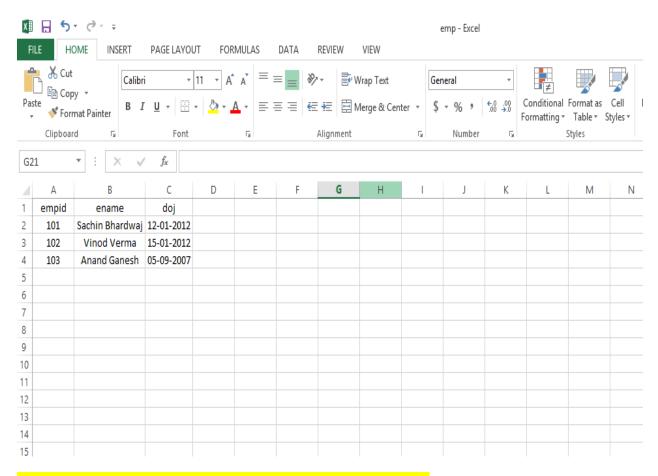
```
id_x Value1_x Value2_x id_y Value1_y Value2_y
    1
                         2
            Α
   2
                        3
1
            C
                    D
                                 Μ
                                         N
2
   3
           Ε
                    F 6
                                 0
3
            G
                        7
    4
                    Н
                                         R
                                 5
   5
            Τ
                    J
                                         Τ
```

# CSV File

A CSV is a comma separated values file, which allows data to be saved in a tabular format. CSV is a simple file such as a spreadsheet or database. Files in the csv format can be imported and exported from programs that store data in tables, such as Microsoft excel or Open Office.

CSV files data fields are most often separated, or delimited by a comma. Here the data in each row are delimited by comma and individual rows are separated by newline.

To create a csv file, first choose your favorite text editor such as- Notepad and open a new file. Then enter the text data you want the file to contain, separating each value with a comma and each row with a new line. Save the file with the extension.csv. You can open the file using MS Excel or another spread sheet program. It will create the table of similar data.



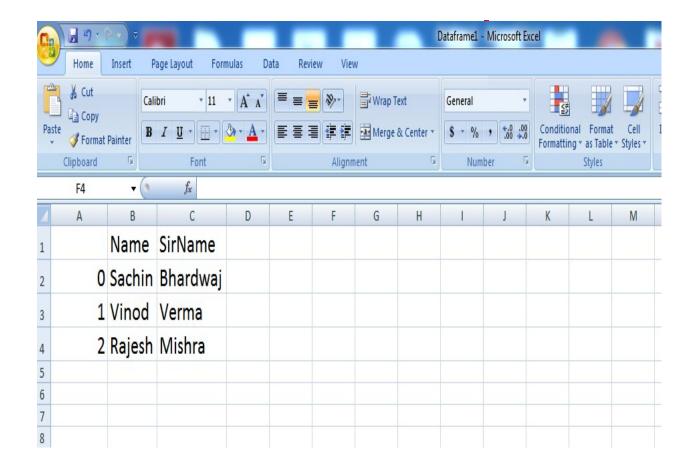
pd.read\_csv() method is used to read a csv file.

```
# importing pandas module
import pandas as pd
# making data frame
df = pd.read_csv("E:\emp.csv")
print(df)

empid ename doj
101 Sachin Bhardwaj 12-01-2012
1 102 Vinod Verma 15-01-2012
2 103 Anand Ganesh 05-09-2007
```

# Exporting data from dataframe to CSV File

To export a data frame into a csv file first of all, we create a data frame say df1 and use dataframe.to\_csv(' E:\Dataframe1.csv') method to export data frame df1 into csv file Dataframe1.csv.



And now the content of df1 is exported to csv file Dataframe1.

# **Data cleaning**

Data cleaning means fixing bad data in your data set.

#### Bad data could be:

- . Empty cells
- . Data in wrong format
- . Wrong data
- Duplicates

## **Empty Cells**

Empty cells can potentially give you a wrong result when you analyze data.

import pandas as pd

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

new\_df = df.dropna()

print(new\_df.to\_string())

#Notice in the result that some rows have been removed (row 18, 22 and 28).

#These rows had cells with empty values.

By default, the dropna() method returns a new DataFrame, and will not change the original.

If you want to change the original DataFrame, use the inplace = True argument:

Remove all rows with NULL values:

import pandas as pd

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

df.dropna(inplace = True)

```
print(df.to_string())
```

# Data of Wrong Format

Cells with data of wrong format can make it difficult, or even impossible, to analyze data.

To fix it, you have two options: remove the rows, or convert all cells in the columns into the same format.

convert all cells in the 'Date' column into dates.

```
import pandas as pd

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

df['Date'] = pd.to_datetime(df['Date'])

print(df.to_string())
        Pandas - Fixing Wrong Data
import pandas as pd

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

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

print(df.to_string())
```