

## ▼ Pandas

### ▼ pandas DataFrame

```
1 import numpy as np
2 import pandas as pd
3 # importing NumPy & pandas
```

```
1 student_values = [["Rayan", 53, "DBDA"], ["Nelly", 56, "DAC"], ["Eminem", 57, "DITISS"]]
2 # creating a 2-D list for DataFrame
3 df = pd.DataFrame(student_values, columns=['Name', 'RollNo', 'Course'])
4 # using 2-D DataFrame to create DataFrame with columns 'Name', 'RollNo', 'Course'
5 df
```

	Name	RollNo	Course
0	Rayan	53	DBDA
1	Nelly	56	DAC
2	Eminem	57	DITISS

```
1 df = pd.DataFrame([['Sudeep', 3545, 'MBA'], ['Ruchika', 1254, 'MCA'], ['Rahul', 6547, 'MS'], ['Purnima', 1389, 'BE']], columns=['Name', 'RollNo', 'Qualification'])
2 # Creating 2nd DataFrame with columns 'Name', 'RollNo', 'Qualification'
3 df
```

	Name	RollNo	Qualification
0	Sudeep	3545	MBA
1	Ruchika	1254	MCA
2	Rahul	6547	MS
3	Purnima	1389	BE

### ▼ DataFrame.iloc[Integer Locator]

```
1 mydf = df.iloc[[0, 2]]
2 # df.iloc(integer location/index or array of indexes or slice of indices)
3 # here, internal bracket is list/array of indices
4 # creating a new DataFrame mydf with rows 0 & 2 from previous DataFrame df
5 mydf
```

**Name RollNo Qualification**

```
1 mydf.iloc[1]
2 # printing 1st row from mydf DataFrame using iloc property
```

```
Name      Rahul
RollNo     6547
Qualification  MS
Name: 2, dtype: object
```

```
1 import os
2 # importing os package
3 os.getcwd()
4 # checking current working directory
```

```
'/content'
```

```
1 from google.colab import files
2 uploaded=files.upload()
3 # ERPData.xlsx
4 # to be used with google colab
5
6 # os.chdir(r'D:/advanced-analytics-files/day03')
7 # to change current working directory to specified path
8 # to be used while running on local system
```

No file chosen

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving ERPData.xlsx to ERPData.xlsx

## ▼ pd.read\_excel('filename.xlsx')

- to open a Workbook as a DataFrame

```
1 df = pd.read_excel('ERPDData.xlsx')
2 # reading excel file to DataFrame
3 df.head()
4 # returns first 5 records from DataFrame
```

	MaterialID	Location	Quantity
0	TMI-43T	MWH-4	34
1	AXCP-78	MWH-1	67
2	LXCV-21	MWH-2	27
3	AXCP-78	MWH-5	65
4	AXCP-78	MWH-4	36

## ▼ DataFrame.Shape

```
1 df.shape
2 # prints the shape in (r, c) format for DataFrame

(50, 3)
```

### ▼ DataFrame['columnName']

```
1 df['Quantity']
2 # method to fetch one column from dataframe as Series Type
3 # df.Quantity    # alt method to fetch one column
```

```
0      34
1      67
2      27
3      65
4      36
5      78
6      31
7      29
8      10
9     120
10     34
11     58
12     76
13     32
14     65
15     87
16     12
17    102
18     34
19     52
20     39
21     75
22     48
23     71
24    152
25    109
26     57
27     83
28     34
29     57
30     27
31     65
32     43
33    112
34     26
35     31
36     28
37     65
38     70
39    145
40     34
41     57
42     31
43     38
44     42
45     69
```

```

46      85
47      29
48      75
49      39
Name: Quantity, dtype: int64

```

```

1 x = list(df['Quantity'])
2 # df['Quantity'] returns a series type which is now converted into list

```

```

1 x[:5]
2 # slicing the list

```

```
[34, 67, 27, 65, 36]
```

```

1 len(x)
2 # finding length of list

```

```
50
```

#### ▼ Q. find the count of records with quantity >=75

```

1 arr1 = np.array(x)
2 # converting list of quantities to ndarray
3 idx = np.where(arr1 >= 75)
4 # using where conditon with ndarray to find quantities >= 75
5 len(idx[0])
6 # finding length of index array
7
8 # alt method
9 # len(np.where(np.array(x)>=75)[0])

```

```
13
```

#### ▼ Q. find the count of records with quantity >=75

```

1 (arr1>=75).sum()
2 # alt method

```

```
13
```

#### ▼ Q. find count of records where location is MWH-3

```

1 arr1 = np.array(df.Location)
2 # creating an ndarray with Location column from DataFrame
3 (arr1=='MWH-3').sum()
4 # using condition on ndarray and then sum() operator to count records

```

```
4
```

```
1 # alt method
2 (df.Location == 'MWH-3').sum()
```

4

```
1 idx
2 # printing the indices that satisfy the condition
```

```
(array([ 5,  9, 12, 15, 17, 21, 24, 25, 27, 33, 39, 46, 48]),)
```

```
1 mydf = df.iloc[idx]
2 # using indices in iloc to print the records that satisfy the condition
3 mydf
4 # note that the indices are improper and according to the previous DataFrame
```

	MaterialID	Location	Quantity
5	TMI-43T	MWH-4	78
9	TMI-43T	MWH-1	120
12	TMI-43T	MWH-4	76
15	GCVB-79	MWH-2	87
17	SDRT-67	MWH-5	102
21	TMI-43T	MWH-4	75
24	AXCP-78	MWH-1	152
25	AXCP-78	MWH-1	109
27	LXCV-21	MWH-5	83
33	TMI-43T	MWH-5	112
39	AXCP-78	MWH-1	145
46	AXCP-78	MWH-2	85
48	AXCP-78	MWH-1	75

```
1 mydf.iloc[[0, 1, 2]]
2 # accessing mydf DataFrame using default indices,
3 # but it prints indices from previous indices
```

	MaterialID	Location	Quantity
5	TMI-43T	MWH-4	78
9	TMI-43T	MWH-1	120
12	TMI-43T	MWH-4	76

▼ Q. find records with quantity > 100 and location == MWH-5

```

1 ind = np.where((df['Quantity'] >= 100) & (df['Location'] == 'MWH-5'))
2 # using two conditions in np.where() using logical '&' operator
3 ind
4 # printing the array of indices satisfying the condition

(array([17, 33]),)

```

```

1 mydf = df.iloc[ind]
2 # using indices in iloc to get values from previous DataFrame
3 mydf
4 # here also, it shows indices as per previous DataFrame

```

	MaterialID	Location	Quantity
17	SDRT-67	MWH-5	102
33	TMI-43T	MWH-5	112

▼ Q. records where either material id is tmi43t or the quantity >=100

```

1 ind = np.where((df['MaterialID'] == 'TMI-43T') | (df['Quantity'] >= 100))
2 # using two conditions in np.where() using logical '|' operator
3 ind
4 # printing the array of indices satisfying the condition

(array([ 0,  5,  6,  7,  9, 12, 13, 14, 17, 18, 19, 20, 21, 24, 25, 31, 32,
        33, 39]),)

```

```

1 len(ind[0])
2 # getting the count of records where condition satisfied

19

```

```

1 mydf = df.iloc[ind]
2 # using indices in iloc to get values from previous DataFrame
3 mydf
4 # here also, it shows indices as per previous DataFrame

```

	MaterialID	Location	Quantity
0	TMI-43T	MWH-4	34
5	TMI-43T	MWH-4	78
6	TMI-43T	MWH-4	31
7	TMI-43T	MWH-2	29
9	TMI-43T	MWH-1	120
12	TMI-43T	MWH-4	76
13	TMI-43T	MWH-4	32
14	TMI-43T	MWH-3	65
17	SDRT-67	MWH-5	102
18	TMI-43T	MWH-4	34
19	TMI-43T	MWH-4	52

▼ Q. find the count of records with quantity <=25

```
1 ind = np.where(df['Quantity']<= 25)
2 # using conditions in np.where() over 'Quantity' column from DataFrame df
3 ind
```

```
(array([ 8, 16]),)
```

```
1 len(ind[0])
2 # count of indices which satisfy the condition
```

```
2
```

▼ DataFrame.drop(index)

```
1 mydf = df.drop(ind[0])
2 # creates a new Dataframe dropping the specifying the index
3 mydf.shape
4 # prints the shape of DataFrame
```

```
(48, 3)
```

```
1 df.iloc[ind]
```

	MaterialID	Location	Quantity
8	GCVB-79	MWH-2	10
16	SDRT-67	MWH-5	12

## ▼ DataFrame.head()

```
1 df.head()
```

	MaterialID	Location	Quantity
0	TMI-43T	MWH-4	34
1	AXCP-78	MWH-1	67
2	LXCV-21	MWH-2	27
3	AXCP-78	MWH-5	65
4	AXCP-78	MWH-4	36

## ▼ Q. add another column 'Status', and put value 'Scrap' if Qty&lt;=50 else put value 'Useful'

```
1 # (df.Quantity<=50)-> df['Status']=Scrap
2 # else df['Status']=useful
3
4 qty_lst = list(df['Quantity'])
5 #qty_lst
6 sts_lst = []
7 for a in qty_lst:
8     if a <=50:
9         sts_lst.append('Scrap')
10    else:
11        sts_lst.append("Useful")
12 # sts_lst
13 df['Status'] = sts_lst
14 df.head()
```

	MaterialID	Location	Quantity	Status
0	TMI-43T	MWH-4	34	Scrap
1	AXCP-78	MWH-1	67	Useful
2	LXCV-21	MWH-2	27	Scrap
3	AXCP-78	MWH-5	65	Useful
4	AXCP-78	MWH-4	36	Scrap

```
1 mydf = df[['Location', 'Status']]
2 # creating a new dataframe with only two columns from previous DataFrame
3 mydf.head()
4 # printing only top 5 records from DataFrame
```



	Location	Status
0	MWH-4	Scrap
1	MWH-1	Useful
2	MWH-2	Scrap

### ▼ DataFrame('ColumnName', axis=1)

```
1 mydf = df.drop('Location', axis=1)
2 # dropping a row/column based on axis
3 mydf.head()
4 # printing top 5 records from DataFrame
```

	MaterialID	Quantity	Status
0	TMI-43T	34	Scrap
1	AXCP-78	67	Useful
2	LXCV-21	27	Scrap
3	AXCP-78	65	Useful
4	AXCP-78	36	Scrap

### ▼ DataFrame.columns

```
1 mydf.columns
2 # prints an index of column names

Index(['MaterialID', 'Quantity', 'Status'], dtype='object')
```

### ▼ renaming column names in DataFrame

```
1 x = list(mydf.columns)
2 # converts index of headers in DataFrame
3 x
```

```
['MaterialID', 'Quantity', 'Status']
```

```
1 x[0] = 'Name'
2 # updating value of first column name in list for index of
3 # column Names in DataFrame
4 mydf.columns = x
5 # updating the column Name index with list containing updated list names
```

```
1 mydf.head()
2 # printing first five records to check if column names have been updated
```

	Name	Quantity	Status
0	TMI-43T	34	Scrap
1	AXCP-78	67	Useful
2	LXCV-21	27	Scrap
3	AXCP-78	65	Useful
4	AXCP-78	36	Scrap

## ▼ Data Cleaning

```
1 import numpy as np
2 import pandas as pd
3 # importing NumPy & pandas
```

```
1 from google.colab import files
2 uploaded=files.upload()
3 # emp_info.xlsx
4 # to be used with google colab
5
6 # os.chdir(r'D:/advanced-analytics-files/day03')
7 # to change current working directory to specified path
8 # to be used while running on local system
```

No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.  
Saving emp\_info.xlsx to emp\_info.xlsx

```
1 df = pd.read_excel('emp_info.xlsx')
2 # loading excel file into dataframe
3 df
```

	Name	EmpID	Deptt	Passport
0	Sudeep	A342	Quality	DF453278
1	Ruchika	C578	Sales	CA567657
2	NaN	J436	Admin	DF453278
3	Arjun	A746	NaN	RF453432
4	Rishi	R475	Prodn	DF453278
5	Chetan	NaN	NaN	ER534867
6	Deepak	NaN	Quality	RF453432

## ▼ DataFrame['ColumnName'].isnull()

```
1 ind = np.where(df['Deptt'].isnull())
2 # using np.where() with isnull() condition
3 ind
4 # printing indices which have null values
```

```
(array([3, 5]),)
```

```
1 # deptt is not missing / data filtered without null values
2 df1 = df.drop(ind[0], axis=0)
3 # creating a new DataFrame while dropping indices which have null values
4 df1
5 # printing new DataFrame that has removed null values
```

	Name	EmpID	Deptt	Passport
0	Sudeep	A342	Quality	DF453278
1	Ruchika	C578	Sales	CA567657
2	NaN	J436	Admin	DF453278
4	Rishi	R475	Prodn	DF453278
6	Deepak	NaN	Quality	RF453432

```
1 # deptt is missing / data filtered with null values
2 df2 = df.iloc[ind]
3 # creating a new DataFrame using indices which have null values
4 df2
5 # printing new DataFrame which contains null values
```

	Name	EmpID	Deptt	Passport
3	Arjun	A746	NaN	RF453432
5	Chetan	NaN	NaN	ER534867

▼ Q. select the rows where either empid is missing or deptt is missing

```
1 ind = np.where( (df['EmpID'].isnull()) | (df['Deptt'].isnull()) )
2 # creating a new DataFrame using where() specifying
3 # either EmpID or Deptt is null
4 ind
5 # printing indices which satisfies condition
```

```
(array([3, 5, 6]),)
```

```

1 # either empid is missing or deptt is missing
2 df2 = df.iloc[ind]
3 # creating a new DataFrame using indices in df.iloc which satisfies
4 # either EmpID or Deptt is null
5 df2
6 # printing Dataframe which satisfies condition

```

	Name	EmpID	Deptt	Passport
3	Arjun	A746	NaN	RF453432
5	Chetan	NaN	NaN	ER534867
6	Deepak	NaN	Quality	RF453432

#### ▼ Q. neither empid is missing nor deptt is missing

```

1 df1 = df.drop(ind[0], axis=0)
2 # dropping indices with null vlaue condition to create
3 # a new DataFrame without null values
4 df1
5 # printing Dataframe which satisfies condition for without null values

```

	Name	EmpID	Deptt	Passport
0	Sudeep	A342	Quality	DF453278
1	Ruchika	C578	Sales	CA567657
2	NaN	J436	Admin	DF453278
4	Rishi	R475	Prodn	DF453278

Q. select the rows where none of the values is missing

#### ▼ DataFrame.dropna()

```

1 mydf = df.dropna()
2 # DataFrame.dropna() removes all the records with any null values
3 mydf
4 # printing new DataFrame without any null values

```

	Name	EmpID	Deptt	Passport
0	Sudeep	A342	Quality	DF453278
1	Ruchika	C578	Sales	CA567657
4	Rishi	R475	Prodn	DF453278

```
1 df
2 # printing original DataFrame
```

	Name	EmpID	Deptt	Passport
0	Sudeep	A342	Quality	DF453278
1	Ruchika	C578	Sales	CA567657
2	NaN	J436	Admin	DF453278
3	Arjun	A746	NaN	RF453432
4	Rishi	R475	Prodn	DF453278
5	Chetan	NaN	NaN	ER534867
6	Deepak	NaN	Quality	RF453432

Q. remove duplicate values on Passport column

▼ DataFrame['columnName'].duplicated(keep=False)

```
1 ind = np.where(df['Passport'].duplicated(keep=False))
2 # create index of column 'Passport' from DataFrame, without duplicate values
3 # .duplicated(keep=False) will mark all the duplicates as True
4 ind
5 # printing indices which satisfy the condition
```

```
(array([0, 2, 3, 4, 6]),)
```

```
1 df1 = df.drop(ind[0], axis=0)
2 # unique values on Passport column
3 # using index of column 'Passport' from DataFrame to drop records without
4 # duplicate values or unique values
5 df1
6 # printing new DataFrame which contains unique values
```

	Name	EmpID	Deptt	Passport
1	Ruchika	C578	Sales	CA567657
5	Chetan	NaN	NaN	ER534867

```
1 # non-unique values on Passport column
2 df2 = df.iloc[ind]
3 # creating a new DataFrame with iloc using indices with duplicate entries
4 df2
5 # printing the new DataFrame with duplicate entries
```

	Name	EmpID	Deptt	Passport
0	Sudeep	A342	Quality	DF453278
2	NaN	J436	Admin	DF453278
3	Arjun	A746	NaN	RF453432
4	Rishi	R475	Prodn	DF453278

### ▼ DataFrame[columnName].duplicated(keep='last')

```
1 ind = np.where(df['Passport'].duplicated(keep='last'))
2 ind
3 # keep='last' marks all duplicates, but not the last one
4 # considers last duplicate occurrence as correct
5 # while other duplicate values are removed
```

```
(array([0, 2, 3]),)
```

```
1 keep_last_mydf = df.drop(ind[0], axis=0)
2 #
3 keep_last_mydf
4 #
```

	Name	EmpID	Deptt	Passport
1	Ruchika	C578	Sales	CA567657
4	Rishi	R475	Prodn	DF453278
5	Chetan	NaN	NaN	ER534867
6	Deepak	NaN	Quality	RF453432

### ▼ DataFrame[columnName].duplicated(keep='first')

```
1 ind = np.where(df['Passport'].duplicated(keep='first'))
2 ind
3 # keep='first' considers last duplicate occurrence is considered correct
4 # while other duplicate values are removed
```

```
(array([2, 4, 6]),)
```

### ▼ DataFrame.drop(index, axis=0)

```
1 keep_first_mydf = df.drop(ind[0], axis=0)
2 # dropping indices of duplicates keeping 'first' as original value
3 # but all other as duplicate values
4 keep_first_mydf
5 # printing DataFrame with duplicates as 'first' as removed
```

	Name	EmpID	Deptt	Passport
1	Ruchika	C578	Sales	CA567657
4	Rishi	R475	Prodn	DF453278
5	Chetan	NaN	NaN	ER534867
6	Deepak	NaN	Quality	RF453432

## ▼ DataFrame GroupBy

```
1 import numpy as np
2 import pandas as pd
3 # importing NumPy & pandas
```

```
1 from google.colab import files
2 uploaded=files.upload()
3 # ERPData.xlsx
4 # to be used with google colab
5
6 # os.chdir(r'D:/advanced-analytics-files/day03')
7 # to change current working directory to specified path
8 # to be used while running on local system
```

No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.  
Saving ERPData.xlsx to ERPData.xlsx

```
1 df = pd.read_excel('ERPDData.xlsx')
2 # using pd.read_excel('filename') to load excel file to DataFrame
3 df.head()
4 # printing first 5 records from DataFrame
```

	MaterialID	Location	Quantity
0	TMI-43T	MWH-4	34
1	AXCP-78	MWH-1	67
2	LXCV-21	MWH-2	27
3	AXCP-78	MWH-5	65
4	AXCP-78	MWH-4	36

## ▼ DataFrame.groupby('ColumnName')

- groupby: Grouped information based on columns
- groups created can be used to perform aggregate functions

```
1 grp1 = df.groupby('Location')
2 # creates a DataFrame GroupBy on Location column
3 grp1
4 # prints object information for DataFrame GroupBy
```

<pandas.core.groupby.generic.DataFrameGroupBy object at 0x7f0b49ad8eb0>

### ▼ DataFrameGroupBy.groups

```
1 grp1.groups
2 # gruoups property lsts down all the groups in a DataFrame Group
```

```
{'MWH-1': [1, 9, 10, 24, 25, 26, 39, 40, 41, 48, 49], 'MWH-2': [2, 7, 8, 15, 35, 36, 42, 45, 46], 'MWH-3': [14, 43, 44, 47], 'MWH-4': [0, 4, 5, 6, 11, 12, 13, 18, 19, 20, 21, 22, 30, 31, 32],
'MWH-5': [3, 16, 17, 23, 27, 28, 29, 33, 34, 37, 38]}
```

```
1 len(grp1.groups)
2 # prints count of unique groups in the column where GroupBy is created
```

5

```
1 type(grp1)
2 # type of groupby object is DataFrameGroupBy
```

pandas.core.groupby.generic.DataFrameGroupBy

### ▼ DataFrameGroupBy.get\_group('UniqueValue')

```
1 grp1.get_group('MWH-2')
2 # uses DataFrameGroupBy.get_group property to fetch records
3 # for a specified unique value / group only
```

	MaterialID	Location	Quantity
2	LXCV-21	MWH-2	27
7	TMI-43T	MWH-2	29
8	GCVB-79	MWH-2	10
15	GCVB-79	MWH-2	87
35	GCVB-79	MWH-2	31
36	GCVB-79	MWH-2	28
42	SDRT-67	MWH-2	31
45	DDBN-89	MWH-2	69
46	AXCP-78	MWH-2	85

### ▼ DataFrameGroupBy.agg(AggregateFunction)



- works with DataFrameGroupBy, but returns a DataFrame

```
1 grp1.agg(np.mean)
2 # can specify the aggregate function to be used with DataFrame GroupBy
3 # printing the mean of each of the groups created on location column
```

```
<ipython-input-9-1c597bc92aa3>:1: FutureWarning: The operation <function mean at 0x7f0b899
grp1.agg(np.mean)
```

	Quantity
Location	
MWH-1	80.818182
MWH-2	44.111111
MWH-3	43.500000
MWH-4	48.533333
MWH-5	63.363636

```
1 grp1.agg(np.size)
2 # printing the size/length of each of the groups created on location column
3
4 # NOTE that, now it prints size for MaterialID column too,
5 # because it can count all type of values,
6 # but for numeric aggregate functions like mean, etc. ,
7 # it takes up column which have numeric values only
```

	MaterialID	Quantity
Location		
MWH-1	11	11
MWH-2	9	9
MWH-3	4	4
MWH-4	15	15
MWH-5	11	11

```
1 grp1.agg(np.sum)
2 # printing the sum of each of the groups created on location column
3
4 # Here again, it'll take columns which have numeric values only
5 # for numeric functions like sum, mean
```

```
<ipython-input-11-70f497a5407a>:1: FutureWarning: The operation <function sum at 0x7f0b899
grp1.agg(np.sum)
      Quantity
```

Location

MWH-1 889

MWH-2 397

### ▼ DataFrameGroupBy.agg([ListOfAggregateFunctions])

```
MWH-5 697
1 grp1.agg([np.sum, np.mean])
2 # printing the sum & mean of each of the groups created on location column
3 # but it'll work on numeric columns only
```

```
<ipython-input-12-61ce2952447e>:1: FutureWarning: ['MaterialID'] did not aggregate success
grp1.agg([np.sum, np.mean])
```

Quantity

sum mean

Location

MWH-1 889 80.818182

MWH-2 397 44.111111

MWH-3 174 43.500000

MWH-4 728 48.533333

MWH-5 697 63.363636

```
1 r = grp1.agg([np.sum, np.mean])
2 # returns a DataFrame from DataFrameGroupBy using aggregate functions
3 r
4 # printing DataFrame r
```

```
FutureWarning: ['MaterialID'] did not aggregate successfully
```

```
1 type(r)
2 # checking type of DataFrame containing aggregate information
3 # returned by DataFrameGroupBy.agg()
```

```
pandas.core.frame.DataFrame
```

```
1 ind = np.where(r.index == 'MWH-3')
2 # using np.where() to find the indices of records of
3 # DataFrame which has index = 'MWH-3'
4 ind
5 # printing indices of records satisfying the condition
```

```
(array([2]),)
```

```
1 r.iloc[ind]
2 # using iloc on DataFrame to print the records on indices
3 # satisfying the condition
```

	Quantity	
	sum	mean
Location		
MWH-3	174	43.5

```
1 grp2 = df.groupby(['MaterialID', 'Location'])
2 # creating another group using DataFrame.groupby() to make
3 # groups by 'MaterialID' column, and then sub-grouping by 'Location' column
4 grp2.groups
```

```
{('AXCP-78', 'MWH-1'): [1, 24, 25, 39, 48], ('AXCP-78', 'MWH-2'): [46], ('AXCP-78', 'MWH-3'): [47], ('AXCP-78', 'MWH-4'): [4], ('AXCP-78', 'MWH-5'): [3, 37, 38], ('DDBN-89', 'MWH-1'): [26, 40, 49], ('DDBN-89', 'MWH-2'): [45], ('DDBN-89', 'MWH-3'): [43], ('DDBN-89', 'MWH-4'): [22, 30], ('DDBN-89', 'MWH-5'): [23], ('GCVB-79', 'MWH-2'): [8, 15, 35, 36], ('GCVB-79', 'MWH-5'): [34], ('LXCV-21', 'MWH-2'): [2], ('LXCV-21', 'MWH-3'): [44], ('LXCV-21', 'MWH-5'): [27, 28, 29], ('SDRT-67', 'MWH-1'): [10, 41], ('SDRT-67', 'MWH-2'): [42], ('SDRT-67', 'MWH-4'): [11], ('SDRT-67', 'MWH-5'): [16, 17], ('TMI-43T', 'MWH-1'): [9], ('TMI-43T', 'MWH-2'): [7], ('TMI-43T', 'MWH-3'): [14], ('TMI-43T', 'MWH-4'): [0, 5, 6, 12, 13, 18, 19, 20, 21, 31, 32], ('TMI-43T', 'MWH-5'): [33]}
```

```
1 grp3 = df.groupby(['Location', 'MaterialID'])
2 # creating another group using DataFrame.groupby() to make
3 # groups by 'Location' column, and then sub-grouping by 'MaterialID' column
4 grp3.groups
```

```
{('MWH-1', 'AXCP-78'): [1, 24, 25, 39, 48], ('MWH-1', 'DDBN-89'): [26, 40, 49], ('MWH-1', 'SDRT-67'): [10, 41], ('MWH-1', 'TMI-43T'): [9], ('MWH-2', 'AXCP-78'): [46], ('MWH-2', 'DDBN-89'): [45], ('MWH-2', 'GCVB-79'): [8, 15, 35, 36], ('MWH-2', 'LXCV-21'): [2], ('MWH-2', 'SDRT-67'): [42], ('MWH-2', 'TMI-43T'): [7], ('MWH-3', 'AXCP-78'): [47], ('MWH-3', 'DDBN-89'): [43], ('MWH-3', 'LXCV-21'): [44], ('MWH-3', 'TMI-43T'): [14], ('MWH-4', 'AXCP-78'): [4], ('MWH-4', 'DDBN-89'): [22, 30], ('MWH-4', 'SDRT-67'): [11], ('MWH-4', 'TMI-43T'): [0, 5, 6, 12, 13, 18, 19, 20, 21, 31, 32], ('MWH-5', 'AXCP-78'): [3, 37, 38], ('MWH-5', 'DDBN-89'): [23], ('MWH-5', 'GCVB-79'): [34], ('MWH-5', 'LXCV-21'): [27, 28, 29], ('MWH-5', 'SDRT-67'): [16, 17], ('MWH-5', 'TMI-43T'): [33]}
```

```
1 grp2.agg([np.sum, np.size])
2 # using aggregate functions on DataFrameGroupBy which is grouped
3 # on two columns of DataFrame,
4 # so there will be sub-groups for groups,
5 # and then it'll execute aggregate functions for each of the sub-groups
6
7 # Note that, here the order of grouping is 'MaterialID' then 'Location'
```

```
8 # groups are arranged alphabetically, so access time is reduced
9 # and DataFrame becomes optimized
```

MaterialID	Location	Quantity	
		sum	size
AXCP-78	MWH-1	548	5
	MWH-2	85	1
	MWH-3	29	1
	MWH-4	36	1
	MWH-5	200	3
DDBN-89	MWH-1	130	3
	MWH-2	69	1
	MWH-3	38	1
	MWH-4	75	2
	MWH-5	71	1
GCVB-79	MWH-2	156	4
	MWH-5	26	1
LXCV-21	MWH-2	27	1
	MWH-3	42	1
	MWH-5	174	3
SDRT-67	MWH-1	91	2
	MWH-2	31	1
	MWH-4	58	1
	MWH-5	114	2
TMI-43T	MWH-1	120	1
	MWH-2	29	1
	MWH-3	65	1
	MWH-4	559	11
	MWH-5	112	1

```
1 grp3.agg([np.sum, np.mean])
2 # Note that, here the order of grouping is 'Location' then 'MaterialID'
```

Location	MaterialID	Quantity	
		sum	mean
MWH-1	AXCP-78	548	109.600000
	DDBN-89	130	43.333333
	SDRT-67	91	45.500000
	TMI-43T	120	120.000000
MWH-2	AXCP-78	85	85.000000
	DDBN-89	69	69.000000
	GCVB-79	156	39.000000
	LXCV-21	27	27.000000
	SDRT-67	31	31.000000
	TMI-43T	29	29.000000
MWH-3	AXCP-78	29	29.000000
	DDBN-89	38	38.000000
	LXCV-21	42	42.000000
	TMI-43T	65	65.000000
MWH-4	AXCP-78	36	36.000000
	DDBN-89	75	37.500000
	SDRT-67	58	58.000000
	TMI-43T	559	50.818182
MWH-5	AXCP-78	200	66.666667
	DDBN-89	71	71.000000
	GCVB-79	26	26.000000
	LXCV-21	174	58.000000

```
1 df.head()
2 # first 5 records from original DataFrame
```

	MaterialID	Location	Quantity
0	TMI-43T	MWH-4	34
1	AXCP-78	MWH-1	67
2	LXCV-21	MWH-2	27
3	AXCP-78	MWH-5	65
4	AXCP-78	MWH-4	36

## ▼ DataFrame.sort\_values('columnName')

```
1 mydf = df.sort_values('Quantity')
2 # creating a new DataFrame while sorting on
3 # column 'Quantity' from original DataFrame
```

## ▼ DataFrame.head(n)

- DataFrame.head() : by default prints first 5 records
- DataFrame.head(n) : prints first n records

```
1 mydf.head()
2 # printing first 5 records from DataFrame
```

	MaterialID	Location	Quantity
8	GCVB-79	MWH-2	10
16	SDRT-67	MWH-5	12
34	GCVB-79	MWH-5	26
2	LXCV-21	MWH-2	27
30	DDBN-89	MWH-4	27

## ▼ DataFrame.tail(n)

- DataFrame.tail() : by default prints last 5 records
- DataFrame.tail(n) : prints last n records

```
1 mydf.tail()
2 # printing last 5 records from DataFrame
```

	MaterialID	Location	Quantity
25	AXCP-78	MWH-1	109
33	TMI-43T	MWH-5	112
9	TMI-43T	MWH-1	120
39	AXCP-78	MWH-1	145
24	AXCP-78	MWH-1	152

## ▼ DataFrame.sort\_values([listOfColumns])

```
1 mydf = df.sort_values(['Location', 'MaterialID', 'Quantity'])
2 # creating a new Dataframe with sorting original DataFrame on three columns
3 mydf
4 # printing sorted DataFrame
```

MaterialID	Location	Quantity
39	AXCP-78	MWH-1
145		

## ▼ Pandas merge

- used to merge two DataFrames
- acts as SQL Join from

MaterialID	Location	Quantity
39	AXCP-78	MWH-1
145		

## ▼ pd.merge with 'Pune' & 'Mumbai'

MaterialID	Location	Quantity
39	AXCP-78	MWH-1
145		

```
1 import numpy as np
2 import pandas as pd
```

```
1 from google.colab import files
2 uploaded=files.upload()
3 # locations.xlsx
4 # to be used with google colab
5
6 # os.chdir(r'D:/advanced-analytics-files/day03')
7 # to change current working directory to specified path
8 # to be used while running on local system
```

No file chosen      Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.  
Saving locations.xlsx to locations.xlsx

MaterialID	Location	Quantity
39	AXCP-78	MWH-1
145		

## ▼ pd.read\_excel('filename.xlsx', sheet\_name='WorkSheet')

MaterialID	Location	Quantity
39	AXCP-78	MWH-1
145		

```
1 df1 = pd.read_excel('locations.xlsx', sheet_name='Pune')
2 # taking workSheet 'Pune' from WorkBook 'locations.xlsx' into DataFrame
3 df2 = pd.read_excel('locations.xlsx', sheet_name='Mumbai')
4 # taking workSheet 'Mumbai' from WorkBook 'locations.xlsx' into DataFrame
```

```
1 df1
2 # printing DataFrame df1
```

	Name	Subject	Grade
0	Rakesh	Python	A
1	Manoj	MLPython	C
2	Vaibhav	Statistics	B
3	Hitesh	CommSkills	A
4	Suyash	ProjectMgmt	B

MaterialID	Location	Quantity
39	AXCP-78	MWH-1
145		

```
1 df2
2 # printing DataFrame df2
```



	Name	Subject	Grade
0	Vaibhav	Six Sigma	A
1	Deepika	Statistics	B
2	Arjun	CommSkills	A
3	Chetan	Python	A
4	Abhishek	MLPython	B

### ▼ `pd.merge(DataFrame1, DataFrame2, on='CommonColumnName')`

- produces inner join of DataFrame1 & DataFrame2 on common column

```
1 pd.merge(df1, df2, on='Subject')
2 # returns resultant DataFrame which joins df1 with df2 on column subject
3 # produces inner join by default
4
```

	Name_x	Subject	Grade_x	Name_y	Grade_y
0	Rakesh	Python	A	Chetan	A
1	Manoj	MLPython	C	Abhishek	B
2	Vaibhav	Statistics	B	Deepika	B
3	Hitesh	CommSkills	A	Arjun	A

### ▼ `pd.merge(DataFrame1, DataFrame2, on='CommonColumnName', how='JoinType')`

```
1 pd.merge(df1, df2, on='Subject', how='inner')
2 # returns resultant DataFrame which joins df1 with df2
3 # on column 'subject' by 'inner join'
4
5 # performs intersection of two DataFrames
```

	Name_x	Subject	Grade_x	Name_y	Grade_y
0	Rakesh	Python	A	Chetan	A
1	Manoj	MLPython	C	Abhishek	B
2	Vaibhav	Statistics	B	Deepika	B
3	Hitesh	CommSkills	A	Arjun	A

```

1 pd.merge(df1, df2, on='Subject', how='outer')
2 # returns resultant DataFrame which joins df1 with df2
3 # on column 'subject' by 'outer join'
4
5 # performs union of two DataFrames

```

	Name_x	Subject	Grade_x	Name_y	Grade_y
0	Rakesh	Python	A	Chetan	A
1	Manoj	MLPython	C	Abhishek	B
2	Vaibhav	Statistics	B	Deepika	B
3	Hitesh	CommSkills	A	Arjun	A
4	Suyash	ProjectMgmt	B	NaN	NaN
5	NaN	Six Sigma	NaN	Vaibhav	A

```

1 pd.merge(df1, df2, on='Subject', how='left')
2 # returns resultant DataFrame which joins df1 with df2
3 # on column 'subject' by 'left join'

```

	Name_x	Subject	Grade_x	Name_y	Grade_y
0	Rakesh	Python	A	Chetan	A
1	Manoj	MLPython	C	Abhishek	B
2	Vaibhav	Statistics	B	Deepika	B
3	Hitesh	CommSkills	A	Arjun	A
4	Suyash	ProjectMgmt	B	NaN	NaN

```

1 pd.merge(df1, df2, on='Subject', how='right')
2 # returns resultant DataFrame which maps df1 with df2
3 # on column 'subject' by 'right join'

```

	Name_x	Subject	Grade_x	Name_y	Grade_y
0	NaN	Six Sigma	NaN	Vaibhav	A
1	Vaibhav	Statistics	B	Deepika	B
2	Hitesh	CommSkills	A	Arjun	A
3	Rakesh	Python	A	Chetan	A
4	Manoj	MLPython	C	Abhishek	B

```

1 df1
2 # printing actual DataFrame df1

```

	Name	Subject	Grade
0	Rakesh	Python	A
1	Manoj	MLPython	C
2	Vaibhav	Statistics	B
3	Hitesh	CommSkills	A

```
1 df2
2 # printing actual DataFrame df2
```

	Name	Subject	Grade
0	Vaibhav	Six Sigma	A
1	Deepika	Statistics	B
2	Arjun	CommSkills	A
3	Chetan	Python	A
4	Abhishek	MLPython	B

▼ `pd.merge(DataFrame1, DataFrame2, on=[listOfColumns], how='JoinType')`

```
1 pd.merge(df1, df2, on=['Subject', 'Grade'], how='inner')
2 # returns resultant DataFrame which joins df1 with df2
3 # on two columns 'Subject' & 'Grade' by 'inner join'
```

	Name_x	Subject	Grade	Name_y
0	Rakesh	Python	A	Chetan
1	Vaibhav	Statistics	B	Deepika
2	Hitesh	CommSkills	A	Arjun

```
1 pd.merge(df1, df2, on=['Subject', 'Grade'], how='outer')
2 # returns resultant DataFrame which joins df1 with df2
3 # on two columns 'Subject' & 'Grade' by 'outer join'
```

	Name_x	Subject	Grade	Name_y
0	Rakesh	Python	A	Chetan
1	Manoj	MLPython	C	NaN
2	Vaibhav	Statistics	B	Deepika
3	Hitesh	CommSkills	A	Arjun
4	Suyash	ProjectMgmt	B	NaN
5	NaN	Six Sigma	A	Vaibhav
6	NaN	MLPython	B	Abhishek

```
1 pd.merge(df1, df2, on=['Subject', 'Grade'], how='left')
2 # returns resultant DataFrame which joins df1 with df2
3 # on two columns 'Subject' & 'Grade' by 'left join'
```

	Name_x	Subject	Grade	Name_y
0	Rakesh	Python	A	Chetan
1	Manoj	MLPython	C	NaN
2	Vaibhav	Statistics	B	Deepika
3	Hitesh	CommSkills	A	Arjun
4	Suyash	ProjectMgmt	B	NaN

```
1 pd.merge(df1, df2, on=['Subject', 'Grade'], how='right')
2 # returns resultant DataFrame which joins df1 with df2
3 # on two columns 'Subject' & 'Grade' by 'right join'
```

	Name_x	Subject	Grade	Name_y
0	NaN	Six Sigma	A	Vaibhav
1	Vaibhav	Statistics	B	Deepika
2	Hitesh	CommSkills	A	Arjun
3	Rakesh	Python	A	Chetan
4	NaN	MLPython	B	Abhishek

#### ▼ pd.merge with 'Salary' & 'Deptt'

```
1 import numpy as np
2 import pandas as pd
```

```
1 from google.colab import files
2 uploaded=files.upload()
3 # locations.xlsx
4 # to be used with google colab
5
6 # os.chdir(r'D:/advanced-analytics-files/day03')
7 # to change current working directory to specified path
8 # to be used while running on local system
```

No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.  
Saving locations.xlsx to locations (1).xlsx

```
1 df1 = pd.read_excel('locations.xlsx', sheet_name='Salary')
2 # taking workSheet 'Salary' from WorkBook 'locations.xlsx' into DataFrame
```

```
3 df2 = pd.read_excel('locations.xlsx', sheet_name='Deptt')
4 # taking workSheet 'Deptt' from WorkBook 'locations.xlsx' into DataFrame
```

```
1 df1
2 # printing DataFrame df1
```

	Name	EmpID	Salary
0	Sudeep	A342	56
1	Deepika	J436	45
2	Chetan	B435	76
3	Abhishek	C234	47
4	Rishi	R475	38

```
1 df2
2 # printing DataFrame df2
```

	Name	EmpID	Deptt
0	Sudeep	A342	Quality
1	Ruchika	C578	Sales
2	Deepika	J436	Admin
3	Arjun	A342	Procurement
4	Rishi	R475	Prodn

```
1 pd.merge(df1, df2)
2 # same as inner merge / join
```

	Name	EmpID	Salary	Deptt
0	Sudeep	A342	56	Quality
1	Deepika	J436	45	Admin
2	Rishi	R475	38	Prodn

```
1 pd.merge(df1, df2, how='inner')
2 # inner merge/join acts like intersection
```

	Name	EmpID	Salary	Deptt
0	Sudeep	A342	56	Quality
1	Deepika	J436	45	Admin
2	Rishi	R475	38	Prodn

```
1 pd.merge(df1, df2, how='outer')
2 # outer merge/join acts like union
```

	Name	EmpID	Salary	Deptt
0	Sudeep	A342	56.0	Quality
1	Deepika	J436	45.0	Admin
2	Chetan	B435	76.0	NaN
3	Abhishek	C234	47.0	NaN
4	Rishi	R475	38.0	Prodn
5	Ruchika	C578	NaN	Sales
6	Arjun	A342	NaN	Procurement

```
1 df1
2 # printing DataFrame df1
```

	Name	EmpID	Salary
0	Sudeep	A342	56
1	Deepika	J436	45
2	Chetan	B435	76
3	Abhishek	C234	47
4	Rishi	R475	38

```
1 df2
2 # printing DataFrame df2
```

	Name	EmpID	Deptt
0	Sudeep	A342	Quality
1	Ruchika	C578	Sales
2	Deepika	J436	Admin
3	Arjun	A342	Procurement
4	Rishi	R475	Prodn

```
1 pd.merge(df1, df2, how='left')
2 # Note that, indexing based on 'left' DataFrame in case of left merge / join
```

	Name	EmpID	Salary	Deptt
0	Sudeep	A342	56	Quality

4      Deepika      J436      45      Admin

```
1 pd.merge(df1, df2, how='right')
```

```
2 # Note that, indexing based on 'right' DataFrame in case of right merge / join
```

	Name	EmpID	Salary	Deptt
0	Sudeep	A342	56.0	Quality
1	Ruchika	C578	NaN	Sales
2	Deepika	J436	45.0	Admin
3	Arjun	A342	NaN	Procurement
4	Rishi	R475	38.0	Prodn

### ▼ pd.merge with 'inventory' & 'Info'

```
1 import numpy as np
2 import pandas as pd
```

```
1 from google.colab import files
2 uploaded=files.upload()
3 # MaterialInfo.xlsx
4 # to be used with google colab
5
6 # os.chdir(r'D:/advanced-analytics-files/day03')
7 # to change current working directory to specified path
8 # to be used while running on local system
```

No file chosen

Upload widget is only available when the cell has been executed in

the current browser session. Please rerun this cell to enable.

Saving MaterialInfo.xlsx to MaterialInfo.xlsx

```
1 df1 = pd.read_excel('MaterialInfo.xlsx', sheet_name='inventory')
2 # taking workSheet 'inventory' from WorkBook 'MaterialInfo.xlsx' into DataFrame
3 df2 = pd.read_excel('MaterialInfo.xlsx', sheet_name='Info')
4 # taking workSheet 'Info' from WorkBook 'MaterialInfo.xlsx' into DataFrame
```

```
1 df1
2 # printing DataFrame df1
```

	MaterialID	Quantity
0	A	34
1	B	23
2	A	16
3	A	52
4	B	34
5	C	41
6	B	27
7	C	63

```
1 df2
2 # printing DataFrame df2
```

	MaterialID	Rate
0	A	1200
1	B	1500
2	C	2400

```
1 pd.merge(df1, df2)
2 # sorting will take place automatically as per first column MaterialID
```

	MaterialID	Quantity	Rate
0	A	34	1200
1	A	16	1200
2	A	52	1200
3	A	45	1200
4	A	38	1200
5	B	23	1500
6	B	34	1500
7	B	27	1500
8	B	36	1500
9	C	41	2400
10	C	63	2400
11	C	31	2400

```
1 pd.merge(df1, df2, how='left')
2 # indexing / sequencing as per left DataFrame
```



	MaterialID	Quantity	Rate
0	A	34	1200
1	B	23	1500
2	A	16	1200
3	A	52	1200
4	B	34	1500
5	C	41	2400
6	B	27	1500
7	C	63	2400
8	A	45	1200
9	B	36	1500
10	C	31	2400
11	A	38	1200

```
1 pd.merge(df1, df2, how='right')
2 # indexing / sequencing as per right DataFrame
```

	MaterialID	Quantity	Rate
0	A	34	1200
1	A	16	1200
2	A	52	1200
3	A	45	1200
4	A	38	1200
5	B	23	1500
6	B	34	1500
7	B	27	1500
8	B	36	1500
9	C	41	2400
10	C	63	2400
11	C	31	2400

1

[Colab paid products](#) - [Cancel contracts here](#)

