

Chapter 1
Data Handling
using Pandas

Informatics Practices
Class XII (As per
CBSE Board)

Visit : python.mykvs.in for regular updates

Data Handling using Pandas

Python Library – Pandas

It is a most famous Python package for data science, which offers powerful and flexible data structures that make data analysis and manipulation easy. Pandas makes data importing and data analyzing much easier. Pandas builds on packages like [NumPy](#) and [matplotlib](#) to give us a single & convenient place for data analysis and visualization work.



Visit : python.mykvs.in for regular updates

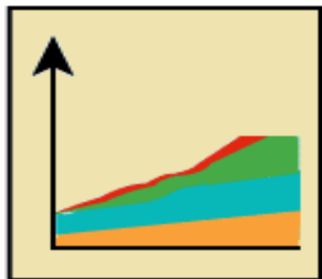
Data Handling using Pandas

Python Library – Matplotlib

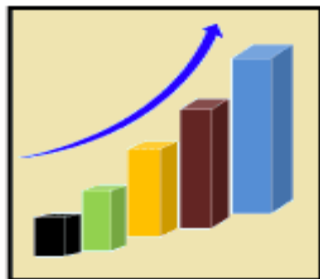
Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python. It is used to create

1. Develop publication quality plots with just a few lines of code
2. Use interactive figures that can zoom, pan, update...

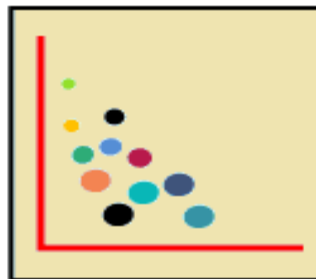
We can customize and Take full control of line styles, font properties, axes properties... as well as export and embed to a number of file formats and interactive environments



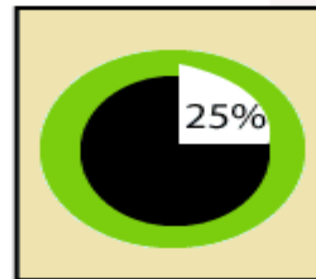
Area plot



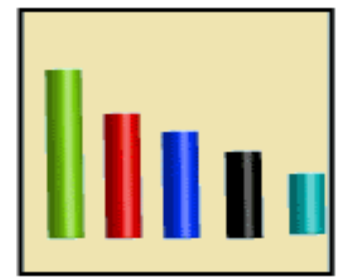
Histogram



Scatter plot



Pie Plot



Bar Graph

Visit : python.mykvs.in for regular updates

Data Handling using Pandas

Basic Features of Pandas

help a lot in keeping track of our data.

2. With a pandas dataframe, we can have different data types (float, int, string, datetime, etc) all in one place
3. Pandas has built in functionality for like easy grouping & easy joins of data, rolling windows
4. Good IO capabilities; Easily pull data from a MySQL database directly into a data frame
5. With pandas, you can use patsy for R-style syntax in doing regressions.
6. Tools for loading data into in-memory data objects from different file formats.
7. Data alignment and integrated handling of missing data.
8. Reshaping and pivoting of data sets.
9. Label-based slicing, indexing and subsetting of large data sets.

Data Handling using Pandas

Pandas – Installation/Environment Setup

Pandas module doesn't come bundled with Standard Python.

1

If we install Anaconda Python package Pandas will be installed by default.

Steps for Anaconda installation & Use

1. visit the site <https://www.anaconda.com/download/>
2. Download appropriate anaconda installer
3. After download install it.
4. During installation check for set path and all user
5. After installation start spyder utility of anaconda from start menu
6. Type `import pandas as pd` in left pane(temp.py)
7. Then run it.
8. If no error is show then it shows pandas is installed.
9. Like default temp.py we can create another .py file from new window option of file menu for new program.

Data Handling using Pandas

Pandas – Installation/Environment Setup

2

Pandas installation can be done in Standard Python distribution, using following steps.

1. There must be service pack installed on our computer if we are using windows. If it is not installed then we will not be able to install pandas in existing Standard Python (which is already installed). So install it first (google it).
2. We can check it through properties option of my computer icon.



3. Now install latest version (any one above 3.4) of python.

Visit : python.mykvs.in for regular updates

2

4. Now move to script folder of python distribution in command prompt (through cmd command of windows).

5. Execute following commands in command prompt serially.

```
>pip install numpy
```

```
>pip install six
```

```
>pip install pandas
```

Wait after each command for installation

Now we will be able to use pandas in standard python distribution.

6. Type `import pandas as pd` in python (IDLE) shell.

7. If it executed without error(it means pandas is installed on your system)

Data Handling using Pandas

Data Structures in Pandas

Two important data structures of pandas are–Series, DataFrame

1. Series

Series is like a one-dimensional array like structure with homogeneous data. For example, the following series is a collection of integers.

78	45	12	89	56
----	----	----	----	----

Basic feature of series are

- ❖ Homogeneous data
- ❖ Size Immutable
- ❖ Values of Data Mutable

2. DataFrame

DataFrame is like a two-dimensional array with heterogeneous data.

SR. No.	Admn No	Student Name	Class	Section	Gender	Date Of Birth
1	001284	NIDHI MANDAL	I	A	Girl	07/08/2010
2	001285	SOUMYADIP BHATTACHARYA	I	A	Boy	24/02/2011
3	001286	SHREYAANG SHANDILYA	I	A	Boy	29/12/2010

Basic feature of DataFrame are

- ❖ Heterogeneous data
- ❖ Size Mutable
- ❖ Data Mutable

Pandas Series

It is like one-dimensional array capable of holding data of any type (integer, string, float, python objects, etc.). Series can be created using constructor.

Syntax :- **pandas.Series(data, index, dtype, copy)**

Creation of Series is also possible from – ndarray, dictionary, scalar value.

Series can be created using

1. Array
2. Dict
3. Scalar value or constant

Pandas Series

Create an Empty Series

e.g.

```
import pandas as pseries
s = pseries.Series()
print(s)
```

Output

```
Series([], dtype: float64)
```

Data Handling using Pandas

Pandas Series

Create a Series from ndarray

Without index

e.g.

```
import pandas as pd1
import numpy as np1
data = np1.array(['a','b','c','d'])
s = pd1.Series(data)
print(s)
```

Output

```
1  a
2  b
3  c
4  d
```

dtype: object

Note : default index is starting from 0

With index position

e.g.

```
import pandas as p1
import numpy as np1
data = np1.array(['a','b','c','d'])
s = p1.Series(data,index=[100,101,102,103])
print(s)
```

Output

```
100  a
101  b
102  c
103d  dtype:
object
```

Note : index is starting from 100

Data Handling using Pandas

Pandas Series

Create a Series from dict

Eg.1(without index)

```
import pandas as pd1
import numpy as np1
data = {'a' : 0., 'b' : 1., 'c' : 2.}
s = pd1.Series(data)
print(s)
```

Output

```
a    0.0
b    1.0
c    2.0
dtype: float64
```

Eg.2 (with index)

```
import pandas as pd1
import numpy as np1
data = {'a' : 0., 'b' : 1., 'c' : 2.}
s = pd1.Series(data,index=['b','c','d','a'])
print(s)
```

Output

```
b    1.0
c    2.0
d    NaN
a    0.0
dtype: float64
```


Create a Series from Scalar

e.g

```
import pandas as pd1
import numpy as np1
s = pd1.Series(5, index=[0, 1, 2, 3])
print(s)
```

Output

```
0    5
1    5
2    5
3    5
```

dtype: int64

Note :- here 5 is repeated for 4 times (as per no of index)

Pandas Series

Maths operations with Series

e.g.

```
import pandas as pd
```

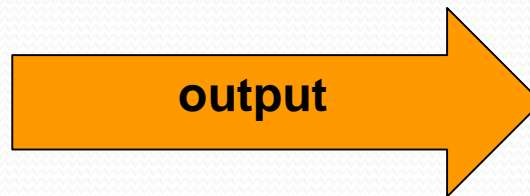
```
s = pd.Series([1,2,3])
```

```
t = pd.Series([1,2,4])
```

```
u=s+t #addition operation print (u)
```

```
u=s*t # multiplication operation
```

```
print (u)
```



```
0  2
1  4
2  7
dtype: int64
```

```
0  1
1  4
2  12
dtype: int64
```

Pandas Series

Head function

e.g

```
import pandas as pd1
s = pd1.Series([1,2,3,4,5],index = ['a','b','c','d','e'])
print (s.head(3))
```

Output

a 1

b. 2

c. 3

dtype: int64

Return first 3 elements

Pandas Series

tail function

e.g

```
import pandas as pd1
s = pd1.Series([1,2,3,4,5],index = ['a','b','c','d','e'])
print (s.tail(3))
```

Output

c 3

d. 4

e. 5

dtype: int64

Return last 3 elements

Accessing Data from Series with indexing and slicing

e.g.

```
import pandas as pd1
s = pd1.Series([1,2,3,4,5],index = ['a','b','c','d','e'])
print (s[0])# for 0 index position
print (s[:3]) #for first 3 index values
print (s[-3:]) #slicing for last 3 index values
```

Output

```
1
a. 1
b. 2
c. 3
dtype: int64 c      3
d. 4
e. 5
dtype: int64
```


Pandas Series

Retrieve Data Using Label as (Index)

e.g.

```
import pandas as pd1
s = pd1.Series([1,2,3,4,5],index = ['a','b','c','d','e'])
print (s[['c','d']])
```

Output c

3

d 4

dtype: int64

Pandas Series

Retrieve Data from selection

There are three methods for data selection:

- loc gets rows (or columns) with particular labels from the index.
- iloc gets rows (or columns) at particular positions in the index (so it only takes integers).
- ix usually tries to behave like loc but falls back to behaving like iloc if a label is not present in the index.

ix is deprecated and the use of loc and iloc is encouraged instead

Pandas Series

Retrieve Data from selection

e.g.

```
>>> s = pd.Series(np.nan,  
index=[49,48,47,46,45, 1, 2, 3, 4, 5])
```

```
>>> s.iloc[:3] # slice the first three rows
```

```
49 NaN
```

```
48 NaN
```

```
47 NaN
```

```
>>> s.loc[:3] # slice up to and including  
label 3
```

```
49 NaN
```

```
48 NaN
```

```
47 NaN
```

```
46 NaN
```

```
45 NaN
```

```
1 NaN
```

```
2 NaN
```

```
3 NaN
```

```
>>> s.ix[:3] # the integer is in the index so  
s.ix[:3] works like loc
```

```
49 NaN
```

```
48 NaN
```

```
47 NaN
```

```
46 NaN
```

```
45 NaN
```

```
1 NaN
```

```
2 NaN
```

```
3 NaN
```

Data Handling using Pandas

Pandas DataFrame

It is a two-dimensional data structure, just like any table (with rows & columns).

Basic Features of DataFrame

- Columns may be of different types
- Size can be changed(Mutable)
- Labeled axes (rows / columns)
- Arithmetic operations on rows and columns

Structure

SR. No.	Admn No	Student Name	Class	Section	Gender	Date Of Birth
1	001284	NIDHI MANDAL	I	A	Girl	07/08/2010
2	001285	SOUMYADIP BHATTACHARYA	I	A	Boy	24/02/2011
3	001286	SHREYAANG SHANDILYA	I	A	Boy	29/12/2010



columns



Rows

It can be created using constructor

pandas.DataFrame(data, index, columns, dtype, copy)

Visit : python.mykvs.in for regular updates

Data Handling using Pandas

Pandas DataFrame

Create DataFrame

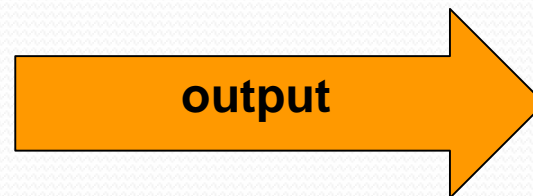
It can be created with followings

- Lists
- dict
- Series
- Numpy ndarrays
- Another DataFrame

Create an Empty DataFrame

e.g.

```
import pandas as pd1
df1 = pd1.DataFrame()
print(df1)
```



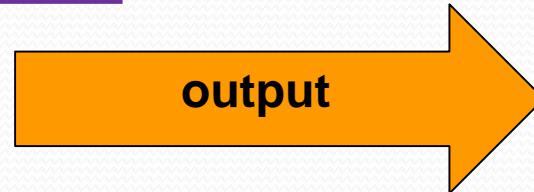
```
Empty
DataFrame
Columns: []
Index: []
```


Pandas DataFrame

Create a DataFrame from Lists

e.g.1

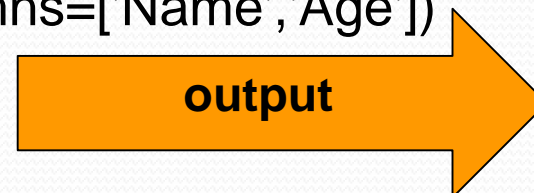
```
import pandas as pd1
data1 = [1,2,3,4,5]
df1 = pd1.DataFrame(data1)
print (df1)
```



```
0
0 1
1 2
2 3
3 4
4 5
```

e.g.2

```
import pandas as pd1
data1 = [['Freya',10],['Mohak',12],['Dwivedi',13]]
df1 = pd1.DataFrame(data1,columns=['Name','Age'])
print (df1)
```



```
      Name Age
1  Freya  10
2  Mohak  12
2 Dwivedi  13
```

Write below for numeric value as float

```
df1 = pd1.DataFrame(data,columns=['Name','Age'],dtype=float)
```

Pandas DataFrame

Create a DataFrame from Dict of ndarrays / Lists

e.g.1

```
import pandas as pd1
data1 = {'Name':['Freya', 'Mohak'],'Age':[9,10]}
df1 = pd1.DataFrame(data1)
print (df1)
```

Output

	Name	Age
1	Freya	9
2	Mohak	10

Write below as 3rd statement in above prog for indexing

```
df1 = pd1.DataFrame(data1, index=['rank1','rank2','rank3','rank4'])
```

Pandas DataFrame

Create a DataFrame from List of Dicts

e.g.1

```
import pandas as pd1
data1 = [{'x': 1, 'y': 2},{'x': 5, 'y': 4, 'z': 5}]
df1 = pd1.DataFrame(data1)
print (df1)
```

Output

	x	y	z
0	1	2	NaN
1	5	4	5.0

Write below as 3rd stmt in above program for indexing

```
df = pd.DataFrame(data, index=['first', 'second'])
```

Pandas DataFrame

Create a DataFrame from Dict of Series

e.g.1

```
import pandas as pd1
```

```
d1 = {'one' : pd1.Series([1, 2, 3], index=['a', 'b', 'c']),  
      'two' : pd1.Series([1, 2, 3, 4], index=['a', 'b', 'c', 'd'])}
```

```
df1 = pd1.DataFrame(d1)
```

```
print (df1)
```

Output

	one	two
a	1.0	1
b	2.0	2
c	3.0	3
d	NaN	4

Column Selection -> `print (df ['one'])`

Adding a new column by passing as Series: ->

```
df1['three']=pd1.Series([10,20,30],index=['a','b','c'])
```

Adding a new column using the existing columns values

```
df1['four']=df1['one']+df1['three']
```

Create a DataFrame from .txt file

Having a text file './inputs/dist.txt' as:

1	1	12.92
1	2	90.75
1	3	60.90
2	1	71.34

Pandas is shipped with built-in reader methods. For example the `pandas.read_table` method seems to be a good way to read (also in chunks) a tabular data file.

```
import pandas
```

```
df = pandas.read_table('./input/dists.txt', delim_whitespace=True,  
names=('A', 'B', 'C'))
```

will create a DataFrame objects with column named A made of data of type int64, B of int64 and C of float64

Create a DataFrame from csv(comma separated value) file / import data from cvs file

e.g.

Suppose filename.csv file contains following data

```
Date,"price","factor_1","factor_2"
```

```
2012-06-11,1600.20,1.255,1.548
```

```
2012-06-12,1610.02,1.258,1.554
```

```
import pandas as pd
```

```
# Read data from file 'filename.csv'
```

```
# (in the same directory that your python program is based)
```

```
# Control delimiters, rows, column names with read_csv
```

```
data = pd.read_csv("filename.csv")
```

```
# Preview the first 1 line of the loaded data
```

```
data.head(1)
```

Pandas DataFrame

Column addition

```
df = pd.DataFrame({"A": [1, 2, 3], "B": [4, 5, 6]})
```

```
c = [7,8,9]
```

```
df['C'] = c
```

Column Deletion

```
del df1['one'] # Deleting the first column using DEL function
```

```
df.pop('two') #Deleting another column using POP function
```

Rename columns

```
df = pd.DataFrame({"A": [1, 2, 3], "B": [4, 5, 6]})
```

```
>>> df.rename(columns={"A": "a", "B": "c"})
```

```
  a  c
```

```
0  1  4
```

```
1  2  5
```

```
2  3  6
```

Pandas DataFrame

Row Selection, Addition, and Deletion

#Selection by Label

```
import pandas as pd1
d1 = {'one' : pd1.Series([1, 2, 3], index=['a', 'b', 'c']),
      'two' : pd1.Series([1, 2, 3, 4], index=['a', 'b', 'c', 'd'])} df1
= pd1.DataFrame(d1)
print (df1.loc['b'])
```

Output

one 2.0

two 2.0

Name: b, dtype: float64

Pandas DataFrame

#Selection by integer location

```
import pandas as pd
d1 = {'one' : pd.Series([1, 2, 3], index=['a', 'b', 'c']),
      'two' : pd.Series([1, 2, 3, 4], index=['a', 'b', 'c', 'd'])}
df1 = pd.DataFrame(d1)
print (df1.iloc[2])
```

Output

```
one  3.0
two  3.0
Name: c, dtype: float64
```

Slice Rows : Multiple rows can be selected using ' : ' operator.

```
print (df1[2:4])
```

Pandas DataFrame

Addition of Rows

```
import pandas as pd1
```

```
df1 = pd1.DataFrame([[1, 2], [3, 4]], columns = ['a','b'])
```

```
df2 = pd1.DataFrame([[5, 6], [7, 8]], columns = ['a','b'])
```

```
df1 = df1.append(df2)
```

```
print (df1)
```

Deletion of Rows

```
# Drop rows with label 0
```

```
df1 = df1.drop(0)
```

Pandas DataFrame

Iterate over rows in a dataframe

e.g.

```
import pandas as pd1
import numpy as np1
raw_data1 = {'name': ['freya', 'mohak'],
             'age': [10, 1],
             'favorite_color': ['pink', 'blue'],
             'grade': [88, 92]}
df1 = pd1.DataFrame(raw_data1, columns = ['name', 'age',
'favorite_color', 'grade'])
for index, row in df1.iterrows():
    print (row["name"], row["age"])
```

Output

```
freya 10
mohak 1
```

Pandas DataFrame

Head & Tail

`head()` returns the first `n` rows (observe the index values). The default number of elements to display is five, but you may pass a custom number. `tail()` returns the last `n` rows .e.g.

```
import pandas as pd
import numpy as np
#Create a Dictionary of series
d = {'Name':pd.Series(['Tom','James','Ricky','Vin','Steve','Smith','Jack']),
     'Age':pd.Series([25,26,25,23,30,29,23]),
     'Rating':pd.Series([4.23,3.24,3.98,2.56,3.20,4.6,3.8])}
#Create a DataFrame
df = pd.DataFrame(d)
print ("Our data frame is:")
print df
print ("The first two rows of the data frame is:")
print df.head(2)
```

Pandas DataFrame

Indexing a DataFrame using `.loc[]` :

This function selects data by the label of the rows and columns.

```
#import the pandas library and aliasing as pd
import pandas as pd
import numpy as np
```

```
df = pd.DataFrame(np.random.randn(8, 4),
index = ['a','b','c','d','e','f','g','h'], columns = ['A', 'B', 'C', 'D'])
```

```
#select all rows for a specific column
print df.loc[:, 'A']
```

Accessing a DataFrame with a boolean index :

In order to access a dataframe with a boolean index, we have to create a dataframe in which index of dataframe contains a boolean value that is “True” or “False”.

```
# importing pandas as pd
import pandas as pd
```

```
# dictionary of lists
```

```
dict = {'name':['Mohak', 'Freya', 'Roshni'],
        'degree': ["MBA", "BCA", "M.Tech"],
        'score':[90, 40, 80]}
```

```
# creating a dataframe with boolean index
```

```
df = pd.DataFrame(dict, index = [True, False, True])
```

```
# accessing a dataframe using .loc[] function
```

```
print(df.loc[True]) #it will return rows of Mohak and Roshni only(matching true only)
```

import csv file in Pandas DataFrame

e.g.

```
import pandas as pd
```

```
# Takes the file's folder  
filepath = r"csv file path"
```

```
# read the CSV file  
df = pd.read_csv(filepath)
```

```
# print the first five rows  
print(df.head())
```

Export Pandas DataFrame to a CSV File

e.g.

```
import pandas as pd
```

```
cars = {'Brand': ['Honda Civic', 'Toyota Corolla', 'Ford Focus', 'Audi A4'],  
        'Price': [22000, 25000, 27000, 35000]  
        }
```

```
df = pd.DataFrame(cars, columns= ['Brand', 'Price'])
```

```
df.to_csv(r'C:\export_dataframe.csv', index = False, header=True)
```

```
print (df)
```