# LAB – 2
# PYTHON BASIC PRACTICE – II

**Introduction to PANDAS**

```python
import pandas as pd
import numpy as np


s=pd.Series([3,9,-2,10,5])
s.sum()
s.min()
s.max()
```
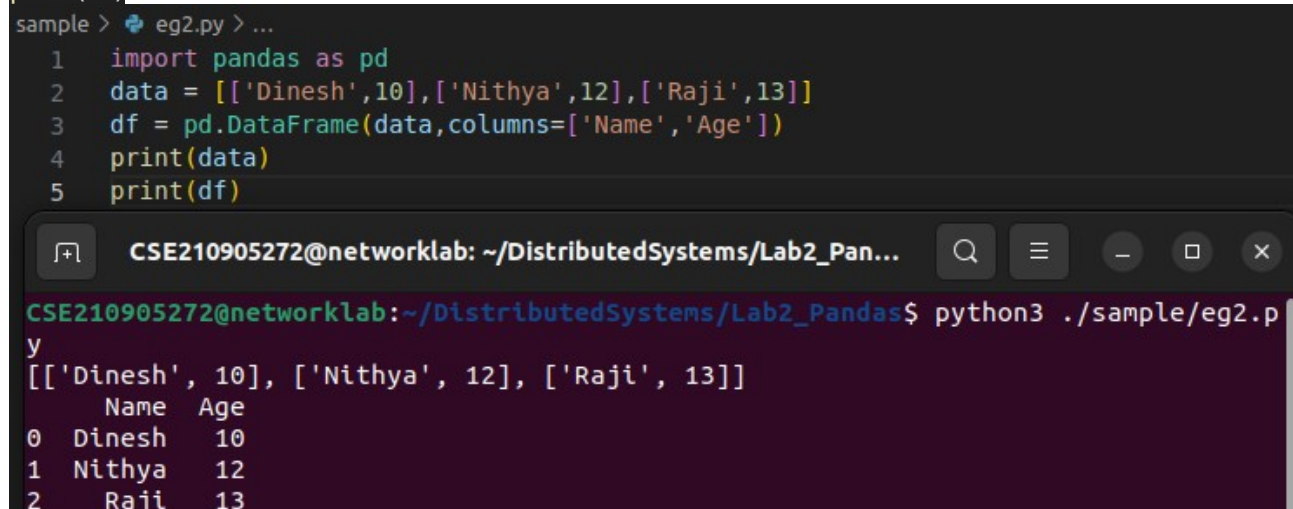
```
sample > 🐍 eg1.py > ...
    1    import pandas as pd
    2    import numpy as np
    3
    4    s=pd.Series([3,9,-2,10,5])
    5    print(s.sum())
    6    print(s.min())
    7    print(s.max())
    8
```

```
CSE210905272@networklab: ~/DistributedSystems/Lab2_Pan...

CSE210905272@networklab:~/DistributedSystems/Lab2_Pandas$ pyth
y
25
-2
10
CSE210905272@networklab:~/DistributedSystems/Lab2_Pandas$
```

**Creating a Data Frame**

```python
import pandas as pd
data = [['Dinesh',10],['Nithya',12],['Raji',13]]
df = pd.DataFrame(data,columns=['Name','Age'])
print(data)
print(df)
```

```
sample > 🐍 eg2.py > ...
    1    import pandas as pd
    2    data = [['Dinesh',10],['Nithya',12],['Raji',13]]
    3    df = pd.DataFrame(data,columns=['Name','Age'])
    4    print(data)
    5    print(df)
```

```
CSE210905272@networklab: ~/DistributedSystems/Lab2_Pan...

CSE210905272@networklab:~/DistributedSystems/Lab2_Pandas$ python3 ./sample/eg2.p
y
[['Dinesh', 10], ['Nithya', 12], ['Raji', 13]]
      Name  Age
0   Dinesh   10
1   Nithya   12
2     Raji   13
```

**Indexed Data Frame**

```python
import pandas as pd


data = {'Name':['Kavitha', 'Sudha', 'Raju','Vignesh'],'Age':[28,34,29,42]}
df = pd.DataFrame(data, index=['rank1','rank2','rank3','rank4'])


print(data)
```

```
print(df)
```



## Creating a DataFrame using Dictionary

```python
import pandas as pd
import numpy as np

df1=pd.DataFrame({'A':pd.Timestamp('20130102'),'B':np.array([3]*4,dtype='int32'),'C':pd.Categorical(['Male','Female','Male','Female'])})

print(df1.shape)
print(df1.dtypes)
print(df1.head())
print(df1.tail())
# print(df1.summary())
print(df1.T)
```



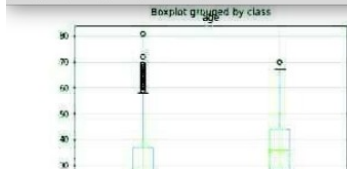.

## Creating a DataFrame using Dictionary

```python
import pandas as pd
import numpy as np
```

```python
dates=pd.date_range('20130101', periods=100)
df = pd.DataFrame(np.random.randn(100,4), index=dates, columns=list('ABCD'))
print(df.head()) #To view first 5 records
print(df.tail()) #To view last 5 records
print(df.index) #To view the index
print(df.columns) #To view the column names
print(df.T) #To transpose the df
print(df.sort_index(axis=1, ascending=False)) #Sorting by Axis
print(df.sort_values(by='B')) #Sorting by Values
print(df[0:3]) #Slicing the rows
print(df['20130105':'20130110']) #Slicing with index name
print(df.iloc[0]) #slicing with row and column index (like 2D Matrix)
print(df.iloc[0,:2]) #will fetch 1st row, first 2 columns
print(df.iloc[0,0]) #will fetch 1st row, 1st column element (single element)
print(df['A'])#which yields a Series
print(df['A','B']) #Selecting more than one column
print(df[['A','B']][:5]) #Selecting more than one column, with selected number of records
```

**Boolean Indexing**
**df[df.A>0], will fetch all positive values of A column**
**Include a 6th column (a categorical) character data**
**df['F']=['Male','Female','Female','Male','Female','Female']**
**Setting by assigning with a numpy array**
**df.loc[:,'D']=np.array([5]*len(df))**
**Deleting a row or column**
**df.drop ('col_name', axis =1, inplace=True)**
**will drop the column name specified in col_name**
**df.drop ('row_index', axis =0, inplace=True)**
**Df_new= pd.concat (df1, df2, axis=1)**
**Df_new.shape**
**D= pd.concat (A, B, axis=0)**
**D.shape**



.

**I/O operations**

Terminal output:

```
     1   85   66   29     0   26.6   0.351   31   0
2    8  183   64    0     0   23.3   0.672   32   1
3    1   89   66   23    94   28.1   0.167   21   0
4    0  137   40   35   168   43.1   2.288   33   1
        0    1    2    3     4      5       6    7   8
763  10  101   76   48   180   32.9   0.171   63   0
764   2  122   70   27     0   36.8   0.340   27   0
765   5  121   72   23   112   26.2   0.245   30   0
766   1  126   60    0     0   30.1   0.349   47   1
767   1   93   70   31     0   30.4   0.315   23   0
CSE210905272@networklab:~/DistributedSystems/Lab2_P
/eg6.py
        0    1    2    3     4      5       6    7   8
0    6  148   72   35     0   33.6   0.627   50   1
1    1   85   66   29     0   26.6   0.351   31   0
2    8  183   64    0     0   23.3   0.672   32   1
3    1   89   66   23    94   28.1   0.167   21   0
4    0  137   40   35   168   43.1   2.288   33   1
        0    1    2    3     4      5       6    7   8
763  10  101   76   48   180   32.9   0.171   63   0
764   2  122   70   27     0   36.8   0.340   27   0
765   5  121   72   23   112   26.2   0.245   30   0
766   1  126   60    0     0   30.1   0.349   47   1
767   1   93   70   31     0   30.4   0.315   23   0
CSE210905272@networklab:~/DistributedSystems/Lab2_P
/eg6.py
        0    1    2    3     4      5       6    7   8
0    6  148   72   35     0   33.6   0.627   50   1
1    1   85   66   29     0   26.6   0.351   31   0
2    8  183   64    0     0   23.3   0.672   32   1
3    1   89   66   23    94   28.1   0.167   21   0
4    0  137   40   35   168   43.1   2.288   33   1
        0    1    2    3     4      5       6    7   8
763  10  101   76   48   180   32.9   0.171   63   0
764   2  122   70   27     0   36.8   0.340   27   0
765   5  121   72   23   112   26.2   0.245   30   0
766   1  126   60    0     0   30.1   0.349   47   1
767   1   93   70   31     0   30.4   0.315   23   0
```



```python
import pandas as pd
import matplotlib.pyplot as plt
# %matplotlib inline

df=pd.read_csv('././lab2_req_files/xyz.csv',header=None)
print(df.head())
print(df.tail())
#attach header
df.columns=['preg','glu','bp','sft','ins','bmi','dpf','age','class']

#Let us visualize the scatter plot of two continuous variable

plt.scatter(df['bmi'],df['glu'])
plt.xlabel('bmi')
plt.ylabel('Glucose')
plt.title('bmi vs glucose')
plt.show()

#Let us visualize the histogram of another continuous variable 'Age'
plt.hist(df['age'])
plt.show()

#box plot
plt.boxplot(df['age'])
plt.show()
```
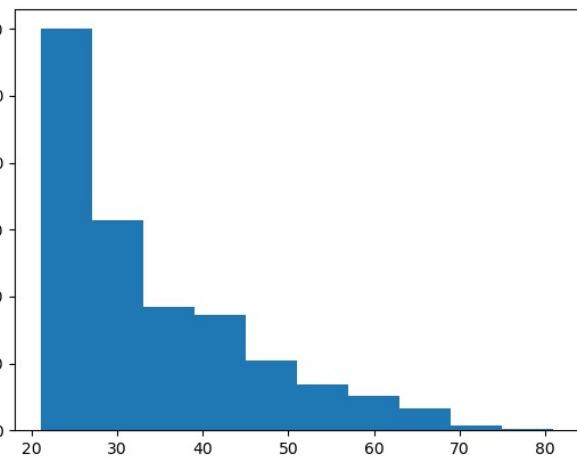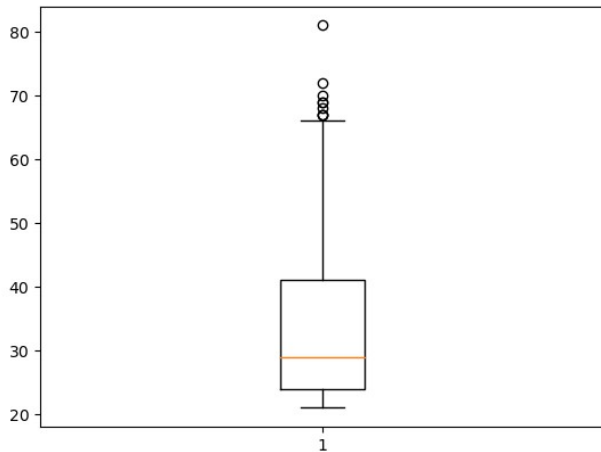
```python
import pandas as pd
import matplotlib.pyplot as plt
# %matplotlib inline

df=pd.read_csv('./../lab2_req_files/xyz.csv',header=None)
print(df.head())
print(df.tail())
#attach header
df.columns=['preg','glu','bp','sft','ins','bmi','dpf','age','class']

#Let us visualize the scatter plot of two continuous variable

plt.scatter(df['bmi'],df['glu'])
plt.xlabel('bmi')
plt.ylabel('Glucose')
plt.title('bmi vs glucose')
plt.show()

#Let us visualize the histogram of another continuous variable 'Age'
plt.hist(df['age'])
plt.show()

#box plot
plt.boxplot(df['age'])
plt.show()
```

#W = pd.read_csv('xyz.xls',header=None)
#W.head() #XLS file format also, we can read using pd.read_csv
#D= np.loadtxt('xyz.data',delimiter=",")
#D[:5,:] # this file is loaded in Numpy 2D array format
# Reading a XLSX file format
#G=pd.read_excel(xyz.xlsx',sheet_name='Sheet1')
#Pandas can read table tabs off of html. For example:
B = pd.read_html('./../lab2_req_files/Test runs-1.html')
for df in B:
print(df.head())

**Reading a TXT file format**

**H = pd.read_table('HR_for_week2.txt')**

.

# LAB EXERCISES

**Q1)Write a program to demonstrate while loop with else**
**->**
**i=0;**
**while(i<3):**
**        print("Value of i is:",i)**
**        i+=1**
**else:**
**        print("now value of i is>3")**



.

**Q2)Write a program to print negative Numbers in a List using while loop.**
**->**

**from icecream import ic**
**list1 = [1,2,3,-1,2,-3, -9]**
**i=0**

```python
print("Negative numbers: ")
while i < len(list1):
    if list1[i]<0:
        print(list1[i])
    i+=1
```



Q3Define a dictionary containing Students data {Name, Height, Qualification}.
a) Convert the dictionary into DataFrame
b) Declare a list that is to be converted into a new column (Address}
c) Using 'Address' as the column name and equate it to the list and display the result.)

->
```python
import pandas as pd
from icecream import ic
import numpy as np
dict1={'Name':['A','B','C'], 'Height':[180,183,190],
'Qualification':'UG'}
studentData= pd.DataFrame(dict1)
print(studentData)

list1=['Shivaji Peth', 'Budhwar Peth', 'Mangalwar Peth']
studentData['Address']=list1
print(studentData)
```

.

**Q4) Define a dictionary containing Students data {Name, Height, Qualification}.**
**a) Convert the dictionary into DataFrame**
**b) Use DataFrame.insert() to add a column and display the result.**

->
```python
import pandas as pd
from icecream import ic
import numpy as np
dict1={'Name':['A','B','C'], 'Height':[180,183,190], 'Qualification':'UG'}
studentData= pd.DataFrame(dict1)

studentData.insert(1,"Address",['Shivaji Peth', 'Budhwar Peth', 'Mangalwar Peth'])
print(studentData)
```



.

**Q5)**
**a) Create two data frames df1 and df2. df1 contains one column 'Name' and df2 contains 4 columns 'Maths', 'Physics', 'Chemistry' and 'Biology' .**
**b) Concatenate two data frames df1 and df2. Now insert one column 'Total' to the new data frame df_new and find the sum of all marks.**
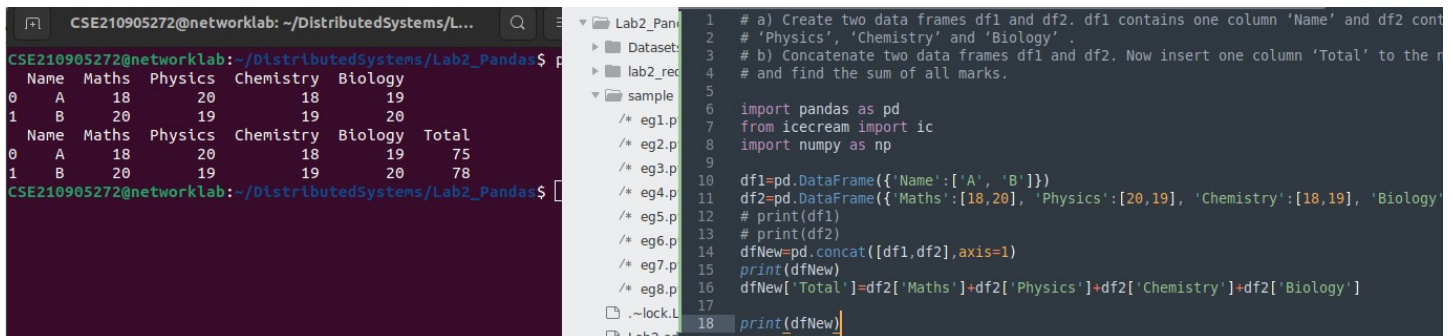
->
**import pandas as pd**
**from icecream import ic**
**import numpy as np**

**df1=pd.DataFrame({'Name':['A', 'B']})**
**df2=pd.DataFrame({'Maths':[18,20], 'Physics':[20,19], 'Chemistry':**
**[18,19], 'Biology':[19,20]})**
**# print(df1)**
**# print(df2)**
**dfNew=pd.concat([df1,df2],axis=1)**
**print(dfNew)**
**dfNew['Total']=df2['Maths']+df2['Physics']+df2['Chemistry']**
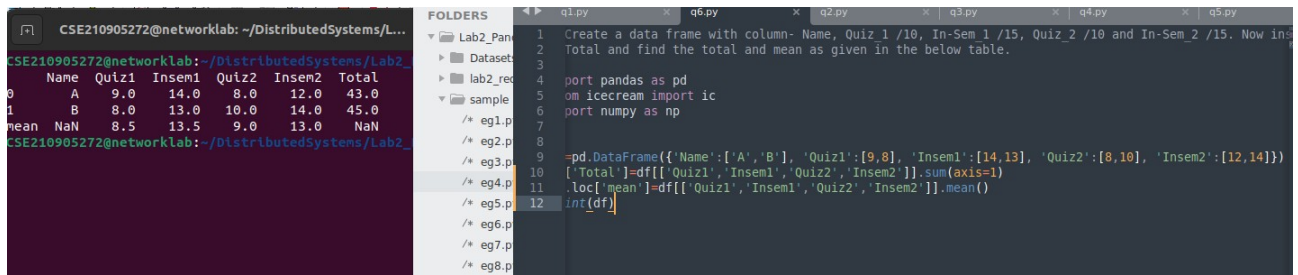**+df2['Biology']**

**print(dfNew)**



.

**Q6) Create a data frame with column- Name, Quiz_1 /10, In-Sem_1**
**/15, Quiz_2 /10 and In-Sem_2 /15. Now insert a column**
**Total and find the total and mean as given in the below table.**

->
**import pandas as pd**
**from icecream import ic**
**import numpy as np**

**df=pd.DataFrame({'Name':['A','B'], 'Quiz1':[9,8], 'Insem1':[14,13],**
**'Quiz2':[8,10], 'Insem2':[12,14]})**
**df['Total']=df[['Quiz1','Insem1','Quiz2','Insem2']].sum()**

# df.loc['mean']=df[['Quiz1','Insem1','Quiz2','Insem2']].mean()
# print(df)

```
CSE210905272@networklab:~/DistributedSystems/Lab2_
      Name  Quiz1  Insem1  Quiz2  Insem2  Total
0        A    9.0    14.0    8.0    12.0   43.0
1        B    8.0    13.0   10.0    14.0   45.0
mean   NaN    8.5    13.5    9.0    13.0    NaN
CSE210905272@networklab:~/DistributedSystems/Lab2_
```

```python
1  Create a data frame with column- Name, Quiz_1 /10, In-Sem_1 /15, Quiz_2 /10 and In-Sem_2 /15. Now ins
2  Total and find the total and mean as given in the below table.
3
4  port pandas as pd
5  om icecream import ic
6  port numpy as np
7
8
9  =pd.DataFrame({'Name':['A','B'], 'Quiz1':[9,8], 'Insem1':[14,13], 'Quiz2':[8,10], 'Insem2':[12,14]})
10  ['Total']=df[['Quiz1','Insem1','Quiz2','Insem2']].sum(axis=1)
11  .loc['mean']=df[['Quiz1','Insem1','Quiz2','Insem2']].mean()
12  int(df)
```