

Name:- Rohini Janardan Devkar

Roll no:- 23272

PRN no:- 72030818G

Class :- TE2

Problem Statement:-

Perform the following operations using Python on any open source dataset (eg. data.csv)

- Scan all variables for missing values and inconsistencies. If there are missing values and/or inconsistencies, use any of the suitable techniques to deal with them.
- Scan all numeric variables for outliers. If there are outliers, use any of the suitable techniques to deal with them.
- Apply data transformations on at least one of the variables. The purpose of this transformation should be one of the following reasons: to change the scale for better understanding of the variable, to convert a non-linear relation into a linear one, or to decrease the skewness and convert the distribution into a normal distribution. Reason and document your approach properly

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

```
In [2]: url_link = 'pract2.csv'
df=pd.read_csv(url_link)
```

```
In [3]: df.head(28)
```

	Rollno	Name	Gender	Branch	Attendance	Phy_marks	Che_marks	EMI_marks	PPS_marks	SME_marks
Out[3]:	0	1	Mohammed	M	Comp	72.0	62.0	98.0	83.0	89.0
	1	2	Reyanash	M	IT	58.0	62.0	83.0	63.0	58.0
	2	3	Aarav	M	IT	57.0	-20.0	100.0	NaN	96.0
	3	4	Atharv	M	IT	60.0	89.0	83.0	70.0	23.0
	4	5	Vivaan	M	Comp	85.0	90.0	NaN	78.0	23.0
	5	6	Advik	M	ENTC	94.0	99.0	84.0	100.0	96.0
	6	7	Ansh	M	ENTC	98.0	88.0	95.0	81.0	78.0
	7	8	Ishaan	M	ENTC	75.0	66.0	66.0	83.0	99.0
	8	9	Dhruv	M	ENTC	63.0	NaN	NaN	97.0	56.0
	9	10	Siddharth	M	ENTC	96.0	67.0	78.0	95.0	NaN
	10	11	Vivaan	M	ENTC	82.0	54.0	70.0	88.0	55.0
	11	12	Ajvan	M	IT	75.0	64.0	67.0	71.0	66.0
	12	13	Aarush	M	IT	67.0	56.0	81.0	NaN	90.0
	13	14	Leo	M	IT	98.0	-34.0	70.0	94.0	77.0
	14	15	Mayam	F	IT	64.0	87.0	60.0	90.0	65.0
	15	16	Sakshi	F	Comp	66.0	90.0	95.0	67.0	99.0
	16	17	Zaranew	F	Comp	93.0	54.0	NaN	75.0	90.0
	17	18	Inaya	F	Comp	74.0	67.0	93.0	93.0	87.0
	18	19	Aarya	F	Comp	72.0	88.0	84.0	81.0	86.0
	19	20	Pari	F	Comp	53.0	76.0	81.0	93.0	65.0

```
In [4]: df.shape
```

```
Out[4]: (28, 10)
```

```
In [5]: df.dtypes.value_counts()
```

```
Out[5]: float64    6
object      3
int64       1
dtype: int64
```

```
In [6]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 28 entries, 0 to 19
Data columns (total 10 columns):
#   Column          Non-Null Count  Dtype
---  --
0   Rollno          28 non-null    int64
1   Name            28 non-null    object
2   Gender          28 non-null    object
3   Branch          28 non-null    object
4   Attendance      28 non-null    float64
5   Phy_marks       19 non-null    float64
6   Che_marks       17 non-null    float64
7   EMI_marks       18 non-null    float64
8   PPS_marks       19 non-null    float64
9   SME_marks       18 non-null    float64
dtypes: float64(6), int64(1), object(3)
memory usage: 1.7+ KB
```

```
In [7]: df.describe()
```

	Rollno	Attendance	Phy_marks	Che_marks	EMI_marks	PPS_marks	SME_marks
count	20.00000	20.00000	19.00000	17.00000	18.00000	19.00000	18.00000
mean	10.50000	75.100000	63.421053	81.647059	83.444444	61.052632	52.055556
std	5.91608	14.660724	34.940133	12.036098	11.078449	43.767661	37.805185
min	1.00000	53.000000	-34.900000	60.000000	63.000000	-99.000000	-23.000000
25%	5.75000	63.750000	59.000000	70.000000	75.750000	56.000000	38.000000
50%	10.50000	73.000000	67.000000	83.000000	83.000000	66.000000	61.000000
75%	15.25000	87.000000	88.000000	93.000000	93.000000	87.500000	77.500000
max	20.00000	98.000000	99.000000	100.000000	100.000000	99.000000	100.000000

```
In [8]: df.isnull().sum()
```

```
Out[8]: Rollno      0
Name            0
Gender          0
Branch          0
Attendance      0
Phy_marks       9
Che_marks       1
EMI_marks       1
PPS_marks       1
SME_marks       2
dtype: int64
```

```
In [9]: data=df
col=[
miss=[]
data=data[data.columns]
pprint(miss)
for i in col:
    t=data[i].isnull().sum()
    if t!=0:
        miss.append(i)
    else:
        continue
print(miss)
```

```
[ 'Phy_marks', 'Che_marks', 'EMI_marks', 'PPS_marks', 'SME_marks' ]
```

```
In [10]: pd.options.mode.chained_assignment = None
for j in miss:
    q=data[j].dtypes
    if (q=="int64" or q=="float64") :
        Fdata[i]
        for k in range(data.shape[0]):
            if (f[k]<0 or f[k]>100):
                f[k]=(np.nan)
            else:
                continue
```

```
In [11]: for j in miss:
    q=data[j].dtypes
    if (q=="int64" or q=="float64") :
        data[j].fillna(data[j].mean(),inplace=True)
    else:
        data.fillna(method="bfill")
data.head(18)
```

	Rollno	Name	Gender	Branch	Attendance	Phy_marks	Che_marks	EMI_marks	PPS_marks	SME_marks
Out[11]:	0	1	Mohammed	M	Comp	72.0	62.000000	98.000000	63.000000	89.000000
	1	2	Reyanash	M	IT	58.0	62.000000	83.000000	63.000000	58.000000
	2	3	Aarav	M	IT	57.0	74.058824	100.000000	83.444444	56.000000
	3	4	Atharv	M	IT	60.0	89.000000	83.000000	70.000000	23.000000
	4	5	Vivaan	M	Comp	85.0	90.000000	81.647059	78.000000	23.000000
	5	6	Advik	M	ENTC	94.0	99.000000	84.000000	100.000000	96.000000
	6	7	Ansh	M	ENTC	98.0	88.000000	95.000000	81.000000	78.000000
	7	8	Ishaan	M	ENTC	75.0	66.000000	66.000000	83.000000	99.000000
	8	9	Dhruv	M	ENTC	63.0	74.058824	81.647059	97.000000	56.000000
	9	10	Siddharth	M	ENTC	96.0	67.000000	78.000000	96.000000	69.944444
	10	11	Vivaan	M	ENTC	82.0	54.000000	70.000000	88.000000	55.000000
	11	12	Ajvan	M	IT	75.0	64.000000	67.000000	71.000000	66.000000
	12	13	Aarush	M	IT	67.0	56.000000	81.000000	83.444444	90.000000
	13	14	Leo	M	IT	98.0	74.058824	70.000000	94.000000	77.000000
	14	15	Mayyam	F	IT	64.0	87.000000	60.000000	90.000000	65.000000
	15	16	Sakshi	F	Comp	66.0	90.000000	95.000000	67.000000	99.000000
	16	17	Zaranew	F	Comp	93.0	54.000000	81.647059	75.000000	90.000000
	17	18	Inaya	F	Comp	74.0	67.000000	93.000000	93.000000	87.000000
	18	19	Aarya	F	Comp	72.0	88.000000	84.000000	81.000000	86.000000
	19	20	Pari	F	Comp	53.0	76.000000	81.000000	93.000000	65.000000

```
In [12]: data['Total Marks']=data['Phy_marks']+data['Che_marks']+data['EMI_marks']+data['PPS_marks']+data['SME_marks']
data['Percentage']=data['Total Marks']/5
```

```
In [13]: data
```

	Rollno	Name	Gender	Branch	Attendance	Phy_marks	Che_marks	EMI_marks	PPS_marks	SME_marks	Total Marks	Percentage
Out[13]:	0	1	Mohammed	M	Comp	72.0	62.000000	98.000000	63.000000	89.000000	62.0	374.000000
	1	2	Reyanash	M	IT	58.0	62.000000	83.000000	63.000000	58.000000	58.0	374.000000
	2	3	Aarav	M	IT	57.0	74.058824	100.000000	83.444444	56.000000	100.0	413.503268
	3	4	Atharv	M	IT	60.0	89.000000	83.000000	70.000000	23.000000	23.0	341.708883
	4	5	Vivaan	M	Comp	85.0	90.000000	81.647059	78.000000	23.000000	60.0	332.647059
	5	6	Advik	M	ENTC	94.0	99.000000	84.000000	100.000000	96.000000	66.2	405.200000
	6	7	Ansh	M	ENTC	98.0	88.000000	95.000000	81.000000	78.000000	66.2	408.200000
	7	8	Ishaan	M	ENTC	75.0	66.000000	66.000000	83.000000	99.000000	56.0	340.944444
	8	9	Dhruv	M	ENTC	63.0	74.058824	81.647059	97.000000	56.000000	23.0	341.708883
	9	10	Siddharth	M	ENTC	96.0	67.000000	78.000000	96.000000	69.944444	23.0	332.944444
	10	11	Vivaan	M	ENTC	82.0	54.000000	70.000000	88.000000	55.000000	77.0	344.000000
	11	12	Ajvan	M	IT	75.0	64.000000	67.000000	71.000000	66.000000	65.0	333.000000
	12	13	Aarush	M	IT	67.0	56.000000	81.000000	83.444444	90.000000	86.0	396.444444
	13	14	Leo	M	IT	98.0	74.058824	70.000000	94.000000	77.000000	66.2	381.258824
	14	15	Mayyam	F	IT	64.0	87.000000	60.000000	90.000000	65.000000	64.0	366.000000
	15	16	Sakshi	F	Comp	66.0	90.000000	95.000000	67.000000	99.000000	93.0	444.000000
	16	17	Zaranew	F	Comp	93.0	54.000000	81.647059	75.000000	90.000000	53.0	353.647059
	17	18	Inaya	F	Comp	74.0	67.000000	93.000000	93.000000	87.000000	66.2	406.200000
	18	19	Aarya	F	Comp	72.0	88.000000	84.000000	81.000000	86.000000	78.0	417.000000
	19	20	Pari	F	Comp	53.0	76.000000	81.000000	93.000000	65.000000	66.2	381.200000

```
In [14]: data['Attendance'].plot(kind='box')
```

```
<AxesSubplot:~>
```



```
In [15]: data['Phy_marks'].plot(kind='box')
```

```
<AxesSubplot:~>
```



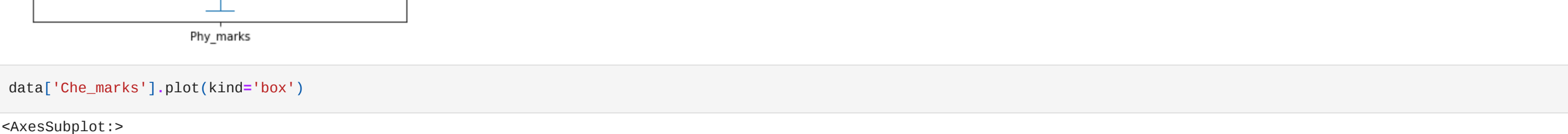
```
In [16]: data['Che_marks'].plot(kind='box')
```

```
<AxesSubplot:~>
```



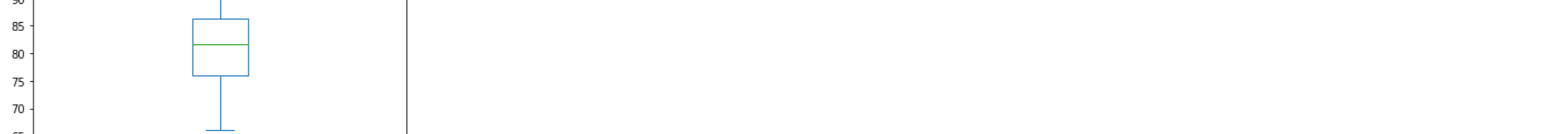
```
In [17]: data['EMI_marks'].plot(kind='box')
```

```
<AxesSubplot:~>
```



```
In [18]: data['PPS_marks'].plot(kind='box')
```

```
<AxesSubplot:~>
```



```
In [19]: data['SME_marks'].plot(kind='box')
```

```
<AxesSubplot:~>
```



```
In [20]: data['Total Marks'].plot(kind='box')
```

```
<AxesSubplot:~>
```



```
In [21]: data['Percentage'].plot(kind='box')
```

```
<AxesSubplot:~>
```



```
In [22]: Q1 = data['Attendance'].quantile(0.25)
Q3 = data['Attendance'].quantile(0.75)
IQR = Q3 - Q1
Lower_limit = Q1 - 1.5*IQR
Upper_limit = Q3 + 1.5*IQR
print("Q1 :",Q1,"nQ3 :",Q3,"nIQR :",IQR,"nLower_limit :",Lower_limit,"nUpper_limit:",Upper_limit)
```

```
Q1 : 63.75
Q3 : 87.8
IQR : 24.05
Lower_limit : 28.875
Upper_limit : 121.875
```

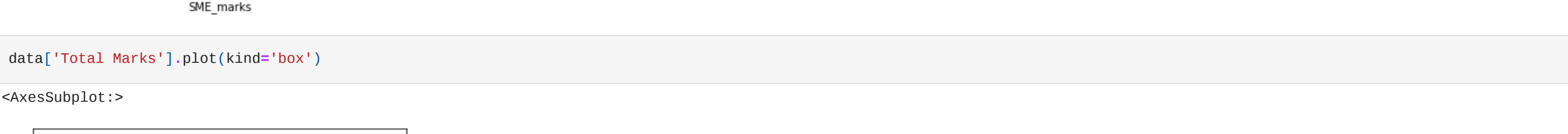
```
Out[23]: data[(data['Attendance']>Lower_limit)&(data['Attendance']<Upper_limit)]
```

```
In [24]: data[data[(data['Attendance']>Lower_limit)&(data['Attendance']<Upper_limit)]]
```

```
Out[24]: Rollno Name Gender Branch Attendance Phy_marks Che_marks EMI_marks PPS_marks SME_marks Total Marks Percentage
```

```
In [25]: data['Attendance'].plot(kind='box')
```

```
<AxesSubplot:~>
```



```
In [26]: Q1 = data['Che_marks'].quantile(0.25)
Q3 = data['Che_marks'].quantile(0.75)
IQR = Q3 - Q1
Lower_limit = Q1 - 1.5*IQR
Upper_limit = Q3 + 1.5*IQR
print("Q1 :",Q1,"nQ3 :",Q3,"nIQR :",IQR,"nLower_limit :",Lower_limit,"nUpper_limit:",Upper_limit)
```

```
Q1 : 76.8
Q3 : 86.25
IQR : 9.45
Lower_limit : 60.625
Upper_limit : 101.625
```

```
In [27]: data[(data['Che_marks']>Lower_limit)&(data['Che_marks']<Upper_limit)]
```

```
Out[27]: Rollno Name Gender Branch Attendance Phy_marks Che_marks EMI_marks PPS_marks SME_marks Total Marks Percentage
```

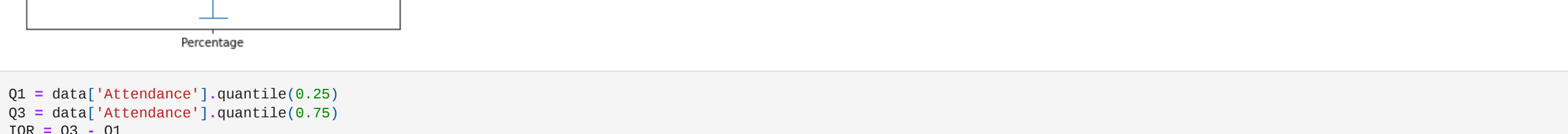
```
14 15 Mayam F IT 64.0 87.0 60.0 90.0 65.0 64.0 366.0 73.200000
```

```
In [28]: data[data[(data['Che_marks']>Lower_limit)&(data['Che_marks']<Upper_limit)]]
```

```
Out[28]: Rollno Name Gender Branch Attendance Phy_marks Che_marks EMI_marks PPS_marks SME_marks Total Marks Percentage
```

```
In [29]: data['Che_marks'].plot(kind='box')
```

```
<AxesSubplot:~>
```



```
In [30]: Q1 = data['EMI_marks'].quantile(0.25)
Q3 = data['EMI_marks'].quantile(0.75)
IQR = Q3 - Q1
Lower_limit = Q1 - 1.5*IQR
Upper_limit = Q3 + 1.5*IQR
print("Q1 :",Q1,"nQ3 :",Q3,"nIQR :",IQR,"nLower_limit :",Lower_limit,"nUpper_limit:",Upper_limit)
```

```
Q1 : 76.8
Q3 : 93.0
IQR : 16.2
Lower_limit : 51.75
Upper_limit : 117.75
```

```
In [31]: data[(data['EMI_marks']>Lower_limit)&(data['EMI_marks']<Upper_limit)]
```

```
Out[31]: Rollno Name Gender Branch Attendance Phy_marks Che_marks EMI_marks PPS_marks SME_marks Total Marks Percentage
```

```
In [32]: data[data[(data['EMI_marks']>Lower_limit)&(data
```