

Data Wrangling II

Create an "Academic performance" dataset of students and perform the following operations using Python.

1. 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.
2. Scan all numeric variables for outliers. If there are outliers, use any of the suitable techniques to deal with them.
3. 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.

IMPORTING ALL REQUIRED LIBRARIES

```
import pandas as pd

from google.colab import files
files.upload()

 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 om2.csv to om2.csv
{'om2.csv': b'Roll
No.Name.DSBDAAI.WT.CC.TOTAL.PERCENTAGE.RESULT\r\n1.Om.44.47.65.27.183.45.75.PASS\r\n'}
```

Load the Dataset into pandas dataframe.

```
om=pd.read_csv("/content/om2.csv")
```

```
om
```

Roll No	Name	DSBDA	AI	WT	CC	TOTAL	PERCENTAGE	RESULT
0	1	Om	44	47.0	65.0	27.0	183	45.75
1	2	NaN	47	45.0	53.0	41.0	186	46.50
2	3	NaN	70	70.0	75.0	50.0	265	66.25
3	4	NaN	50	50.0	40.0	30.0	160	40.00
4	5	NaN	50	50.0	40.0	30.0	160	40.00
5	6	NaN	50	50.0	40.0	30.0	160	40.00
6	7	NaN	50	50.0	40.0	30.0	160	40.00
7	8	NaN	50	50.0	40.0	30.0	160	40.00
8	9	NaN	50	50.0	40.0	30.0	160	40.00
9	10	NaN	50	50.0	40.0	30.0	160	40.00
10	11	NaN	50	50.0	40.0	30.0	160	40.00
11	12	NaN	50	50.0	40.0	30.0	160	40.00
12	13	NaN	50	50.0	40.0	30.0	160	40.00
13	14	NaN	50	50.0	40.0	30.0	160	40.00
14	15	NaN	50	50.0	40.0	30.0	160	40.00
15	16	NaN	50	50.0	40.0	30.0	160	40.00
16	17	NaN	50	50.0	40.0	30.0	160	40.00
17	18	NaN	50	50.0	40.0	30.0	160	40.00
18	19	NaN	50	50.0	40.0	30.0	160	40.00
19	20	NaN	50	50.0	40.0	30.0	160	40.00
20	21	NaN	50	50.0	40.0	30.0	160	40.00
21	22	NaN	50	50.0	40.0	30.0	160	40.00
22	23	NaN	50	50.0	40.0	30.0	160	40.00
23	24	NaN	50	50.0	40.0	30.0	160	40.00
24	25	NaN	50	50.0	40.0	30.0	160	40.00
25	26	NaN	50	50.0	40.0	30.0	160	40.00
26	27	NaN	50	50.0	40.0	30.0	160	40.00
27	28	NaN	50	50.0	40.0	30.0	160	40.00
28	29	NaN	50	50.0	40.0	30.0	160	40.00
29	30	NaN	50	50.0	40.0	30.0	160	40.00

Data Preprocessing

```
om.isnull()
```

Roll No	Name	DSBDA	AI	WT	CC	TOTAL	PERCENTAGE	RESULT
0	False	False						
1	False	True	False	False	False	False	False	False
2	False	True	False	False	False	False	False	False
3	False	True	False	True	False	False	False	False
4	False	True	False	False	False	False	False	False
5	False	True	False	False	False	True	False	False
6	False	True	False	False	False	False	False	False
7	False	True	False	False	True	False	False	False
8	False	True	False	False	False	False	False	False
9	False	True	False	False	False	False	False	False
10	False	True	False	False	False	False	False	False
11	False	True	False	True	False	False	False	False
12	False	True	False	False	False	False	False	False
13	False	True	False	False	False	False	False	False
14	False	True	False	False	False	False	False	False
15	False	True	False	False	False	False	False	False
16	False	True	False	False	True	False	False	False
17	False	True	False	False	False	False	False	False
18	False	True	False	False	False	False	False	False
19	False	True	False	False	False	False	False	False
20	False	True	False	False	False	False	False	False
21	False	True	False	False	False	False	False	False
22	False	True	False	False	False	False	False	False
23	False	True	False	False	False	False	False	False
24	False	True	False	False	False	False	False	False
25	False	True	False	False	False	False	False	False
26	False	True	False	False	False	False	False	False
27	False	True	False	False	False	False	False	False
28	False	True	False	False	False	False	False	False
29	False	True	False	False	False	False	False	False

CHECKING FOR NULL VALUES

```
pd.isnull(om['WT'])
```

0	False
1	False
2	False
3	False
4	False
5	False
6	False
7	True
8	False
9	False
10	False
11	False

```

12  False
13  False
14  False
15  False
16  True
17  False
18  False
19  False
20  False
21  False
22  False
23  False
24  False
25  False
26  False
27  False
28  False
29  False
Name: WT, dtype: bool

```

```
om.notnull()
```

	Roll No	Name	DSBDA	AI	WT	CC	TOTAL	PERCENTAGE	RESULT
0	True	True	True	True	True	True	True	True	True
1	True	False	True	True	True	True	True	True	True
2	True	False	True	True	True	True	True	True	True
3	True	False	True	False	True	True	True	True	True
4	True	False	True	True	True	True	True	True	True
5	True	False	True	True	True	False	True	True	True
6	True	False	True	True	True	True	True	True	True
7	True	False	True	True	False	True	True	True	True
8	True	False	True	True	True	True	True	True	True
9	True	False	True	True	True	True	True	True	True
10	True	False	True	True	True	True	True	True	True
11	True	False	True	False	True	True	True	True	True
12	True	False	True	True	True	True	True	True	True
13	True	False	True	True	True	True	True	True	True
14	True	False	True	True	True	True	True	True	True
15	True	False	True	True	True	True	True	True	True
16	True	False	True	True	False	True	True	True	True
17	True	False	True	True	True	True	True	True	True
18	True	False	True	True	True	True	True	True	True
19	True	False	True	True	True	True	True	True	True
20	True	False	True	True	True	True	True	True	True
21	True	False	True	True	True	True	True	True	True
22	True	False	True	True	True	True	True	True	True
23	True	False	True	True	True	True	True	True	True
24	True	False	True	True	True	True	True	True	True
25	True	False	True	True	True	True	True	True	True
26	True	False	True	True	True	True	True	True	True
27	True	False	True	True	True	True	True	True	True
28	True	False	True	True	True	True	True	True	True
29	True	False	True	True	True	True	True	True	True

FILLING NULL VALUES

```
om.fillna(0)
```

Roll No	Name	DSBDA	AI	WT	CC	TOTAL	PERCENTAGE	RESULT	
0	1	Om	44	47.0	65.0	27.0	183	45.75	PASS
1	2	0	47	45.0	53.0	41.0	186	46.50	PASS
2	3	0	70	70.0	75.0	50.0	265	66.25	PASS
3	4	0	58	0.0	42.0	34.0	134	33.50	FAIL
4	5	0	56	40.0	41.0	78.0	215	53.75	PASS
5	6	0	29	32.0	69.0	0.0	130	32.50	FAIL
6	7	0	68	33.0	51.0	42.0	194	48.50	PASS
7	8	0	29	67.0	0.0	58.0	154	38.50	FAIL
8	9	0	43	63.0	59.0	55.0	220	55.00	PASS
9	10	0	45	78.0	55.0	50.0	228	57.00	PASS
10	11	0	68	67.0	72.0	38.0	245	61.25	PASS
11	12	0	24	0.0	58.0	58.0	140	35.00	FAIL
12	13	0	38	67.0	43.0	77.0	225	56.25	PASS
13	14	0	75	54.0	44.0	42.0	215	53.75	PASS
14	15	0	41	70.0	39.0	63.0	213	53.25	PASS
15	16	0	27	55.0	48.0	74.0	204	51.00	PASS
16	17	0	33	63.0	0.0	68.0	164	41.00	PASS
17	18	0	34	61.0	63.0	29.0	187	46.75	PASS
18	19	0	28	29.0	41.0	26.0	124	31.00	FAIL
19	20	0	37	67.0	76.0	23.0	203	50.75	PASS
20	21	0	41	31.0	34.0	55.0	161	40.25	PASS
21	22	0	69	45.0	77.0	59.0	250	62.50	PASS
22	23	0	73	47.0	39.0	55.0	214	53.50	PASS
23	24	0	44	50.0	51.0	65.0	210	52.50	PASS
24	25	0	38	43.0	61.0	52.0	194	48.50	PASS
25	26	0	62	48.0	42.0	68.0	220	55.00	PASS
26	27	0	41	74.0	37.0	40.0	192	48.00	PASS
27	28	0	64	33.0	44.0	60.0	201	50.25	PASS
28	29	0	67	45.0	38.0	61.0	211	52.75	PASS

```
cols_with_na=[]
for col in om.columns:
    if om[col].isna().any():
        cols_with_na.append(col)
cols_with_na
['Name', 'AI', 'WT', 'CC']
```

```
om.fillna(1)
```

	Roll No	Name	DSBDA	AI	WT	CC	TOTAL	PERCENTAGE	RESULT
0	1	Om	44	47.000000	65.000000	27.000000	183	45.75	PASS
1	2	1	47	45.000000	53.000000	41.000000	186	46.50	PASS
2	3	1	70	70.000000	75.000000	50.000000	265	66.25	PASS
3	4	1	58	53.535714	42.000000	34.000000	134	33.50	FAIL
4	5	1	56	40.000000	41.000000	78.000000	215	53.75	PASS
5	6	1	29	32.000000	69.000000	51.551724	130	32.50	FAIL
6	7	1	68	33.000000	51.000000	42.000000	194	48.50	PASS
7	8	1	29	67.000000	52.428571	58.000000	154	38.50	FAIL
8	9	1	43	63.000000	59.000000	55.000000	220	55.00	PASS
9	10	1	45	78.000000	55.000000	50.000000	228	57.00	PASS
10	11	1	68	67.000000	72.000000	38.000000	245	61.25	PASS
11	12	1	24	53.535714	58.000000	58.000000	140	35.00	FAIL
12	13	1	38	67.000000	43.000000	77.000000	225	56.25	PASS
13	14	1	75	54.000000	44.000000	42.000000	215	53.75	PASS
14	15	1	41	70.000000	39.000000	63.000000	213	53.25	PASS
15	16	1	27	55.000000	48.000000	74.000000	204	51.00	PASS
16	17	1	33	63.000000	52.428571	68.000000	164	41.00	PASS

FILLING NULL VALUES WITH MEAN

```
20 21 1 41 31.000000 34.000000 55.000000 161 40.25 PAss
```

```
om['DSBDA']=om['DSBDA'].fillna(om['DSBDA'].mean())
```

```
om['WT']=om['WT'].fillna(om['WT'].mean())
```

```
om['CC']=om['CC'].fillna(om['CC'].mean())
```

```
25 26 1 62 18.000000 12.000000 68.000000 220 55.00 PAss
```

```
om['AI']=om['AI'].fillna(om['AI'].mean())
```

```
om['PERCENTAGE']=om['PERCENTAGE'].fillna(om['PERCENTAGE'].mean())
```

```
28 29 1 67 45.000000 38.000000 61.000000 211 52.75 PAss
```

```
om['TOTAL']=om['TOTAL'].fillna(om['TOTAL'].mean())
```

```
om
```

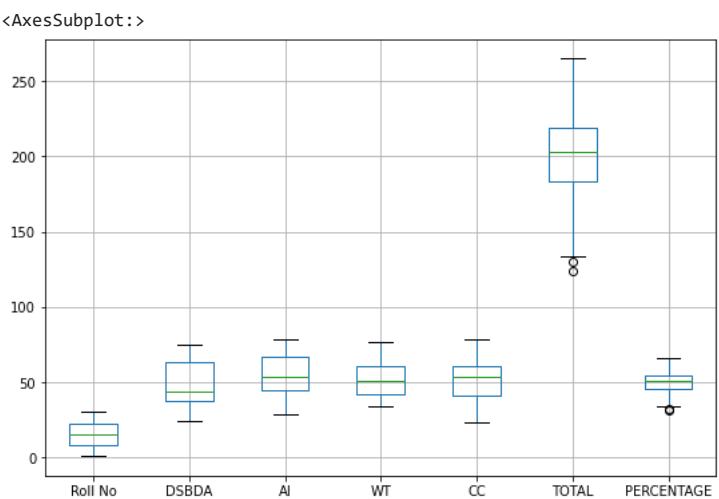
Roll No	Name	DSBDA	AI	WT	CC	TOTAL	PERCENTAGE	RESULT	
0	1	Om	44	47.000000	65.000000	27.000000	183	45.75	PASS
1	2	NaN	47	45.000000	53.000000	41.000000	186	46.50	PASS
2	3	NaN	70	70.000000	75.000000	50.000000	265	66.25	PASS
3	4	NaN	58	53.535714	42.000000	34.000000	134	33.50	FAIL
4	5	NaN	56	40.000000	41.000000	78.000000	215	53.75	PASS
5	6	NaN	29	32.000000	69.000000	51.551724	130	32.50	FAIL
6	7	NaN	68	33.000000	51.000000	42.000000	194	48.50	PASS
7	8	NaN	29	67.000000	52.428571	58.000000	154	38.50	FAIL
8	9	NaN	43	63.000000	59.000000	55.000000	220	55.00	PASS
9	10	NaN	45	78.000000	55.000000	50.000000	228	57.00	PASS
10	11	NaN	68	67.000000	72.000000	38.000000	245	61.25	PASS
11	12	NaN	24	53.535714	58.000000	58.000000	140	35.00	FAIL
12	13	NaN	38	67.000000	43.000000	77.000000	225	56.25	PASS
13	14	NaN	75	54.000000	44.000000	42.000000	215	53.75	PASS

DATA VISUALIZATION

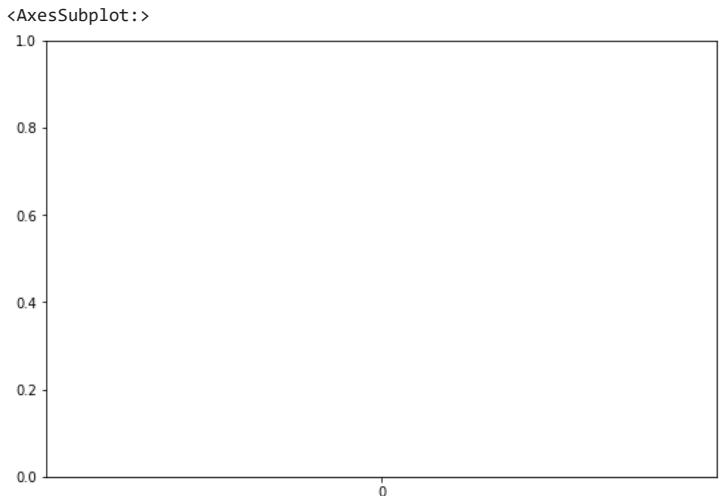
```
import seaborn as sns  
import matplotlib.pyplot as plt
```

CHECKING FOR OUTLIERS

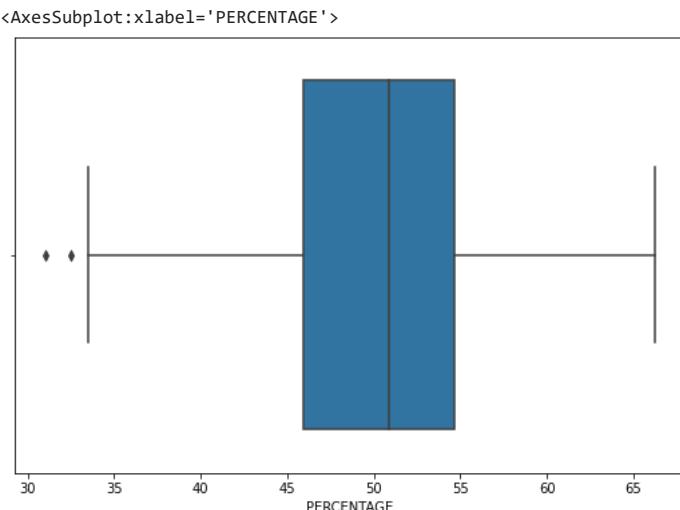
```
om_boxplot()
```



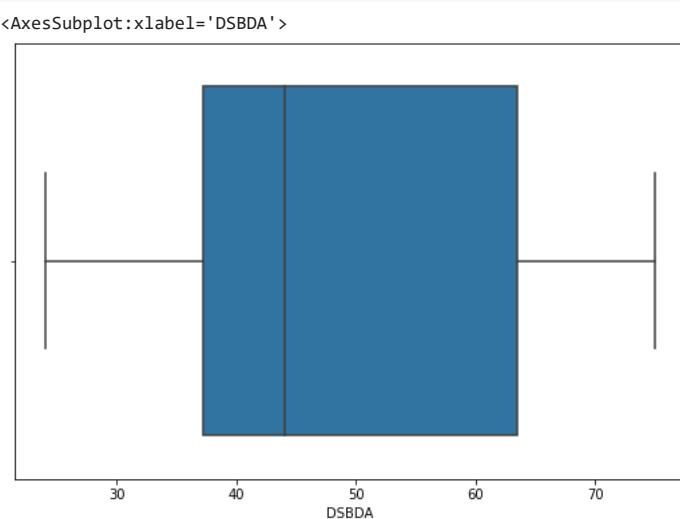
```
sns.boxplot()
```



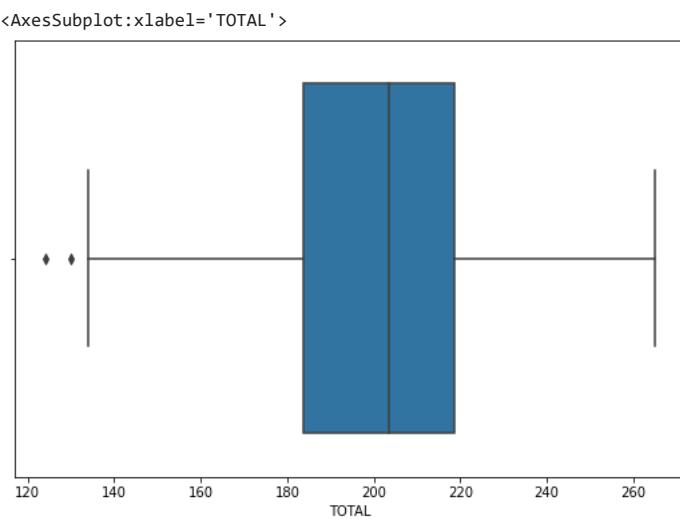
```
sns.boxplot(x=om.PERCENTAGE)
```



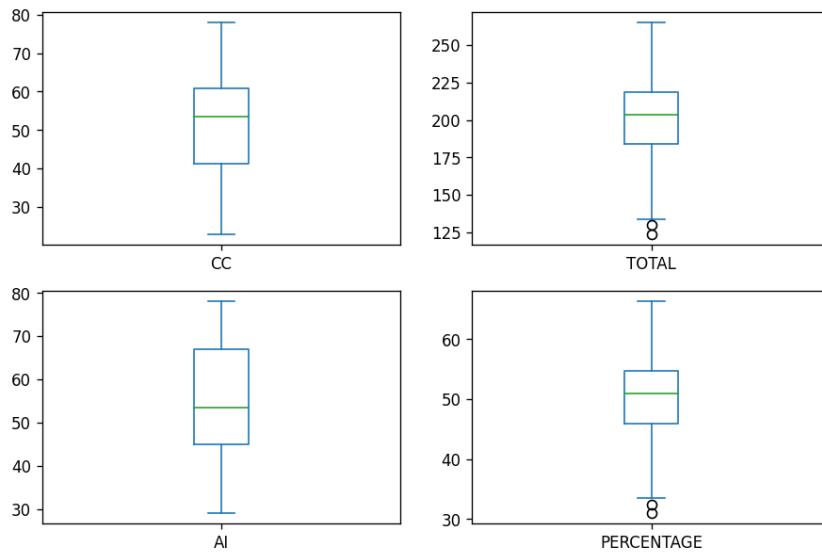
```
sns.boxplot(x=om.DSBDA)
```



```
sns.boxplot(x=om.TOTAL)
```



```
plt.rcParams["figure.figsize"]=(9,6)
om_list=['CC','TOTAL','AI','PERCENTAGE']
fig,axes=plt.subplots(2,2)
fig.set_dpi(120)
count=0
for r in range(2):
    for c in range(2):
        _=om[om_list[count]].plot(kind='box',ax=axes[r,c])
        count+=1
```



REMOVING OUTLIERS

```

Q1=om['PERCENTAGE'].quantile(0.25)
Q3=om['PERCENTAGE'].quantile(0.75)
IQR=Q3-Q1
Lower_limit=Q1 - 1.5*IQR
Upper_limit=Q1 + 1.5*IQR
print(f'Q1={Q1},Q3={Q3},IQR={IQR},LOWER LIMIT={Lower_limit}, UPPER LIMIT={Upper_limit}')

```

Q1=45.9375,Q3=54.6875,IQR=8.75,LOWER LIMIT=32.8125, UPPER LIMIT=59.0625

```
om[(om['PERCENTAGE']<Lower_limit) | (om['PERCENTAGE']>Upper_limit)]
```

	Roll No	Name	DSBDA	AI	WT	CC	TOTAL	PERCENTAGE	RESULT
2	3	NaN	70	70.0	75.0	50.000000	265	66.25	PASS
5	6	NaN	29	32.0	69.0	51.551724	130	32.50	FAIL
10	11	NaN	68	67.0	72.0	38.000000	245	61.25	PASS
18	19	NaN	28	29.0	41.0	26.000000	124	31.00	FAIL
21	22	NaN	69	45.0	77.0	59.000000	250	62.50	PASS

```

Q1=om['TOTAL'].quantile(0.25)
Q3=om['TOTAL'].quantile(0.75)
IQR=Q3-Q1
Lower_limit=Q1 - 1.5*IQR
Upper_limit=Q1 + 1.5*IQR
print(f'Q1={Q1},Q3={Q3},IQR={IQR},LOWER LIMIT={Lower_limit}, UPPER LIMIT={Upper_limit}')

```

Q1=183.75,Q3=218.75,IQR=35.0,LOWER LIMIT=131.25, UPPER LIMIT=236.25

```
om[(om['TOTAL']<Lower_limit) | (om['TOTAL']>Upper_limit)]
```

	Roll No	Name	DSBDA	AI	WT	CC	TOTAL	PERCENTAGE	RESULT
2	3	NaN	70	70.0	75.0	50.000000	265	66.25	PASS
5	6	NaN	29	32.0	69.0	51.551724	130	32.50	FAIL
10	11	NaN	68	67.0	72.0	38.000000	245	61.25	PASS
18	19	NaN	28	29.0	41.0	26.000000	124	31.00	FAIL
21	22	NaN	69	45.0	77.0	59.000000	250	62.50	PASS

```

outliers=[]
for i in om.PERCENTAGE:
    if i<Lower_limit or i>Upper_limit:
        outliers.append(i)
print("outliers are",outliers)

outliers are [45.75, 46.5, 66.25, 33.5, 53.75, 32.5, 48.5, 38.5, 55.0, 57.0, 61.25, 35.0, 56.25, 53.75, 53.25, 51.0, 41.0, 46.75, 3

```

◀ ▶

```

outliers=[]
for i in om.TOTAL:
    if i<Lower_limit or i>Upper_limit:
        outliers.append(i)
print("outliers are",outliers)

outliers are [265, 130, 245, 124, 250]

```

Upper_limit

236.25

Lower_limit

131.25

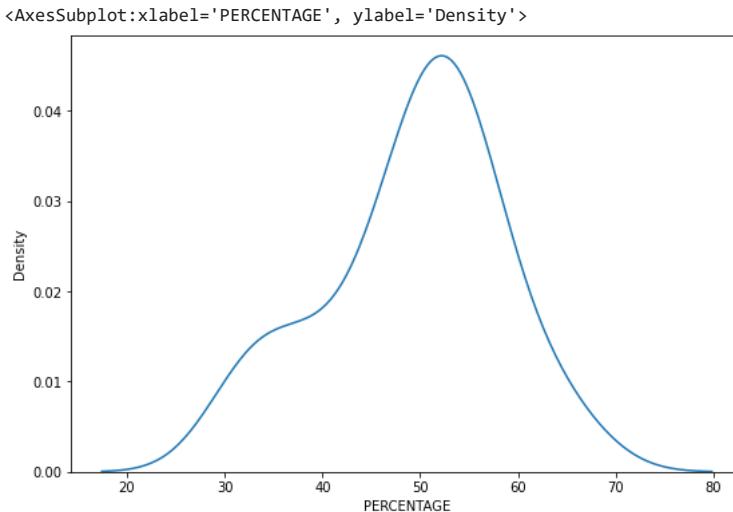
```

om2=om[om.PERCENTAGE<Lower_limit]
om2

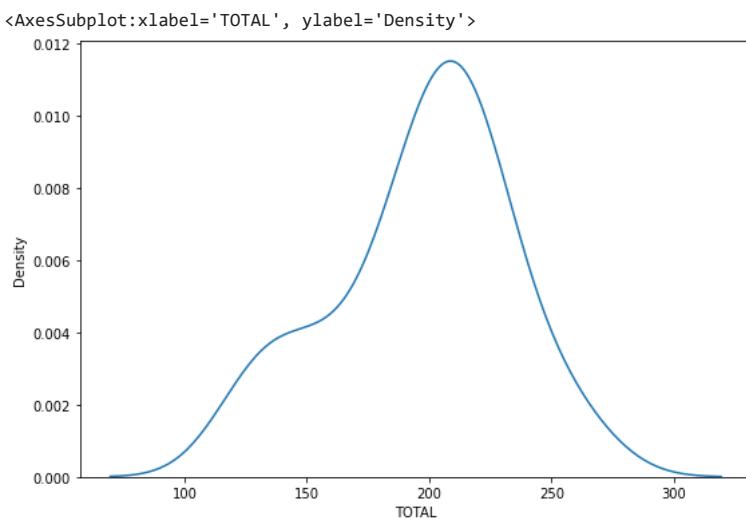
```

	Roll No	Name	DSBDA	AI	WT	CC	TOTAL	PERCENTAGE	RESULT
0	1	Om	44	47.000000	65.000000	27.000000	183	45.75	PASS
1	2	NaN	47	45.000000	53.000000	41.000000	186	46.50	PASS
2	3	NaN	70	70.000000	75.000000	50.000000	265	66.25	PASS
3	4	NaN	58	53.535714	42.000000	34.000000	134	33.50	FAIL
4	5	NaN	56	40.000000	41.000000	78.000000	215	53.75	PASS
5	6	NaN	29	32.000000	69.000000	51.551724	130	32.50	FAIL
6	7	NaN	68	33.000000	51.000000	42.000000	194	48.50	PASS
7	8	NaN	29	67.000000	52.428571	58.000000	154	38.50	FAIL
8	9	NaN	43	63.000000	59.000000	55.000000	220	55.00	PASS
9	10	NaN	45	78.000000	55.000000	50.000000	228	57.00	PASS
10	11	NaN	68	67.000000	72.000000	38.000000	245	61.25	PASS
11	12	NaN	24	53.535714	58.000000	58.000000	140	35.00	FAIL
12	13	NaN	38	67.000000	43.000000	77.000000	225	56.25	PASS
13	14	NaN	75	54.000000	44.000000	42.000000	215	53.75	PASS
14	15	NaN	41	70.000000	39.000000	63.000000	213	53.25	PASS
15	16	NaN	27	55.000000	48.000000	74.000000	204	51.00	PASS
16	17	NaN	33	63.000000	52.428571	68.000000	164	41.00	PASS
17	18	NaN	34	61.000000	63.000000	29.000000	187	46.75	PASS
18	19	NaN	28	29.000000	41.000000	26.000000	124	31.00	FAIL
19	20	NaN	37	67.000000	76.000000	23.000000	203	50.75	PASS
20	21	NaN	41	31.000000	34.000000	55.000000	161	40.25	PASS
21	22	NaN	69	45.000000	77.000000	59.000000	250	62.50	PASS
22	23	NaN	73	47.000000	39.000000	55.000000	214	53.50	PASS
23	24	NaN	44	50.000000	51.000000	65.000000	210	52.50	PASS
24	25	NaN	38	43.000000	61.000000	52.000000	194	48.50	PASS
25	26	NaN	62	48.000000	42.000000	68.000000	220	55.00	PASS
26	27	NaN	41	74.000000	37.000000	40.000000	192	48.00	PASS
27	28	NaN	64	33.000000	44.000000	60.000000	201	50.25	PASS
28	29	NaN	67	45.000000	38.000000	61.000000	211	52.75	PASS
29	30	NaN	58	75.000000	51.000000	47.000000	231	57.75	PASS

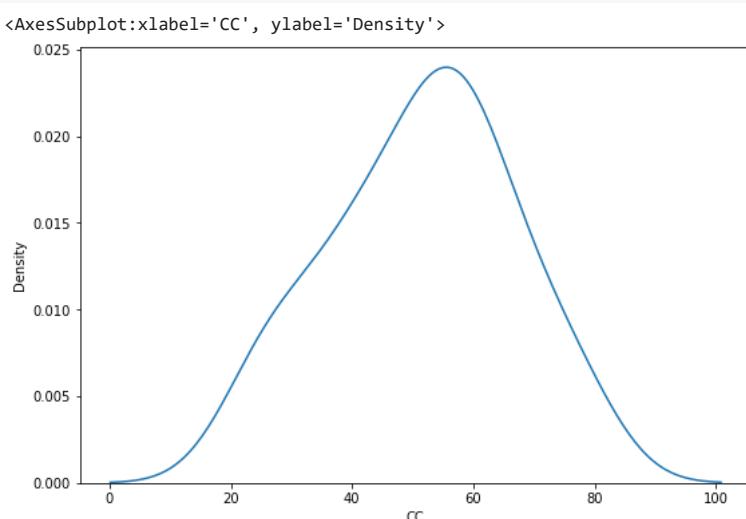
```
sns.kdeplot(om.PERCENTAGE)
```



```
sns.kdeplot(om.TOTAL)
```

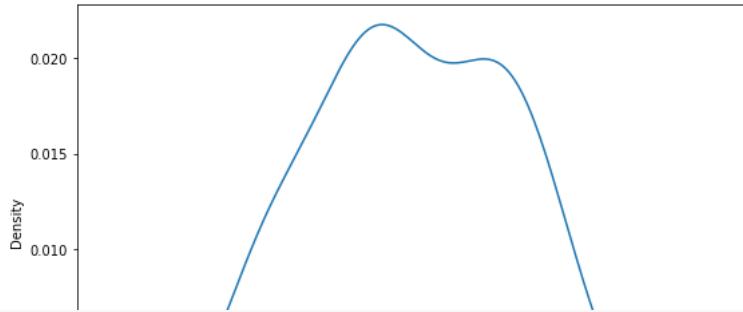


```
sns.kdeplot(om.CC)
```



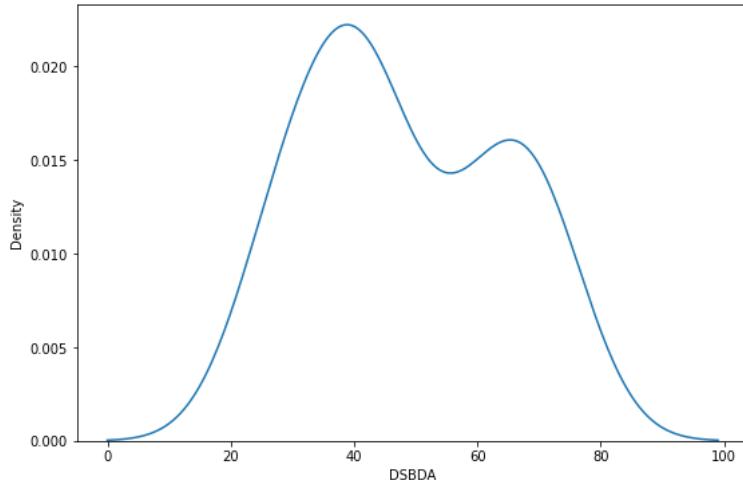
```
sns.kdeplot(om.AI)
```

```
<AxesSubplot:xlabel='AI', ylabel='Density'>
```



```
sns.kdeplot(om.DSBDA)
```

```
<AxesSubplot:xlabel='DSBDA', ylabel='Density'>
```



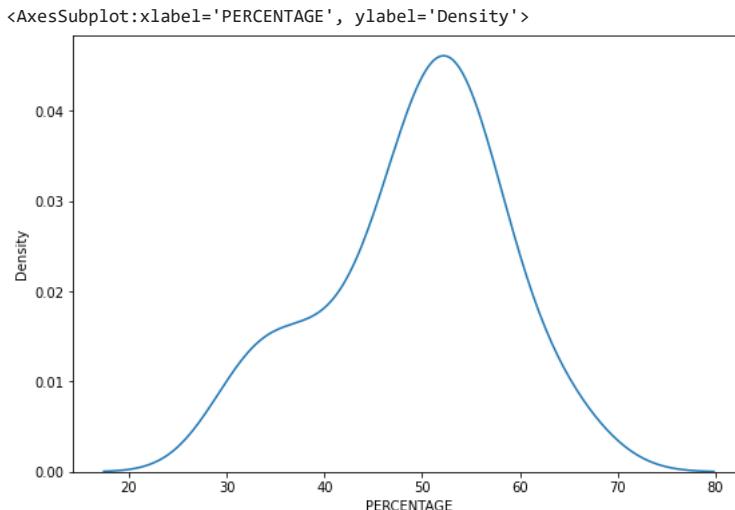
```
om.PERCENTAGE
```

```
0    45.75
1    46.50
2    66.25
3    33.50
4    53.75
5    32.50
6    48.50
7    38.50
8    55.00
9    57.00
10   61.25
11   35.00
12   56.25
13   53.75
14   53.25
15   51.00
16   41.00
17   46.75
18   31.00
19   50.75
20   40.25
21   62.50
22   53.50
23   52.50
24   48.50
25   55.00
26   48.00
27   50.25
28   52.75
29   57.75
Name: PERCENTAGE, dtype: float64
```

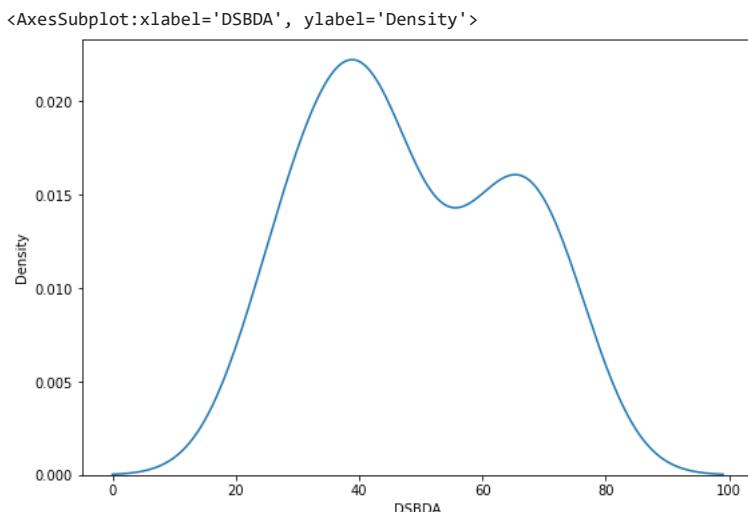
```
om.skew()
```

```
<ipython-input-154-54adde609add>:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None')
  om.skew()
Roll No      0.000000
DSBDA       0.223449
AI        -0.072577
WT         0.559667
CC        -0.184889
TOTAL     -0.475043
PERCENTAGE -0.475043
dtype: float64
```

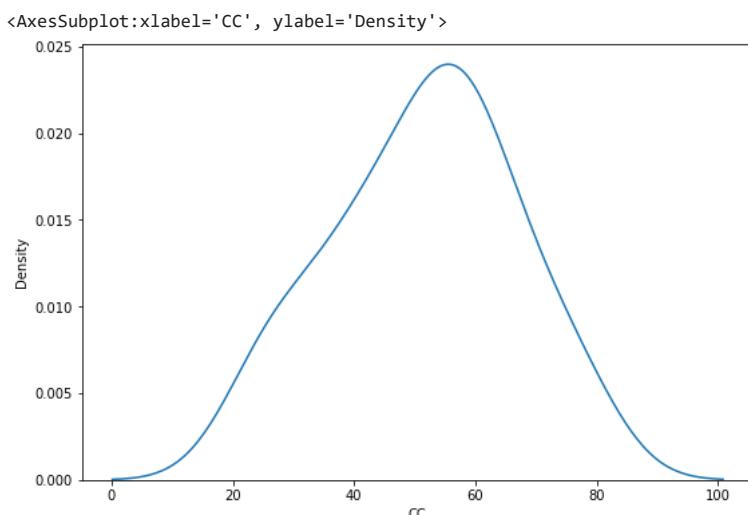
```
sns.kdeplot(om.PERCENTAGE)
```



```
sns.kdeplot(om.DSBDA)
```



```
sns.kdeplot(om.CC)
```



Descriptive Statistics -

Measures of Central Tendency and variability Perform the following operations on any open source dataset (e.g., data.csv)

1. Provide summary statistics (mean, median, minimum, maximum, standard deviation) for a dataset (age, income etc.) with numeric variables grouped by one of the qualitative (categorical) variable. For example, if your categorical variable is age groups and quantitative variable is income, then provide summary statistics of income grouped by the age groups. Create a list that contains a numeric value for each response to the categorical variable.
 2. Write a Python program to display some basic statistical details like percentile, mean, standard deviation etc. of the species of 'Iris-setosa', 'Iris-versicolor' and 'Iris-versicolor' of iris.csv dataset

IMPORTING ALL REQUIRED LIBRARIES

LOADING DATASET

```
iris=pd.read_csv("/content/Iris.csv")
```

iris

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa
3	4	4.6	3.1	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa

DATA PREPROCESSING

iris.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 6 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   Id          150 non-null    int64  
 1   SepalLengthCm 150 non-null   float64 
 2   SepalWidthCm  150 non-null   float64 
 3   PetalLengthCm 150 non-null   float64 
 4   PetalWidthCm  150 non-null   float64 
 5   Species      150 non-null   object  
dtypes: float64(4), int64(1), object(1)
memory usage: 7.2+ KB
```

iris.describe(include='all')

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Spec:
count	150.000000	150.000000	150.000000	150.000000	150.000000	150.000000
unique	Nan	Nan	Nan	Nan	Nan	Nan
top	Nan	Nan	Nan	Nan	Nan	Iris-setosa
freq	Nan	Nan	Nan	Nan	Nan	Nan
mean	75.500000	5.843333	3.054000	3.758667	1.198667	N
std	43.445368	0.828066	0.433594	1.764420	0.763161	N
min	1.000000	4.300000	2.000000	1.000000	0.100000	N
25%	38.250000	5.100000	2.800000	1.600000	0.300000	N
50%	75.500000	5.800000	3.000000	4.350000	1.300000	N
75%	112.750000	6.400000	3.300000	5.100000	1.800000	N

iris.head(10)

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa
3	4	4.6	3.1	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa
5	6	5.4	3.9	1.7	0.4	Iris-setosa
6	7	4.6	3.4	1.4	0.3	Iris-setosa
7	8	5.0	3.4	1.5	0.2	Iris-setosa
8	9	4.4	2.9	1.4	0.2	Iris-setosa
9	10	4.9	3.1	1.5	0.1	Iris-setosa

MEAN

iris.mean()

```
<ipython-input-12-7eed97565d6e>:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') in
iris.mean()
Id      75.500000
SepalLengthCm   5.843333
SepalWidthCm    3.054000
PetalLengthCm   3.758667
PetalWidthCm    1.198667
dtype: float64
```

```
iris.loc[:, 'PetalLengthCm'].mean()
```

```
3.758666666666666
```

```
iris.mean(axis=1)[0:3]
```

```
<ipython-input-15-ac32167c8676>:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') in
iris.mean(axis=1)[0:3]
0    2.24
1    2.30
2    2.48
dtype: float64
```

MEDIAN

```
iris.median()
```

```
<ipython-input-16-7f9b4600b8ea>:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') in
iris.median()
Id      75.50
SepalLengthCm   5.80
SepalWidthCm    3.00
PetalLengthCm   4.35
PetalWidthCm    1.30
dtype: float64
```

```
iris.loc[:, 'SepalWidthCm'].median()
```

```
3.0
```

MODE

```
iris.mode()
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.0	3.0	1.5	0.2	Iris-setosa
1	2	NaN	NaN	NaN	NaN	Iris-versicolor
2	3	NaN	NaN	NaN	NaN	Iris-virginica
3	4	NaN	NaN	NaN	NaN	NaN
4	5	NaN	NaN	NaN	NaN	NaN
...
145	146	NaN	NaN	NaN	NaN	NaN
146	147	NaN	NaN	NaN	NaN	NaN
147	148	NaN	NaN	NaN	NaN	NaN
148	149	NaN	NaN	NaN	NaN	NaN
149	150	NaN	NaN	NaN	NaN	NaN

150 rows × 6 columns

```
iris.loc[:, 'SepalWidthCm'].mode()
```

```
0    3.0
dtype: float64
```

MIN AND MAX

```
iris.min()
```

```

Id           1
SepalLengthCm   4.3
SepalWidthCm    2.0
PetalLengthCm   1.0
PetalWidthCm    0.1
Species      Iris-setosa
dtype: object

```

```
iris.max()
```

```

Id           150
SepalLengthCm   7.9
SepalWidthCm    4.4
PetalLengthCm   6.9
PetalWidthCm    2.5
Species      Iris-virginica
dtype: object

```

STANDARD DEVIATION

```
iris.std()
```

```

<ipython-input-22-c5ab3f85284a>:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated, and will be removed in a future version.
  iris.std()
Id          43.445368
SepalLengthCm  0.828066
SepalWidthCm   0.433594
PetalLengthCm  1.764420
PetalWidthCm   0.763161
dtype: float64

```

```
iris.std(axis=1)[0:3]
```

```

<ipython-input-23-3fedfc1ee20e>:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated, and will be removed in a future version.
  iris.std(axis=1)[0:3]
0    2.010721
1    1.772005
2    1.754138
dtype: float64

```

PRINTING STATISTICAL DATA

```
grouped_data = iris.groupby('Species')
grouped_data
```

```
<pandas.core.groupby.generic.DataFrameGroupBy object at 0x7fe6d68bb730>
```

```

setosa_data = grouped_data.get_group('Iris-setosa')
versicolor_data = grouped_data.get_group('Iris-versicolor')
virginica_data = grouped_data.get_group('Iris-virginica')

```

```

print("Statistical details for Iris-setosa:")
print(setosa_data.describe())
print("\nStatistical details for Iris-versicolor:")
print(versicolor_data.describe())
print("\nStatistical details for Iris-virginica:")
print(virginica_data.describe())

```

```

Statistical details for Iris-setosa:
   Id  SepalLengthCm  SepalWidthCm  PetalLengthCm  PetalWidthCm
count  50.00000     50.00000     50.00000     50.00000
mean   25.50000      5.00000     3.418000    1.464000    0.244000
std    14.57738     0.35249     0.381024    0.173511    0.10721
min    1.00000     4.30000     2.300000    1.000000    0.10000
25%   13.25000     4.80000     3.125000    1.400000    0.20000
50%   25.50000     5.00000     3.400000    1.500000    0.20000
75%   37.75000     5.20000     3.675000    1.575000    0.30000
max   50.00000     5.80000     4.400000    1.900000    0.60000

```

```

Statistical details for Iris-versicolor:
   Id  SepalLengthCm  SepalWidthCm  PetalLengthCm  PetalWidthCm
count  50.00000     50.00000     50.00000     50.00000
mean   75.50000     5.936000    2.770000    4.260000    1.326000
std    14.57738     0.516171    0.313798    0.469911    0.197753
min    51.00000     4.900000    2.000000    3.000000    1.000000
25%   63.25000     5.600000    2.525000    4.000000    1.200000
50%   75.50000     5.900000    2.800000    4.350000    1.300000

```

```
75%    87.75000    6.300000   3.000000   4.600000   1.500000
max    100.00000   7.000000   3.400000   5.100000   1.800000
```

Statistical details for Iris-virginica:

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
count	50.00000	50.00000	50.000000	50.000000	50.00000
mean	125.50000	6.58800	2.974000	5.552000	2.02600
std	14.57738	0.63588	0.322497	0.551895	0.27465
min	101.00000	4.90000	2.200000	4.500000	1.40000
25%	113.25000	6.22500	2.800000	5.100000	1.80000
50%	125.50000	6.50000	3.000000	5.550000	2.00000
75%	137.75000	6.90000	3.175000	5.875000	2.30000
max	150.00000	7.90000	3.800000	6.900000	2.50000

```
print('Iris-setosa')
setosa = iris['Species'] == 'Iris-setosa'
print(iris[setosa].describe())
```

Iris-setosa

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
count	50.00000	50.00000	50.000000	50.000000	50.00000
mean	25.50000	5.00600	3.418000	1.464000	0.24400
std	14.57738	0.35249	0.381024	0.173511	0.10721
min	1.00000	4.30000	2.300000	1.000000	0.10000
25%	13.25000	4.80000	3.125000	1.400000	0.20000
50%	25.50000	5.00000	3.400000	1.500000	0.20000
75%	37.75000	5.20000	3.675000	1.575000	0.30000
max	50.00000	5.80000	4.400000	1.900000	0.60000

```
print('\nIris-versicolor')
versicolor = iris['Species'] == 'Iris-versicolor'
print(iris[versicolor].describe())
```

Iris-versicolor

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
count	50.00000	50.000000	50.000000	50.000000	50.000000
mean	75.50000	5.936000	2.770000	4.260000	1.326000
std	14.57738	0.516171	0.313798	0.469911	0.197753
min	51.00000	4.900000	2.000000	3.000000	1.000000
25%	63.25000	5.600000	2.525000	4.000000	1.200000
50%	75.50000	5.900000	2.800000	4.350000	1.300000
75%	87.75000	6.300000	3.000000	4.600000	1.500000
max	100.00000	7.000000	3.400000	5.100000	1.800000

```
print('\nIris-virginica')
virginica = iris['Species'] == 'Iris-virginica'
print(iris[virginica].describe())
```

Iris-virginica

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
count	50.00000	50.00000	50.000000	50.000000	50.00000
mean	125.50000	6.58800	2.974000	5.552000	2.02600
std	14.57738	0.63588	0.322497	0.551895	0.27465
min	101.00000	4.90000	2.200000	4.500000	1.40000
25%	113.25000	6.22500	2.800000	5.100000	1.80000
50%	125.50000	6.50000	3.000000	5.550000	2.00000
75%	137.75000	6.90000	3.175000	5.875000	2.30000
max	150.00000	7.90000	3.800000	6.900000	2.50000

```
import pandas as pd
```

```
from google.colab import files
files.upload()
```

LOADING DATASET

```
\x00\x86\x9A\xC7\x05\xD0\x89\xB1\xFF\xD4\xB3\xAA\xD4:\x01\xB1\xE7:\xA7\xA4\x1:  
om=pd.read_excel("/content/dsba A3.xlsx")  
\x00\x05\xB7\x8D\x1C\x0A\xDD\xC7\x00\xF0\x83\xD7\x3A\xA7\xB2\xD1\xF8\x1+\x3C\x11  
om
```

DATA PREPROCESSING

```
om.info
```

```
<bound method DataFrame.info of      Id  Age  Gender  Income
 0   1   40    male   32879
 1   2   32  Female  29976
 2   3   47    male  39440
 3   4   49  Female  39453
 4   5   32  Female  24008
 5   6   30  Female  22621
 6   7   50    male  27621
 7   8   40  Female  35625
 8   9   45    male  29312
 9  10   43  Female  36535
10  11   40  Female  25391
11  12   41    male  28793
12  13   39  Female  23070
13  14   32  Female  38285
14  15   35    male  34431
15  16   40  Female  23633
16  17   42    male  21205
17  18   46    male  21464
18  19   45  Female  23773
19  20   31    male  38252>
   Id  Age  Gender  Income
 0   1   40    male   32879
 1   2   32  Female  29976
 2   3   47    male  39440
 3   4   49  Female  39453
 4   5   32  Female  24008
 5   6   30  Female  22621
 6   7   50    male  27621
 7   8   40  Female  35625
 8   9   45    male  29312
 9  10   43  Female  36535
10  11   40  Female  25391
11  12   41    male  28793
12  13   39  Female  23070
13  14   32  Female  38285
14  15   35    male  34431
15  16   40  Female  23633
16  17   42    male  21205
17  18   46    male  21464
18  19   45  Female  23773
19  20   31    male  38252>
```

```
om.Id
```

```
0    1
1    2
2    3
3    4
4    5
5    6
6    7
7    8
8    9
9   10
10  11
11  12
12  13
13  14
14  15
15  16
16  17
17  18
18  19
19  20
Name: Id, dtype: int64
```

```
om.describe(include="all")
```

	Id	Age	Gender	Income
count	20.00000	20.00000	20	20.00000
unique	NaN	NaN	2	NaN
top	NaN	NaN	Female	NaN
freq	NaN	NaN	11	NaN
mean	10.50000	39.95000	NaN	29788.35000
std	5.91608	6.194012	NaN	6542.670912
min	1.00000	30.00000	NaN	21205.00000
25%	5.75000	34.25000	NaN	23738.00000
50%	10.50000	40.00000	NaN	29052.50000
75%	15.25000	45.00000	NaN	35852.50000
max	20.00000	50.00000	NaN	39453.00000

```
summary = om.groupby("Gender")["Income"].describe()
summary
```

```
count      mean      std     min    25%    50%    75%    max
counts = om["Gender"].value_counts().tolist()
counts
```

```
[11, 9]
```

```
om.groupby("Gender")["Income"].mean()
```

```
Gender
Female    29306.363636
male      30377.444444
Name: Income, dtype: float64
```

```
om.groupby("Gender")["Income"].median()
```

```
Gender
Female    25391.0
male      29312.0
Name: Income, dtype: float64
```

```
om.groupby("Gender")["Income"].std()
```

```
Gender
Female    6826.354214
male      6535.044780
Name: Income, dtype: float64
```

```
om.groupby("Gender")["Income"].quantile(0.25)
```

```
Gender
Female    23703.0
male      27621.0
Name: Income, dtype: float64
```

```
om.groupby("Gender")["Income"].quantile(0.50)
```

```
Gender
Female    25391.0
male      29312.0
Name: Income, dtype: float64
```

```
om.groupby("Gender")["Income"].quantile(0.75)
```

```
Gender
Female    36080.0
male      34431.0
Name: Income, dtype: float64
```

```
om.groupby("Gender")["Income"].min()
```

```
Gender
Female    22621
male      21205
Name: Income, dtype: int64
```

```
om.groupby("Gender")["Income"].max()
```

```
Gender
Female    39453
male      39440
Name: Income, dtype: int64
```


Data Analytics I

Create a Linear Regression Model using Python/R to predict home prices using Boston Housing Dataset (<https://www.kaggle.com/c/boston-housing>). The Boston Housing dataset contains information about various houses in Boston through different parameters. There are 506 samples and 14 feature variables in this dataset.

IMPORTING REQUIRED LIBRARIES

```
import pandas as pd
```

```
from google.colab import files  
files.upload()
```

LOADING DATASET

```
om=pd.read_csv("/content/HousingData.csv")
```

```
om
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	L
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1	296	15.3	396.90	
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2	242	17.8	396.90	
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2	242	17.8	392.83	
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63	
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3	222	18.7	396.90	
...	
501	0.06263	0.0	11.93	0.0	0.573	6.593	69.1	2.4786	1	273	21.0	391.99	
502	0.04527	0.0	11.93	0.0	0.573	6.120	76.7	2.2875	1	273	21.0	396.90	
503	0.06076	0.0	11.93	0.0	0.573	6.976	91.0	2.1675	1	273	21.0	396.90	
504	0.10959	0.0	11.93	0.0	0.573	6.794	89.3	2.3889	1	273	21.0	393.45	
505	0.04741	0.0	11.93	0.0	0.573	6.030	NaN	2.5050	1	273	21.0	396.90	

506 rows × 14 columns

DATA PREPROCESSING

```
om.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 506 entries, 0 to 505  
Data columns (total 14 columns):  
 #   Column   Non-Null Count  Dtype     
---  --    
 0   CRIM    486 non-null   float64  
 1   ZN      486 non-null   float64  
 2   INDUS   486 non-null   float64  
 3   CHAS    486 non-null   float64  
 4   NOX     506 non-null   float64  
 5   RM      506 non-null   float64  
 6   AGE     486 non-null   float64  
 7   DIS     506 non-null   float64  
 8   RAD     506 non-null   int64  
 9   TAX     506 non-null   int64  
 10  PTRATIO 506 non-null   float64  
 11  B       506 non-null   float64  
 12  LSTAT   486 non-null   float64  
 13  MEDV    506 non-null   float64  
dtypes: float64(12), int64(2)  
memory usage: 55.5 KB
```

```
om.describe()
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	
count	486.000000	486.000000	486.000000	486.000000	506.000000	506.000000	486.00
mean	3.611874	11.211934	11.083992	0.069959	0.554695	6.284634	68.51
std	8.720192	23.388876	6.835896	0.255340	0.115878	0.702617	27.99
min	0.006320	0.000000	0.460000	0.000000	0.385000	3.561000	2.90
25%	0.081900	0.000000	5.190000	0.000000	0.449000	5.885500	45.17
50%	0.253715	0.000000	9.690000	0.000000	0.538000	6.208500	76.80
75%	3.560263	12.500000	18.100000	0.000000	0.624000	6.623500	93.97

```
om.isnull()
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	E
0	False	False										
1	False	False										
2	False	False										
3	False	False										
4	False	False										
...
501	False	False										
502	False	False										
503	False	False										
504	False	False										
505	False	False	False	False	False	False	True	False	False	False	False	False

506 rows × 14 columns

```
om.fillna(0)
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1	296	15.3	396.90
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2	242	17.8	396.90
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2	242	17.8	392.83
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3	222	18.7	396.90
...
501	0.06263	0.0	11.93	0.0	0.573	6.593	69.1	2.4786	1	273	21.0	391.99
502	0.04527	0.0	11.93	0.0	0.573	6.120	76.7	2.2875	1	273	21.0	396.90
503	0.06076	0.0	11.93	0.0	0.573	6.976	91.0	2.1675	1	273	21.0	396.90
504	0.10959	0.0	11.93	0.0	0.573	6.794	89.3	2.3889	1	273	21.0	393.45
505	0.04741	0.0	11.93	0.0	0.573	6.030	0.0	2.5050	1	273	21.0	396.90

506 rows × 14 columns

CHECKING FOR NULL VALUES

```
om.isnull().sum()
```

CRIM	20
ZN	20
INDUS	20
CHAS	20
NOX	0
RM	0
AGE	20
DIS	0
RAD	0
TAX	0
PTRATIO	0
B	0
LSTAT	20

```
MEDV      0  
dtype: int64
```

FILLING NULL VALUES

```
om['CRIM']=om["CRIM"].fillna(om['CRIM'].mean())
```

```
om.isnull().sum()
```

```
CRIM      0  
ZN        20  
INDUS     20  
CHAS      20  
NOX       0  
RM        0  
AGE       20  
DIS       0  
RAD       0  
TAX       0  
PTRATIO   0  
B         0  
LSTAT     20  
MEDV      0  
dtype: int64
```

```
om['ZN']=om["ZN"].fillna(om['ZN'].mean())
```

```
om['INDUS']=om["INDUS"].fillna(om['INDUS'].mean())
```

```
om['CHAS']=om["CHAS"].fillna(om['CHAS'].mean())
```

```
om['AGE']=om["AGE"].fillna(om['AGE'].mean())
```

```
om['LSTAT']=om["LSTAT"].fillna(om['LSTAT'].mean())
```

```
om.isnull().sum()
```

```
CRIM      0  
ZN        0  
INDUS     0  
CHAS      0  
NOX       0  
RM        0  
AGE       0  
DIS       0  
RAD       0  
TAX       0  
PTRATIO   0  
B         0  
LSTAT     0  
MEDV      0  
dtype: int64
```

DATA VISUALIZATION

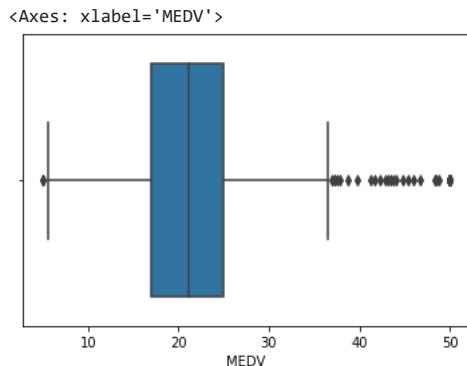
```
import seaborn as sns  
import matplotlib.pyplot as plt
```

CHECKING FOR OUTLIERS

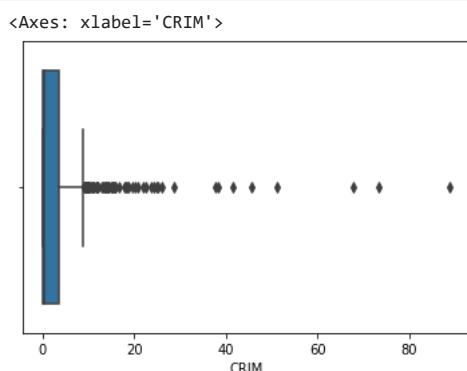
```
om.boxplot()
```

<Axes: >

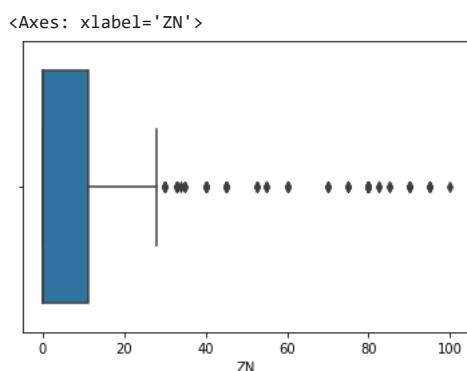
sns.boxplot(x=om.MEDV)



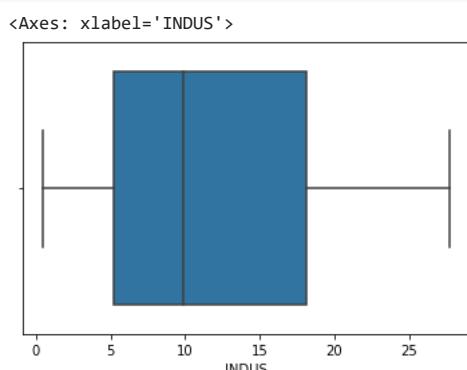
sns.boxplot(x=om.CRIM)



sns.boxplot(x=om.ZN)



sns.boxplot(x=om.INDUS)



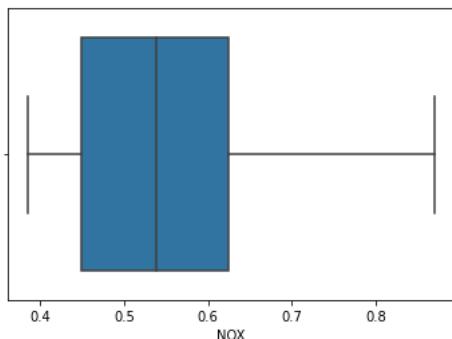
sns.boxplot(x=om.CHAS)

```
<Axes: xlabel='CHAS'>
```



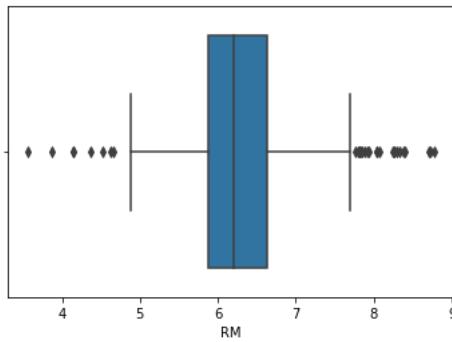
```
sns.boxplot(x=om.NOX)
```

```
<Axes: xlabel='NOX'>
```



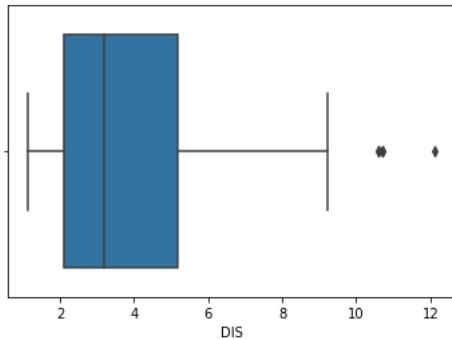
```
sns.boxplot(x=om.RM)
```

```
<Axes: xlabel='RM'>
```



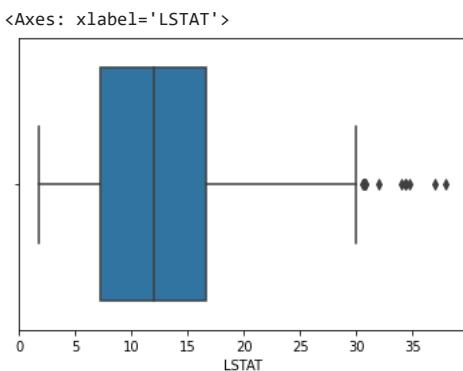
```
sns.boxplot(x=om.DIS)
```

```
<Axes: xlabel='DIS'>
```



```
sns.boxplot(x=om.B)
```

```
<Axes: xlabel='B'>
sns.boxplot(x=om.LSTAT)
```



REMOVING OUTLIERS

```
Q1=om['MEDV'].quantile(0.25)
Q3=om['MEDV'].quantile(0.75)
IQR=Q3-Q1
Lower_limit=Q1 - 1.5*IQR
Upper_limit=Q1 + 1.5*IQR
print(f'Q1={Q1}, Q3={Q3}, IQR={IQR}, LOWER LIMIT={Lower_limit}, UPPER LIMIT={Upper_limit}')
```

```
Q1=17.025, Q3=25.0, IQR=7.97500000000001, LOWER LIMIT=5.062499999999964, UPPER LIMIT=28.9875
```

```
om[(om['MEDV']<Lower_limit) | (om["MEDV"]>Upper_limit)]
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2	242	17.8	392.
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3	222	18.7	394.
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3	222	18.7	396.
39	0.02763	75.0	2.95	0.0	0.428	6.595	21.8	5.4011	3	252	18.3	395.
40	0.03359	75.0	2.95	0.0	0.428	7.024	15.8	5.4011	3	252	18.3	395.
...
371	9.23230	0.0	18.10	0.0	0.631	6.216	100.0	1.1691	24	666	20.2	366.
372	8.26725	0.0	18.10	1.0	0.668	5.875	89.6	1.1296	24	666	20.2	347.
398	38.35180	0.0	18.10	0.0	0.693	5.453	100.0	1.4896	24	666	20.2	396.
405	67.92080	0.0	18.10	0.0	0.693	5.683	100.0	1.4254	24	666	20.2	384.
473	4.64689	0.0	18.10	0.0	0.614	6.980	67.6	2.5329	24	666	20.2	374.

96 rows × 14 columns

```
outliers=[]
for i in om.MEDV:
    if i<Lower_limit or i>Upper_limit:
        outliers.append(i)
print("outliers are",outliers)
```

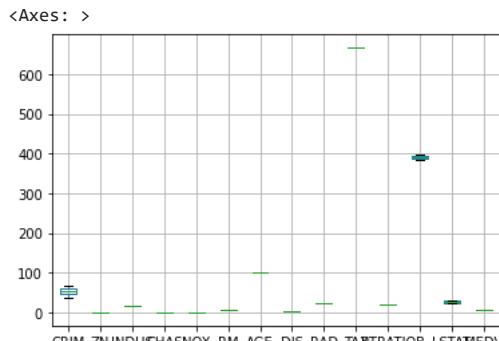
```
outliers are [34.7, 33.4, 36.2, 30.8, 34.9, 35.4, 31.6, 33.0, 38.7, 43.8, 33.2, 41.3, 50.0, 50.0, 50.0, 50.0, 29.4, 29.9, 37.2, 39.
```

```
om1=om.drop(om[om.MEDV<Lower_limit].index)
```

```
om2=om[om.MEDV<Lower_limit]
om2
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B
398	38.3518	0.0	18.1	0.0	0.693	5.453	100.0	1.4896	24	666	20.2	396.90
405	67.9208	0.0	18.1	0.0	0.693	5.683	100.0	1.4254	24	666	20.2	384.97

```
om2.boxplot()
```



Upper_limit

28.9875

Lower_limit

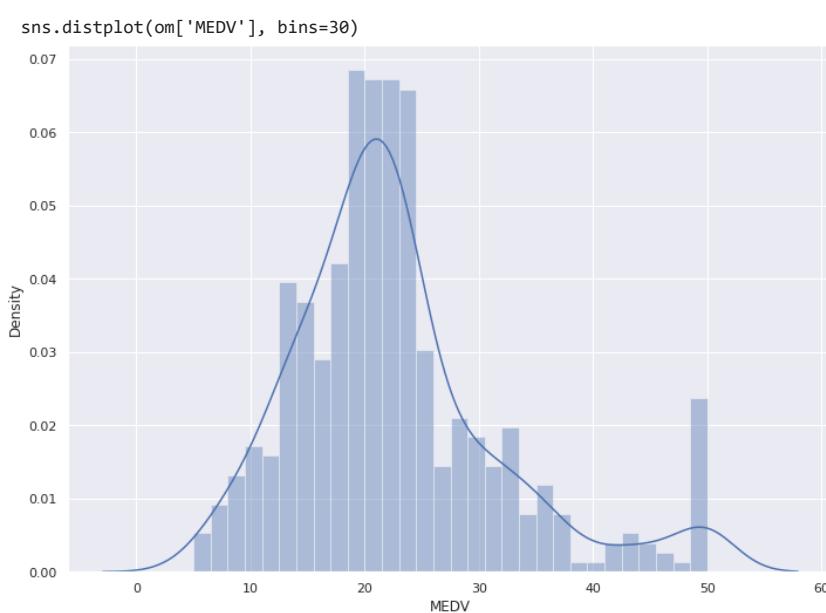
5.0624999999999964

DISPLOT

```
sns.set(rc={'figure.figsize':(11.7,8.27)})
sns.distplot(om['MEDV'], bins=30)
plt.show()
```

<ipython-input-39-d4bd883d1b20>:2: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either `displot` (a figure-level function with
similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see
<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

**CORELATION MATRIX**

```
correlation_matrix = om.corr().round(2)
sns.heatmap(data=correlation_matrix, annot=True)
```

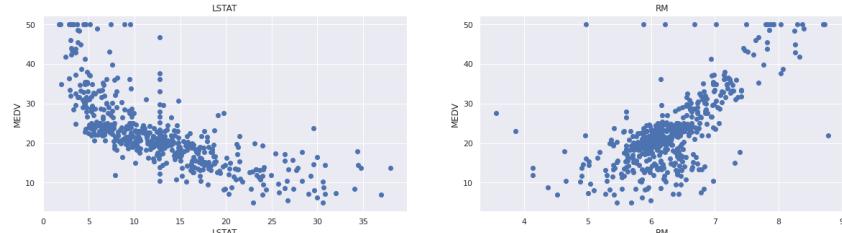
<Axes: >



plt.figure(figsize=(20, 5))

```
features = ['LSTAT', 'RM']
target = om['MEDV']
```

```
for i, col in enumerate(features):
    plt.subplot(1, len(features), i+1)
    x = om[col]
    y = target
    plt.scatter(x, y, marker='o')
    plt.title(col)
    plt.xlabel(col)
    plt.ylabel('MEDV')
```



```
X=om.drop('MEDV',axis=1)
Y=om['MEDV']
```

```
X
Y
```

```
0      24.0
1      21.6
2      34.7
3      33.4
4      36.2
...
501     22.4
502     20.6
503     23.9
504     22.0
505     11.9
Name: MEDV, Length: 506, dtype: float64
```

LINEAR REGRESSION

```
from sklearn.model_selection import train_test_split

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.2, random_state=5)
print(X_train.shape)
print(X_test.shape)
print(Y_train.shape)
print(Y_test.shape)

(404, 13)
(102, 13)
(404,)
(102,)
```

```
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score

lin_model = LinearRegression()
lin_model.fit(X_train, Y_train)
```

LinearRegression
LinearRegression()

lin_model

LinearRegression
LinearRegression()

```
import numpy as np
```

TRAINING AND TESTING

```
y_train_predict = lin_model.predict(X_train)
rmse = (np.sqrt(mean_squared_error(Y_train, y_train_predict)))
r2 = r2_score(Y_train, y_train_predict)

print("The model performance for training set")
print("-----")
print('RMSE is {}'.format(rmse))
print('R2 score is {}'.format(r2))
print("\n")

y_test_predict = lin_model.predict(X_test)
rmse = (np.sqrt(mean_squared_error(Y_test, y_test_predict)))
r2 = r2_score(Y_test, y_test_predict)

print("The model performance for testing set")
print("-----")
print('RMSE is {}'.format(rmse))
print('R2 score is {}'.format(r2))
```

The model performance for training set

RMSE is 4.840940238126834
R2 score is 0.7271916265566591

The model performance for testing set

RMSE is 4.586413823734064
R2 score is 0.7313302802551167

```
mse=mean_squared_error(Y_test,y_test_predict)
print(mse)
mse=mean_squared_error(y_train_predict,Y_train)
print(mse)

21.035191762538922
23.434702389115483
```


Data Analytics II

1. Implement logistic regression using Python/R to perform classification on Social_Network_Ads.csv dataset.
2. Compute Confusion matrix to find TP, FP, TN, FN, Accuracy, Error rate, Precision, Recall on the given dataset

IMPORTING REQUIRED LIBRARIES

Double-click (or enter) to edit

```
import pandas as pd

from google.colab import files
files.upload()

 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 social network ads.csv to social network ads.csv
{'social network ads.csv':
b'Age,EstimatedSalary,Purchased\r\n19,19000,0\r\n35,20000,0\r\n26,43000,0\r\n27,57000,1\r\n...  
395,46,1  
396,51,1  
397,50,1  
398,36,0  
399,49,1'}
```

LOADING DATASET

```
om=pd.read_csv("/content/social network ads.csv")
```

```
om
```

	Age	EstimatedSalary	Purchased
0	19	19000	0
1	35	20000	0
2	26	43000	0
3	27	57000	0
4	19	76000	0
...
395	46	41000	1
396	51	23000	1
397	50	20000	1
398	36	33000	0
399	49	36000	1

400 rows × 3 columns

DATA PREPROCESSING

```
om.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype  
 ---  --          --          --      
 0   Age         400 non-null    int64  
 1   EstimatedSalary 400 non-null    int64  
 2   Purchased    400 non-null    int64  
dtypes: int64(3)
memory usage: 9.5 KB
```

```
om.describe()
```

	Age	EstimatedSalary	Purchased
count	400.000000	400.000000	400.000000
mean	37.655000	69742.500000	0.357500
std	10.482877	34096.960282	0.479864
min	18.000000	15000.000000	0.000000

CHECKING FOR NULL VALUES

om.isNull()

	Age	EstimatedSalary	Purchased
0	False	False	False
1	False	False	False
2	False	False	False
3	False	False	False
4	False	False	False
...
395	False	False	False
396	False	False	False
397	False	False	False
398	False	False	False
399	False	False	False

400 rows × 3 columns

om.isNull().sum

```
<bound method NDFrame._add_numeric_operations.<locals>.sum of           Age  EstimatedSalary  Purchased
0    False            False    False
1    False            False    False
2    False            False    False
3    False            False    False
4    False            False    False
...
395   ...            ...
396   False           False    False
397   False           False    False
398   False           False    False
399   False           False    False

[400 rows x 3 columns]>
```

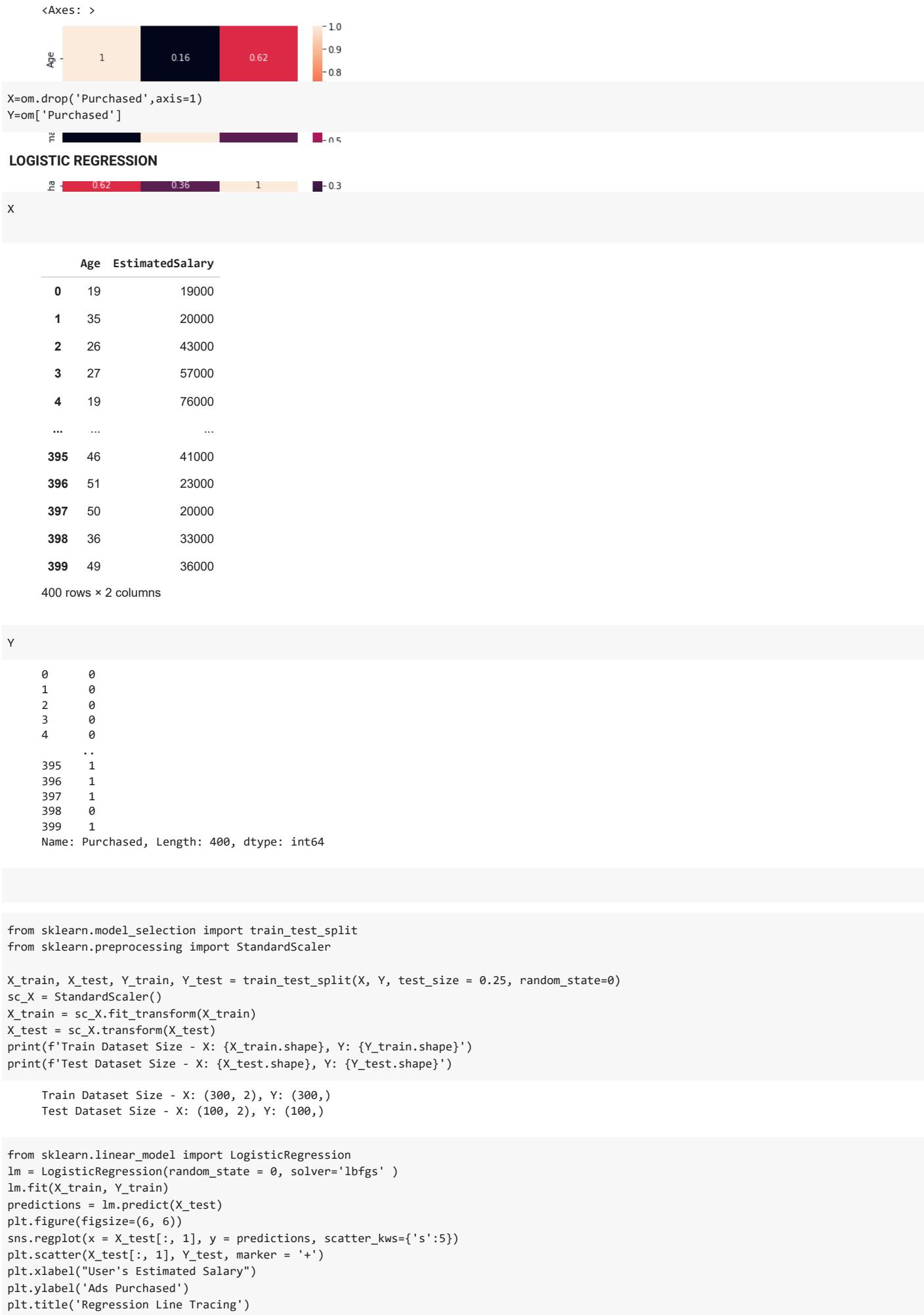
on isnull() sum()

```
Age          0  
EstimatedSalary 0  
Purchased     0  
dtype: int64
```

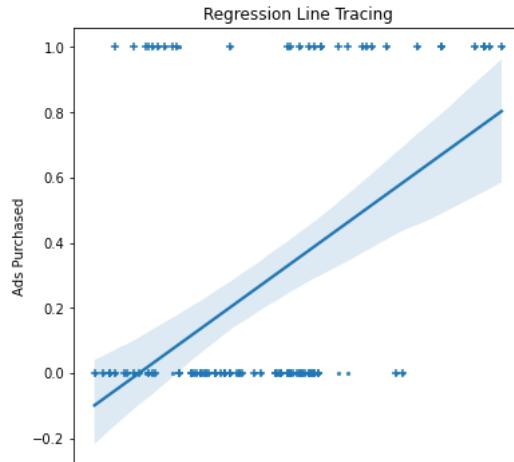
```
import seaborn as sns  
import matplotlib.pyplot as plt
```

CORRELATION MATRIX

```
correlation_matrix = om.corr().round(2)  
sns.heatmap(data=correlation_matrix, annot=True)
```



```
Text(0.5, 1.0, 'Regression Line Tracing')
```



```
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report,precision_score,accuracy_score,recall_score
cm = confusion_matrix(Y_test, predictions)
print(f'''Confusion matrix :\n| Positive Prediction | Negative Prediction\n-----+-----+-----\nPositive Class | True Positive (TP) {cm[0, 0]}\t| False Negative (FN) {cm[0, 1]}\n-----+-----+-----\nNegative Class | False Positive (FP) {cm[1, 0]}\t| True Negative (TN) {cm[1, 1]}\n\n'''')
```

Confusion matrix :

	Positive Prediction	Negative Prediction
Positive Class	True Positive (TP) 65	False Negative (FN) 3
Negative Class	False Positive (FP) 8	True Negative (TN) 24

```
cr = classification_report(Y_test, predictions)
print('Classification report : \n', cr)
```

	precision	recall	f1-score	support
0	0.89	0.96	0.92	68
1	0.89	0.75	0.81	32
accuracy			0.89	100
macro avg	0.89	0.85	0.87	100
weighted avg	0.89	0.89	0.89	100

```
import numpy as np
```

```
precision_score=precision_score(Y_test,predictions)
print(precision_score)
```

0.8888888888888888

```
ac=accuracy_score(Y_test,predictions)
print(ac)
```

0.89

```
rc=recall_score(Y_test,predictions)
print(rc)
```

0.75

```
error_rate=1-ac
print(error_rate)
```

0.1099999999999999



Data Analytics III

1. Implement Simple Naïve Bayes classification algorithm using Python/R on iris.csv dataset.
 2. Compute Confusion matrix to find TP, FP, TN, FN, Accuracy, Error rate, Precision, Recall on the given dataset

IMPORTING REQUIRED LIBRARIES

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

```
from google.colab import files  
files.upload()
```

Choose files 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 Iris.csv to Iris.csv

```
{'Iris.csv': b'Id,SepalLengthCm,SepalWidthCm,PetalLengthCm,PetalWidthCm,Species\n1,5.1,3.5,1.4,0.2,Iris-setosa\n2,4.9,3.0,1.4,0.2,Iris-setosa\n3,4.7,3.2,1.3,0.2,Iris-setosa\n4,4.6,3.1,1.5,0.2,Iris-setosa\n5,5.0,3.6,1.4,0.2,Iris-setosa\n6,5.4,3.9,1.7,0.4,Iris-setosa\n7,4.6,3.4,1.4,0.3,Iris-setosa\n8,5.0,3.4,1.5,0.2,Iris-setosa\n9,4.4,2.9,1.4,0.2,Iris-setosa\n10,4.6,3.1,1.5,0.1,Iris-setosa\n11,5.4,3.7,1.5,0.2,Iris-setosa\n12,4.8,3.4,1.6,0.2,Iris-setosa\n13,4.8,3.0,1.4,0.1,Iris-setosa\n14,4.3,3.0,1.1,0.1,Iris-setosa\n15,5.8,4.0,1.2,0.2,Iris-setosa\n16,5.7,4.4,1.5,0.4,Iris-setosa\n17,5.4,3.9,1.3,0.4,Iris-setosa\n18,5.1,3.5,1.4,0.3,Iris-setosa\n19,5.7,3.8,1.7,0.3,Iris-setosa\n20,5.1,3.8,1.5,0.3,Iris-setosa\n21,5.4,3.4,1.7,0.2,Iris-setosa\n22,5.1,3.7,1.5,0.4,Iris-setosa\n23,4.6,3.6,1.0,0.2,Iris-setosa\n24,5.1,3.3,1.7,0.5,Iris-setosa\n25,4.8,3.4,1.9,0.2,Iris-setosa\n26,5.0,3.0,1.6,0.2,Iris-setosa\n27,5.0,3.4,1.6,0.4,Iris-setosa\n28,5.2,3.5,1.5,0.2,Iris-setosa\n29,5.2,3.4,1.4,0.2,Iris-setosa\n30,4.7,3.2,1.6,0.2,Iris-setosa\n31,4.8,3.1,1.6,0.2,Iris-setosa\n32,5.4,3.4,1.5,0.4,Iris-setosa\n33,5.2,4.1,1.5,0.1,Iris-setosa\n34,5.5,4.2,1.4,0.2,Iris-setosa\n35,4.9,3.1,1.5,0.1,Iris-setosa\n36,5.0,3.2,1.2,0.2,Iris-setosa\n37,5.5,3.5,1.3,0.2,Iris-setosa\n38,4.9,3.1,1.5,0.1,Iris-setosa\n39,4.4,3.0,1.3,0.2,Iris-setosa\n40,5.1,3.4,1.5,0.2,Iris-setosa\n41,5.0,3.5,1.3,0.3,Iris-setosa\n42,4.5,2.3,1.3,0.3,Iris-setosa\n43,4.4,3.2,1.3,0.2,Iris-setosa\n44,5.0,3.5,1.6,0.6,Iris-setosa\n45,5.1,3.8,1.9,0.4,Iris-setosa\n46,4.8,3.0,1.4,0.3,Iris-setosa\n47,5.1,3.8,1.6,0.2,Iris-setosa\n48,4.6,3.2,1.4,0.2,Iris-setosa\n49,5.3,3.7,1.5,0.2,Iris-setosa\n50,5.0,3.3,1.4,0.2,Iris-setosa\n51,7.0,3.2,4.7,1.4,Iris-versicolor\n52,6.4,3.2,4.5,1.5,Iris-versicolor\n53,6.9,3.1,4.9,1.5,Iris-versicolor\n54,5.5,2.3,4.0,1.3,Iris-versicolor\n55,6.5,2.8,4.6,1.5,Iris-versicolor\n56,5.7,2.8,4.5,1.3,Iris-versicolor\n57,6.3,3.3,4.7,1.6,Iris-versicolor\n58,4.9,2.4,3.3,1.0,Iris-versicolor\n59,6.6,2.9,4.6,1.3,Iris-versicolor\n60,5.2,2.7,3.9,1.4,Iris-versicolor\n61,5.0,2.0,3.5,1.0,Iris-versicolor\n62,5.9,3.0,4.2,1.5,Iris-versicolor\n63,6.0,2.2,4.0,1.0,Iris-versicolor\n64,6.1,2.9,4.7,1.4,Iris-versicolor\n65,5.6,2.9,3.6,1.3,Iris-versicolor\n66,6.7,3.1,4.1,1.4,Iris-versicolor\n67,5.6,3.0,4.5,1.5,Iris-versicolor\n68,5.5,2.3,4.0,1.3,Iris-versicolor\n69,5.4,2.5,4.3,1.5,Iris-versicolor\n70,5.7,2.8,4.5,1.6,Iris-versicolor\n71,6.3,3.3,4.7,1.7,Iris-versicolor\n72,6.5,2.9,4.6,1.4,Iris-versicolor\n73,6.7,3.0,4.4,1.5,Iris-versicolor\n74,6.9,3.1,4.8,1.6,Iris-versicolor\n75,7.0,3.2,4.7,1.5,Iris-versicolor\n76,7.2,3.3,4.9,1.7,Iris-versicolor\n77,7.4,3.4,5.0,1.5,Iris-versicolor\n78,7.6,3.5,5.1,1.4,Iris-versicolor\n79,7.8,3.6,5.2,1.3,Iris-versicolor\n80,8.0,3.7,5.3,1.5,Iris-versicolor\n81,8.2,3.8,5.4,1.6,Iris-versicolor\n82,8.4,3.9,5.5,1.7,Iris-versicolor\n83,8.6,4.0,5.6,1.8,Iris-versicolor\n84,8.8,4.1,5.7,1.9,Iris-versicolor\n85,9.0,4.2,5.8,2.0,Iris-versicolor\n86,9.2,4.3,5.9,2.1,Iris-versicolor\n87,9.4,4.4,6.0,2.2,Iris-versicolor\n88,9.6,4.5,6.1,2.3,Iris-versicolor\n89,9.8,4.6,6.2,2.4,Iris-versicolor\n90,10.0,4.7,6.3,2.5,Iris-versicolor\n91,10.2,4.8,6.4,2.6,Iris-versicolor\n92,10.4,4.9,6.5,2.7,Iris-versicolor\n93,10.6,5.0,6.6,2.8,Iris-versicolor\n94,10.8,5.1,6.7,2.9,Iris-versicolor\n95,11.0,5.2,6.8,3.0,Iris-versicolor\n96,11.2,5.3,6.9,3.1,Iris-versicolor\n97,11.4,5.4,7.0,3.2,Iris-versicolor\n98,11.6,5.5,7.1,3.3,Iris-versicolor\n99,11.8,5.6,7.2,3.4,Iris-versicolor\n100,12.0,5.7,7.3,3.5,Iris-versicolor\n101,12.2,5.8,7.4,3.6,Iris-versicolor\n102,12.4,5.9,7.5,3.7,Iris-versicolor\n103,12.6,6.0,7.6,3.8,Iris-versicolor\n104,12.8,6.1,7.7,3.9,Iris-versicolor\n105,13.0,6.2,7.8,4.0,Iris-versicolor\n106,13.2,6.3,7.9,4.1,Iris-versicolor\n107,13.4,6.4,8.0,4.2,Iris-versicolor\n108,13.6,6.5,8.1,4.3,Iris-versicolor\n109,13.8,6.6,8.2,4.4,Iris-versicolor\n110,14.0,6.7,8.3,4.5,Iris-versicolor\n111,14.2,6.8,8.4,4.6,Iris-versicolor\n112,14.4,6.9,8.5,4.7,Iris-versicolor\n113,14.6,7.0,8.6,4.8,Iris-versicolor\n114,14.8,7.1,8.7,4.9,Iris-versicolor\n115,15.0,7.2,8.8,5.0,Iris-versicolor\n116,15.2,7.3,8.9,5.1,Iris-versicolor\n117,15.4,7.4,9.0,5.2,Iris-versicolor\n118,15.6,7.5,9.1,5.3,Iris-versicolor\n119,15.8,7.6,9.2,5.4,Iris-versicolor\n120,16.0,7.7,9.3,5.5,Iris-versicolor\n121,16.2,7.8,9.4,5.6,Iris-versicolor\n122,16.4,7.9,9.5,5.7,Iris-versicolor\n123,16.6,8.0,9.6,5.8,Iris-versicolor\n124,16.8,8.1,9.7,5.9,Iris-versicolor\n125,17.0,8.2,9.8,6.0,Iris-versicolor\n126,17.2,8.3,9.9,6.1,Iris-versicolor\n127,17.4,8.4,10.0,6.2,Iris-versicolor\n128,17.6,8.5,10.1,6.3,Iris-versicolor\n129,17.8,8.6,10.2,6.4,Iris-versicolor\n130,18.0,8.7,10.3,6.5,Iris-versicolor\n131,18.2,8.8,10.4,6.6,Iris-versicolor\n132,18.4,8.9,10.5,6.7,Iris-versicolor\n133,18.6,9.0,10.6,6.8,Iris-versicolor\n134,18.8,9.1,10.7,6.9,Iris-versicolor\n135,19.0,9.2,10.8,7.0,Iris-versicolor\n136,19.2,9.3,10.9,7.1,Iris-versicolor\n137,19.4,9.4,11.0,7.2,Iris-versicolor\n138,19.6,9.5,11.1,7.3,Iris-versicolor\n139,19.8,9.6,11.2,7.4,Iris-versicolor\n140,20.0,9.7,11.3,7.5,Iris-versicolor\n141,20.2,9.8,11.4,7.6,Iris-versicolor\n142,20.4,9.9,11.5,7.7,Iris-versicolor\n143,20.6,10.0,11.6,7.8,Iris-versicolor\n144,20.8,10.1,11.7,7.9,Iris-versicolor\n145,21.0,10.2,11.8,8.0,Iris-versicolor\n146,21.2,10.3,11.9,8.1,Iris-versicolor\n147,21.4,10.4,12.0,8.2,Iris-versicolor\n148,21.6,10.5,12.1,8.3,Iris-versicolor\n149,21.8,10.6,12.2,8.4,Iris-versicolor\n150,22.0,10.7,12.3,8.5,Iris-versicolor\n151,22.2,10.8,12.4,8.6,Iris-versicolor\n152,22.4,10.9,12.5,8.7,Iris-versicolor\n153,22.6,11.0,12.6,8.8,Iris-versicolor\n154,22.8,11.1,12.7,8.9,Iris-versicolor\n155,23.0,11.2,12.8,9.0,Iris-versicolor\n156,23.2,11.3,12.9,9.1,Iris-versicolor\n157,23.4,11.4,13.0,9.2,Iris-versicolor\n158,23.6,11.5,13.1,9.3,Iris-versicolor\n159,23.8,11.6,13.2,9.4,Iris-versicolor\n160,24.0,11.7,13.3,9.5,Iris-versicolor\n161,24.2,11.8,13.4,9.6,Iris-versicolor\n162,24.4,11.9,13.5,9.7,Iris-versicolor\n163,24.6,12.0,13.6,9.8,Iris-versicolor\n164,24.8,12.1,13.7,9.9,Iris-versicolor\n165,25.0,12.2,13.8,10.0,Iris-versicolor\n166,25.2,12.3,13.9,10.1,Iris-versicolor\n167,25.4,12.4,14.0,10.2,Iris-versicolor\n168,25.6,12.5,14.1,10.3,Iris-versicolor\n169,25.8,12.6,14.2,10.4,Iris-versicolor\n170,26.0,12.7,14.3,10.5,Iris-versicolor\n171,26.2,12.8,14.4,10.6,Iris-versicolor\n172,26.4,12.9,14.5,10.7,Iris-versicolor\n173,26.6,13.0,14.6,10.8,Iris-versicolor\n174,26.8,13.1,14.7,10.9,Iris-versicolor\n175,27.0,13.2,14.8,11.0,Iris-versicolor\n176,27.2,13.3,14.9,11.1,Iris-versicolor\n177,27.4,13.4,15.0,11.2,Iris-versicolor\n178,27.6,13.5,15.1,11.3,Iris-versicolor\n179,27.8,13.6,15.2,11.4,Iris-versicolor\n180,28.0,13.7,15.3,11.5,Iris-versicolor\n181,28.2,13.8,15.4,11.6,Iris-versicolor\n182,28.4,13.9,15.5,11.7,Iris-versicolor\n183,28.6,14.0,15.6,11.8,Iris-versicolor\n184,28.8,14.1,15.7,11.9,Iris-versicolor\n185,29.0,14.2,15.8,12.0,Iris-versicolor\n186,29.2,14.3,15.9,12.1,Iris-versicolor\n187,29.4,14.4,16.0,12.2,Iris-versicolor\n188,29.6,14.5,16.1,12.3,Iris-versicolor\n189,29.8,14.6,16.2,12.4,Iris-versicolor\n190,30.0,14.7,16.3,12.5,Iris-versicolor\n191,30.2,14.8,16.4,12.6,Iris-versicolor\n192,30.4,14.9,16.5,12.7,Iris-versicolor\n193,30.6,15.0,16.6,12.8,Iris-versicolor\n194,30.8,15.1,16.7,12.9,Iris-versicolor\n195,31.0,15.2,16.8,13.0,Iris-versicolor\n196,31.2,15.3,16.9,13.1,Iris-versicolor\n197,31.4,15.4,17.0,13.2,Iris-versicolor\n198,31.6,15.5,17.1,13.3,Iris-versicolor\n199,31.8,15.6,17.2,13.4,Iris-versicolor\n200,32.0,15.7,17.3,13.5,Iris-versicolor\n201,32.2,15.8,17.4,13.6,Iris-versicolor\n202,32.4,15.9,17.5,13.7,Iris-versicolor\n203,32.6,16.0,17.6,13.8,Iris-versicolor\n204,32.8,16.1,17.7,13.9,Iris-versicolor\n205,33.0,16.2,17.8,14.0,Iris-versicolor\n206,33.2,16.3,17.9,14.1,Iris-versicolor\n207,33.4,16.4,18.0,14.2,Iris-versicolor\n208,33.6,16.5,18.1,14.3,Iris-versicolor\n209,33.8,16.6,18.2,14.4,Iris-versicolor\n210,34.0,16.7,18.3,14.5,Iris-versicolor\n211,34.2,16.8,18.4,14.6,Iris-versicolor\n212,34.4,16.9,18.5,14.7,Iris-versicolor\n213,34.6,17.0,18.6,14.8,Iris-versicolor\n214,34.8,17.1,18.7,14.9,Iris-versicolor\n215,35.0,17.2,18.8,15.0,Iris-versicolor\n216,35.2,17.3,18.9,15.1,Iris-versicolor\n217,35.4,17.4,19.0,15.2,Iris-versicolor\n218,35.6,17.5,19.1,15.3,Iris-versicolor\n219,35.8,17.6,19.2,15.4,Iris-versicolor\n220,36.0,17.7,19.3,15.5,Iris-versicolor\n221,36.2,17.8,19.4,15.6,Iris-versicolor\n222,36.4,17.9,19.5,15.7,Iris-versicolor\n223,36.6,18.0,19.6,15.8,Iris-versicolor\n224,36.8,18.1,19.7,15.9,Iris-versicolor\n225,37.0,18.2,19.8,16.0,Iris-versicolor\n226,37.2,18.3,19.9,16.1,Iris-versicolor\n227,37.4,18.4,20.0,16.2,Iris-versicolor\n228,37.6,18.5,20.1,16.3,Iris-versicolor\n229,37.8,18.6,20.2,16.4,Iris-versicolor\n230,38.0,18.7,20.3,16.5,Iris-versicolor\n231,38.2,18.8,20.4,16.6,Iris-versicolor\n232,38.4,18.9,20.5,16.7,Iris-versicolor\n233,38.6,19.0,20.6,16.8,Iris-versicolor\n234,38.8,19.1,20.7,16.9,Iris-versicolor\n235,39.0,19.2,20.8,17.0,Iris-versicolor\n236,39.2,19.3,20.9,17.1,Iris-versicolor\n237,39.4,19.4,21.0,17.2,Iris-versicolor\n238,39.6,19.5,21.1,17.3,Iris-versicolor\n239,39.8,19.6,21.2,17.4,Iris-versicolor\n240,40.0,19.7,21.3,17.5,Iris-versicolor\n241,40.2,19.8,21.4,17.6,Iris-versicolor\n242,40.4,19.9,21.5,17.7,Iris-versicolor\n243,40.6,19.8,21.6,17.8,Iris-versicolor\n244,40.8,19.7,21.7,17.9,Iris-versicolor\n245,41.0,19.6,21.8,18.0,Iris-versicolor\n246,41.2,19.5,21.9,18.1,Iris-versicolor\n247,41.4,19.4,22.0,18.2,Iris-versicolor\n248,41.6,19.3,22.1,18.3,Iris-versicolor\n249,41.8,19.2,22.2,18.4,Iris-versicolor\n250,42.0,19.1,22.3,18.5,Iris-versicolor\n251,42.2,19.0,22.4,18.6,Iris-versicolor\n252,42.4,18.9,22.5,18.7,Iris-versicolor\n253,42.6,18.8,22.6,18.8,Iris-versicolor\n254,42.8,18.7,22.7,18.9,Iris-versicolor\n255,43.0,18.6,22.8,19.0,Iris-versicolor\n256,43.2,18.5,22.9,19.1,Iris-versicolor\n257,43.4,18.4,23.0,19.2,Iris-versicolor\n258,43.6,18.3,23.1,19.3,Iris-versicolor\n259,43.8,18.2,23.2,19.4,Iris-versicolor\n260,44.0,18.1,23.3,19.5,Iris-versicolor\n261,44.2,18.0,23.4,19.6,Iris-versicolor\n262,44.4,17.9,23.5,19.7,Iris-versicolor\n263,44.6,17.8,23.6,19.8,Iris-versicolor\n264,44.8,17.7,23.7,19.9,Iris-versicolor\n265,45.0,17.6,23.8,20.0,Iris-versicolor\n266,45.2,17.5,23.9,20.1,Iris-versicolor\n267,45.4,17.4,24.0,20.2,Iris-versicolor\n268,45.6,17.3,24.1,20.3,Iris-versicolor\n269,45.8,17.2,24.2,20.4,Iris-versicolor\n270,46.0,17.1,24.3,20.5,Iris-versicolor\n271,46.2,17.0,24.4,20.6,Iris-versicolor\n272,46.4,16.9,24.5,20.7,Iris-versicolor\n273,46.6,16.8,24.6,20.8,Iris-versicolor\n274,46.8,16.7,24.7,20.9,Iris-versicolor\n275,47.0,16.6,24.8,21.0,Iris-versicolor\n276,47.2,16.5,24.9,21.1,Iris-versicolor\n277,47.4,16.4,25.0,21.2,Iris-versicolor\n278,47.6,16.3,25.1,21.3,Iris-versicolor\n279,47.8,16.2,25.2,21.4,Iris-versicolor\n280,48.0,16.1,25.3,21.5,Iris-versicolor\n281,48.2,16.0,25.4,21.6,Iris-versicolor\n282,48.4,15.9,25.5,21.7,Iris-versicolor\n283,48.6,15.8,25.6,21.8,Iris-versicolor\n284,48.8,15.7,25.7,21.9,Iris-versicolor\n285,49.0,15.6,25.8,22.0,Iris-versicolor\n286,49.2,15.5,25.9,22.1,Iris-versicolor\n287,49.4,15.4,26.0,22.2,Iris-versicolor\n288,49.6,15.3,26.1,22.3,Iris-versicolor\n289,49.8,15.2,26.2,22.4,Iris-versicolor\n290,50.0,15.1,26.3,22.5,Iris-versicolor\n291,50.2,15.0,26.4,22.6,Iris-versicolor\n292,50.4,14.9,26.5,22.7,Iris-versicolor\n293,50.6,14.8,26.6,22.8,Iris-versicolor\n294,50.8,14.7,26.7,22.9,Iris-versicolor\n295,51.0,14.6,26.8,23.0,Iris-versicolor\n296,51.2,14.5,26.9,23.1,Iris-versicolor\n297,51.4,14.4,27.0,23.2,Iris-versicolor\n298,51.6,14.3,27.1,23.3,Iris-versicolor\n299,51.8,14.2,27.2,23.4,Iris-versicolor\n300,52.0,14.1,27.3,23.5,Iris-versicolor\n301,52.2,14.0,27.4,23.6,Iris-versicolor\n302,52.4,13.9,27.5,23.7,Iris-versicolor\n303,52.6,13.8,27.6,23.8,Iris-versicolor\n304,52.8,13.7,27.7,23.9,Iris-versicolor\n305,53.0,13.6,27.8,24.0,Iris-versicolor\n306,53.2,13.5,27.9,24.1,Iris-versicolor\n307,53.4,13.4,28.0,24.2,Iris-versicolor\n308,53.6,13.3,28.1,24.3,Iris-versicolor\n309,53.8,13.2,28.2,24.4,Iris-versicolor\n310,54.0,13.1,28.3,24.5,Iris-versicolor\n311,54.2,13.0,28.4,24.6,Iris-versicolor\n312,54.4,12.9,28.5,24.7,Iris-versicolor\n313,54.6,12.8,28.6,24.8,Iris-versicolor\n314,54.8,12.7,28.7,24.9,Iris-versicolor\n315,55.0,12.6,28.8,25.0,Iris-versicolor\n316,55.2,12.5,28.9,25.1,Iris-versicolor\n317,55.4,12.4,29.0,25.2,Iris-versicolor\n318,55.6,12.3,29.1,25.3,Iris-versicolor\n319,55.8,12.2,29.2,25.4,Iris-versicolor\n320,56.0,12.1,29.3,25.5,Iris-versicolor\n321,56.2,12.0,29.4,25.6,Iris-versicolor\n322,56.4,11.9,29.5,25.7,Iris-versicolor\n323,56.6,11.8,29.6,25.8,Iris-versicolor\n324,56.8,11.7,29.7,25.9,Iris-versicolor\n325,57.0,11.6,29.8,26.0,Iris-versicolor\n326,57.2,11.5,29.9,26.1,Iris-versicolor\n327,57.4,11.4,30.0,26.2,Iris-versicolor\n328,57.6,11.3,30.1,26.3,Iris-versicolor\n329,57.8,11.2,30.2,26.4,Iris-versicolor\n330,58.0,11.1,30.3,26.5,Iris-versicolor\n331,58.2,11.0,30.4,26.6,Iris-versicolor\n332,58.4,10.9,30.5,26.7,Iris-versicolor\n333,58.6,10.8,30.6,26.8,Iris-versicolor\n334,58.8,10.7,30.7,26.9,Iris-versicolor\n335,59.0,10.6,30.8,27.0,Iris-versicolor\n336,59.2,10.5,30.9,27.1,Iris-versicolor\n337,59.4,10.4,31.0,27.2,Iris-versicolor\n338,59.6,10.3,31.1,27.3,Iris-versicolor\n339,59.8,10.2,31.2,27.4,Iris-versicolor\n340,60.0,10.1,31.3,27.5,Iris-versicolor\n341,60.2,10.0,31.4,27.6,Iris-versicolor\n342,60.4,9.9,31.5,27.7,Iris-versicolor\n343,60.6,9.8,31.6,27.8,Iris-versicolor\n344,60.8,9.7,31.7,27.9,Iris-versicolor\n345,61.0,9.6,31.8,28.0,Iris-versicolor\n346,61.2,9.5,31.9,28.1,Iris-versicolor\n347,61.4,9.4,32.0,28.2,Iris-versicolor\n348,61.6,9.3,32.1,28.3,Iris-versicolor\n349,61.8,9.2,32.2,28.4,Iris-versicolor\n350,62.0,9.1,32.3,28.5,Iris-versicolor\n351,62.2,9.0,32.4,28.6,Iris-versicolor\n352,62.4,8.9,32.5,28.7,Iris-versicolor\n353,62.6,8.8,32.6,28.8,Iris-versicolor\n354,62.8,8.7,32.7,28.9,Iris-versicolor\n355,63.0,8.6,32.8,29.0,Iris-versicolor\n356,63.2,8.5,32.9,29.1,Iris-versicolor\n357,63.4,8.4,33.0,29.2,Iris-versicolor\n358,63.6,8.3,33.1,29.3,Iris-versicolor\n359,63.8,8.2,33.2,29.4,Iris-versicolor\n360,64.0,8.1,33.3,29.5,Iris-versicolor\n361,64.2,8.0,33.4,29.6,Iris-versicolor\n362,64.4,7.9,33.5,29.7,Iris-versicolor\n363,64.6,7.8,33.6,29.8,Iris-versicolor\n364,64.8,7.7,33.7,29.9,Iris-versicolor\n365,65.0,7.6,33.8,30.0,Iris-versicolor\n366,65.2,7.5,33.9,30.1,Iris-versicolor\n367,65.4,7.4,34.0,30.2,Iris-versicolor\n368,65.6,7.3,34.1,30.3,Iris-versicolor\n369,65.8,7.2,34.2,30.4,Iris-versicolor\n370,66.0,7.1,34.3,30.5,Iris-versicolor\n371,66.2,7.0,34.4,30.6,Iris-versicolor\n372,66.4,6.9,34.5,30.7,Iris-versicolor\n373,66.6,6.8,34.6,30.8,Iris-versicolor\n374,66.8,6.7,34.7,30.9,Iris-versicolor\n375,67.0,6.6,34.8,31.0,Iris-versicolor\n376,67.2,6.5,34.9,31.1,Iris-versicolor\n377,67.4,6.4,35.0,31.2,Iris-versicolor\n378,67.6,6.3,35.1,31.3,Iris-versicolor\n379,67.8,6.2,35.2,31.4,Iris-versicolor\n380,68.0,6.1,35.3,31.5,Iris-versicolor\n381,68.2,6.0,35.4,31.6,Iris-versicolor\n382,68.4,5.9,35.5,31.7,Iris-versicolor\n383,68.6,5.8,35.6,31.8,Iris-versicolor\n384,68.8,5.7,35.7,31.9,Iris-versicolor\n385,69.0,5.6,35.8,32.0,Iris-versicolor\n386,69.2,5.5,35.9,32.1,Iris-versicolor\n387,69.4,5.4,36.0,32.2,Iris-versicolor\n388,69.6,5.3,36.1,32.3,Iris-versicolor\n389,69.8,5.2,36.2,32.4,Iris-versicolor\n390,70.0,5.1,36.3,32.5,Iris-versicolor\n391,70.2,5.0,36.4,32.6,Iris-versicolor\n392,70.4,4.9,36.5,32.7,Iris-versicolor\n393,70.6,4.8,36.6,32.8,Iris-versicolor\n394,70.8,4.7,36.7,32.9,Iris-versicolor\n395,71.0,4.6,36.8,33.0,Iris-versicolor\n396,71.2,4.5,36.9,33.1,Iris-versicolor\n397,71.4,4.4,37.0,33.2,Iris-versicolor\n398,71.6,4.3,37.1,33.3,Iris-versicolor\n399,71.8,4.2,37.2,33.4,Iris-versicolor\n400,72.0,4.1,37.3,33.5,Iris-versicolor\n401,72.2,4.0,37.4,33.6,Iris-versicolor\n402,72.4,3.9,37.5,33.7,Iris-versicolor\n403,72.6,3.8,37.6,33.8,Iris-versicolor\n404,72.8,3.7,37.7,33.9,Iris-versicolor\n405,73.0,3.6,37.8,34.0,Iris-versicolor\n406,73.2,3.5,37.9,34.1,Iris-versicolor\n407,73.4,3.4,38.0,34.2,Iris-versicolor\n408,73.6,3.3,38.1,34.3,Iris-versicolor\n409,73.8,3.2,38.2,34.4,Iris-versicolor\n410,74.0,3.1,38.3,34.5,Iris-versicolor\n411,74.2,3.0,38.4,34.6,Iris-versicolor\n412,74.4,2.9,38.5,34.7,Iris-versicolor\n413,74.6,2.8,38.6,34.8,Iris-versicolor\n414,74.8,2.7,38.7,34.9,Iris-versicolor\n415,75.0,2.6,38.8,35.0,Iris-versicolor\n416,75.2,2.5,38.9,35.1,Iris-versicolor\n417,75.4,2.4,39.0,35.2,Iris-versicolor\n418,75.6,2.3,39.1,35.3,Iris-versicolor\n419,75.8,2.2,39.2,35.4,Iris-versicolor\n420,76.0,2.1,39.3,35.5,Iris-versicolor\n421,76.2,2.0,39.4,35.6,Iris-versicolor\n422,76.4,1.9,39.5,35.7,Iris-versicolor\n423,76.6,1.8,39.6,35.8,Iris-versicolor\n424,76.8,1.7,39.7,35.9,Iris-versicolor\n425,77.0,1.6,39.8,36.0,Iris-versicolor\n426,77.2,1.5,39.9,36.1,Iris-versicolor\n427,77.4,1.4,40.0,36.2,Iris-versicolor\n428,77.6,1.3,40.1,36.3,Iris-versicolor\n429,77.8,1.2,40.2,36.4,Iris-versicolor\n430,78.0,1.1,40.3,36.5,Iris-versicolor\n431,78.2,1.0,40.4,36.6,Iris-versicolor\n432,78.4,0.9,40.5,36.7,Iris-versicolor\n433,78.6,0.8,40.6,36.8,Iris-versicolor\n434,78.8,0.7,40.7,36.9,Iris-versicolor\n435,79.0,0.6,40.8,37.0,Iris-versicolor\n436,79.2,0.5,40.9,37.1,Iris-versicolor\n437,79.4,0.4,41.0,37.2,Iris-versicolor\n438,79.6,0.3,41.1,37.3,Iris-versicolor\n439,79.8,0.2,41.2,37.4,Iris-versicolor\n440,80.0,0.1,41.3,37.5,Iris-versicolor\n441,80.2,0.0,41.4,37.6,Iris-versicolor\n442,80.4,0.1,41.5,37.7,Iris-versicolor\n443,80.6,0.2,41.6,37.8,Iris-versicolor\n444,80.8,0.3,41.7,37.9,Iris-versicolor\n445,81.0,0.4,41.8,38.0,Iris-versicolor\n446,81.2,0.5,41.9,38.1,Iris-versicolor\n447,81.4,0.6,42.0,38.2,Iris-versicolor\n448,81.6,0.7,42.1,38.3,Iris-versicolor\n449,81.8,0.8,42.2,38.4,Iris-versicolor\n450,82.0,0.9,42.3,38.5,Iris-versicolor\n451,82.2,0.8,42.4,38.6,Iris-versicolor\n452,82.4,0.7,42.5,38.7,Iris-versicolor\n453,82.6,0.6,42.6,38.8,Iris-versicolor\n454,82.8,0.5,42.7,38.9,Iris-versicolor\n455,83.0,0.4,42.8,39.0,Iris-versicolor\n456,83.2,0.3,42.9,39.1,Iris-versicolor\n457,83.4,0.2,43.0,39.2,Iris-versicolor\n458,83.6,0.1,43.1,39.3,Iris-versicolor\n459,83.8,0.0,43.2,39.4,Iris-versicolor\n460,84.0,0.1,43.3,39.5,Iris-versicolor\n461,84.2,0.2,43.4,39.6,Iris-versicolor\n462,84.4,0.3,43.5,39.7,Iris-versicolor\n463,84.6,0.4,43.6,39.8,Iris-versicolor\n464,84.8,0.5,43.7,39.9,Iris-versicolor\n465,85.0,0.6,43.8,40.0,Iris-versicolor\n466,85.2,0.
```

LOADING DATASET

```
df=pd.read_csv("Iris.csv")
```

DATA PREPROCESSING

```
df.head()
```

Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.1	3.5	1.4	0.2 Iris-setosa
1	2	4.9	3.0	1.4	0.2 Iris-setosa
2	3	4.7	3.2	1.3	0.2 Iris-setosa
3	4	4.6	3.1	1.5	0.2 Iris-setosa
4	5	5.0	3.6	1.4	0.2 Iris-setosa

```
df.describe()
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
count	150.000000	150.000000	150.000000	150.000000	150.000000
mean	75.500000	5.843333	3.054000	3.758667	1.198667
std	43.445368	0.828066	0.433594	1.764420	0.763161
min	1.000000	4.300000	2.000000	1.000000	0.100000
25%	38.250000	5.100000	2.800000	1.600000	0.300000
50%	75.500000	5.800000	3.000000	4.350000	1.300000
75%	110.750000	6.100000	3.300000	5.100000	1.800000

CHECKING FOR NULL VALUES

```
df.isnull().sum()
```

```
Id          0
SepalLengthCm  0
SepalWidthCm   0
PetalLengthCm  0
PetalWidthCm   0
Species        0
dtype: int64
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 6 columns):
 #   Column      Non-Null Count  Dtype  
 --- 
 0   Id          150 non-null    int64  
 1   SepalLengthCm 150 non-null    float64 
 2   SepalWidthCm  150 non-null    float64 
 3   PetalLengthCm 150 non-null    float64 
 4   PetalWidthCm  150 non-null    float64 
 5   Species      150 non-null    object  
dtypes: float64(4), int64(1), object(1)
memory usage: 7.2+ KB
```

```
X = df.drop(['Species'], axis = 1)
Y = df['Species']
```

NAIVE BAYES CLASSIFICATION

```
from sklearn.model_selection import train_test_split
xtrain, xtest, ytrain, ytest = train_test_split(X, Y, test_size = 0.2, random_state = 0)
```

```
from sklearn.naive_bayes import GaussianNB
gaussian = GaussianNB()
```

```
gaussian.fit(xtrain, ytrain)
```

```
▼ GaussianNB
GaussianNB()
```

```
y_pred = gaussian.predict(xtest)
```

```
print(xtrain)
print("-----\n")
print(xtest)
print("-----\n")
print(ytrain)
print("-----\n")
print(ytest)
print("-----\n")
print(y_pred)
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
137	138	6.4	3.1	5.5	1.8
84	85	5.4	3.0	4.5	1.5
27	28	5.2	3.5	1.5	0.2
127	128	6.1	3.0	4.9	1.8
132	133	6.4	2.8	5.6	2.2
...

```

9      10      4.9      3.1      1.5      0.1
103    104      6.3      2.9      5.6      1.8
67     68      5.8      2.7      4.1      1.0
117    118      7.7      3.8      6.7      2.2
47     48      4.6      3.2      1.4      0.2

```

[120 rows x 5 columns]

```

Id SepalLengthCm SepalWidthCm PetalLengthCm PetalWidthCm
114 115      5.8      2.8      5.1      2.4
62   63      6.0      2.2      4.0      1.0
33   34      5.5      4.2      1.4      0.2
107 108      7.3      2.9      6.3      1.8
7    8       5.0      3.4      1.5      0.2
100 101      6.3      3.3      6.0      2.5
40   41      5.0      3.5      1.3      0.3
86   87      6.7      3.1      4.7      1.5
76   77      6.8      2.8      4.8      1.4
71   72      6.1      2.8      4.0      1.3
134 135      6.1      2.6      5.6      1.4
51   52      6.4      3.2      4.5      1.5
73   74      6.1      2.8      4.7      1.2
54   55      6.5      2.8      4.6      1.5
63   64      6.1      2.9      4.7      1.4
37   38      4.9      3.1      1.5      0.1
78   79      6.0      2.9      4.5      1.5
90   91      5.5      2.6      4.4      1.2
45   46      4.8      3.0      1.4      0.3
16   17      5.4      3.9      1.3      0.4
121 122      5.6      2.8      4.9      2.0
66   67      5.6      3.0      4.5      1.5
24   25      4.8      3.4      1.9      0.2
8    9       4.4      2.9      1.4      0.2
126 127      6.2      2.8      4.8      1.8
22   23      4.6      3.6      1.0      0.2
44   45      5.1      3.8      1.9      0.4
97   98      6.2      2.9      4.3      1.3
93   94      5.0      2.3      3.3      1.0
26   27      5.0      3.4      1.6      0.4

```

```

137     Iris-virginica
84      Iris-versicolor
27      Iris-setosa
127    Iris-virginica
132    Iris-virginica
...
9       Iris-setosa
103   Iris-virginica
67     Iris-versicolor

```

CONFUSION MATRIX

```

from sklearn.metrics import precision_score,confusion_matrix,accuracy_score,recall_score, classification_report
cm= confusion_matrix(ytest, y_pred)

```

cm

```

array([[11,  0,  0],
 [ 0, 13,  0],
 [ 0,  0,  6]])

```

```

from sklearn.metrics import accuracy_score

```

```

print ("Accuracy : ", accuracy_score(ytest, y_pred))

```

Accuracy : 1.0

```

error_rate = 1- accuracy_score(ytest, y_pred)

```

error_rate

0.0

```

print("classification report: ",classification_report(ytest, y_pred))

```

	precision	recall	f1-score	support
Iris-setosa	1.00	1.00	1.00	11
Iris-versicolor	1.00	1.00	1.00	13
Iris-virginica	1.00	1.00	1.00	6

accuracy		1.00		30
macro avg	1.00	1.00	1.00	30
weighted avg	1.00	1.00	1.00	30

ACCURACY SCORE, PRECISION SCORE, RECALL SCORE

```
accuracy = accuracy_score(ytest,y_pred)
precision =precision_score(ytest, y_pred,average='macro')
recall = recall_score(ytest, y_pred,average='macro')
```

```
accuracy
```

```
1.0
```

```
precision
```

```
1.0
```

```
recall
```

```
1.0
```

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



Text Analytics

1. Extract Sample document and apply following document preprocessing methods: Tokenization, POS Tagging, stop words removal, Stemming and Lemmatization.
2. Create representation of document by calculating Term Frequency and Inverse Document Frequency.

IMPORTING REQUIRED LIBRARIES

```
pip install nltk
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: nltk in /usr/local/lib/python3.9/dist-packages (3.8.1)
Requirement already satisfied: click in /usr/local/lib/python3.9/dist-packages (from nltk) (8.1.3)
Requirement already satisfied: tqdm in /usr/local/lib/python3.9/dist-packages (from nltk) (4.65.0)
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.9/dist-packages (from nltk) (2022.10.31)
Requirement already satisfied: joblib in /usr/local/lib/python3.9/dist-packages (from nltk) (1.2.0)
```

```
import nltk
```

```
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')
nltk.download('averaged_perceptron_tagger')

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]  Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]  Unzipping corpora/stopwords.zip.
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]      /root/nltk_data...
[nltk_data]  Unzipping taggers/averaged_perceptron_tagger.zip.
True
```

TOKENIZATION

```
text= "Tokenization is the first step in text analytics. The process of breaking down a text paragraph into smaller chunks such as words
```

```
from nltk.tokenize import sent_tokenize
tokenized_text= sent_tokenize(text)
print(tokenized_text)
```

```
['Tokenization is the first step in text analytics.', 'The process of breaking down a text paragraph into smaller chunks such as wo
```

```
from nltk.tokenize import word_tokenize
tokenized_word=word_tokenize(text)
print(tokenized_word)
```

```
['Tokenization', 'is', 'the', 'first', 'step', 'in', 'text', 'analytics', '.', 'The', 'process', 'of', 'breaking', 'down', 'a', 'te
```

STOP WORDS REMOVAL

```
from nltk.corpus import stopwords
stop_words=set(stopwords.words("english"))
print(stop_words)
```

```
{'all', 'those', 'weren', 'myself', 'by', 'before', 'further', "didn't", 'while', 'against', 'hers', 'now', "won't", "you'd", 'his'
```

```
import re
re.compile('<title>(.*)</title>')

re.compile(r'<title>(.*)</title>', re.UNICODE)
```

STEMMING

```
text= "How to remove stop words with NLTK library in Python?"
text= re.sub('[^a-zA-Z]', ' ',text)
tokens = word_tokenize(text.lower())
```

```

filtered_text=[]
for w in tokens:
    if w not in stop_words:
        filtered_text.append(w)
print("Tokenized Sentence:",tokens)
print("Filtered Sentence:",filtered_text)

Tokenized Sentence: ['how', 'to', 'remove', 'stop', 'words', 'with', 'nltk', 'library', 'in', 'python']
Filtered Sentence: ['remove', 'stop', 'words', 'nltk', 'library', 'python']

```

```

from nltk.stem import PorterStemmer
e_words= ["wait", "waiting", "waited", "waits"]
ps =PorterStemmer()
for w in e_words:
    rootWord=ps.stem(w)
print(rootWord)

```

wait

LEMMATIZATION

```

from nltk.stem import WordNetLemmatizer
wordnet_lemmatizer = WordNetLemmatizer()
text = "studies studying cries cry"
tokenization = nltk.word_tokenize(text)
for w in tokenization:
    print("Lemma for {} is {}".format(w,wordnet_lemmatizer.lemmatize(w)))

```

```

Lemma for studies is study
Lemma for studying is studying
Lemma for cries is cry
Lemma for cry is cry

```

```
import pandas as pd
```

```
from sklearn.feature_extraction.text import TfidfVectorizer
```

```

documentA = 'Jupiter is the largest Planet'
documentB = 'Mars is the fourth planet from the Sun'

```

```

bagOfWordsA = documentA.split(' ')
bagOfWordsB = documentB.split(' ')

```

```
uniqueWords = set(bagOfWordsA).union(set(bagOfWordsB))
```

```

numOfWordsA = dict.fromkeys(uniqueWords, 0)
for word in bagOfWordsA:
    numOfWordsA[word] += 1
numOfWordsB = dict.fromkeys(uniqueWords, 0)
for word in bagOfWordsB:
    numOfWordsB[word] += 1

```

```

def computeTF(wordDict, bagOfWords):
    tfDict = {}
    bagOfWordsCount = len(bagOfWords)
    for word, count in wordDict.items():
        tfDict[word] = count / float(bagOfWordsCount)
    return tfDict
tfA = computeTF(numOfWordsA, bagOfWordsA)
tfB = computeTF(numOfWordsB, bagOfWordsB)

```

```
def computeIDF(documents):
    import math
    N = len(documents)
    idfDict = dict.fromkeys(documents[0].keys(), 0)
    for document in documents:
        for word, val in document.items():
            if val > 0:
                idfDict[word] += 1
    for word, val in idfDict.items():
        idfDict[word] = math.log(N / float(val))
    return idfDict
idfs = computeIDF([numOfWordsA, numOfWordsB])
idfs
```

```
{'Mars': 0.6931471805599453,
 'Sun': 0.6931471805599453,
 'from': 0.6931471805599453,
 'fourth': 0.6931471805599453,
 'is': 0.0,
 'largest': 0.6931471805599453,
 'Jupiter': 0.6931471805599453,
 'the': 0.0,
 'planet': 0.6931471805599453,
 'Planet': 0.6931471805599453}
```

```
def computeTFIDF(tfBagOfWords, idfs):
    tfidf = {}
    for word, val in tfBagOfWords.items():
        tfidf[word] = val * idfs[word]
    return tfidf
tfidfA = computeTFIDF(tfA, idfs)
tfidfB = computeTFIDF(tfB, idfs)
df = pd.DataFrame([tfidfA, tfidfB])
df
```

	Mars	Sun	from	fourth	is	largest	Jupiter	the	planet	P1
0	0.000000	0.000000	0.000000	0.000000	0.0	0.138629	0.138629	0.0	0.000000	0.13
1	0.086643	0.086643	0.086643	0.086643	0.0	0.000000	0.000000	0.0	0.086643	0.00

Data Visualization I

1. Use the inbuilt dataset 'titanic'. The dataset contains 891 rows and contains information about the passengers who boarded the unfortunate Titanic ship. Use the Seaborn library to see if we can find any patterns in the data.
2. Write a code to check how the price of the ticket (column name: 'fare') for each passenger is distributed by plotting a histogram.

IMPORTING REQUIRED LIBRARIES

```
import pandas as pd  
  
import numpy as np  
  
import matplotlib.pyplot as plt  
  
import seaborn as sns
```

LOADING DATASET

```
dataset=sns.load_dataset('titanic')
```

DATA PREPROCESSING

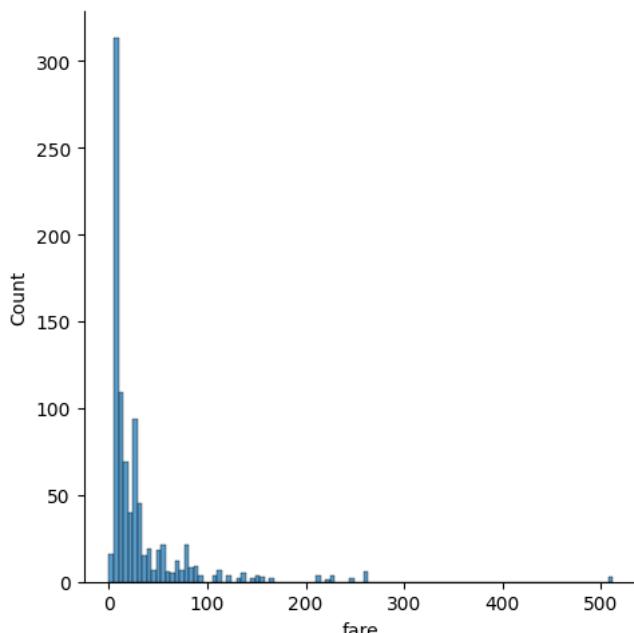
```
dataset.head()
```

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adul
0	0	3	male	22.0	1	0	7.2500	S	Third	man	
1	1	1	female	38.0	1	0	71.2833	C	First	woman	
2	1	3	female	26.0	0	0	7.9250	S	Third	woman	
3	1	1	female	35.0	1	0	53.1000	S	First	woman	
4	0	3	male	35.0	0	0	8.0500	S	Third	man	

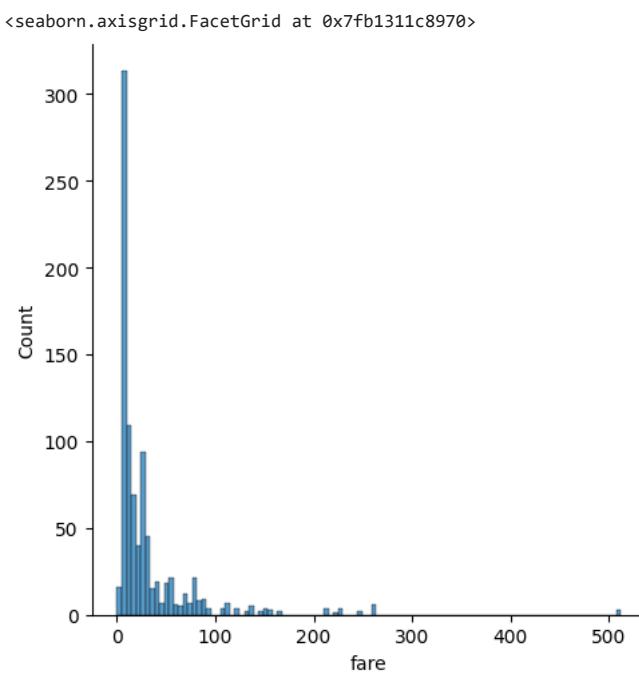
DATA VISUALIZATION

DISPLOT

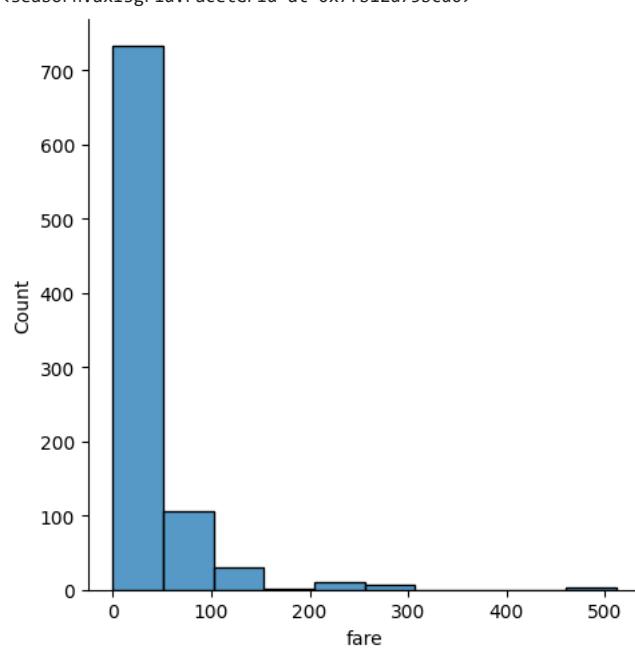
```
sns.displot(dataset['fare'])  
  
<seaborn.axisgrid.FacetGrid at 0x7fb13030a880>
```



```
sns.displot(dataset['fare'],kde=False)
```



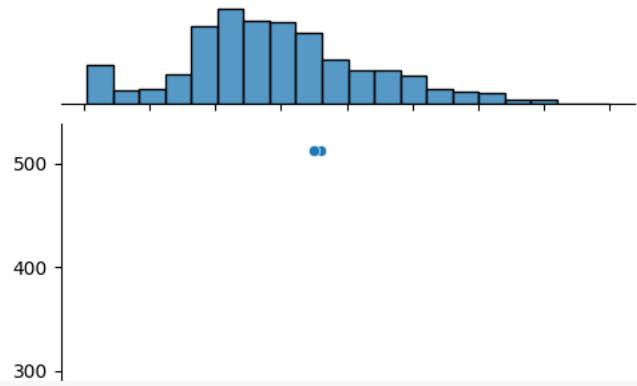
```
sns.displot(dataset['fare'],kde=False,bins=10)
```



JOINTPLOT

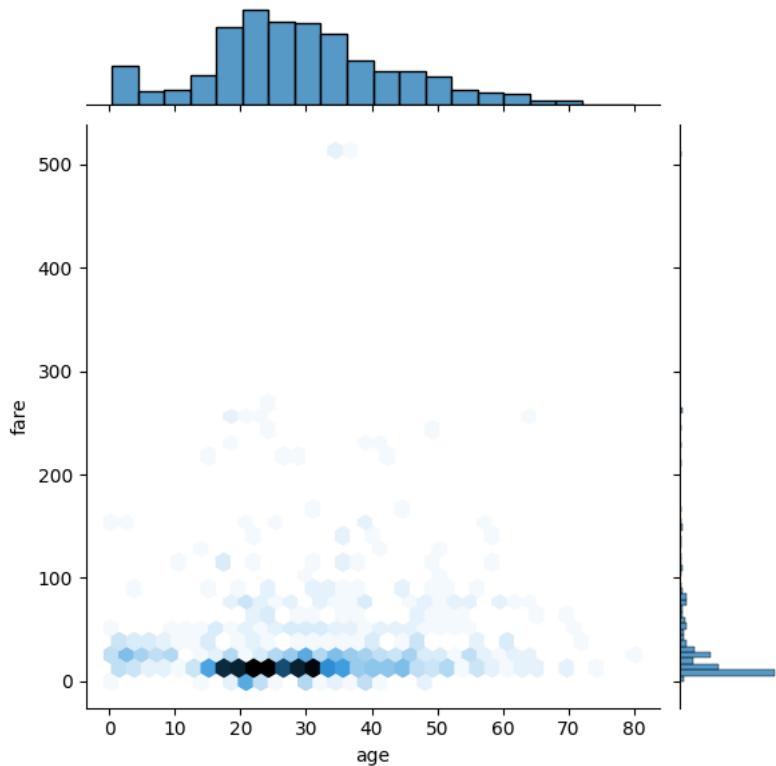
```
sns.jointplot(x='age',y='fare',data=dataset)
```

```
<seaborn.axisgrid.JointGrid at 0x7fb12d5f8e20>
```



```
sns.jointplot(x='age',y='fare',data=dataset,kind='hex')
```

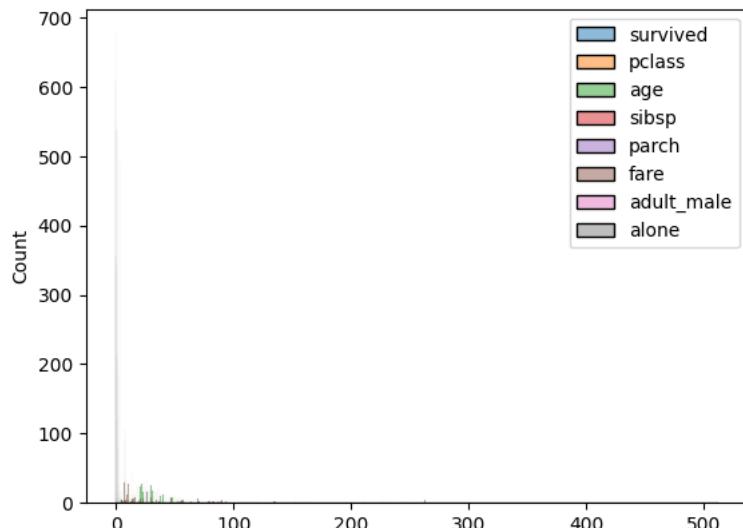
```
<seaborn.axisgrid.JointGrid at 0x7fb12d677ee0>
```



HISTPLOT

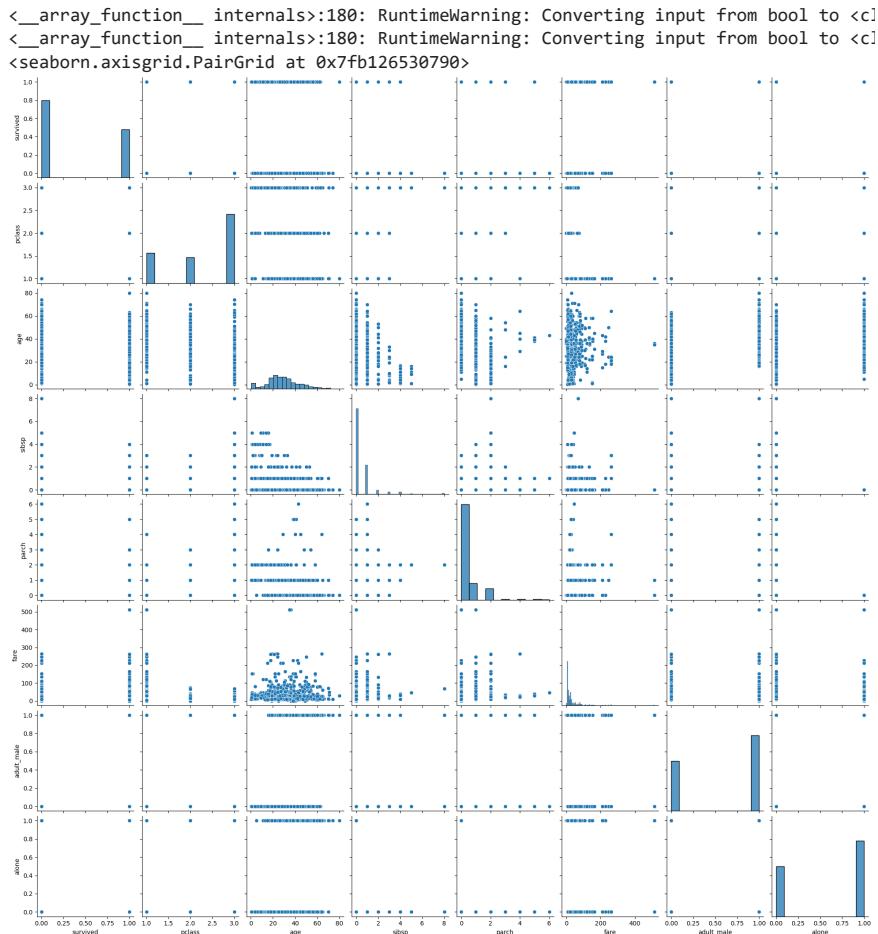
```
sns.histplot(dataset)
```

```
<Axes: ylabel='Count'>
```



PAIRPLOT

```
sns.pairplot(dataset)
```



Data Visualization II

1. Use the inbuilt dataset 'titanic' as used in the above problem. Plot a box plot for distribution of age with respect to each gender along with the information about whether they survived or not. (Column names : 'sex' and 'age')
2. Write observations on the inference from the above statistics.

IMPORTING REQUIRED LIBRARIES

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

from seaborn import load_dataset
```

```
sns.get_dataset_names()
```

```
['anagrams',
 'anscombe',
 'attention',
 'brain_networks',
 'car_crashes',
 'diamonds',
 'dots',
 'dowjones',
 'exercise',
 'flights',
 'fmri',
 'geyser',
 'glue',
 'healthexp',
 'iris',
 'mpg',
 'penguins',
 'planets',
 'seoice',
 'taxis',
 'tips',
 'titanic']
```

LOADING DATASET

```
data=load_dataset('titanic')
```

```
data
```

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	a
0	0	3	male	22.0	1	0	7.2500	S	Third	man	
1	1	1	female	38.0	1	0	71.2833	C	First	woman	
2	1	3	female	26.0	0	0	7.9250	S	Third	woman	
3	1	1	female	35.0	1	0	53.1000	S	First	woman	
4	0	3	male	35.0	0	0	8.0500	S	Third	man	
...
886	0	2	male	27.0	0	0	13.0000	S	Second	man	
887	1	1	female	19.0	0	0	30.0000	S	First	woman	
888	0	3	female	NaN	1	2	23.4500	S	Third	woman	
889	1	1	male	26.0	0	0	30.0000	C	First	man	
890	0	3	male	32.0	0	0	7.7500	Q	Third	man	

891 rows × 15 columns

DATA PREPROCESSING

```
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
```

```
Data columns (total 15 columns):
 #   Column      Non-Null Count Dtype  
 --- 
 0   survived    891 non-null   int64  
 1   pclass      891 non-null   int64  
 2   sex         891 non-null   object  
 3   age         714 non-null   float64 
 4   sibsp       891 non-null   int64  
 5   parch       891 non-null   int64  
 6   fare        891 non-null   float64 
 7   embarked    889 non-null   object  
 8   class       891 non-null   category 
 9   who         891 non-null   object  
 10  adult_male  891 non-null   bool    
 11  deck        203 non-null   category 
 12  embark_town 889 non-null   object  
 13  alive        891 non-null   object  
 14  alone        891 non-null   bool    
dtypes: bool(2), category(2), float64(2), int64(4), object(5)
memory usage: 80.7+ KB
```

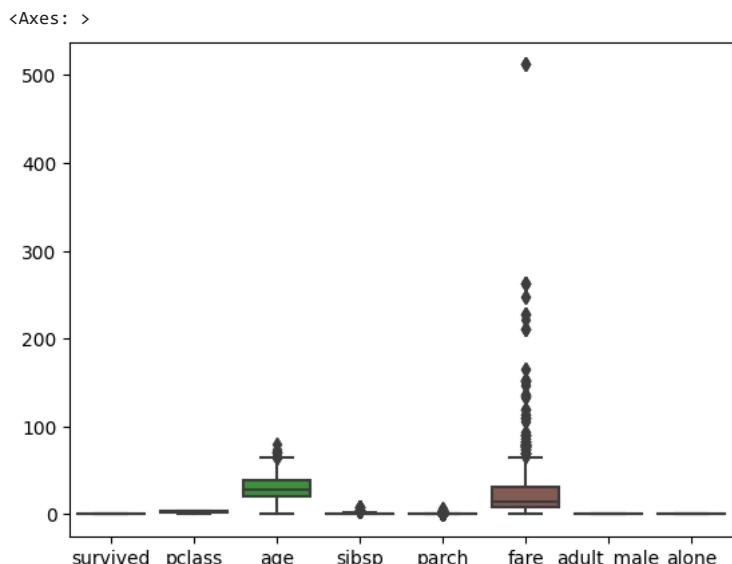
```
data.describe()
```

	survived	pclass	age	sibsp	parch	fare
count	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000
mean	0.383838	2.308642	29.699118	0.523008	0.381594	32.204208
std	0.486592	0.836071	14.526497	1.102743	0.806057	49.693429
min	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
25%	0.000000	2.000000	20.125000	0.000000	0.000000	7.910400
50%	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
75%	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000
max	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

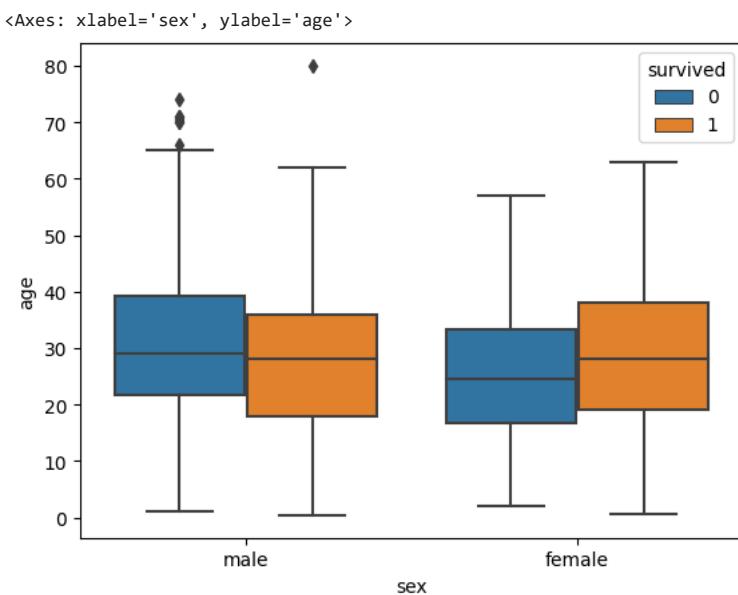
DATA VISUALIZATION

BOXPLOT

```
sns.boxplot(data)
```



```
sns.boxplot(x=data['sex'],y=data['age'],hue=data['survived'])
```



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



Data Visualization III Download the Iris flower dataset or any other dataset into a DataFrame. (e.g., <https://archive.ics.uci.edu/ml/datasets/Iris>). Scan the dataset and give the inference as:

1. List down the features and their types (e.g., numeric, nominal) available in the dataset.
2. Create a histogram for each feature in the dataset to illustrate the feature distributions.
3. Create a boxplot for each feature in the dataset.
4. Compare distributions and identify outliers

IMPORTING REQUIRED LIBRARIES

```
import seaborn as sns

sns.get_dataset_names()

['anagrams',
 'anscombe',
 'attention',
 'brain_networks',
 'car_crashes',
 'diamonds',
 'dots',
 'dowjones',
 'exercise',
 'flights',
 'fmri',
 'geyser',
 'glue',
 'healthexp',
 'iris',
 'mpg',
 'penguins',
 'planets',
 'seice',
 'taxis',
 'tips',
 'titanic']
```

LOADING DATASET

```
dataset=sns.load_dataset('iris')
```

DATA PREPROCESSING

```
dataset.head()
```

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa

```
dataset.describe()
```

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa
..
145	6.7	3.0	5.2	2.3	virginica
146	6.3	2.5	5.0	1.9	virginica
147	6.5	3.0	5.2	2.0	virginica
148	6.2	3.4	5.4	2.3	virginica
149	5.9	3.0	5.1	1.8	virginica

```
[150 rows x 5 columns]>
```

```
dataset.info
```

```
<bound method DataFrame.info of
0      5.1      3.5      1.4      0.2    setosa
1      4.9      3.0      1.4      0.2    setosa
2      4.7      3.2      1.3      0.2    setosa
3      4.6      3.1      1.5      0.2    setosa
4      5.0      3.6      1.4      0.2    setosa
..     ...
145     6.7      3.0      5.2      2.3  virginica
146     6.3      2.5      5.0      1.9  virginica
147     6.5      3.0      5.2      2.0  virginica
148     6.2      3.4      5.4      2.3  virginica
149     5.9      3.0      5.1      1.8  virginica
```

[150 rows x 5 columns]>

CHECKING FOR NULL VALUES

```
dataset.isnull().any()
```

```
sepal_length    False
sepal_width     False
petal_length    False
petal_width     False
species         False
dtype: bool
```

```
dataset.isnull().sum()
```

```
sepal_length    0
sepal_width     0
petal_length    0
petal_width     0
species         0
dtype: int64
```

MAX AND MIN

```
dataset.max()
```

```
sepal_length      7.9
sepal_width       4.4
petal_length      6.9
petal_width       2.5
species          virginica
dtype: object
```

```
dataset.min()
```

```
sepal_length      4.3
sepal_width       2.0
petal_length      1.0
petal_width       0.1
species          setosa
dtype: object
```

MEAN

```
dataset.mean()
```

```
<ipython-input-16-e55bc0ed4499>:1: FutureWarning: The default value of numeric_only in DataFrame.mean is deprecated. In a future version of pandas, this will return a Series.
dataset.mean()
sepal_length      5.843333
sepal_width       3.057333
petal_length      3.758000
petal_width       1.199333
dtype: float64
```

STANDARD DEVIATION

```
a=dataset.std()
a
```

```
<ipython-input-19-e1e8b4cdc687>:1: FutureWarning: The default value of numeric_only in DataFrame.std is deprecated. In a future version of pandas, this will return a Series.
a=dataset.std()
sepal_length      0.828066
sepal_width       0.435866
petal_length      1.765298
```

```
petal_width      0.762238
dtype: float64
```

VARIANCE

```
variance = a*a
variance

sepal_length    0.685694
sepal_width     0.189979
petal_length    3.116278
petal_width     0.581006
dtype: float64
```

```
import numpy as np
```

PERCENTILE

```
np.percentile(dataset['sepal_length'],90)
```

```
6.9
```

```
np.percentile(dataset['sepal_width'],90)
```

```
3.6099999999999994
```

```
np.percentile(dataset['petal_length'],90)
```

```
5.8
```

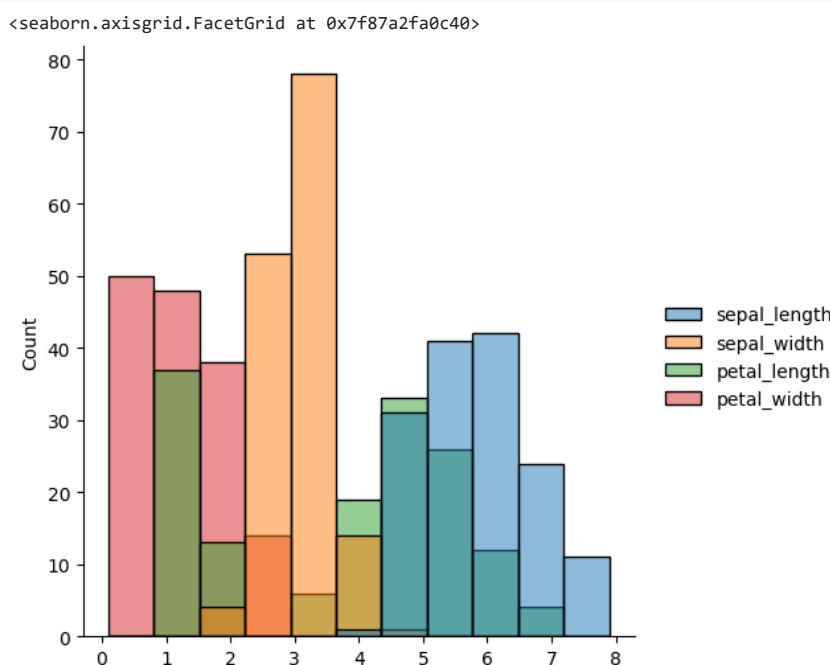
```
np.percentile(dataset['petal_width'],90)
```

```
2.2
```

DATA VISUALIZATION

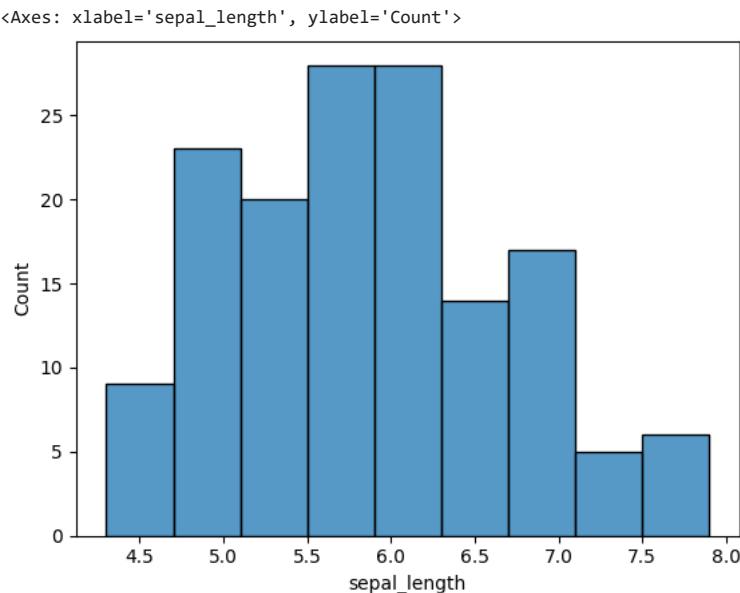
DISPLOT

```
sns.displot(dataset)
```

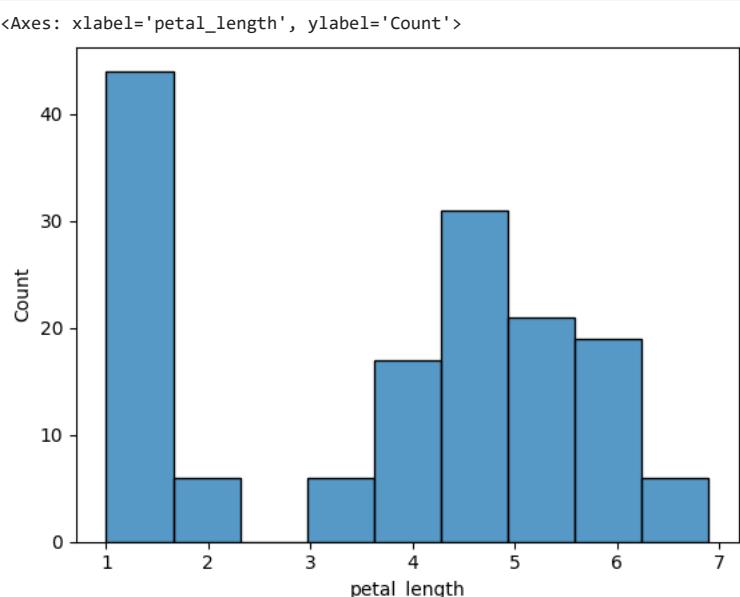


HISTPLOT

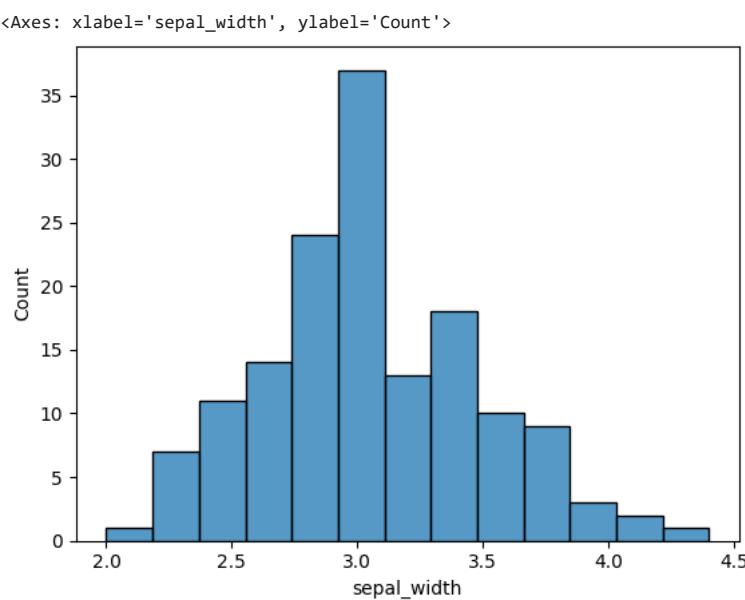
```
sns.histplot(dataset['sepal_length'])
```



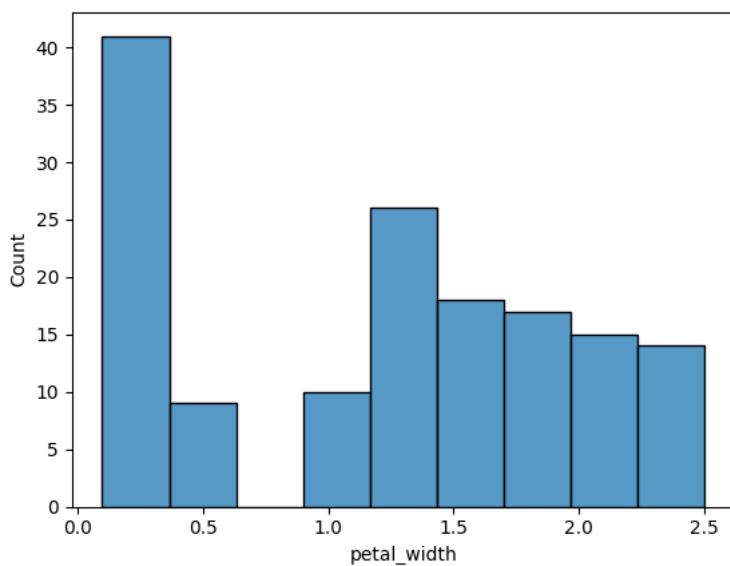
```
sns.histplot(dataset['petal_length'])
```



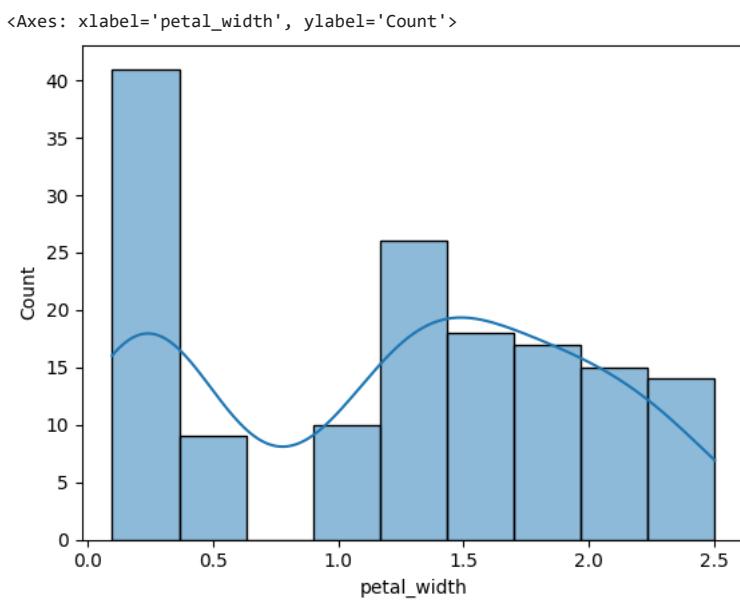
```
sns.histplot(dataset['sepal_width'])
```



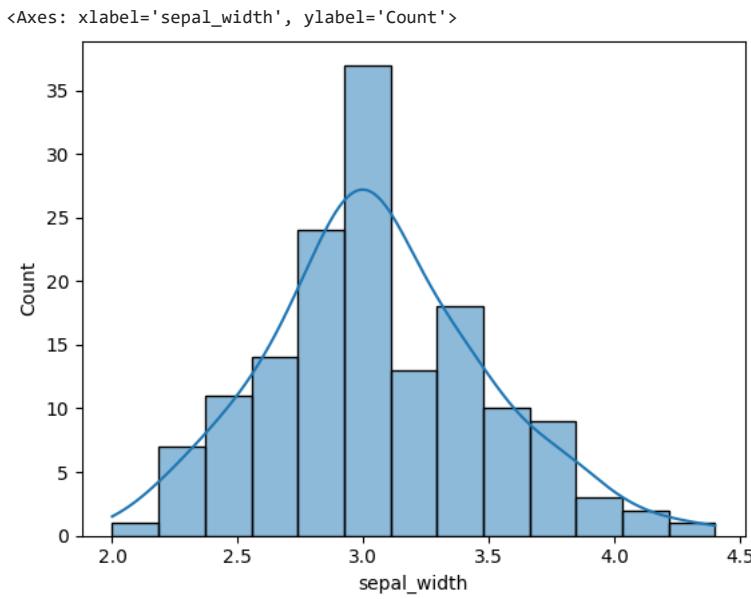
```
sns.histplot(dataset['petal_width'])  
  
<Axes: xlabel='petal_width', ylabel='Count'>
```



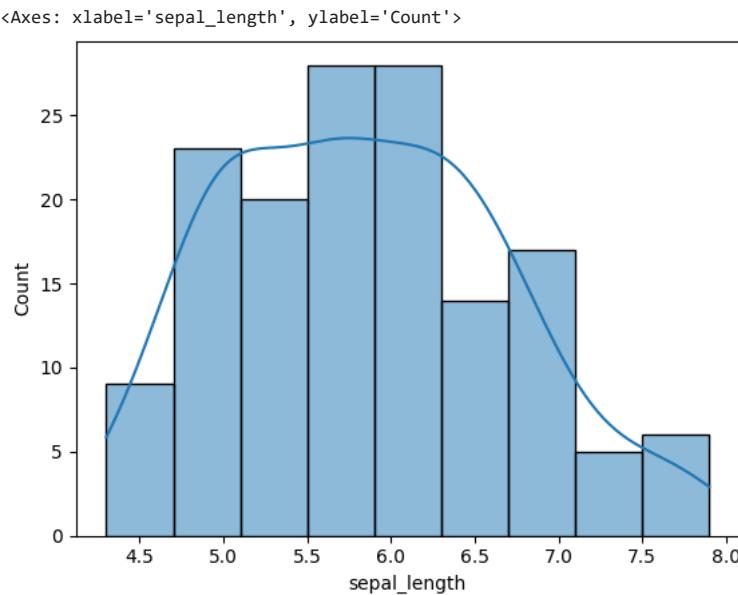
```
sns.histplot(dataset['petal_width'], kde=True)
```



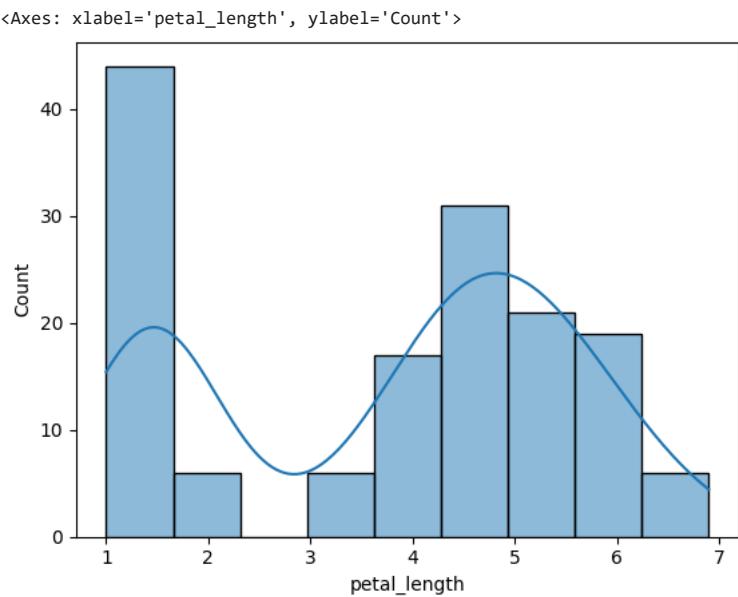
```
sns.histplot(dataset['sepal_width'], kde=True)
```



```
sns.histplot(dataset['sepal_length'], kde=True)
```

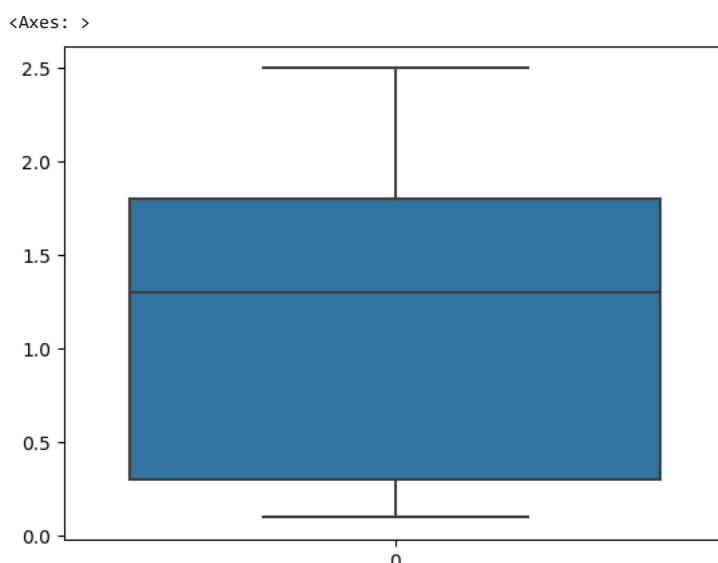
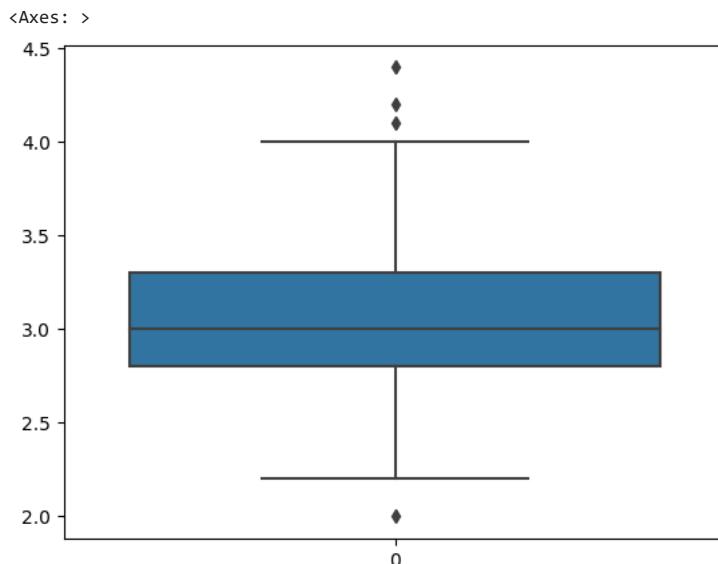
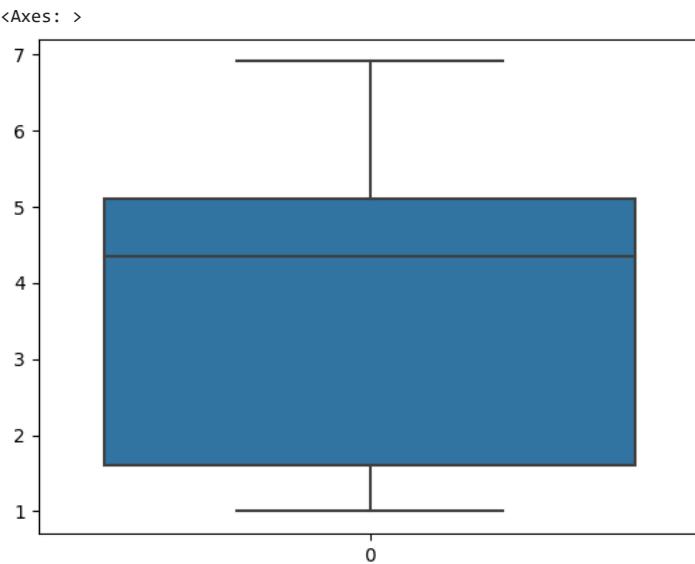
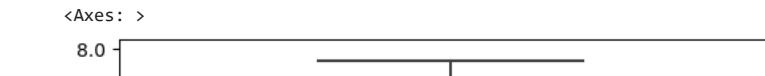


```
sns.histplot(dataset['petal_length'], kde=True)
```

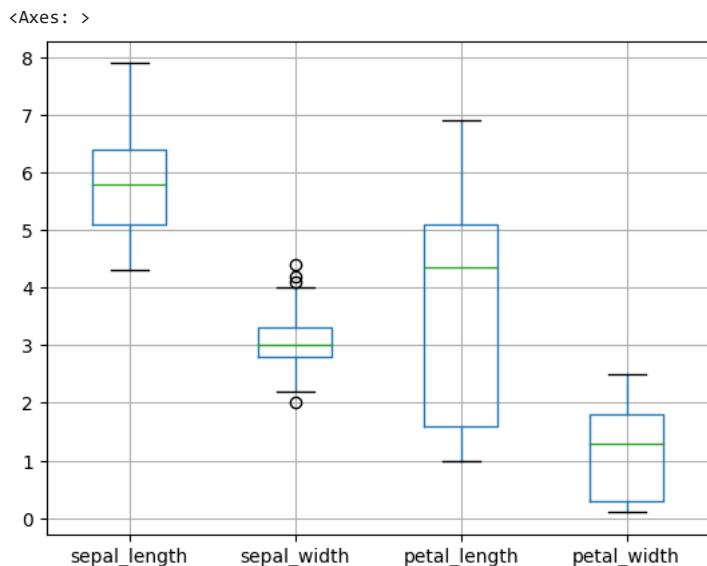


BOXPLOT

```
sns.boxplot(dataset['sepal_length'])
```



```
dataset.boxplot()
```



**Data Wrangling, I

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

1. Import all the required Python Libraries.
2. Locate an open source data from the web (e.g., <https://www.kaggle.com>). Provide a clear description of the data and its source (i.e., URL of the web site).
3. Load the Dataset into pandas dataframe.
4. Data Preprocessing: check for missing values in the data using pandas isnull(), describe() function to get some initial statistics. Provide variable descriptions. Types of variables etc. Check the dimensions of the data frame.
5. Data Formatting and Data Normalization: Summarize the types of variables by checking the data types (i.e., character, numeric, integer, factor, and logical) of the variables in the data set. If variables are not in the correct data type, apply proper type conversions.
6. Turn categorical variables into quantitative variables in Python

Import all the required Python Libraries

```
import pandas as pd
```

```
!pip install -q kaggle
```

```
from google.colab import files
files.upload()
```

Choose files | 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 Iris.csv to Iris.csv
{'Iris.csv':
b'Id,SepalLengthCm,SepalWidthCm,PetalLengthCm,PetalWidthCm,Species\n1,5.1,3.5,1.4,
setosa\n2,4.9,3.0,1.4,0.2,Iris-setosa\n3,4.7,3.2,1.3,0.2,Iris-
setosa\n4,4.6,3.1,1.5,0.2,Iris-setosa\n5,5.0,3.6,1.4,0.2,Iris-
setosa\n6,5.4,3.9,1.7,0.4,Iris-setosa\n7,4.6,3.4,1.4,0.3,Iris-
setosa\n8,5.0,3.4,1.5,0.2,Iris-setosa\n9,4.4,2.9,1.4,0.2,Iris-
setosa\n10,4.9,3.1,1.5,0.1,Iris-setosa\n11,5.4,3.7,1.5,0.2,Iris-
setosa\n12,4.8,3.4,1.6,0.2,Iris-setosa\n13,4.8,3.0,1.4,0.1,Iris-
setosa\n14,4.3,3.0,1.1,0.1,Iris-setosa\n15,5.8,4.0,1.2,0.2,Iris-
setosa\n16,5.7,4.4,1.5,0.4,Iris-setosa\n17,5.4,3.9,1.3,0.4,Iris-
setosa\n18,5.1,3.5,1.4,0.3,Iris-setosa\n19,5.7,3.8,1.7,0.3,Iris-
setosa\n20,5.1,3.8,1.5,0.3,Iris-setosa\n21,5.4,3.4,1.7,0.2,Iris-
setosa\n22,5.1,3.7,1.5,0.4,Iris-setosa\n23,4.6,3.6,1.0,0.2,Iris-
setosa\n24,5.1,3.3,1.7,0.5,Iris-setosa\n25,4.8,3.4,1.9,0.2,Iris-
setosa\n26,5.0,3.0,1.6,0.2,Iris-setosa\n27,5.0,3.4,1.6,0.4,Iris-
setosa\n28,5.2,3.5,1.5,0.2,Iris-setosa\n29,5.2,3.4,1.4,0.2,Iris-
setosa\n30,4.7,3.2,1.6,0.2,Iris-setosa\n31,4.8,3.1,1.6,0.2,Iris-
setosa\n32,5.4,3.4,1.5,0.4,Iris-setosa\n33,5.2,4.1,1.5,0.1,Iris-
setosa\n34,5.5,4.2,1.4,0.2,Iris-setosa\n35,4.9,3.1,1.5,0.1,Iris-
setosa\n36,5.0,3.2,1.2,0.2,Iris-setosa\n37,5.5,3.5,1.3,0.2,Iris-
setosa\n38,4.9,3.1,1.5,0.1,Iris-setosa\n39,4.4,3.0,1.3,0.2,Iris-
setosa\n40,5.1,3.4,1.5,0.2,Iris-setosa\n41,5.0,3.5,1.3,0.3,Iris-
setosa\n42,4.5,2.3,1.3,0.3,Iris-setosa\n43,4.4,3.2,1.3,0.2,Iris-
setosa\n44,5.0,3.5,1.6,0.6,Iris-setosa\n45,5.1,3.8,1.9,0.4,Iris-
setosa\n46,4.8,3.0,1.4,0.3,Iris-setosa\n47,5.1,3.8,1.6,0.2,Iris-
setosa\n48,4.6,3.2,1.4,0.2,Iris-setosa\n49,5.3,3.7,1.5,0.2,Iris-
setosa\n50,5.0,3.3,1.4,0.2,Iris-setosa\n51,7.0,3.2,4.7,1.4,Iris-
versicolor\n52,6.4,3.2,4.5,1.5,Iris-versicolor\n53,6.9,3.1,4.9,1.5,Iris-
versicolor\n54,5.5,2.3,4.0,1.3,Iris-versicolor\n55,6.5,2.8,4.6,1.5,Iris-
versicolor\n56,5.7,2.8,4.5,1.3,Iris-versicolor\n57,6.3,3.3,4.7,1.6,Iris-
versicolor\n58,4.9,2.4,3.3,1.0,Iris-versicolor\n59,6.6,2.9,4.6,1.3,Iris-
versicolor\n60,5.2,2.7,3.9,1.4,Iris-versicolor\n61,5.0,2.0,3.5,1.0,Iris-
versicolor\n62,5.9,3.0,4.2,1.5,Iris-versicolor\n63,6.0,2.2,4.0,1.0,Iris-
versicolor\n64,6.1,2.9,4.7,1.4,Iris-versicolor\n65,5.6,2.9,3.6,1.3,Iris-
versicolor\n66,6.7,3.1,1.6,1.1,Iris-versicolor\n67,5.6,3.0,1.5,1.5,Iris-

Load the Dataset into pandas dataframe.

Indented block

```
g,mial
```

Data Preprocessing

```
iris.head()
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa
3	4	4.6	3.1	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa

```
iris.tail()
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
145	146	6.7	3.0	5.2	2.3	Iris-virginica
146	147	6.3	2.5	5.0	1.9	Iris-virginica
147	148	6.5	3.0	5.2	2.0	Iris-virginica
148	149	6.2	3.4	5.4	2.3	Iris-virginica
149	150	5.9	3.0	5.1	1.8	Iris-virginica

```
iris.describe(include="all")
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Spec:
count	150.000000	150.000000	150.000000	150.000000	150.000000	'
unique	NaN	NaN	NaN	NaN	NaN	
top	NaN	NaN	NaN	NaN	NaN	I seto
freq	NaN	NaN	NaN	NaN	NaN	
mean	75.500000	5.843333	3.054000	3.758667	1.198667	N
std	43.445368	0.828066	0.433594	1.764420	0.763161	N
min	1.000000	4.300000	2.000000	1.000000	0.100000	N
25%	38.250000	5.100000	2.800000	1.600000	0.300000	N
50%	75.500000	5.800000	3.000000	4.350000	1.300000	N
75%	112.750000	6.400000	3.300000	5.100000	1.800000	N

```
iris.shape
```

```
(150, 6)
```

```
iris.columns
```

```
Index(['Id', 'SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', 'PetalWidthCm', 'Species'],
      dtype='object')
```

```
iris['Id']
```

```
0      1
1      2
2      3
3      4
4      5
...
145   146
146   147
147   148
148   149
149   150
Name: Id, Length: 150, dtype: int64
```

```
iris[0:3]
```

```
iris.loc[0:2]
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa

```
iris.loc[0:2,'Id':'PetalWidthCm']
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
0	1	5.1	3.5	1.4	0.2
1	2	4.9	3.0	1.4	0.2
2	3	4.7	3.2	1.3	0.2

```
iris.iloc[1:3]
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa

```
iris.iloc[1:5,1:5]
```

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2

Check for missing values in the data using pandas isnull()

```
iris.isnull()
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	False	False	False	False	False	False
1	False	False	False	False	False	False
2	False	False	False	False	False	False
3	False	False	False	False	False	False
4	False	False	False	False	False	False
...
145	False	False	False	False	False	False
146	False	False	False	False	False	False
147	False	False	False	False	False	False
148	False	False	False	False	False	False
149	False	False	False	False	False	False

150 rows × 6 columns

```
iris.isna()
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	False	False	False	False	False	False
1	False	False	False	False	False	False
2	False	False	False	False	False	False
3	False	False	False	False	False	False
4	False	False	False	False	False	False
...

```
iris.isnull().any()
```

```
Id      False
SepalLengthCm  False
SepalWidthCm  False
PetalLengthCm  False
PetalWidthCm  False
Species     False
dtype: bool
```

```
iris.head
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa
3	4	4.6	3.1	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa
..
145	146	6.7	3.0	5.2	2.3	Iris-virginica
146	147	6.3	2.5	5.0	1.9	Iris-virginica
147	148	6.5	3.0	5.2	2.0	Iris-virginica
148	149	6.2	3.4	5.4	2.3	Iris-virginica
149	150	5.9	3.0	5.1	1.8	Iris-virginica

```
[150 rows x 6 columns]>
```

```
iris.isnull().sum()
```

```
Id      0
SepalLengthCm  0
SepalWidthCm  0
PetalLengthCm  0
PetalWidthCm  0
Species     0
dtype: int64
```

```
iris.head
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa
3	4	4.6	3.1	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa
..
145	146	6.7	3.0	5.2	2.3	Iris-virginica
146	147	6.3	2.5	5.0	1.9	Iris-virginica
147	148	6.5	3.0	5.2	2.0	Iris-virginica
148	149	6.2	3.4	5.4	2.3	Iris-virginica
149	150	5.9	3.0	5.1	1.8	Iris-virginica

```
146 Iris-virginica  
147 Iris-virginica  
148 Iris-virginica  
149 Iris-virginica
```

[150 rows x 6 columns]>

Count of missing values of a specific column

```
iris.SepalLengthCm.isnull().sum()
```

8

iris.tail

<bound	method	NDFrame.tail of	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	\
0	1	5.1	3.5	1.4	0.2			
1	2	4.9	3.0	1.4	0.2			
2	3	4.7	3.2	1.3	0.2			
3	4	4.6	3.1	1.5	0.2			
4	5	5.0	3.6	1.4	0.2			
..			
145	146	6.7	3.0	5.2	2.3			
146	147	6.3	2.5	5.0	1.9			
147	148	6.5	3.0	5.2	2.0			
148	149	6.2	3.4	5.4	2.3			
149	150	5.9	3.0	5.1	1.8			

```
      Species
0    Iris-setosa
1    Iris-setosa
2    Iris-setosa
3    Iris-setosa
4    Iris-setosa
...
145   Iris-virginica
146   Iris-virginica
147   Iris-virginica
148   Iris-virginica
149   Iris-virginica
```

[150 rows x 6 columns]>

```
iris.head
```

<bound	method	NDFrame.head of						
0	1	5.1	3.5	1.4	0.2			
1	2	4.9	3.0	1.4	0.2			
2	3	4.7	3.2	1.3	0.2			
3	4	4.6	3.1	1.5	0.2			
4	5	5.0	3.6	1.4	0.2			
..			
145	146	6.7	3.0	5.2	2.3			
146	147	6.3	2.5	5.0	1.9			
147	148	6.5	3.0	5.2	2.0			
148	149	6.2	3.4	5.4	2.3			
149	150	5.9	3.0	5.1	1.8			

```

          Species
0    Iris-setosa
1    Iris-setosa
2    Iris-setosa
3    Iris-setosa
4    Iris-setosa
...
145   Iris-virginica
146   Iris-virginica
147   Iris-virginica
148   Iris-virginica
149   Iris-virginica

```

[150 rows x 6 columns]>

Data Formatting

```
iris.dtypes
```

```
Id          int64
SepalLengthCm  float64
SepalWidthCm   float64
PetalLengthCm  float64
```

```
PetalWidthCm      float64
Species          object
dtype: object
```

```
iris.SepalLengthCm=iris.SepalLengthCm.astype("int")
```

```
iris.dtypes
```

```
Id            int64
SepalLengthCm    int64
SepalWidthCm     float64
PetalLengthCm    float64
PetalWidthCm     float64
Species          object
dtype: object
```

Data Normalization

```
from sklearn import preprocessing
```

```
iris.head()
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5	3.5	1.4	0.2	Iris-setosa
1	2	4	3.0	1.4	0.2	Iris-setosa
2	3	4	3.2	1.3	0.2	Iris-setosa
3	4	4	3.1	1.5	0.2	Iris-setosa
4	5	5	3.6	1.4	0.2	Iris-setosa

```
min_max_scaler = preprocessing.MinMaxScaler() #min-max scalar
```

```
x=iris.iloc[:, :4] #separate input from dataset
```

```
iris.head
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5	3.5	1.4	0.2	Iris-setosa
1	2	4	3.0	1.4	0.2	Iris-setosa
2	3	4	3.2	1.3	0.2	Iris-setosa
3	4	4	3.1	1.5	0.2	Iris-setosa
4	5	5	3.6	1.4	0.2	Iris-setosa
..
145	146	6	3.0	5.2	2.3	Iris-virginica
146	147	6	2.5	5.0	1.9	Iris-virginica
147	148	6	3.0	5.2	2.0	Iris-virginica
148	149	6	3.4	5.4	2.3	Iris-virginica
149	150	5	3.0	5.1	1.8	Iris-virginica

```
Species
0    Iris-setosa
1    Iris-setosa
2    Iris-setosa
3    Iris-setosa
4    Iris-setosa
..
145   Iris-virginica
146   Iris-virginica
147   Iris-virginica
148   Iris-virginica
149   Iris-virginica
```

```
[150 rows x 6 columns]>
```

```
x_scaled = min_max_scaler.fit_transform(x)
```

```
iris.head
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
0	1	5	3.5	1.4	0.2
1	2	4	3.0	1.4	0.2
2	3	4	3.2	1.3	0.2
3	4	4	3.1	1.5	0.2
4	5	5	3.6	1.4	0.2
..

```

145 146      6      3.0      5.2      2.3
146 147      6      2.5      5.0      1.9
147 148      6      3.0      5.2      2.0
148 149      6      3.4      5.4      2.3
149 150      5      3.0      5.1      1.8

```

```

Species
0    Iris-setosa
1    Iris-setosa
2    Iris-setosa
3    Iris-setosa
4    Iris-setosa
...
145 Iris-virginica
146 Iris-virginica
147 Iris-virginica
148 Iris-virginica
149 Iris-virginica

```

[150 rows x 6 columns]>

```
df_normalized = pd.DataFrame(x_scaled)
```

```
iris.tail
```

```

<bound method NDFrame.tail of      Id  SepalLengthCm  SepalWidthCm  PetalLengthCm  PetalWidthCm  \
0      1      5      3.5      1.4      0.2
1      2      4      3.0      1.4      0.2
2      3      4      3.2      1.3      0.2
3      4      4      3.1      1.5      0.2
4      5      5      3.6      1.4      0.2
... ...
145 146      6      3.0      5.2      2.3
146 147      6      2.5      5.0      1.9
147 148      6      3.0      5.2      2.0
148 149      6      3.4      5.4      2.3
149 150      5      3.0      5.1      1.8

Species
0    Iris-setosa
1    Iris-setosa
2    Iris-setosa
3    Iris-setosa
4    Iris-setosa
...
145 Iris-virginica
146 Iris-virginica
147 Iris-virginica
148 Iris-virginica
149 Iris-virginica

```

[150 rows x 6 columns]>

```
iris.Id
```

```

0      1
1      2
2      3
3      4
4      5
...
145 146
146 147
147 148
148 149
149 150
Name: Id, Length: 150, dtype: int64

```

```
df_normalized
```

```
    0      1      2      3
0  0.000000  0.333333  0.625000  0.067797
1  0.006711  0.000000  0.416667  0.067797
2  0.013423  0.000000  0.500000  0.050847
3  0.020134  0.000000  0.458333  0.084746
```

Label Encoding

```
...     ...     ...     ...     ...
from sklearn import preprocessing

iris['Species'].unique()
array(['Iris-setosa', 'Iris-versicolor', 'Iris-virginica'], dtype=object)

label_encoder = preprocessing.LabelEncoder()
150 rows × 4 columns
iris['Species']= label_encoder.fit_transform(iris['Species'])

iris['Species'].unique()
array([0, 1, 2])
```

```
iris.Id
0      1
1      2
2      3
3      4
4      5
...
145    146
146    147
147    148
148    149
149    150
Name: Id, Length: 150, dtype: int64
```

```
iris.Species
0      0
1      0
2      0
3      0
4      0
...
145    2
146    2
147    2
148    2
149    2
Name: Species, Length: 150, dtype: int64
```

```
iris.PetalLengthCm
0      1.4
1      1.4
2      1.3
3      1.5
4      1.4
...
145    5.2
146    5.0
147    5.2
148    5.4
149    5.1
Name: PetalLengthCm, Length: 150, dtype: float64
```

