

kaviyadevi 20106064

```
In [1]: #to import Libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [49]: #to import dataset
data=pd.read_csv(r"C:\Users\user\Downloads\4_drug200 - 4_drug200.csv")
data
```

Out[49]:

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	F	HIGH	HIGH	25.355	drugY
1	47	M	LOW	HIGH	13.093	drugC
2	47	M	LOW	HIGH	10.114	drugC
3	28	F	NORMAL	HIGH	7.798	drugX
4	61	F	LOW	HIGH	18.043	drugY
...	...	...	...	...	...	...
195	56	F	LOW	HIGH	11.567	drugC
196	16	M	LOW	HIGH	12.006	drugC
197	52	M	NORMAL	HIGH	9.894	drugX
198	23	M	NORMAL	NORMAL	14.020	drugX
199	40	F	LOW	NORMAL	11.349	drugX

200 rows × 6 columns

```
In [50]: #to display top 5 rows
data.head()
```

Out[50]:

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	F	HIGH	HIGH	25.355	drugY
1	47	M	LOW	HIGH	13.093	drugC
2	47	M	LOW	HIGH	10.114	drugC
3	28	F	NORMAL	HIGH	7.798	drugX
4	61	F	LOW	HIGH	18.043	drugY

## DATA CLEANING AND PREPROCESSING

In [51]: data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 6 columns):
 #   Column          Non-Null Count  Dtype  
---  -
 0   Age             200 non-null   int64  
 1   Sex             200 non-null   object  
 2   BP              200 non-null   object  
 3   Cholesterol     200 non-null   object  
 4   Na_to_K         200 non-null   float64  
 5   Drug            200 non-null   object  
dtypes: float64(1), int64(1), object(4)
memory usage: 9.5+ KB
```

In [52]: *#to display summary of statistics*  
data.describe()

Out[52]:

	Age	Na_to_K
<b>count</b>	200.000000	200.000000
<b>mean</b>	44.315000	16.084485
<b>std</b>	16.544315	7.223956
<b>min</b>	15.000000	6.269000
<b>25%</b>	31.000000	10.445500
<b>50%</b>	45.000000	13.936500
<b>75%</b>	58.000000	19.380000
<b>max</b>	74.000000	38.247000

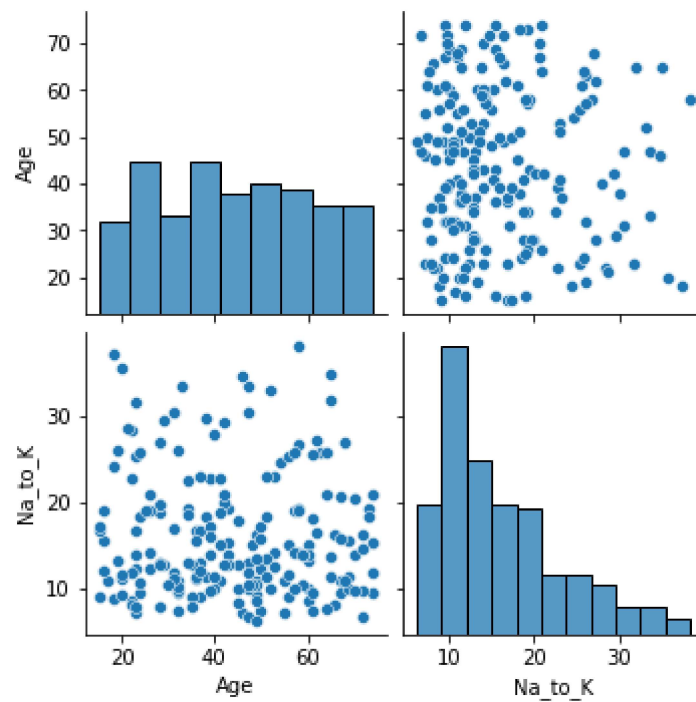
In [53]: *#to display the column heading*  
data.columns

Out[53]: Index(['Age', 'Sex', 'BP', 'Cholesterol', 'Na\_to\_K', 'Drug'], dtype='object')

## EDA and DATA VISUALIZATION

```
In [54]: sns.pairplot(data)
```

```
Out[54]: <seaborn.axisgrid.PairGrid at 0x1e5af7fc0a0>
```

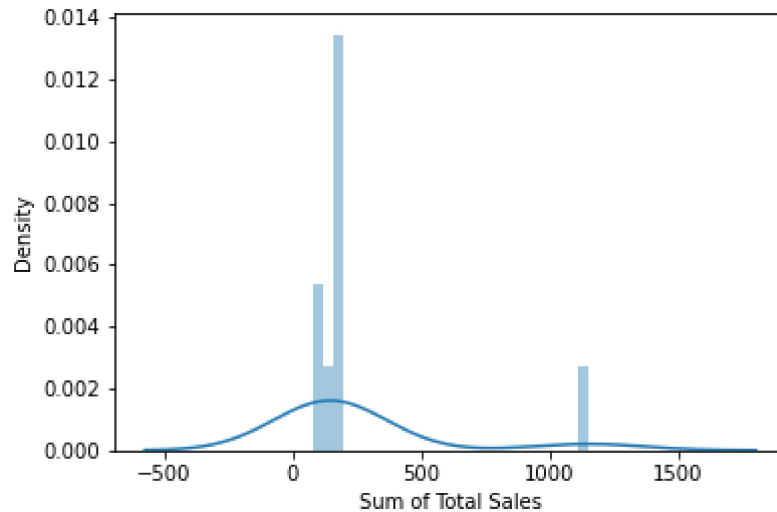


```
In [9]: sns.distplot(data["Sum of Total Sales"])
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

```
warnings.warn(msg, FutureWarning)
```

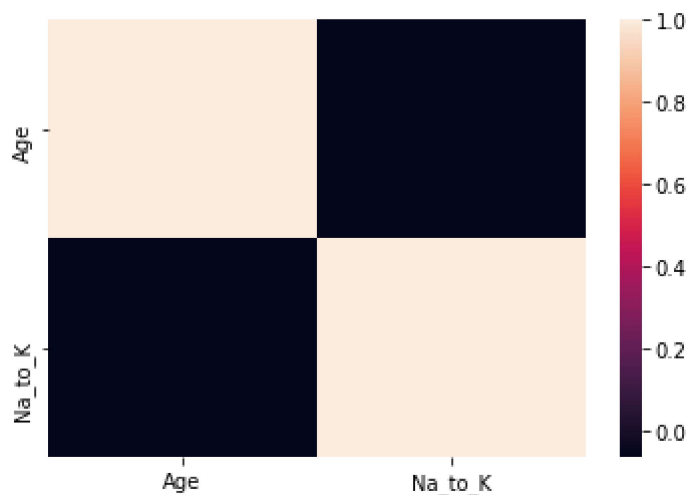
```
Out[9]: <AxesSubplot:xlabel='Sum of Total Sales', ylabel='Density'>
```



```
In [56]: df=data[['Age', 'Sex', 'BP', 'Cholesterol', 'Na_to_K', 'Drug']]
```

```
In [57]: sns.heatmap(df.corr())
```

```
Out[57]: <AxesSubplot:>
```



## MODEL TRAINING

```
In [74]: x=df[['Age']]
         y=df[['Na_to_K']]
```

```
In [75]: #to split my dataset into training and test

from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

```
In [76]: from sklearn.linear_model import LinearRegression

lr=LinearRegression()
lr.fit(x_train,y_train)
```

```
Out[76]: LinearRegression()
```

```
In [77]: #to find intercept
         print(lr.intercept_)
```

```
[18.04819595]
```

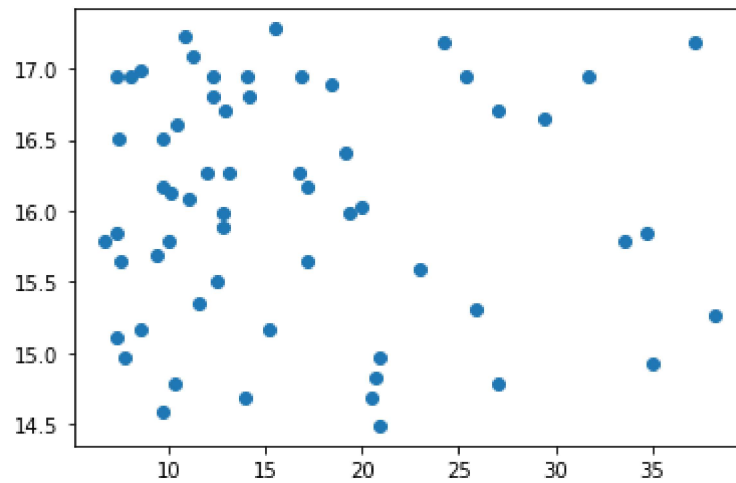
```
In [78]: coeff = pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
         coeff
```

```
Out[78]:
```

	Co-efficient
Age	-0.048056

```
In [79]: prediction = lr.predict(x_test)
plt.scatter(y_test, prediction)
```

Out[79]: <matplotlib.collections.PathCollection at 0x1e5b2a50910>



```
In [80]: print(lr.score(x_test, y_test))
```

-0.022059454830758662