# HEAMNATH

# 20104028

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

In [2]:
```python
df=pd.read_csv("4_drug200.csv")
df
```

Out[2]:

|     | 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 [3]:
```python
df.head()
```

Out[3]:

|     | 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 DATA PREPROCESSING

In [4]:
```python
df.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 [5]:
```python
df.describe()
```

Out[5]:

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

In [6]:
```python
df.columns
```

Out[6]: Index(['Age', 'Sex', 'BP', 'Cholesterol', 'Na_to_K', 'Drug'], dtype='object')

In [7]:
```python
df1=df.dropna(axis=1)
df1
```

Out[7]:

|   | 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 |

| | Age | Sex | BP | Cholesterol | Na_to_K | Drug |
|---|---|---|---|---|---|---|
| **...** | ... | ... | ... | ... | ... | ... |
| **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 [8]:
```python
df1.columns
```

Out[8]: Index(['Age', 'Sex', 'BP', 'Cholesterol', 'Na_to_K', 'Drug'], dtype='object')

In [9]:
```python
df1=df1[['Age','Na_to_K']]
```

# EDA AND VISUALIZATION

In [10]:
```python
sns.pairplot(df1)
```
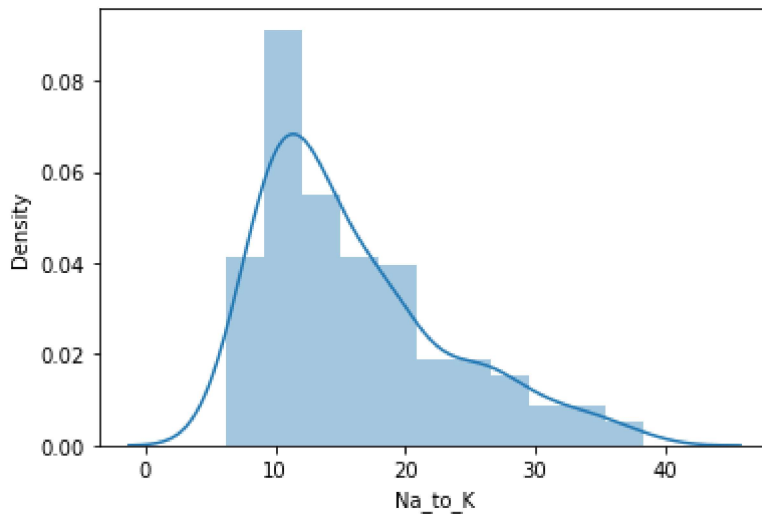
Out[10]: <seaborn.axisgrid.PairGrid at 0x27aa3162b20>

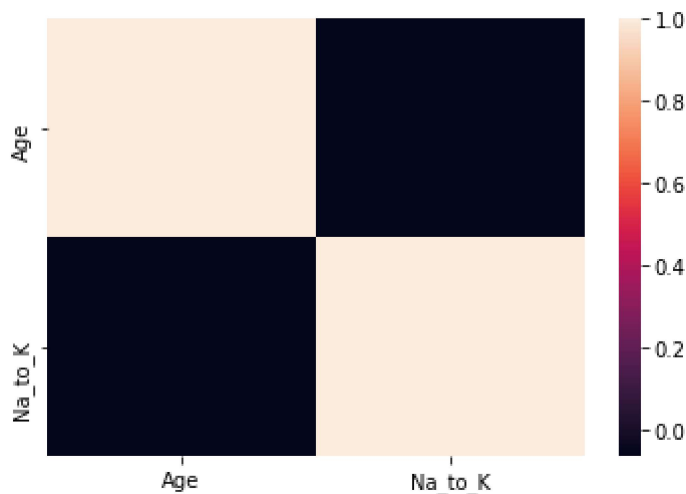

In [11]:
```python
sns.distplot(df1['Na_to_K'])
```

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 adap
t your code to use either `displot` (a figure-level function with similar flexibility) o
r `histplot` (an axes-level function for histograms).
  warnings.warn(msg, FutureWarning)

Out[11]: <AxesSubplot:xlabel='Na_to_K', ylabel='Density'>



In [12]:
```python
sns.heatmap(df1.corr())
```

Out[12]: <AxesSubplot:>



# TO TRAIN THE MODEL AND MODEL BULDING

In [13]:
```python
x=df[['Age','Na_to_K']]
y=df['Na_to_K']
```

In [14]:
```python
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)
```

In [15]:
```python
from sklearn.linear_model import LinearRegression
lr=LinearRegression()
```

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

Out[15]: LinearRegression()

In [16]:
```
lr.intercept_
```
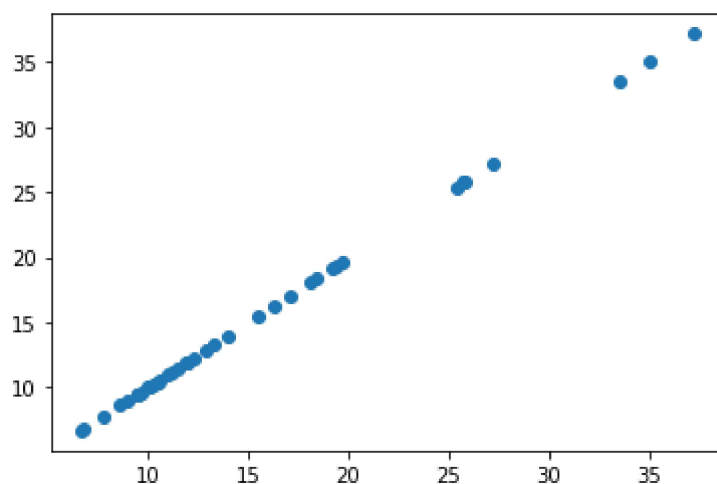
Out[16]: -1.0658141036401503e-14

In [17]:
```
coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
coeff
```

Out[17]:

|         | Co-efficient |
|---------|--------------|
| Age     | 0.0          |
| Na_to_K | 1.0          |

In [18]:
```
prediction =lr.predict(x_test)
plt.scatter(y_test,prediction)
```

Out[18]: <matplotlib.collections.PathCollection at 0x27aa542a670>



# ACCURACY

In [19]:
```
lr.score(x_test,y_test)
```

Out[19]: 1.0