

In [1]:

```
import numpy as np
import pandas as pd
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
```

In [2]:

```
df=pd.read_csv(r"C:\Users\monim\Downloads\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]:

```
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 |

In [4]:

```
df.tail()
```

Out[4]:

| | 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 |

In [5]:

```
df.shape
```

Out[5]:

```
(200, 6)
```

In [6]:

```
df.describe()
```

Out[6]:

| | 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 [7]:

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

Out[7]:

```
Age          0
Sex          0
BP           0
Cholesterol  0
Na_to_K      0
Drug         0
dtype: int64
```

In [8]:

```
df.isnull().any()
```

Out[8]:

```
Age           False
Sex           False
BP            False
Cholesterol   False
Na_to_K       False
Drug          False
dtype: bool
```

In [9]:

```
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 [10]:

```
df['Drug'].value_counts()
```

Out[10]:

```
Drug
drugY    91
drugX    54
drugA    23
drugC    16
drugB    16
Name: count, dtype: int64
```

In [11]:

```
df['Cholesterol'].value_counts()
```

Out[11]:

```
Cholesterol
HIGH      103
NORMAL    97
Name: count, dtype: int64
```

In [12]:

```
convert={"BP":{"LOW":1,"HIGH":2,"NORMAL":0}}
df=df.replace(convert)
df
```

Out[12]:

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

200 rows × 6 columns

In [13]:

```
convert={"Cholesterol":{"HIGH":1,"NORMAL":0}}
df=df.replace(convert)
df
```

Out[13]:

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

200 rows × 6 columns

In [32]:

```
x=["BP","Cholesterol"]  
y=["drugX","drugC","drugY","drugA","drugB"]  
all_inputs=df[x]  
all_classes=df["Drug"]
```

In [33]:

```
(x_train,x_test,y_train,y_test)=train_test_split(all_inputs,all_classes,test_size=0.2)
```

In [34]:

```
clf=DecisionTreeClassifier(random_state=0)
```

In [35]:

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

Out[35]:

```
DecisionTreeClassifier  
DecisionTreeClassifier(random_state=0)
```

In [30]:

```
score=clf.score(x_test,y_test)  
print(score)
```

0.525

In []: