

In [1]:

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

In [2]:

```
diabetes = pd.read_csv("D:/PGDAI - lec/Machine learning/Dataset/diabetes.csv")
diabetes.head(15)
```

Out[2]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction
0	6	148	72	35	0	33.6	0.625
1	1	85	66	29	0	26.6	0.351
2	8	183	64	0	0	23.3	0.672
3	1	89	66	23	94	28.1	0.167
4	0	137	40	35	168	43.1	2.278
5	5	116	74	0	0	25.6	0.258
6	3	78	50	32	88	31.0	0.254
7	10	115	0	0	0	35.3	0.161
8	2	197	70	45	543	30.5	0.161
9	8	125	96	0	0	0.0	0.232
10	4	110	92	0	0	37.6	0.161
11	10	168	74	0	0	38.0	0.516
12	10	139	80	0	0	27.1	1.472
13	1	189	60	23	846	30.1	0.349
14	5	166	72	19	175	25.8	0.516

In [3]:

diabetes.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Pregnancies                          768 non-null    int64
1   Glucose                              768 non-null    int64
2   BloodPressure                        768 non-null    int64
3   SkinThickness                        768 non-null    int64
4   Insulin                              768 non-null    int64
5   BMI                                  768 non-null    float64
6   DiabetesPedigreeFunction             768 non-null    float64
7   Age                                  768 non-null    int64
8   Outcome                              768 non-null    int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
```

In [4]:

diabetes.describe()

Out[4]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Diabetes
<b>count</b>	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	
<b>mean</b>	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	
<b>std</b>	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	
<b>min</b>	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
<b>25%</b>	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	
<b>50%</b>	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	
<b>75%</b>	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	
<b>max</b>	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	

In [5]:

diabetes.isnull().sum()

Out[5]:

```
Pregnancies      0
Glucose           0
BloodPressure     0
SkinThickness     0
Insulin           0
BMI               0
DiabetesPedigreeFunction  0
Age               0
Outcome           0
dtype: int64
```

In [6]:

```
#zeros or missing values will be replaced by the mean of that particular column
cols_clean = ['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'DiabetesPedigreeFu

for i in cols_clean:
    diabetes[i] = diabetes[i].replace(0,np.NaN)
    cols_mean = int(diabetes[i].mean(skipna=True))
    diabetes[i] = diabetes[i].replace(np.NaN,cols_mean)
data = diabetes
data.head(15)
```

Out[6]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction
0	6	148.0	72.0	35.0	155.0	33.6	0.627
1	1	85.0	66.0	29.0	155.0	26.6	0.351
2	8	183.0	64.0	29.0	155.0	23.3	0.672
3	1	89.0	66.0	23.0	94.0	28.1	0.167
4	0	137.0	40.0	35.0	168.0	43.1	2.278
5	5	116.0	74.0	29.0	155.0	25.6	0.245
6	3	78.0	50.0	32.0	88.0	31.0	0.245
7	10	115.0	72.0	29.0	155.0	35.3	0.167
8	2	197.0	70.0	45.0	543.0	30.5	0.167
9	8	125.0	96.0	29.0	155.0	32.0	0.245
10	4	110.0	92.0	29.0	155.0	37.6	0.167
11	10	168.0	74.0	29.0	155.0	38.0	0.516
12	10	139.0	80.0	29.0	155.0	27.1	1.472
13	1	189.0	60.0	23.0	846.0	30.1	0.351
14	5	166.0	72.0	19.0	175.0	25.8	0.516

In [7]:

```
x = diabetes.iloc[:, :-1].values
y = diabetes.iloc[:, 8].values
```

In [8]:

```
print(x)
```

```
[[ 6.  148.  72.  ...  33.6  0.627  50.  ]
 [ 1.   85.  66.  ...  26.6  0.351  31.  ]
 [ 8.  183.  64.  ...  23.3  0.672  32.  ]
 ...
 [ 5.  121.  72.  ...  26.2  0.245  30.  ]
 [ 1.  126.  60.  ...  30.1  0.349  47.  ]
 [ 1.   93.  70.  ...  30.4  0.315  23.  ]]
```

In [9]:

```
print(y)
```

```
[1 0 1 0 1 0 1 0 1 1 0 1 0 1 1 1 1 1 0 1 0 0 1 1 1 1 1 0 0 0 0 1 0 0 0 0
 1 1 1 0 0 0 1 0 1 0 0 1 0 0 0 0 1 0 0 1 0 0 0 0 1 0 0 1 0 1 0 0 0 0 1 0 1 0
 0 0 0 0 1 0 0 0 0 0 1 0 0 0 1 0 0 0 0 1 0 0 0 0 0 1 1 0 0 0 0 0 0 0 0 0 1 1
 1 0 0 1 1 1 0 0 0 1 0 0 0 1 1 0 0 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0
 0 0 0 0 1 0 1 1 0 0 0 1 0 0 0 0 1 1 0 0 0 0 1 1 0 0 0 1 0 1 0 1 0 0 0 0 0
 1 1 1 1 1 0 0 1 1 0 1 0 1 1 1 0 0 0 0 0 0 1 1 0 1 0 0 0 1 1 1 1 0 1 1 1 1
 0 0 0 0 0 1 0 0 1 1 0 0 0 1 1 1 1 0 0 0 1 1 0 1 0 0 0 0 0 0 0 0 0 1 1 0 0 0
 1 0 1 0 0 1 0 1 0 0 1 1 0 0 0 0 0 1 0 0 0 1 0 0 1 1 0 0 1 0 0 0 1 1 1 0 0
 1 0 1 0 1 1 0 1 0 0 1 0 1 1 0 0 1 0 1 0 0 1 0 1 0 1 1 1 0 0 1 0 1 0 0 0 1
 0 0 0 0 1 1 1 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 1 1 1 0 1 1 0 0 1 0 0 1 0 0 1
 1 0 0 0 0 1 0 0 1 0 0 0 0 0 0 0 1 1 1 0 0 1 0 0 1 0 0 1 0 1 1 0 1 0 1 0 1
 0 1 1 0 0 0 0 1 1 0 1 0 1 0 0 0 0 1 1 0 1 0 1 0 0 0 0 0 1 0 0 0 0 1 0 0 1
 1 1 0 0 1 0 0 1 0 0 0 1 0 0 1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 0 0 1
 0 0 0 1 1 0 0 0 0 0 0 0 1 0 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 0 1 1 0
 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 1 1 1 1 0 0 1 1 0 0 0 0 0 0 0 0
 0 0 0 0 0 1 1 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 1 1 0 0 0 1 0 1 0 1 0 1 0
 1 0 0 1 0 0 1 0 0 0 0 1 1 0 1 0 0 0 0 1 1 0 1 0 0 0 1 1 0 0 0 0 0 0 0 0
 0 1 0 0 0 0 1 0 0 1 0 0 0 1 0 0 0 1 1 1 0 0 0 0 0 0 1 0 0 0 1 0 1 1 1 1 0
 1 1 0 0 0 0 0 0 0 1 1 0 1 0 0 1 0 1 0 0 0 0 0 1 0 1 0 1 0 1 1 0 0 0 0 1 1
 0 0 0 1 0 1 1 0 0 1 0 0 1 1 0 0 1 0 0 1 0 0 0 0 0 0 0 1 1 1 0 0 0 0 0 0 1
 1 0 0 1 0 0 1 0 1 1 1 0 0 1 1 1 0 1 0 1 0 1 0 0 0 0 1 0]
```

In [10]:

```
# Splitting the dataset into the Training set and Test set
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 [11]:

```
# Feature Scaling to bring the variable in a single scale
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
```

In [13]:

```
# Fitting Naive Bayes to the Training set
from sklearn.naive_bayes import GaussianNB
classifier = GaussianNB()
classifier.fit(X_train, y_train)
```

Out[13]:

GaussianNB()

In [15]:

```
classifier.score(X_test,y_test)
```

Out[15]:

0.7662337662337663

In [17]:

```
#Predicting the Test set results
y_pred = classifier.predict(X_test)
y_pred
```

Out[17]:

```
array([0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1,
        0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0,
        0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0,
        1, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0,
        0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0,
        0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0],
      dtype=int64)
```

In [18]:

```
# confusion matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test,y_pred)
cm
```

Out[18]:

```
array([[83, 12],
       [24, 35]], dtype=int64)
```

In [19]:

```
from sklearn.metrics import accuracy_score
acc = accuracy_score(y_test,y_pred)
print("Accuracy : {0:.2f}%".format(acc*100))
```

Accuracy : 76.62%

In [20]:

```
from sklearn import metrics
print(metrics.classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	0.78	0.87	0.82	95
1	0.74	0.59	0.66	59
accuracy			0.77	154
macro avg	0.76	0.73	0.74	154
weighted avg	0.76	0.77	0.76	154

In [ ]:

