

Importing libraries and Loading Data set

```
In [15]: # Load the diabetes dataset
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn import svm, metrics, model_selection

path = './diabetes.csv'

data = pd.read_csv(path)
print(data.shape)
data
```

(768, 9)

Out[15]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1
...
763	10	101	76	48	180	32.9	0.171	63	0
764	2	122	70	27	0	36.8	0.340	27	0
765	5	121	72	23	112	26.2	0.245	30	0
766	1	126	60	0	0	30.1	0.349	47	1
767	1	93	70	31	0	30.4	0.315	23	0

768 rows × 9 columns

Data Preprocessing

In [4]: data.head()

Out[4]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

In [2]: data.tail()

Out[2]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
763	10	101	76	48	180	32.9	0.171	63	0
764	2	122	70	27	0	36.8	0.340	27	0
765	5	121	72	23	112	26.2	0.245	30	0
766	1	126	60	0	0	30.1	0.349	47	1
767	1	93	70	31	0	30.4	0.315	23	0

In [7]: data.dtypes

Out[7]: Pregnancies int64
Glucose int64
BloodPressure int64
SkinThickness int64
Insulin int64
BMI float64
DiabetesPedigreeFunction float64
Age int64
Outcome int64
dtype: object

```
In [19]: data.columns
```

```
Out[19]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',  
              'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],  
              dtype='object')
```

```
In [10]: data.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 [18]: data.describe()
```

```
Out[18]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	0.471876	33.240885	0.348958
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.331329	11.760232	0.476951
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	21.000000	0.000000
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.243750	24.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	29.000000	0.000000
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.626250	41.000000	1.000000
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	81.000000	1.000000

checking for missing value

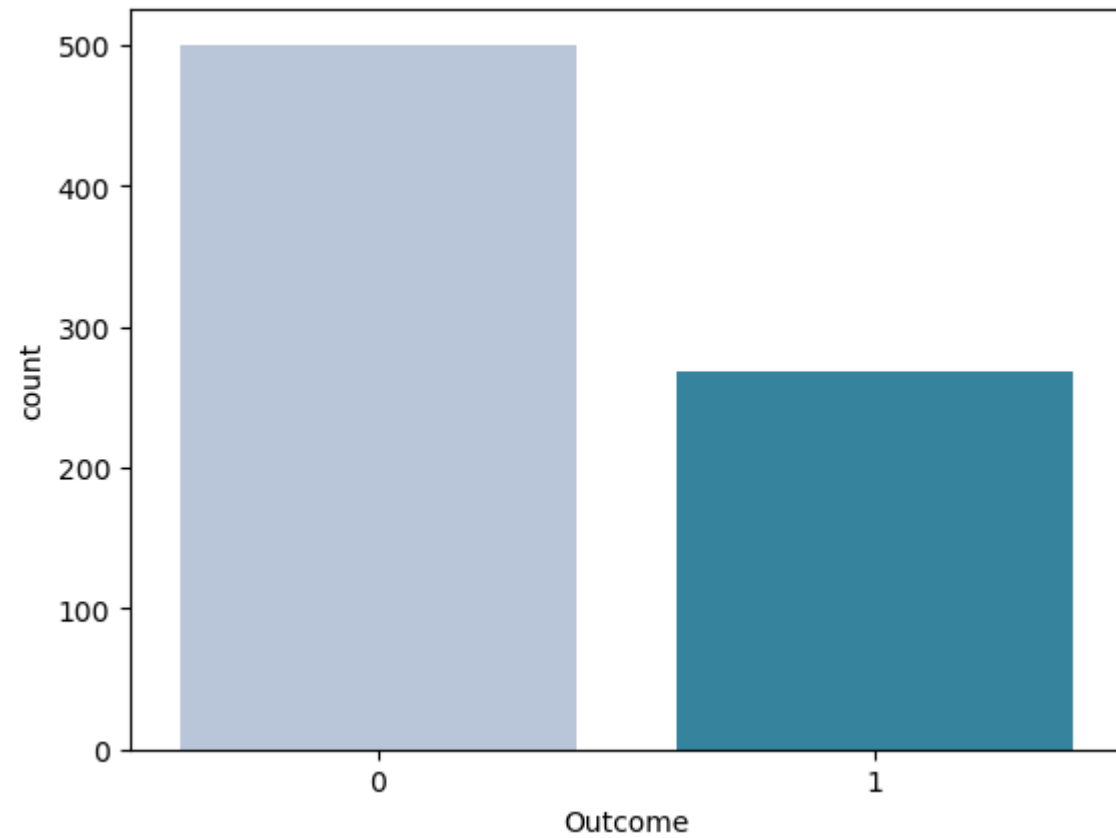
```
In [13]: data.isnull().sum()
```

```
Out[13]: Pregnancies      0
          Glucose         0
          BloodPressure   0
          SkinThickness   0
          Insulin         0
          BMI             0
          DiabetesPedigreeFunction  0
          Age             0
          Outcome         0
          dtype: int64
```

Data visualization

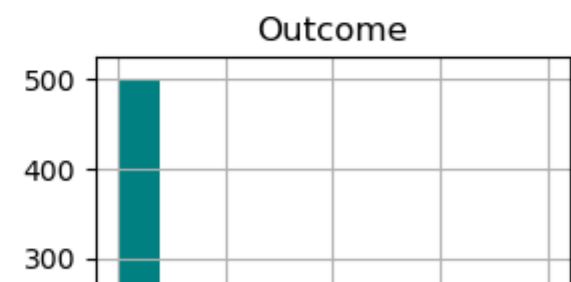
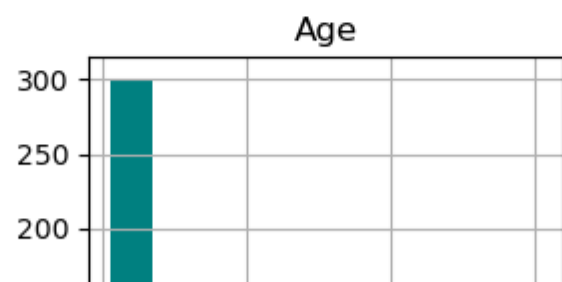
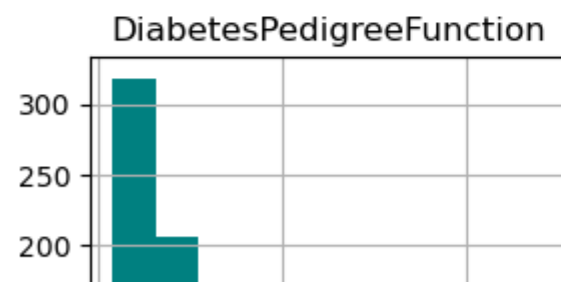
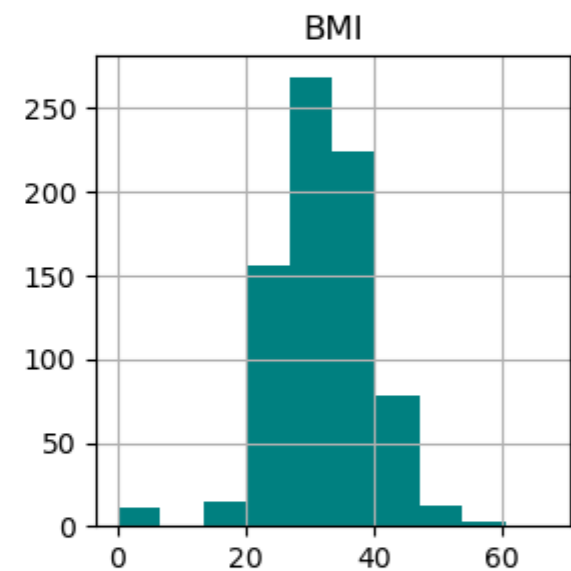
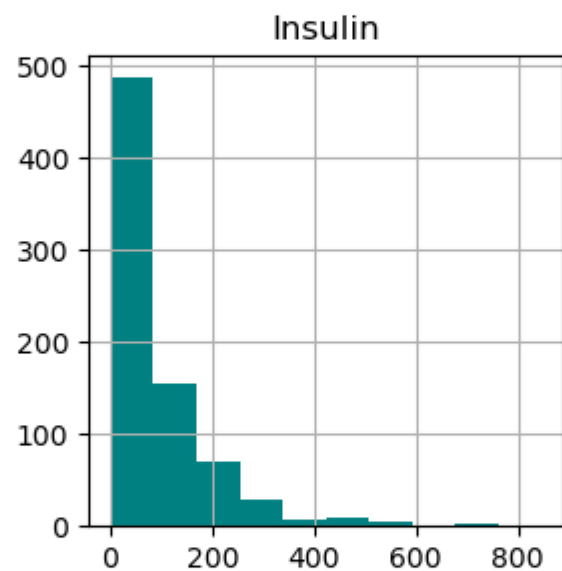
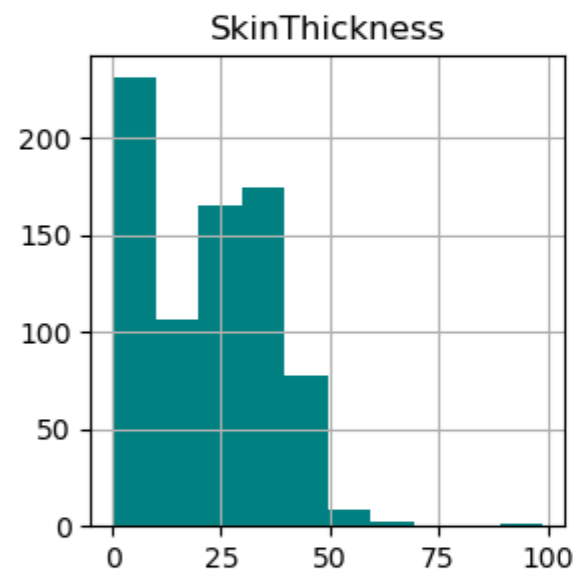
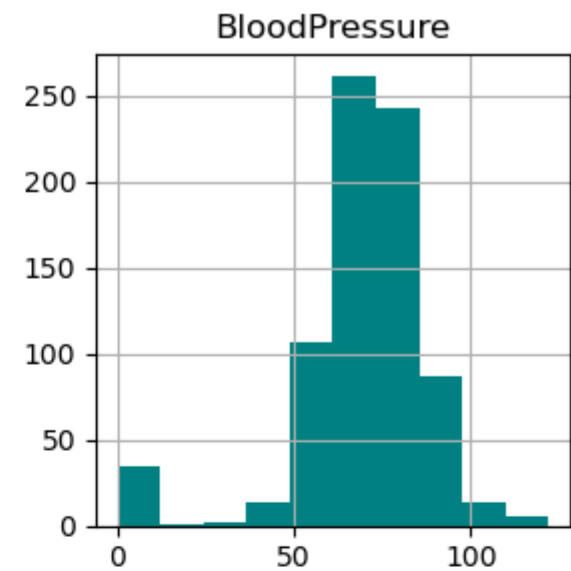
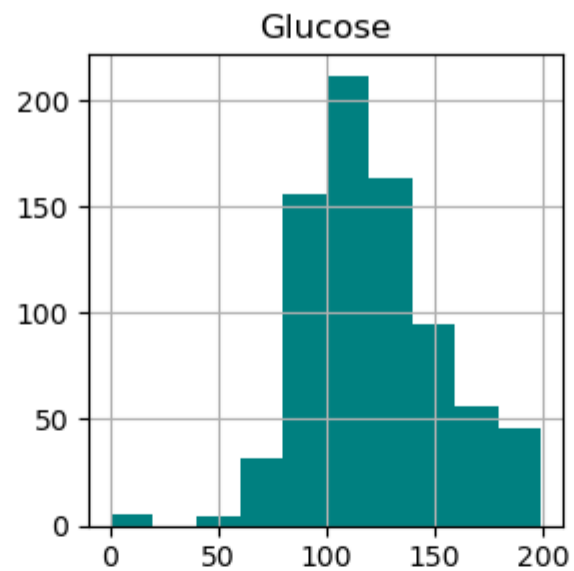
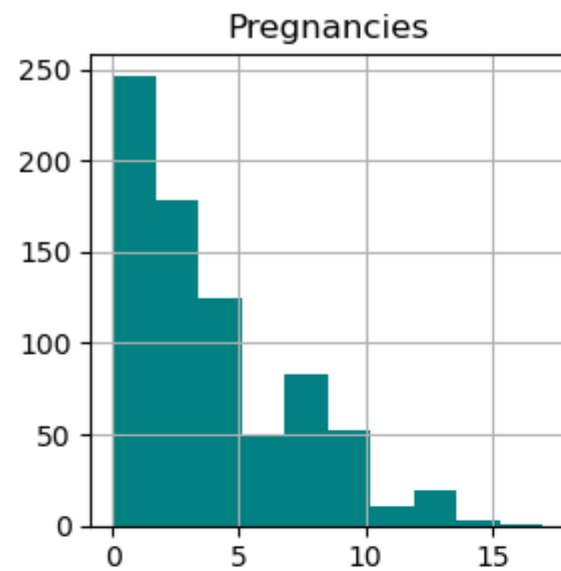
```
In [16]: sns.countplot(data=data,x="Outcome",palette='PuBuGn')
```

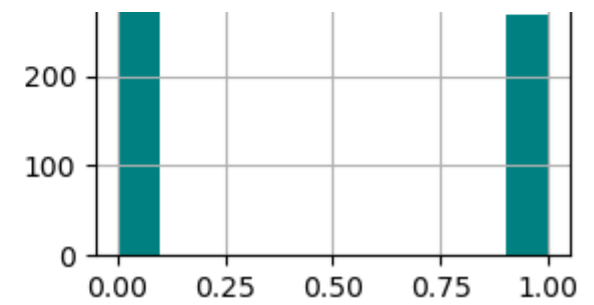
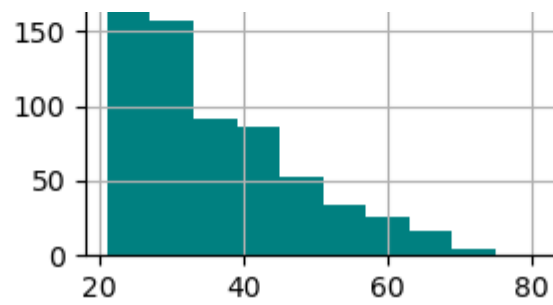
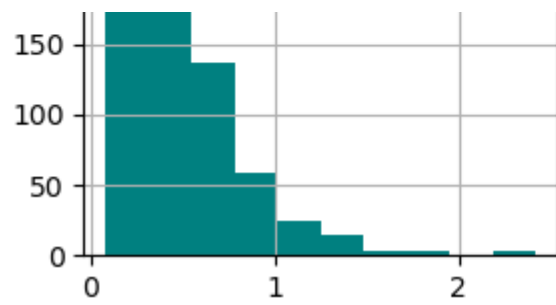
```
Out[16]: <AxesSubplot:xlabel='Outcome', ylabel='count'>
```



```
In [17]: data.hist(figsize=(11,11),color='teal')
```

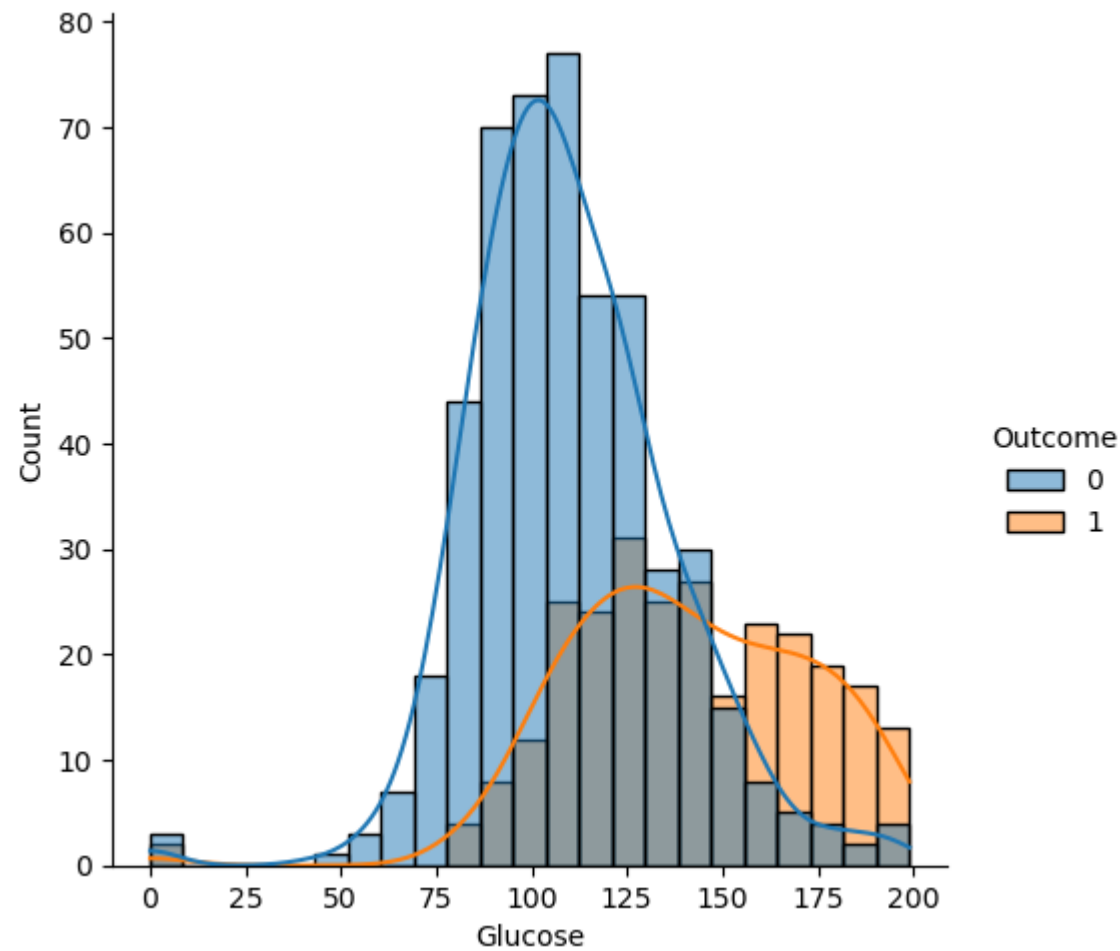
```
Out[17]: array([[<AxesSubplot:title={'center':'Pregnancies'}>,
                  <AxesSubplot:title={'center':'Glucose'}>,
                  <AxesSubplot:title={'center':'BloodPressure'}>],
                [<AxesSubplot:title={'center':'SkinThickness'}>,
                  <AxesSubplot:title={'center':'Insulin'}>,
                  <AxesSubplot:title={'center':'BMI'}>],
                [<AxesSubplot:title={'center':'DiabetesPedigreeFunction'}>,
                  <AxesSubplot:title={'center':'Age'}>,
                  <AxesSubplot:title={'center':'Outcome'}>]], dtype=object)
```



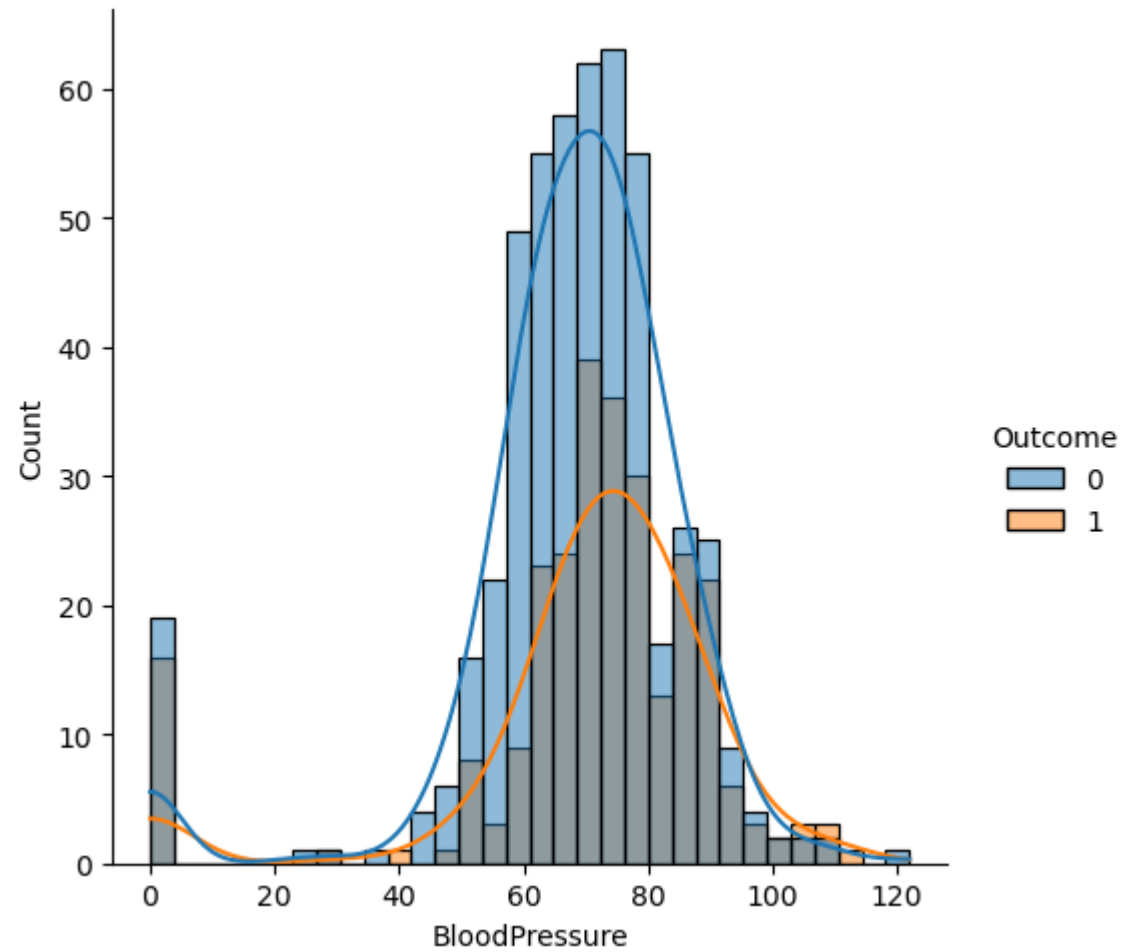
```
In [21]: sns.displot(x='Glucose',hue="Outcome",kde=True,data=data,color="PuBuGn")
```

```
Out[21]: <seaborn.axisgrid.FacetGrid at 0x1800186b370>
```



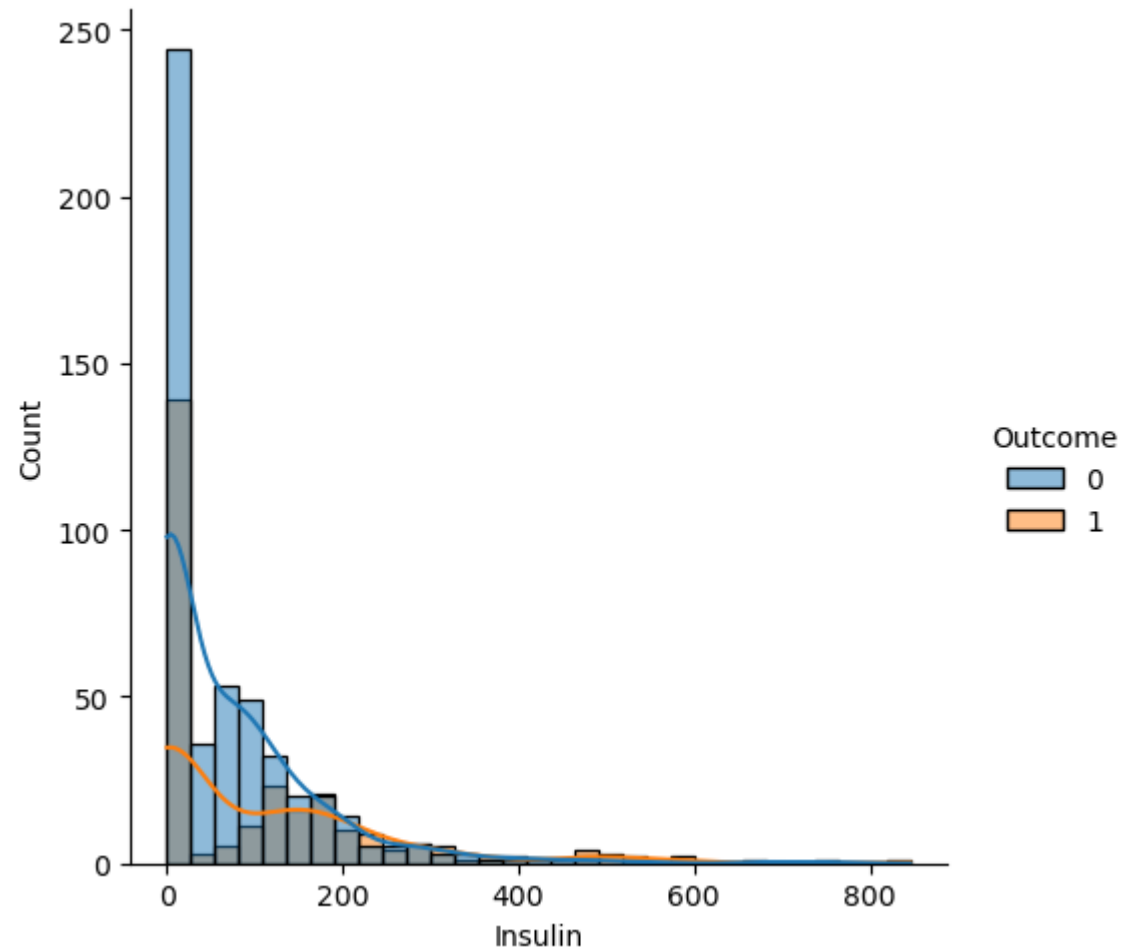
```
In [22]: sns.displot(x='BloodPressure',hue="Outcome",kde=True,data=data,color="PuBuGn")
```

```
Out[22]: <seaborn.axisgrid.FacetGrid at 0x1800156e6a0>
```



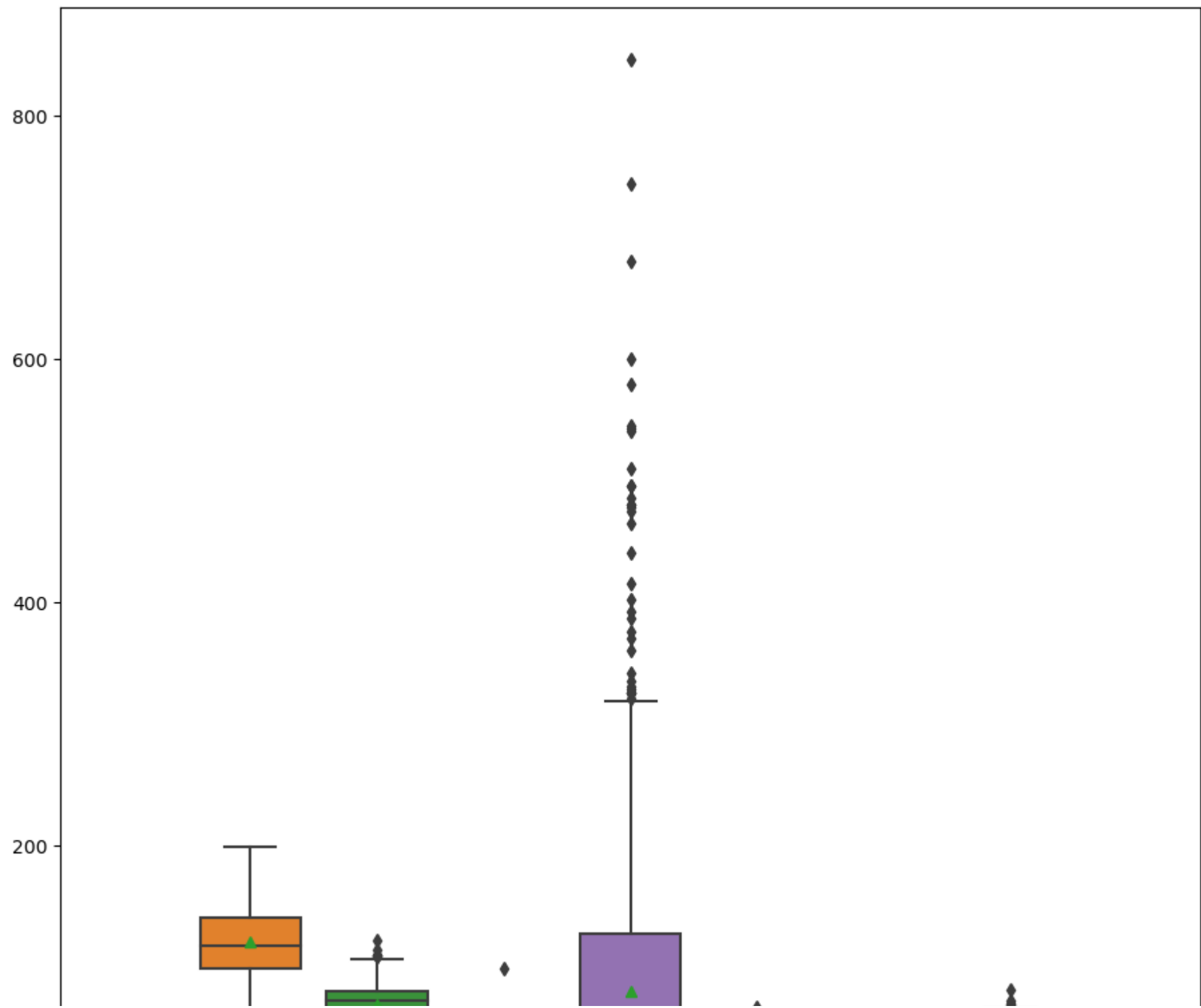
```
In [23]: sns.displot(x='Insulin',hue="Outcome",kde=True,data=data,color="PuBuGn")
```

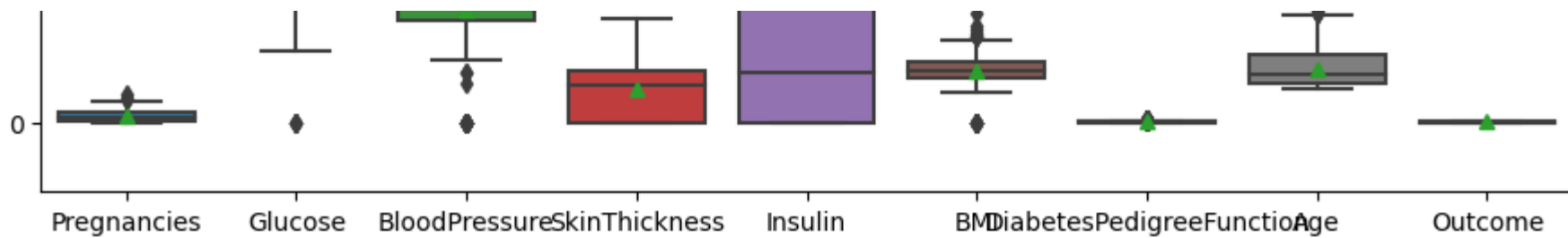
```
Out[23]: <seaborn.axisgrid.FacetGrid at 0x18002e9ffa0>
```



```
In [29]: plt.figure(figsize=(11,11))  
sns.boxplot(data=data, showmeans=True)
```

```
Out[29]: <AxesSubplot:>
```



Outlier Treatment

```
In [30]: from scipy import stats
z_score_threshold=3.0
z_score=np.abs(stats.zscore(data))
outlier_mask=(z_score >z_score_threshold).any(axis=1)
data1=data[ outlier_mask]
```

Input Output sepration

```
In [32]: x=data.iloc[:, :-1].values
x
```

```
Out[32]: array([[ 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 [34]: y=data.iloc[:, -1].values  
y
```

```
Out[34]: array([1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0,  
1, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 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, 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, 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, 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, 1, 1, 0, 0, 0, 1, 0, 0,  
1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 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, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0,  
1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0,  
1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 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, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 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, 0, 0, 1, 0, 0, 0, 0, 1, 0,  
1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 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, 0, 1, 1, 0, 0, 1, 0, 0, 1,  
1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0],  
dtype=int64)
```


Training and Testing

```
In [37]: from sklearn import model_selection
xtrain,xtest,ytrain,ytest = model_selection.train_test_split(x,y,test_size=0.3,random_state=10)
xtrain
```

```
Out[37]: array([[2.00e+00, 8.90e+01, 9.00e+01, ..., 3.35e+01, 2.92e-01, 4.20e+01],
 [4.00e+00, 1.46e+02, 8.50e+01, ..., 2.89e+01, 1.89e-01, 2.70e+01],
 [1.00e+01, 1.11e+02, 7.00e+01, ..., 2.75e+01, 1.41e-01, 4.00e+01],
 ...,
 [3.00e+00, 1.16e+02, 7.40e+01, ..., 2.63e+01, 1.07e-01, 2.40e+01],
 [1.00e+00, 8.80e+01, 3.00e+01, ..., 5.50e+01, 4.96e-01, 2.60e+01],
 [5.00e+00, 9.60e+01, 7.40e+01, ..., 3.36e+01, 9.97e-01, 4.30e+01]])
```

```
In [38]: xtest
```

```
Out[38]: array([[ 4.   , 154.   , 72.   , ..., 31.3   , 0.338, 37.   ],
 [ 2.   , 112.   , 86.   , ..., 38.4   , 0.246, 28.   ],
 [ 1.   , 135.   , 54.   , ..., 26.7   , 0.687, 62.   ],
 ...,
 [ 3.   , 150.   , 76.   , ..., 21.   , 0.207, 37.   ],
 [ 3.   , 130.   , 64.   , ..., 23.1   , 0.314, 22.   ],
 [ 0.   , 108.   , 68.   , ..., 27.3   , 0.787, 32.   ]])
```

```
In [39]: ytrain
```

```
Out[39]: array([0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1,
                1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1,
                0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1,
                0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0, 0,
                0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0,
                0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1,
                0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 1, 1, 0, 1, 1, 0, 0, 0,
                1, 1, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1,
                0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0,
                0, 0, 0, 1, 1, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0,
                1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0, 0, 1, 0,
                0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0,
                1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0,
                1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0,
                0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0,
                0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0,
                0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 1, 0,
                1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 0, 0, 1,
                0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 0, 0, 1, 0,
                1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0,
                1, 0, 0, 0, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0,
                1, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1,
                0, 0, 0, 0, 1, 0, 0, 1, 0], dtype=int64)
```

```
In [40]: ytest
```

```
Out[40]: array([0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0,
                0, 1, 0, 1, 0, 1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 1, 1, 0,
                1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0,
                1, 0, 0, 1, 1, 1, 0, 1, 0, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1,
                1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0,
                0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 1,
                0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 0,
                1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 0,
                0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 0,
                0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1,
                0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0], dtype=int64)
```

Normalization

```
In [41]: from sklearn.preprocessing import StandardScaler
scaler=StandardScaler()
scaler.fit(xtrain)
xtrain=scaler.transform(xtrain)
xtest=scaler.transform(xtest)
xtrain
```

```
Out[41]: array([[ -0.53788077, -0.98266127,  1.07584045, ...,  0.19662751,
                -0.51280558,  0.74882677],
                [ 0.06137644,  0.74886132,  0.81461495, ..., -0.39218237,
                -0.8218938 , -0.52209603],
                [ 1.85914807, -0.31435431,  0.03093844, ..., -0.57138538,
                -0.96593491,  0.57937039],
                ...,
                [-0.23825217, -0.16246636,  0.23991884, ..., -0.72498795,
                -1.06796403, -0.77628059],
                [-0.83750938, -1.01303886, -2.05886558, ...,  2.94867368,
                 0.09936914, -0.60682422],
                [ 0.36100504, -0.77001815,  0.23991884, ...,  0.20942772,
                 1.60279824,  0.83355495]])
```

```
In [42]: xtest
```

```
Out[42]: array([[ 0.06137644,  0.99188203,  0.13542864, ..., -0.08497722,
                -0.37476618,  0.32518583],
                [-0.53788077, -0.28397672,  0.86686005, ...,  0.82383803,
                -0.65084498, -0.43736785],
                [-0.83750938,  0.41470783, -0.80498317, ..., -0.6737871 ,
                 0.67253273,  2.44339049],
                ...,
                [-0.23825217,  0.87037167,  0.34440904, ..., -1.40339934,
                -0.76787838,  0.32518583],
                [-0.23825217,  0.26281989, -0.28253216, ..., -1.13459483,
                -0.44678674, -0.94573696],
                [-1.13713798, -0.40548708, -0.07355176, ..., -0.59698581,
                 0.97261838, -0.0984551 ]])
```

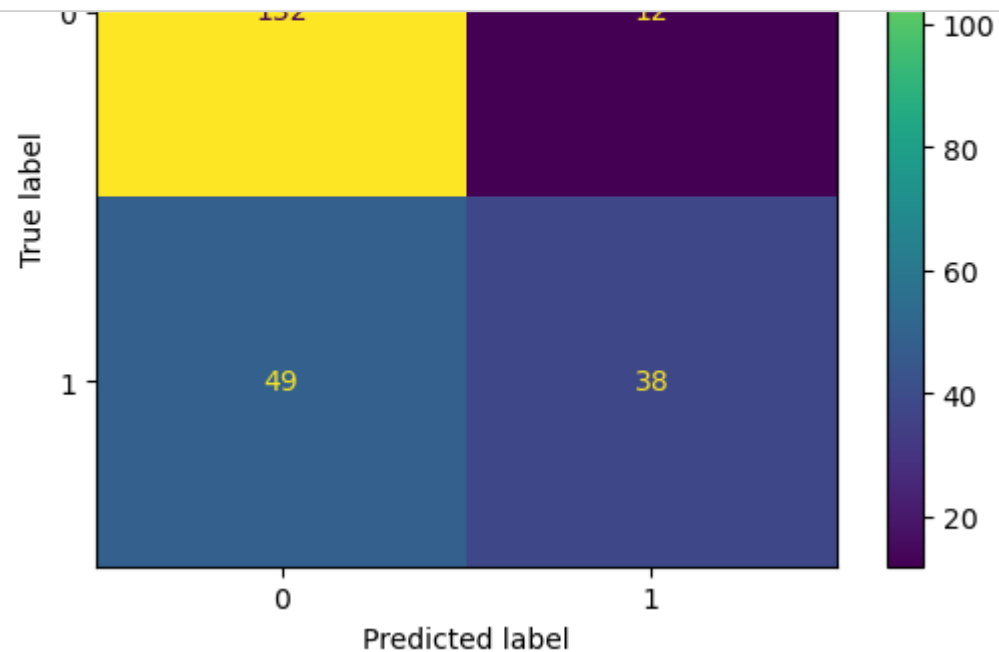
Modelselection using Classification Algorithms

Algorithms
K-Nearest Neighbors
Naive Bayes
Support Vector Machine
Decision Tree
Random Forest

```
In [47]: from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import ConfusionMatrixDisplay
knn=KNeighborsClassifier()
svm=SVC()
nb=GaussianNB()
dt=DecisionTreeClassifier(criterion='entropy')
rf=RandomForestClassifier(n_estimators=10,criterion='entropy')
lst=[knn,svm,nb,dt,rf]
```

Perfomance Evaluvation

```
In [49]: from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
for i in lst:
    print(i)
    i.fit(xtrain, ytrain)
    y_pred = i.predict(xtest)
    print(accuracy_score(ytest, y_pred))
    print('*****')
    print(confusion_matrix(ytest, y_pred))
    print('*****')
    print(classification_report(ytest, y_pred))
    print('*****')
    result = confusion_matrix(ytest, y_pred)
    labels = [0, 1]
    cmd = ConfusionMatrixDisplay(result, display_labels=labels)
    cmd.plot()
```



Model prediction

```
In [53]: input_data=nb.predict(scaler.transform([[2,148,66,27,219,28.1,0.313,50]]))
if input_data==0:
    print('The person is NOT Diabetic')
else:
    print('The person is Diabetic')
```

The person is Diabetic

Conclusion

It is concluded that out of the 5 models Random Forest is the best performing model as it has the highest accuracy of 76% as compared to all the models