

Project : Diabetes Prediction

Parameter which are invol

- 1. Pregnancies: Number of times pregnant
- 2. Glucose: Plasma glucose concentration a 2 hours in an oral glucose tolerance test
- 3. BloodPressure: Diastolic blood pressure (mm Hg)
- 4. SkinThickness: Triceps skin fold thickness (mm)
- 5. Insulin: 2-Hour serum insulin (mu U/ml)
- 6. BMI: Body mass index (weight in kg/(height in m)^2)
- 7. DiabetesPedigreeFunction: Diabetes pedigree function
- 8. Age: Age (years)
- 9. Outcome: Class variable (0 or 1)

Importing Requried Python Libraries

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

Uploading Dataset

We have our data saved in a CSV file called Diabetes_1.csv. We first read our dataset into a pandas dataframe called data, and then use the head() function to show the first five records from our dataset.

In [2]:

```
1 # Loading the data using pandas
2
3 Data = pd.read_csv('Diabetes_1.csv')
4
5 # Displaying top 5 rows
6
7 Data.head()
```

Out[2]:

	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

Exploratory Data Analysis

In [3]:

```
1 # Checking for Null Values in Dataset
2
3 Data.isnull().sum()
```

Out[3]:

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

In [4]:

```
1 # Summery of Dataframe
2
3 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 [5]:

```
1 # Displaying total rows and columns
2
3 Data.shape
```

Out[5]: (768, 9)

In [6]:

```
1 # Display Distinct values
2
3 Data['Outcome'].unique()
```

Out[6]: array([1, 0], dtype=int64)

In [7]:

```
1 # Display Distinct values
2
3 Data['Age'].unique()
```

Out[7]: array([50, 31, 32, 21, 33, 30, 26, 29, 53, 54, 34, 57, 59, 51, 27, 41, 43, 22, 38, 60, 28, 45, 35, 46, 56, 37, 48, 40, 25, 24, 58, 42, 44, 39, 36, 23, 61, 69, 62, 55, 65, 47, 52, 66, 49, 63, 67, 72, 81, 64, 70, 68], dtype=int64)

In [8]:

```
1 # Classified Age into Young , Middle , Siniors
2
3 def Age_Class(Age):
4
5     if Age <= 35:
6         return 'Young Age'
7
8     elif Age <= 60:
9         return 'Middle Age'
10
11    else:
12        return 'Siniors'
13
14 Data['Age Classification'] = Data['Age'].apply(Age_Class)
15 Data.head()
```

Out[8]:

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

```
In [9]: 1 # Classified BMI into Under Weight , Normal Weight , Over Weight
2
3 def BMI_Class(BMI):
4
5     if BMI <= 45:
6         return 'Under Weight'
7
8     elif BMI <= 60:
9         return 'Normal Weight'
10
11    else:
12        return 'Over Weight'
13
14 Data['BMI Classification'] = Data['BMI'].apply(BMI_Class)
15 Data.head()
```

Out[9]:

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

```
In [10]: 1 # Display Descriptive Statistics of Data
2
3 Data.describe().T
```

Out[10]:

	count	mean	std	min	25%	50%	75%	max
Pregnancies	768.0	3.845052	3.369578	0.000	1.00000	3.0000	6.00000	17.00
Glucose	768.0	120.894531	31.972618	0.000	99.00000	117.0000	140.25000	199.00
BloodPressure	768.0	69.105469	19.355807	0.000	62.00000	72.0000	80.00000	122.00
SkinThickness	768.0	20.536458	15.952218	0.000	0.00000	23.0000	32.00000	99.00
Insulin	768.0	79.799479	115.244002	0.000	0.00000	30.5000	127.25000	846.00
BMI	768.0	31.992578	7.884160	0.000	27.30000	32.0000	36.60000	67.10
DiabetesPedigreeFunction	768.0	0.471876	0.331329	0.078	0.24375	0.3725	0.62625	2.42
Age	768.0	33.240885	11.760232	21.000	24.00000	29.0000	41.00000	81.00
Outcome	768.0	0.348958	0.476951	0.000	0.00000	0.0000	1.00000	1.00

```
In [11]: 1 # Display Average value of each Parameter from Table
2
3 round(Data.mean() , 2)
```

C:\Users\admin\AppData\Local\Temp\ipykernel_12336\3818815197.py:3: FutureWarning: The default value of numeric_only in DataFrame.mean is deprecated. In a future version, it will default to False. In addition, specifying 'numeric_only=None' is deprecated. Select only valid columns or specify the value of numeric_only to silence this warning.

```
round(Data.mean() , 2)
```

```
Out[11]: Pregnancies      3.85
Glucose      120.89
BloodPressure  69.11
SkinThickness  20.54
Insulin      79.80
BMI          31.99
DiabetesPedigreeFunction  0.47
Age          33.24
Outcome      0.35
dtype: float64
```

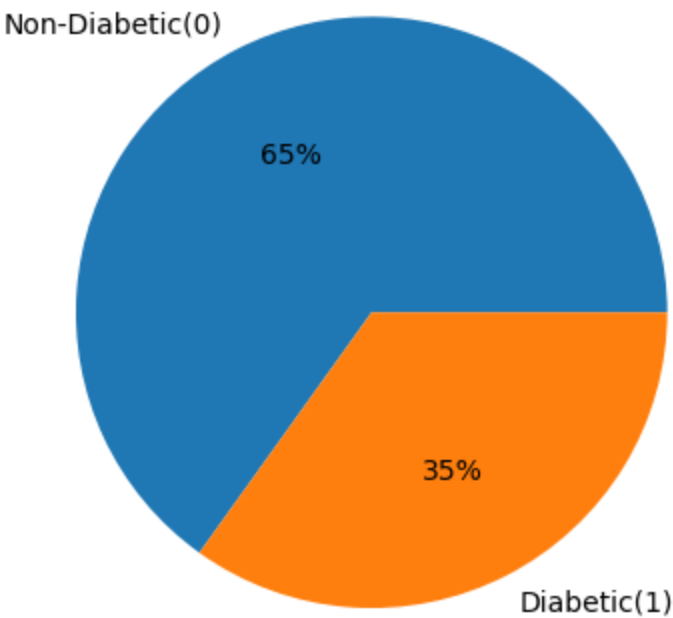
Data Visualization

```
In [12]: 1 # Display Count Diabetic Patient and Not Diabetic Patients
2
3 Data['Outcome'].value_counts(normalize = True) * 100
```

```
Out[12]: 0    65.104167
1    34.895833
Name: Outcome, dtype: float64
```

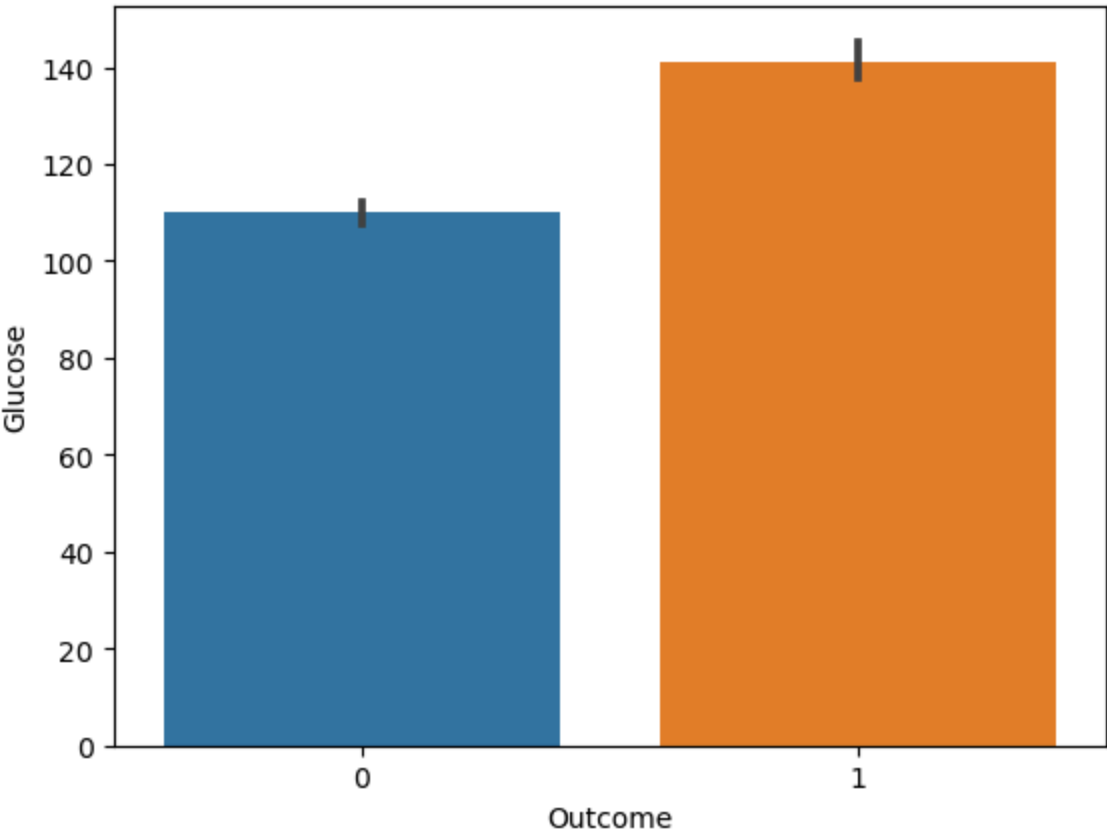
```
In [13]: 1 plt.pie(Data['Outcome'].value_counts(normalize = True) * 100 , labels = ['Non-Diabetic(0)' , 'Diabetic(1)']
```

Out[13]: ([<matplotlib.patches.Wedge at 0x13e781066d0>, <matplotlib.patches.Wedge at 0x13e78118490>], [Text(-0.5025941410457841, 0.9784677457057288, 'Non-Diabetic(0)'), Text(0.502593957824251, -0.9784678398182308, 'Diabetic(1)')], [Text(-0.27414225875224585, 0.533709679475852, '65%'), Text(0.2741421588132278, -0.533709730809944, '35%')])



```
In [14]: 1 # Plot for Glucose and Insulin
2
3 sns.barplot(y = 'Glucose' , x = 'Outcome' , data = Data)
```

Out[14]: <Axes: xlabel='Outcome', ylabel='Glucose'>



```
In [15]: 1 Age_C = Data.groupby(['Age Classification', 'Outcome']).size().unstack(fill_value = 0)
2 Age_C
```

Out[15]:

Outcome	0	1
Age Classification		
Middle Age	113	130
Seniors	20	7
Young Age	367	131

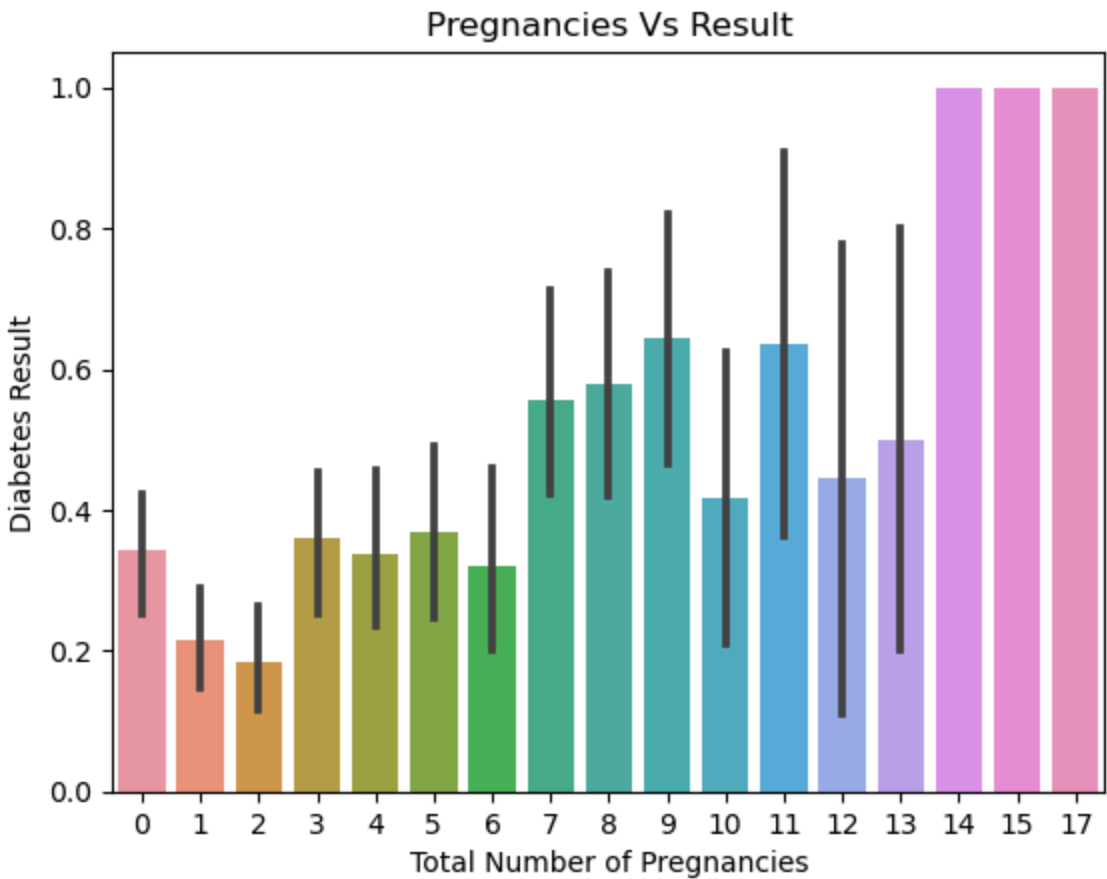
```
In [16]: 1 BMI_C = Data.groupby(['BMI Classification', 'Outcome']).size().unstack(fill_value = 0)
2 BMI_C
```

Out[16]:

Outcome	0	1
BMI Classification		
Normal Weight	13	21
Over Weight	0	1
Under Weight	487	246

```
In [17]: 1 sns.barplot(x = 'Pregnancies' , y = 'Outcome' , data = Data)
2 plt.xlabel("Total Number of Pregnancies")
3 plt.ylabel("Diabetes Result")
4 plt.title("Pregnancies Vs Result")
```

Out[17]: Text(0.5, 1.0, 'Pregnancies Vs Result')



```
In [18]: 1 plt.figure(figsize = (20,5))
2
3 plt.subplot(1,2,1)
4 plt.title('Counter Plot')
5 sns.countplot(x = 'Pregnancies',data = Data)
6
7 plt.subplot(1,2,2)
8 plt.title('Distribution Plot')
9 sns.distplot(Data["Pregnancies"])
```

C:\Users\admin\AppData\Local\Temp\ipykernel_12336\913136878.py:9: UserWarning:

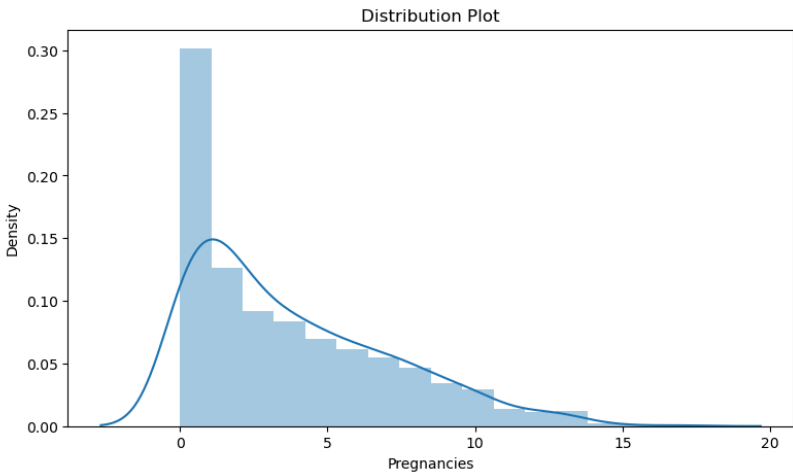
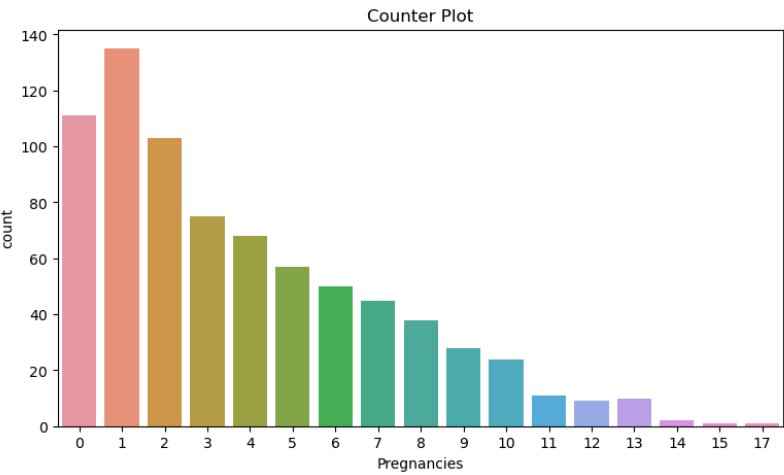
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

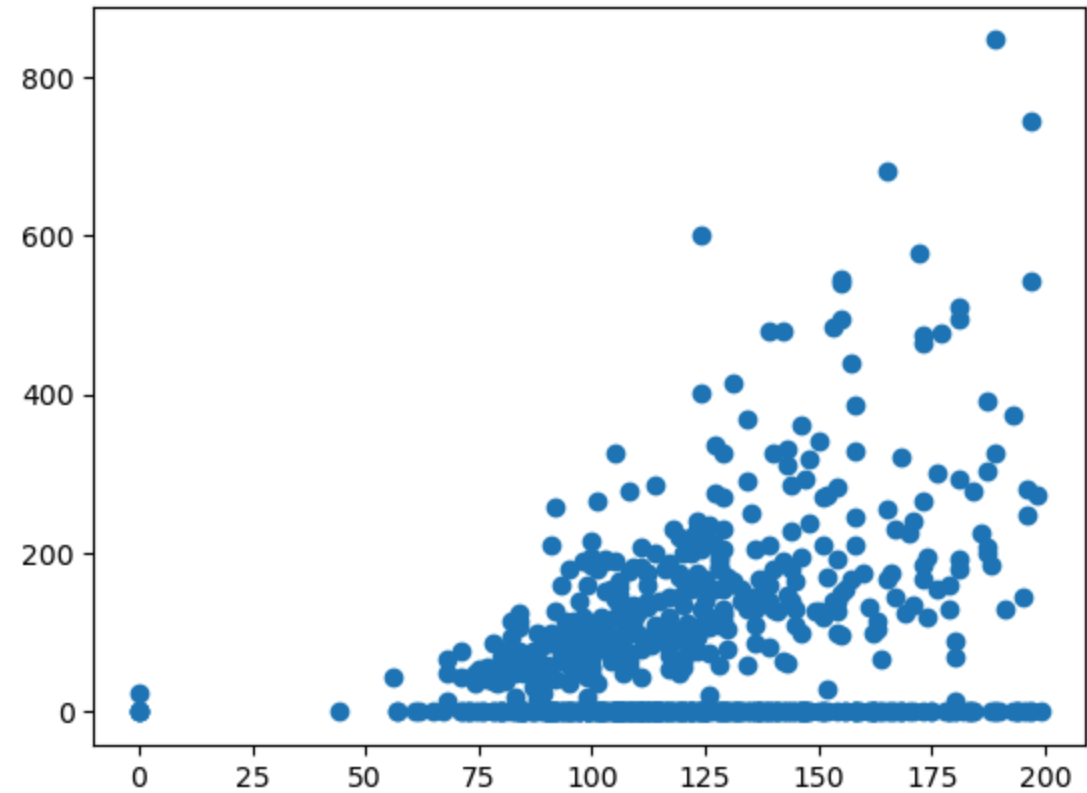
```
sns.distplot(Data["Pregnancies"])
```

Out[18]: <Axes: title={'center': 'Distribution Plot'}, xlabel='Pregnancies', ylabel='Density'>



```
In [19]: 1 plt.scatter(Data['Glucose'] , Data['Insulin'])
```

Out[19]: <matplotlib.collections.PathCollection at 0x13e78b2ee50>



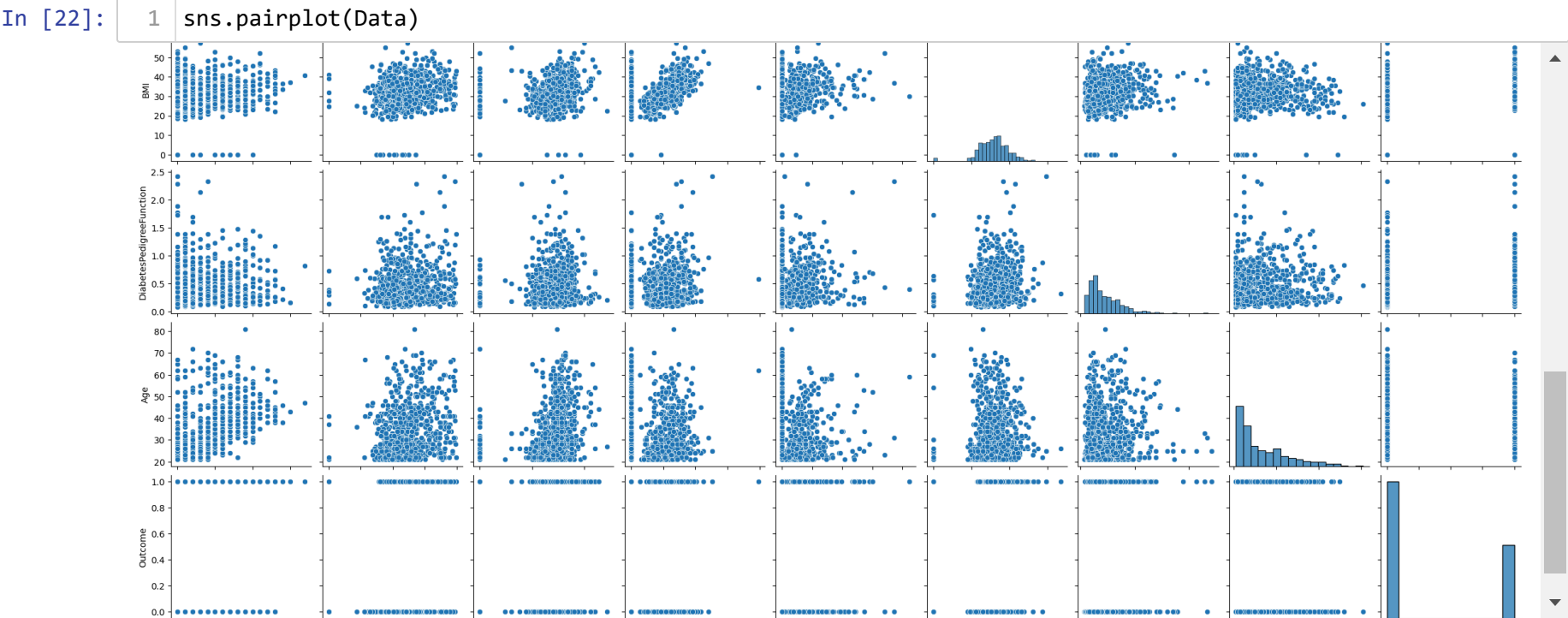
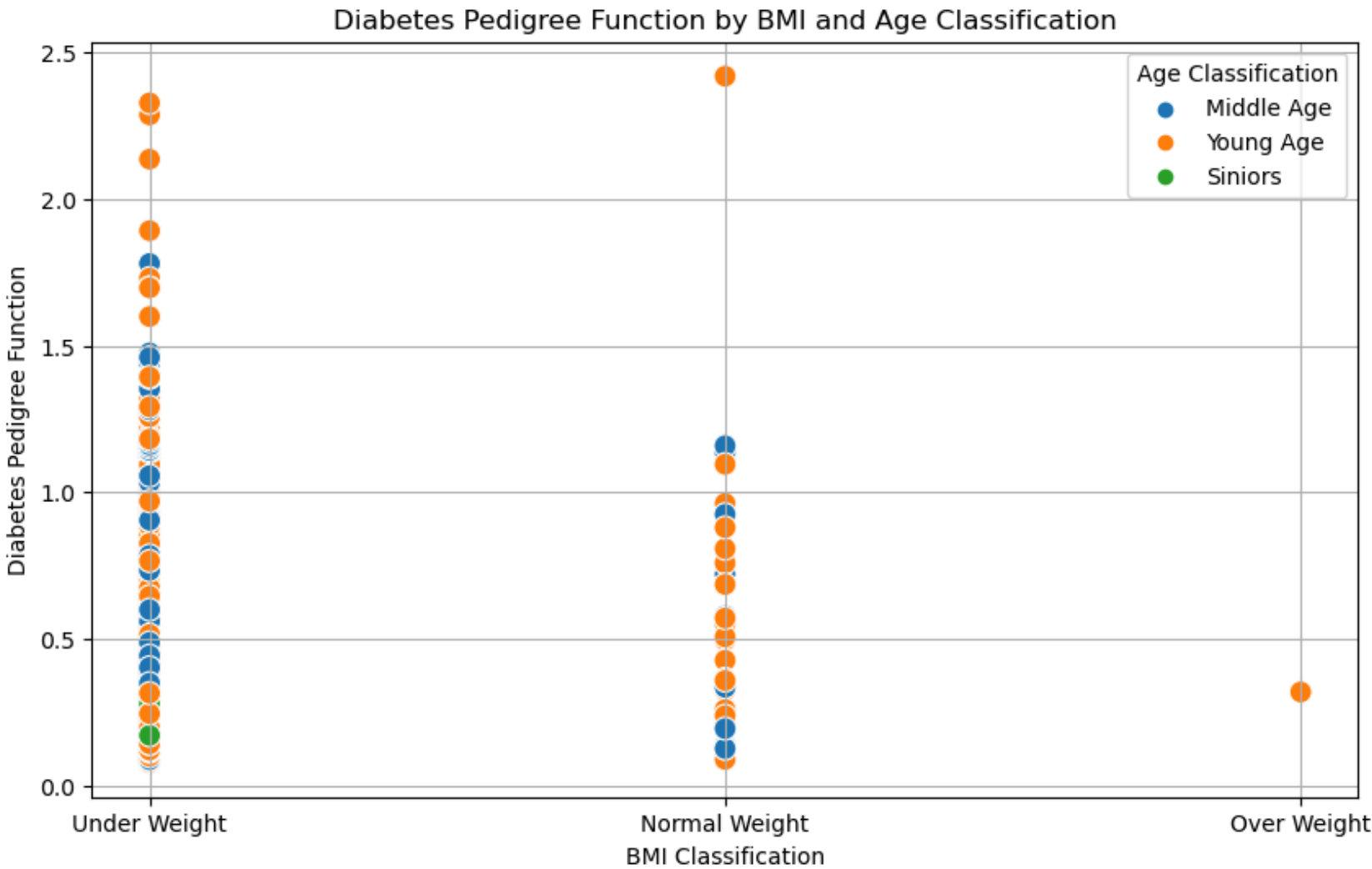
```
In [20]: 1 B = Data.groupby(['BMI Classification', 'Age Classification' , 'DiabetesPedigreeFunction']).size().unstack()
2 B
```

Out[20]:

BMI Classification	Normal Weight		Over Weight		Under Weight	
Age Classification	Middle Age	Young Age	Young Age	Middle Age	Seniors	Young Age
DiabetesPedigreeFunction						
0.078	0	0	0	0	0	1
0.084	0	0	0	0	0	1
0.085	0	0	0	1	0	1
0.088	0	0	0	1	0	1
0.089	0	1	0	0	0	0
...
1.893	0	0	0	0	0	1
2.137	0	0	0	0	0	1
2.288	0	0	0	0	0	1
2.329	0	0	0	0	0	1
2.420	0	1	0	0	0	0

517 rows × 6 columns

```
In [21]: 1 # Plotting
2 plt.figure(figsize=(10, 6))
3 sns.scatterplot(x='BMI Classification', y='DiabetesPedigreeFunction', data=Data , hue='Age Classification')
4
5 # Adding Labels and title
6 plt.xlabel('BMI Classification')
7 plt.ylabel('Diabetes Pedigree Function')
8 plt.title('Diabetes Pedigree Function by BMI and Age Classification')
9
10 # Show the plot
11 plt.legend(title='Age Classification')
12 plt.grid(True)
13 plt.show()
```



```
In [23]: 1 # Checking Data Type of all Paramerters
2
3 Data.dtypes
```

```
Out[23]: Pregnancies      int64
Glucose      int64
BloodPressure  int64
SkinThickness int64
Insulin       int64
BMI           float64
DiabetesPedigreeFunction float64
Age          int64
Outcome      int64
Age Classification object
BMI Classification object
dtype: object
```

Converting Object Data Type into Integer

In [24]:

```
1 from sklearn.preprocessing import LabelEncoder
2
3 L = LabelEncoder()
4
5 Data['Age Classification']= L.fit_transform(Data['Age Classification'])
6 Data['BMI Classification'] = L.fit_transform(Data['BMI Classification'])
7
8 Data.info()
```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 11 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
9 Age Classification 768 non-null int32
10 BMI Classification 768 non-null int32
dtypes: float64(2), int32(2), int64(7)
memory usage: 60.1 KB

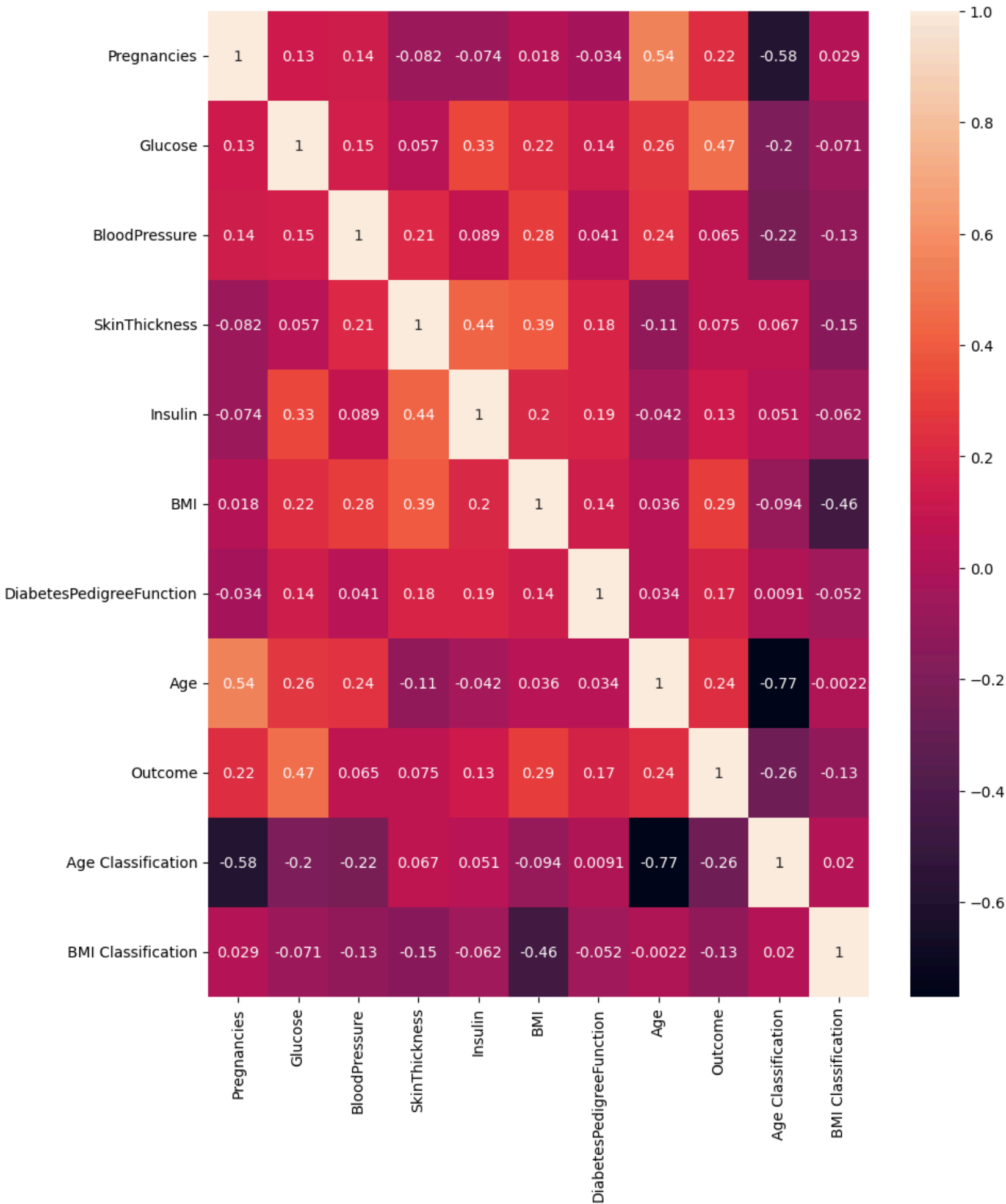
In [25]:

```
1 # Correlation between Columns
2
3 Corr = Data.corr()
```



```
In [26]: 1 plt.figure(figsize = (10 , 12))
2         sns.heatmap(Corr , annot = True)
```

Out[26]: <Axes: >



Machine Learning Models

```
In [27]: 1 # Splitting Columns into X(Independent Columns) and Y(Dependent Columns)
2
3 X = Data.drop(columns = ['Outcome'] , axis = 1)
4 Y = Data['Outcome']
```

```
In [28]: 1 # Splitting Splitting Training and Testing Data
2
3 from sklearn.model_selection import train_test_split
4
5 X_train , X_test , Y_train , Y_test = train_test_split(X , Y , test_size=0.2 , random_state=42 , stratif
```

```
In [29]: 1 Y.value_counts()
```

Out[29]: 0 500
1 268
Name: Outcome, dtype: int64

```
In [30]: 1 # # Train Test Split
2
3 print(" X\tY ".center(66))
4 print ('Train set\t:\t', X_train.shape, Y_train.shape)
5 print ('Test set\t:\t', X_test.shape, Y_test.shape)
```

		X	Y
Train set	:	(614, 10)	(614,)
Test set	:	(154, 10)	(154,)

Importing Machine Learning Algorithms

The following algorithms will be tested on the given dataset and the one with the best performance would be declared suitable:

1. Logistic Regression
2. Random Forest Classifier
3. Decision Tree Classifier
4. K-Nearest Neighbors Classifier

```
In [31]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.tree import DecisionTreeClassifier
4 from sklearn.neighbors import KNeighborsClassifier
5 from sklearn.metrics import accuracy_score , f1_score , classification_report , confusion_matrix , roc_auc_score
```

Logistic Regression

```
In [32]: 1 LR = LogisticRegression()
2
3 LR.fit(X_train , Y_train)
4
5 Y_true , Y_pred = Y_test , LR.predict(X_test)
```

C:\Users\admin\anaconda3\Lib\site-packages\sklearn\linear_model_logistic.py:458: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)
n_iter_i = _check_optimize_result(

```
In [33]: 1 print("Accuracy Score\t:\t" , accuracy_score(Y_true , Y_pred) * 100)
2 print('\n')
3 print("Classification Report\t:\n" , classification_report(Y_true , Y_pred))
```

Accuracy Score : 71.42857142857143

Classification Report :					
	precision	recall	f1-score	support	
0	0.75	0.83	0.79	100	
1	0.61	0.50	0.55	54	
accuracy			0.71	154	
macro avg	0.68	0.67	0.67	154	
weighted avg	0.71	0.71	0.71	154	

Random Forest Classifier

```
In [34]: 1 RF = RandomForestClassifier(n_estimators=100, random_state=42)
2
3 RF.fit(X_train, Y_train)
4
5 Y_pred = RF.predict(X_test)
```

```
In [35]: 1 print("Accuracy Score\t:\t" , accuracy_score(Y_true , Y_pred) * 100)
2 print('\n')
3 print("Classification Report\t:\n" , classification_report(Y_true , Y_pred))
```

Accuracy Score : 72.72727272727273

Classification Report :

	precision	recall	f1-score	support
0	0.77	0.82	0.80	100
1	0.62	0.56	0.59	54
accuracy			0.73	154
macro avg	0.70	0.69	0.69	154
weighted avg	0.72	0.73	0.72	154

Decision Tree Classifier

```
In [36]: 1 DF = DecisionTreeClassifier()
2
3 DF.fit(X_train, Y_train)
4
5 Y_pred = DF.predict(X_test)
```

```
In [37]: 1 print("Accuracy Score\t:\t" , accuracy_score(Y_true , Y_pred) * 100)
2 print('\n')
3 print("Classification Report\t:\n" , classification_report(Y_true , Y_pred))
```

Accuracy Score : 74.67532467532467

Classification Report :

	precision	recall	f1-score	support
0	0.77	0.87	0.82	100
1	0.68	0.52	0.59	54
accuracy			0.75	154
macro avg	0.73	0.69	0.70	154
weighted avg	0.74	0.75	0.74	154

K-Nearest Neighbors Classifier

```
In [38]: 1 KNN = KNeighborsClassifier(n_neighbors=5)
2
3 KNN.fit(X_train , Y_train)
4
5 Y_pred = KNN.predict(X_test)
```

```
In [39]: 1 print("Accuracy Score\t:\t" , accuracy_score(Y_true , Y_pred) * 100)
2 print('\n')
3 print("Classification Report\t:\n" , classification_report(Y_true , Y_pred))
```

Accuracy Score : 66.23376623376623

Classification Report :

	precision	recall	f1-score	support
0	0.73	0.76	0.75	100
1	0.52	0.48	0.50	54
accuracy			0.66	154
macro avg	0.63	0.62	0.62	154
weighted avg	0.66	0.66	0.66	154

Confusion Matrix

```
In [40]: 1 confusion_matrix(Y_true , Y_pred) * 100
```

Out[40]: array([[7600, 2400],
[2800, 2600]], dtype=int64)

In [41]:

1

sns.heatmap(confusion_matrix(Y_true , Y_pred) , annot = True)

2

plt.title("Confusion Matrix")

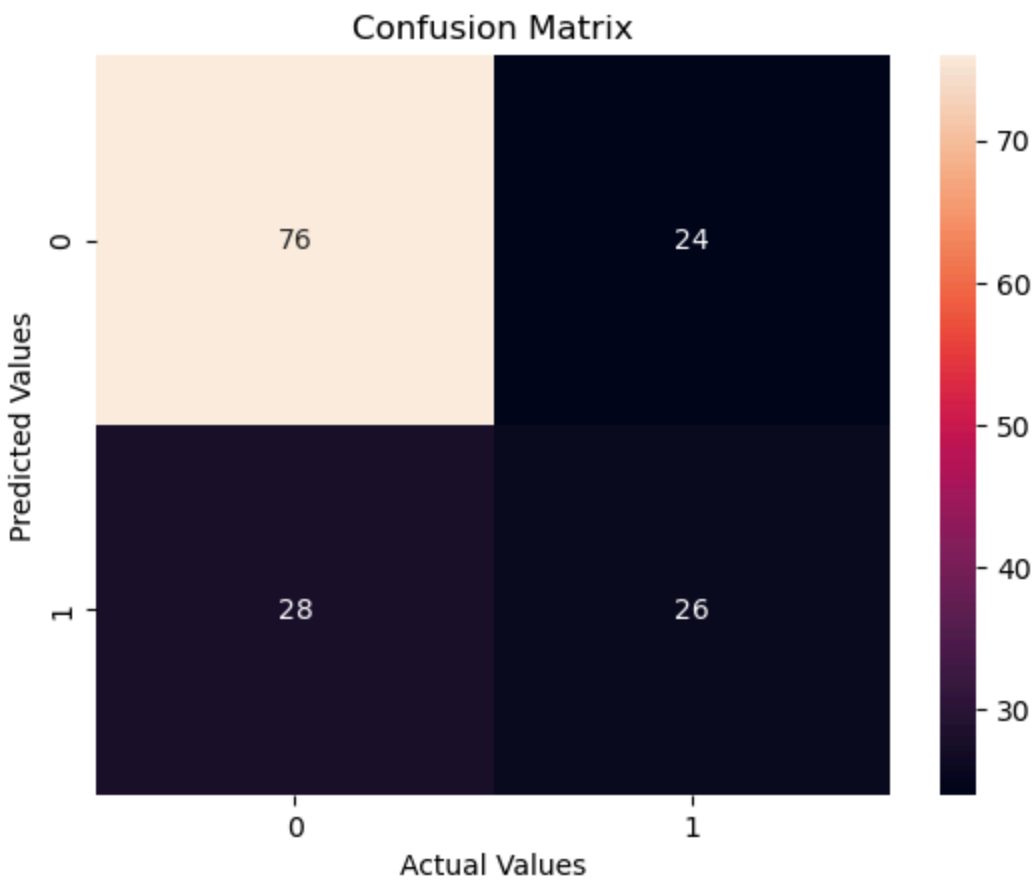
3

plt.xlabel('Actual Values')

4

plt.ylabel('Predicted Values')

Out[41]: Text(50.72222222222214, 0.5, 'Predicted Values')



In [42]:

1

tnr , tpr ,_ = roc_curve(Y_test , Y_pred)

2

plt.plot(tnr , tpr)

3

plt.title('ROC - AUC Curve')

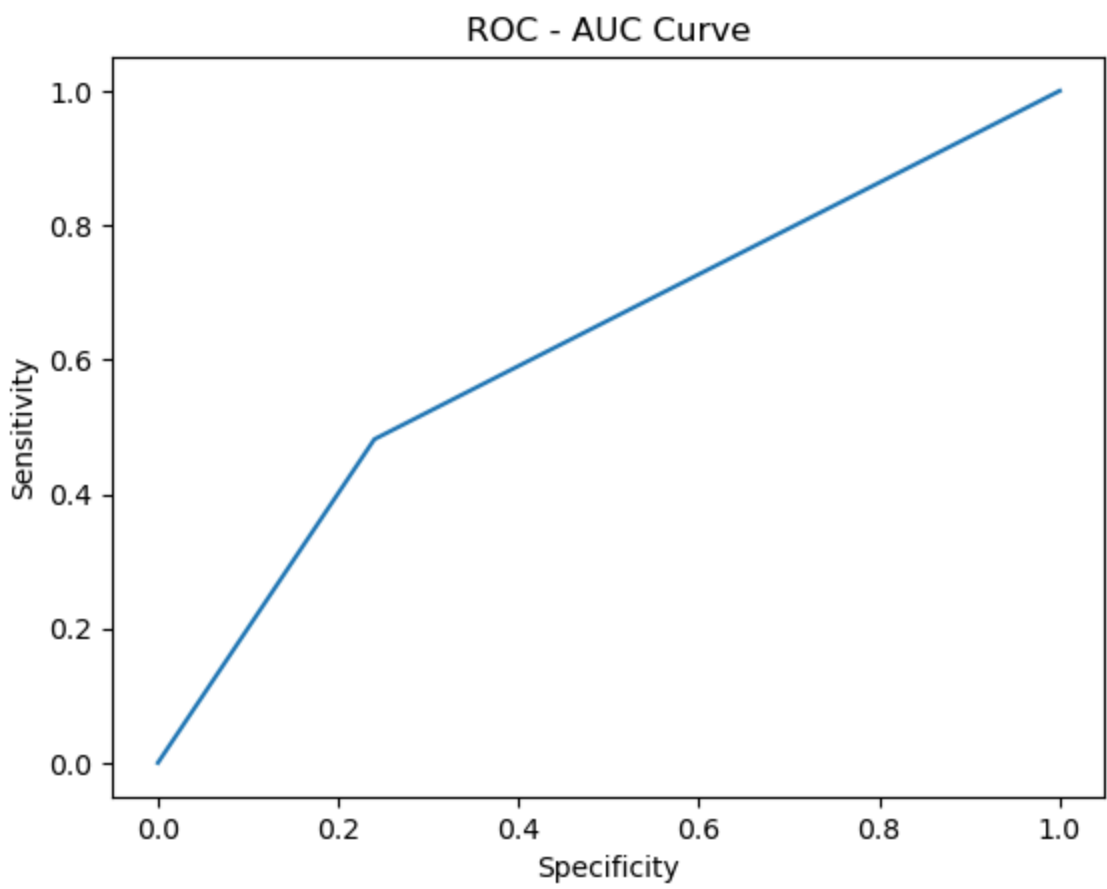
4

plt.xlabel('Specificity')

5

plt.ylabel('Sensitivity')

Out[42]: Text(0, 0.5, 'Sensitivity')



In []:

1