

# Diabetes Prediction using LogisticRegression

```
In [1]: #Let's start with importing necessary Libraries
import pandas as pd
import numpy as np
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import Ridge,Lasso,RidgeCV, LassoCV, ElasticNet, EL
from sklearn.model_selection import train_test_split
from statsmodels.stats.outliers_influence import variance_inflation_factor
from sklearn.metrics import accuracy_score, confusion_matrix, roc_curve, roc_a
import matplotlib.pyplot as plt
import seaborn as sns
sns.set()
```

```
In [2]: #read the data file
data = pd.read_csv("diabetes.csv")
data.head()
```

Out[2]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction
0	6	148	72	35	0	33.6	0.627
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	2.288

```
In [3]: data.describe()
```

Out[3]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Diabetes
count	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	
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	

In [4]: `data.describe().T`

Out[4]:

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

In [5]: `data.isnull().sum()`

Out[5]:

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

Seems like there is no missing values in our data. Great, let's see the distribution of data:

```
In [6]: import matplotlib.pyplot as plt
import seaborn as sns

# Let's see how data is distributed for every column
plt.figure(figsize=(15, 15), facecolor='white')
plotnumber = 1

colors = ['blue', 'green', 'red', 'orange', 'purple', 'brown', 'pink', 'gray',

# Loop through each column in the dataset
for i, column in enumerate(data):
    if plotnumber <= 9:
        plt.subplot(3, 3, plotnumber)
        sns.distplot(data[column], color=colors[i % len(colors)]) # Use modul
        plt.xlabel(column, fontsize=12)
        plotnumber += 1

plt.tight_layout()
plt.show()
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

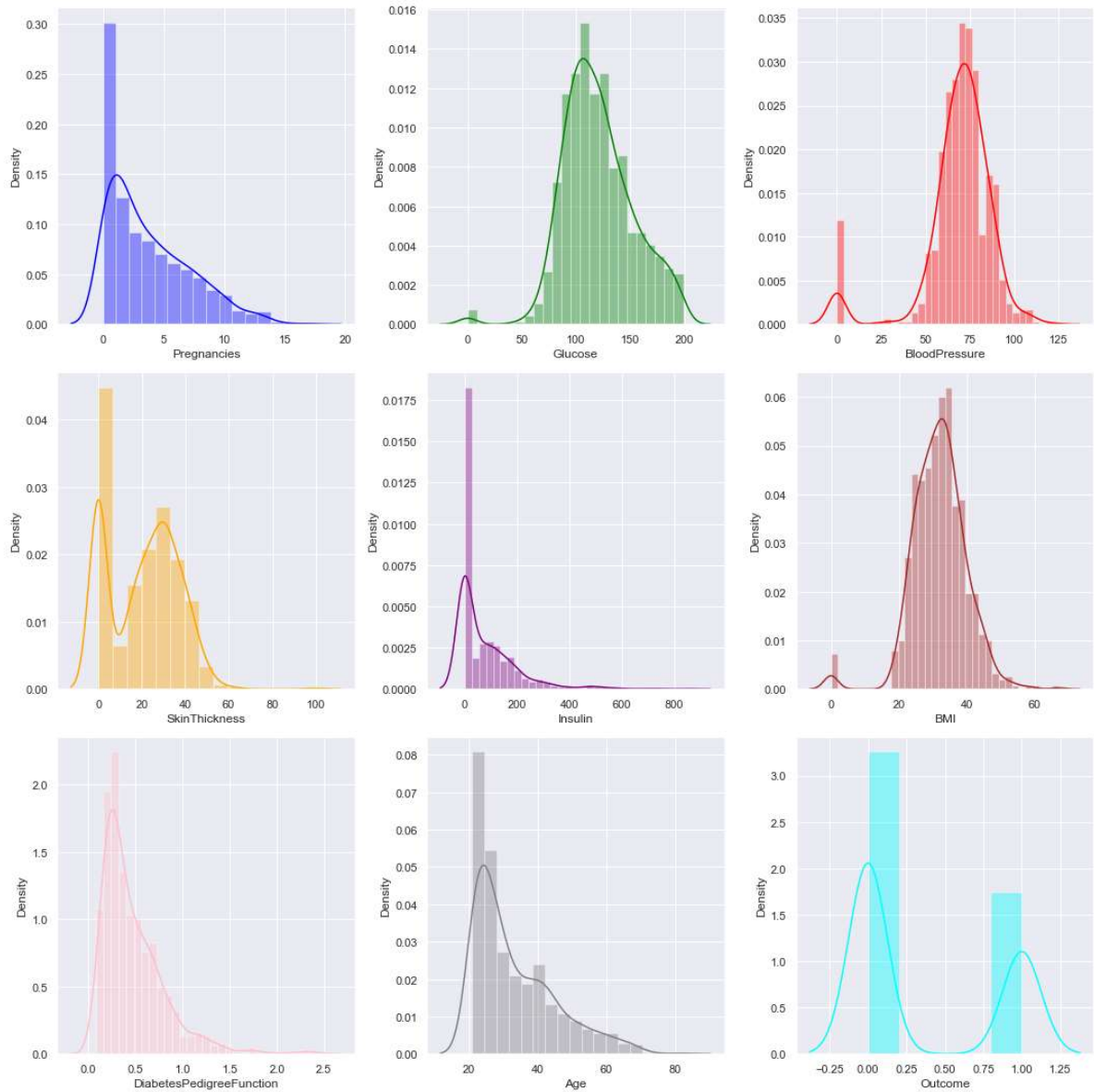
```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```



We can see there is some skewness in the data, let's deal with data.

Also, we can see there few data for columns Glucose , Insulin, skin thickenss, BMI and Blood Pressure which have value as 0. That's not possible,right? you can do a quick search to see that one cannot have 0 values for these. Let's deal with that. we can either remove such data or simply replace it with their respective mean values. Let's do the latter.

```
In [7]: #here few misconception is there lke BMI can not be zero, BP can't be zero, gl  
# now replacing zero values with the mean of the column  
data['BMI'] = data['BMI'].replace(0,data['BMI'].mean())  
data['BloodPressure'] = data['BloodPressure'].replace(0,data['BloodPressure'].  
data['Glucose'] = data['Glucose'].replace(0,data['Glucose'].mean())  
data['Insulin'] = data['Insulin'].replace(0,data['Insulin'].mean())  
data['SkinThickness'] = data['SkinThickness'].replace(0,data['SkinThickness'].  
#pregnancies data also look skewed towards left because of some outliers, let  
# q = data['Pregnancies'].quantile(0.95)  
# data_cleaned = data[data['Pregnancies']<q]
```

```
In [8]: import matplotlib.pyplot as plt
import seaborn as sns

# Let's see how data is distributed for every column
plt.figure(figsize=(15, 15), facecolor='white')
plotnumber = 1

colors = ['blue', 'green', 'red', 'orange', 'purple', 'brown', 'pink', 'gray',

# Loop through each column in the dataset
for i, column in enumerate(data):
    if plotnumber <= 9:
        plt.subplot(3, 3, plotnumber)
        sns.distplot(data[column], color=colors[i % len(colors)]) # Use modul
        plt.xlabel(column, fontsize=12)
        plotnumber += 1

plt.tight_layout()
plt.show()
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

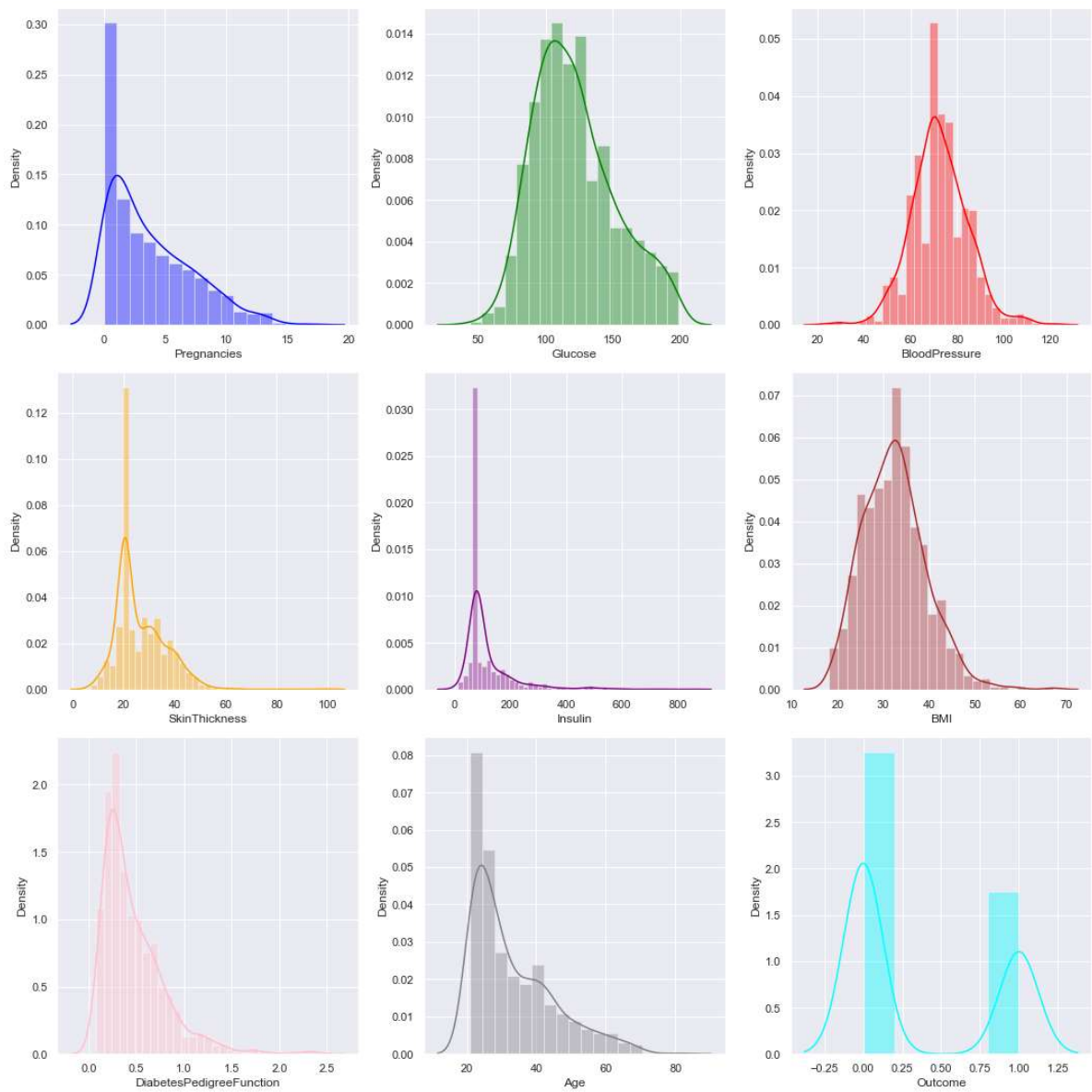
```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

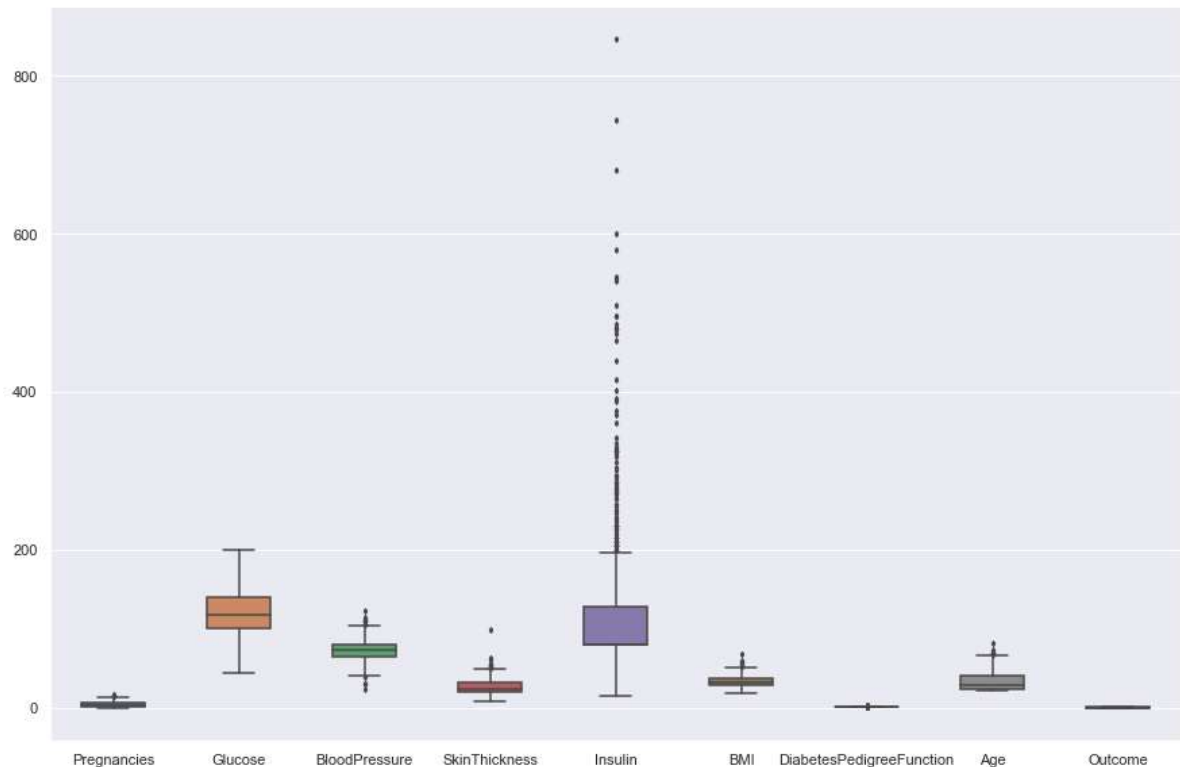
```
warnings.warn(msg, FutureWarning)
```





```
In [9]: #now we have dealt with the 0 values and data looks better. But, there still a
fig, ax = plt.subplots(figsize=(15,10))
sns.boxplot(data=data, width= 0.5,ax=ax, fliersize=3)
```

Out[9]: <AxesSubplot:>



```
In [10]: #Let's deal with the outliers.
q = data['Pregnancies'].quantile(0.98)
# we are removing the top 2% data from the Pregnancies column
data_cleaned = data[data['Pregnancies']<q]
q = data_cleaned['BMI'].quantile(0.99)
# we are removing the top 1% data from the BMI column
data_cleaned = data_cleaned[data_cleaned['BMI']<q]
q = data_cleaned['SkinThickness'].quantile(0.99)
# we are removing the top 1% data from the SkinThickness column
data_cleaned = data_cleaned[data_cleaned['SkinThickness']<q]
q = data_cleaned['Insulin'].quantile(0.95)
# we are removing the top 5% data from the Insulin column
data_cleaned = data_cleaned[data_cleaned['Insulin']<q]
q = data_cleaned['DiabetesPedigreeFunction'].quantile(0.99)
# we are removing the top 1% data from the DiabetesPedigreeFunction column
data_cleaned = data_cleaned[data_cleaned['DiabetesPedigreeFunction']<q]
q = data_cleaned['Age'].quantile(0.99)
# we are removing the top 1% data from the Age column
data_cleaned = data_cleaned[data_cleaned['Age']<q]
```

The data looks much better now than before. We will start our analysis with this data now as we don't want to lose important information. If our model doesn't work with accuracy, we will come back for more preprocessing.

```
In [19]: import matplotlib.pyplot as plt
import seaborn as sns

# Let's see how data is distributed for every column
plt.figure(figsize=(15, 15), facecolor='white')
plotnumber = 1

colors = ['blue', 'green', 'red', 'orange', 'purple', 'brown', 'pink', 'gray',

# Loop through each column in the dataset
for i, column in enumerate(data_cleaned):
    if plotnumber <= 9:
        plt.subplot(3, 3, plotnumber)
        sns.stripplot(data[column], color=colors[i % len(colors)]) # Use modu
        plt.xlabel(column, fontsize=12)
        plotnumber += 1

plt.tight_layout()
plt.show()
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
```

```
warnings.warn(
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
```

```
warnings.warn(
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
```

```
warnings.warn(
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
```

```
warnings.warn(
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
```

```
warnings.warn(
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
```

```
warnings.warn(
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
```

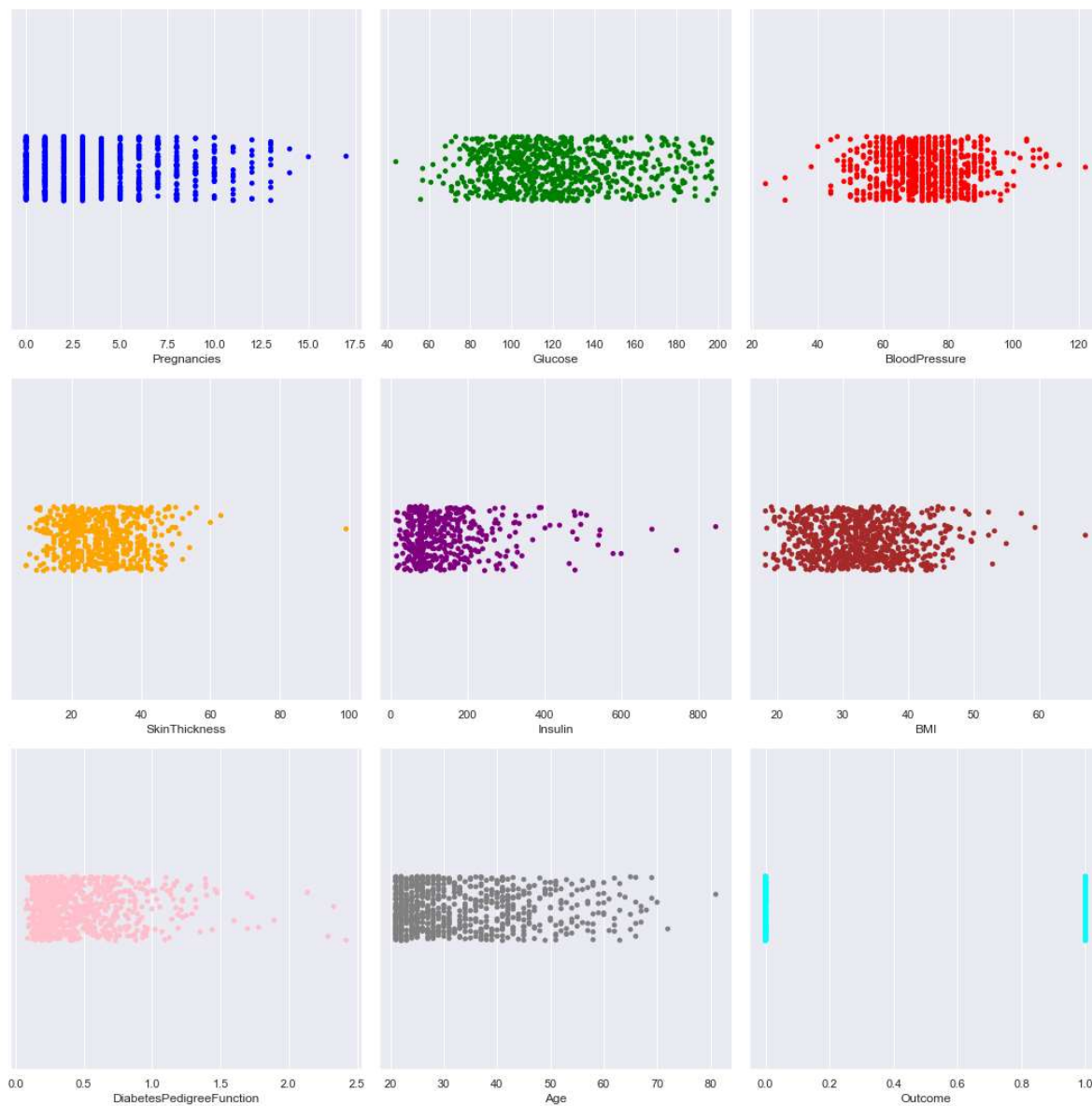
```
warnings.warn(
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
```

```
warnings.warn(
```

```
C:\Users\kishu\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
```

```
warnings.warn(
```



```
In [20]: #segregate the dependent and independent variable
X = data.drop(columns = ['Outcome'])
y = data['Outcome']
```

```
In [21]: # separate dataset into train and test
X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.25,random_
X_train.shape, X_test.shape
```

```
Out[21]: ((576, 8), (192, 8))
```

```
In [22]: import bz2,pickle
def scaler_standard(X_train, X_test):
    #scaling the data
    scaler = StandardScaler()
    X_train_scaled = scaler.fit_transform(X_train)
    X_test_scaled = scaler.transform(X_test)

    #saving the model
    file = bz2.BZ2File('standardScalar.pkl', 'wb')
    pickle.dump(scaler,file)
    file.close()

    return X_train_scaled, X_test_scaled
```

```
In [23]: X_train_scaled, X_test_scaled = scaler_standard(X_train, X_test)
```

This is how our data looks now after scaling. Great, now we will check for multicollinearity using VIF(Variance Inflation factor)

```
In [24]: X_train_scaled
```

```
Out[24]: array([[ 1.50755225, -1.09947934, -0.89942504, ..., -1.45561965,
                -0.98325882, -0.04863985],
                [-0.82986389, -0.1331471 , -1.23618124, ...,  0.09272955,
                -0.62493647, -0.88246592],
                [-1.12204091, -1.03283573,  0.61597784, ..., -0.03629955,
                 0.39884168, -0.5489355 ],
                ...,
                [ 0.04666716, -0.93287033, -0.64685789, ..., -1.14021518,
                -0.96519215, -1.04923114],
                [ 2.09190629, -1.23276654,  0.11084355, ..., -0.36604058,
                -0.5075031 ,  0.11812536],
                [ 0.33884418,  0.46664532,  0.78435594, ..., -0.09470985,
                 0.51627505,  2.953134  ]])
```

```
In [25]: log_reg = LogisticRegression()

log_reg.fit(X_train_scaled,y_train)
```

```
Out[25]: LogisticRegression()
```

```
In [26]: # r2 score
log_reg.score(X_train_scaled,y_train)
```

```
Out[26]: 0.7708333333333334
```

```
In [27]: # Let's use the handy function we created
def adj_r2(x,y,r2):
    n = x.shape[0]
    p = x.shape[1]
    adjusted_r2 = 1-(1-r2)*(n-1)/(n-p-1)
    return adjusted_r2
```

```
In [28]: # adj_r2 score

adj_r2(X_train_scaled,y_train,log_reg.score(X_train_scaled,y_train))
```

Out[28]: 0.7675999412110524

Great, our adjusted r2 score is almost same as r2 score, thus we are not being penalized for use of many features.

let's see how well our model performs on the test data set.

```
In [29]: y_pred = log_reg.predict(X_test_scaled)
```

```
accuracy = accuracy_score(y_test,y_pred) accuracy
```

```
In [30]: conf_mat = confusion_matrix(y_test,y_pred)
conf_mat
```

Out[30]: array([[117, 13],  
[ 26, 36]], dtype=int64)

```
In [31]: true_positive = conf_mat[0][0]
false_positive = conf_mat[0][1]
false_negative = conf_mat[1][0]
true_negative = conf_mat[1][1]
```

```
In [32]: Accuracy = (true_positive + true_negative) / (true_positive +false_positive +
Accuracy
```

Out[32]: 0.796875

```
In [33]: Precision = true_positive/(true_positive+false_positive)
Precision
```

Out[33]: 0.9

```
In [34]: Recall = true_positive/(true_positive+false_negative)
Recall
```

Out[34]: 0.8181818181818182

```
In [35]: F1_Score = 2*(Recall * Precision) / (Recall + Precision)
F1_Score
```

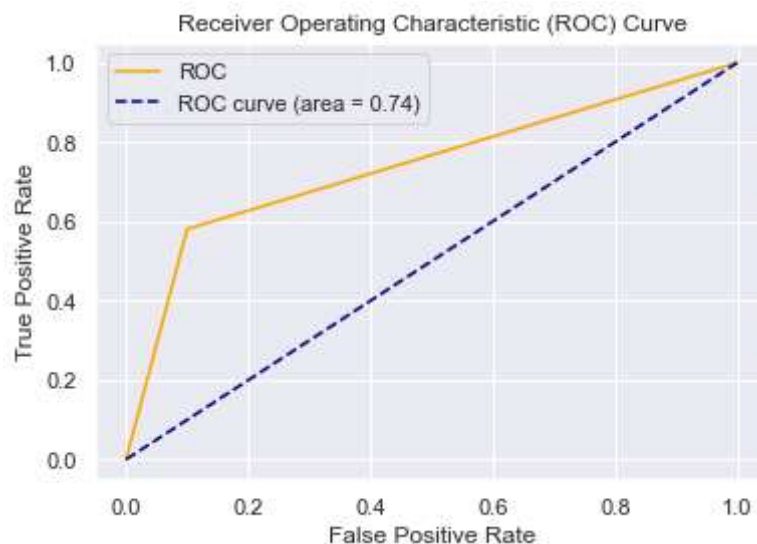
```
Out[35]: 0.8571428571428572
```

```
In [36]: auc = roc_auc_score(y_test, y_pred)
auc
```

```
Out[36]: 0.7403225806451613
```

```
In [37]: fpr, tpr, thresholds = roc_curve(y_test, y_pred)
```

```
In [38]: plt.plot(fpr, tpr, color='orange', label='ROC')
plt.plot([0, 1], [0, 1], color='darkblue', linestyle='--', label='ROC curve (ar
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC) Curve')
plt.legend()
plt.show()
```



```
In [39]: import bz2,pickle
file = bz2.BZ2File('modelForPrediction.pkl','wb')
pickle.dump(log_reg,file)
file.close()
```

```
In [ ]:
```