

Problem Statement :-

Build a model to accurately predict whether the patients in the dataset have diabetes or not?

Data Description :-

The datasets consists of several medical predictor variables and one target variable, Outcome. Predictor variables includes the number of pregnancies the patient has had, their BMI, insulin level, age, and so on.

Pregnancies: Number of times pregnant

Glucose: Plasma glucose concentration a 2 hours in an oral glucose tolerance test

BloodPressure: Diastolic blood pressure (mm Hg)

SkinThickness: Triceps skin fold thickness (mm)

Insulin: 2-Hour serum insulin (mu U/ml)

BMI: Body mass index (weight in kg/(height in m)^2)

DiabetesPedigreeFunction: Diabetes pedigree function

Age: Age (years)

Outcome: Class variable (0 or 1) 268 of 768 are 1, the others are 0

```
In [653... # importing basic libraries
import numpy as np
import pandas as pd
from matplotlib import pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")
```

Importing the data

```
In [654... data1=pd.read_csv('C:/Harsha/Capstone Projects/Project_2/Project 2/Healthcare - Diabetes/health care diabetes.csv')
```

Overview of the data

```
In [655... data1.shape
```

```
Out[655... (768, 9)
```

```
In [656... data1.columns
```

```
Out[656... Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
        'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
        dtype='object')
```

```
In [657... data1.head(2)
```

```
Out[657...
```

	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

```
In [658... data1.tail(2)
```

```
Out[658...
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
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 [659... data1.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 [660... data1.describe()

Out[660...
      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
```

```
In [661... data1['Outcome'].unique()

Out[661... array([1, 0], dtype=int64)
```

```
In [662... # Last column of dataset i.e., Outcome is categorical because it has only two
# unique items that are 0(not have diabetes) and 1(have diabetes) so I am changing int64 to category
```

```
In [663... data1['Outcome']=data1['Outcome'].astype('category')
```

```
In [664... data1.info() # now Outcome is changes to category

<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   category
dtypes: category(1), float64(2), int64(6)
memory usage: 49.0 KB
```

Checking null values in dataset

```
In [665... data1.isnull().sum()

Out[665... Pregnancies      0
Glucose          0
BloodPressure    0
SkinThickness    0
```

```
Insulin          0
BMI              0
DiabetesPedigreeFunction  0
Age              0
Outcome          0
dtype: int64
```

```
In [666... # By this we clearly known that there are no null values in our dataset
# But in some columns we have Zero values means missing values
# Why because see in Glucose,BloodPressure,SkinThickness,Insulin,BMI zero means
# As per domain knowledge Zero not makes sense so these Zero's are missing values
```

Checking Missing Values In Dataset

```
In [667... (data1['Glucose']==0).sum()
```

```
Out[667... 5
```

```
In [668... (data1['BloodPressure']==0).sum()
```

```
Out[668... 35
```

```
In [669... (data1['SkinThickness']==0).sum()
```

```
Out[669... 227
```

```
In [670... (data1['Insulin']==0).sum()
```

```
Out[670... 374
```

```
In [671... (data1['BMI']==0).sum()
```

```
Out[671... 11
```

```
In [672... missing=[5,35,227,374,11]
missing_values = pd.DataFrame(missing, index=['Glucose','BloodPressure','SkinThickness','Insulin','BMI'],columns=
```

```
In [673... missing_values
```

```
Out[673...
      Missing values
Glucose          5
BloodPressure    35
SkinThickness    227
Insulin          374
BMI              11
```

```
In [674... # By this we can see that missing values in Glucose,BloodPressure and BMI are negligible
# But in SkinThickness and Insulin we have more missing values
# I am not removing these missing values because our dataset is small dataset if we remove these
# data points data becomes not sufficient so I am filing Zero's with Mean or Median
```

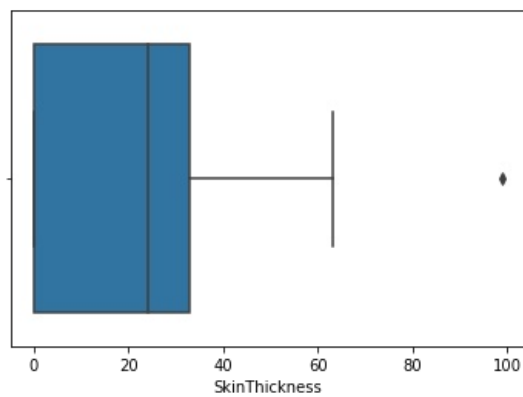
```
In [675... # And Glucose=5, BloodPressure=35 and BMI=11 these have small missing values
# so I am removing these rows with missing values
```

```
In [676... data2 = data1[(data1["Glucose"] >0) & (data1["BloodPressure"] >0) & (data1['BMI']>0)]
```

```
In [677... # First see outliers in these (SkinThickness and Insulin) two columns and decide which is best to fill
```

```
In [678... sns.boxplot(data2['SkinThickness'])
```

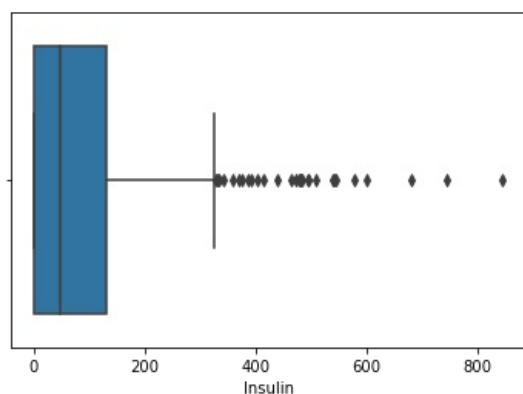
```
Out[678... <AxesSubplot:xlabel='SkinThickness'>
```



```
In [679... # It has less outliers so fill Zero's with mean
```

```
In [680... sns.boxplot(data2['Insulin'])
```

```
Out[680... <AxesSubplot:xlabel='Insulin'>
```



```
In [681... # It has more outliers so fill Zero's with Median
# Because mean is affected by outliers and Median is not affected by outliers
```

```
In [682... data2['SkinThickness'].mean()
```

```
Out[682... 21.443370165745858
```

```
In [683... data2['Insulin'].median()
```

```
Out[683... 48.0
```

```
In [684... data2['SkinThickness']=data2['SkinThickness'].replace(0,21.443370165745858)
```

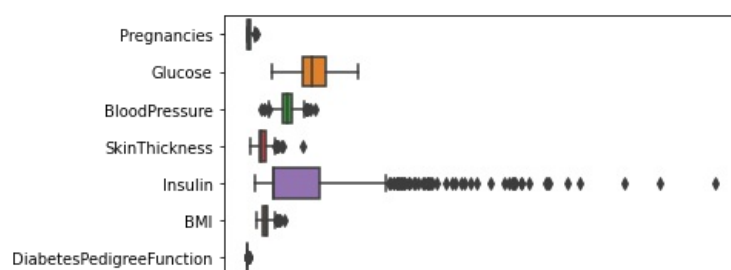
```
In [685... data2['Insulin']=data2['Insulin'].replace(0,48.0)
```

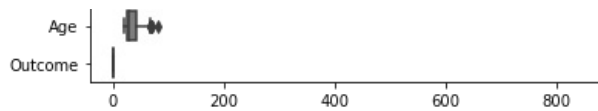
```
In [686... # Now we filled missing values
```

Detecting and Treating of outliers in dataset

```
In [687... sns.boxplot(data=data2,orient='h')
```

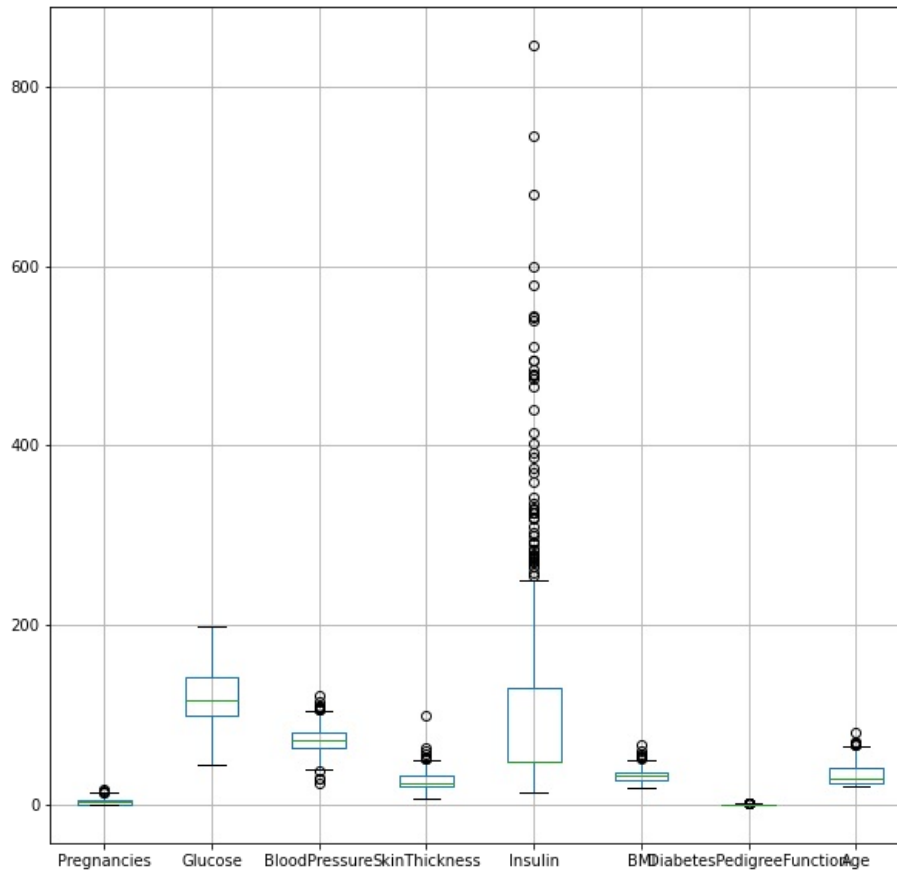
```
Out[687... <AxesSubplot:>
```





```
In [688.. data2.boxplot(figsize=(10, 10))
```

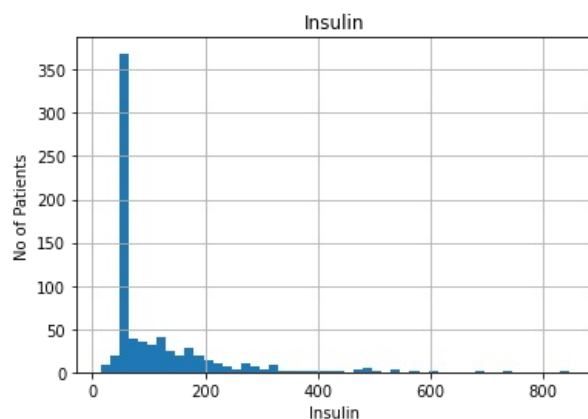
```
Out[688.. <AxesSubplot:>
```



```
In [689.. # By above boxplot Insulin have more outliers so we have to treat these outliers
# Removing outliers is not an option to our dataset because our dataset is small in size
```

```
In [690.. figure=data2['Insulin'].hist(bins=50)
figure.set_title('Insulin')
figure.set_xlabel('Insulin')
figure.set_ylabel('No of Patients')
```

```
Out[690.. Text(0, 0.5, 'No of Patients')
```



```
In [691.. # Here the Data is not Normally distributed and is following the Right Skewed distribution
# that means we can use the Interquartile Range to measure the boundaries for outliers
# IQR = Q3 - Q1 (quartile 3 - quartile 1)
```

```
In [692.. IQR= data2['Insulin'].quantile(0.75) - data2['Insulin'].quantile(0.25)
```

IQR

Out[692...] 82.5

```
In [693...] ## Calculating the boundaries  
lower_bridge= data2['Insulin'].quantile(0.25)-(IQR*1.5)  
upper_bridge= data2['Insulin'].quantile(0.75)+(IQR*1.5)  
print(lower_bridge), print(upper_bridge)
```

-75.75

254.25

Out[693...] (None, None)

```
In [694...] data2['Insulin'].describe()
```

```
Out[694...] count      724.000000  
mean       106.505525  
std        102.669035  
min         14.000000  
25%         48.000000  
50%         48.000000  
75%        130.500000  
max         846.000000  
Name: Insulin, dtype: float64
```

```
In [695...] # Here the maximum value of outliers is very high compare to upper boundary that indicates  
# we need to calculate the extreme outliers boundaries
```

```
In [696...] ## Calculating the extreme boundaries  
lower_bridge= data2['Insulin'].quantile(0.25)-(IQR*3)  
upper_bridge= data2['Insulin'].quantile(0.75)+(IQR*3)  
print(lower_bridge), print(upper_bridge)
```

-199.5

378.0

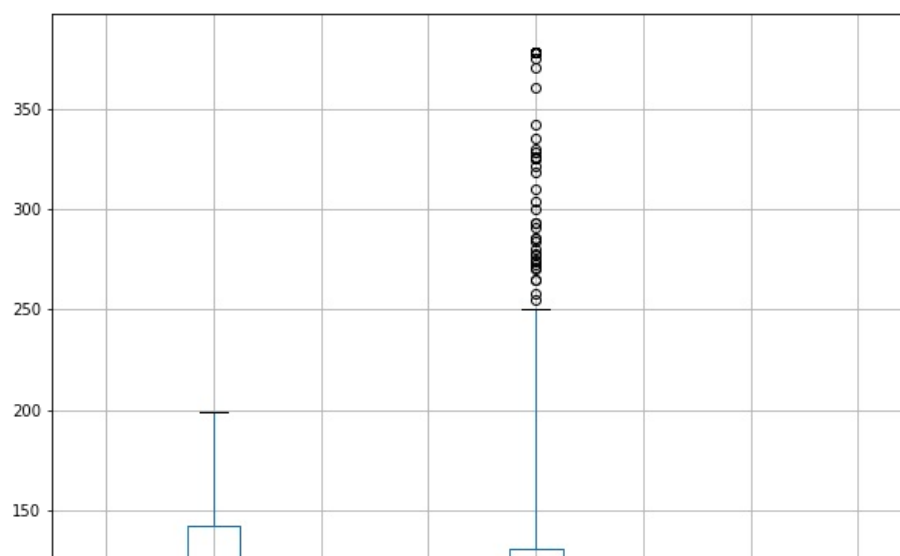
Out[696...] (None, None)

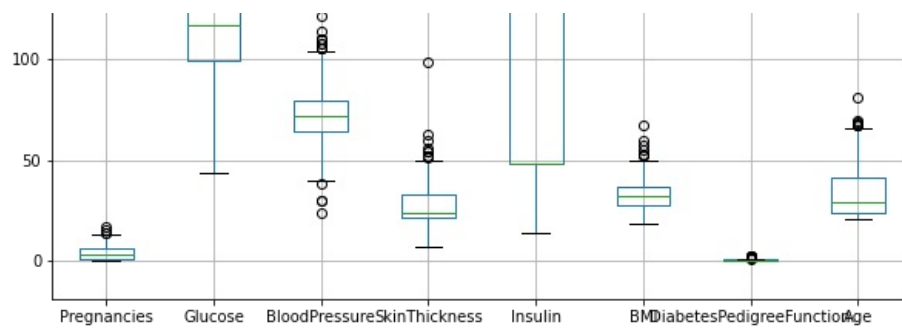
```
In [697...] # Replacing the outliers  
# Since the lower boundaries of Insulin column is negative value  
# we do not need to consider the lower boundary because as per the domain knowledge  
# there won't be any negative values exists for Insulin Column  
# Insulin upper bridge is 378
```

```
In [698...] data2.loc[data2['Insulin']>=378,'Insulin']=378
```

```
In [699...] data2.boxplot(figsize=(10, 10))
```

Out[699...] <AxesSubplot:>





In [700...] *# By comparing above boxplot outliers are only reduced not completely gone*

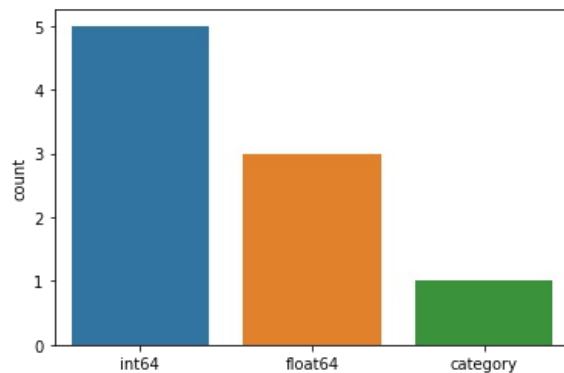
In [897...] `data2.to_csv('C:/Harsha/Capstone Projects/Project_2/Project 2/Healthcare - Diabetes/data.csv')`

Create a count (frequency) plot describing the data types and the count of variables.

In [701...] `data2.dtypes.value_counts()`

Out[701...] `int64 5`
`float64 3`
`category 1`
`dtype: int64`

In [702...] `sns.countplot(data2.dtypes.map(str))`
`plt.show()`

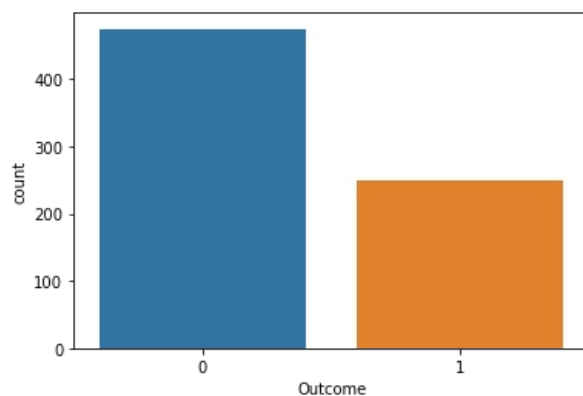


Check the balance of the data by plotting the count of outcomes by their value. Describe your findings and plan future course of actions.

In [703...] `data2['Outcome'].value_counts()`

Out[703...] `0 475`
`1 249`
`Name: Outcome, dtype: int64`

In [704...] `sns.countplot(data2['Outcome'].map(str))`
`plt.show()`



```
In [705... # By above we have clear evidence that our dataset has class imbalance issue
# There are different methods to solve class imbalance issue
# 1.Over Sampling Technique
# 2.Under sampling Technique
# 3.Cross validation Technique
# But see in both under and over sampling we are adding and deleting data points.
# High Accuracy will come but we cannot assure that that accuracy is biased or not
# But with cross validation in this technique we train our model using the subset of the data-set
# so we can assure that accuracy got by this technique is unbiased.
```

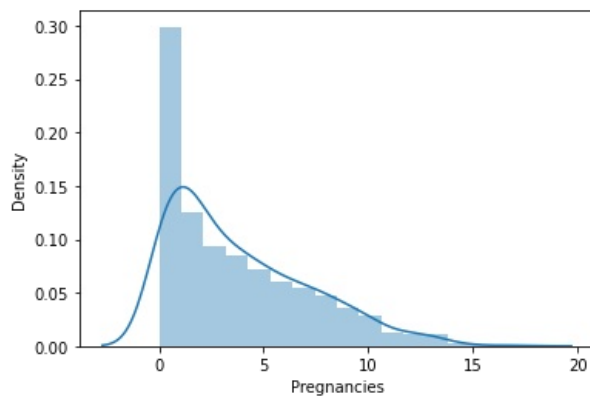
Checking variables normal or not if not normal making to normal

```
In [706... data2.columns
```

```
Out[706... Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
        'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
        dtype='object')
```

```
In [707... sns.distplot(data2['Pregnancies'])
```

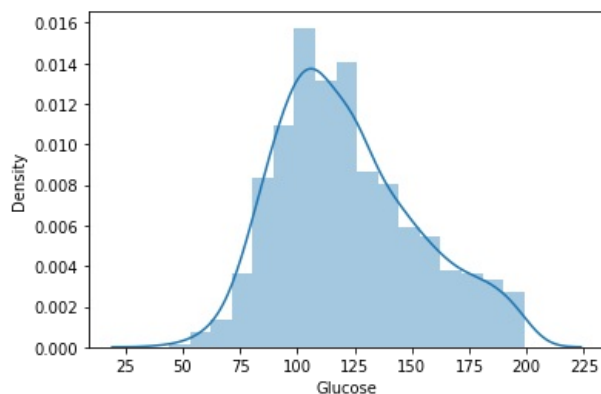
```
Out[707... <AxesSubplot:xlabel='Pregnancies', ylabel='Density'>
```



```
In [708... # it is not normal it is right skewed
```

```
In [709... sns.distplot(data2['Glucose'])
```

```
Out[709... <AxesSubplot:xlabel='Glucose', ylabel='Density'>
```



```
In [710... # it is visually normal but slight skewed so we can declare it is skewed or not by summary statistics
```

```
In [711... data2['Glucose'].describe()
```

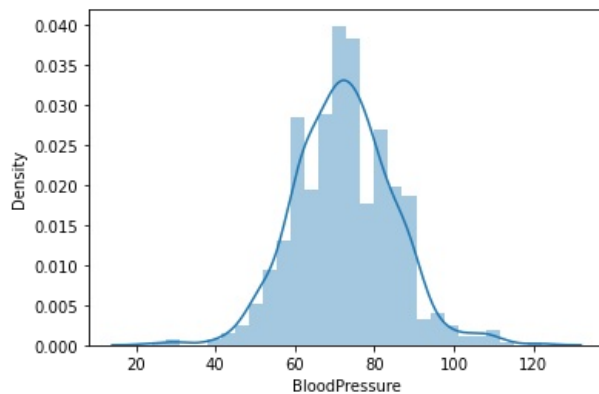
```
Out[711... count    724.000000
mean      121.882597
std       30.750030
min       44.000000
25%      99.750000
50%     117.000000
75%     142.000000
max     199.000000
```


Name: Glucose, dtype: float64

```
In [712...] # Here there is 4 difference between mean and median so it is not normal
```

```
In [713...] sns.distplot(data2['BloodPressure'])
```

```
Out[713...] <AxesSubplot:xlabel='BloodPressure', ylabel='Density'>
```



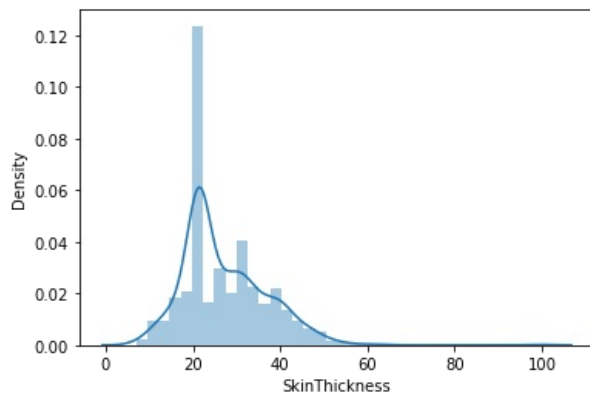
```
In [714...] data2['BloodPressure'].describe()
```

```
Out[714...] count      724.000000
mean        72.400552
std         12.379870
min         24.000000
25%         64.000000
50%         72.000000
75%         80.000000
max         122.000000
Name: BloodPressure, dtype: float64
```

```
In [715...] # it is normal
```

```
In [716...] sns.distplot(data2['SkinThickness'])
```

```
Out[716...] <AxesSubplot:xlabel='SkinThickness', ylabel='Density'>
```



```
In [717...] data2['SkinThickness'].describe()
```

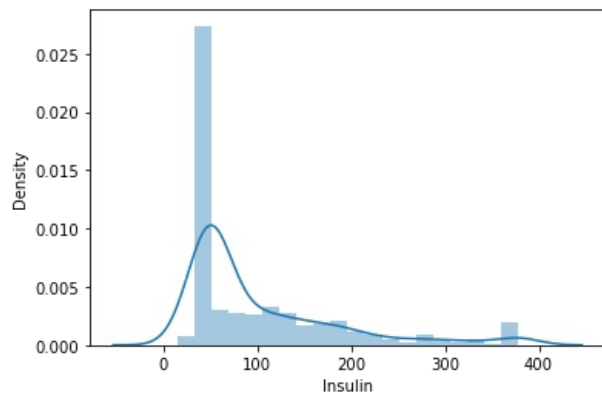
```
Out[717...] count      724.000000
mean        27.130010
std          9.645083
min          7.000000
25%         21.443370
50%         24.000000
75%         33.000000
max         99.000000
Name: SkinThickness, dtype: float64
```

```
In [718...] # Not normal right skewed
```

```
In [718... # Not normal and right skewed
```

```
In [719... sns.distplot(data2['Insulin'])
```

```
Out[719... <AxesSubplot:xlabel='Insulin', ylabel='Density'>
```



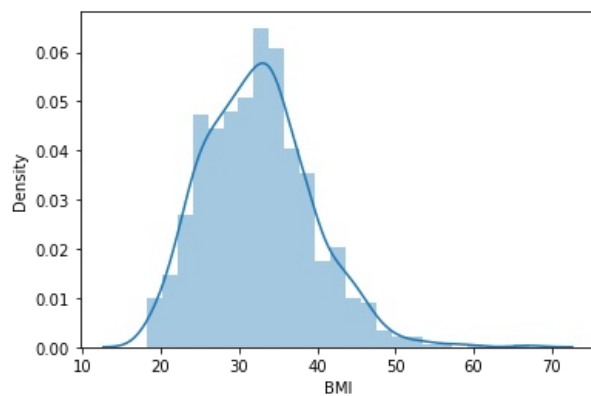
```
In [720... data2['Insulin'].describe()
```

```
Out[720... count      724.000000
mean       102.142265
std        84.536572
min        14.000000
25%        48.000000
50%        48.000000
75%       130.500000
max       378.000000
Name: Insulin, dtype: float64
```

```
In [721... # Not normal and right skewed
```

```
In [722... sns.distplot(data2['BMI'])
```

```
Out[722... <AxesSubplot:xlabel='BMI', ylabel='Density'>
```



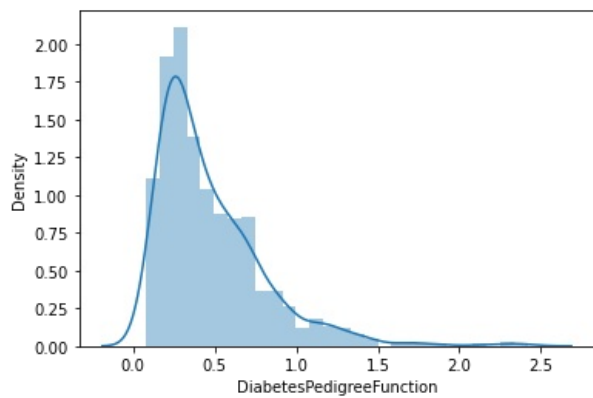
```
In [723... data2['BMI'].describe()
```

```
Out[723... count      724.000000
mean       32.467127
std        6.888941
min       18.200000
25%       27.500000
50%       32.400000
75%       36.600000
max       67.100000
Name: BMI, dtype: float64
```

```
In [724... # It is normal distribution
```

```
In [725... sns.distplot(data2['DiabetesPedigreeFunction'])
```

```
Out[725... <AxesSubplot:xlabel='DiabetesPedigreeFunction', ylabel='Density'>
```



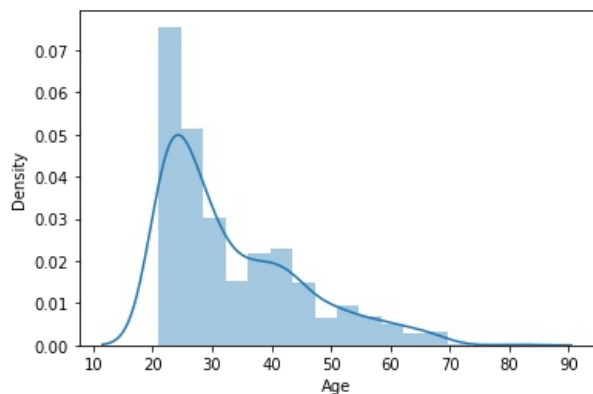
```
In [726...] data2['DiabetesPedigreeFunction'].describe()
```

```
Out[726...] count      724.000000
mean         0.474765
std          0.332315
min          0.078000
25%          0.245000
50%          0.379000
75%          0.627500
max          2.420000
Name: DiabetesPedigreeFunction, dtype: float64
```

```
In [727...] # It is normal difference between mean and median is negligble
```

```
In [728...] sns.distplot(data2['Age'])
```

```
Out[728...] <AxesSubplot:xlabel='Age', ylabel='Density'>
```



```
In [729...] data2['Age'].describe()
```

```
Out[729...] count      724.000000
mean         33.350829
std          11.765393
min          21.000000
25%          24.000000
50%          29.000000
75%          41.000000
max          81.000000
Name: Age, dtype: float64
```

```
In [730...] # It is not normal
```

```
In [731...] # So here Pregnancies,Glucose,SkinThickness,Insulin,Age are not normal we have to make it to normal
```

```
In [732...] # using log transformation technique to make variables normal
data2['Pregnancies']=np.log1p(data2['Pregnancies'])
data2['Glucose']=np.log1p(data2['Glucose'])
data2['SkinThickness']=np.log1p(data2['SkinThickness'])
data2['Insulin']=np.log1p(data2['Insulin'])
data2['Age']=np.log1p(data2['Age'])
```

In [733]	data2.describe()								
Out [733]		Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
	count	724.000000	724.000000	724.000000	724.000000	724.000000	724.000000	724.000000	724.000000
	mean	1.319311	4.780277	72.400552	3.280856	4.387084	32.467127	0.474765	3.484730
	std	0.762929	0.249781	12.379870	0.336677	0.662871	6.888941	0.332315	0.313971
	min	0.000000	3.806662	24.000000	2.079442	2.708050	18.200000	0.078000	3.091042
	25%	0.693147	4.612633	64.000000	3.110995	3.891820	27.500000	0.245000	3.218876
	50%	1.386294	4.770685	72.000000	3.218876	3.891820	32.400000	0.379000	3.401197
	75%	1.945910	4.962845	80.000000	3.526361	4.878985	36.600000	0.627500	3.737670
	max	2.890372	5.298317	122.000000	4.605170	5.937536	67.100000	2.420000	4.406719

In [734]

By above summary statistics for all variables mean is approximately equals to median so all
variables are normal now

Checking correlation between variables

First Doing Scaling to all dataset

In [735]	<pre>from sklearn.preprocessing import MinMaxScaler,StandardScaler m=MinMaxScaler()</pre>								
In [736]	# Here x is independent variable and y is dependent variable (Outcome)								
In [737]	<pre>x=data2.drop(['Outcome'],axis=1) y=data2['Outcome']</pre>								
In [738]	<pre>from sklearn.preprocessing import MinMaxScaler,StandardScaler m=MinMaxScaler()</pre>								
In [739]	data_sc=m.fit_transform(x)								
In [740]	data_sc_df=pd.DataFrame(data_sc,columns=x.columns,index=x.index)								
In [741]	data_sc_df.head(5)								
Out [741]		Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
	0	0.673239	0.802655	0.489796	0.595502	0.366551	0.314928	0.234415	0.639050
	1	0.239812	0.434206	0.428571	0.523317	0.366551	0.171779	0.116567	0.284791
	2	0.760188	0.944101	0.408163	0.408418	0.366551	0.104294	0.253629	0.308180
	3	0.239812	0.464683	0.428571	0.434968	0.571554	0.202454	0.038002	0.000000
	4	0.000000	0.751240	0.163265	0.595502	0.749918	0.509202	0.943638	0.330870

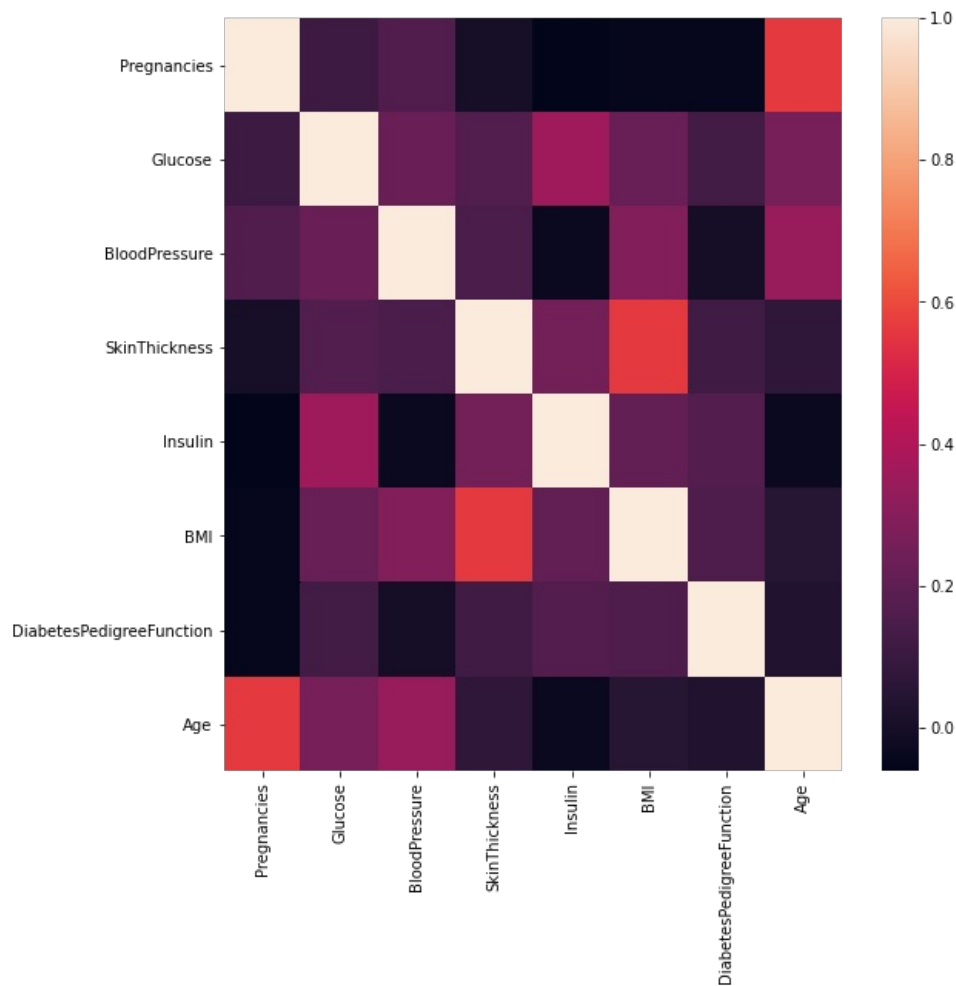
Finding correlation between variables

In [742]	data_sc_df.corr()								
Out [742]		Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
	Pregnancies	1.000000	0.110236	0.159556	0.000653	-0.061777	-0.046226	-0.046698	0.560507
	Glucose	0.110236	1.000000	0.226538	0.164075	0.354076	0.220961	0.125562	0.257609
	BloodPressure	0.159556	0.226538	1.000000	0.147366	-0.035981	0.287403	-0.000075	0.340852
	SkinThickness	0.000653	0.164075	0.147366	1.000000	0.249853	0.561603	0.119397	0.069261
	Insulin	-0.061777	0.354076	-0.035981	0.249853	1.000000	0.205610	0.167646	-0.035048
	BMI	-0.046226	0.220961	0.287403	0.561603	0.205610	1.000000	0.154858	0.049484
	DiabetesPedigreeFunction	-0.046698	0.125562	-0.000075	0.119397	0.167646	0.154858	1.000000	0.032301
	Age	0.560507	0.257609	0.340852	0.069261	-0.035048	0.049484	0.032301	1.000000

In [743]

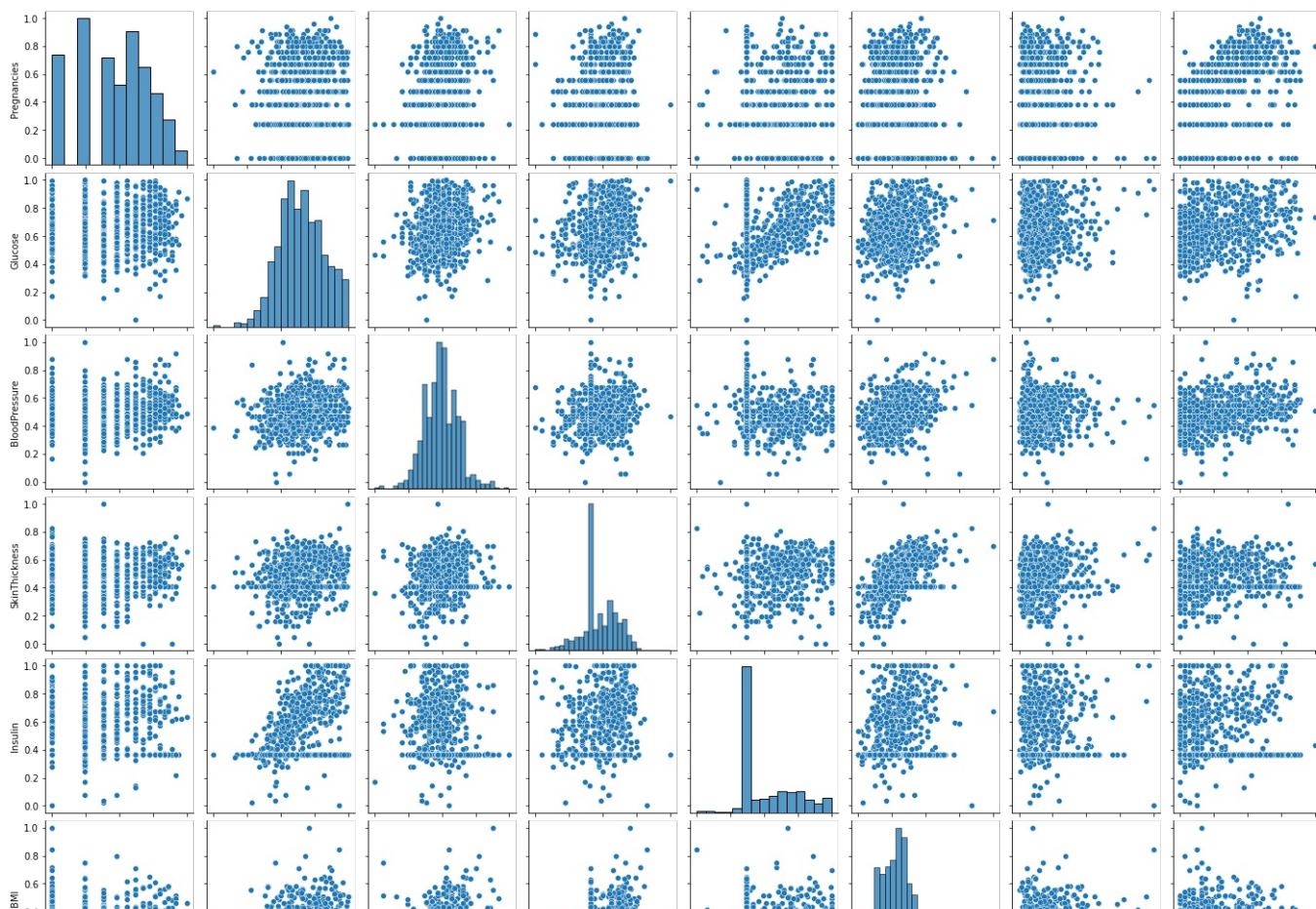
plt.figure(figsize = (9,9))
sns.heatmap(data_sc_df.corr())

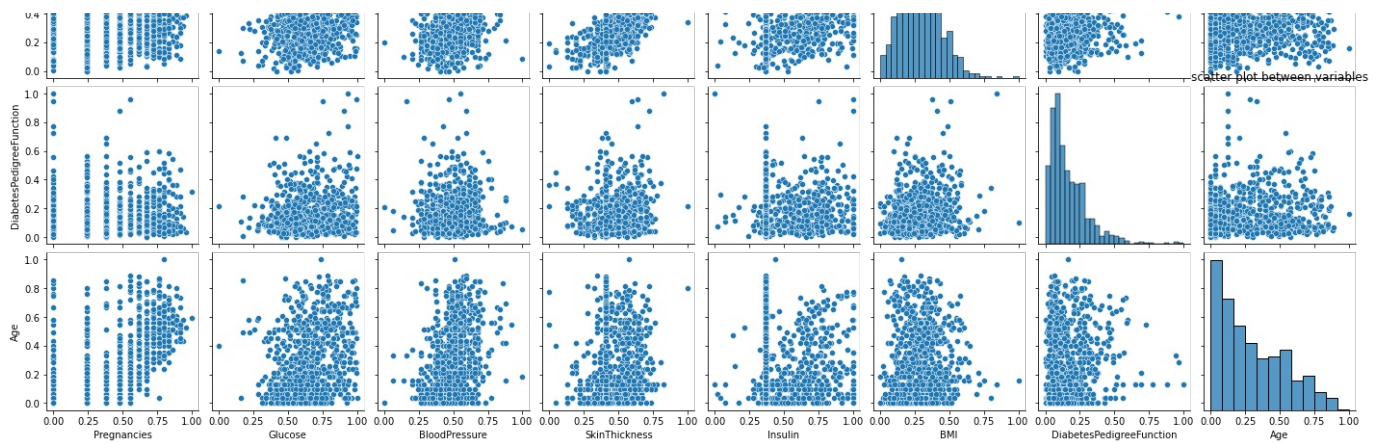
Out [743... <AxesSubplot:>



```
In [899... sns.pairplot(data_sc_df)
plt.title('scatter plot between variables')
```

Out [899... Text(0.5, 1.0, 'scatter plot between variables')





```
In [900... # Clearly multi-collinearity exists in variables in our data set
# We can see from scatter plot that there is no strong multicollinearity among features
# But between skin thickness and BMI, Pregnancies and age it looks like there positive correlation.
# In heat map we got cleared there is multi-collinearity
```

```
In [744... # For Logistic Regression multi-collinearity should not be in data
# But here there is clearly multi-collinearity exists between variables
# Multi-collinearity means presence of relation between independent variables
# Multi-collinearity is a problem we should avoid
# There are two types o avoid multi-collinearity
# 1.) Remove independent variables which are corelated
# 2.) Use dimensionality reduction technique to merge variables
```

```
In [745... # Variance Inflation Factor (VIF)
# equal to the ratio of the overall model variance to the variance of a model
# that includes only that single independent variable.
# VIF=1/1-R square (R-square is correlation coefficient of that variable)
```

```
In [746... # VIF ~ 1: Negligible
# 1<VIF<5 : Moderate
# VIF>5 : Extreme
# if a variable with extreme VIF is should be removed
```

```
In [747... from statsmodels.stats.outliers_influence import variance_inflation_factor
```

```
In [748... # VIF dataframe
vif_data = pd.DataFrame()
vif_data["feature"] = data_sc_df.columns
```

```
In [749... # calculating VIF for each feature
vif_data["VIF"] = [variance_inflation_factor(data_sc_df.values, i)
                    for i in range(len(data_sc_df.columns))]
```

```
In [750... print(vif_data)
```

	feature	VIF
0	Pregnancies	5.512882
1	Glucose	18.220714
2	BloodPressure	15.132647
3	SkinThickness	17.090657
4	Insulin	9.077077
5	BMI	8.266541
6	DiabetesPedigreeFunction	2.534800
7	Age	4.104236

```
In [751... # Here five variables have high VIF value so we have to remove these variables
# That leads heavy data loss in our scenario so by using PCA echnique.
# we are going to merge there all variables.
```

PCA Technique

```
In [752... # using PCA technique
from sklearn.decomposition import PCA
```

```
In [753... pca=PCA(n_components=6)
var_x=pca.fit_transform(data_sc_df)
```

```
In [754... new_data=pd.DataFrame(var_x,columns=['pca1','pca2','pca3','pca4','pca5','pca6'])
```

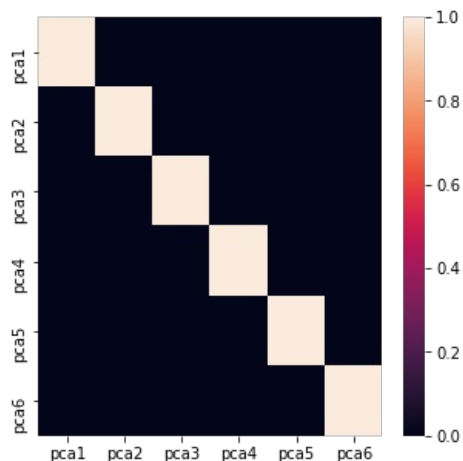
```
In [755... b=pca.explained_variance_ratio_
b*100

Out[755... array([37.81312702, 20.79637505, 11.53703291,  8.99408893,  7.07498587,
        6.66814181])
```

```
In [756... # Here 88% data captured by these pca variables
```

```
In [757... plt.figure(figsize = (5,5))
sns.heatmap(new_data.corr())
```

```
Out[757... <AxesSubplot:>
```



```
In [758... # so finally no multicollinearity so we can proceed with logistic regression
```

Logistic Regression Model Building

```
In [759... # so now
x_var=new_data
y_var=y
```

```
In [760... # splitting our data
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score,classification_report
```

```
In [761... x_train,x_test,y_train,y_test=train_test_split(x_var,y_var,test_size=0.2,random_state=50)
```

```
In [762... from sklearn.linear_model import LogisticRegression
lr=LogisticRegression()
lr.fit(x_train,y_train)
```

```
Out[762... LogisticRegression()
```

```
In [763... pred=lr.predict(x_test)
```

```
In [764... accuracy_score(y_test,pred)
```

```
Out[764... 0.7793103448275862
```

```
In [765... accuracy_score(y_train,lr.predict(x_train))
```

```
Out[765... 0.7651122625215889
```

```
In [766... # see clearly our model is not over fitted so see for other reports
```

```
In [767... print(classification_report(y_test,pred))
```

```
precision    recall  f1-score   support
```


0	0.76	0.94	0.84	89
1	0.85	0.52	0.64	56

accuracy			0.78	145
macro avg	0.80	0.73	0.74	145
weighted avg	0.79	0.78	0.76	145

```
In [893]_ #1.Accuracy is over all 78% percent it is good because our data predicts 78% correctly
#2.But here our outcome is having diabetes (1) and not having diabetes (0) means both 1 and 0 are
# important here so we have to consider macro average because it gives equal importance to both
# 1 and 0 so our macro average is 74%
```

AUC and ROC curve

AUC Curve

```
In [768]_ # The Area Under the Curve (AUC)
# Is the measure of the ability of a classifier to distinguish between classes and is used as a summary of the ROC
# The higher the AUC, the better the performance of the model at distinguishing between the positive and negative
```

ROC curve

```
In [769]_ # ROC (Reciever Operating Characteristic Curve)
# In ROC curve we plot pairs of the true positive (Sensitivity) rate vs the false positive rate (Specificity).
# for every possible decision threshold of a logistic regression model.
# By ROC curve we can find optimal Threshold value for our model.
```

```
In [770]_ # plot roc and auc curve
from sklearn.metrics import roc_auc_score,roc_curve
```

```
In [771]_ # Here first find predicted probabilities of test data
# Then created dataframe of actual and predicted probabilities of data
```

```
In [772]_ y_prob=lr.predict_proba(x_test)
df_pred_prob=pd.DataFrame({'Actual':y_test,'Predicted_prob':y_prob[:,1]})
```

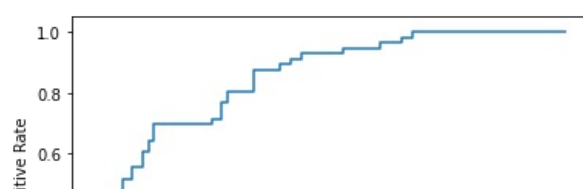
```
In [773]_ df_pred_prob
```

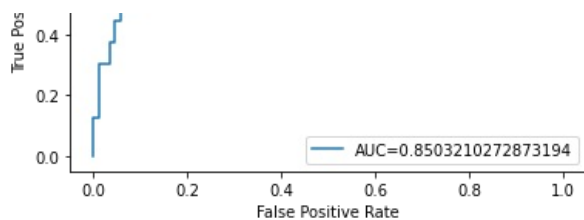
```
Out[773]_
```

	Actual	Predicted_prob
19	1	0.274795
750	1	0.396869
213	1	0.318996
424	1	0.713071
501	0	0.144898
...
647	1	0.405455
108	0	0.088859
730	1	0.315705
761	1	0.726563
45	1	0.704235

145 rows × 2 columns

```
In [774]_ y_prob[:,1][0:10]
fpr, tpr, thresholds=roc_curve(y_test,y_prob[:,1])
auc_curve=roc_auc_score(y_test,y_prob[:,1])
plt.plot(fpr,tpr,label="AUC="+str(auc_curve))
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.legend(loc=4)
plt.show()
```





```
In [775... # Here our AUC is 85% means our model distinguished 85% correctly
```

```
In [776... from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
confusion_matrix(y_test, pred)
```

```
Out[776... array([[84,  5],
        [27, 29]], dtype=int64)
```

```
In [777... # BY this confusion matrix we can say that actually 5 patients have no diabetes
# but our model says have diabetes
# Actually 27 patients have diabetes but our predicted model says
# not have diabetes this is false negative
# here our false negative=27 we have to reduce our False Negative
# we have to find threshold value to reduce our error(false negative) by ROC curve we can know that
```

```
In [778... # put threshold at 0.4 means probabilities more than 0.4 are 1 (have diabetes) rest are 0 (not have diabetes)
```

```
In [779... df_pred_prob['prediction_0.4']=np.where(df_pred_prob['Predicted_prob']>0.4,1,0)
confusion_matrix(df_pred_prob['Actual'],df_pred_prob['prediction_0.4'])
```

```
Out[779... array([[78, 11],
        [18, 38]], dtype=int64)
```

```
In [780... # Here our false negative is reduced I think it's fair so I am fixing threshold value at 0.4
# By domain knowledge if we want min (least) false negative we can reduce threshold value below 0.4
```

```
In [781... # put threshold at 0.2
```

```
In [782... df_pred_prob['prediction_0.2']=np.where(df_pred_prob['Predicted_prob']>0.2,1,0)
confusion_matrix(df_pred_prob['Actual'],df_pred_prob['prediction_0.2'])
```

```
Out[782... array([[46, 43],
        [ 4, 52]], dtype=int64)
```

```
In [783... # Here see false negative(FN) is 4 means FN is reduced but see False positive(FP) is increased
```

Other Classification Models Except Logistic Model

```
In [784... # For Tree models preprocessing steps are not necessary but we will use above data
```

Decision Trees

```
In [785... # Here I am using Scaling Data not PCA data
x_dtree=data_sc_df
y_dtree=y
```

```
In [786... x_dtree.head(5)
```

```
Out[786...
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
0	0.673239	0.802655	0.489796	0.595502	0.366551	0.314928	0.234415	0.639050
1	0.239812	0.434206	0.428571	0.523317	0.366551	0.171779	0.116567	0.284791
2	0.760188	0.944101	0.408163	0.408418	0.366551	0.104294	0.253629	0.308180
3	0.239812	0.464683	0.428571	0.434968	0.571554	0.202454	0.038002	0.000000
4	0.000000	0.751240	0.163265	0.595502	0.749918	0.509202	0.943638	0.330870

```
In [787... y_dtree.head(5)
```

```
Out[787... 0    1
1    0
2    1
3    0
4    1
Name: Outcome, dtype: category
Categories (2, int64): [0, 1]
```

```
In [788... x_dtree_train,x_dtree_test,y_dtree_train,y_dtree_test=train_test_split(x_dtree,y_dtree,test_size=0.2,random_state=50)
```

```
In [789... from sklearn.tree import DecisionTreeClassifier
dt=DecisionTreeClassifier()
dt.fit(x_dtree_train,y_dtree_train)
dpred=dt.predict(x_dtree_test)
```

```
In [790... dpred
```

```
Out[790... array([0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0,
        0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1,
        1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1,
        0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0,
        0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1,
        0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1,
        1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 1], dtype=int64)
```

```
In [791... from sklearn.metrics import accuracy_score,confusion_matrix,classification_report
```

```
In [792... accuracy_score(y_dtree_test,dpred)
```

```
Out[792... 0.7034482758620689
```

```
In [793... accuracy_score(y_dtree_train,dt.predict(x_dtree_train))
```

```
Out[793... 1.0
```

```
In [794... # Now use PCA independent variables and check accuracy
```

```
In [795... x_var.head(5)
```

```
Out[795...
```

	pca1	pca2	pca3	pca4	pca5	pca6
0	0.407391	0.001247	0.171459	-0.012504	0.031812	0.080258
1	-0.205940	0.217559	0.105583	0.029990	0.076450	-0.148843
2	0.260472	0.098875	-0.134737	0.106131	-0.066869	0.350525
3	-0.391812	0.105149	-0.140500	0.006706	-0.048653	-0.130832
4	-0.348571	-0.508681	0.185572	0.031150	0.696183	0.253937

```
In [796... y_var.head(5)
```

```
Out[796... 0    1
1    0
2    1
3    0
4    1
Name: Outcome, dtype: category
Categories (2, int64): [0, 1]
```

```
In [797... x_train_p,x_test_p,y_train_p,y_test_p=train_test_split(x_var,y_var,test_size=0.2,random_state=50)
```

```
In [798... from sklearn.tree import DecisionTreeClassifier
pdt=DecisionTreeClassifier()
```

```
In [799... pdt.fit(x_train_p,y_train_p)
```

```
Out[799... DecisionTreeClassifier()
```

```
In [800...] pdpred=pdt.predict(x_test_p)
```

```
In [801...] accuracy_score(y_test_p,pdpred)
```

```
Out[801...] 0.6827586206896552
```

```
In [802...] accuracy_score(y_train_p,pdt.predict(x_train_p))
```

```
Out[802...] 1.0
```

```
In [803...] # There is no difference in accuracy between scaling data and PCA data  
# So we can use PCA independent variables for further models because our PCA data is multicollinearity  
# free and it reduces overfitting.  
# Here our model is clearly over fitted it is I think due to class imbalance issue  
# First we build all classification models and finally we see every model  
# with cross validation technique and finalize our model based on accuracy  
# got by cross validation technique.
```

Random Forest

```
In [804...] x_rforest=x_var  
           y_rforest=y_var
```

```
In [805...] x_rforest.head(5)
```

```
Out[805...]      pca1    pca2    pca3    pca4    pca5    pca6  
0    0.407391  0.001247  0.171459 -0.012504  0.031812  0.080258  
1   -0.205940  0.217559  0.105583  0.029990  0.076450 -0.148843  
2    0.260472  0.098875 -0.134737  0.106131 -0.066869  0.350525  
3   -0.391812  0.105149 -0.140500  0.006706 -0.048653 -0.130832  
4   -0.348571 -0.508681  0.185572  0.031150  0.696183  0.253937
```

```
In [806...] y_rforest.head(5)
```

```
Out[806...] 0    1  
           1    0  
           2    1  
           3    0  
           4    1  
Name: Outcome, dtype: category  
Categories (2, int64): [0, 1]
```

```
In [807...] from sklearn.ensemble import RandomForestClassifier  
           rf=RandomForestClassifier()
```

```
In [808...] x_rforest_train,x_rforest_test,y_rforest_train,y_rforest_test=train_test_split(x_rforest,y_rforest,test_size=0.2,
```

```
In [809...] rf.fit(x_rforest_train,y_rforest_train)
```

```
Out[809...] RandomForestClassifier()
```

```
In [810...] rf_pred=rf.predict(x_rforest_test)
```

```
In [811...] accuracy_score(y_rforest_test,rf_pred)
```

```
Out[811...] 0.7655172413793103
```

```
In [812...] accuracy_score(y_rforest_train,rf.predict(x_rforest_train))
```

```
Out[812...] 1.0
```

```
In [813... print(classification_report(y_rforest_test,rf_pred))
```

	precision	recall	f1-score	support
0	0.78	0.88	0.83	92
1	0.73	0.57	0.64	53
accuracy			0.77	145
macro avg	0.76	0.72	0.73	145
weighted avg	0.76	0.77	0.76	145

Support Vector Machine Model

```
In [814... x_svm=x_var  
y_svm=y_var
```

```
In [815... x_svm.head(5)
```

```
Out[815...      pca1    pca2    pca3    pca4    pca5    pca6  
0  0.407391  0.001247  0.171459 -0.012504  0.031812  0.080258  
1 -0.205940  0.217559  0.105583  0.029990  0.076450 -0.148843  
2  0.260472  0.098875 -0.134737  0.106131 -0.066869  0.350525  
3 -0.391812  0.105149 -0.140500  0.006706 -0.048653 -0.130832  
4 -0.348571 -0.508681  0.185572  0.031150  0.696183  0.253937
```

```
In [816... y_svm.head(5)
```

```
Out[816... 0    1  
1    0  
2    1  
3    0  
4    1  
Name: Outcome, dtype: category  
Categories (2, int64): [0, 1]
```

```
In [817... x_svm_train,x_svm_test,y_svm_train,y_svm_test=train_test_split(x_svm,y_svm,test_size=0.2,random_state=50)
```

```
In [818... from sklearn.svm import SVC
```

```
In [819... svm=SVC()
```

```
In [820... svm.fit(x_svm_train,y_svm_train)
```

```
Out[820... SVC()
```

```
In [821... svm_pred=svm.predict(x_svm_test)
```

```
In [822... accuracy_score(y_svm_test,svm_pred)
```

```
Out[822... 0.7241379310344828
```

```
In [823... accuracy_score(y_svm_train,svm.predict(x_svm_train))
```

```
Out[823... 0.7962003454231433
```

```
In [824... print(classification_report(y_svm_test,svm_pred))
```

	precision	recall	f1-score	support
0	0.72	0.91	0.80	89
1	0.75	0.43	0.55	56

accuracy			0.72	145
macro avg	0.73	0.67	0.67	145
weighted avg	0.73	0.72	0.70	145

Naive Bayes Classification Algorithm

```
In [825... x_train_nb,x_test_nb,y_train_nb,y_test_nb=train_test_split(data_sc_df,y_var,test_size=0.2,random_state=50)
```

```
In [826... from sklearn.naive_bayes import MultinomialNB
nb=MultinomialNB()
```

```
In [827... nb.fit(x_train_nb,y_train_nb)
```

```
Out[827... MultinomialNB()
```

```
In [828... pred_nb=nb.predict(x_test_nb)
```

```
In [829... accuracy_score(y_test_nb,pred_nb)
```

```
Out[829... 0.6137931034482759
```

```
In [830... accuracy_score(y_train_nb,nb.predict(x_train_nb))
```

```
Out[830... 0.6666666666666666
```

```
In [831... print(classification_report(y_test_nb,pred_nb))
```

	precision	recall	f1-score	support
0	0.61	1.00	0.76	89
1	0.00	0.00	0.00	56
accuracy			0.61	145
macro avg	0.31	0.50	0.38	145
weighted avg	0.38	0.61	0.47	145

K-Nearest Neighbors Algorithm

```
In [832... # I am using PCA independent variables
```

```
In [833... x_var.head(5)
```

```
Out[833...
```

	pca1	pca2	pca3	pca4	pca5	pca6
0	0.407391	0.001247	0.171459	-0.012504	0.031812	0.080258
1	-0.205940	0.217559	0.105583	0.029990	0.076450	-0.148843
2	0.260472	0.098875	-0.134737	0.106131	-0.066869	0.350525
3	-0.391812	0.105149	-0.140500	0.006706	-0.048653	-0.130832
4	-0.348571	-0.508681	0.185572	0.031150	0.696183	0.253937

```
In [834... from sklearn.cluster import KMeans
```

```
In [835... km = KMeans(n_clusters=5)
```

```
In [836... km.fit(x_var)
```

```
Out[836... KMeans(n_clusters=5)
```

```
In [837... km.labels_
```

```
Out[837...] array([[1, 3, 1, 3, 2, 3, 3, 0, 1, 1, 1, 0, 0, 2, 1, 4, 4, 2, 1, 1, 1, 0,
0, 1, 3, 1, 1, 1, 2, 3, 3, 1, 0, 1, 1, 3, 0, 3, 1, 1, 0, 1, 4, 4,
3, 1, 3, 3, 3, 0, 0, 3, 0, 4, 4, 2, 1, 3, 2, 1, 1, 4, 1, 3, 3, 3,
0, 1, 2, 3, 1, 3, 3, 3, 1, 4, 1, 2, 1, 3, 0, 3, 3, 0, 1, 1, 3, 0,
3, 3, 3, 2, 4, 3, 4, 3, 3, 2, 4, 0, 3, 4, 2, 0, 4, 3, 0, 1, 1, 3,
3, 3, 2, 3, 3, 1, 4, 2, 2, 2, 2, 4, 0, 1, 2, 1, 3, 2, 4, 4, 4, 0,
1, 1, 3, 1, 2, 1, 2, 1, 3, 2, 1, 0, 2, 1, 1, 3, 3, 3, 0, 1, 1, 2,
3, 4, 0, 3, 1, 3, 3, 1, 0, 4, 3, 0, 1, 2, 1, 1, 1, 4, 3, 1, 1, 0,
4, 0, 0, 3, 1, 1, 1, 0, 3, 3, 3, 0, 4, 4, 4, 3, 0, 3, 0, 1, 3, 1,
3, 4, 1, 2, 0, 0, 2, 0, 3, 1, 2, 1, 0, 3, 4, 4, 3, 0, 4, 3, 0, 3,
3, 3, 3, 0, 4, 1, 4, 3, 3, 3, 0, 2, 1, 1, 2, 0, 3, 1, 3, 3, 4, 0,
4, 3, 3, 2, 0, 0, 3, 1, 1, 1, 4, 4, 1, 3, 1, 4, 1, 4, 3, 4, 1, 2,
4, 0, 0, 1, 1, 0, 0, 2, 3, 1, 4, 4, 2, 2, 4, 0, 2, 2, 0, 1, 2, 1,
1, 1, 2, 0, 4, 2, 2, 1, 2, 3, 3, 1, 3, 3, 3, 2, 1, 0, 3, 4, 1, 3,
2, 2, 1, 2, 1, 1, 3, 1, 3, 2, 1, 0, 1, 3, 4, 1, 1, 0, 3, 3, 3, 1,
1, 3, 3, 1, 2, 1, 2, 0, 1, 1, 1, 0, 3, 1, 4, 3, 2, 2, 4, 3, 2, 0,
4, 4, 1, 4, 4, 4, 3, 3, 3, 3, 1, 1, 0, 3, 2, 1, 2, 1, 1, 2, 0, 4,
3, 3, 3, 1, 1, 1, 1, 2, 1, 4, 1, 2, 1, 2, 2, 3, 2, 2, 3, 1, 3, 3,
2, 3, 4, 3, 0, 0, 2, 2, 2, 3, 3, 3, 3, 1, 3, 3, 1, 4, 3, 3, 1, 3,
4, 3, 4, 4, 4, 3, 3, 2, 3, 1, 1, 3, 0, 1, 1, 3, 1, 1, 1, 4, 4, 4,
0, 4, 4, 4, 1, 3, 4, 2, 1, 1, 1, 0, 4, 3, 4, 2, 2, 2, 3, 1, 3, 1,
1, 0, 1, 3, 3, 0, 0, 3, 3, 1, 1, 1, 4, 2, 3, 1, 1, 2, 1, 3, 3, 3,
0, 1, 1, 0, 3, 2, 1, 3, 3, 3, 3, 2, 4, 3, 4, 4, 4, 4, 4, 2, 2, 0,
2, 1, 3, 4, 0, 0, 0, 1, 1, 3, 3, 1, 3, 4, 0, 4, 1, 1, 1, 1, 2, 4,
3, 4, 3, 4, 1, 0, 2, 1, 3, 3, 3, 2, 4, 0, 3, 1, 1, 4, 3, 1, 1, 0,
3, 1, 3, 0, 1, 2, 1, 3, 0, 2, 4, 3, 4, 3, 3, 4, 0, 4, 2, 3, 2, 3,
3, 0, 0, 3, 0, 3, 1, 3, 1, 2, 3, 1, 4, 3, 1, 4, 4, 1, 3, 1, 4, 3,
2, 1, 1, 1, 3, 1, 3, 4, 3, 1, 2, 2, 2, 2, 0, 4, 3, 2, 3, 3, 2, 2,
3, 2, 1, 3, 1, 4, 0, 0, 1, 2, 1, 1, 0, 0, 0, 3, 1, 2, 1, 1, 1, 4,
3, 2, 3, 4, 4, 3, 2, 3, 3, 2, 2, 1, 1, 2, 0, 3, 0, 0, 0, 3, 2, 1,
1, 3, 1, 2, 3, 1, 2, 2, 1, 1, 2, 3, 0, 2, 1, 2, 1, 1, 2, 2, 1, 4, 1,
2, 4, 3, 3, 1, 3, 2, 2, 1, 3, 4, 1, 3, 4, 0, 3, 3, 1, 0, 1, 4, 4,
0, 1, 3, 4, 3, 2, 1, 2, 1, 4, 4, 1, 3, 1, 1, 0, 3, 3, 4, 3])
```

```
In [838...] d=x_var
d["labels"] =km.labels_
```

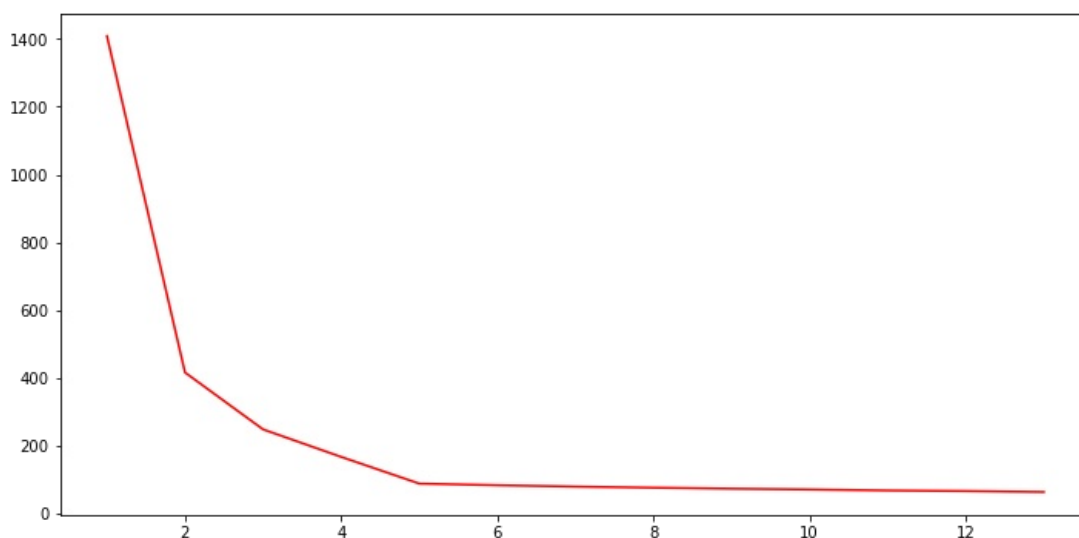
```
In [839...] ## elbow technique to decide the number of clusters
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [840...] wcss=[]

for i in range(1,14):
    km = KMeans(n_clusters=i, init="k-means++")
    km.fit(x_var)
    wcss.append(km.inertia_)

plt.figure(figsize=(12,6))
sns.lineplot(range(1,14), wcss, color="red")
```

```
Out[840...] <AxesSubplot:>
```



```
In [841...] km.inertia_
```

```
Out[841...] 62.186828070049316
```

```
In [842... # See by inertia at 5 clusters we got low inertia that is 62% so our data distinguished
# into 5 clusters
```

Model Selection

Using Cross Validation Technique to finalize model

```
In [843... #1.I am using cross validation technique other than over and under sampling.
#2.Because In over and under sampling there is a chance that accyracy is biased because we are adding
# data points in over sampling and removing some data points in under samplaing.
#3.But in cross validation technique this algorithm takes random sample subsets in dataset as test data and
# goes through further analysis so in this process there is less chance of overfitting and we can get
# unbiased accuracy.
#4.And in this technique it calculates accuracies of all subsets and takes its average so means it
# considers all values in data even data is in imbalance.
```

Logistic Regression Cross Validation

```
In [844... # Cross validation technique for logistic regression where
# we used PCA variables as independent variables.
```

```
In [845... from sklearn.model_selection import cross_validate
```

```
In [846... # I am taking 5 subsets
```

```
In [847... res_lr=cross_validate(lr,x_var,y_var,cv=5,return_train_score=True)
```

```
In [848... # Test Score
```

```
In [849... test_lr=np.average(res_lr['test_score'])
```

```
In [850... test_lr
```

```
Out[850... 0.7652107279693485
```

```
In [851... train_lr=np.average(res_lr['train_score'])
```

```
In [852... train_lr
```

```
Out[852... 0.7765880531236972
```

Decision Tree Cross Validation

```
In [853... x_var.head(5)
```

```
Out[853...
```

	pca1	pca2	pca3	pca4	pca5	pca6	labels
0	0.407391	0.001247	0.171459	-0.012504	0.031812	0.080258	1
1	-0.205940	0.217559	0.105583	0.029990	0.076450	-0.148843	3
2	0.260472	0.098875	-0.134737	0.106131	-0.066869	0.350525	1
3	-0.391812	0.105149	-0.140500	0.006706	-0.048653	-0.130832	3
4	-0.348571	-0.508681	0.185572	0.031150	0.696183	0.253937	2

```
In [854... y_var.head(5)
```

```
Out[854... 0    1
1    0
2    1
3    0
4    1
Name: Outcome, dtype: category
Categories (2, int64): [0, 1]
```

```
- from sklearn.model_selection import cross_validate
```

```
In [855... # I am taking 5 subsets
```

```
In [856... res_dt=cross_validate(pdt,x_var,y_var,cv=5,return_train_score=True)
```

```
In [857... # Test Score
```

```
In [858... test_dt=np.average(res_dt['test_score'])
```

```
In [859... test_dt
```

```
Out[859... 0.653227969348659
```

```
In [860... train_dt=np.average(res_dt['train_score'])
```

```
In [861... train_dt
```

```
Out[861... 1.0
```

Random Forest Cross Validation

```
In [862... # I am taking 4 subsets
```

```
In [863... res_rf=cross_validate(rf,x_rforest,y_rforest,cv=4,return_train_score=True)
```

```
In [864... # Test Score
```

```
In [865... test_rf=np.average(res_rf['test_score'])
```

```
In [866... test_rf
```

```
Out[866... 0.7569060773480663
```

```
In [867... train_rf=np.average(res_rf['train_score'])
```

```
In [868... train_rf
```

```
Out[868... 1.0
```

Support Vector Machine Model Cross Validation

```
In [869... # I am taking 5 subsets
```

```
In [870... res_svm=cross_validate(svm,x_svm,y_svm,cv=5,return_train_score=True)
```

```
In [871... # Test Score
```

```
In [872... test_svm=np.average(res_svm['test_score'])
```

```
In [873... test_svm
```

```
Out[873... 0.7416954022988506
```

```
In [874... train_svm=np.average(res_svm['train_score'])
```

```
In [875... train_svm
```

```
Out[875... 0.744131379905902
```

```
In [876... # I am taking 4 subsets
```

```
In [877... res_svm=cross_validate(svm,x_svm,y_svm,cv=4,return_train_score=True)
```

```
In [878... # Test Score
```



```
In [879... test_svm=np.average(res_svm['test_score'])
```

```
In [880... test_svm
```

```
Out[880... 0.7375690607734806
```

```
In [881... train_svm=np.average(res_svm['train_score'])
```

```
In [882... train_svm
```

```
Out[882... 0.7430939226519336
```

Naive Bayes Classification Algorithm Model Cross Validation

```
In [883... # I am taking 5 subsets
```

```
In [884... res_nb=cross_validate(nb,data_sc_df,y_var,cv=5,return_train_score=True)
```

```
In [885... # Test Score
```

```
In [886... test_nb=np.average(res_nb['test_score'])
```

```
In [887... test_nb
```

```
Out[887... 0.6560823754789272
```

```
In [888... # Train Score
```

```
In [889... train_nb=np.average(res_nb['train_score'])
```

```
In [890... train_nb
```

```
Out[890... 0.6560776606515395
```

Final Model Based On Cross Validation Technique

```
In [891... #1.There is only 1% difference in accuracy of train and test data in both Logistic Regression and  
# Naive Bayes Algorithm.  
#2.So these two models are finalized models.  
#3.But in these two we will take logistic regression as our final model with threshold 0.4  
#4.Because limitation of Naive Bayes is the assumption of independent predictors.  
# In real life it is almost impossible that we get a set of predictors.  
# which are completely independent.  
#5.Our data
```

```
In [ ]:
```

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js