

```
In [23]: import pandas as pd
import numpy as np
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn import svm
from sklearn.metrics import accuracy_score

In [2]: df=pd.read_csv("C:\\Users\\Pranav\\Desktop\\DATA SCIENCE DATA\\CVC file\\diabetes.csv")

In [3]: df.head()

Out[3]:      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 [6]: #number of rows and columns
df.shape

Out[6]: (768, 9)

In [5]: #describe mathamatical data df['']
df.describe()

Out[5]:      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 [8]: df['Outcome'].value_counts()
0      500
1      268
Name: Outcome, dtype: int64

In [ ]: #0--->NO Diabetes      1---> Diabetes

In [9]: df.groupby('Outcome').mean()

Out[9]:      Pregnancies      Glucose  BloodPressure  SkinThickness      Insulin      BMI  DiabetesPedigreeFunction      Age
Outcome
0      3.298000   109.980000      68.184000      19.664000   68.792000   30.304200              0.429734   31.190000
1      4.865672   141.257463      70.824627      22.164179  100.335821   35.142537              0.550500   37.067164

In [10]: #separting data and labels
X=df.drop(columns='Outcome',axis=1)
y=df['Outcome']

In [11]: print(X)
print(y)

      Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin   BMI  \
0            6         148                72              35         0   33.6
1            1          85                66              29         0   26.6
2            8         183                64              0         0   23.3
3            1          89                66              23         94  28.1
4            0         137                40              35        168  43.1
..          ...         ...              ...          ...      ...      ...
763         10         101                76              48        180   32.9
764          2         122                70              27         0   36.8
765          5         121                72              23        112   26.2
766          1         126                60              0         0   30.1
767          1          93                70              31         0   30.4

      DiabetesPedigreeFunction  Age
0              0.627    50
1              0.351    31
2              0.672    32
3              0.167    21
4              2.288    33
..          ...      ...
763             0.171    63
764             0.340    27
765             0.245    30
766             0.349    47
767             0.315    23

[768 rows x 8 columns]
0      1
1      0
2      1
3      0
4      1
..
763    0
764    0
765    0
766    1
767    0
Name: Outcome, Length: 768, dtype: int64

In [15]: #DATA STANDARDIZATION
scaler=StandardScaler()
scaler.fit(X)

Out[15]: ▼ StandardScaler
StandardScaler()

In [16]: standardized_data=scaler.transform(X)

In [18]: print(standardized_data)

[[ 0.63994726  0.84832379  0.14964075 ...  0.20401277  0.46849198
  1.4259954 ]
 [-0.84488505 -1.12339636 -0.16054575 ... -0.68442195 -0.36506078
 -0.19067191]
 [ 1.23388019  1.94372388 -0.26394125 ... -1.10325546  0.60439732
 -0.10558415]
 ...
 [ 0.3429808  0.00330087  0.14964075 ... -0.73518964 -0.68519336
 -0.27575966]
 [-0.84488505  0.1597866  -0.47073225 ... -0.24020459 -0.37110101
  1.17073215]
 [-0.84488505 -0.8730192  0.04624525 ... -0.20212881 -0.47378505
 -0.87137393]]

In [19]: X=standardized_data
y=y=df['Outcome']

In [20]: #training and test data
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=2,stratify=y)

In [21]: print("shape of X_train= ",X_train.shape)
print("shape of X_test= ",X_test.shape)
print("shape of y_train= ",y_train.shape)
print("shape of y_test= ",y_test.shape)

shape of X_train= (614, 8)
shape of X_test= (154, 8)
shape of y_train= (614,)
shape of y_test= (154,)

In [27]: #Training model
classifier=svm.SVC(kernel='linear')

In [28]: # trainin the support vector Machine classifier
classifier.fit(X_train,y_train)

Out[28]: ▼ SVC
SVC(kernel='linear')

In [30]: #model evaluation
#Accuracy score
#accarcy of training data
X_train_prediction=classifier.predict(X_train)
training_data_accuracy=accuracy_score(X_train_prediction,y_train)

In [31]: print('Accuracy on training data:',training_data_accuracy)

Accuracy on training data: 0.7886449511400652

In [33]: #accuracy of testing data
X_test_prediction=classifier.predict(X_test)
testing_data_accuracy=accuracy_score(X_test_prediction,y_test)
print('Accuracy on test data:',testing_data_accuracy)

Accuracy on test data: 0.7727272727272727

In [34]: #Making a Predictive system
#changing the input _data to a numpy array
input_data1=(5,116,74,0,0,25.6,0.201,30)
input_data_numpy_array1=np.asarray(input_data1)

In [35]: #reshape the np array as we are predicting for one instance
input_data_resaped1=input_data_numpy_array1.reshape(1,-1)

In [36]: input_data_resaped1

Out[36]: array([[ 5. , 116. , 74. , 0. , 0. , 25.6 , 0.201,
 30. ]])

In [61]: #standard input data
std_data=scaler.transform(input_data_resaped1)

C:\Users\Pranav\.continuum\ppp\lib\site-packages\sklearn\base.py:450: UserWarning: X does not have valid feature names, but StandardScaler was fitted with feature names
warnings.warn(

In [62]: prediction=classifier.predict(std_data)
print(prediction)

[0]

In [65]: if(prediction[0]==0):
print('The patient is Non-Diabetes')
else:
print('The patient is Diabetes')

The patient is Non-Diabetes

In [74]: input_data2=(3,158,76,36,245,31.6,0.851,28)
input_data_numpy_array2=np.asarray(input_data2)

In [75]: #reshape the np array as we are predicting for one instance
input_data_resaped2=input_data_numpy_array2.reshape(1,-1)

In [76]: input_data_resaped2

Out[76]: array([[ 3. , 158. , 76. , 36. , 245. , 31.6 , 0.851,
 28. ]])

In [77]: #standard input data
std_data1=scaler.transform(input_data_resaped2)

C:\Users\Pranav\.continuum\ppp\lib\site-packages\sklearn\base.py:450: UserWarning: X does not have valid feature names, but StandardScaler was fitted with feature names
warnings.warn(

In [78]: prediction=classifier.predict(std_data1)
print(prediction)

[1]

In [79]: if(prediction[0]==0):
print('The patient is Non-Diabetes')
else:
print('The patient is Diabetes')

The patient is Diabetes

In [ ]:

In [ ]: 
```