

Importing the Libraries

```
In [1]: import numpy as np # used for array
import pandas as pd # used for data processing (data loading, data manipulation, etc)
from sklearn.preprocessing import StandardScaler # standardize data to fit in a common range
from sklearn.model_selection import train_test_split # used for splitting the data into training data and testing data
from sklearn import svm #importing support vector machine
from sklearn.metrics import accuracy_score
```

Data Collection

```
In [2]: # loading the data from csv file to Pandas DataFrame

diabetes_dataset = pd.read_csv('diabetes.csv')
```

```
In [3]: # .head() displays the first 5 rows, and first n columns

# python indexes the column and row number starting from 0.

diabetes_dataset.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 [4]: # number of data points & number of features
# .shape() gives number of rows x columns (38523,12)

diabetes_dataset.shape
```

```
Out[4]: (768, 9)
```

```
In [5]: # getting some information about the data

diabetes_dataset.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 [6]: # getting the statistical measures of the data

diabetes_dataset.describe()
```

```
Out[6]:
```

	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 [7]: # label column's different values and their number of occurrences
# spread isn't good, the closer the value counts, better the data
diabetes_dataset['Outcome'].value_counts()
```

```
Out[7]: 0    500
        1    268
        Name: Outcome, dtype: int64

0 --> Non-Diabetic

1 --> Diabetic
```

```
In [8]: diabetes_dataset.groupby('Outcome').mean()
```

```
Out[8]:
```

	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 [9]: # separating the data and labels
X = diabetes_dataset.drop(columns = 'Outcome', axis=1)
Y = diabetes_dataset['Outcome']
```

```
In [10]: print(X)
```

```
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              0              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]

```
In [11]: print(Y)
```

```
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
```

Data Standardization

```
In [12]: scaler = StandardScaler()
```

```
In [13]: scaler.fit(X)
```

```
Out[13]: StandardScaler()
```

```
In [14]: standardized_data = scaler.transform(X)
```

```
In [15]: 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 [16]: X = standardized_data
Y = diabetes_dataset['Outcome']
```

```
In [17]: print(X)
print(Y)
```

```
[[ 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]]
```

```
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
```

Train Test Split

```
In [18]: X_train, X_test, Y_train, Y_test = train_test_split(X,Y, test_size = 0.2, stratify=Y, random_state=2)

# split the into 4 variables, X_train, X_test, Y_train, Y_test.

# test_size is 20%, therefore 80% of the data and labels will be in X_train, Y_train, and remaining 20% will be in X_test, Y_test.

# stratify = Y references that we split the data according to Y, i.e, 1 or 0
```

```
In [19]: print(X.shape, X_train.shape, X_test.shape)

(768, 8) (614, 8) (154, 8)
```

Training the Model

```
In [20]: classifier = svm.SVC(kernel='linear')

# load the svm.SVC function into variable classifier

# svm - support vector machine

# SVC - support vector classifier function

# kernel = 'linear' implies the type of SVC
```

```
In [21]: #training the support vector Machine Classifier
classifier.fit(X_train, Y_train)
```

```
Out[21]: SVC(kernel='linear')
```

Model Evaluation

Accuracy Score

```
In [22]: # accuracy score on the training data
X_train_prediction = classifier.predict(X_train)
training_data_accuracy = accuracy_score(X_train_prediction, Y_train)
```

```
In [23]: print('Accuracy score of the training data : ', training_data_accuracy)

Accuracy score of the training data :  0.7866449511400652
```

```
In [24]: # accuracy score on the test data
X_test_prediction = classifier.predict(X_test)
test_data_accuracy = accuracy_score(X_test_prediction, Y_test)
```

```
In [25]: print('Accuracy score of the test data : ', test_data_accuracy)

Accuracy score of the test data :  0.7727272727272727
```

Making a Predictive System

```
In [26]: input_data = (5,166,72,19,175,25.8,0.587,51)
# input_data = (Pregnancies,Glucose,BP,Skin Thickness, Insulin, BMI, Diabetes Pedigree Function, Age)

# changing the input_data to numpy array
input_data_as_numpy_array = np.asarray(input_data)

# reshape the array as we are predicting for one instance
input_data_reshaped = input_data_as_numpy_array.reshape(1,-1)

# standardize the input data
std_data = scaler.transform(input_data_reshaped)
print(std_data)

prediction = classifier.predict(std_data)
print(prediction)

if (prediction[0] == 0):
    print('The person is not diabetic')
else:
    print('The person is diabetic')
```

```
[[ 0.3429808  1.41167241  0.14964075 -0.09637905  0.82661621 -0.78595734
  0.34768723  1.51108316]]
```

[1]  
The person is diabetic