

Import libraries

```
In [2]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.tree import DecisionTreeClassifier
```

```
In [3]: df = pd.read_csv('diabetes.csv')
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 [4]:

```
df.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]:

```
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 [6]: df.shape
```

```
Out[6]: (768, 9)
```

```
In [7]: # seperate out features and target value from dataset
```

```
X = df.drop(['Outcome'],axis = 1).values  
y = df['Outcome'].values
```

```
In [8]: X.shape
```

```
Out[8]: (768, 8)
```

```
In [9]: y.shape
```

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

```
In [10]: # split the data in training and testing set
```

```
X_train, X_test, y_train,y_test = train_test_split(X,y, test_size = 0.25, random_state = 42)
```

```
In [11]: print("X_train shape : " , X_train.shape)  
print("X_test shape : " , X_test.shape)  
print("y_train shape : " , y_train.shape)  
print("y_test shape : " , y_test.shape)
```

```
X_train shape : (576, 8)  
X_test shape : (192, 8)  
y_train shape : (576,)  
y_test shape : (192,)
```

```
In [12]: # Model

clf = DecisionTreeClassifier()

# fitting
clf.fit(X_train,y_train)
```

```
Out[12]: DecisionTreeClassifier()
```

```
In [13]: # predicting
y_pred = clf.predict(X_test)
y_pred
```

```
Out[13]: array([0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0,
                0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 0,
                0, 1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 1,
                0, 1, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 1,
                1, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1,
                0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1, 1, 1,
                0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0,
                0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1,
                0, 0, 1, 0, 1, 1, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0])
```

```
In [14]: acc = metrics.accuracy_score(y_test,y_pred)
print("Accuracy : ",acc)
```

```
Accuracy : 0.703125
```

```
In [18]: y_pred_df = pd.DataFrame(y_pred)
```

```
In [22]: y_pred_df["Actual"] = y_test
```

Out[22]:

	0	Actual
0	0	0
1	0	0
2	0	0
3	0	0
4	0	0
...
187	0	0
188	1	1
189	1	0
190	0	1
191	0	0

192 rows × 2 columns

```
In [23]: y_pred_df.columns = ['Predcited', 'Actual']
```

```
In [24]: y_pred_df
```

```
Out[24]:
```

	Predcited	Actual
0	0	0
1	0	0
2	0	0
3	0	0
4	0	0
...
187	0	0
188	1	1
189	1	0
190	0	1
191	0	0

192 rows × 2 columns

```
In [26]: # saving results to csv  
y_pred_df.to_csv("Actual vs predicted DT.csv")
```

```
In [ ]:
```