

Untitled108

July 2, 2023

```
[1]: ## import the required library
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

1 Module 8: Decision Tree

Assignment Contact us: support@intellipaat.com / © Copyright Intellipaat / All rights reserved
Intel iPaat Python for Data Science Certification Course Problem Statement: You work in XYZ Company as a Python Data Scientist. The company officials have collected some data on salaries based on year of experience and wish for you to create a model from it. Dataset: diabetes.csv
Tasks To Be Performed: 1. Load the dataset using pandas 2. Extract data from outcome column is a variable named Y 3. Extract data from every column except outcome column in a variable named X 4. Divide the dataset into two parts for training and testing in 70% and 30% proportion 5. Create and train Decision Tree Model on training set 6. Make predictions based on the testing set using the trained model 7. Check the performance by calculating the confusion matrix and accuracy score of the model

```
[2]: ##Load the dataset using pandas
data=pd.read_csv(r'C:/Users/Vikas/Desktop/diabetes-1.csv')
```

```
[3]: data
```

```
[3]:
```

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

[768 rows x 9 columns]

```
[31]: data.head(5)
```

```
[31]:   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
```

	DiabetesPedigreeFunction	Age	Outcome
0	0.627	50	1
1	0.351	31	0
2	0.672	32	1
3	0.167	21	0
4	2.288	33	1

```
[32]: data.tail(5)
```

```
[32]:   Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin   BMI  \
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	Outcome
763	0.171	63	0
764	0.340	27	0
765	0.245	30	0
766	0.349	47	1
767	0.315	23	0

```
[33]: data.shape
```

```
[33]: (768, 9)
```

```
[28]: data.columns
```

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

```
[29]: data.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
```

```
[30]: data.isnull().sum()
```

```
[30]: Pregnancies           0  
      Glucose             0  
      BloodPressure        0  
      SkinThickness        0  
      Insulin              0  
      BMI                  0  
      DiabetesPedigreeFunction 0  
      Age                  0  
      Outcome              0  
      dtype: int64
```

```
[12]: ##Extract data from outcome column is a variable named Y  
      y=pd.DataFrame(data.iloc[:,-1])
```

```
[ ]:
```

```
[13]: y
```

```
[13]: Outcome
0      1
1      0
2      1
3      0
4      1
..     ...
763    0
764    0
765    0
766    1
767    0

[768 rows x 1 columns]
```

```
[15]: ##Extract data from every column except outcome column in a variable named X
x=pd.DataFrame(data.iloc[:, :-1])
```

```
[16]: x
```

```
[16]: 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]
```

```
[19]: ##Divide the dataset into two parts for training and testing in 70% and 30%
      ↪proportion
      from sklearn.model_selection import train_test_split
      x_train, x_test, y_train, y_test = train_test_split(
      x, y, test_size=0.30, random_state=0)
```

```
[20]: from sklearn.tree import DecisionTreeClassifier
      clf=DecisionTreeClassifier()
      clf=clf.fit(x_train,y_train)
```

```
[22]: clf.get_params()
```

```
[22]: {'ccp_alpha': 0.0,
      'class_weight': None,
      'criterion': 'gini',
      'max_depth': None,
      'max_features': None,
      'max_leaf_nodes': None,
      'min_impurity_decrease': 0.0,
      'min_samples_leaf': 1,
      'min_samples_split': 2,
      'min_weight_fraction_leaf': 0.0,
      'random_state': None,
      'splitter': 'best'}
```

```
[25]: ##Make predictions based on the testing set using the trained model
      y_pred=clf.predict(x_test)
      y_pred
```

```
[25]: array([0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0,
            0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1,
            0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1,
            0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
            1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0,
            0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0,
            0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0,
            1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0,
            1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 1, 1, 0,
            0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0,
            0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0], dtype=int64)
```

```
[26]: ## confusion matrix
      from sklearn.metrics import confusion_matrix
      cm=confusion_matrix(y_test,y_pred)
      print('confusionmatrix:\n',cm)
```

```
confusionmatrix:
[[123  34]
```

[26 48]]

```
[27]: ##accuracy score
      from sklearn.metrics import accuracy_score
      print("Accuracy:",accuracy_score(y_test,y_pred))
```

Accuracy: 0.7402597402597403

[]: