DecessionTreeClassifier

Akshat Kumar (242211003) 12-02-2025

```
# Libraries
import pandas as pd
from time import time,ctime
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn import metrics
import matplotlib.pyplot as plt
from sklearn.tree import plot_tree
import seaborn as sns
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

print("Timestamp: "+ctime(time()))

Timestamp: Wed Feb 12 06:27:22 2025

df = pd.read_csv('diabetes.csv')
df.head()
```

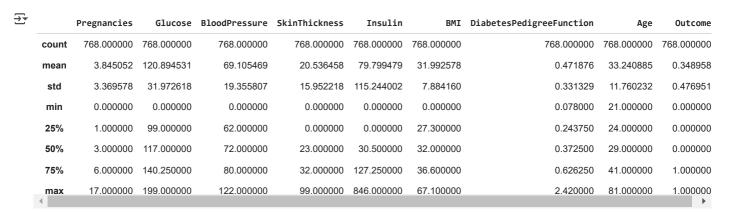
→		Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
	0	6	148	72	35	0	33.6	0.627	50	1
	1	1	85	66	29	0	26.6	0.351	31	0
	2	8	183	64	0	0	23.3	0.672	32	1
	3	1	89	66	23	94	28.1	0.167	21	0
	4	0	137	40	35	168	43.1	2.288	33	1

df.info()

</pr RangeIndex: 768 entries, 0 to 767 Data columns (total 9 columns): # Column Non-Null Count Dtype 0 Pregnancies 768 non-null int64 Glucose 768 non-null int64 BloodPressure 768 non-null int64 SkinThickness 768 non-null int64 Insulin 768 non-null int64 BMI 768 non-null float64 DiabetesPedigreeFunction 768 non-null float64 768 non-null int64 Age 8 Outcome 768 non-null

dtypes: float64(2), int64(7)
memory usage: 54.1 KB

Data Exploration and Preprocessing:



df.Outcome.value_counts()*100/len(df)



dtype: float64

1

df.groupby('Outcome').mean()

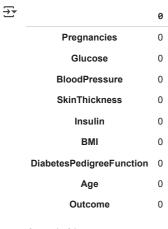
_		Pregnancies Glucose		BloodPressure	SkinThickness	inThickness Insulin		DiabetesPedigreeFunction	Age
	Outcome								
	0	3.298000	109.980000	68.184000	19.664000	68.792000	30.304200	0.429734	31.190000
	1	4.865672	141.257463	70.824627	22.164179	100.335821	35.142537	0.550500	37.067164

df.groupby('Outcome').agg(['mean','median'])

₹		Pregnancies		Glucose	cose BloodPressure SkinThickness		Insulin		BMI		DiabetesPed			
		mean	median	mean	median	mean	median	mean	median	mean	median	mean	median	mean
	Outcome													
	0	3.298000	2.0	109.980000	107.0	68.184000	70.0	19.664000	21.0	68.792000	39.0	30.304200	30.05	0.42973
	1	4.865672	4.0	141.257463	140.0	70.824627	74.0	22.164179	27.0	100.335821	0.0	35.142537	34.25	0.55050

2. Handle missing values (if any) appropriately.

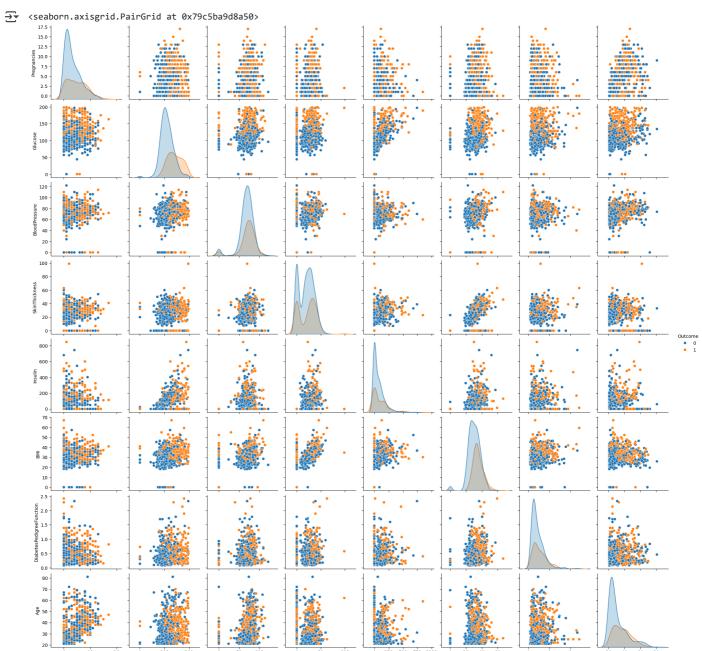
df.isnull().sum()



dtype: int64

sns.pairplot(df,hue='Outcome')







	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
Pregnancies	1.000000	0.129459	0.141282	-0.081672	-0.073535	0.017683	-0.033523	0.544341
Glucose	0.129459	1.000000	0.152590	0.057328	0.331357	0.221071	0.137337	0.263514
BloodPressure	0.141282	0.152590	1.000000	0.207371	0.088933	0.281805	0.041265	0.239528
SkinThickness	-0.081672	0.057328	0.207371	1.000000	0.436783	0.392573	0.183928	-0.113970
Insulin	-0.073535	0.331357	0.088933	0.436783	1.000000	0.197859	0.185071	-0.042163
ВМІ	0.017683	0.221071	0.281805	0.392573	0.197859	1.000000	0.140647	0.036242
DiabetesPedigreeFunction	-0.033523	0.137337	0.041265	0.183928	0.185071	0.140647	1.000000	0.033561
Age	0.544341	0.263514	0.239528	-0.113970	-0.042163	0.036242	0.033561	1.000000
Outcome	0.221898	0.466581	0.065068	0.074752	0.130548	0.292695	0.173844	0.238356

Model Implementation:

sns.heatmap(mat, annot=True)

```
from sklearn.neighbors import KNeighborsClassifier
X=df.drop('Outcome',axis=1)
X.head()
y=df['Outcome']
X_train, X_test, y_train, y_test = train_test_split(X, y,test_size=0.2,random_state=0)
from sklearn.preprocessing import StandardScaler
sc_x=StandardScaler()
X_train=sc_x.fit_transform(X_train)
X_test=sc_x.transform(X_test)
knn=KNeighborsClassifier(n_neighbors=5,metric='euclidean',p=2)
knn.fit(X_train,y_train)
             KNeighborsClassifier
     KNeighborsClassifier(metric='euclidean')
y_pred=knn.predict(X_test)
y_pred
0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1,
           1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 1,
           1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0,
           1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1,
           0, 0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0,
           0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0])
metrics.accuracy_score(y_test,y_pred)
mat = confusion_matrix(y_test, y_pred)
target_names = ['Diabetes', 'Normal']
print(classification_report(y_test, y_pred, target_names=target_names))
                             recall f1-score
                  precision
                                                support
                                0.87
        Diabetes
                       0.85
                                         0.86
                                                    107
          Normal
                       0.68
                                0.64
                                         0.66
                                                     47
        accuracy
                                          0.80
                                                    154
                       0.76
                                0.75
                                          0.76
                                                    154
       macro avg
    weighted avg
                       0.80
                                0.80
                                         0.80
                                                    154
knn=KNeighborsClassifier(n_neighbors=18,metric='euclidean',p=2)
knn.fit(X_train,y_train)
y_pred=knn.predict(X_test)
y_pred
knn.score(X_test,y_test)
from sklearn.metrics import confusion_matrix
mat = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(10, 8))
```

		precision	recall	f1-score	support
	Diabetes	0.83	0.92	0.87	107
	Normal	0.75	0.57	0.65	47
	accuracy			0.81	154
	macro avg	0.79	0.75	0.76	154
	weighted avg	0.81	0.81	0.80	154

