

Decision Tree

import dataset

```
import pandas as pd
df = pd.read_csv("drug200.csv")
df.head()
```

```
↗
```

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	F	HIGH	HIGH	25.355	drugY
1	47	M	LOW	HIGH	13.093	drugC
2	47	M	LOW	HIGH	10.114	drugC
3	28	F	NORMAL	HIGH	7.798	drugX
4	61	F	LOW	HIGH	18.043	drugY

```
df['Sex'].unique()
```

```
↗ array(['F', 'M'], dtype=object)
```

cleaning

```
# clean the data
```

encoding

```
from sklearn.preprocessing import LabelEncoder
```

```
le_sex = LabelEncoder().fit(df['Sex'])
df['Sex'] = le_sex.transform(df['Sex'])
```

```
le_BP = LabelEncoder().fit(df['BP'])
df['BP'] = le_BP.transform(df['BP'])
```

```
le_Cholesterol = LabelEncoder().fit(df['Cholesterol'])
df['Cholesterol'] = le_Cholesterol.fit_transform(df['Cholesterol'])
```

```
df.head()
```

```
↗
```

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	0	0	0	25.355	drugY
1	47	1	1	0	13.093	drugC
2	47	1	1	0	10.114	drugC
3	28	0	2	0	7.798	drugX
4	61	0	1	0	18.043	drugY

define x and y

```
import numpy as np
x = df[['Age', 'Sex', 'BP', 'Cholesterol', 'Na_to_K']].values
x[0:5]
```

```
↗ array([[23.,  0.,  0.,  0., 25.355],
        [47.,  1.,  1.,  0., 13.093],
        [47.,  1.,  1.,  0., 10.114],
```

```
[28. , 0. , 2. , 0. , 7.798],  
[61. , 0. , 1. , 0. , 18.043]])
```

```
y = df["Drug"].values  
y[0:5]
```

```
array(['drugY', 'drugC', 'drugC', 'drugX', 'drugY'], dtype=object)
```

✖ splitting

```
### finding best random state
```

```
# from sklearn.model_selection import train_test_split  
# from sklearn.preprocessing import StandardScaler  
# from sklearn.tree import DecisionTreeClassifier  
# from sklearn.metrics import accuracy_score  
  
# import time  
# t1 = time.time()  
# lst = []  
# for i in range(1,10):  
#     x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=i)  
#     sc = StandardScaler().fit(x_train)  
#     x_train = sc.transform(x_train)  
#     x_test = sc.transform(x_test)  
#     dtc = DecisionTreeClassifier(criterion="entropy", max_depth = 4)  
#     dtc.fit(x_train,y_train)  
#     yhat_test = dtc.predict(x_test)  
#     acc = accuracy_score(y_test, yhat_test)  
#     lst.append(acc)  
# t2 = time.time()  
# print(f"run time: {round((t2 - t1) / 60 , 0)} min")  
# print(f"accuracy_score = {round(max(lst),2)}")  
# print(f"random_state = {np.argmax(lst) + 1}")
```

```
from sklearn.model_selection import train_test_split  
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=7)
```

✖ scaling

```
from sklearn.preprocessing import StandardScaler  
sc = StandardScaler().fit(x_train)  
x_train = sc.transform(x_train)  
x_test = sc.transform(x_test)
```

✖ fit train data

```
from sklearn.tree import DecisionTreeClassifier  
dtc = DecisionTreeClassifier(criterion="entropy", max_depth = 4)  
dtc.fit(x_train,y_train)
```

```
DecisionTreeClassifier  
DecisionTreeClassifier(criterion='entropy', max_depth=4)
```

✖ predict test data

```
yhat_test = dtc.predict(x_test)  
  
print (yhat_test [0:5])  
print (y_test [0:5])
```

```
➦ ['drugX' 'drugY' 'drugY' 'drugY' 'drugC']  
   ['drugX' 'drugY' 'drugY' 'drugY' 'drugC']
```

✓ evaluate the model

```
from sklearn.metrics import accuracy_score  
print("Accuracy_score (train data): ", accuracy_score(y_train, dtc.predict(x_train)))  
print("Accuracy_score (test data): ", accuracy_score(y_test, yhat_test))
```

```
➦ Accuracy_score (train data):  1.0  
   Accuracy_score (test data):  0.96
```

✓ predict new data

```
dtc.predict([[23, 0, 0, 0, 25.355]])
```

```
➦ array(['drugY'], dtype=object)
```

```
dtc.predict([[23, le_sex.transform(['F'])[0], le_BP.transform(['HIGH'])[0], le_Chol.transform(['HIGH'])[0], le_HbA1c.transform(['HIGH'])[0]])])
```

```
➦ array(['drugY'], dtype=object)
```

✓ save the model

```
# import joblib  
# joblib.dump(dtc, 'dtc_model.pkl')
```

✓ load the model

```
# import joblib  
# dtc = joblib.load('dtc_model.pkl')
```