

Decision Tree Classification

Import Libraries

In [1]:

```
# import libraries
import numpy as np
import pandas as pd
```

Load Dataset

In [2]:

```
#load dataset
from sklearn.datasets import load_breast_cancer
data = load_breast_cancer()
```

data

```
{ 'data': array([[1.799e+01, 1.038e+01, 1.228e+02, ..., 2.654e-01, 4.601e-01,
1.189e-01],
[2.057e+01, 1.777e+01, 1.329e+02, ..., 1.860e-01, 2.750e-01,
8.902e-02],
[1.969e+01, 2.125e+01, 1.300e+02, ..., 2.430e-01, 3.613e-01,
8.758e-02],
...,
[1.660e+01, 2.808e+01, 1.083e+02, ..., 1.418e-01, 2.218e-01,
7.820e-02],
[2.060e+01, 2.933e+01, 1.401e+02, ..., 2.650e-01, 4.087e-01,
1.240e-01],
[7.760e+00, 2.454e+01, 4.792e+01, ..., 0.000e+00, 2.871e-01,
7.039e-02]]),
'target': array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
1, 1, 1,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0,
1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0,
1, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 1,
1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0,
0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1,
1, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0,
0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0,
1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1,
1, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0,
0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0,
0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0,
1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1,
1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1,
1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1,
1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1])),
'frame': None,
'target_names': array(['malignant', 'benign'], dtype='<U9'),
'DESCR': '.. _breast_cancer_dataset:\n\nBreast cancer wisconsin (diagnosti
c) dataset\n-----\n\n**Data Set Characteristics:**\n\n      :Number of Instances: 569\n\n      :Number of Attributes: 30 numeric, predictive attributes and the class\n\n      :Attribute Information:\n      - radius (mean of distances from center to points on the perimeter)\n      - texture (standard deviation of gray-scale values)\n      - perimeter\n      - area\n      - smoothness (local variation in radius lengths)\n      - compactness (perimeter^2 / area - 1.0)\n      - concavity (severity of concave portions of the contour)\n      - concave points (number of concave portions of the contour)\n      - symmetry\n      - fractal dimension ("coastline approximation" - 1)\n\n      The mean, standard error, and "worst" or largest (mean of the three worst/largest values) of these features were computed for each image, resulting in
```

```

30 features. For instance, field 0 is Mean Radius, field 10 is Rad
ius SE, field 20 is Worst Radius.\n\n      - class:\n      - WDB
C-Malignant\n      - WDBC-Benign\n\n      :Summary Statistics:\n\n
===== \n
Min      Max\n      ===== \n      rad
ius (mean):          6.981 28.11\n      texture (mean):
9.71 39.28\n      perimeter (mean):          43.79 188.5\n      ar
ea (mean):          143.5 2501.0\n      smoothness (mean):
0.053 0.163\n      compactness (mean):          0.019 0.345\n      co
ncavity (mean):          0.0 0.427\n      concave points (mean):
0.0 0.201\n      symmetry (mean):          0.106 0.304\n      fr
actal dimension (mean):          0.05 0.097\n      radius (standard erro
r):          0.112 2.873\n      texture (standard error):          0.3
6 4.885\n      perimeter (standard error):          0.757 21.98\n      area
(standard error):          6.802 542.2\n      smoothness (standard erro
r):          0.002 0.031\n      compactness (standard error):          0.002
0.135\n      concavity (standard error):          0.0 0.396\n      concave p
oints (standard error):          0.0 0.053\n      symmetry (standard error):
0.008 0.079\n      fractal dimension (standard error):          0.001 0.03\n      rad
ius (worst):          7.93 36.04\n      texture (worst):
12.02 49.54\n      perimeter (worst):          50.41 251.2\n      ar
ea (worst):          185.2 4254.0\n      smoothness (worst):
0.071 0.223\n      compactness (worst):          0.027 1.058\n      co
ncavity (worst):          0.0 1.252\n      concave points (wors
t):          0.0 0.291\n      symmetry (worst):          0.
156 0.664\n      fractal dimension (worst):          0.055 0.208\n      ====
===== \n\n      :Missing Attribute Va
lues: None\n\n      :Class Distribution: 212 - Malignant, 357 - Benign\n\n
:Creator: Dr. William H. Wolberg, W. Nick Street, Olvi L. Mangasarian\n\n
:Donor: Nick Street\n\n      :Date: November, 1995\n\nThis is a copy of UCI ML
Breast Cancer Wisconsin (Diagnostic) datasets.\nhttps://goo.gl/U2Uwz2\n\nFea
tures are computed from a digitized image of a fine needle\naspirate (FNA) o
f a breast mass. They describe\ncharacteristics of the cell nuclei present
in the image.\n\nSeparating plane described above was obtained using\nMultis
urface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree\nConstruction Via
Linear Programming." Proceedings of the 4th\nMidwest Artificial Intelligence
and Cognitive Science Society,\npp. 97-101, 1992], a classification method w
hich uses linear\nprogramming to construct a decision tree. Relevant featur
es\nwere selected using an exhaustive search in the space of 1-4\nfeatures a
nd 1-3 separating planes.\n\nThe actual linear program used to obtain the se
parating plane\nin the 3-dimensional space is that described in:\n[K. P. Ben
nett and O. L. Mangasarian: "Robust Linear\nProgramming Discrimination of Tw
o Linearly Inseparable Sets",\nOptimization Methods and Software 1, 1992, 23
-34].\n\nThis database is also available through the UW CS ftp server:\n\nft
p ftp.cs.wisc.edu\ncd math-prog/cpo-dataset/machine-learn/WDBC/\n\n.. topi
c:: References\n\n      - W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nucle
ar feature extraction \n      for breast tumor diagnosis. IS&T/SPIE 1993 Inte
rnational Symposium on \n      Electronic Imaging: Science and Technology, vo
lume 1905, pages 861-870,\n      San Jose, CA, 1993.\n      - O.L. Mangasarian,
W.N. Street and W.H. Wolberg. Breast cancer diagnosis and \n      prognosis v
ia linear programming. Operations Research, 43(4), pages 570-577, \n      Jul
y-August 1995.\n      - W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machin
e learning techniques\n      to diagnose breast cancer from fine-needle aspir
ates. Cancer Letters 77 (1994) \n      163-171.',
'feature_names': array(['mean radius', 'mean texture', 'mean perimeter', 'm
ean area',
      'mean smoothness', 'mean compactness', 'mean concavity',
      'mean concave points', 'mean symmetry', 'mean fractal dimension',
      'radius error', 'texture error', 'perimeter error', 'area error',
      'smoothness error', 'compactness error', 'concavity error',
      'concave points error', 'symmetry error',

```

```
'fractal dimension error', 'worst radius', 'worst texture',
'worst perimeter', 'worst area', 'worst smoothness',
'worst compactness', 'worst concavity', 'worst concave points',
'worst symmetry', 'worst fractal dimension'], dtype='<U23'),
'filename': 'C:\\Users\\avira\\anaconda3\\lib\\site-packages\\sklearn\\data
sets\\data\\breast_cancer.csv'}
```

In [4]:

```
data.feature_names
```

Out[4]:

```
array(['mean radius', 'mean texture', 'mean perimeter', 'mean area',
'mean smoothness', 'mean compactness', 'mean concavity',
'mean concave points', 'mean symmetry', 'mean fractal dimension',
'radius error', 'texture error', 'perimeter error', 'area error',
'smoothness error', 'compactness error', 'concavity error',
'concave points error', 'symmetry error',
'fractal dimension error', 'worst radius', 'worst texture',
'worst perimeter', 'worst area', 'worst smoothness',
'worst compactness', 'worst concavity', 'worst concave points',
'worst symmetry', 'worst fractal dimension'], dtype='<U23')
```

In [5]:

```
data.target
```

Out[5]:

```
array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0,
1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0,
1, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 1,
1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0,
0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1,
1, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0,
0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0,
1, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1,
1, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1,
0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0,
0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0,
1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 1,
1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 0,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1,
1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1,
1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1,
1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1])
```

In [6]:

```
data.target_names
```

Out[6]:

```
array(['malignant', 'benign'], dtype='<U9')
```

In [7]:

```
# create dataframe
df = pd.DataFrame(np.c_[data.data, data.target], columns=[list(data.feature_names)+['target']]
df.head()
```

Out[7]:

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2069
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809

5 rows × 31 columns

In [8]:

```
df.tail()
```

Out[8]:

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symme
564	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890	0.17
565	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.09791	0.17
566	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.05302	0.15
567	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200	0.25
568	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000	0.15

5 rows × 31 columns

In [9]:

```
df.shape
```

Out[9]:

```
(569, 31)
```

Split Data

In [10]:

```
X = df.iloc[:, 0:-1]
y = df.iloc[:, -1]
```

In [11]:

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=2020)

print('Shape of X_train = ', X_train.shape)
print('Shape of y_train = ', y_train.shape)
print('Shape of X_test = ', X_test.shape)
print('Shape of y_test = ', y_test.shape)
```

```
Shape of X_train = (455, 30)
```

```
Shape of y_train = (455,)
```

```
Shape of X_test = (114, 30)
```

```
Shape of y_test = (114,)
```

Train Decision Tree Classification Model

In [12]:

```
from sklearn.tree import DecisionTreeClassifier
```

In [13]:

```
classifier = DecisionTreeClassifier(criterion='gini')
classifier.fit(X_train, y_train)
```

Out[13]:

```
DecisionTreeClassifier()
```

In [14]:

```
classifier.score(X_test, y_test)
```

Out[14]:

```
0.9385964912280702
```

In [15]:

```
classifier_entropy = DecisionTreeClassifier(criterion='entropy')  
classifier_entropy.fit(X_train, y_train)
```

Out[15]:

```
DecisionTreeClassifier(criterion='entropy')
```

In [16]:

```
classifier_entropy.score(X_test, y_test)
```

Out[16]:

```
0.9210526315789473
```

Feature Scaling

In [17]:

```
from sklearn.preprocessing import StandardScaler  
sc = StandardScaler()
```

In [18]:

```
sc.fit(X_train)
```

Out[18]:

```
StandardScaler()
```

In [19]:

```
X_train_sc = sc.transform(X_train)  
X_test_sc = sc.transform(X_test)
```

In [20]:

```
classifier_sc = DecisionTreeClassifier(criterion='gini')  
classifier_sc.fit(X_train_sc, y_train)  
  
classifier_sc.score(X_test_sc, y_test)
```

Out[20]:

```
0.9298245614035088
```

Predict Cancer

In [21]:

```
patient1 = [17.99,  
10.38,  
122.8,  
1001.0,  
0.1184,  
0.2776,  
0.3001,  
0.1471,  
0.2419,  
0.07871,  
1.095,  
0.9053,  
8.589,  
153.4,  
0.006399,  
0.04904,  
0.05373,  
0.01587,  
0.03003,  
0.006193,  
25.38,  
17.33,  
184.6,  
2019.0,  
0.1622,  
0.6656,  
0.7119,  
0.2654,  
0.4601,  
0.1189]
```

In [22]:

```
patient1 = np.array([patient1])  
patient1
```

Out[22]:

```
array([[1.799e+01, 1.038e+01, 1.228e+02, 1.001e+03, 1.184e-01, 2.776e-01,  
3.001e-01, 1.471e-01, 2.419e-01, 7.871e-02, 1.095e+00, 9.053e-01,  
8.589e+00, 1.534e+02, 6.399e-03, 4.904e-02, 5.373e-02, 1.587e-02,  
3.003e-02, 6.193e-03, 2.538e+01, 1.733e+01, 1.846e+02, 2.019e+03,  
1.622e-01, 6.656e-01, 7.119e-01, 2.654e-01, 4.601e-01, 1.189e-01]])
```

In [23]:

```
classifier.predict(patient1)
```

Out[23]:

```
array([0.])
```


In [24]:

```
data.target_names
```

Out[24]:

```
array(['malignant', 'benign'], dtype='<U9')
```

In [25]:

```
pred = classifier.predict(patient1)
```

In [26]:

```
if pred[0] == 0:  
    print('Patient has Cancer (malignant tumor)')  
else:  
    print('Patient has no Cancer (malignant benign)')
```

```
Patient has Cancer (malignant tumor)
```