# Decision Tree:

Decision Tree set up a tree structure on training data which helps make a decision based on rules

We can apply Decision tree for both classification and regression problems, so this decision tree also called as CART(Classification and Regression Tree)

Rules:

```
* Entropy(Information gain)
* Gini
```

In [2]:
```python
import pandas as pd
import numpy as np
```

In [4]:
```python
from sklearn.datasets import load_iris
iris=load_iris()
iris
```

. . .

In [6]:
```python
# conver iris dictionary into dataframe
data1=iris.data
c1=iris.feature_names
df=pd.DataFrame(data1,columns=c1)
df.head()
```

Out[6]:

|   | sepal length (cm) | sepal width (cm) | petal length (cm) | petal width (cm) |
|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 |

In [39]: 
```python
#now set the target to dataframe
df['Target']=iris.target
df.sample(5)
```

Out[39]:

|     | sepal length (cm) | sepal width (cm) | petal length (cm) | petal width (cm) | Target |
|-----|-------------------|------------------|-------------------|------------------|--------|
| 102 | 7.1               | 3.0              | 5.9               | 2.1              | 2      |
| 92  | 5.8               | 2.6              | 4.0               | 1.2              | 1      |
| 76  | 6.8               | 2.8              | 4.8               | 1.4              | 1      |
| 114 | 5.8               | 2.8              | 5.1               | 2.4              | 2      |
| 75  | 6.6               | 3.0              | 4.4               | 1.4              | 1      |

In [8]: 
```python
# find the null values
df.isna().sum()
```

Out[8]: 
```
sepal length (cm)    0
sepal width (cm)     0
petal length (cm)    0
petal width (cm)     0
Target               0
dtype: int64
```

In [15]: 
```python
# Now take the feature names into X and target values into y
X=df.drop('Target',axis=1)
y=df[['Target']]
```

In [25]: 
```python
# split the data into train test split
from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.3,random_state=42
X_train.head()
```

Out[25]:

|     | sepal length (cm) | sepal width (cm) | petal length (cm) | petal width (cm) |
|-----|-------------------|------------------|-------------------|------------------|
| 81  | 5.5               | 2.4              | 3.7               | 1.0              |
| 133 | 6.3               | 2.8              | 5.1               | 1.5              |
| 137 | 6.4               | 3.1              | 5.5               | 1.8              |
| 75  | 6.6               | 3.0              | 4.4               | 1.4              |
| 109 | 7.2               | 3.6              | 6.1               | 2.5              |

In [26]: 
```python
X_test.shape
```

Out[26]: (45, 4)

In [27]: 
```python
df.shape
```

Out[27]: (150, 5)

```
In [31]:  # import Decision tree module
          from sklearn.tree import DecisionTreeClassifier
          model=DecisionTreeClassifier(criterion="entropy")
          model.fit(X_train,y_train)
```

```
Out[31]:  DecisionTreeClassifier(class_weight=None, criterion='entropy', max_depth=None,
                                 max_features=None, max_leaf_nodes=None,
                                 min_impurity_decrease=0.0, min_impurity_split=None,
                                 min_samples_leaf=1, min_samples_split=2,
                                 min_weight_fraction_leaf=0.0, presort=False,
                                 random_state=None, splitter='best')
```

```
In [32]:  # find accuracy scorer
          from sklearn.metrics import accuracy_score
          y_pred=model.predict(X_test)
          y_pred
```
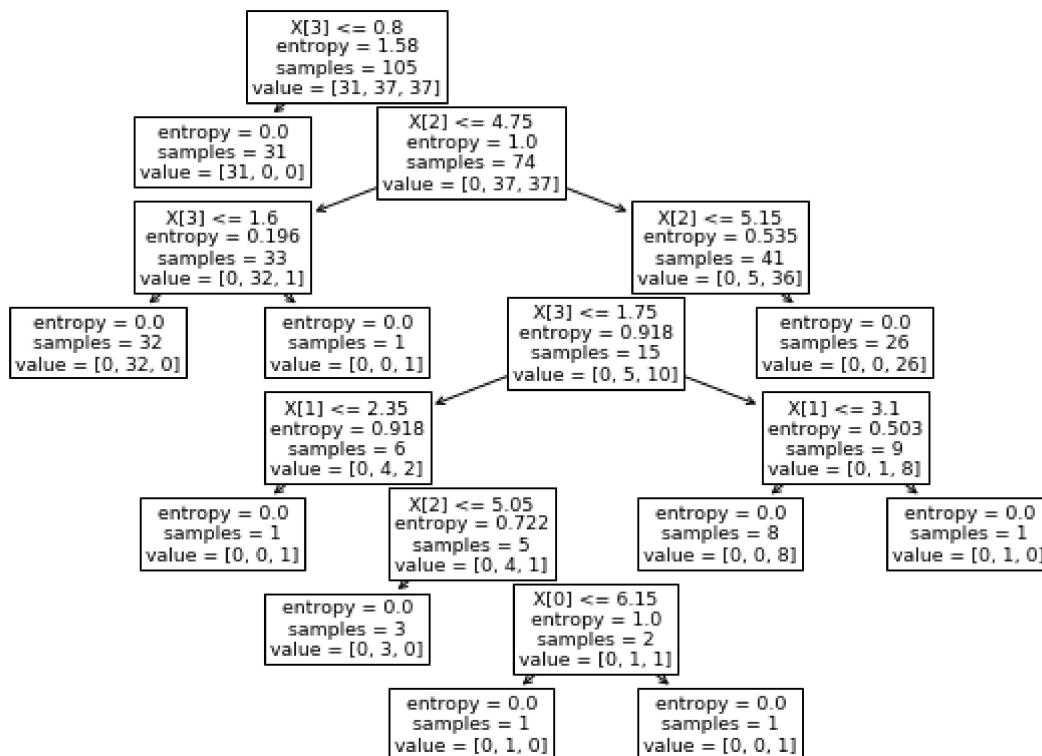
```
Out[32]:  array([1, 0, 2, 1, 1, 0, 1, 2, 1, 1, 1, 0, 0, 0, 0, 1, 2, 1, 1, 2, 0, 2,
                 0, 2, 2, 2, 2, 2, 0, 0, 0, 0, 1, 0, 0, 2, 1, 0, 0, 0, 2, 1, 1, 0,
                 0])
```

```
In [33]:  accuracy_score(y_test,y_pred)
```

```
Out[33]:  0.9777777777777777
```

In [36]:
```python
# Decistion tree visulization
import matplotlib.pyplot as plt
from sklearn import tree
plt.figure(figsize=(10,7))
tree.plot_tree(model)
plt.show()
```

X[3] <= 0.8
entropy = 1.58
samples = 105
value = [31, 37, 37]

entropy = 0.0
samples = 31
value = [31, 0, 0]

X[2] <= 4.75
entropy = 1.0
samples = 74
value = [0, 37, 37]

X[3] <= 1.6
entropy = 0.196
samples = 33
value = [0, 32, 1]

X[2] <= 5.15
entropy = 0.535
samples = 41
value = [0, 5, 36]

entropy = 0.0
samples = 32
value = [0, 32, 0]

entropy = 0.0
samples = 1
value = [0, 0, 1]

X[3] <= 1.75
entropy = 0.918
samples = 15
value = [0, 5, 10]

entropy = 0.0
samples = 26
value = [0, 0, 26]

X[1] <= 2.35
entropy = 0.918
samples = 6
value = [0, 4, 2]

X[1] <= 3.1
entropy = 0.503
samples = 9
value = [0, 1, 8]

entropy = 0.0
samples = 1
value = [0, 0, 1]

X[2] <= 5.05
entropy = 0.722
samples = 5
value = [0, 4, 1]

entropy = 0.0
samples = 8
value = [0, 0, 8]

entropy = 0.0
samples = 1
value = [0, 1, 0]

entropy = 0.0
samples = 3
value = [0, 3, 0]

X[0] <= 6.15
entropy = 1.0
samples = 2
value = [0, 1, 1]

entropy = 0.0
samples = 1
value = [0, 1, 0]

entropy = 0.0
samples = 1
value = [0, 0, 1]

In [40]:
```python
model.predict([[5.8,2.6,4.0,1.2]])
```

Out[40]: array([1])

## Task: Apply decision tree classification for breast cancer dataset

In [38]:
```python
from sklearn.datasets import load_breast_cancer
c=load_breast_cancer()
c
```

Out[38]: {'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, 1,
        1, 1,
           0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 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, 1, 0, 1, 1, 0, 1,
           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, 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, 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, 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, 0, 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, 0, 0, 1, 0, 0, 0, 1, 0, 0,
           1, 1, 1, 1, 1, 0, 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, 0, 1, 1, 1, 1, 1, 0, 1, 1, 0,
           1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1,
           1, 0, 1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 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, 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, 0, 0, 0, 0, 0, 0, 1]),
    'target_names': array(['malignant', 'benign'], dtype='<U9'),
    'DESCR': '.. _breast_cancer_dataset:\n\nBreast cancer wisconsin (diagnostic)
    dataset\n---------------------------------------------\n\n**Data Set Character
    istics:**\n\n    :Number of Instances: 569\n\n    :Number of Attributes: 30 n
    umeric, 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        - c
    ompactness (perimeter^2 / area - 1.0)\n        - concavity (severity of conca
    ve portions of the contour)\n        - concave points (number of concave port
    ions of the contour)\n        - symmetry \n        - fractal dimension ("coas
    tline approximation" - 1)\n\n        The mean, standard error, and "worst" or

```
largest (mean of the three\n          largest values) of these features were co
mputed for each image,\n          resulting in 30 features.  For instance, fiel
d 3 is Mean Radius, field\n          13 is Radius SE, field 23 is Worst Radiu
s.\n\n        - class:\n                - WDBC-Malignant\n                    - W
DBC-Benign\n\n    :Summary Statistics:\n    ================================
====== ====== ======\n                                           Min    Max\n
=================================== ====== ======\n    radius (mean):
6.981  28.11\n    texture (mean):                        9.71   39.28\n    per
imeter (mean):                     43.79  188.5\n    area (mean):
143.5  2501.0\n    smoothness (mean):                   0.053  0.163\n    co
mpactness (mean):                  0.019  0.345\n    concavity (mean):
0.0    0.427\n    concave points (mean):               0.0    0.201\n    sym
metry (mean):                        0.106  0.304\n    fractal dimension (mea
n):            0.05   0.097\n    radius (standard error):             0.112
2.873\n    texture (standard error):            0.36   4.885\n    perimeter
(standard error):           0.757  21.98\n    area (standard error):
6.802  542.2\n    smoothness (standard error):         0.002  0.031\n    com
pactness (standard error):          0.002  0.135\n    concavity (standard erro
r):             0.0    0.396\n    concave points (standard error):     0.0
0.053\n    symmetry (standard error):           0.008  0.079\n    fractal di
mension (standard error):   0.001  0.03\n    radius (worst):
7.93   36.04\n    texture (worst):                     12.02  49.54\n    per
imeter (worst):                    50.41  251.2\n    area (worst):
185.2  4254.0\n    smoothness (worst):                  0.071  0.223\n    co
mpactness (worst):                 0.027  1.058\n    concavity (worst):
0.0    1.252\n    concave points (worst):              0.0    0.291\n    sym
metry (worst):                       0.156  0.664\n    fractal dimension (wors
t):            0.055  0.208\n    =================================== ======
======\n\n    :Missing Attribute Values: None\n\n    :Class Distribution: 212
- Malignant, 357 - Benign\n\n    :Creator:  Dr. William H. Wolberg, W. Nick S
treet, 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) dataset
s.\nhttps://goo.gl/U2Uwz2\n\nFeatures are computed from a digitized image of
a fine needle\naspirate (FNA) of a breast mass.  They describe\ncharacteristi
cs of the cell nuclei present in the image.\n\nSeparating plane described abo
ve was obtained using\nMultisurface Method-Tree (MSM-T) [K. P. Bennett, "Deci
sion Tree\nConstruction Via Linear Programming." Proceedings of the 4th\nMidw
est Artificial Intelligence and Cognitive Science Society,\npp. 97-101, 199
2], a classification method which uses linear\nprogramming to construct a dec
ision tree.  Relevant features\nwere selected using an exhaustive search in t
he space of 1-4\nfeatures and 1-3 separating planes.\n\nThe actual linear pro
gram used to obtain the separating plane\nin the 3-dimensional space is that
described in:\n[K. P. Bennett and O. L. Mangasarian: "Robust Linear\nProgramm
ing Discrimination of Two Linearly Inseparable Sets",\nOptimization Methods a
nd Software 1, 1992, 23-34].\n\nThis database is also available through the U
W CS ftp server:\n\nftp ftp.cs.wisc.edu\ncd math-prog/cpo-dataset/machine-lea
rn/WDBC/\n\n.. topic:: References\n\n   - W.N. Street, W.H. Wolberg and O.L.
Mangasarian. Nuclear feature extraction \n     for breast tumor diagnosis. IS
&T/SPIE 1993 International Symposium on \n     Electronic Imaging: Science an
d Technology, volume 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 via linear programming. Operations Research, 43(4), pages 570-577,
\n     July-August 1995.\n   - W.H. Wolberg, W.N. Street, and O.L. Mangasaria
n. Machine learning techniques\n     to diagnose breast cancer from fine-need
le aspirates. Cancer Letters 77 (1994) \n     163-171.',
 'feature_names': array(['mean radius', 'mean texture', 'mean perimeter', 'me
an 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\\Kanakamma\\Anaconda3\\lib\\site-packages\\sklearn\\d
     atasets\\data\\breast_cancer.csv'}
```

In [ ]: