

## Importing Libraries Required

```
In [1]: 1 import numpy as np
        2 import pandas as pd
        3 import pydotplus
        4 from sklearn import datasets
        5 import math
```

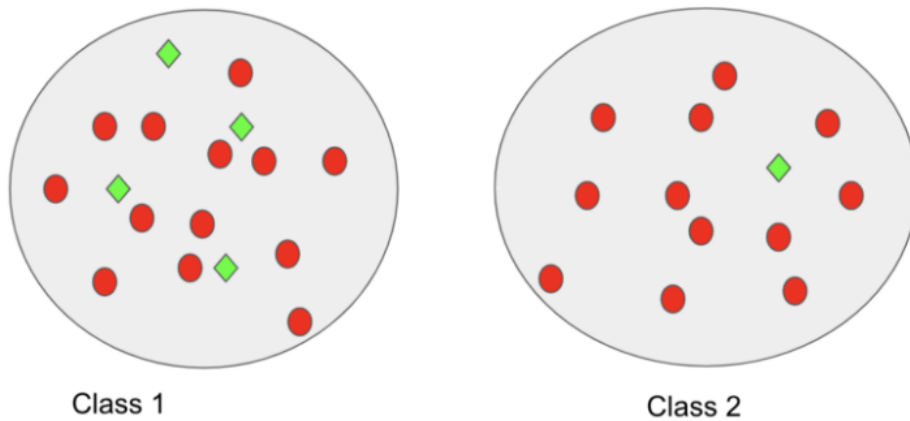
## Importing the 'Iris Dataset'

```
In [2]: 1 data = datasets.load_iris()
```

## Images for better clarity

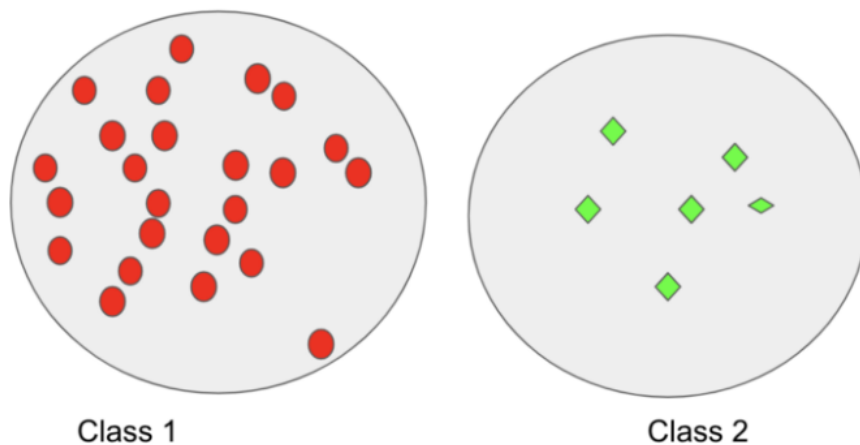
$$\text{Entropy} = \sum_i -p_i \log_2 p_i$$

Splitting can be of 2 types



Not so perfect split

and



Perfect split

## Creating class TreeNode

```
In [3]: 1 class TreeNode:
2
3     def __init__(self,data,output):
4
5
6         '''
7         This returns the root of our decision tree when its built.
8         'TreeNode' is a class which stores data, children, output, index
9
10        1. data: data represents the feature upon which the node was split when fitting the training data, data = None
11           for leaf node.
12        2. children : This is a dictionary which stores the key: value pairs where key represents the feature to be split
13           and value represents the corresponding child TreeNode
14        3. output : output represents the class with current majority at this instance of the decision tree
15        4. index : index is a unique index given to every node
16        '''
17        self.data = data
18        self.children = {}
19        self.output = output
20        self.index = -1
21
22        # We will have another functionality of adding a child to this TreeNode
23        def add_child(self,feature_value,obj):
24
25            # this functionality is used to add children to the dictionary
26            self.children[feature_value] = obj
```

## Creating Class of our classifier

In [4]:

```
1 class DecisionTreeClassifier:
2
3     # Defining the __init__ function
4     def __init__(self):
5         self.__root = None                # Here root represents the root node after building our tree
6
7     def __Count_unique(self, y):
8         # returns a dictionary with keys as unique values of
9         # Y(i.e no of classes) and the corresponding value as its frequency
10        d = {}
11        for i in y:
12            if i in d:
13                d[i] += 1
14            else:
15                d[i] = 1
16        return d
17
18    def __entropy(self, y):
19        # Returns the entropy at current y passed
20
21        freq_map = self.__Count_unique(y)
22
23        # Assigning current entropy to be zero
24        entropy_ = 0
25
26        # Total = how many values are there in y
27        total = len(y)
28
29        # Iterating through freq_map where key is the distinct class and its value is the freq of it
30        for i in freq_map:
31            prob = freq_map[i]/total
32            entropy_ += (-prob)*math.log2(prob)
33        return entropy_
34
35
36
37    def __Gain_ratio(self,x, y, f):
38        '''
39        Step 1: Find original Info
40        Step 2: Find Info(f)
41            - Info(f) is the summation of all (D(i)/D)*(entropy of D(i))
42        Step 3: Information Gain (Info(g)) = Info(origianl) - Info(f)
43        Step 4: Split info (split(info)) = summation of all (D(i)/D)*(math.log2(current_size/initial_size))
44        '''
45
46        info_original = self.__entropy(y) # Before splitting
47        info_f = 0                         # after splitting (initial value = 0)
48        split_info = 0                    # Initial value = 0
49        values = set(x[:,f])
50
51        # Creating a dataframe in X
52        df = pd.DataFrame(x)
53
54        #adding y as the last column
55        df[df.shape[1]] = y
56
57        initial_size = df.shape[0]
58
59        # Iterating through all unique values in our 'f' feature:
60        for uniq_val in values:
61            df1 = df[df[f] == uniq_val] # Creating another DataFrame with just our uniq value of our selected feature
62            current_size = df1.shape[0]
63            info_f += (current_size/initial_size)*self.__entropy(df1[df1.shape[1]-1])
64            split_info += (-current_size/initial_size)*math.log2(current_size/initial_size)
65
66        # In case our split_info is 0, to avoid divide by zero error
67        if split_info == 0:
68            return math.inf
69
70        info_gain = info_original - info_f
71        gain_ratio = info_gain/split_info
72
73        return gain_ratio
74
75
76
77    def __Gini_index(self,y):
78        # returns the gini index
79        freq_map = self.__Count_unique(y)
80        gini_index_ = 1
81        total = len(y)
82        for i in freq_map:
83            p = freq_map[i]/total
84            gini_index_ -= p**2
85        return gini_index_
86
87    def __Gini_gain(self,x, y, f):
88        gini_orig = self.__Gini_index(y)
89        gini_split_f = 0
90        values = set(x[:,f])
91        df = pd.DataFrame(x)
92        # Adding Y values as the last column in the dataframe
93        df[df.shape[1]] = y
```

```

94     initial_size = df.shape[0]
95
96     for i in values:
97         df1 = df[df[f] == i]
98         current_size = df1.shape[0]
99         gini_split_f += (current_size/initial_size)*self.__Gini_index(df1[df1.shape[1]-1])
100     gini_gain_ = gini_orig - gini_split_f
101     return gini_gain_
102
103
104 # Defining __decision_tree function
105 def __decision_tree(self,x, y, features, metric, classes, level):
106
107     '''
108     This is a recursive function so we will implement the main case first and then take a look at our base cases
109     '''
110
111     # If the node consists of only 1 class
112     if len(set(y)) == 1:
113         print("Level",level)
114         output = None
115         for i in classes:
116             if i in y:
117                 output = i
118                 print("Count of",i,"=",len(y))
119             else :
120                 print("Count of",i,"=",0)
121         if metric == "gain_ratio":
122             print("Current Entropy is = 0.0")
123         elif metric == "gini_index":
124             print("Current Gini Index is = 0.0")
125
126         print("Reached leaf Node")
127         print()
128         return TreeNode(None,output)
129
130     # If we have run out of features to split upon
131     # In this case we will output the class with maximum count
132     if len(features) == 0:
133         print("Level",level)
134         freq_map = self.__count_unique(y)
135         output = None
136         max_count = -math.inf
137         for i in classes:
138             if i not in freq_map:
139                 print("Count of",i,"=",0)
140             else :
141                 if freq_map[i] > max_count :
142                     output = i
143                     max_count = freq_map[i]
144                 print("Count of",i,"=",freq_map[i])
145
146         if metric == "gain_ratio":
147             print("Current Entropy is =",self.__entropy(y))
148         elif metric == "gini_index":
149             print("Current Gini Index is =",self.__Gini_index(y))
150
151         print("Reached leaf Node")
152         print()
153         return TreeNode(None,output)
154
155     # Main Loop: Finding the best feature to split upon
156     max_gain = -math.inf
157     final_feature = None
158
159     #Iterating through all features
160     for f in features:
161         if metric == 'Gain_ratio':
162             current_gain = self.__Gain_ratio(x, y, f)
163         elif metric == 'Gini_index':
164             current_gain = self.__Gini_gain(x, y, f)
165
166         # Now we have our gain for one feature f
167         if current_gain > max_gain:
168             max_gain = current_gain
169             final_feature = f
170
171     print("Level", level)
172     freq_map = self.__Count_unique(y)
173     output = None
174     max_count = -math.inf
175
176     for i in classes:
177         if i not in freq_map:
178             print("Count of",i,"= 0")
179         else:
180             if freq_map[i] > max_count:
181                 output = i
182                 max_count = freq_map[i]
183             print("Count of",i,"=",freq_map[i])
184
185     if metric == 'Gain_ratio':
186         print("Current Entropy is",self.__entropy(y))
187         print("Splitting on feature X[" ,final_feature,"] with Gain ratio " ,max_gain,sep = '')

```

```

188     print()
189 elif metric == 'Gini_index':
190     print("Current Gini Index is =",self.__Gini_index(y))
191     print("Splitting on feature X[,final_feature,'] with gini gain ",max_gain,sep="")
192     print()
193
194     unique_values = set(x[:,final_feature]) # unique_values represents the unique values of the feature selected
195     df = pd.DataFrame(x)
196     # Adding Y values as the last column in the dataframe
197     df[df.shape[1]] = y
198
199     current_node = TreeNode(final_feature,output)
200
201     # Now removing the selected feature from the List as we do not want to
202     # split on one feature more than once(in a given root to leaf node path)
203     index = features.index(final_feature)
204     features.remove(final_feature)
205
206     for i in unique_values:
207         # Creating a new dataframe with value of selected feature = i
208         df1 = df[df[final_feature] == i]
209         # Segregating the X and Y values and recursively calling on the splits
210         node = self.__decision_tree(df1.iloc[:,0:df1.shape[1]-1].values,df1.iloc[:,df1.shape[1]-1].values,
211                                     features,metric,classes,level+1)
212         current_node.add_child(i,node)
213
214     # Add the removed feature
215     features.insert(index,final_feature)
216
217     return current_node
218
219
220 # Defining the Main function 'fit':
221 def fit(self,x, y, metric = 'Gini_index'):
222     '''
223     X: This represents the training values upon which we will split
224     Y: This represents the target values or our classes
225     metric: This is the metric to be used while splitting, by default it is Gini Index
226     '''
227     # Creating an array 'features' which has the number of features
228     features = [i for i in range(len(x[0]))]
229
230     # Creating set classes which contains unique values from target y
231     classes = set(y)
232
233     # Because we are root node, initial level = 0
234     level = 0
235
236     # If metric passed is incorrect, assuming default as gain_ratio
237     if metric != 'Gini_index':
238         if metric != 'Gain_ratio':
239             metric = 'Gain_ratio'
240
241     # Now we will split and fit the function on root node
242     self.__root = self.__decision_tree(x, y, features, metric, classes, level)
243
244
245 def __predict_for(self,data,node):
246     # predicts the class for a given testing point and returns the answer
247
248     # We have reached a leaf node
249     if len(node.children) == 0 :
250         return node.output
251
252     val = data[node.data] # represents the value of feature on which the split was made
253     if val not in node.children :
254         return node.output
255
256     # Recursively call on the splits
257     return self.__predict_for(data,node.children[val])
258
259 def predict(self,x):
260     # This function returns Y predicted
261     # X should be a 2-D np array
262     Y = np.array([0 for i in range(len(x))])
263     for i in range(len(x)):
264         Y[i] = self.__predict_for(x[i],self.__root)
265     return Y
266
267 def score(self,x,y):
268     # returns the mean accuracy
269     Y_pred = self.predict(x)
270     count = 0
271     for i in range(len(Y_pred)):
272         if Y_pred[i] == y[i]:
273             count+=1
274     return count/len(Y_pred)
275

```

## Running Code

In [6]:

```
1 x = data.data
2 y = data.target
3
4
5 clf = DecisionTreeClassifier()
6 clf.fit(x, y)
7 y_pred = clf.predict(x)
8 print(y_pred)
9 score = clf.score(x,y)
10 print(score)
```

```
Level 0
Count of 0 = 50
Count of 1 = 50
Count of 2 = 50
Current Gini Index is = 0.6666666666666665
Splitting on feature X[2] with gini gain 0.6039999999999999
```

```
Level 1
Count of 0 = 4
Count of 1 = 0
Count of 2 = 0
Reached leaf Node
```

```
Level 1
Count of 0 = 13
Count of 1 = 0
Count of 2 = 0
Reached leaf Node
```

```
Level 1
Count of 0 = 7
Count of 1 = 0
Count of 2 = 0
Reached leaf Node
```

```
Level 1
Count of 0 = 7
Count of 1 = 0
Count of 2 = 0
Reached leaf Node
```

```
Level 1
Count of 0 = 13
Count of 1 = 0
Count of 2 = 0
Reached leaf Node
```

```
Level 1
Count of 0 = 1
Count of 1 = 0
Count of 2 = 0
Reached leaf Node
```

```
Level 1
Count of 0 = 2
Count of 1 = 0
Count of 2 = 0
Reached leaf Node
```

```
Level 1
Count of 0 = 1
Count of 1 = 0
Count of 2 = 0
Reached leaf Node
```

```
Level 1
Count of 0 = 2
Count of 1 = 0
Count of 2 = 0
Reached leaf Node
```

```
Level 1
Count of 0 = 0
Count of 1 = 5
Count of 2 = 0
Reached leaf Node
```

```
Level 1
Count of 0 = 0
Count of 1 = 7
Count of 2 = 1
Current Gini Index is = 0.21875
Splitting on feature X[0] with gini gain 0.21875
```

```
Level 2
Count of 0 = 0
Count of 1 = 0
Count of 2 = 1
Reached leaf Node
```

```
Level 2
Count of 0 = 0
```

Count of 1 = 1  
Count of 2 = 0  
Reached leaf Node

Level 2  
Count of 0 = 0  
Count of 1 = 2  
Count of 2 = 0  
Reached leaf Node

Level 2  
Count of 0 = 0  
Count of 1 = 1  
Count of 2 = 0  
Reached leaf Node

Level 2  
Count of 0 = 0  
Count of 1 = 1  
Count of 2 = 0  
Reached leaf Node

Level 2  
Count of 0 = 0  
Count of 1 = 1  
Count of 2 = 0  
Reached leaf Node

Level 2  
Count of 0 = 0  
Count of 1 = 1  
Count of 2 = 0  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 2  
Count of 2 = 3  
Current Gini Index is = 0.48  
Splitting on feature X[1] with gini gain 0.48

Level 2  
Count of 0 = 0  
Count of 1 = 1  
Count of 2 = 0  
Reached leaf Node

Level 2  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 1  
Reached leaf Node

Level 2  
Count of 0 = 0  
Count of 1 = 1  
Count of 2 = 0  
Reached leaf Node

Level 2  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 1  
Reached leaf Node

Level 2  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 1  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 5  
Count of 2 = 0  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 1  
Count of 2 = 3  
Current Gini Index is = 0.375  
Splitting on feature X[0] with gini gain 0.375

Level 2  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 1  
Reached leaf Node

Level 2  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 1

Reached leaf Node

Level 2

Count of 0 = 0

Count of 1 = 1

Count of 2 = 0

Reached leaf Node

Level 2

Count of 0 = 0

Count of 1 = 0

Count of 2 = 1

Reached leaf Node

Level 1

Count of 0 = 0

Count of 1 = 0

Count of 2 = 2

Reached leaf Node

Level 1

Count of 0 = 0

Count of 1 = 2

Count of 2 = 0

Reached leaf Node

Level 1

Count of 0 = 0

Count of 1 = 1

Count of 2 = 0

Reached leaf Node

Level 1

Count of 0 = 0

Count of 1 = 3

Count of 2 = 0

Reached leaf Node

Level 1

Count of 0 = 0

Count of 1 = 4

Count of 2 = 0

Reached leaf Node

Level 1

Count of 0 = 0

Count of 1 = 3

Count of 2 = 0

Reached leaf Node

Level 1

Count of 0 = 0

Count of 1 = 1

Count of 2 = 7

Current Gini Index is = 0.21875

Splitting on feature X[0] with gini gain 0.21875

Level 2

Count of 0 = 0

Count of 1 = 0

Count of 2 = 3

Reached leaf Node

Level 2

Count of 0 = 0

Count of 1 = 0

Count of 2 = 1

Reached leaf Node

Level 2

Count of 0 = 0

Count of 1 = 0

Count of 2 = 1

Reached leaf Node

Level 2

Count of 0 = 0

Count of 1 = 1

Count of 2 = 0

Reached leaf Node

Level 2

Count of 0 = 0

Count of 1 = 0

Count of 2 = 1

Reached leaf Node

Level 2

Count of 0 = 0

Count of 1 = 0

Count of 2 = 1

Reached leaf Node

Level 1



Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 2  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 6  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 3  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 2  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 1  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 3  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 1  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 1  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 1  
Count of 2 = 0  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 2  
Count of 2 = 0  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 1  
Count of 2 = 0  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 1  
Count of 2 = 0  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 4  
Count of 2 = 0  
Reached leaf Node

Level 1  
Count of 0 = 0  
Count of 1 = 2  
Count of 2 = 2  
Current Gini Index is = 0.5  
Splitting on feature X[0] with gini gain 0.5

Level 2  
Count of 0 = 0  
Count of 1 = 0  
Count of 2 = 1  
Reached leaf Node

Level 2  
Count of 0 = 0  
Count of 1 = 1  
Count of 2 = 0

## Comparing with sklearn algorithm

