

# Data Mining Assignment (CSE3132)

## Attribute Selection Algorithm

Srijit Roy, CSE-A, 3rd Year, 2351214

August 17, 2025

### Introduction:-

Decision Trees are one of the most widely used supervised learning methods for classification and regression tasks. They split a dataset into smaller subsets based on attribute values, recursively forming a tree structure where each internal node represents an attribute test, each branch represents an outcome of the test, and each leaf node represents a class label or decision.

An important part of constructing a decision tree is the **Attribute Selection Algorithm**, which determines the best attribute to split the dataset at each step. Different impurity measures such as **Entropy**, **Information Gain**, and **Gini Index** are used for selecting the most informative attribute.

### Formulas Used:-

#### Entropy

Entropy is a measurement of impurity or uncertainty in a dataset. It is represented as:-

$$H(S) = - \sum_{i=1}^n p_i \log_2 p_i$$

where  $p_i$  is the probability of class  $i$ .

#### Information Gain

Information Gain measures the reduction in entropy after splitting on an attribute. It is represented as:-

$$IG(S, A) = H(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} H(S_v)$$

where  $S_v$  is the subset of  $S$  for which attribute  $A$  has value  $v$ .

## Gini Index

The Gini Index measures impurity as:-

$$Gini(S) = 1 - \sum_{i=1}^n p_i^2$$

The weighted Gini Index for an attribute  $A$  is:-

$$GiniIndex(S, A) = \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Gini(S_v)$$

## Python Implementation:-

The following Python program (Tree.py) implements the Decision Tree induction using Entropy and Gini Index as attribute selection measures, followed by a visualization of the tree using Graphviz:-

```
1 import numpy as np #np is the alias for numpy
2 import pandas as pd #pd is the alias for pandas
3 from collections import Counter
4 from graphviz import Digraph
5
6 Data=pd.read_csv("Dataset.csv") #reads the data from CSV file
7 features=list(Data.columns[1:-1])
8 target="Decision" #the last column is the independent column named
   Decision
9
10 def entropy(labels):
11     c=Counter(labels) #counts the no of labels as Yes or No
12     tot=len(labels) #counts total no of labels
13     return -sum((count/tot)*np.log2(count/tot) for count in c.values())
   #Calculates total entropy
14
15 def info_gain(df,feature,target):
16     base_entropy=entropy(df[target])
17     values=df[feature].unique()
18     wei_entropy = sum(
19         (len(subset)/len(df))*entropy(subset[target]) #Calculates the
   individual info gains for each classes
20         for v in values
21         if len((subset:=df[df[feature]==v]))>0
22     )
23     return base_entropy - wei_entropy
24
25 def gini(labels):
26     c=Counter(labels)
27     tot=len(labels)
28     return 1-sum((count/tot)**2 for count in c.values()) #Calculates
   each gini index
29
30 def gini_index(df,feature,target):
31     values=df[feature].unique()
32     wei_gini = sum(
33         (len(subset)/len(df)) * gini(subset[target])
34         for v in values
```

```

35         if len (( subset := df[df[ feature ] == v])) > 0
36     )
37     return wei_gini
38
39 def build_decision_tree(df, features ,target ,method):
40     labels=df[target].tolist()
41     if labels.count(labels[0])==len(labels):
42         return labels[0] # Pure leaf
43     if not features:
44         return Counter(labels).most_common(1)[0][0]
45     if method=="entropy":
46         gains={f:info_gain(df,f,target) for f in features}
47         best_feature=max(gains,key=gains.get)
48     else:
49         ginis={f:gini_index(df,f,target) for f in features}
50         best_feature = min (ginis,key=ginis.get)
51
52     d_tree={best_feature:{}}
53     for v in df[ best_feature ].unique():
54         subset=df[df[best_feature]==v].drop(columns=[best_feature])
55         sub_features=[f for f in features if f!=best_feature]
56         d_tree[best_feature][v]=build_decision_tree(subset,sub_features
57             ,target,method)
58     return d_tree
59 no_of_nodes=0 #counter for no of nodes
60 def node_id():
61     global no_of_nodes
62     no_of_nodes+=1
63     return f"node{no_of_nodes}"
64
65 def visualize_tree(tree,graph=None,parent=None,edge_label=""):
66     if graph is None:
67         graph=Digraph(format="png")
68         graph.attr(rankdir="TB",splines="polyline") # Better layout
69         graph.attr("node",fontname="Helvetica")
70     if isinstance(tree,dict):
71         for feature,branches in tree.items():
72             id_of_node=node_id()
73             graph.node(id_of_node,feature,shape="box",style="rounded,
74                 filled",color="lightblue")
75             if parent:
76                 graph.edge(parent,id_of_node,label=edge_label)
77             for value,subtree in branches.items():
78                 visualize_tree(subtree,graph,id_of_node,str(value))
79     else:
80         id_of_leaf=node_id()
81         graph.node(id_of_leaf,str(tree),shape="ellipse",style="filled",
82             color ="lightgreen")
83         if parent:
84             graph.edge(parent,id_of_leaf,label=edge_label)
85     return graph
86
87 tree_entropy=build_decision_tree(Data,features,target,method="entropy")
88 graph_entropy=visualize_tree(tree_entropy)
89 graph_entropy.render("tree_entropy",view=True )
90 no_of_nodes = 0
91 tree_gini=build_decision_tree(Data,features,target,method="gini")

```

```

90 graph_gini=visualize_tree(tree_gini)
91 graph_gini.render("tree_gini",view=True)

```

Listing 1: Decision Tree Implementation (Tree.py)

## Provided Dataset:-

The dataset used for training the decision tree is shown below:

Weather	Parents	Financial Condition	Decision
Sunny	Yes	Poor	Cinema
Sunny	No	Rich	Play Tennis
Windy	Yes	Poor	Cinema
Windy	No	Poor	Cinema
Windy	No	Rich	Shopping
Rainy	Yes	Poor	Cinema
Rainy	No	Poor	Stay in
Rainy	No	Rich	Shopping

## Induced Graphs:-

The constructed decision trees using different attribute selection measures are as follows:

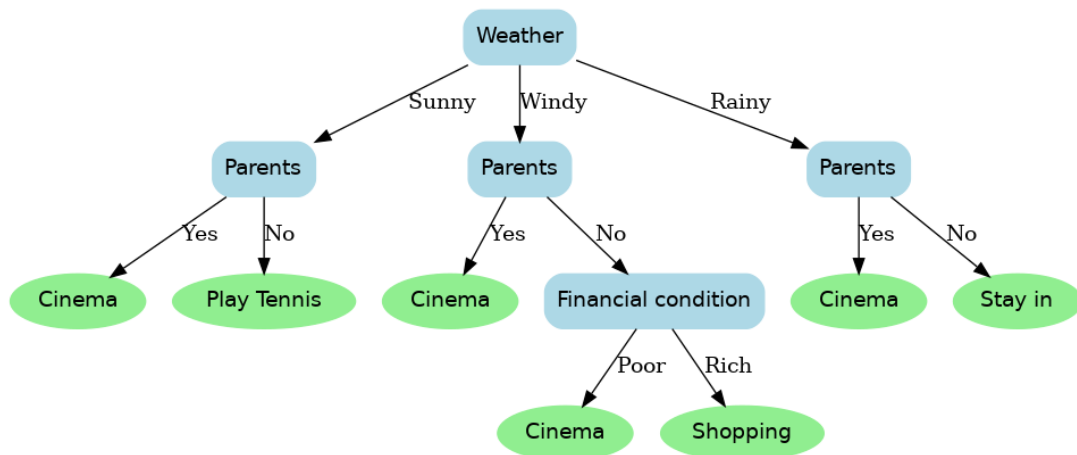


Figure 1: Decision Tree using Entropy (Information Gain)

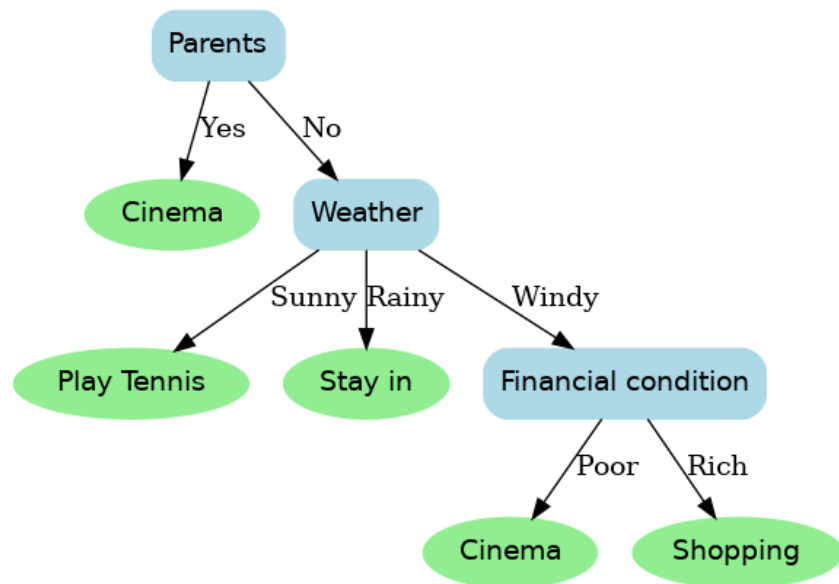


Figure 2: Decision Tree using Gini Index