

CART

# Introduction

- **CART( Classification And Regression Tree)** is a variation of the decision tree algorithm. It can handle both classification and regression tasks.
- Scikit-Learn uses the Classification And Regression Tree (CART) algorithm to train Decision Trees (also called “growing” trees).

## CART Algorithm

- CART is a predictive algorithm used in Machine learning and it explains how the target variable's values can be predicted based on other matters.

- It is a decision tree where each fork is split into a predictor variable and each node has a prediction for the target variable at the end.
- In the decision tree, nodes are split into sub-nodes on the basis of a threshold value of an attribute.
- The root node is taken as the training set and is split into two by considering the best attribute and threshold value. Further, the subsets are also split using the same logic.
- This continues till the last pure sub-set is found in the tree or the maximum number of leaves possible in that growing tree.

- The CART algorithm works via the following process:
  - The best split point of each input is obtained.
  - Based on the best split points of each input in Step 1, the new “best” split point is identified.
  - Split the chosen input according to the “best” split point.
  - Continue splitting until a stopping rule is satisfied or no further desirable splitting is available.

## Gini index/Gini impurity

- The **Gini index** is a metric for the classification tasks in CART. It stores the **sum of squared probabilities of each class**.
- It computes the **degree of probability of a specific variable** that is **wrongly being classified** when chosen randomly and a variation of the Gini coefficient.

- It works on categorical variables, provides outcomes either “successful” or “failure” and hence conducts binary splitting only.
- The degree of the Gini index varies from 0 to 1,
- Where 0 depicts that all the elements are allied to a certain class, or only one class exists there.
- The Gini index of value 1 signifies that all the elements are randomly distributed across various classes, and
- A value of 0.5 denotes the elements are uniformly distributed into some classes.

Mathematically, we can write Gini Impurity as follows:

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

where  $p_i$  is the probability of an object being classified to a particular class.

## Classification tree

- A classification tree is an algorithm where **the target variable is categorical**. The algorithm is then used to identify the “Class” within which the target variable is most likely to fall.
- Classification trees are used when the **dataset needs to be split into classes that belong to the response variable**(like yes or no)

## Regression tree

- A Regression tree is an algorithm where **the target variable is continuous** and the tree is used to predict its value.
- Regression trees are used when the **response variable is continuous**. For example, if the response variable is the temperature of the day.

# CART model representation

- CART models are formed by picking input variables and evaluating split points on those variables until an appropriate tree is produced.

Steps to create a Decision Tree using the CART algorithm:

- ***Greedy algorithm:*** In this the **input space is divided using the Greedy method** which is known as a recursive binary splitting.
- This is a numerical method within which all of the values are **aligned and several other split points** are tried and assessed using a cost function.
- ***Stopping Criterion:*** As it works its way down the tree with the training data, **the recursive binary splitting method** described above must know when to stop splitting.

- The most frequent halting method is to utilize a minimum amount of training data allocated to every leaf node. If the count is smaller than the specified threshold, the split is rejected and also the node is considered the last leaf node.
- ***Tree pruning:*** Decision tree's complexity is defined as the number of splits in the tree. Trees with fewer branches are recommended as they are simple to grasp and less prone to cluster the data.
- ***Data preparation for the CART:*** No special data preparation is required for the CART algorithm.



# Example

## Data set

We will work on same dataset in ID3. There are 14 instances of golf playing decisions based on outlook, temperature, humidity and wind factors.

| Day | Outlook  | Temp. | Humidity | Wind   | Decision |
|-----|----------|-------|----------|--------|----------|
| 1   | Sunny    | Hot   | High     | Weak   | No       |
| 2   | Sunny    | Hot   | High     | Strong | No       |
| 3   | Overcast | Hot   | High     | Weak   | Yes      |
| 4   | Rain     | Mild  | High     | Weak   | Yes      |
| 5   | Rain     | Cool  | Normal   | Weak   | Yes      |
| 6   | Rain     | Cool  | Normal   | Strong | No       |
| 7   | Overcast | Cool  | Normal   | Strong | Yes      |
| 8   | Sunny    | Mild  | High     | Weak   | No       |
| 9   | Sunny    | Cool  | Normal   | Weak   | Yes      |
| 10  | Rain     | Mild  | Normal   | Weak   | Yes      |
| 11  | Sunny    | Mild  | Normal   | Strong | Yes      |
| 12  | Overcast | Mild  | High     | Strong | Yes      |
| 13  | Overcast | Hot   | Normal   | Weak   | Yes      |
| 14  | Rain     | Mild  | High     | Strong | No       |

## Gini index

Gini index is a metric for classification tasks in CART. It stores sum of squared probabilities of each class. We can formulate it as illustrated below.

$$\text{Gini} = 1 - \sum (P_i)^2 \text{ for } i=1 \text{ to number of classes}$$

## Outlook

Outlook is a nominal feature. It can be sunny, overcast or rain. I will summarize the final decisions for outlook feature.

| Outlook  | Yes | No | Number of instances |
|----------|-----|----|---------------------|
| Sunny    | 2   | 3  | 5                   |
| Overcast | 4   | 0  | 4                   |
| Rain     | 3   | 2  | 5                   |

$$\text{Gini}(\text{Outlook}=\text{Sunny}) = 1 - (2/5)^2 - (3/5)^2 = 1 - 0.16 - 0.36 = 0.48$$

$$\text{Gini}(\text{Outlook}=\text{Overcast}) = 1 - (4/4)^2 - (0/4)^2 = 0$$

$$\text{Gini}(\text{Outlook}=\text{Rain}) = 1 - (3/5)^2 - (2/5)^2 = 1 - 0.36 - 0.16 = 0.48$$

Then, we will calculate weighted sum of gini indexes for outlook feature.

$$\text{Gini}(\text{Outlook}) = (5/14) \times 0.48 + (4/14) \times 0 + (5/14) \times 0.48 = 0.171 + 0 + 0.171 = 0.342$$

---

## Temperature

Similarly, temperature is a nominal feature and it could have 3 different values: Cool, Hot and Mild. Let's summarize decisions for temperature feature.

| Temperature | Yes | No | Number of instances |
|-------------|-----|----|---------------------|
| Hot         | 2   | 2  | 4                   |
| Cool        | 3   | 1  | 4                   |
| Mild        | 4   | 2  | 6                   |

$$\text{Gini}(\text{Temp}=\text{Hot}) = 1 - (2/4)^2 - (2/4)^2 = 0.5$$

$$\text{Gini}(\text{Temp}=\text{Cool}) = 1 - (3/4)^2 - (1/4)^2 = 1 - 0.5625 - 0.0625 = 0.375$$

$$\text{Gini}(\text{Temp}=\text{Mild}) = 1 - (4/6)^2 - (2/6)^2 = 1 - 0.444 - 0.111 = 0.445$$

We'll calculate weighted sum of gini index for temperature feature

$$\text{Gini}(\text{Temp}) = (4/14) \times 0.5 + (4/14) \times 0.375 + (6/14) \times 0.445 = 0.142 + 0.107 + 0.190 = 0.439$$

## Humidity

Humidity is a binary class feature. It can be high or normal.

| Humidity | Yes | No | Number of instances |
|----------|-----|----|---------------------|
| High     | 3   | 4  | 7                   |
| Normal   | 6   | 1  | 7                   |

$$\text{Gini}(\text{Humidity}=\text{High}) = 1 - (3/7)^2 - (4/7)^2 = 1 - 0.183 - 0.326 = 0.489$$

$$\text{Gini}(\text{Humidity}=\text{Normal}) = 1 - (6/7)^2 - (1/7)^2 = 1 - 0.734 - 0.02 = 0.244$$

Weighted sum for humidity feature will be calculated next

$$\text{Gini}(\text{Humidity}) = (7/14) \times 0.489 + (7/14) \times 0.244 = 0.367$$

## Wind

Wind is a binary class similar to humidity. It can be weak and strong.

| Wind   | Yes | No | Number of instances |
|--------|-----|----|---------------------|
| Weak   | 6   | 2  | 8                   |
| Strong | 3   | 3  | 6                   |

$$\text{Gini}(\text{Wind}=\text{Weak}) = 1 - (6/8)^2 - (2/8)^2 = 1 - 0.5625 - 0.062 = 0.375$$

$$\text{Gini}(\text{Wind}=\text{Strong}) = 1 - (3/6)^2 - (3/6)^2 = 1 - 0.25 - 0.25 = 0.5$$

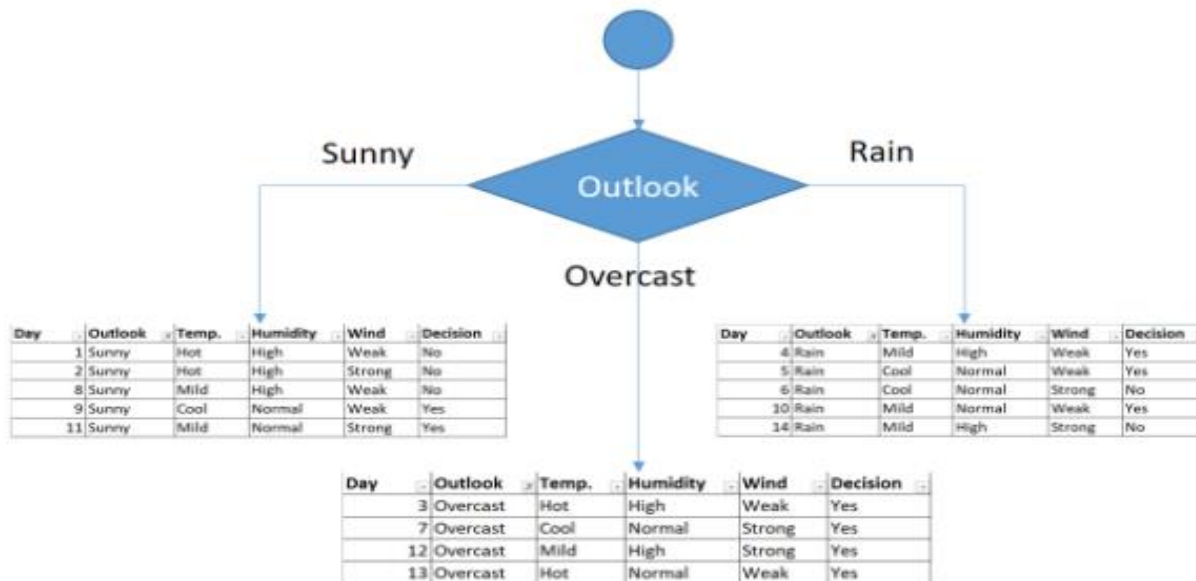
$$\text{Gini}(\text{Wind}) = (8/14) \times 0.375 + (6/14) \times 0.5 = 0.428$$

## Time to decide

We've calculated gini index values for each feature. The winner will be outlook feature because its cost is the lowest.

| Feature     | Gini index |
|-------------|------------|
| Outlook     | 0.342      |
| Temperature | 0.439      |
| Humidity    | 0.367      |
| Wind        | 0.428      |

We'll put outlook decision at the top of the tree.



First decision would be outlook feature

You might realize that sub dataset in the overcast leaf has only yes decisions. This means that overcast leaf is over.



Tree is over for overcast outlook leaf

We will apply same principles to those sub datasets in the following steps.

Focus on the sub dataset for sunny outlook. We need to find the gini index scores for temperature, humidity and wind features respectively.

| Day | Outlook | Temp. | Humidity | Wind   | Decision |
|-----|---------|-------|----------|--------|----------|
| 1   | Sunny   | Hot   | High     | Weak   | No       |
| 2   | Sunny   | Hot   | High     | Strong | No       |
| 8   | Sunny   | Mild  | High     | Weak   | No       |
| 9   | Sunny   | Cool  | Normal   | Weak   | Yes      |
| 11  | Sunny   | Mild  | Normal   | Strong | Yes      |

### Gini of temperature for sunny outlook

| Temperature | Yes | No | Number of instances |
|-------------|-----|----|---------------------|
| Hot         | 0   | 2  | 2                   |
| Cool        | 1   | 0  | 1                   |
| Mild        | 1   | 1  | 2                   |

$$\text{Gini}(\text{Outlook}=\text{Sunny and Temp.}=\text{Hot}) = 1 - (0/2)^2 - (2/2)^2 = 0$$

$$\text{Gini}(\text{Outlook}=\text{Sunny and Temp.}=\text{Cool}) = 1 - (1/1)^2 - (0/1)^2 = 0$$

$$\text{Gini}(\text{Outlook}=\text{Sunny and Temp.}=\text{Mild}) = 1 - (1/2)^2 - (1/2)^2 = 1 - 0.25 - 0.25 = 0.5$$

$$\text{Gini}(\text{Outlook}=\text{Sunny and Temp.}) = (2/5) \times 0 + (1/5) \times 0 + (2/5) \times 0.5 = 0.2$$

## Gini of humidity for sunny outlook

| Humidity | Yes | No | Number of instances |
|----------|-----|----|---------------------|
| High     | 0   | 3  | 3                   |
| Normal   | 2   | 0  | 2                   |

$$\text{Gini}(\text{Outlook}=\text{Sunny and Humidity}=\text{High}) = 1 - (0/3)^2 - (3/3)^2 = 0$$

$$\text{Gini}(\text{Outlook}=\text{Sunny and Humidity}=\text{Normal}) = 1 - (2/2)^2 - (0/2)^2 = 0$$

$$\text{Gini}(\text{Outlook}=\text{Sunny and Humidity}) = (3/5) \times 0 + (2/5) \times 0 = 0$$

## Gini of wind for sunny outlook

| Wind   | Yes | No | Number of instances |
|--------|-----|----|---------------------|
| Weak   | 1   | 2  | 3                   |
| Strong | 1   | 1  | 2                   |

$$\text{Gini}(\text{Outlook}=\text{Sunny and Wind}=\text{Weak}) = 1 - (1/3)^2 - (2/3)^2 = 0.266$$

$$\text{Gini}(\text{Outlook}=\text{Sunny and Wind}=\text{Strong}) = 1 - (1/2)^2 - (1/2)^2 = 0.2$$

$$\text{Gini}(\text{Outlook}=\text{Sunny and Wind}) = (3/5) \times 0.266 + (2/5) \times 0.2 = 0.466$$

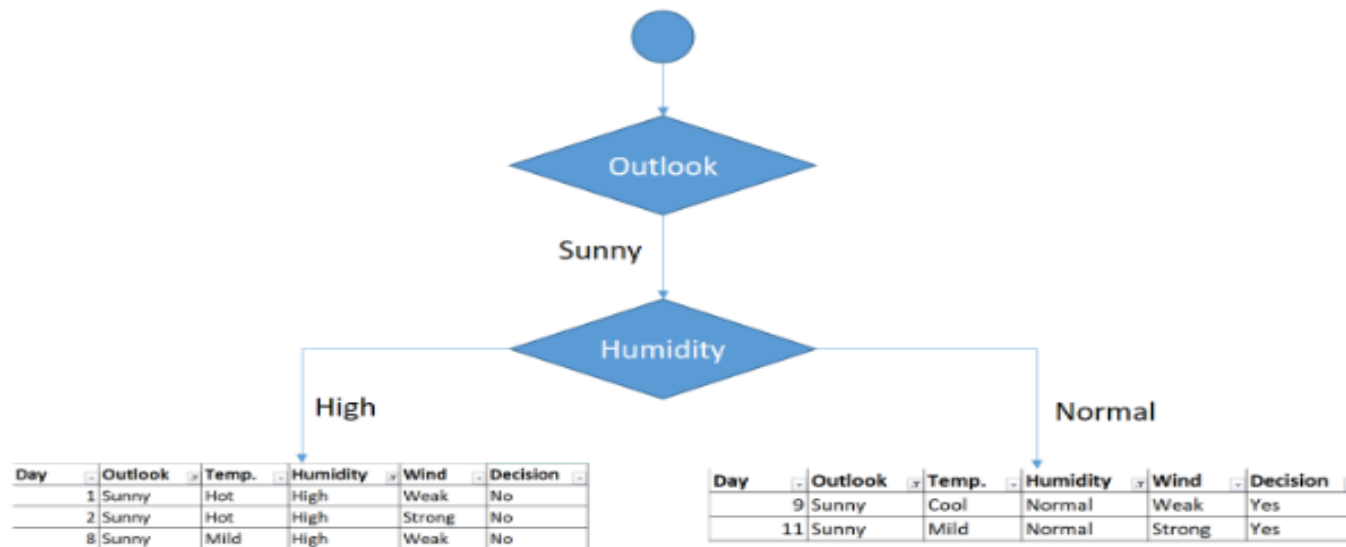


## Decision for sunny outlook

We've calculated gini index scores for feature when outlook is sunny. The winner is humidity because it has the lowest value.

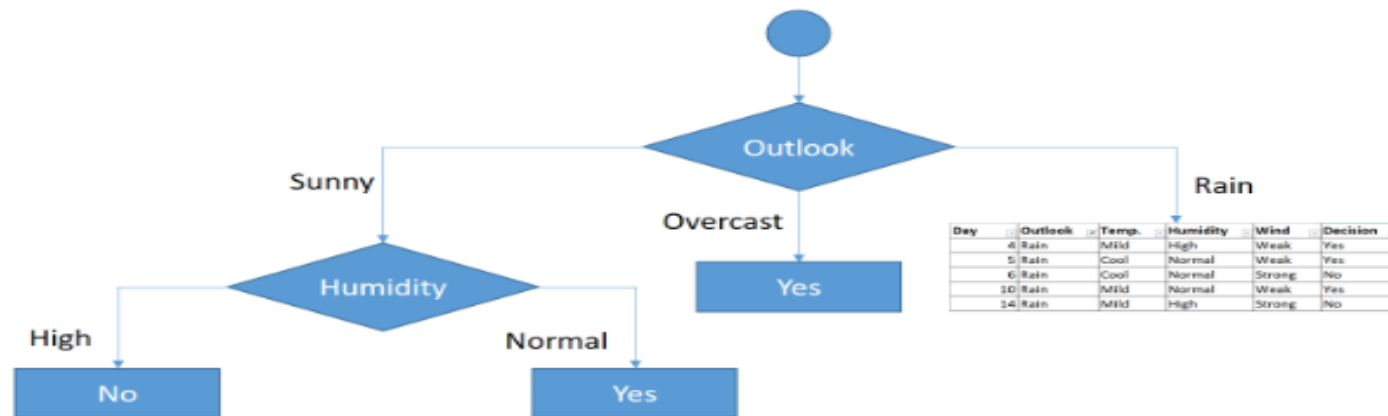
| Feature     | Gini index |
|-------------|------------|
| Temperature | 0.2        |
| Humidity    | 0          |
| Wind        | 0.466      |

We'll put humidity check at the extension of sunny outlook.



Sub datasets for high and normal humidity

As seen, decision is always no for high humidity and sunny outlook. On the other hand, decision will always be yes for normal humidity and sunny outlook. This branch is over.



Decisions for high and normal humidity

Now, we need to focus on rain outlook.

## Rain outlook

| Day | Outlook | Temp. | Humidity | Wind   | Decision |
|-----|---------|-------|----------|--------|----------|
| 4   | Rain    | Mild  | High     | Weak   | Yes      |
| 5   | Rain    | Cool  | Normal   | Weak   | Yes      |
| 6   | Rain    | Cool  | Normal   | Strong | No       |
| 10  | Rain    | Mild  | Normal   | Weak   | Yes      |
| 14  | Rain    | Mild  | High     | Strong | No       |

We'll calculate gini index scores for temperature, humidity and wind features when outlook is rain.

Now, we need to focus on rain outlook.

## Rain outlook

| Day | Outlook | Temp. | Humidity | Wind   | Decision |
|-----|---------|-------|----------|--------|----------|
| 4   | Rain    | Mild  | High     | Weak   | Yes      |
| 5   | Rain    | Cool  | Normal   | Weak   | Yes      |
| 6   | Rain    | Cool  | Normal   | Strong | No       |
| 10  | Rain    | Mild  | Normal   | Weak   | Yes      |
| 14  | Rain    | Mild  | High     | Strong | No       |

We'll calculate gini index scores for temperature, humidity and wind features when outlook is rain.

## Gini of temprature for rain outlook

| Temperature | Yes | No | Number of instances |
|-------------|-----|----|---------------------|
| Cool        | 1   | 1  | 2                   |
| Mild        | 2   | 1  | 3                   |

$$\text{Gini}(\text{Outlook}=\text{Rain and Temp.}=\text{Cool}) = 1 - (1/2)^2 - (1/2)^2 = 0.5$$

$$\text{Gini}(\text{Outlook}=\text{Rain and Temp.}=\text{Mild}) = 1 - (2/3)^2 - (1/3)^2 = 0.444$$

$$\text{Gini}(\text{Outlook}=\text{Rain and Temp.}) = (2/5) \times 0.5 + (3/5) \times 0.444 = 0.466$$

## Gini of humidity for rain outlook

| Humidity | Yes | No | Number of instances |
|----------|-----|----|---------------------|
| High     | 1   | 1  | 2                   |
| Normal   | 2   | 1  | 3                   |

$$\text{Gini}(\text{Outlook}=\text{Rain and Humidity}=\text{High}) = 1 - (1/2)^2 - (1/2)^2 = 0.5$$

$$\text{Gini}(\text{Outlook}=\text{Rain and Humidity}=\text{Normal}) = 1 - (2/3)^2 - (1/3)^2 = 0.444$$

$$\text{Gini}(\text{Outlook}=\text{Rain and Humidity}) = (2/5) \times 0.5 + (3/5) \times 0.444 = 0.466$$

## Gini of wind for rain outlook

| Wind   | Yes | No | Number of instances |
|--------|-----|----|---------------------|
| Weak   | 3   | 0  | 3                   |
| Strong | 0   | 2  | 2                   |

$$\text{Gini}(\text{Outlook}=\text{Rain and Wind}=\text{Weak}) = 1 - (3/3)^2 - (0/3)^2 = 0$$

$$\text{Gini}(\text{Outlook}=\text{Rain and Wind}=\text{Strong}) = 1 - (0/2)^2 - (2/2)^2 = 0$$

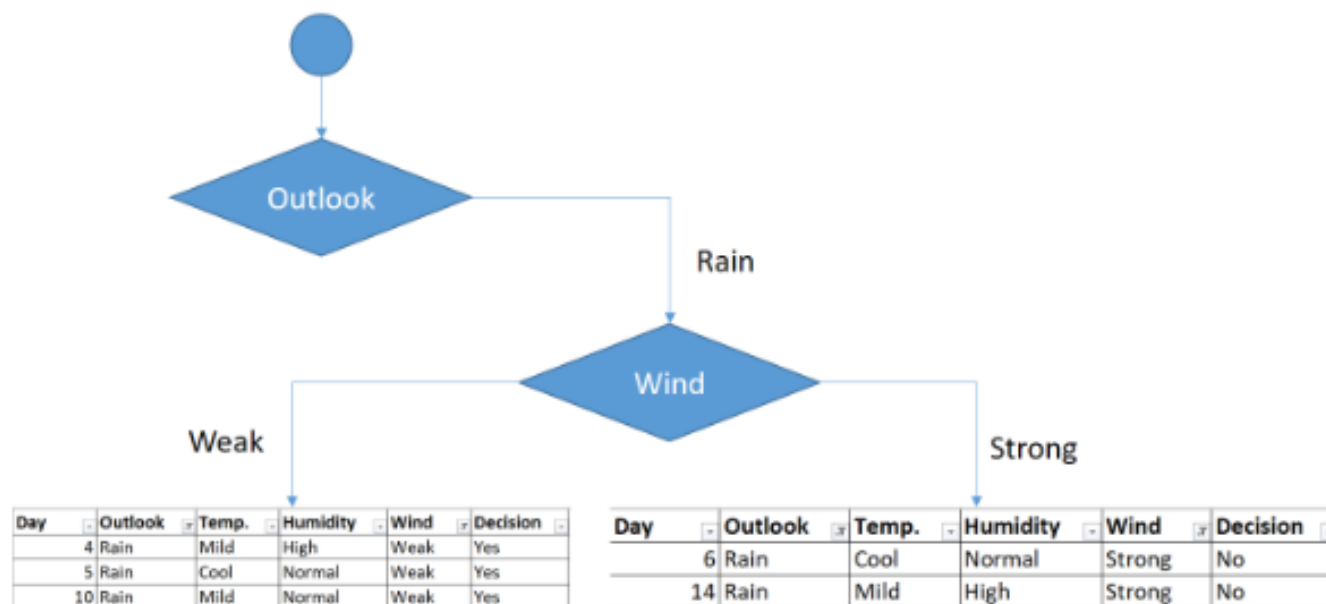
$$\text{Gini}(\text{Outlook}=\text{Rain and Wind}) = (3/5) \times 0 + (2/5) \times 0 = 0$$

## Decision for rain outlook

The winner is wind feature for rain outlook because it has the minimum gini index score in features.

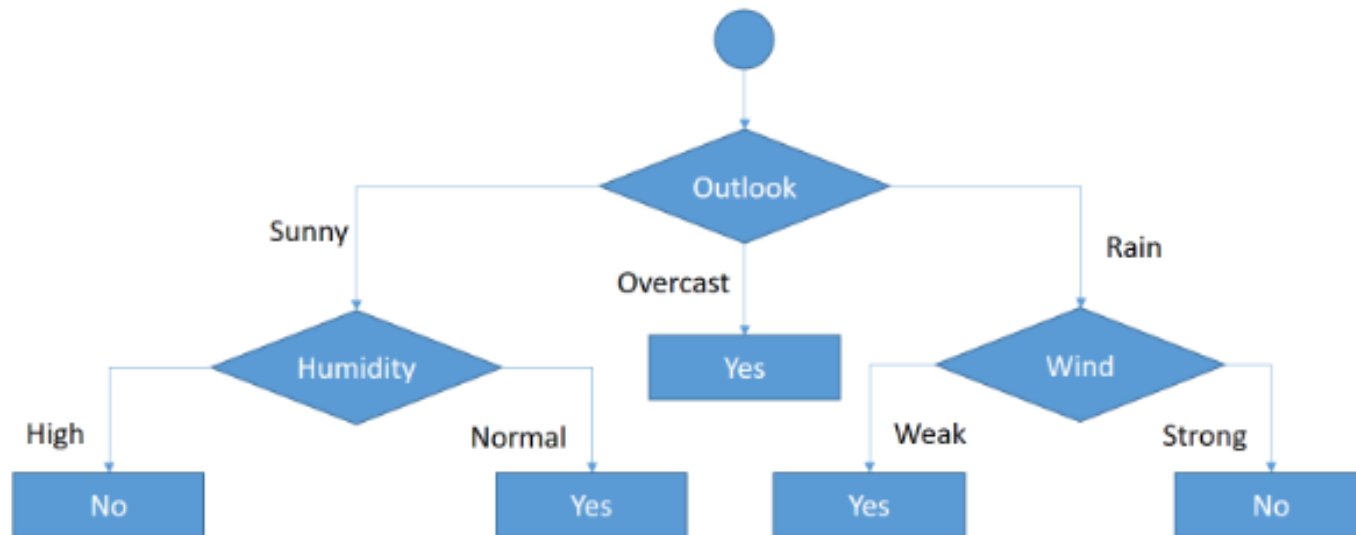
| Feature     | Gini index |
|-------------|------------|
| Temperature | 0.466      |
| Humidity    | 0.466      |
| Wind        | 0          |

Put the wind feature for rain outlook branch and monitor the new sub data sets.



Sub data sets for weak and strong wind and rain outlook

As seen, decision is always yes when wind is weak. On the other hand, decision is always no if wind is strong. This means that this branch is over.



Final form of the decision tree built by CART algorithm