

Decision Trees

AI / ML Workshop
NCIT College

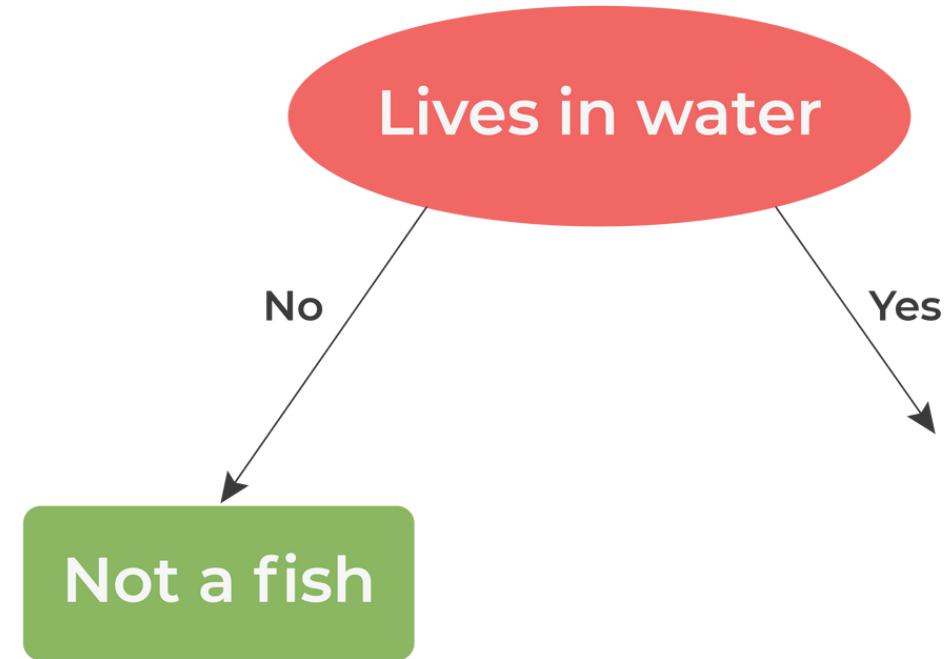
Intuition

- Tree of Decisions (i.e. Every Node represents some decisions)
- Divide and Conquer Approach
- Supervised algorithm for both classification and regression

Fish

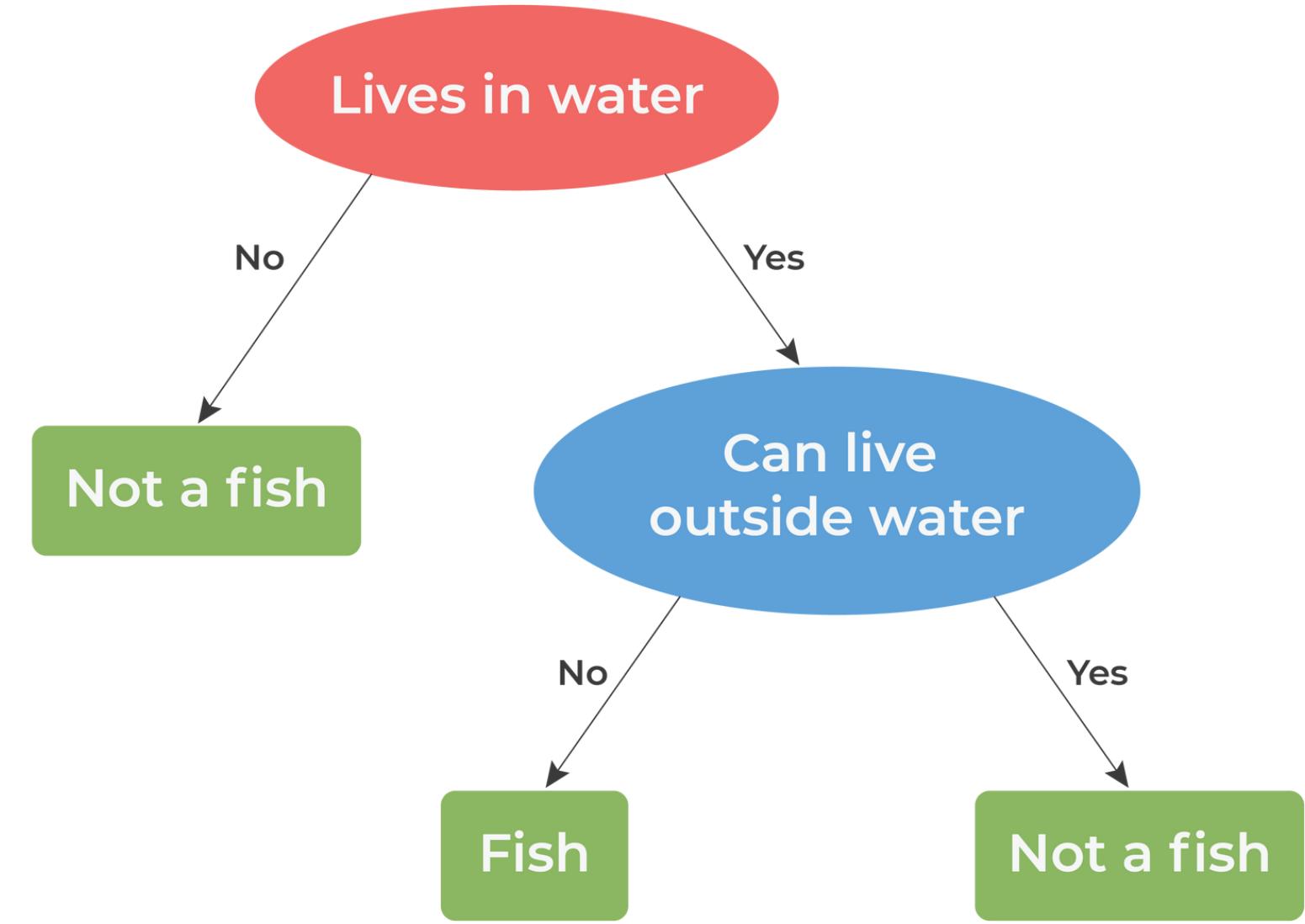
Not Fish

Lives in water

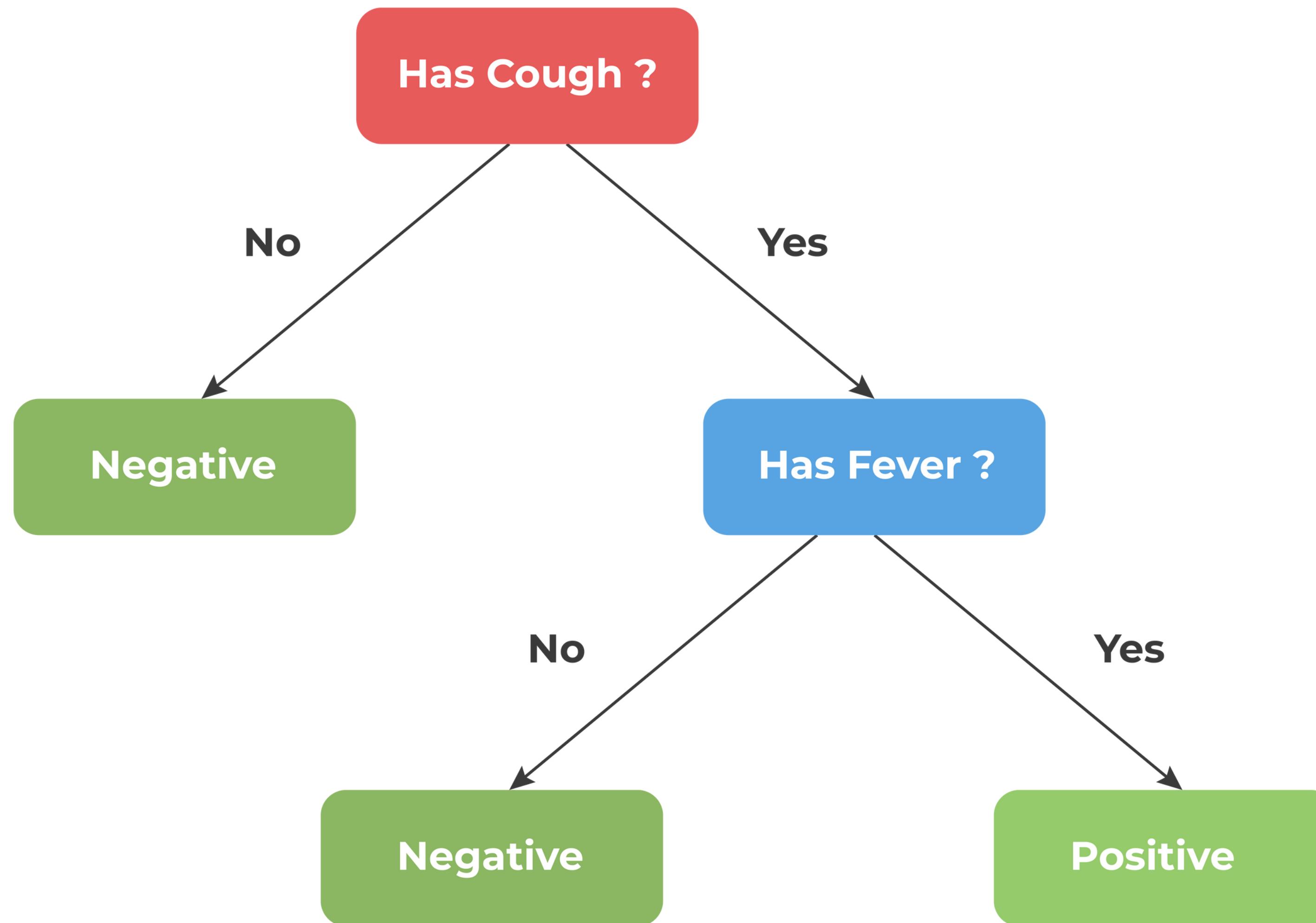


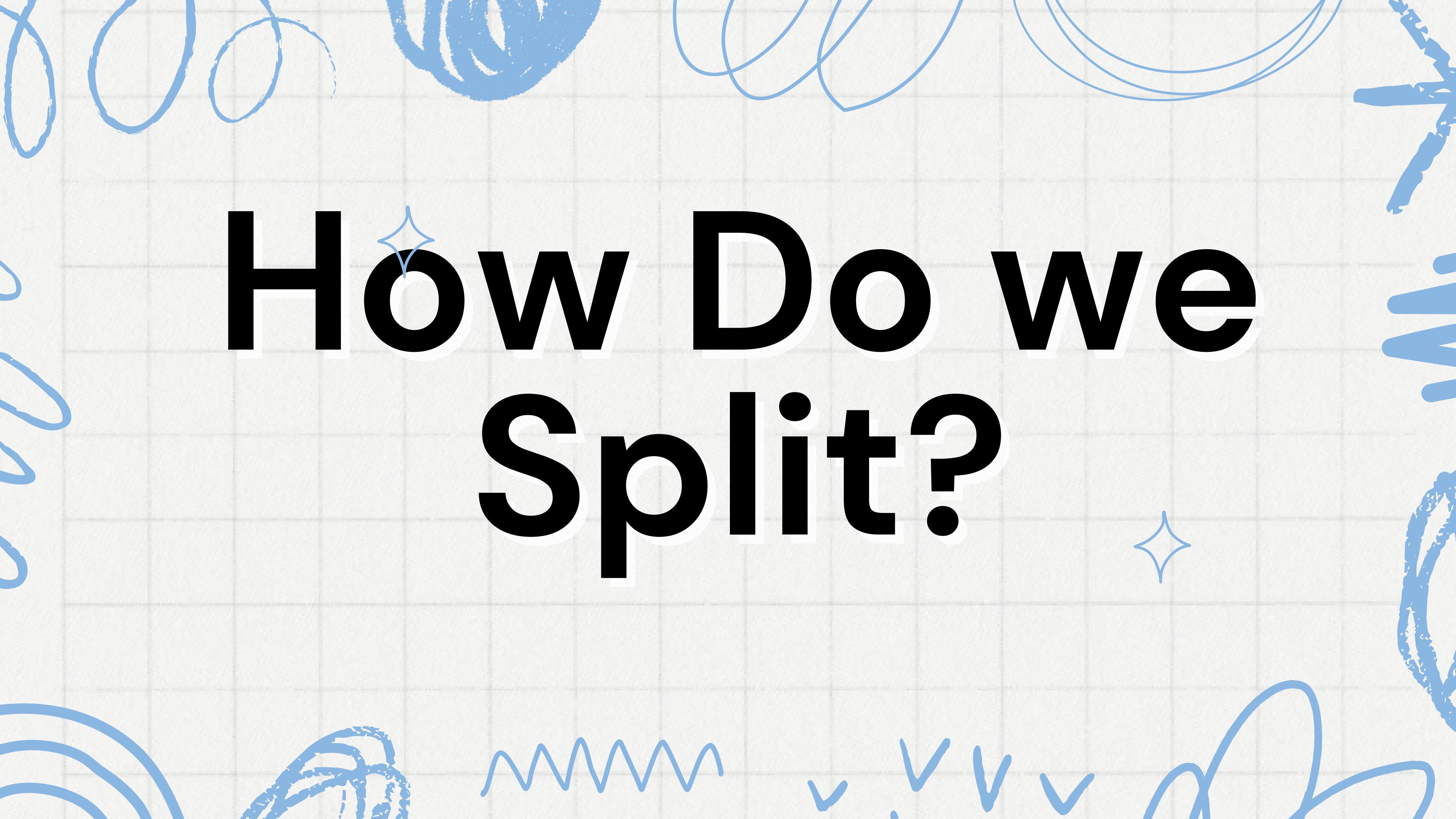
a

b



c





How Do we Split?

Impurity Metrics

Gini

$$\text{Gini} = 1 - \sum_{i=1}^n P_i^2$$

- Probability of Misclassifying a Data point

Entropy

It is a measure of disorder or randomness in data.

$$\text{Entropy} = - \sum_{i=1}^c P_i \log P_i$$

Here,

c = Total number of labels/classes.

p_i = Probability of an item belonging to class i .

Information Gain

Information gain is the decrease in entropy after making a split.

$$\text{Information Gain}(G) = \text{Entropy}(L) - \text{Entropy}(L, A)$$

where,

$\text{Entropy}(L)$ = Entropy of parent node

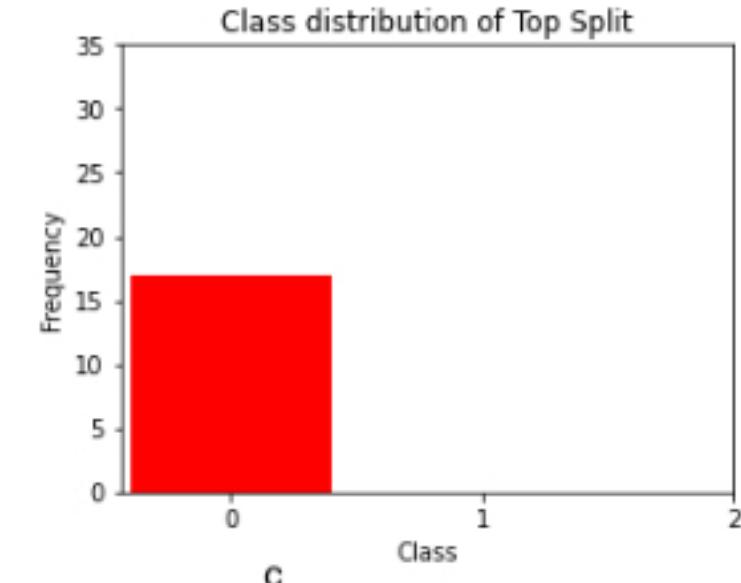
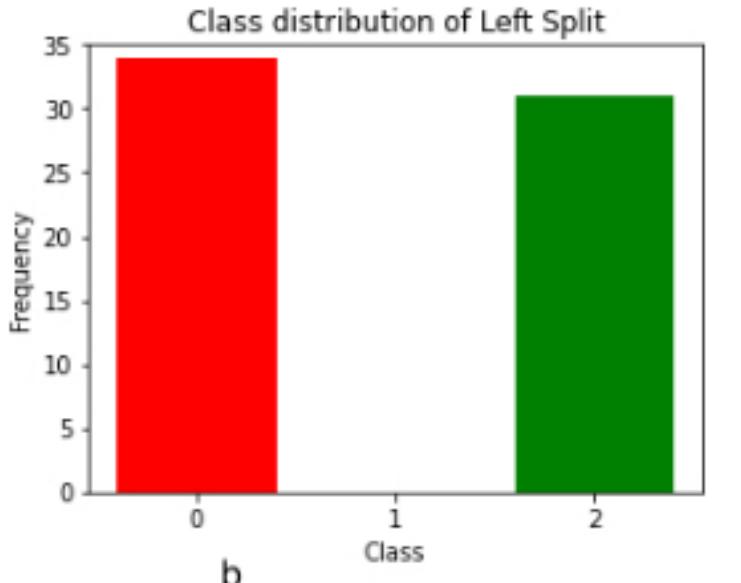
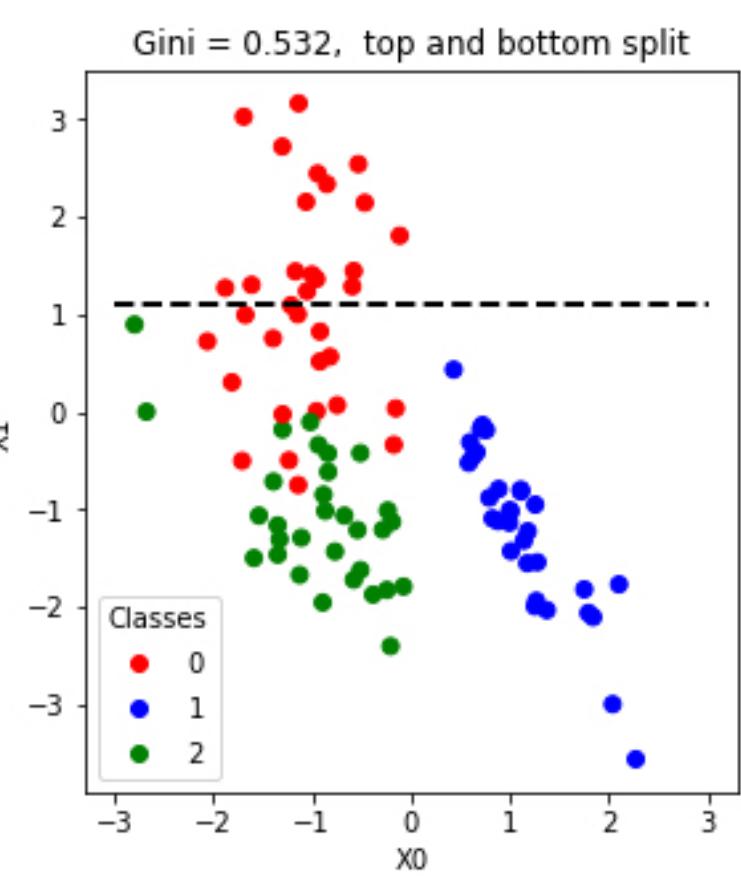
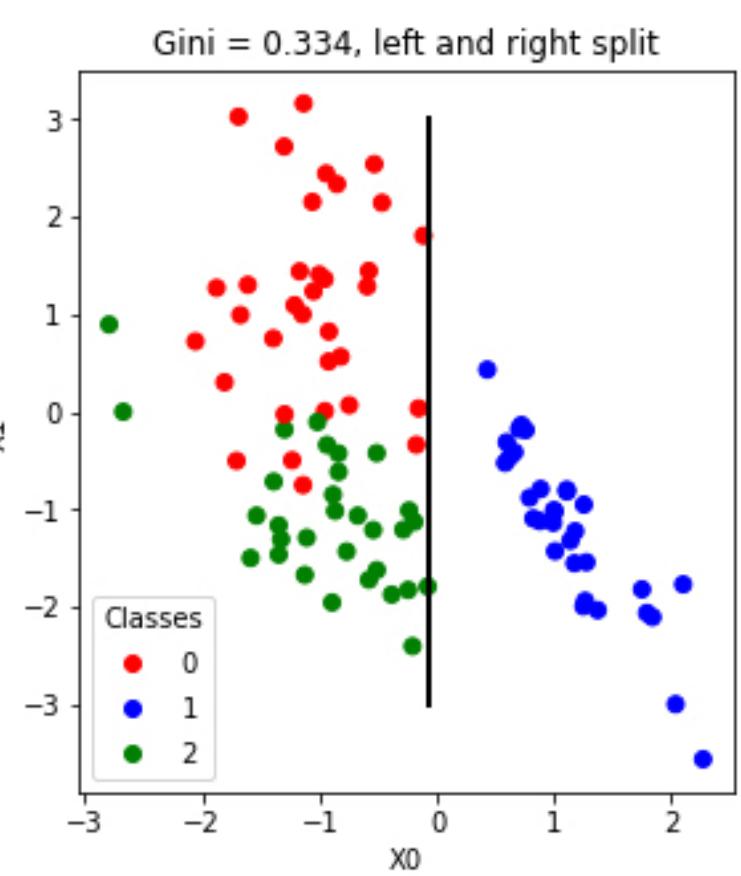
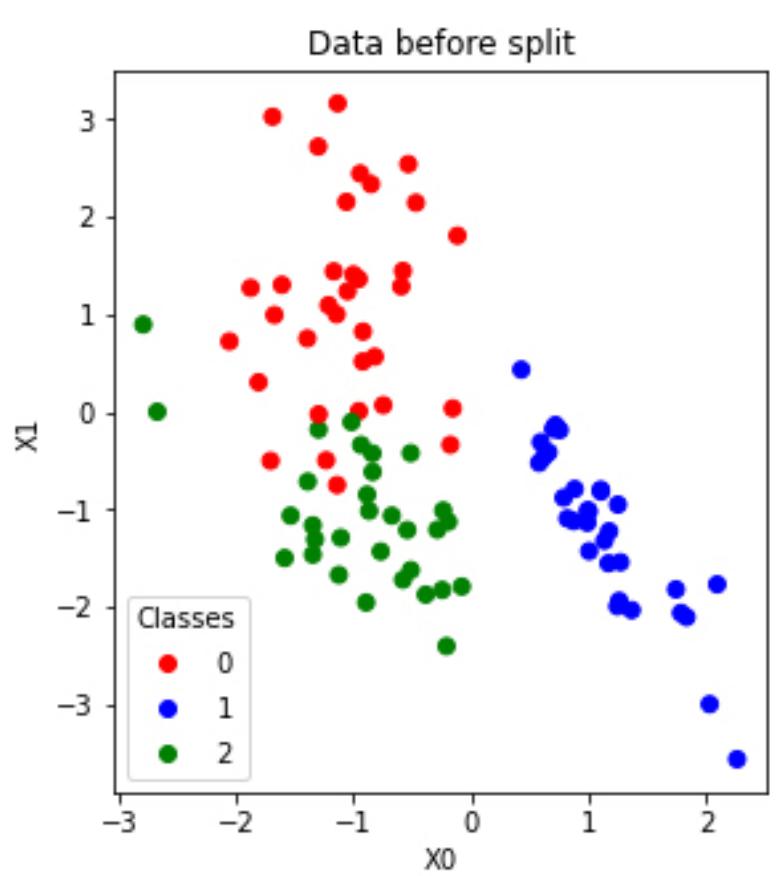
$\text{Entropy}(L, A)$ = Entropy of child nodes after splitting parent node using attribute A

L = Label of the dataset

Gain Ratio

Mostly Used in Decision Tree Problems

$$\text{Gain Ratio}(L, A) = \frac{\text{Information Gain}(G)}{\text{Split Entropy}(L, A)}$$



a

b

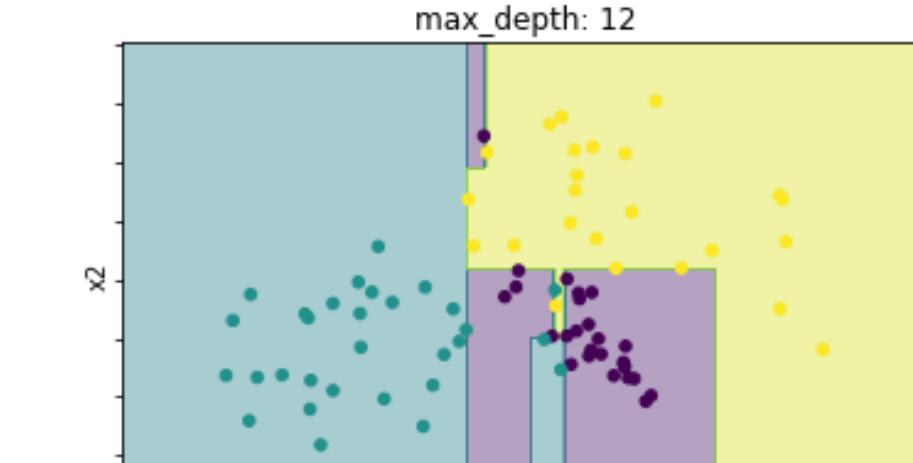
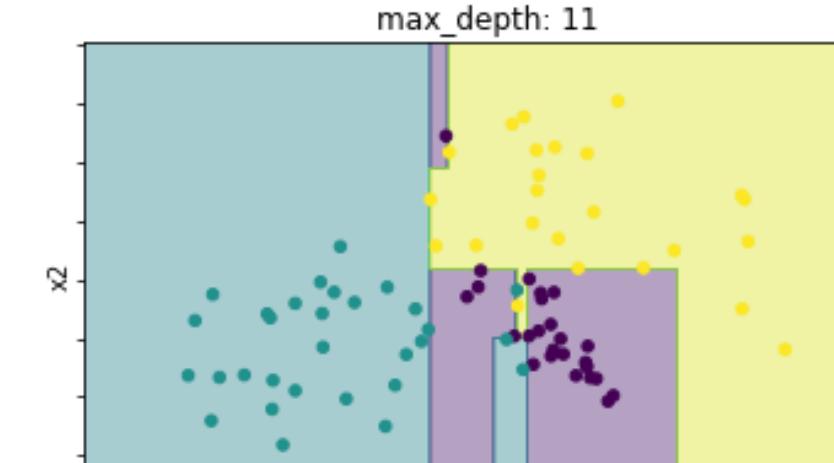
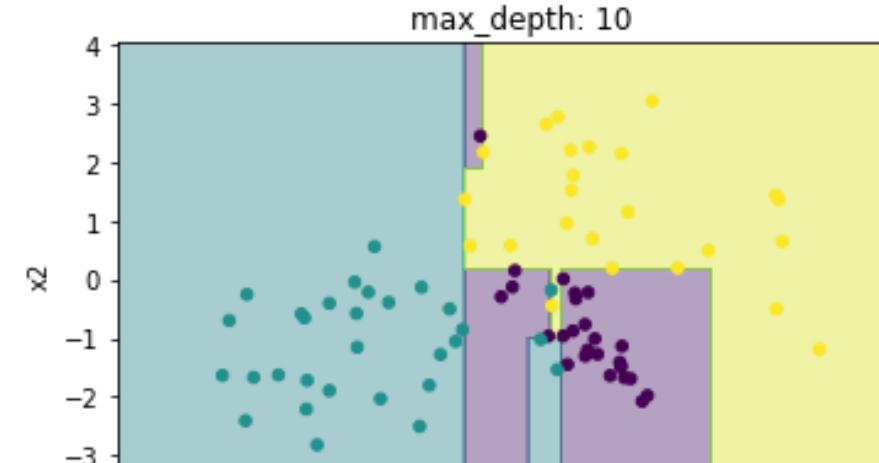
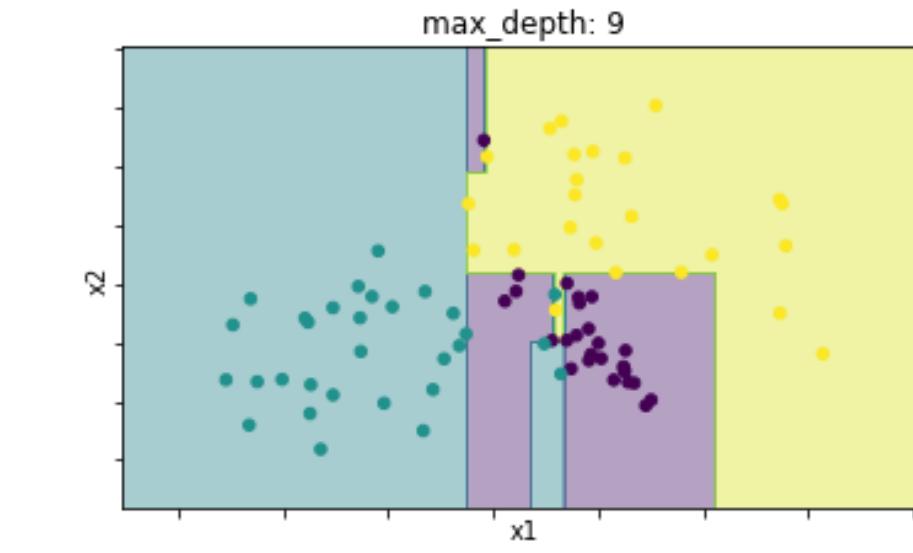
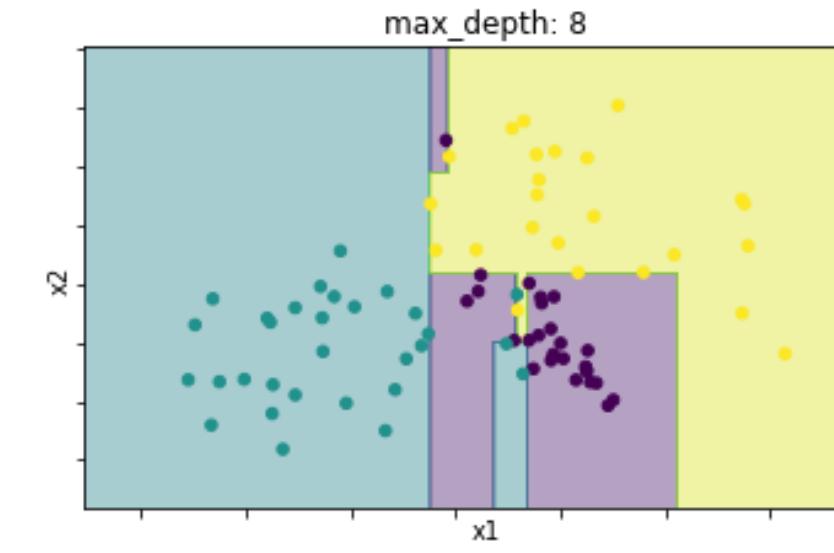
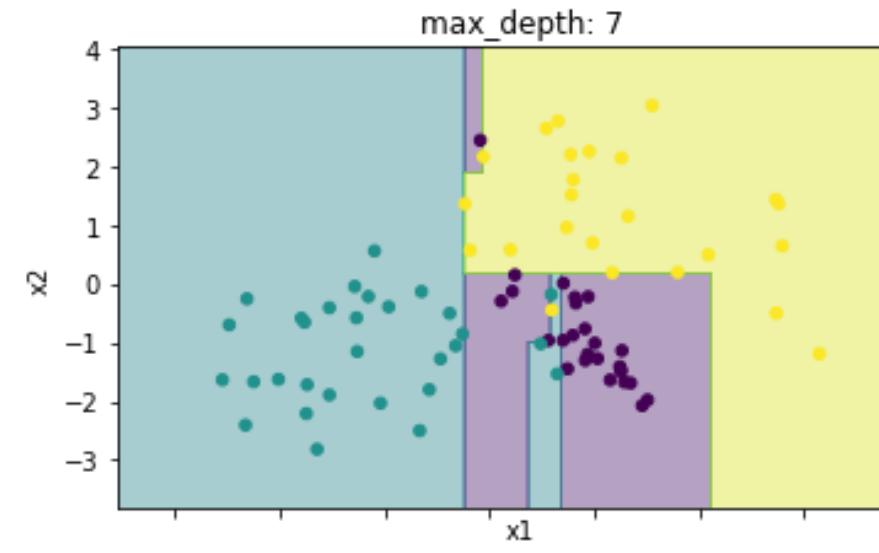
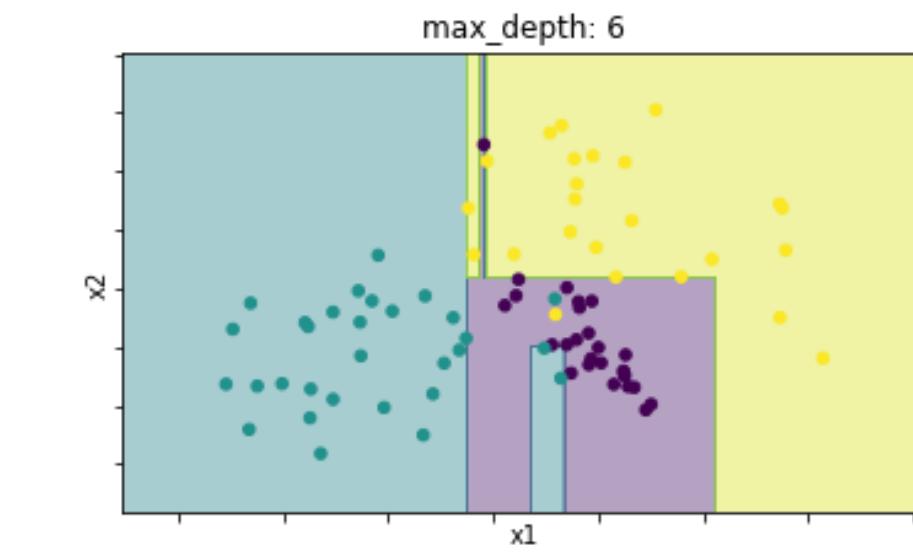
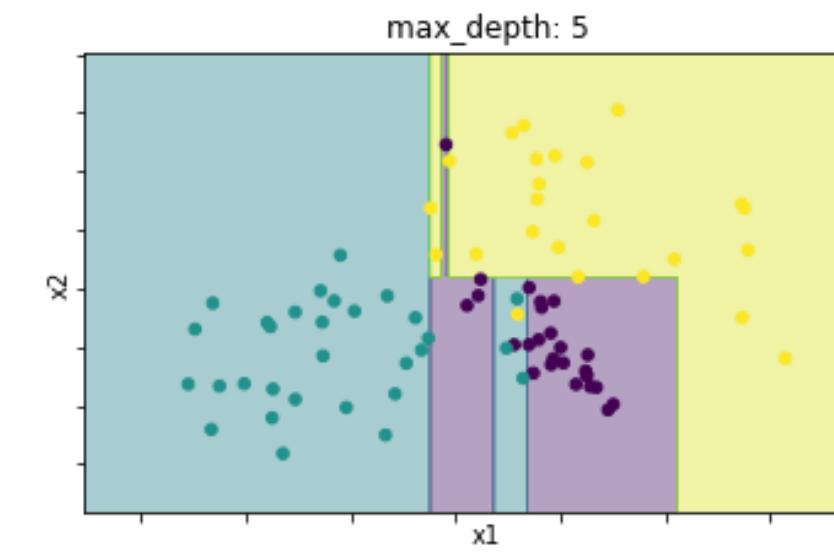
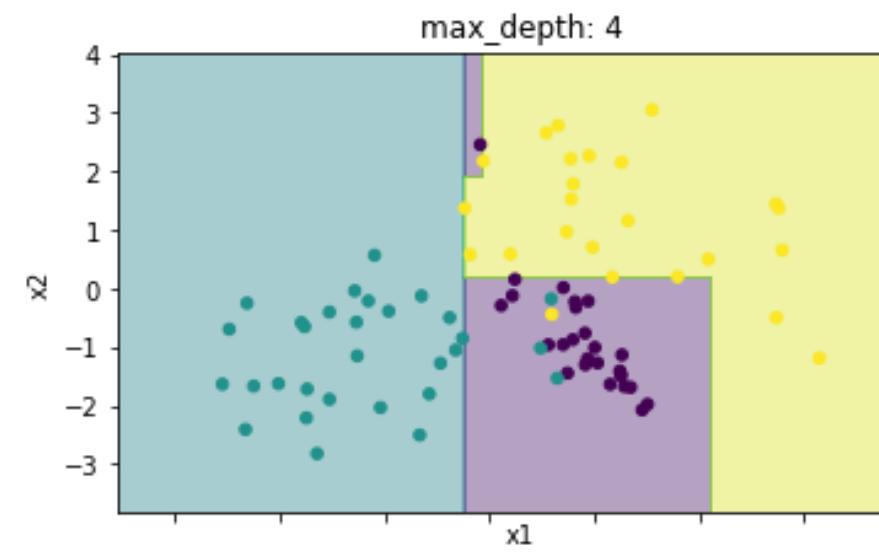
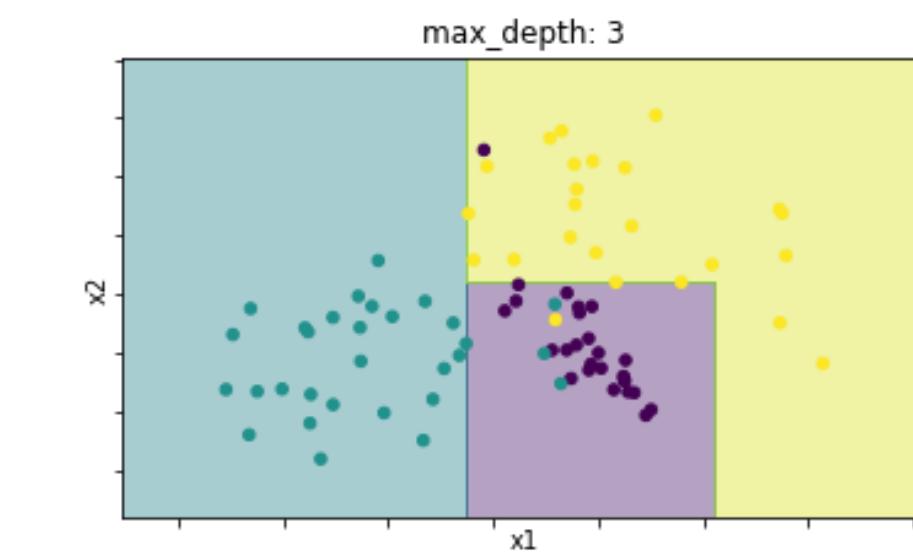
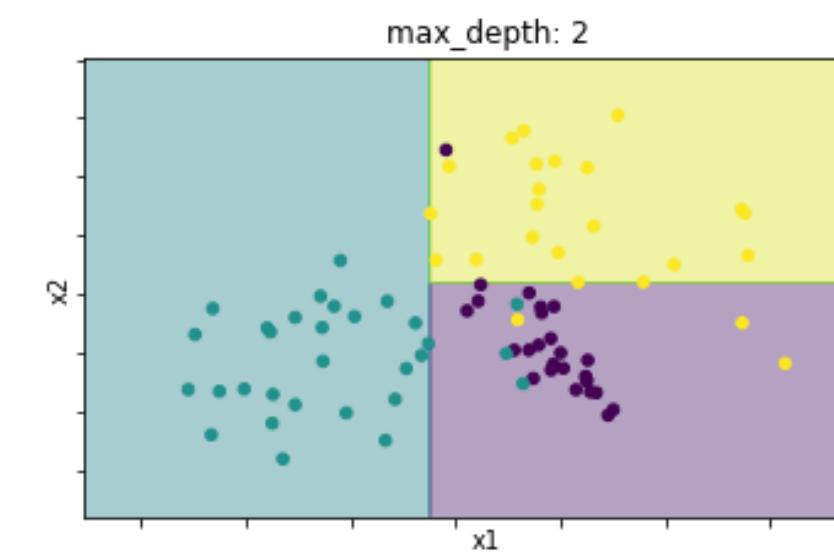
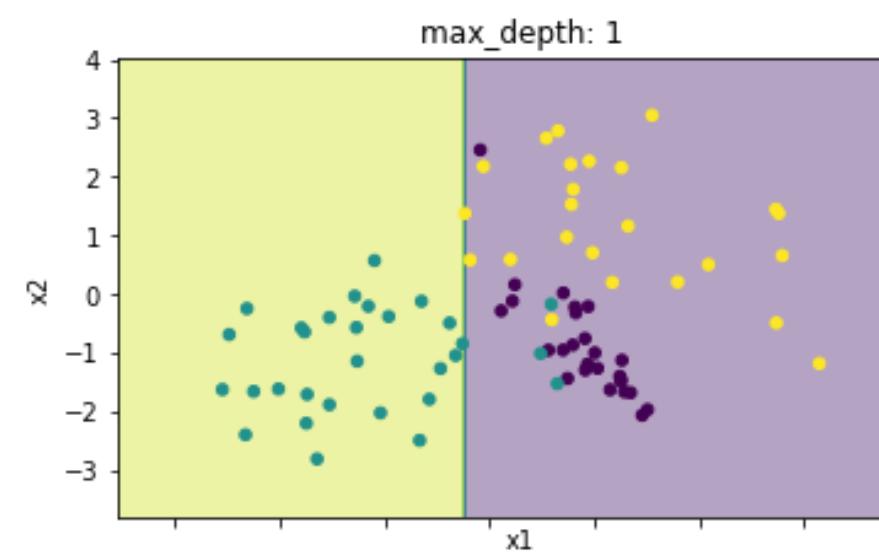
c

Algorithm

- If
 - all records belong to same class. Label as leaf Node
- else
 - Calculate **Impurity metrics** for each attribute i.e **Gini, Entropy, MSE**
 - Select the **best attribute**
 - Split data into subsets
 - Repeat



Should We Worry about Overfitting the Data?



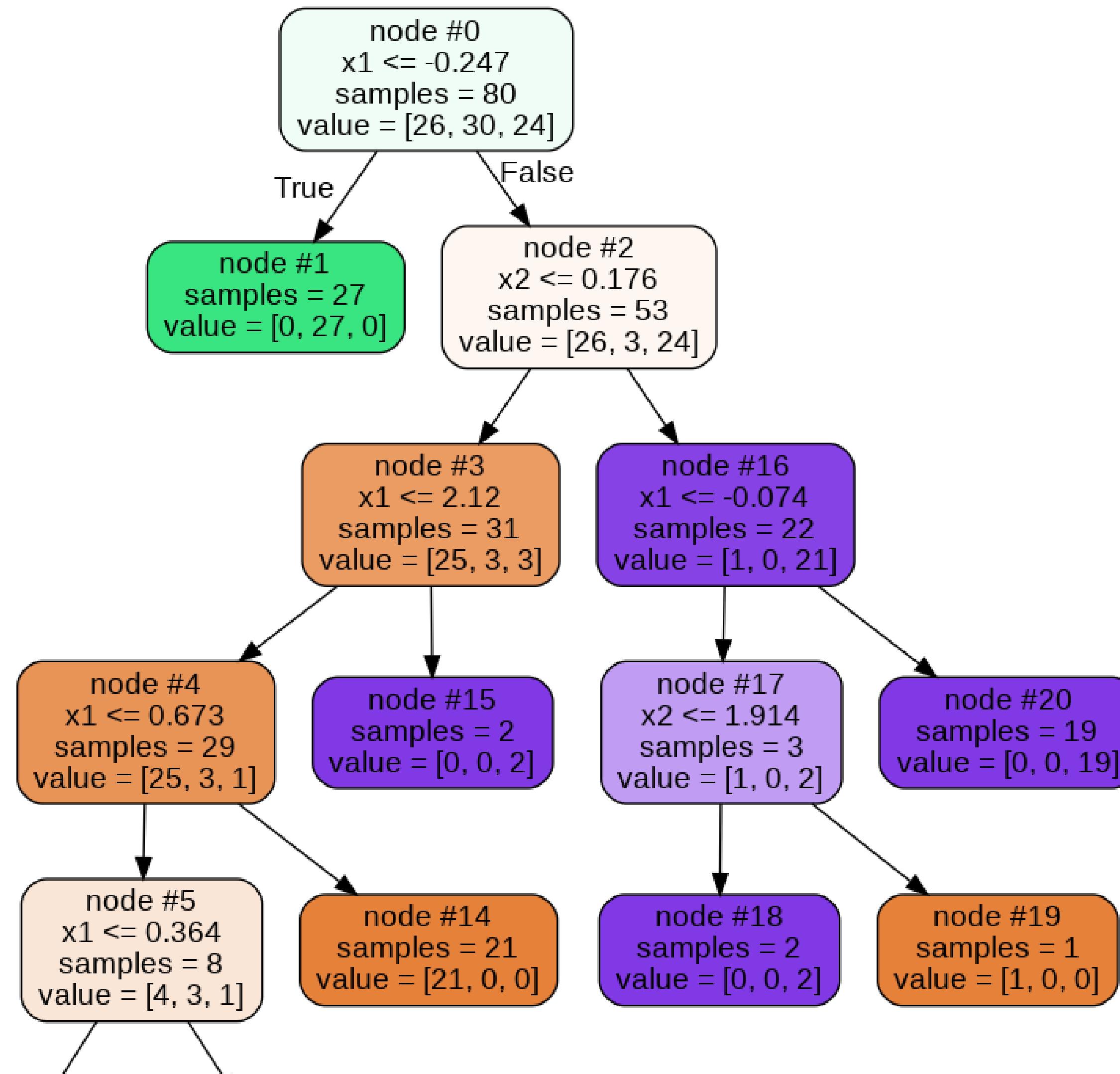
OCCAM'S RAZOR PRiNCiPLE

If there exist multiple solutions to a problem, the simplest one is usually the best.

TECHNIQUES TO OVERCOME

Early Stopping

Pruning



**Thank you
very much!**