# Decision Tree

Abdul Haseeb

BS(AI)-III

# Outline

- What is Decision Tree?

- Terminologies related to Decision Tree

- Different Splitting criterion for Decision Tree

- Pros/Cons of Decision Tree

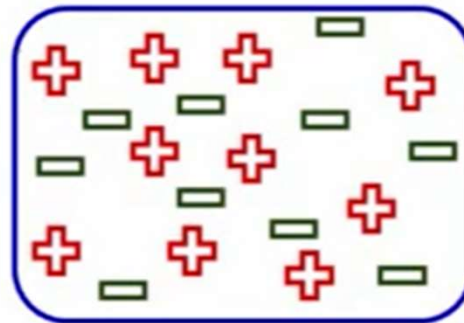- Implementation of Decision Tree

# What is Decision Tree

- Decision Tree is a supervised Learning Algorithm which uses tree like structure to classify data or make predictions

- **Characteristics:**
  - Tree Structure
  - Supervised Learning

# What is Decision Tree

- **Types:**
  - Classification Trees
  - Regression Trees

- **Applications:**
  - **Classification:**
    - Spam classification, Image Classification
  - **Regression:**
    - Predicting stock prices
  - **Feature Selection:**
    - Identifying relevant features

# Decision Tree

- Height

- Performance in class

- Class

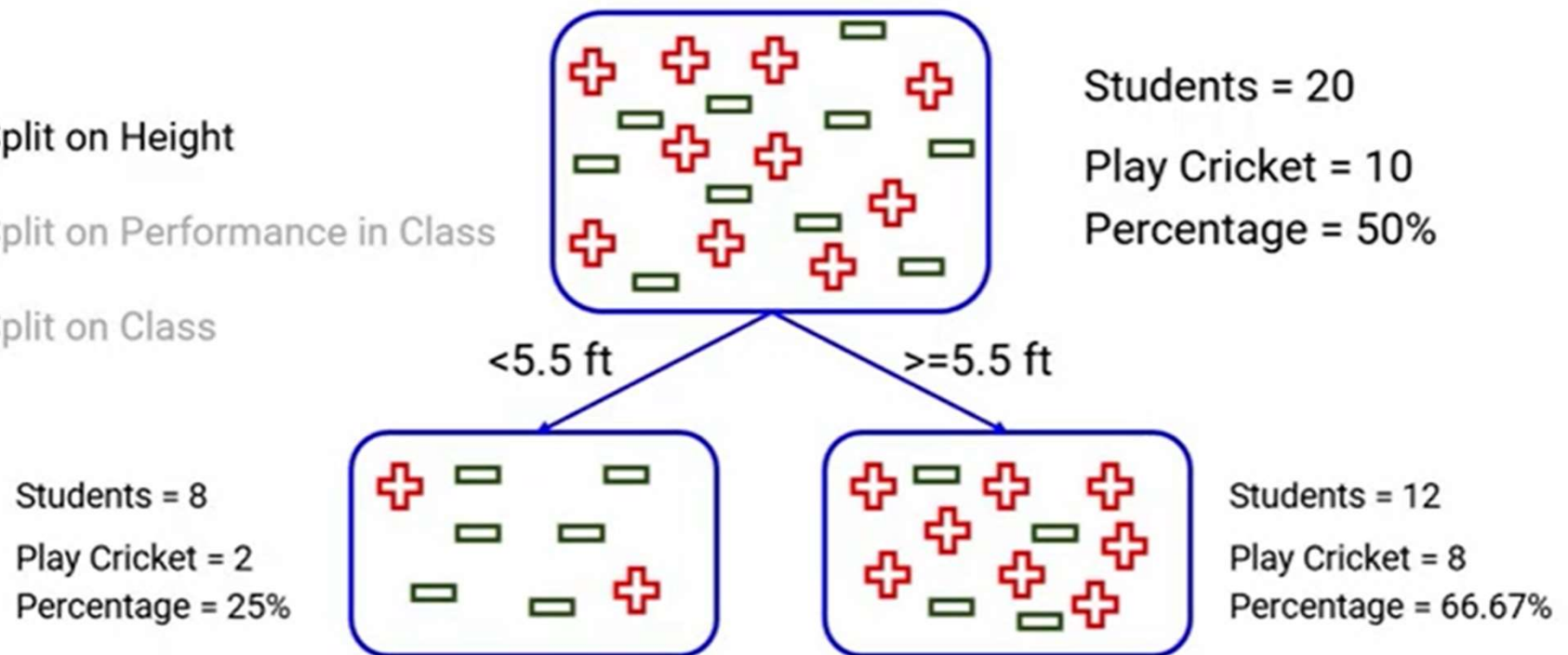

Total number of students = 20

Play cricket = 10

Do not play cricket = 10

# Decision Tree

- **Teacher** wants to identify subgroups, that these subgroups are very much familiar with playing/not playing the cricket with the help of given attributes
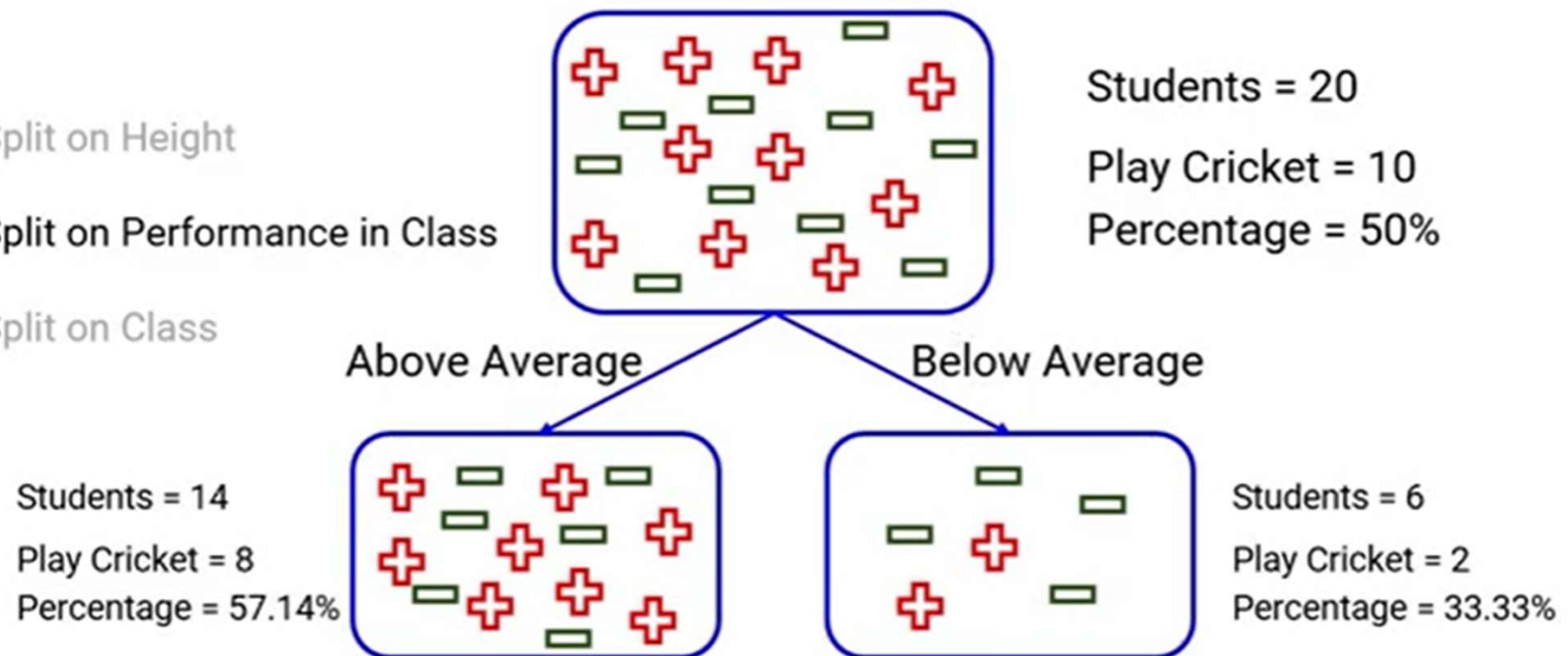
# What is Decision Tree



- Split on Height
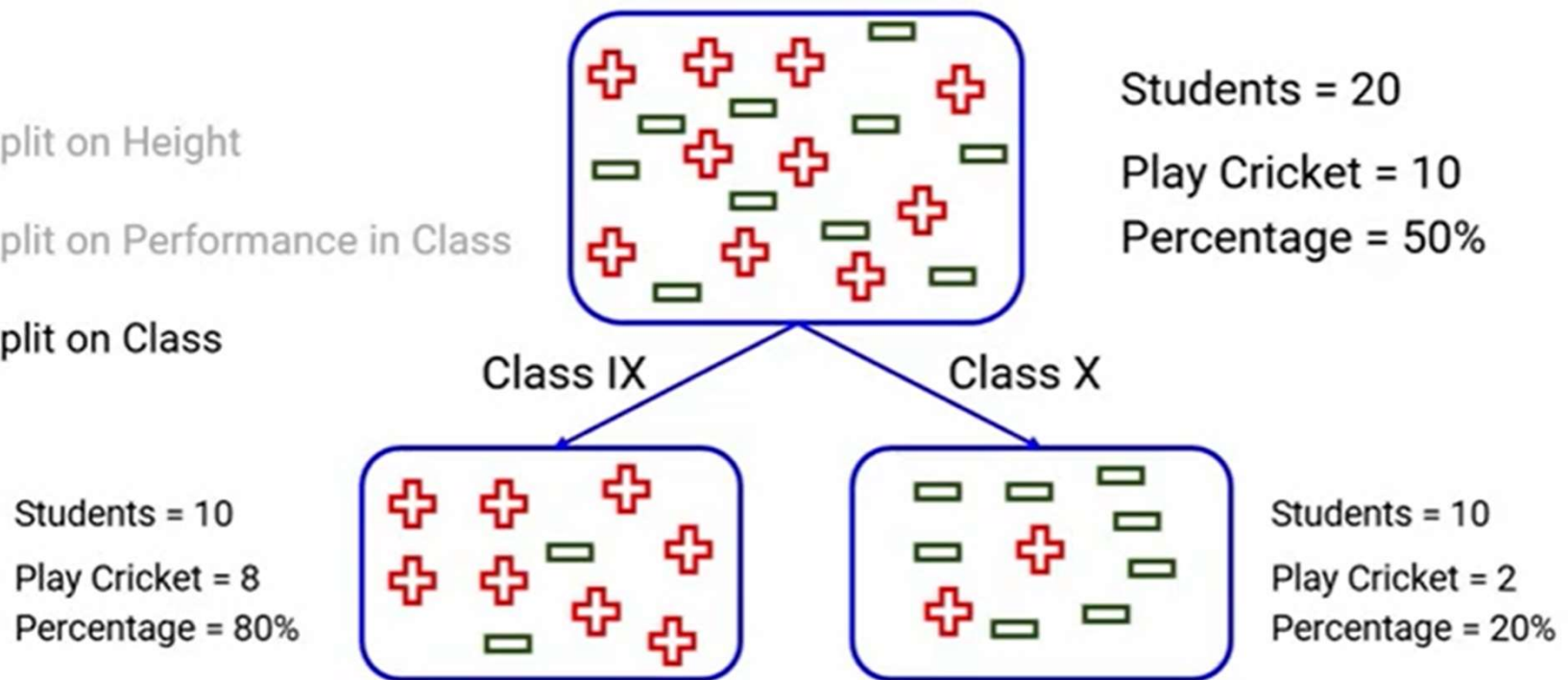- Split on Performance in Class
- Split on Class

Students = 20

Play Cricket = 10

Percentage = 50%

<5.5 ft        >=5.5 ft

Students = 8
Play Cricket = 2
Percentage = 25%

Students = 12
Play Cricket = 8
Percentage = 66.67%

- Split on Height
- Split on Performance in Class
- **Split on Class**

Students = 20

Play Cricket = 10

Percentage = 50%

Class IX

Class X

Students = 10

Play Cricket = 8

Percentage = 80%

Students = 10

Play Cricket = 2

Percentage = 20%

# What is Decision Tree



<5.5 ft  >=5.5 ft

25%  66.67%

Split on Height

Above Average  Below Average

57.14%  33.33%

Split on Performance in Class

Class IX  Class X
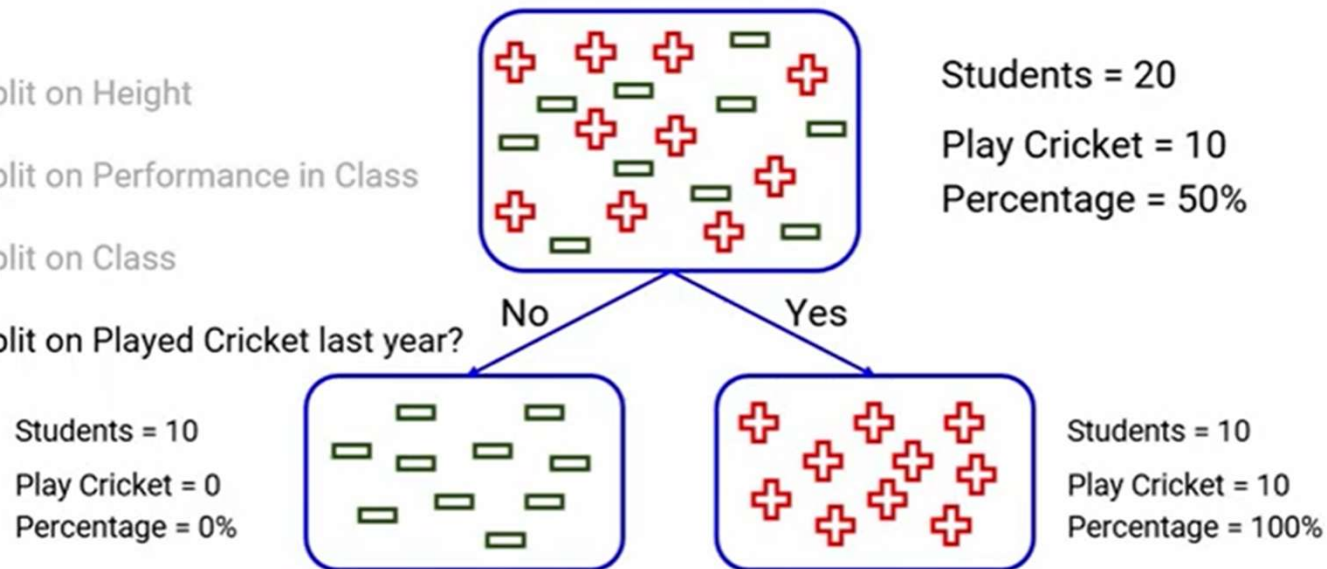
80%  20%

Split on Class

Look at the split on class:

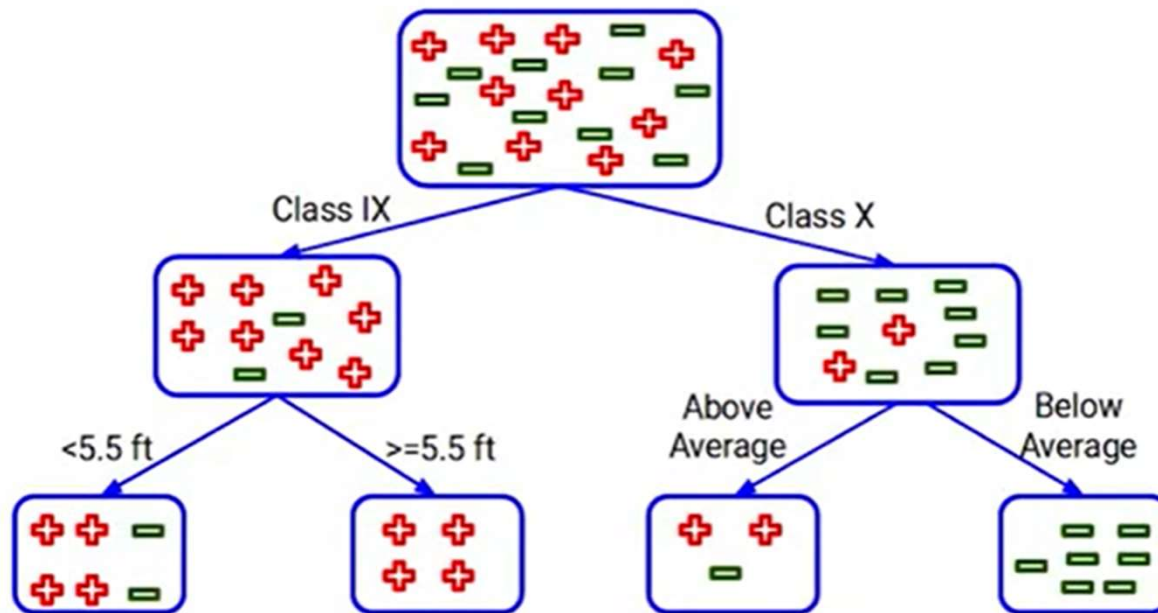It looks the best split as it segregates the most of students

# Purity in Decision Tree

- Split on Height

- Split on Performance in Class

- Split on Class

- **Split on Played Cricket last year?**

Students = 20

Play Cricket = 10

Percentage = 50%

No / Yes

Students = 10

Play Cricket = 0

Percentage = 0%

Students = 10

Play Cricket = 10

Percentage = 100%

Objective of Decision Tree is to produce pure nodes

In practical scenario we rarely will have such features which can produce best split like this

# Purity in Decision Tree



We will have multiple splits and multiple decisions in a decision tree.
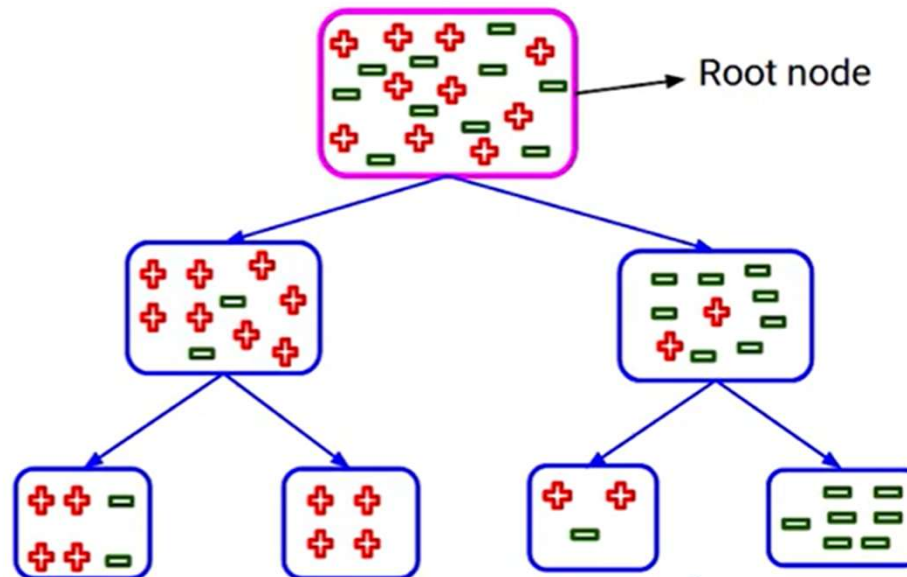
# Purity in Decision Tree

- But wait...
- What should be root node for split?
- What should be the sequence of split?
- How to select the node for split?

# Purity in Decision Tree

- There are techniques to decide purity of nodes: the node which is purest among the others will be taken for the splitting
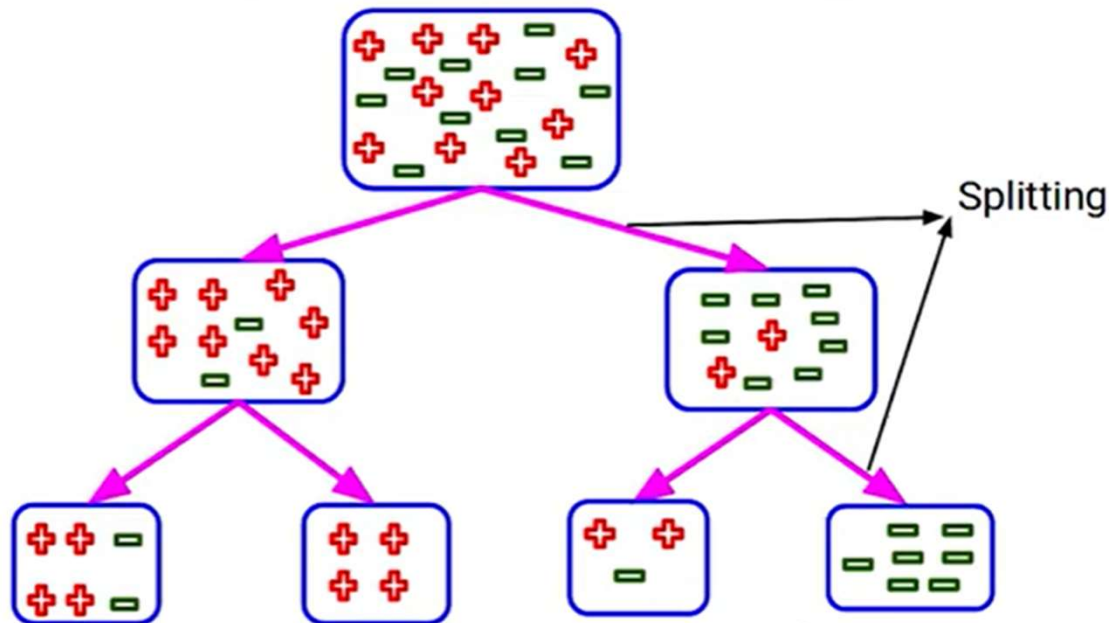
- Root Node



Root node → Root node

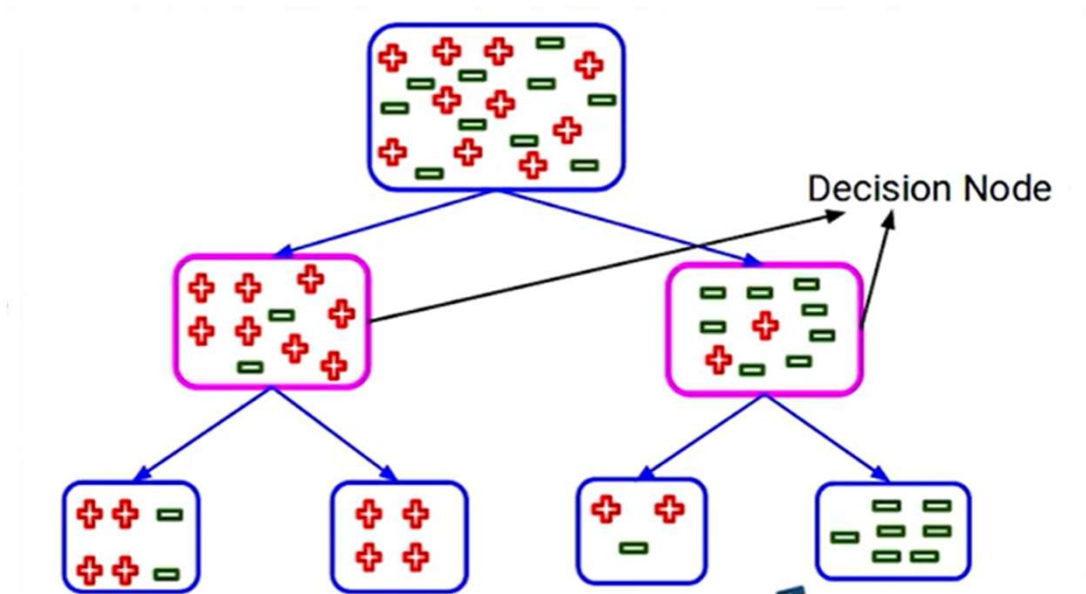Root node describes the entire population

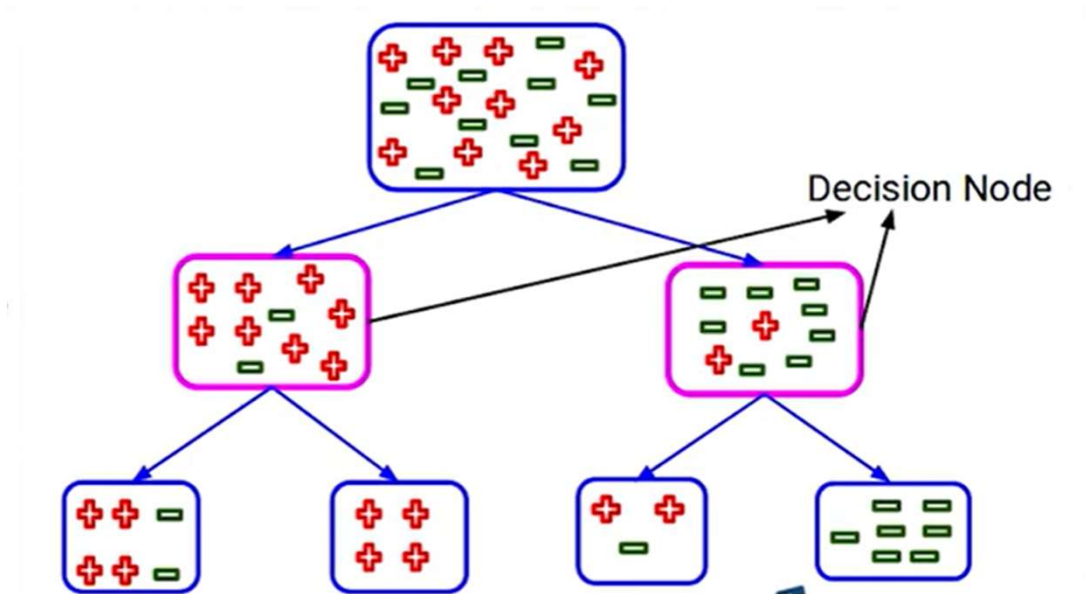- Splitting is dividing a node into further sub nodes

Decision Node

- Decision Nodes are those on which a split is performed

Decision Node

- Leaf nodes don't split further

Decision Node

- Leaf nodes don't split further

- Subtree is a subset/Part of original tree

- Subtree is a subset/Part of original tree

- Parent/Child Nodes

Depth = 2

- Depth is the longest path from root to leaf

Depth = 3

- How many leaf nodes in this tree??

Depth = 3

- How many leaf nodes in this tree??

# How to select the best split point in Decision Tree?

- Decision tree splits the nodes on all available variable

- Select the split which results in most homogenous sub-nodes

- Decision Tree Algorithms measure node Impurity, and two most common techniques for measuring node impurity are:
  - Gini
  - Entropy

# How to select the best split point in Decision Tree?

- Gini Impurity:
    - 1-Gini

- Gini Impurity states:



Same class

If we select two items from a population at random, they
must be of same class

- Gini Impurity:



Probability that randomly picked points belong to same class?

- Gini Impurity:



Probability = 1

Probability is one, as all the samples belong to same class

Node is pure...

Gini ranges from 0-1, Highest the Gini: Highest the Node Purity

- Gini Impurity:



Probability = 1

Probability is one, as all the samples belong to same class

Node is pure...

Gini ranges from 0-1, Highest the Gini: Highest the Node Purity

# Properties of Gini Impurity

- Node split is decided based on the gini impurity

- Lower the gini impurity, higher the homogeneity of the nodes

- Works only with categorical data

- Only performs binary splits

# Steps to calculate Gini Impurity for a split

- Calculate gini impurity for subnode

- Gini=Sum of square probability of each class
  - *Gini=(p1$^2$+p2$^2$+p3….+pn$^2$)*

- To calculate gini impurity of a split, take weighted gini impurity of both sub-nodes of that split

## Split on Performance in Class

- Gini Impurity: sub-node Above Average:
  $1 - [(0.57)*(0.57) + (0.43)*(0.43)] = 0.49$

- Gini Impurity: sub-node Below Average:
  $1 - [(0.33)*(0.33) + (0.67)*(0.67)] = 0.44$

Students = 20
Play Cricket = 10
Percentage = 50%

Above Average

Below Average

Students = 14
Play Cricket = 8
Do not play = 6
Prob. play = 0.57
Prob. Not play = 0.43

Students = 6
Play Cricket = 2
Do not play = 4
Prob. play = 0.33
Prob. Not play = 0.67

## Split on Performance in Class

- Gini Impurity: sub-node Above Average:
  $1 - [(0.57)*(0.57) + (0.43)*(0.43)] = 0.49$

- Gini Impurity: sub-node Below Average:
  $1 - [(0.33)*(0.33) + (0.67)*(0.67)] = 0.44$

- Weighted Gini Impurity: Performance in Class:
  $(14/20)*0.49 + (6/20)*0.44 = 0.475$

Students = 20
Play Cricket = 10
Percentage = 50%

14/20    Above Average          Below Average    6/20

Students = 14
Play Cricket = 8
Do not play = 6
Prob. play = 0.57
Prob. Not play = 0.43

Students = 6
Play Cricket = 2
Do not play = 4
Prob. play = 0.33
Prob. Not play = 0.67

**Split on Class**

- Gini Impurity: sub-node Class IX:
  1 - [(0.8)*(0.8) + (0.2)*(0.2)] = 0.32

- Gini Impurity: sub-node Class X:
  1 - [(0.2)*(0.2) + (0.8)*(0.8)] = 0.32

- Weighted Gini Impurity: Class:
  (10/20)*0.32 + (10/20)*0.32 = 0.32



Students = 20
Play Cricket = 10
Percentage = 50%

10/20    Class IX

Class X    10/20

Students = 10
Play Cricket = 8
Do not play = 2
Prob. play = 0.8
Prob. Not play = 0.2

Students = 10
Play Cricket = 2
Do not play = 8
Prob. play = 0.2
Prob. Not play = 0.8

Weight of the node*Gini Impurity of the node

Students = 20

Play Cricket = 10

Percentage = 50%

10/20   Class IX          Class X   10/20

Calculate the Gini Impurity and Weighted Gini Impurity

# Steps to calculate Gini Impurity for a split
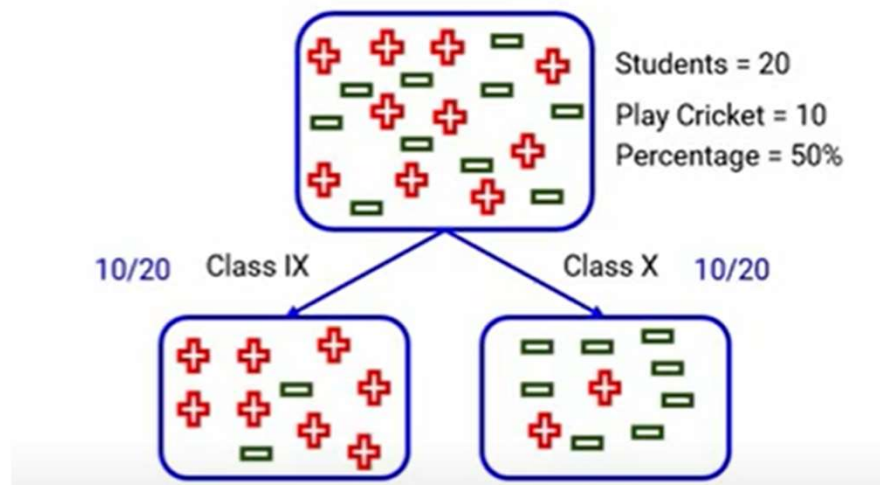
- Gini Impurity: sub-node Class IX:
  1 - [(0.8)*(0.8) + (0.2)*(0.2)] = 0.32

- Gini Impurity: sub-node Class X:
  1 - [(0.2)*(0.2) + (0.8)*(0.8)] = 0.32

- Weighted Gini Impurity: Class:
  (10/20)*0.32 + (10/20)*0.32 = 0.32

# Steps to calculate Gini Impurity for a split

| Split | Weighted Gini Impurity |
|---|---|
| Performance in Class | 0.475 |
| Class | 0.32 |

Node Producing Minimum Weighted Gini Impurity will be selected as the Split

# Another Algorithm for Deciding the Best Split

**Information Gain:**



Node 1        Node 2        Node 3

Which Node will require more explanation?

Which is the purest of the Nodes?

**Information Gain:**



Node 1 > Node 2 > Node 3

Information required to describe the node

**Information Gain:**



What can you infer from this?

**Information Gain:**

- The split on the right is giving less information gain

- So, we can easily say: "***Higher the Information gain Higher the Homogeneity and lesser the impurity***"

**Information Gain:**

- The split on the right is giving less information gain

- So, we can easily say: "***Higher the Information gain Higher the Homogeneity and lesser the impurity***"

# Formula for Information Gain

Information Gain = 1 - Entropy

# Entropy

**Entropy**

$$-p_1 * \log_2 p_1 - p_2 * \log_2 p_2 - p_3 * \log_2 p_3 - \ldots - p_n * \log_2 p_n$$

p refers to percentage of each class in the Node

**Entropy**

$$- p_1 * \log_2 p_1 - p_2 * \log_2 p_2 - p_3 * \log_2 p_3 - ..... - p_n * \log_2 p_n$$



% Play = 0.50
% Not play = 0.50

Entropy = - (0.5) * log2(0.5) - (0.5) * log2(0.5)

= 1

% Play = 0
% Not play = 1

% Play = 0
% Not play = 1

Entropy = - (0) * log2(0) - (1) * log2(1)

= 0

# Entropy



% Play = 0.50
% Not play = 0.50

Entropy = 1

% Play = 0
% Not play = 1

Entropy = 0

Lower the Entropy means?

Higher the Entropy means?

# Properties of Entropy

- Works only with categorical Targets

- Lesser the entropy, higher the homogeneity of Nodes

# Steps to Calculate Entropy of Nodes

- Calculate the entropy of the parent node

- Calculate the entropy of each child node

- Calculate the weighted average entropy of the split

- If weighted entropy of child node is greater than parent node, then we will ignore that node as it is returning more impure node than the parent

# Steps to Calculate Entropy of Nodes



Above Average — Below Average

Split on Performance in Class

Class IX — Class X

Split on Class

**Split on Performance in Class**

- Entropy for Parent node:

$-(0.5)*\log_2(0.5) -(0.5)*\log_2(0.5) = 1$

- Entropy for sub-node Above Average:

$-(0.57)*\log_2(0.57) -(0.43)*\log_2(0.43) = 0.98$

- Entropy for sub-node Below Average:

$-(0.33)*\log_2(0.33) -(0.67)*\log_2(0.67) = 0.91$

- Weighted Entropy: Performance in Class:

$(14/20)*0.98 + (6/20)*0.91 = 0.959$

Students = 20
Play Cricket = 10
Do not play = 10
% play = 0.5
% Not play = 0.5

Above Average

Below Average

Students = 14
Play Cricket = 8
Do not play = 6
% play = 0.57
% Not play = 0.43

Students = 6
Play Cricket = 2
Do not play = 4
% play = 0.33
% Not play = 0.67

# Steps to Calculate Entropy of Nodes



Students = 20
Play Cricket = 10
Do not play = 10
% play = 0.5
% Not play = 0.5

Class IX

Class X

Students = 10
Play Cricket = 8
Do not play = 2
% play = 0.8
% Not play = 0.2

Students = 10
Play Cricket = 2
Do not play = 8
% play = 0.2
% Not play = 0.8

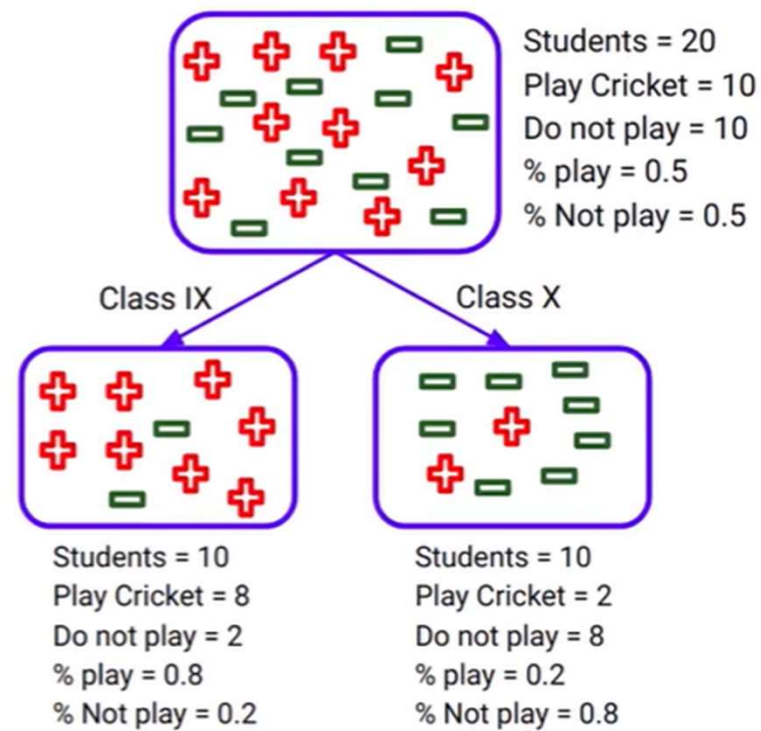# Steps to Calculate Entropy of Nodes

**Split on Class**

- Entropy for Parent node:
  $-(0.5)*\log_2(0,5) -(0.5)*\log_2(0.5) = 1$

- Entropy for sub-node Class IX:
  $-(0.8)*\log_2(0.8) -(0.2)*\log_2(0.2) = 0.722$

- Entropy for sub-node Class X:
  $-(0.2)*\log_2(0.2) -(0.8)*\log_2(0.8) = 0.722$

- Weighted Entropy: Class:
  $(10/20)*0.722 + (10/20)*0.722 = 0.722$

# Steps to Calculate Entropy of Nodes

| Split | Entropy | Information Gain |
|---|---|---|
| Performance in Class | 0.959 | 0.041 |
| Class | 0.722 | 0.278 |

Higher Information Gain is Good or Lower Information Gain is Good?

# Steps to Calculate Entropy of Nodes



Class IX          Class X

# Continuous Values!!

- So far, we dealt with Categorical Values….

- What about continuous data?

# Reduction in Variance

- **Formula:**

$$\text{Variance} = \Sigma \, [(X - \mu)^2] / n$$

- **Formula:**

$$\text{Variance} = \Sigma \, [(X - \mu)^2] \, / \, n$$

# Reduction in Variance

2  6  7
4  7  9

Variance ~ 6

1  1  1
1  1  1

Variance = 0

Lower Value of Variance is
Good or Higher?

# Properties of Variance

- Used when Target variable is Continuous

- Split with lower variance is selected

# Steps to Calculate Variance

- Calculate the variance of each child node

- Variance = $\Sigma \left[ (X - \mu)^2 \right] / n$

- Calculate the variance of each split as weighted average variance of each child node

# Steps to Calculate Variance

- Calculate the variance of each child node

- Variance = $\Sigma\, [(X - \mu)^2] / n$

- Calculate the variance of each split as weighted average variance of each child node
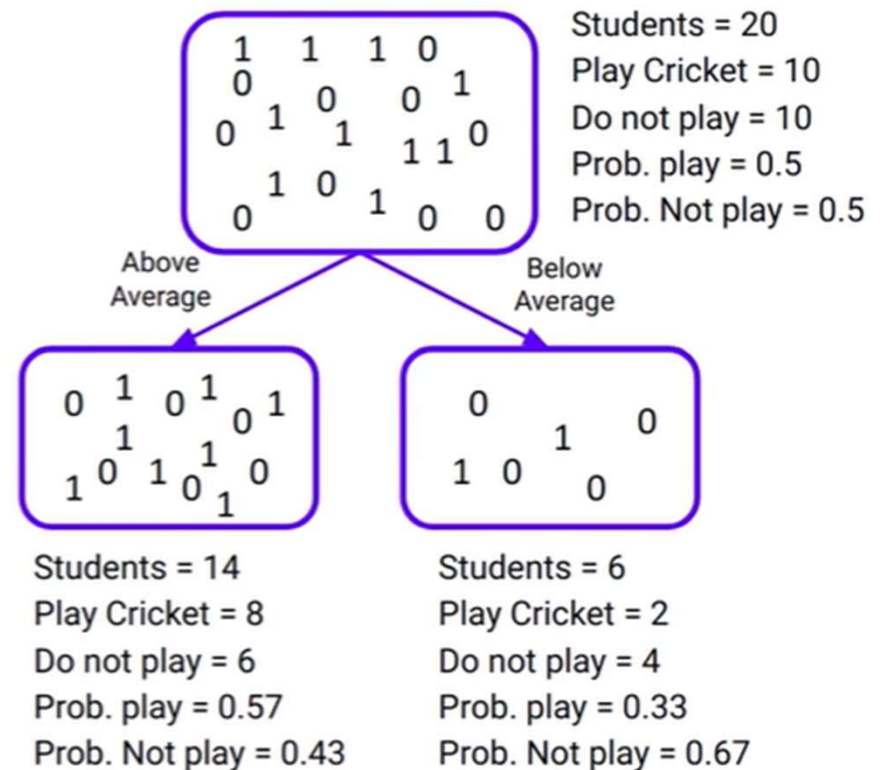
# Steps to Calculate Variance

- Plays Cricket = 1

- Do not play Cricket = 0

- **Above Average node:**
  - Mean = (8*1 + 6*0) / 14 = 0.57
  - Variance =
    $[8*(1-0.57)^2 + 6*(0-0.57)^2] / 14 = 0.245$

- **Below Average node:**
  - Mean = (2*1 + 4*0) / 6 = 0.33
  - Variance =
    $[2*(1-0.33)^2 + 4*(0-0.33)^2] / 6 = 0.222$

- **Variance: Performance in Class:**
  (14/20)*0.245 + (6/20)*0.222 = 0.238



Students = 20
Play Cricket = 10
Do not play = 10
Prob. play = 0.5
Prob. Not play = 0.5

Above Average          Below Average

Students = 14
Play Cricket = 8
Do not play = 6
Prob. play = 0.57
Prob. Not play = 0.43

Students = 6
Play Cricket = 2
Do not play = 4
Prob. play = 0.33
Prob. Not play = 0.67

# Steps to Calculate Variance



Students = 20
Play Cricket = 10
Do not play = 10
Prob. play = 0.5
Prob. Not play = 0.5

Class IX

Class X

Students = 10
Play Cricket = 8
Do not play = 2
Prob. play = 0.8
Prob. Not play = 0.2

Students = 10
Play Cricket = 2
Do not play = 8
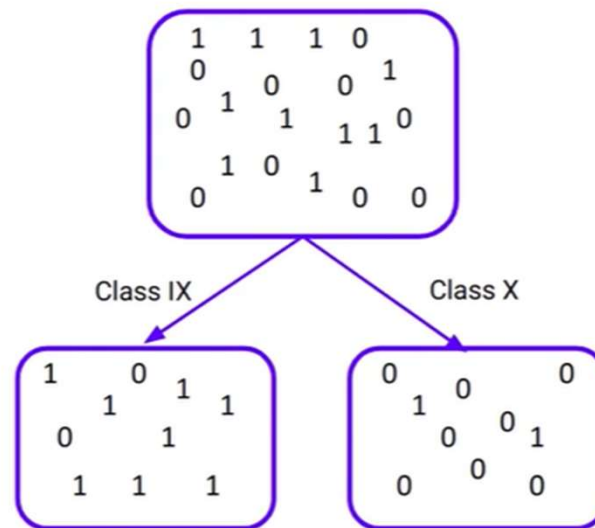Prob. play = 0.2
Prob. Not play = 0.8

# Steps to Calculate Variance

- Class IX node:
  - Mean = (8*1 + 2*0) / 10 = 0.8
  - Variance =
    $[8*(1-0.8)^2 + 2*(0-0.8)^2] / 10 = 0.16$

- Class X node:
  - Mean = (2*1 + 8*0) / 10 = 0.2
  - Variance =
    $[2*(1-0.2)^2 + 8*(0-0.2)^2] / 10 = 0.16$

- Variance: Class:
  (10/20)*0.16 + (10/20)*0.16 = 0.16

# Steps to Calculate Variance

| Split | Variance |
|---|---|
| Performance in Class | 0.238 |
| Class | 0.16 |

Split on Class

- Don't grow too much deeper trees, because it leads to **overfitting (Applying pre-pruning and post pruning can save you)**

- Don't grow too much shallow trees, because it leads to **underfitting (Increase Tree Depth, Use More Additional Features)**

- **CART (Classification and Regression Trees):**
    - Used for both classification and regression tasks
    - CART constructs binary trees by splitting data using the
        - Gini impurity for classification
        - Mean squared error for regression.

# Decision Tree Algorithms

- **ID3 (Iterative Dichotomiser 3):**
  - Primarily used for classification
  - ID3 employs a top-down approach to select the best attribute at each node using information gain, which measures how well an attribute separates the classes.

- **C4.5:An extension of ID3:**
  - C4.5 handles both categorical and continuous data
  - It utilizes the gain ratio to select splits and allowing for the handling of missing values and pruning to avoid overfitting.