# Machine Learning

## Decision Tree – 1

**Prof. Dr. Fazlul Hasan Siddiqui**
Head, Dept. of CSE, DUET, Gazipur
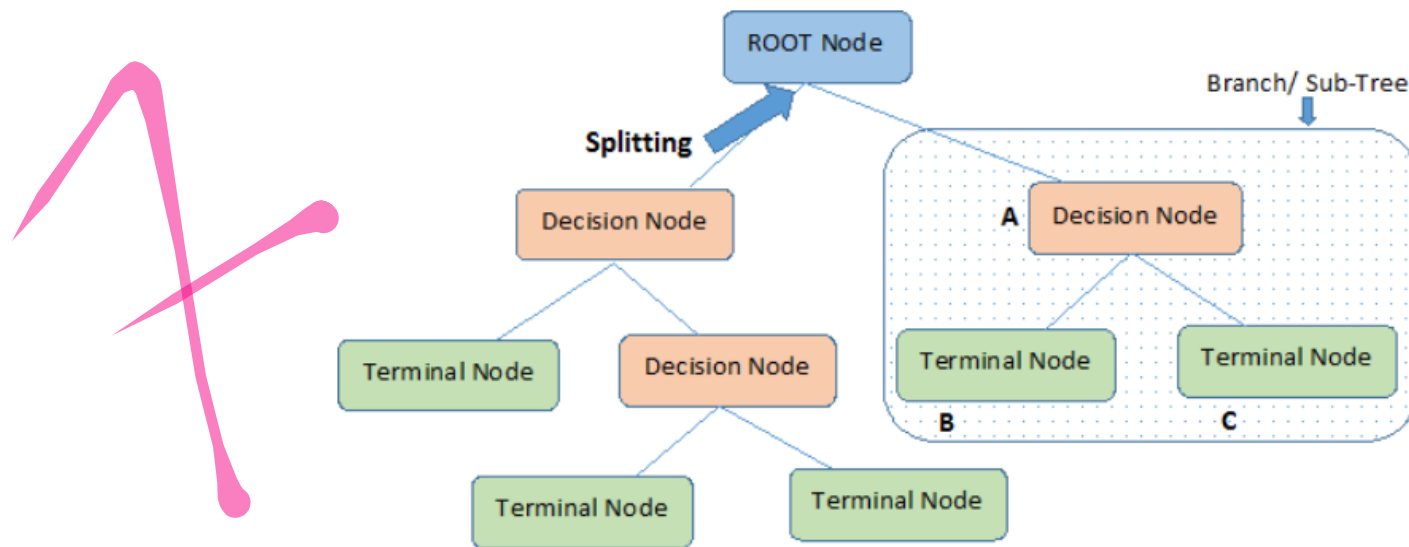BSc:IUT; MSc:BUET; PhD:ANU (Australia)
siddiqui@duet.ac.bd

# Decision Tree | Terminologies

1. **Root Node:** It represents entire population or sample and this further gets divided into two or more homogeneous sets.
2. **Splitting:** It is a process of dividing a node into two or more sub-nodes.
3. **Decision Node:** When a sub-node splits into further sub-nodes, then it is called decision node.
4. **Leaf/ Terminal Node:** Nodes do not split is called Leaf or Terminal node.
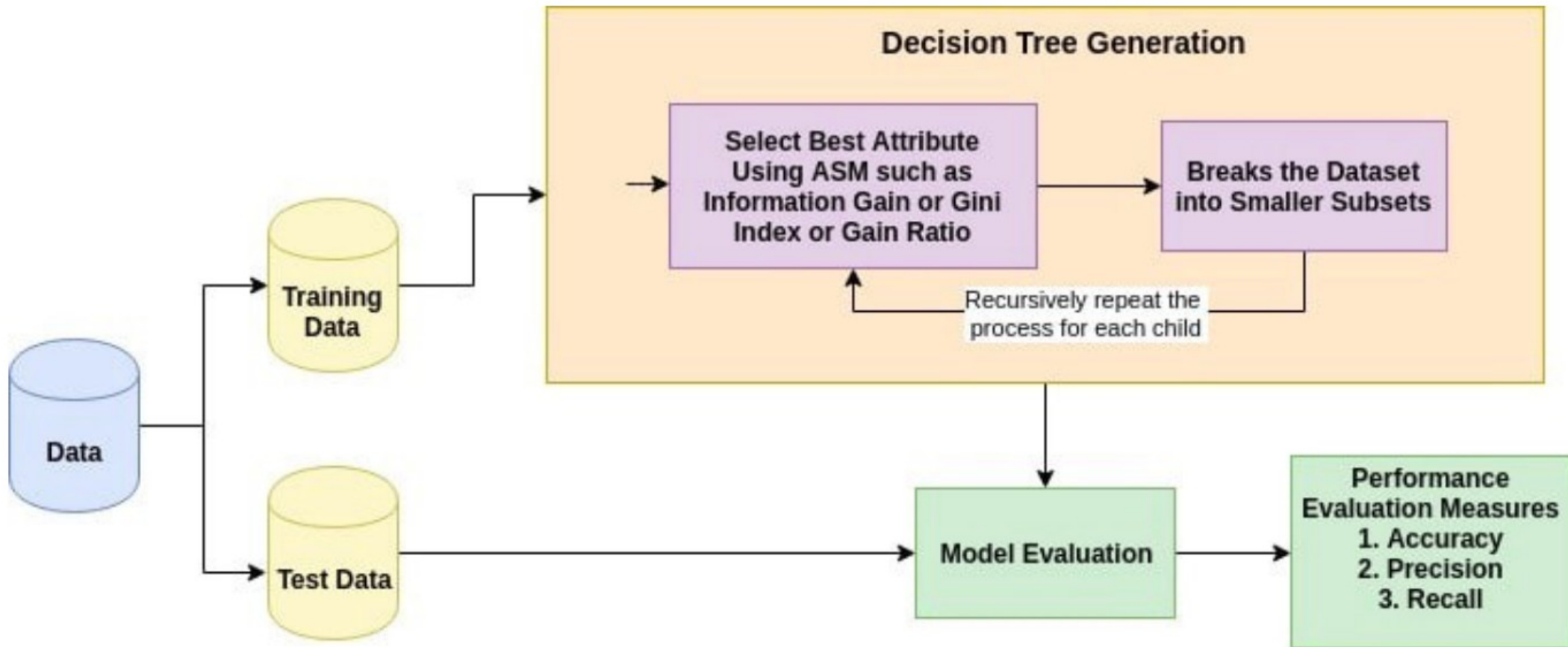5.



Note:- A is parent node of B and C.

**Pruning:** When we remove sub-nodes of a decision node, this process is called pruning. You can say opposite process of splitting.

6. **Branch / Sub-Tree:** A sub section of entire tree is called branch or sub-tree.
7. **Parent and Child Node:** A node, which is divided into sub-nodes is called parent node of sub-nodes where as sub-nodes are the child of parent node.

# Decision Tree | How It Works

1. Select the best attribute using Attribute Selection Measures(ASM) to split the records.

2. Make that attribute a decision node and breaks the dataset into smaller subsets.

3. Starts tree building by repeating this process recursively for each child until one of the condition will match:

   o All the tuples belong to the same attribute value.

   o There are no more remaining attributes.

   o There are no more instances.

# Decision Tree Generation

## Decision Tree Generation

**Select Best Attribute Using ASM such as Information Gain or Gini Index or Gain Ratio** → **Breaks the Dataset into Smaller Subsets**

Recursively repeat the process for each child

**Data** → **Training Data** → Decision Tree Generation

**Data** → **Test Data** → **Model Evaluation** → **Performance Evaluation Measures**
1. Accuracy
2. Precision
3. Recall

# Types of Decision Tree

Types of decision tree is based on the type of target variable we have. It can be of two types:

1. **Categorical Variable Decision Tree:** Decision Tree which has categorical target variable then it called as categorical variable decision tree. Example:- In above scenario of student problem, where the target variable was "Student will play cricket or not" i.e. YES or NO.
2. **Continuous Variable Decision Tree:** Decision Tree has continuous target variable then it is called as Continuous Variable Decision Tree.

**Example:-** Let's say we have a problem to predict whether a customer will pay his renewal premium with an insurance company (yes/ no). Here we know that income of customer is a significant variable but insurance company does not have income details for all customers. Now, as we know this is an important variable, then we can build a decision tree to predict customer income based on occupation, product and various other variables. In this case, we are predicting values for continuous variable.

# Types of Decision Tree

1. Regression trees are used when dependent variable is continuous. Classification trees are used when dependent variable is categorical.
2. In case of regression tree, the value obtained by terminal nodes in the training data is the mean response of observation falling in that region. Thus, if an unseen data observation falls in that region, we'll make its prediction with mean value.
3. In case of classification tree, the value (class) obtained by terminal node in the training data is the mode of observations falling in that region. Thus, if an unseen data observation falls in that region, we'll make its prediction with mode value.
4. Both the trees divide the predictor space (independent variables) into distinct and non-overlapping regions. For the sake of simplicity, you can think of these regions as high dimensional boxes or boxes.
5. Both the trees follow a top-down greedy approach known as recursive binary splitting. We call it as 'top-down' because it begins from the top of tree when all the observations are available in a single region and successively splits the predictor space into two new branches down the tree. It is known as 'greedy' because, the algorithm cares (looks for best variable available) about only the current split, and not about future splits which will lead to a better tree.
6. This splitting process is continued until a user defined stopping criteria is reached. For example: we can tell the the algorithm to stop once the number of observations per node becomes less than 50.
7. In both the cases, the splitting process results in fully grown trees until the stopping criteria is reached. But, the fully grown tree is likely to overfit data, leading to poor accuracy on unseen data. This bring 'pruning'. Pruning is one of the technique used tackle overfitting. We'll learn more about it in following section.
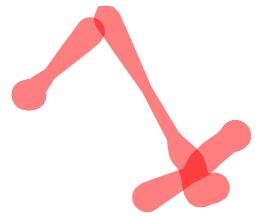
# Decision Tree | Advantages & Disadvantages

## Advantages

1. **Easy to Understand**: Decision tree output is very easy to understand even for people from non-analytical background. It does not require any statistical knowledge to read and interpret them. Its graphical representation is very intuitive and users can easily relate their hypothesis.

2. **Useful in Data exploration:** Decision tree is one of the fastest way to identify most significant variables and relation between two or more variables. With the help of decision trees, we can create new variables / features that has better power to predict target variable. You can refer article ([Trick to enhance power of regression model](#)) for one such trick.  It can also be used in data exploration stage. For example, we are working on a problem where we have information available in hundreds of variables, there decision tree will help to identify most significant variable.

3. **Less data cleaning required:** It requires less data cleaning compared to some other modeling techniques. It is not influenced by outliers and missing values to a fair degree.

4. **Data type is not a constraint:** It can handle both numerical and categorical variables.

5. **Non Parametric Method:** Decision tree is considered to be a non-parametric method. This means that decision trees have no assumptions about the space distribution and the classifier structure.

## Disadvantages

1. **Over fitting:** Over fitting is one of the most practical difficulty for decision tree models. This problem gets solved by setting constraints on model parameters and pruning (discussed in detailed below).

2. **Not fit for continuous variables**: While working with continuous numerical variables, decision tree looses information when it categorizes variables in different categories.

# Information Gain | Decision Tree

We should calculate the Entropy and Information Gain for the right decision tree.

Entropy

Entropy is the measure of randomness or *'impurity'* in the dataset
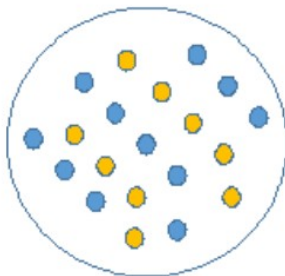
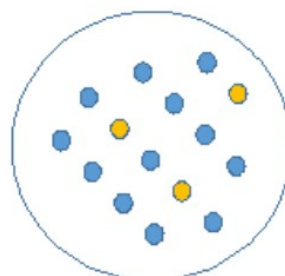Entropy should be low!

Information Gain

It is the measure of decrease in entropy after the dataset is split
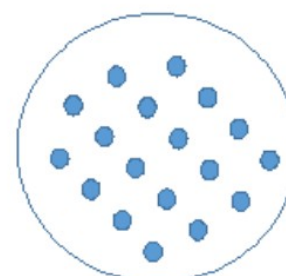
Also known as Entropy Reduction
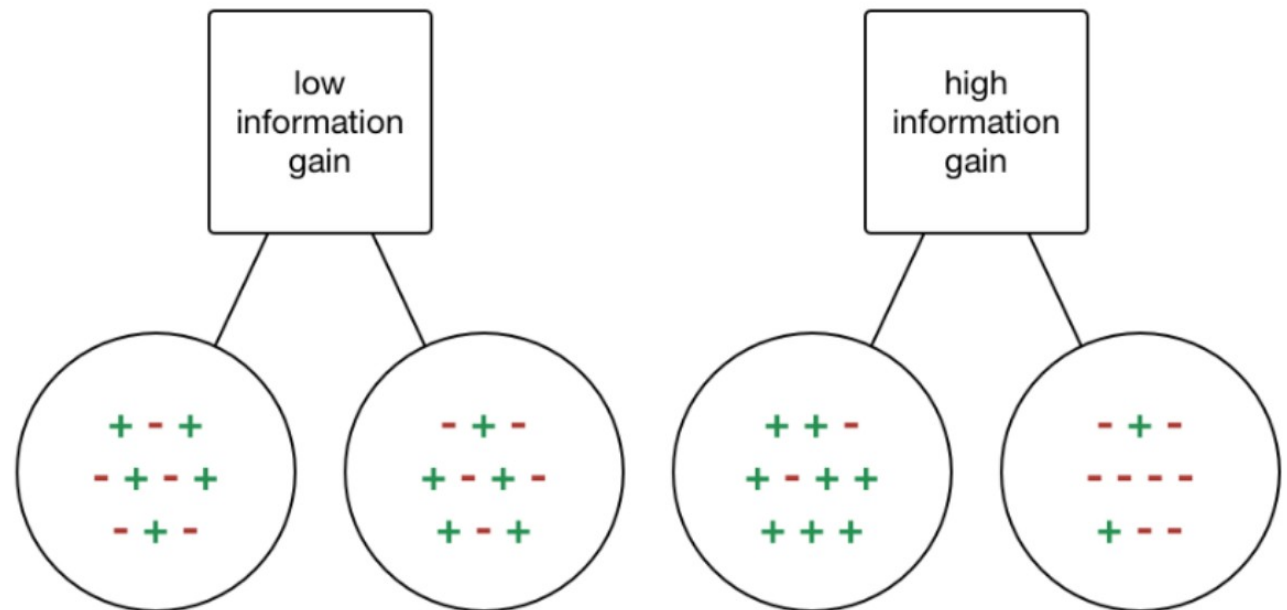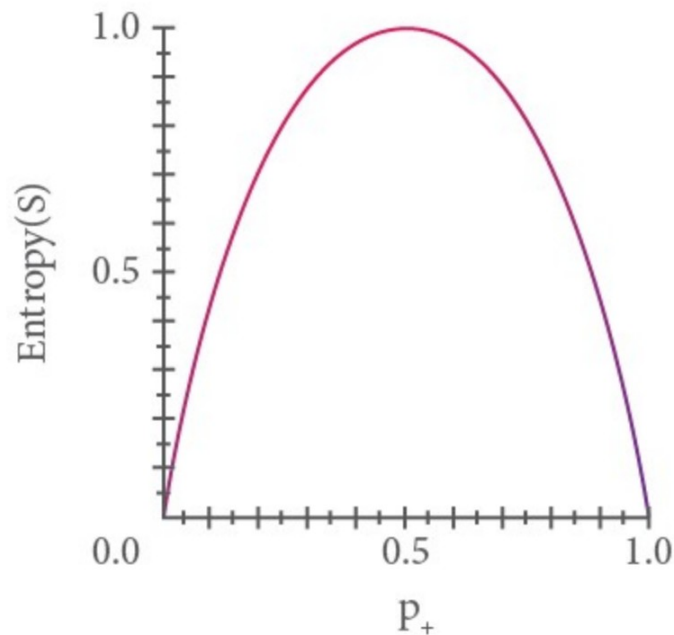
Information Gain should be high!

A

B

C

# Information Gain | Decision Tree

Target: Low Entropy and High Information Gain

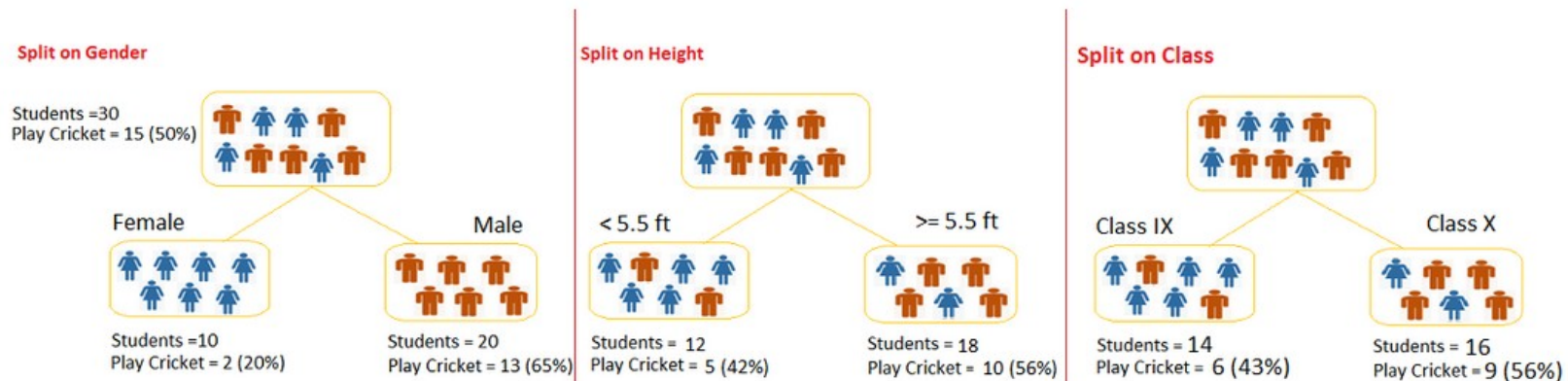The entropy (eg, flipping a coin) is zero when the probability is either 0 or 1.

The Entropy is maximum when the probability is 0.5 because it projects perfect randomness in the data.

# Information Gain Example

Let's say we have a sample of 30 students with three variables Gender (Boy/ Girl), Class( IX/ X) and Height (5 to 6 ft). 15 out of these 30 play cricket in leisure time. Now, I want to create a model to predict who will play cricket during leisure period? In this problem, we need to segregate students who play cricket in their leisure time based on highly significant input variable among all three.

This is where decision tree helps, it will segregate the students based on all values of three variable and identify the variable, which creates the best homogeneous sets of students (which are heterogeneous to each other). In the snapshot below, you can see that variable Gender is able to identify best homogeneous sets compared to the other two variables.

**Split on Gender**

Students =30
Play Cricket = 15 (50%)

Female
Students =10
Play Cricket = 2 (20%)

Male
Students = 20
Play Cricket = 13 (65%)

**Split on Height**

< 5.5 ft
Students = 12
Play Cricket = 5 (42%)

>= 5.5 ft
Students = 18
Play Cricket = 10 (56%)

**Split on Class**

Class IX
Students = 14
Play Cricket = 6 (43%)

Class X
Students = 16
Play Cricket = 9 (56%)

As mentioned above, decision tree identifies the most significant variable and it's value that gives best homogeneous sets of population. Now the question which arises is, how does it identify the variable and the split? To do this, decision tree uses various algorithms, which we will discuss in the following section.

**I.G.**

Entropy can be calculated using formula:-

$$\text{Entropy} = -p \log_2 p - q \log_2 q$$

Here p and q is probability of success and failure respectively in that node. Entropy is also used with categorical target variable. It chooses the split which has lowest entropy compared to parent node and other splits. The lesser the entropy, the better it is.

## Steps to calculate entropy for a split:

1. Calculate entropy of parent node
2. Calculate entropy of each individual node of split and calculate weighted average of all sub-nodes available in split.

**Example:** Let's use this method to identify best split for student example.

1. Entropy for parent node = $-(15/30) \log2 (15/30) - (15/30) \log2 (15/30)$ = **1**. Here 1 shows that it is a impure node.
2. Entropy for Female node = $-(2/10) \log2 (2/10) - (8/10) \log2 (8/10)$ = 0.72 and for male node, $-(13/20) \log2 (13/20) - (7/20) \log2 (7/20)$ = **0.93**
3. Entropy for split Gender = Weighted entropy of sub-nodes = $(10/30)*0.72 + (20/30)*0.93$ = **0.86**
4. Entropy for Class IX node, $-(6/14) \log2 (6/14) - (8/14) \log2 (8/14)$ = 0.99 and for Class X node, $-(9/16) \log2 (9/16) - (7/16) \log2 (7/16)$ = 0.99.
5. Entropy for split Class = $(14/30)*0.99 + (16/30)*0.99$ = **0.99**

Above, you can see that entropy for *Split on Gender* is the lowest among all, so the tree will split on *Gender*. We can derive information gain from entropy as **1- Entropy.**

# Gini Index

You can understand the Gini index as a cost function used to evaluate splits in the dataset. It is calculated by subtracting the sum of the squared probabilities of each class from one. It favors larger partitions and easy to implement whereas information gain favors smaller partitions with distinct values.

Gini says, if we select two items from a population at random then they must be of same class and probability for this is 1 if population is pure.

1. It works with categorical target variable "Success" or "Failure".
2. It performs only Binary splits
3. Higher the value of Gini higher the homogeneity.
4. CART (Classification and Regression Tree) uses Gini method to create binary splits.
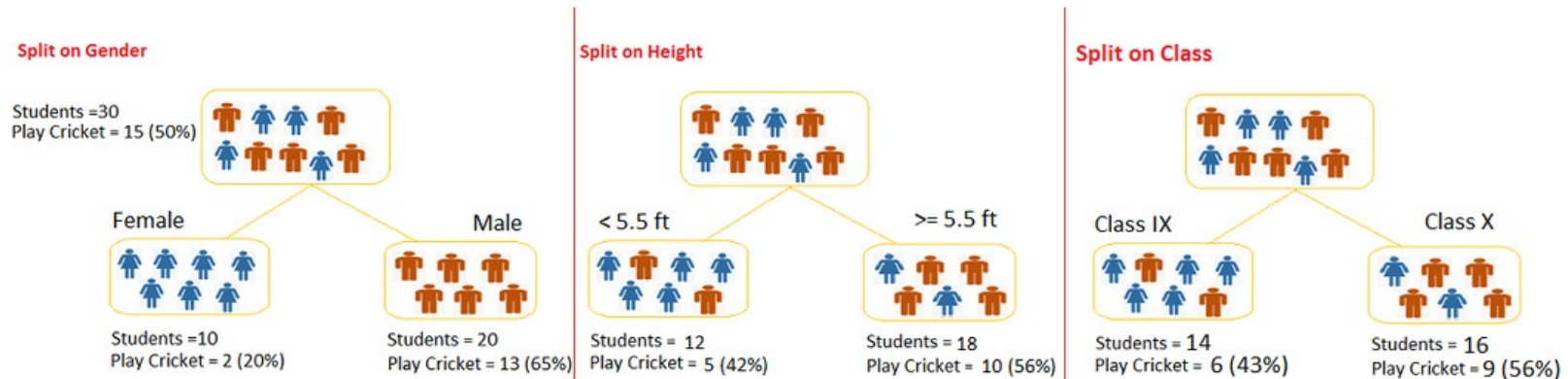
**Steps to Calculate Gini for a split**

1. Calculate Gini for sub-nodes, using formula sum of square of probability for success and failure (p^2+q^2).
2. Calculate Gini for split using weighted Gini score of each node of that split

# Gini Index Example

Let's say we have a sample of 30 students with three variables Gender (Boy/ Girl), Class( IX/ X) and Height (5 to 6 ft). 15 out of these 30 play cricket in leisure time. Now, I want to create a model to predict who will play cricket during leisure period? In this problem, we need to segregate students who play cricket in their leisure time based on highly significant input variable among all three.

This is where decision tree helps, it will segregate the students based on all values of three variable and identify the variable, which creates the best homogeneous sets of students (which are heterogeneous to each other). In the snapshot below, you can see that variable Gender is able to identify best homogeneous sets compared to the other two variables.



**Split on Gender**

Students =30
Play Cricket = 15 (50%)

Female
Students =10
Play Cricket = 2 (20%)

Male
Students = 20
Play Cricket = 13 (65%)

**Split on Height**

< 5.5 ft
Students = 12
Play Cricket = 5 (42%)

>= 5.5 ft
Students = 18
Play Cricket = 10 (56%)

**Split on Class**

Class IX
Students = 14
Play Cricket = 6 (43%)

Class X
Students = 16
Play Cricket = 9 (56%)

As mentioned above, decision tree identifies the most significant variable and it's value that gives best homogeneous sets of population. Now the question which arises is, how does it identify the variable and the split? To do this, decision tree uses various algorithms, which we will discuss in the following section.
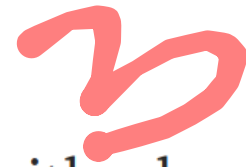
# Gini Index Example

**Split on Gender:**

1. Calculate, Gini for sub-node Female = $(0.2)*(0.2)+(0.8)*(0.8)=0.68$

2. Gini for sub-node Male = $(0.65)*(0.65)+(0.35)*(0.35)=0.55$

3. Calculate weighted Gini for Split Gender = $(10/30)*0.68+(20/30)*0.55 = \mathbf{0.59}$

**Similar for Split on Class:**

1. Gini for sub-node Class IX = $(0.43)*(0.43)+(0.57)*(0.57)=0.51$

2. Gini for sub-node Class X = $(0.56)*(0.56)+(0.44)*(0.44)=0.51$

3. Calculate weighted Gini for Split Class = $(14/30)*0.51+(16/30)*0.51 = \mathbf{0.51}$

Above, you can see that Gini score for *Split on Gender* is higher than *Split on Class*, hence, the node split will take place on Gender.
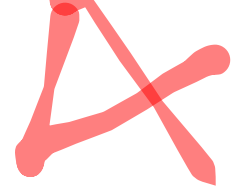
# Gain Ratio

Information gain is biased towards choosing attributes with a large number of values as root nodes. It means it prefers the attribute with a large number of distinct values.

C4.5, an improvement of ID3, uses Gain ratio which is a modification of Information gain that reduces its bias and is usually the best option. Gain ratio overcomes the problem with information gain by taking into account the number of branches that would result before making the split. It corrects information gain by taking the intrinsic information of a split into account.

$$Gain\ Ratio = \frac{Information\ Gain}{SplitInfo} = \frac{Entropy\ (before) - \sum_{j=1}^{K} Entropy(j,\ after)}{\sum_{j=1}^{K} w_j\ log_2\ w_j}$$

Where "before" is the dataset before the split, K is the number of subsets generated by the split, and (j, after) is subset j after the split.

# Reduction in Variance

**Reduction in variance** is an algorithm used for continuous target variables (regression problems). This algorithm uses the standard formula of variance to choose the best split. The split with lower variance is selected as the criteria to split the population:
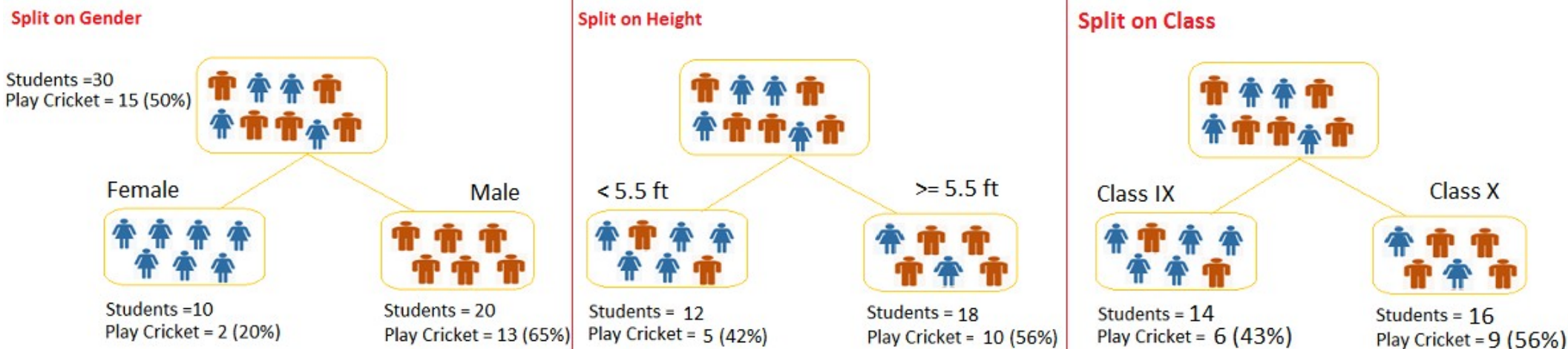
$$\text{Variance} = \frac{\Sigma(X - \overline{X})^2}{n}$$

Above X-bar is the mean of the values, X is actual and n is the number of values.

**Steps to calculate Variance:**

1. Calculate variance for each node.

2. Calculate variance for each split as the weighted average of each node variance.

# Reduction in Variance | Example



**Split on Gender**

Students =30
Play Cricket = 15 (50%)

Female

Students =10
Play Cricket = 2 (20%)

Male

Students = 20
Play Cricket = 13 (65%)

**Split on Height**

< 5.5 ft

>= 5.5 ft

Students = 12
Play Cricket = 5 (42%)

Students = 18
Play Cricket = 10 (56%)

**Split on Class**

Class IX

Class X

Students = 14
Play Cricket = 6 (43%)

Students = 16
Play Cricket = 9 (56%)

**Example:-** Let's assign numerical value 1 for play cricket and 0 for not playing cricket. Now follow the steps to identify the right split:

1. Variance for Root node, here mean value is $(15*1 + 15*0)/30 = 0.5$ and we have 15 one and 15 zero. Now variance would be $((1-0.5)^2+(1-0.5)^2+....15$ times$+(0-0.5)^2+(0-0.5)^2+...15$ times$) / 30$, this can be written as $(15*(1-0.5)^2+15*(0-0.5)^2) / 30 =$ **0.25**

2. Mean of Female node = $(2*1+8*0)/10=0.2$ and Variance = $(2*(1-0.2)^2+8*(0-0.2)^2) / 10 = 0.16$

3. Mean of Male Node = $(13*1+7*0)/20=0.65$ and Variance = $(13*(1-0.65)^2+7*(0-0.65)^2) / 20 = 0.23$

4. Variance for Split Gender = Weighted Variance of Sub-nodes = $(10/30)*0.16 + (20/30) *0.23 =$ **0.21**

5. Mean of Class IX node = $(6*1+8*0)/14=0.43$ and Variance = $(6*(1-0.43)^2+8*(0-0.43)^2) / 14= 0.24$

6. Mean of Class X node = $(9*1+7*0)/16=0.56$ and Variance = $(9*(1-0.56)^2+7*(0-0.56)^2) / 16 = 0.25$

7. Variance for Split Gender = $(14/30)*0.24 + (16/30) *0.25 =$ **0.25**

# Chi-Sqare

It is an algorithm to find out the statistical significance between the differences between sub-nodes and parent node. We measure it by sum of squares of standardized differences between observed and expected frequencies of target variable.

1. It works with categorical target variable "Success" or "Failure".
2. It can perform two or more splits.
3. Higher the value of Chi-Square higher the statistical significance of differences between sub-node and Parent node.
4. Chi-Square of each node is calculated using formula,
5. Chi-square = ((Actual − Expected)^2 / Expected)^1/2
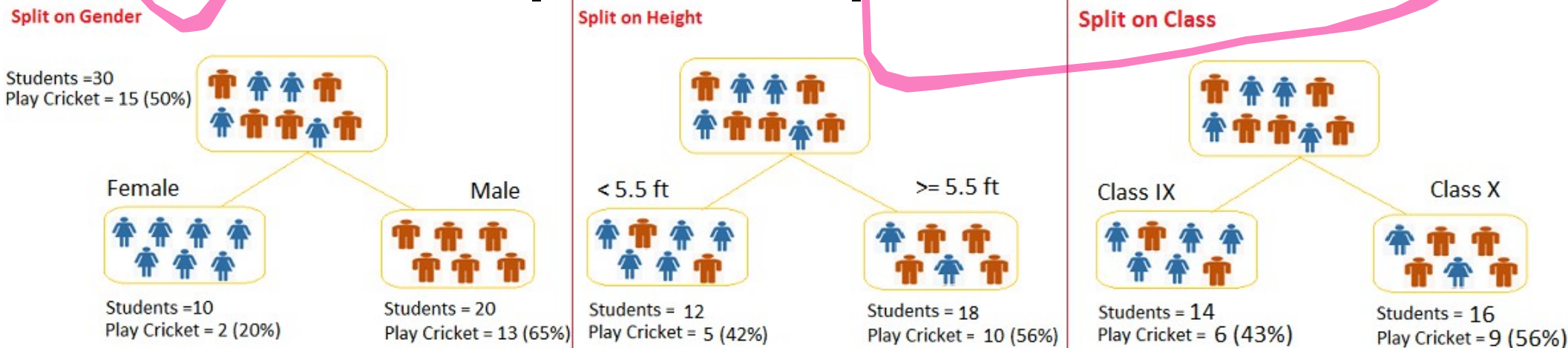6. It generates tree called CHAID (Chi-square Automatic Interaction Detector)

**Steps to Calculate Chi-square for a split:**

1. Calculate Chi-square for an individual node by calculating the deviation for Success and Failure both

2. Calculated Chi-square of Split using Sum of all Chi-square of success and Failure of each node of the split

$$\chi^2 = \sum \frac{(O-E)^2}{E}$$

$\chi^2$ = Chi Square obtained
$\sum$ = the sum of
$O$ = observed score
$E$ = expected score

1. First we are populating for node Female, Populate the actual value for "**Play Cricket**" and "**Not Play Cricket**", here these are 2 and 8 respectively.
2. Calculate expected value for "**Play Cricket**" and "**Not Play Cricket**", here it would be 5 for both because parent node has probability of 50% and we have applied same probability on Female count(10).
3. Calculate deviations by using formula, Actual – Expected. It is for "**Play Cricket**" (2 – 5 = -3) and for "**Not play cricket**" ( 8 – 5 = 3).
4. Calculate Chi-square of node for "**Play Cricket**" and "**Not Play Cricket**" using formula with formula, = $((Actual - Expected)^2 / Expected)^{1/2}$. You can refer below table for calculation.
5. Follow similar steps for calculating Chi-square value for Male node.
6. Now add all Chi-square values to calculate Chi-square for split Gender.

| Node | Play Cricket | Not Play Cricket | Total | Expected Play Cricket | Expected Not Play Cricket | Deviation Play Cricket | Deviation Not Play Cricket | Chi-Square Play Cricket | Chi-Square Not Play Cricket |
|---|---|---|---|---|---|---|---|---|---|
| Female | 2 | 8 | 10 | 5 | 5 | -3 | 3 | 1.34 | 1.34 |
| Male | 13 | 7 | 20 | 10 | 10 | 3 | -3 | 0.95 | 0.95 |
| | | | | | | | Total Chi-Square | 4.58 | |

# Chi-Sqare | Example | Split on Class

**Split on Gender**

Students =30
Play Cricket = 15 (50%)

Female

Students =10
Play Cricket = 2 (20%)

Male

Students = 20
Play Cricket = 13 (65%)

**Split on Height**

< 5.5 ft

Students = 12
Play Cricket = 5 (42%)

>= 5.5 ft

Students = 18
Play Cricket = 10 (56%)

**Split on Class**

Class IX

Students = 14
Play Cricket = 6 (43%)

Class X

Students = 16
Play Cricket = 9 (56%)

Perform similar steps of calculation for split on Class and you will come up with below table.

| Node | Play Cricket | Not Play Cricket | Total | Expected Play Cricket | Expected Not Play Cricket | Deviation Play Cricket | Deviation Not Play Cricket | Chi-Square | |
|------|--------------|------------------|-------|-----------------------|---------------------------|------------------------|----------------------------|------------|------------|
| | | | | | | | | Play Cricket | Not Play Cricket |
| IX | 6 | 8 | 14 | 7 | 7 | -1 | 1 | 0.38 | 0.38 |
| X | 9 | 7 | 16 | 8 | 8 | 1 | -1 | 0.35 | 0.35 |
| | | | | | | | Total Chi-Square | 1.46 | |

Above, you can see that Chi-square also identify the Gender split is more significant compare to Class.