

## Decision Trees

This can be used for both classification and regression tasks.

Decision tree is called as non metric methods for classification.

Decision tree are non metric classifier because they perform classification through rule-based splitting rather than distance or similarity computation.

## Metric Based classification Methods

A metric based classification technique that relies on:

- A distance measure a
- A similarity measure  
(dot product)

## Non metric methods for classification

1) Problem with some datasets

→ in many real world problems, features are not numeric

→ These features are called as nominal features

2) Nominal features:

→ nominal features are descriptive in nature.

→ They describe the characteristics of a class, not quantities

→ Example.

- \* colour
- \* shape
- \* Taste

3) If they Numeric Representation is not meaningful

→ These features do not have a natural order.

→ Converting them into numbers can create wrong assumptions.

→ Example.

- \* Red = 1, Green = 2, Yellow = 3
- \* This wrongly suggests Yellow > Green > Red.

4) Metric cannot be computed

→ Since features are not numeric, we cannot calculate distance or similarity

→ Therefore, metric based methods cannot be used.

5) Non-metric models:

→ Classification models built using nominal features are called as

Non metric method based models.

→ Example:

\* Decision Tree

\* Random Forest.

## Nominal feature : Illustration

- \* Nominal features are non numeric features.
- \* They describe an object in a natural, human readable way.

Features	Values
color	Green, yellow, Red
size	Small, medium, big
shape,	Round, non-Round
Taste	Sweet, non-Sweet

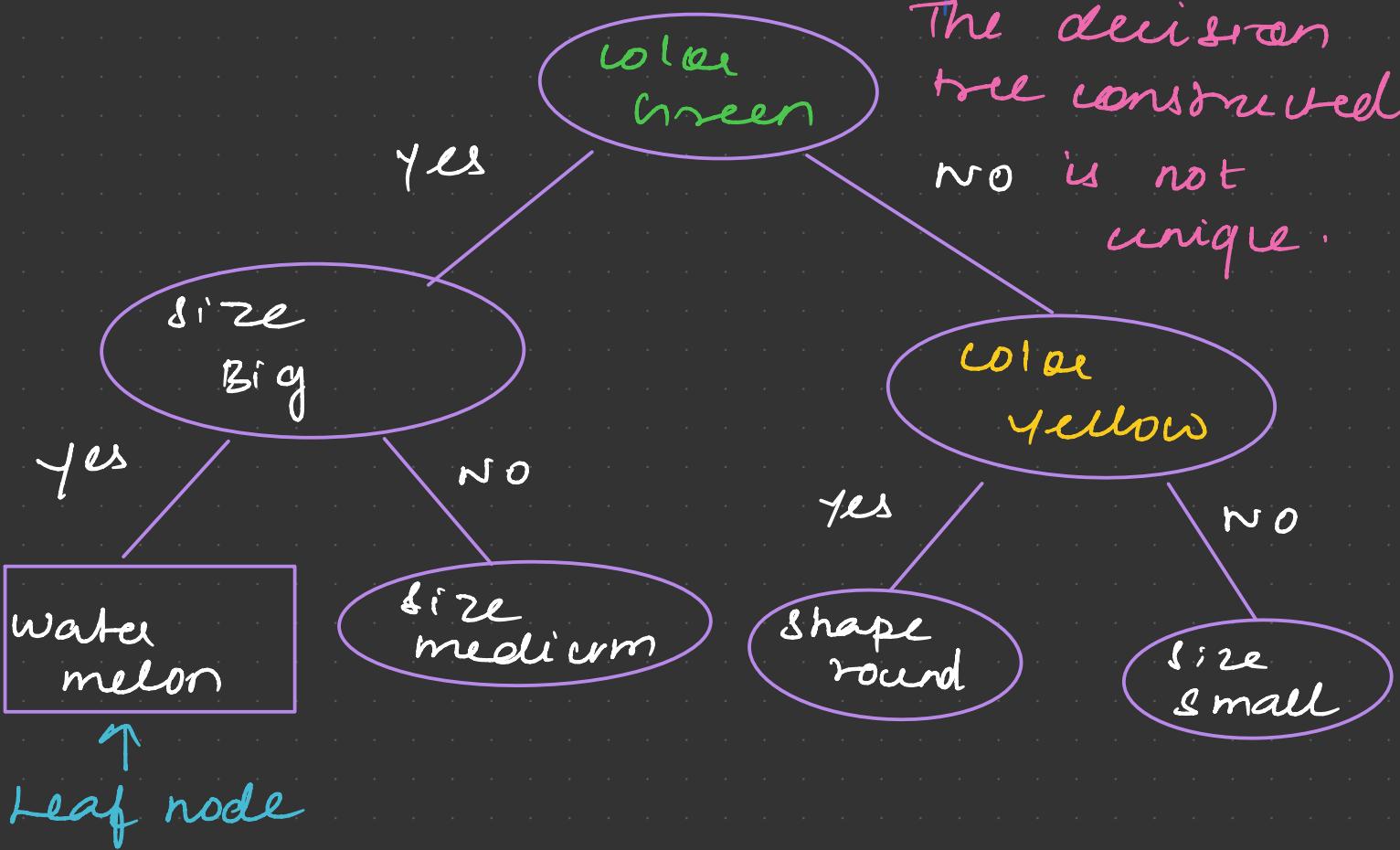
classes :

watermelon, Apple, grape, Banana,  
Lemon, Cherry, grapefruit.

- \* each fruit is described using words, not numbers
- \* Not all the features are used for every fruit
- \* The same fruit can have multiple description

example.

- \* watermelon → green, big
- \* Banana → yellow, non-round
- \* Lemon → yellow, round, not big (small or medium)
- \* Grapefruit → yellow, round, big
- \* Apple → green, red, medium
- \* cherry → small, red, sweet
- \* grape → small, black, not-sweet.
- \* This creates vagueness and variation in data.



\* \* This tree grows until its depth.

In decision tree

- At every node we ask a query or question
- Based on the answer tree is traversed down to reach a leaf node.
- Leaf node is associated with a class.

→ Decision Tree can also be built for numerical feature.

\* no limit to ask question on one feature node.

→ can combine more than one feature.

\* number of questions asked is very large in number

\* set of several questions, which questions to be asked at a particular node is the issue to be addressed.

\* this issue to be addressed is constraining a decision tree.

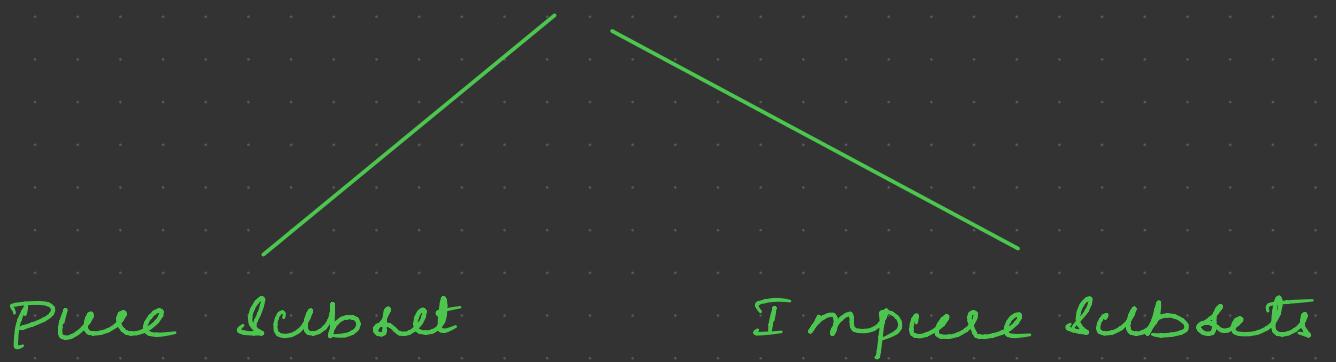
# Classification and Regression Tree (CART)

- \* CART provides framework to construct a decision tree.
- \* works for
  - Classification (class labels)
  - Regression (numerical output)

## How CART Builds the Tree

- \* The decision tree starts with all training data at the root node.
- \* At each node:
  - A question (query) is asked
  - The data is splitted into smaller subsets
- \* Each subset sent to a child node.

## Subsets



### Pure subsets:

All examples in the subsets have the same class label.

Pure subsets are called as pure nodes.

### Impure subset

It includes examples of more than one classes

A node with impure subsets are called as impure nodes.

## \* Quantification of impurity :

Impurity measure.

Impurity measure, will be used to construct the decision tree to decide which query need to be processed at each of the node.

splitting at one node should decrease the impurity from one level to the next level.

## Goal of CART

- \* keep splitting the data such that
  - \* Nodes becomes more pure
  - \* Tree becomes simple and compact

To summarize, CART is a framework used to construct decision trees by recursively splitting the training dataset into smaller subsets to form a pure nodes for classification or regression.

Main issues in CART approach:

(1) Number of split at a node.

→ A node can split into

- \* Two Branches (Binary split)

$$\hookrightarrow B=2$$

- \* Three Branches (Ternary split)

$$\hookrightarrow B=3$$

- \* or more (Multi-way split)

→ Branching Factor (B)

\* Number of branches coming out of a node.

→ In practice:

\* Binary split is most commonly used

\* Makes the tree simpler and easier to control

2) Query selection at a Node

→ Many possible questions can be asked at a node.

→ The challenge is

> which question should be selected.

→ Rule used:

\* Occam's Razor

> The simplest tree that explains the data is preferred.

→ goal

- \* choose the question that
  - > makes child nodes less impure
  - > leads to a small and compact tree.

3) when should we stop splitting

- \* If we keep splitting
  - > Tree becomes very large
  - > model may overfit
  - > classification becomes slower
- \* so we need a stopping condition

4) Assigning class label to impure leaf nodes

- \* sometimes a leaf nodes still has
  - > samples from multiple class
- \* Then we must decide:
  - > which class label to assign to that leaf itself

CART must decide

- \* How many branches to create
- \* which question to ask
- \* when to stop growing the tree
- \* what label to give at the end.

## CART : Measure of Impurity

- \* Impurity measure how mixed the classes are at a node
- \* It helps to decide
  - which question (query) should be asked at a node
- \* the goal is to build a simple and compact decision tree.

### Notation

- \* let a node be called  $N_n$  ( $N_d$ )
- \* The impurity at the node is written as

$$i(N_d)$$

## Meaning of impurity values

$$i(N_d) = 0$$

- \* Node is pure
- \* All samples belong to the same class

$$i(N_d) \text{ is large}$$

- \* Node is impure
- \* samples are evenly distributed among multiple classes.

## Why Impurity is important?

when a node is split

- \* The impurity of child nodes should be lower
- \* This means split is good

CART always tries to:

- \* Reduce impurity from parent node to child nodes.

## Entropy (Information Impurity)

Entropy is the most commonly used impurity measure

It measures the uncertainty or randomness in the distribution

### Formula

$$H(\text{nd}) = - \sum_{j=1}^M P(C_j) \log_2 P(C_j)$$

### Terms,

$M \rightarrow$  number of classes

$P(C_j) \rightarrow$  fraction of samples at node  $\text{nd}$  that belongs to  $C_j$

### Interpretation

\* more class mixing  $\rightarrow$  higher entropy

\* more class purity  $\rightarrow$  lower entropy.

## Summary

- \* If all data at a node belongs to one class  $\rightarrow$  Entropy = 0
- \* If data is evenly split across classes  $\rightarrow$  Entropy is maximum
- \* CART chooses splits that reduce entropy the most.

## Gini Impurity

- \* Gini Impurity is a measure of how impure the classes are at a node
- \* It is an alternative to Entropy
- \* commonly used in CART and Practical implementation

Gini impurity measures the probability of incorrect classification of randomly chosen samples from the node, if its labeled according to the class distribution at the node

### Formula

$$i(N_d) = 1 - \sum_{j=1}^M P(C_j)^2$$

## Entropy v Gini

- \* Both measures Impurity
- \* Both aims to
  - select the split that makes child nodes more pure
- \* Gini is:
  - faster to compute
  - often used in real-world tree implementation

## Summary

- \* If most samples belongs to one class → Gini is low
- \* If samples are evenly spread across classes → Gini is high
- \* CART chooses the split that reduces Gini impurity the most.

CART: Query selection at a node

and splitting - construction of decision tree

- \* At each node, many possible questions (queries) can be asked
- \* CART evaluates all possible queries
- \* Goal is to select the query that:

Reduces the impurity the most.

- \* generally, impurity level of the child nodes should be smaller.
- \* The difference between the impurity level at a node (node) and its child node is important.

## Basic Idea

- \* A node contains a subset of training data
- \* when a query is applied
  - data is split into child nodes
- \* The child nodes should be more pure than the parent node.

## Binary split

This is a most common case

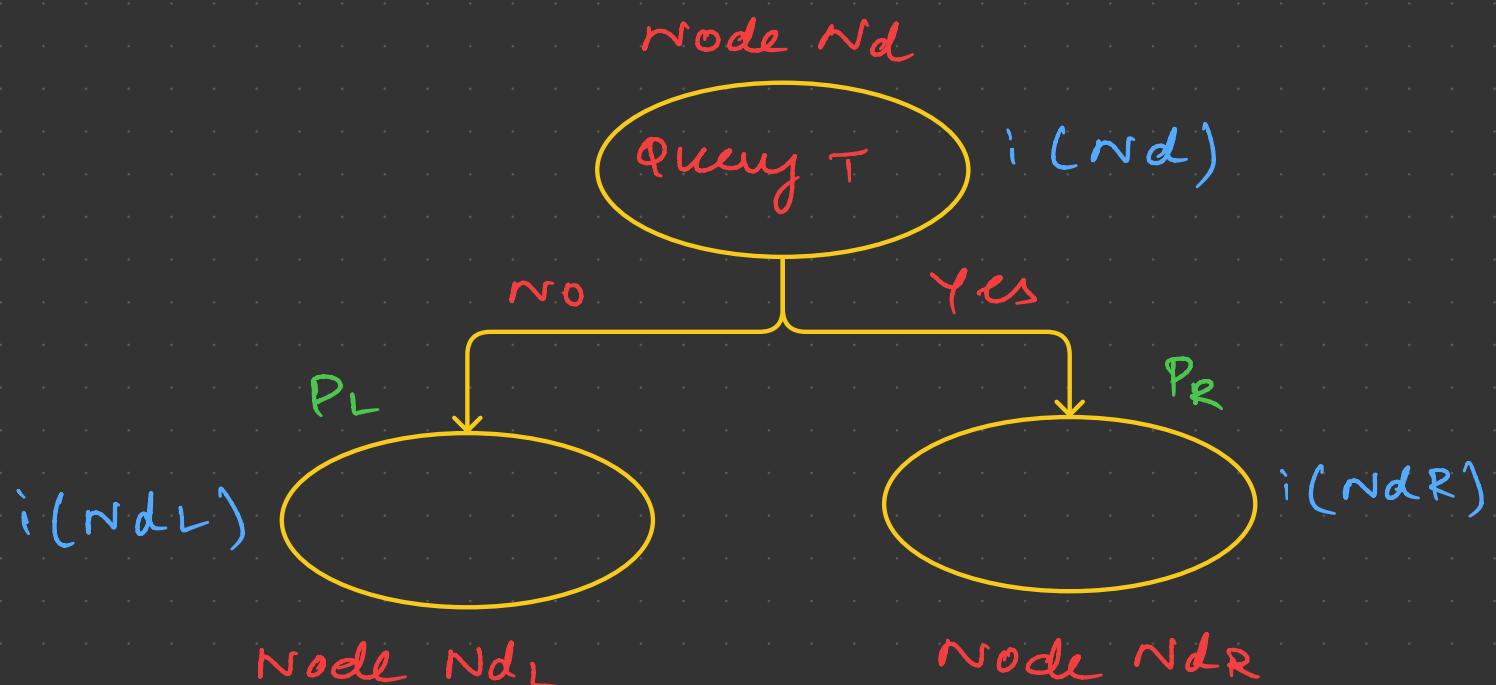
- \* Let the current node be  $nd$
- \* After applying the query  $T$ , it splits into:
  - $nd\ L$  → Left child
  - $nd\ R$  → Right child.

## Fraction of Data

- \*  $P_L \rightarrow$  Fraction of samples that goes to  $N_{dL}$
- \*  $P_R \rightarrow$  Fraction of samples that goes to  $N_{dR}$
- \* Relationship of  $P_R$

Since we have only two nodes.

$$P_R = 1 - P_L$$



## Impurity terms

- \*  $i(\text{nd}) \rightarrow$  Impurity of Parent node
- \*  $i(\text{ndL}) \rightarrow$  Impurity of left child
- \*  $i(\text{ndR}) \rightarrow$  Impurity of right child

## Change in Impurity (Information Gain)

The decrease in impurity caused by query T is

$$\Delta i(\text{nd}) = i(\text{nd}) - (P_L \cdot i(\text{ndL}) + P_R \cdot i(\text{ndR}))$$

## Decision Rule

- \* compute  $\Delta i(\text{nd})$  for all possible queries
- \* select the query which has maximum  $\Delta i(\text{nd})$
- \* This query gives the best split.

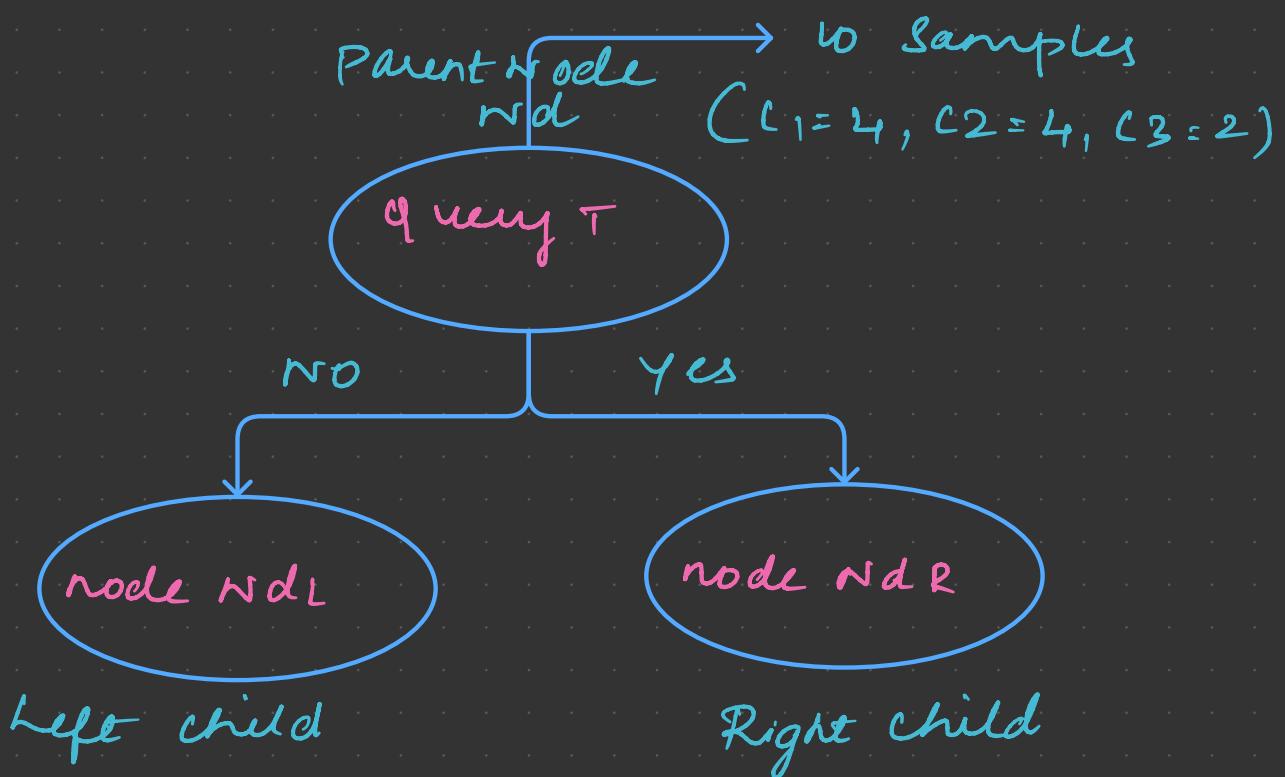
Change in the impurity level is called as Information Gain  $\rightarrow$  when entropy impurity is used.

## Illustration

How to calculate impurity and information gain using numbers and how CART chooses the best query.

Given data at node  $nd$

- \* The total number of samples = 10
- \* Class distribution
  - \* C1 → 4 samples
  - \* C2 → 4 samples
  - \* C3 = 2 samples.

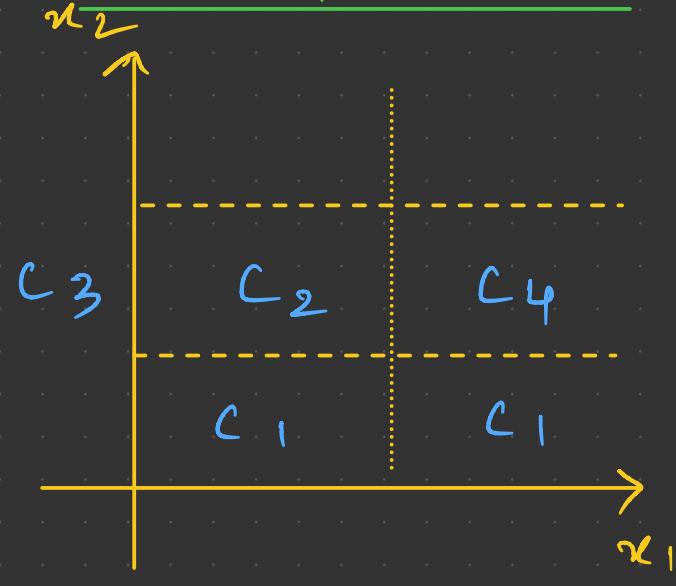


## with tree structure

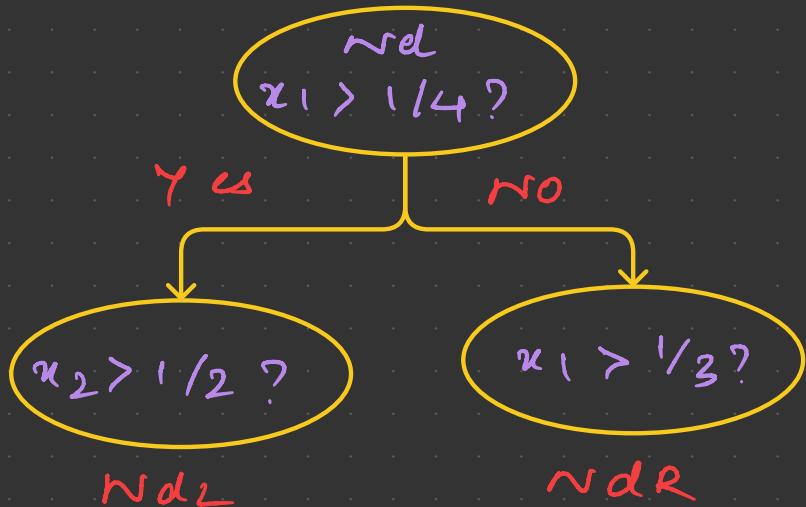
with both :

- \* The tree structure ( $nd \rightarrow ndL/ndR$ )
- \* feature space ( $x_1, x_2 \rightarrow \text{Region split}$ )

feature space view:



Tree View:



Step 1 - Impurity at parent node ( $nd$ )

using Entropy

$$i(nd) = - \sum P(C_j) \log_2 P(C_j)$$

## Probabilities

$$* P(C_1) = 4/10$$

$$* P(C_2) = 4/10$$

$$* P(C_3) = 2/10$$

here,

$$* \text{number of classes } M = 3$$

\* we must compute entropy for  
 $C_1, C_2, C_3$

## compute probabilities

$$P(C_1) = 4/10 = 0.4$$

$$P(C_2) = 4/10 = 0.4$$

$$P(C_3) = 2/10 = 0.2$$

## Substitute in the formula

$$i(\text{rnd}) = -[0.4 \log_2(0.4) + 0.4 \log_2(0.4) + 0.2 \log_2(0.2)]$$

## compute log values

$$\log_2(0.4) \approx -1.3219$$

$$\log_2(0.2) \approx -2.3219$$

## Multiply Terms

$$0.4 (-1.3219) = -0.5288$$

$$0.2 (-2.3219) = -0.4644$$

therefore,

$$i(\text{Nd}) = -0.5288 + (-0.5288) + (-0.4644)$$

$$i(\text{Nd}) = -(-1.522)$$

$$i(\text{Nd}) \approx 1.522$$

→ This node is impure because sample belongs to more classes.

## Query 1 - First possible split

Left child ( $\text{Nd}_L$ )

\* Total samples = 4

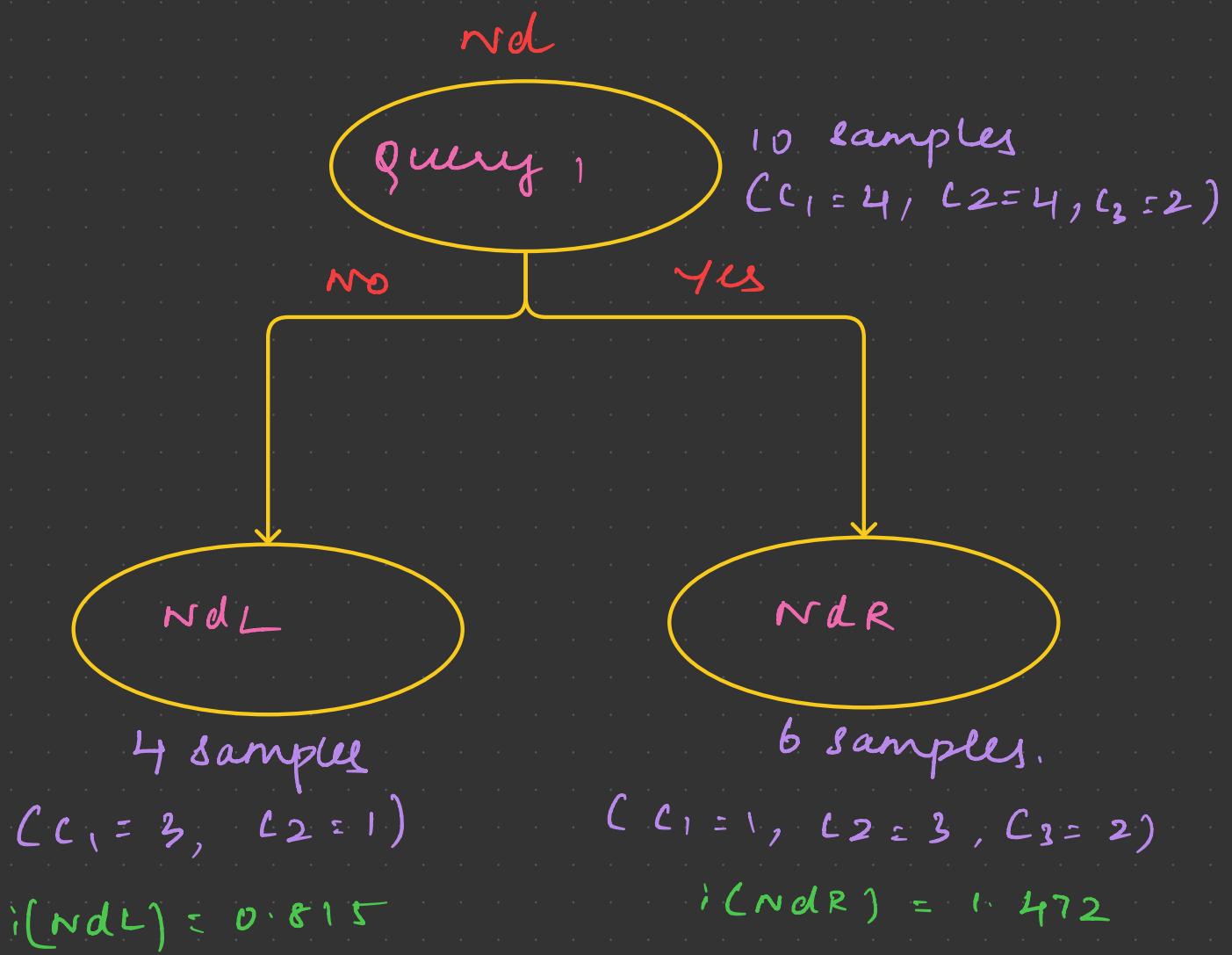
\*  $C_1 = 3, C_2 = 1$

$$i(\text{Nd}_L) = 0.815$$

## Right child ( $ndR$ )

- \* Total samples = 6
- \*  $C_1 = 1, C_2 = 3, C_3 = 2$

$$i(ndR) = 1.472$$



## Step 2 - Information gain for query 1

Fractions,

$$\star PL = 4/10$$

$$\star PR = 6/10$$

Formula

$$\Delta i(\text{nd}) = -i(\text{nd}) - (PL \cdot i(\text{ndl}) + PR \cdot i(\text{ndr}))$$

Result

$$i(\text{nd}) = 0.315$$

## Query 2 - Second possible split

Left child (ndl)

$$\star \text{Total samples} = 3$$

$$\star C_1 = 2, C_3 = 1$$

$$i(\text{ndl}) = 0.918$$

Right child (ndr)

$$\star \text{Total samples} = 7$$

$$\star C_1 = 1, C_2 = 4, C_3 = 2$$

$$i(\text{ndr}) = 1.379$$

### Step 3 - Information Gain for Query 2

$$\Delta_i(N_d) = 0.280$$

#### Decision

Compare,

$$\rightarrow \text{Query 1} = 0.135$$

$$\rightarrow \text{Query 2} = 0.280$$

from the observation Query 1 will be chosen because it gives the maximum impurity reduction.

#### Multimay split

- + In multimay split, a node can split into more than two branches.
- + Each possible answer to a query creates one child node.
- + The number of branches is called as Branching Factor ( $B$ )

## Branching Factor ( $B$ )

$$B = 3$$

node  $n_d$

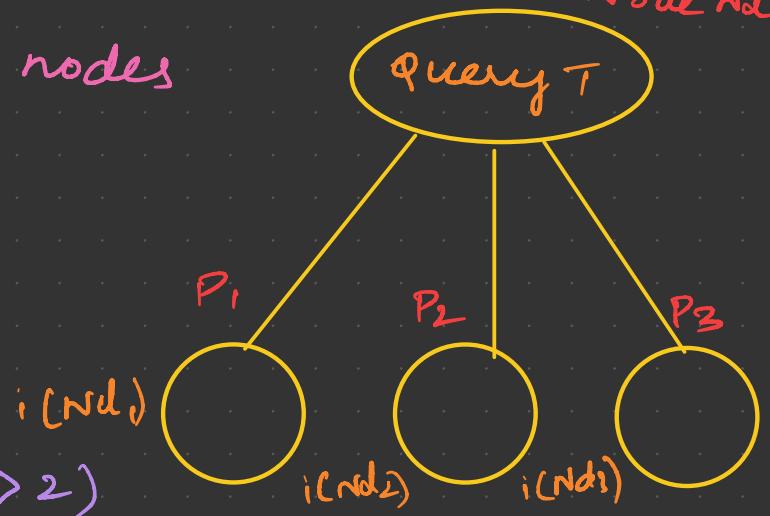
- \*  $B$  = number of child nodes

### \* Examples

→ Binary split ( $B=2$ )

→ Ternary split ( $B=3$ )

→ Multiway split ( $B > 2$ )



## How the split works

### Parent node

- \* Let the parent node  $n_d$
- \* It contains a subset of training samples.

### After query $T$

- \* The nodes splits into

→  $n_d_1, n_d_2, n_d_3, \dots, n_d_n$

- \* Each child gets a fraction of samples.

## Fraction of samples

$P_K \rightarrow$  fraction of samples that can go to the  $K^{th}$  child node

sum of all fractions

$$P_1 + P_2 + \dots + P_n = 1$$

# Impurity Reduction

## Impurity terms

$i(Nd)$  → Impurity of parent node

$i(N_{dk})$  → Impurity of K-th child node

## Formula

$$\Delta_i(Nd) = i(Nd) - \sum_{k=1}^B p_k \cdot i(N_{dk})$$

## Decision Rule

- \* Try all possible queries
- \* For each query
  - compute  $\Delta_i(Nd)$
- \* Select the query with the maximum  $\Delta_i(Nd)$

## Splitting process

- \* Tree construction starts from the root node
- \* At each node, CART :
  - Test all possible queries
  - Measures impurity reduction
  - chooses the best query.

## Step by step process

### Step 1 - Start at Root Node

- \* All training samples are placed at the root node ( $nd$ )
- \* This node usually impure.

### Step 2 - Select Best Query

- \* Try many queries (questions)
- \* For each query
  - Split the data
  - compute  $\Delta_i(nd)$  (impurity reduction)
- \* choose the query with maximum  $\Delta_i(nd)$

### Step 3 - Create child nodes

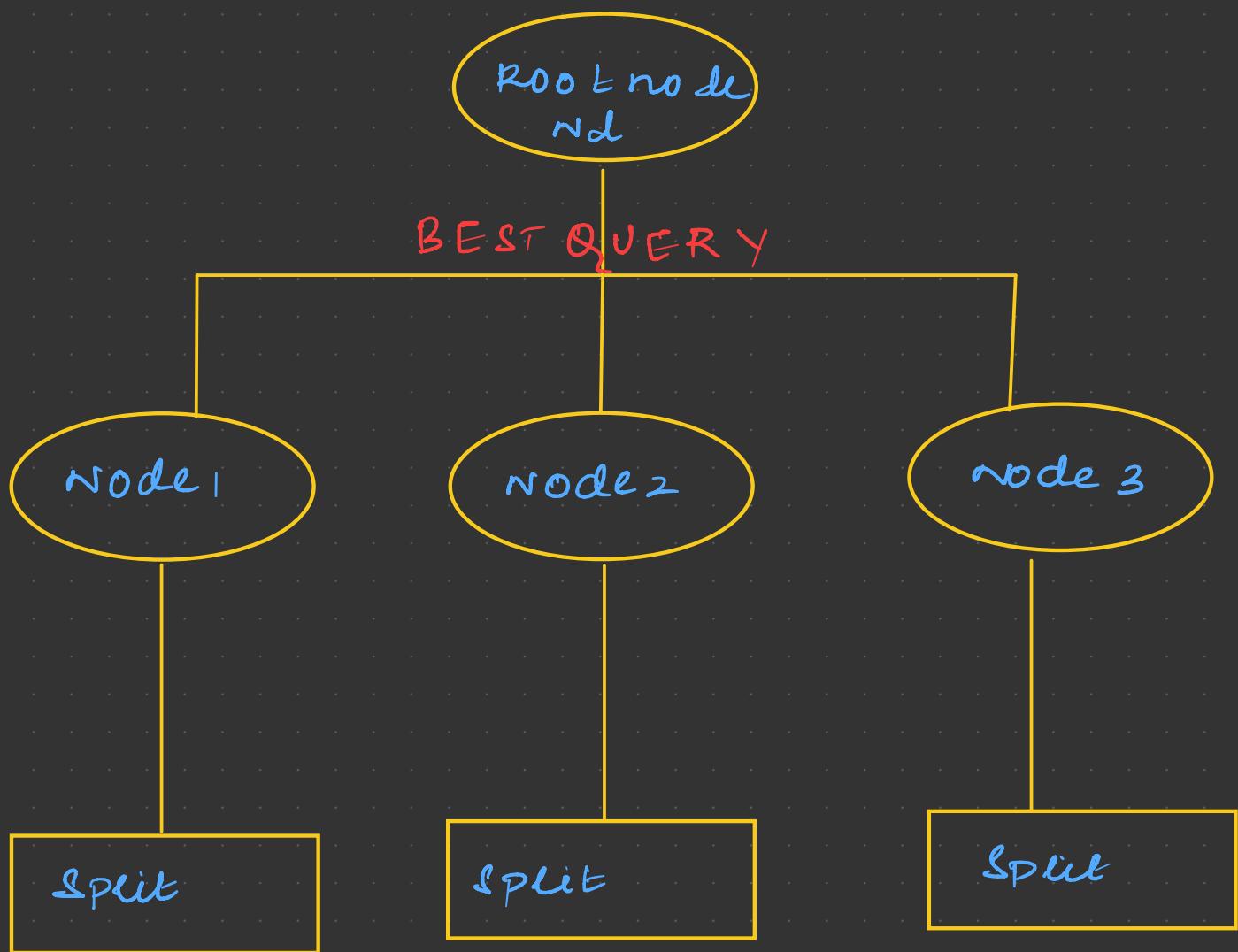
- \* The selected query splits into
  - $Nd_1, Nd_2, Nd_3, \dots, Nd_n$  (child nodes)
- \* Each child gets subset of data

### Step 4 - Repeat Recursively

- \* For each child node:
  - Treat it as a new parent
  - Again select the best query
  - Split it further.

### Step 5 - Stop at Leaf nodes

- \* The process continues until:
  - node becomes pure
  - or
  - A stopping condition is met



## Summary

- Start at the top
- Ask the best question
- Split the data
- Repeat for every branch
- Stop when no split is possible or stopping condition is met.

## criterion to stop splitting

- \* CART must decide when to stop growing the decision tree so that:
  - The tree does not become too large
  - The model does not overfit
  - The tree remains simple and compact

## stopping criterion 1 - Impurity reduction threshold

### threshold

#### Rule

- \* At a node  $nd$ , compute the maximum impurity reduction $\Delta_i(nd)$
- \* If the value is below a small threshold, then stop splitting

#### Typical threshold values

$$\rightarrow 10^{-2} \text{ or } 10^{-3}$$

#### Meaning,

If no query can reduce impurity significantly — further splitting is not useful

## Stopping criteria - 2 - Tree size control (Regularization)

### Rule

stop splitting if

$$\sum_{\text{node} \in S_L} i(\text{nd}) + \alpha \times (\text{Number of nodes in tree})$$

$S_L \rightarrow$  set of all leaf nodes

$i(\text{nd}) \rightarrow$  Impurity of each leaf node

$\alpha \rightarrow$  Hyperparameter that controls  
Penalty for tree nodes

### Meaning

- \* First term  $\rightarrow$  Measures total impurity
- \* Second term  $\rightarrow$  Penalizes large trees
- \* Together  $\rightarrow$  Prefer small, low impurity trees

### Why this is Important

$\rightarrow$  Prevents the tree from growing too deep

$\rightarrow$  Reduce overfitting

$\rightarrow$  Improves generalization on new data.

## Assignment of class label to leaf node

After the decision tree is constructed  
CART must decide what class label  
to assign to each leaf node

### case 1 - Pure leaf node

A leaf node is pure if

→ All training samples that  
reach this node belongs to one class  
only

### Rule

Assign the leaf node the same  
class label

for example,

If leaf node :

$$c_1 = 5, c_2 = 0, c_3 = 0$$

then

$$\text{Label} = c_1$$

## case 2 - Impure Leaf Node

A leaf node is impure if -

→ samples belongs to more than one classes

### Rule

- \* Find the class with the largest number of samples
- \* Assign that class as the leaf node label

for example,

If at a leaf node

$$C_1 = 2, C_2 = 5, C_3 = 3$$

Then

label =  $C_2$  (majority class)

### Why this needed?

- \* sometimes CART stops splitting early
- \* so leaf nodes may still contain mixed classes
- \* we still need a final decision rule.