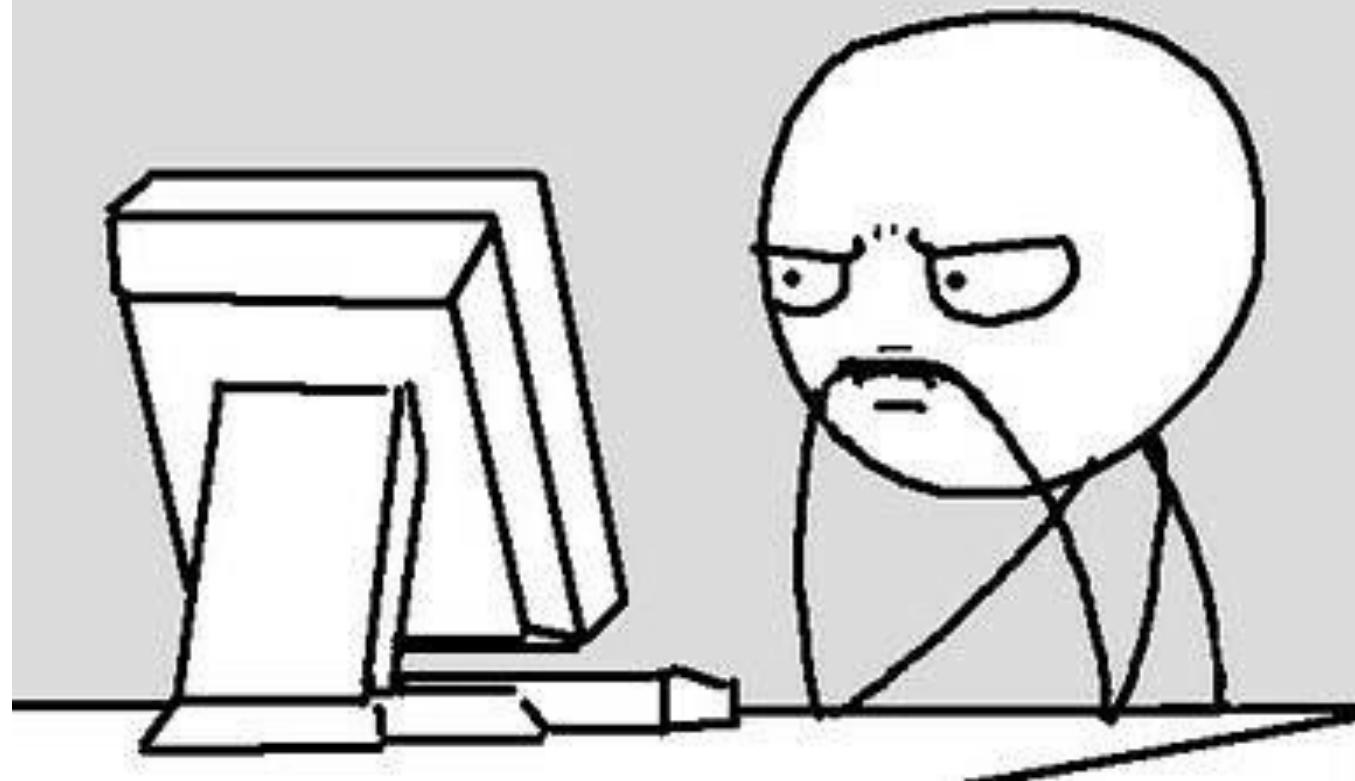
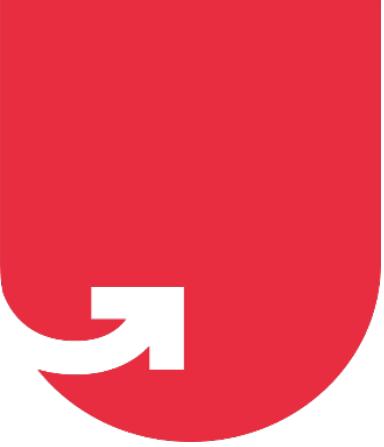


LET'S WAIT





#LifeKoKaroLift

Data Science Certification Program

Course : Machine Learning

Lecture On : Tree Models

Instructor : Shivam Garg



AGENDA

- Why Tree Models?
- Decision Trees
- ✓ Gini Index Calculation – Example
- Doubt Session

linear models

linear reg

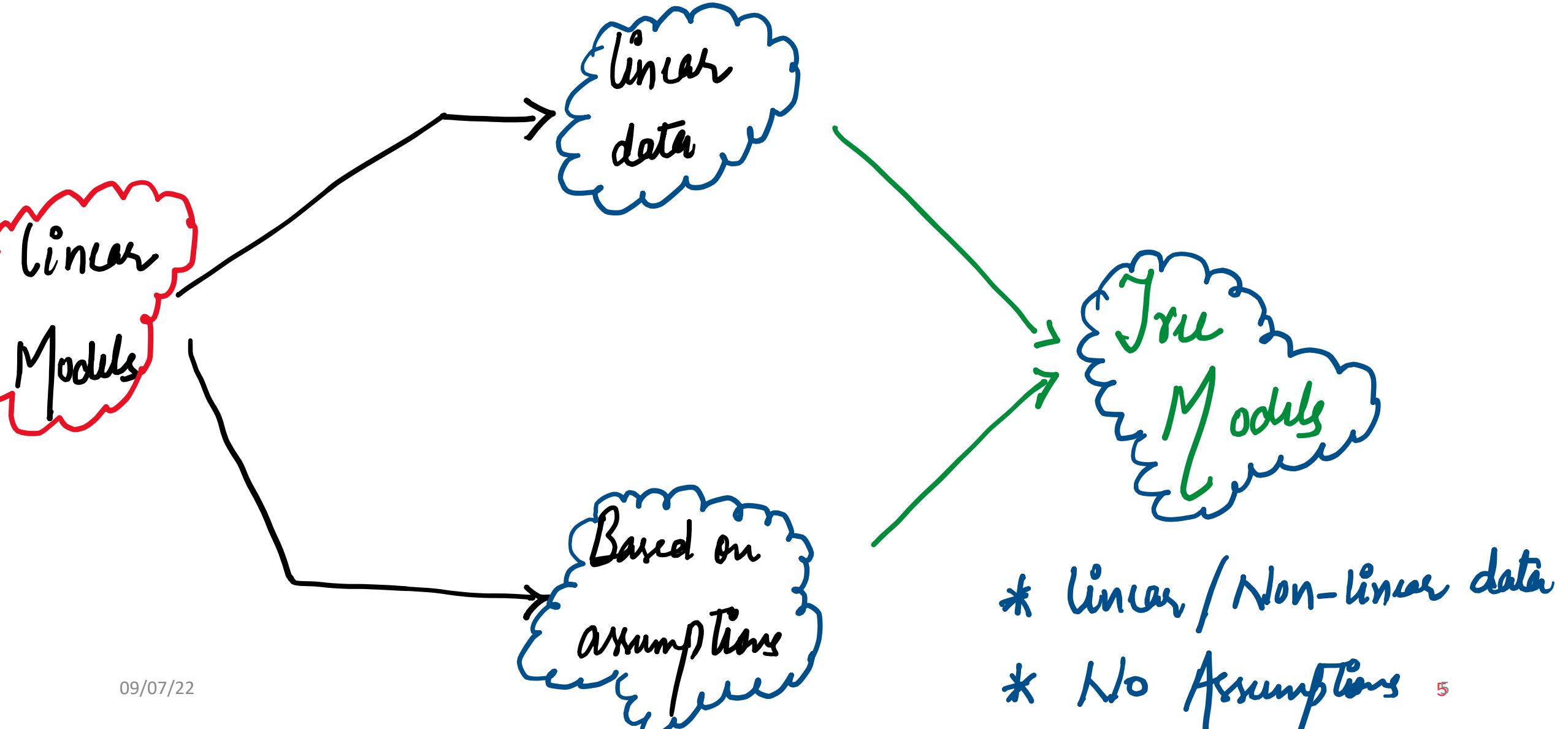
logistic reg

architecture

interpret

constraint

Why Tree Models?:



✓ ✓
Linear
Model

Linear
data

* High interpretability
* less computation

✓
Tree
Models

linear/
Non-linear
data

* partially interpretable
* More computation

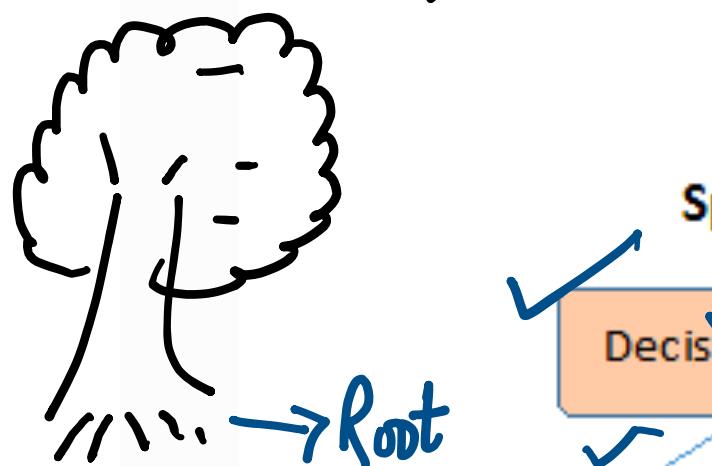
Neural

Networks

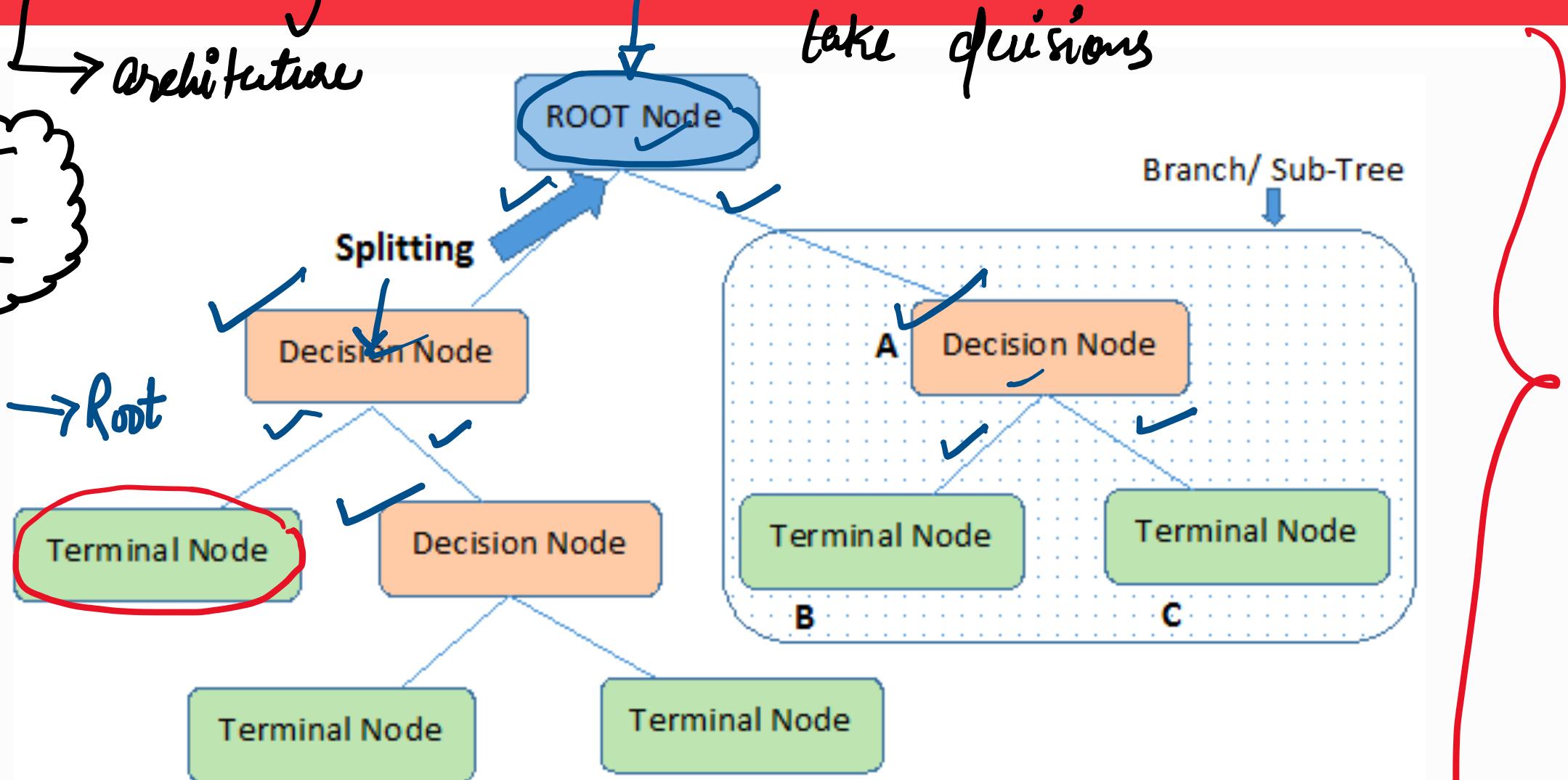
Non-linear

* No interpretability
* High computational cost

Decision Tree: } Tree like Structure which is used to take decisions upGrad

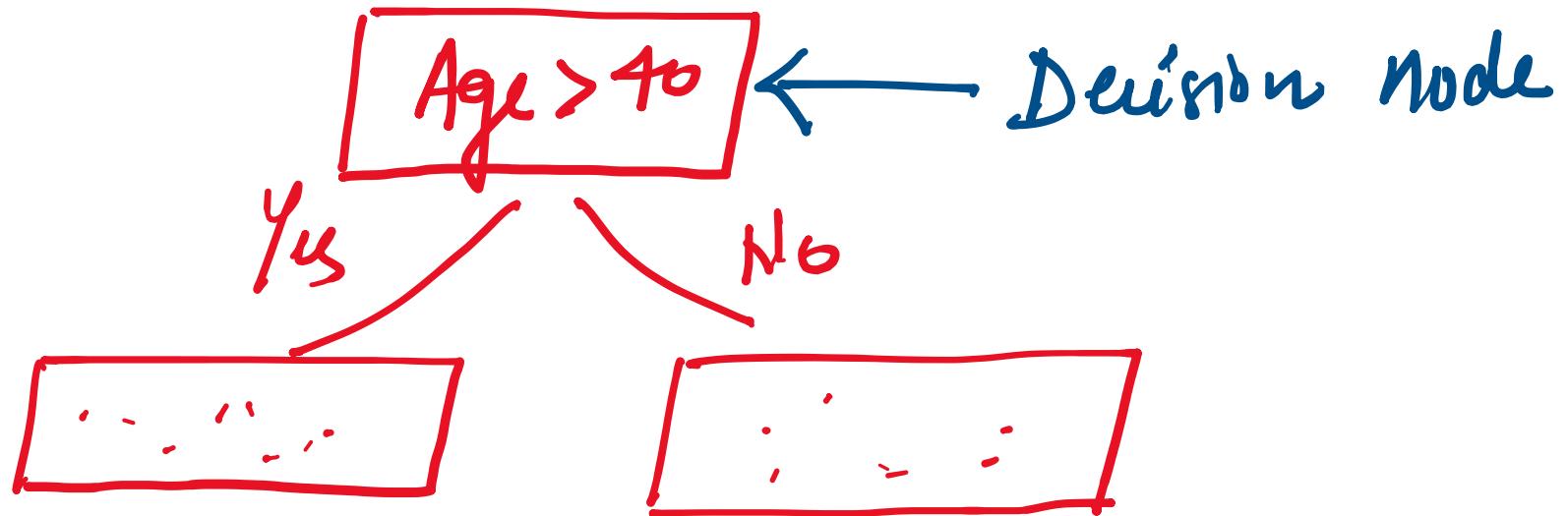


→ Architecture



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

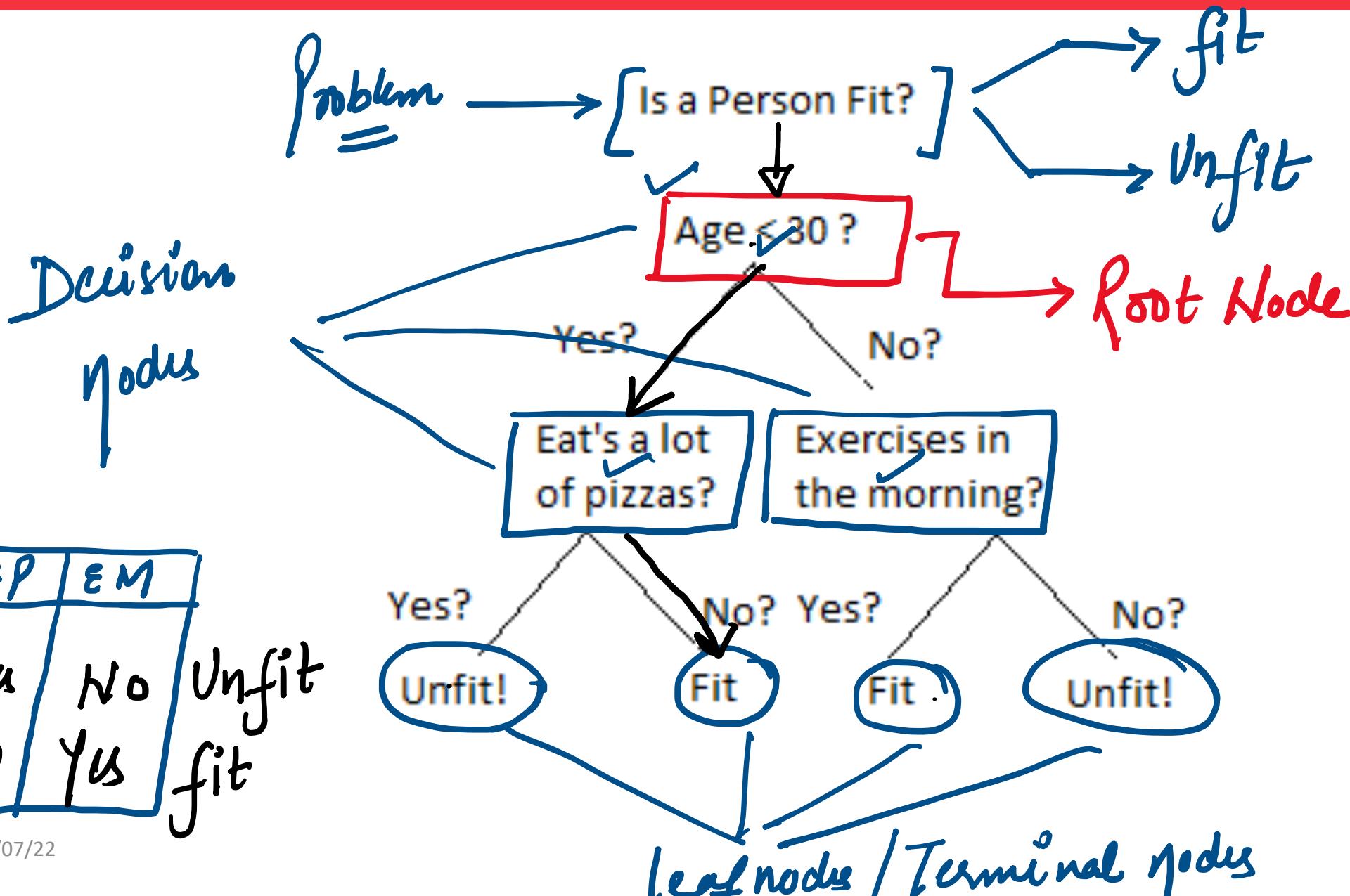
→ Decision Node :- This is the node where we split the data based on some condition on a feature



→ Root Node :- It is very first decision node.

→ Leaf Node or Terminal node :- from where you terminates from the tree and these nodes contain the final output from the tree.

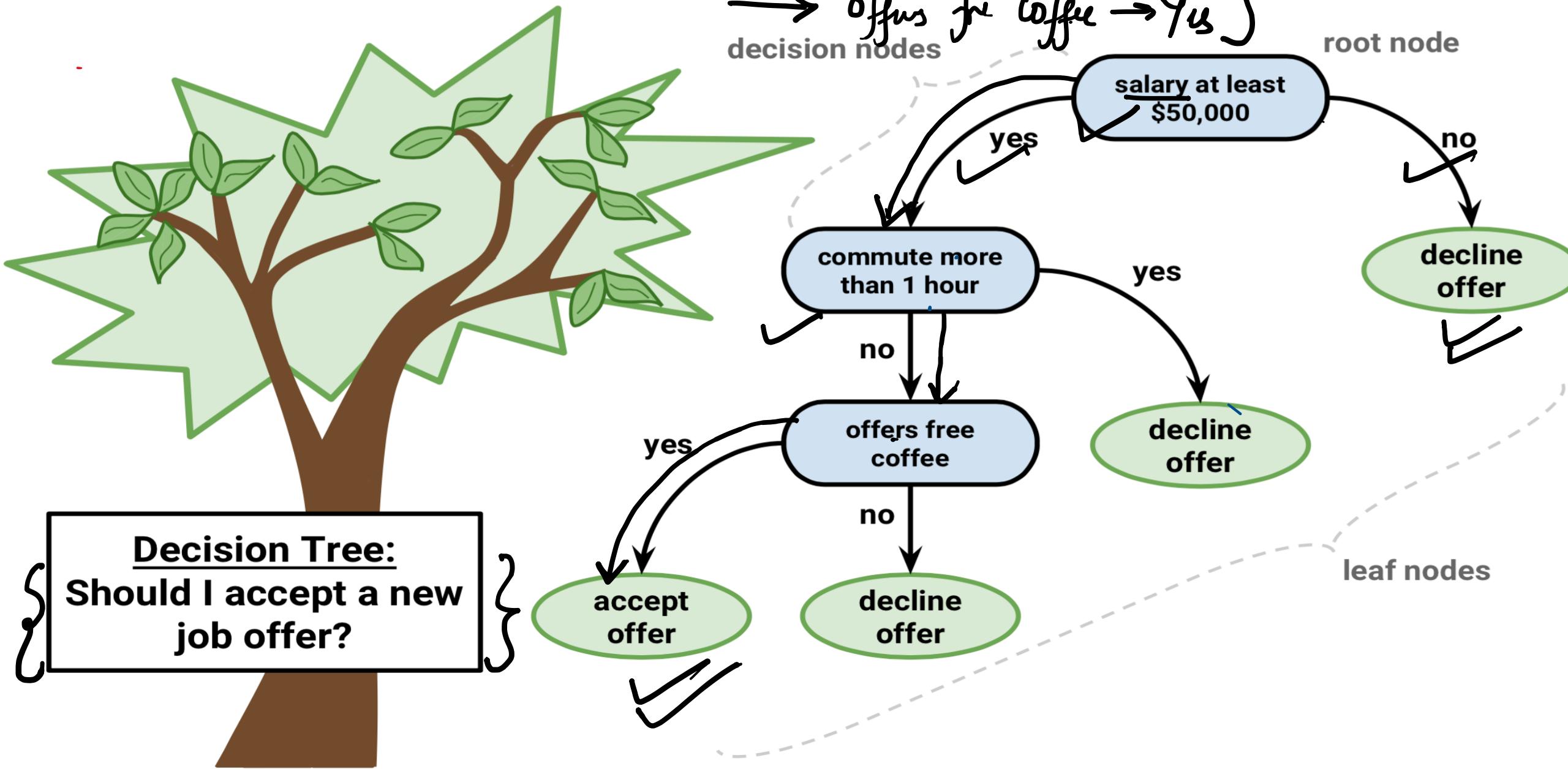
Interpretation of Decision Tree:



Decision Tree Classification:

→ Salary → \$50,000
→ Commute → 35 min
→ Offers free coffee → Yes

upGrad

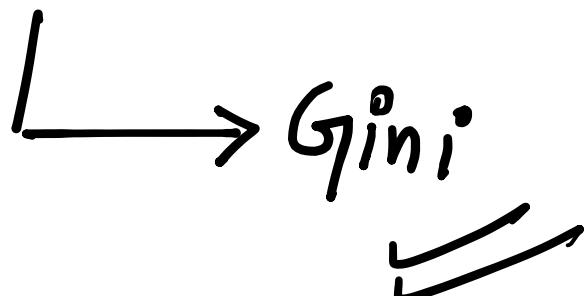


1. Recursive binary splitting/partitioning the data into smaller subsets
2. Selecting the best rule from a variable/ attribute for the split
3. Applying the split based on the rules obtained from the attributes
4. Repeating the process for the subsets obtained
5. Continuing the process until the stopping criterion is reached
6. Assigning the majority class/average value as the prediction

Decision Tree algorithm

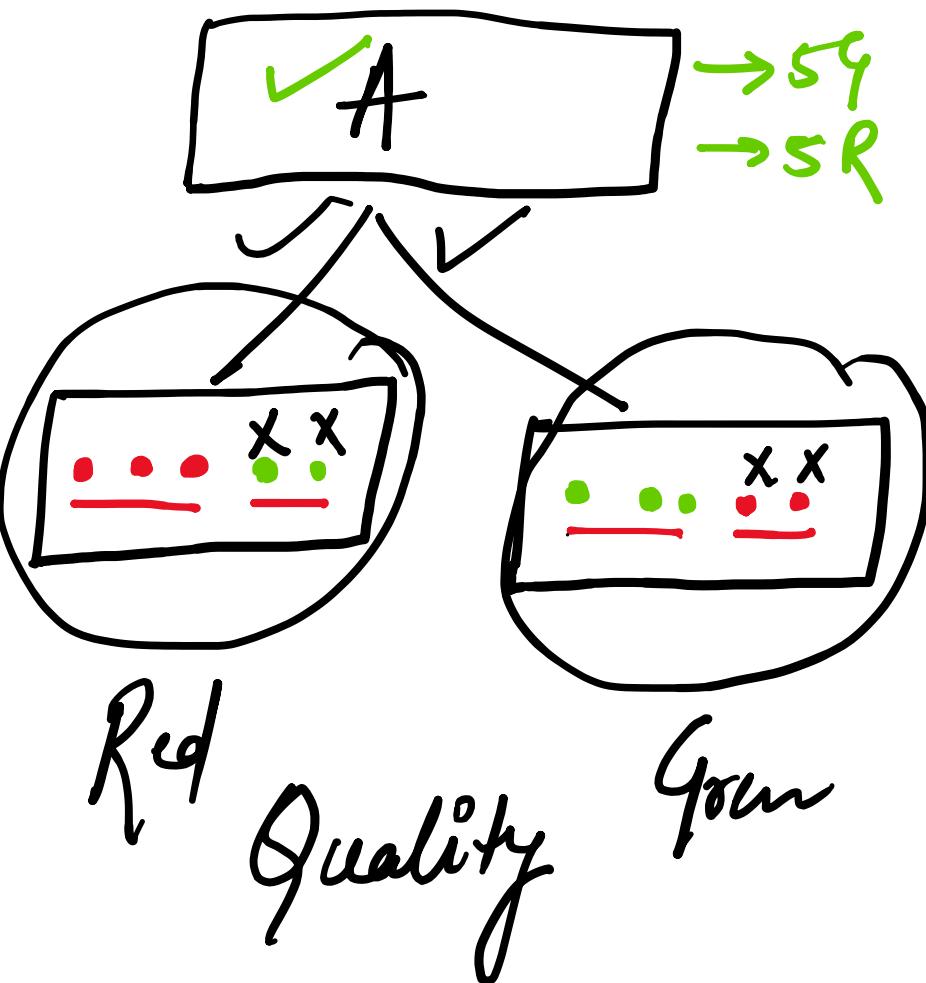
* Recursive selection of features / variables at every decision node using splitting condition and

Selection criteria



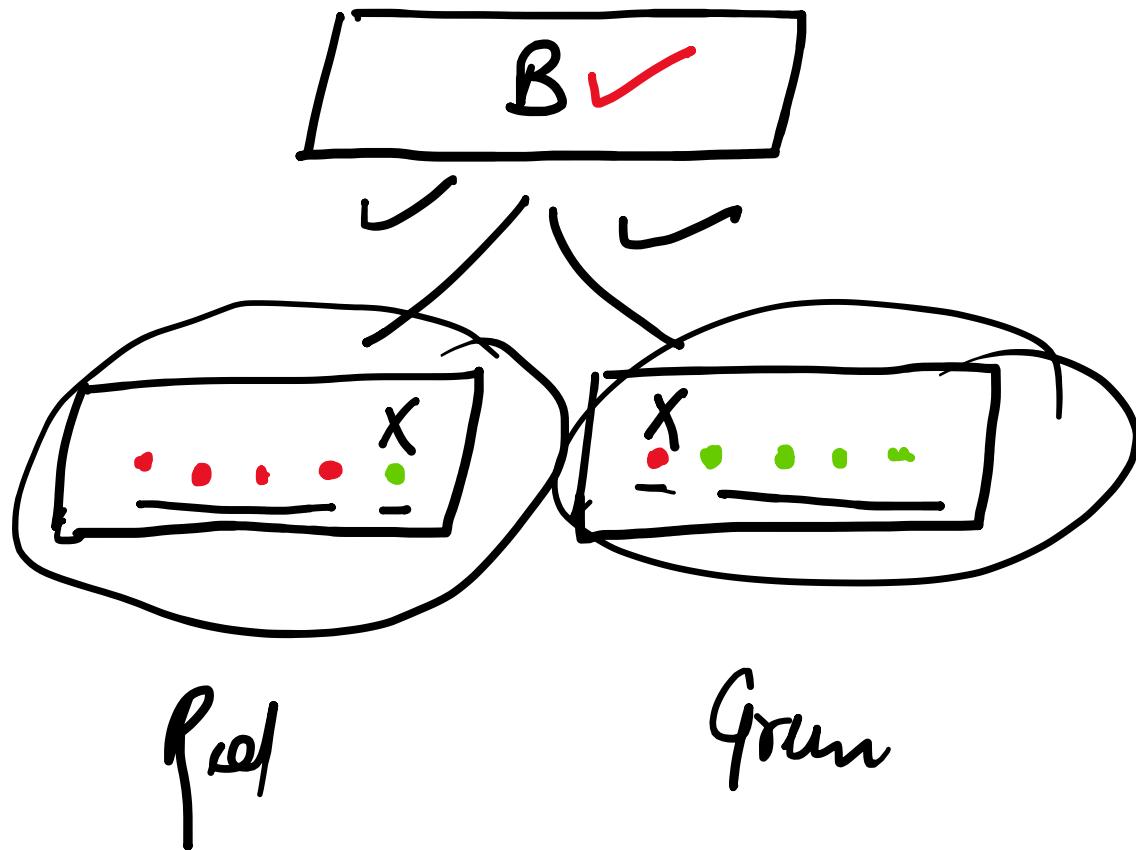
10

4-Misclenfⁿ



10

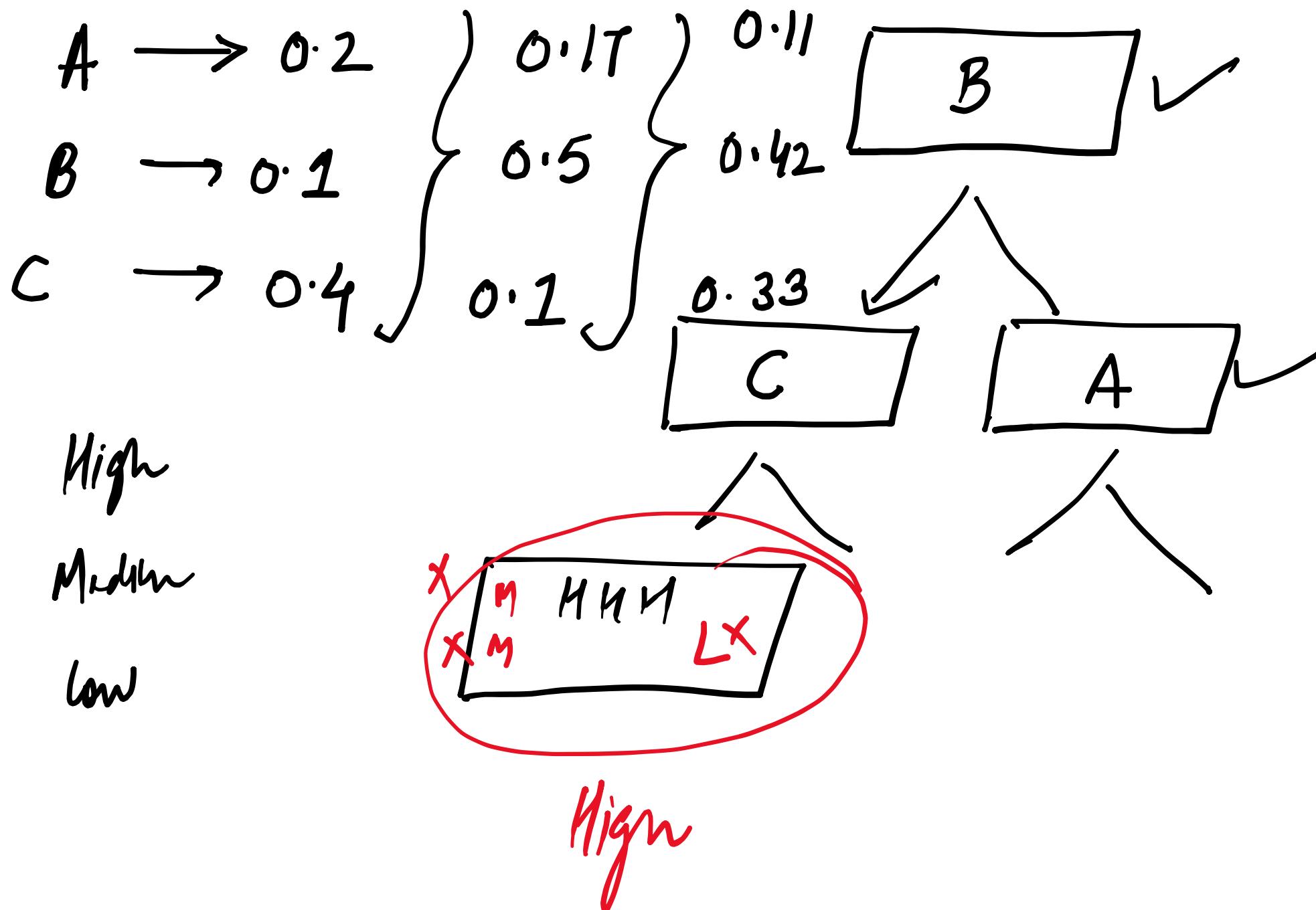
2-Misclenfⁿ



Gini → Statistical method which is used to judge the quality of split.

$Gini \uparrow \rightarrow \text{Quality} \uparrow$

We will be selecting the feature which is giving lower Gini



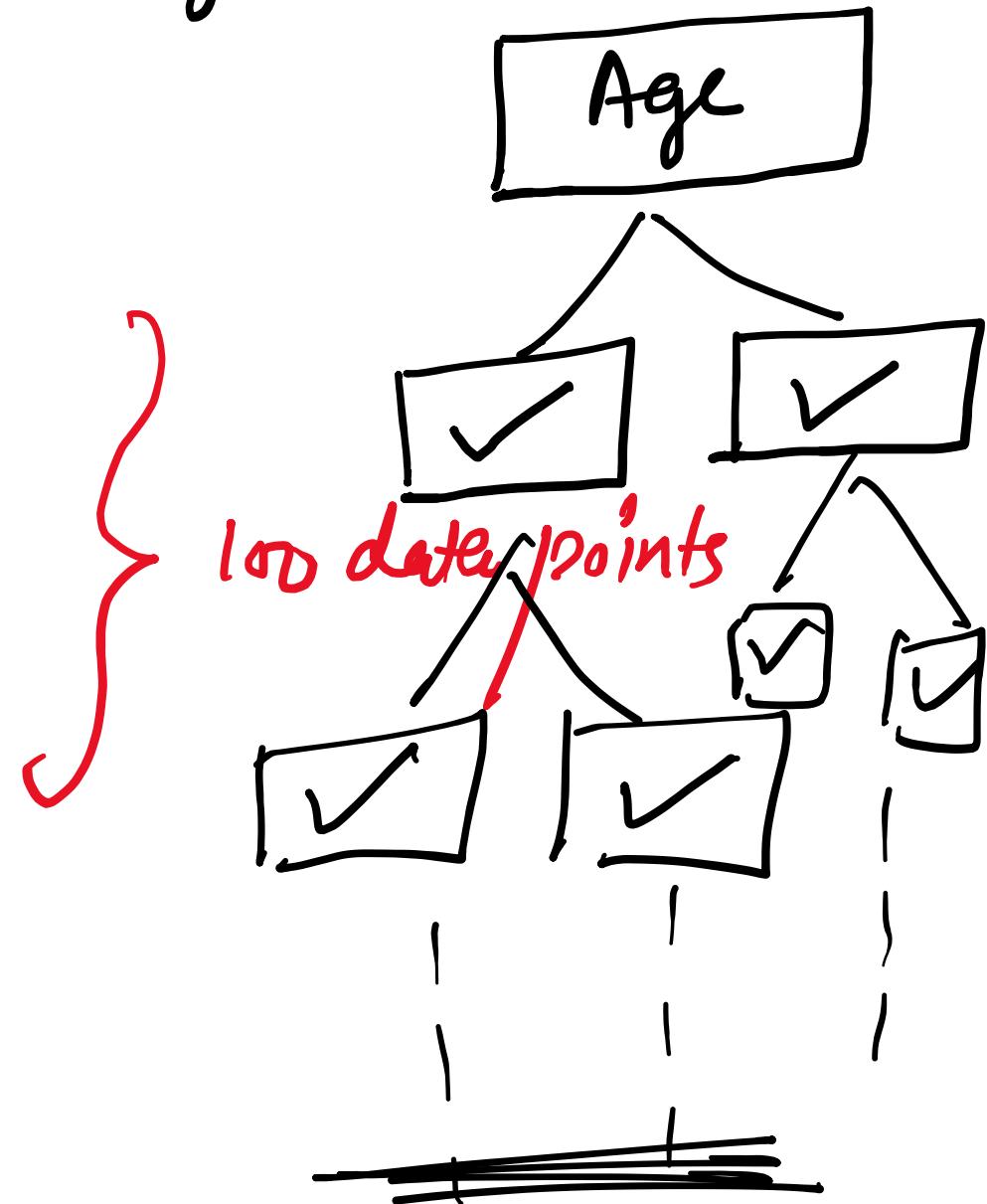
Recall that the Gini index is calculated as follows:

- $G = \underline{\sum_{i=1}^k p_i(1 - p_i)} = 1 - \sum_{i=1}^k p_i^2$

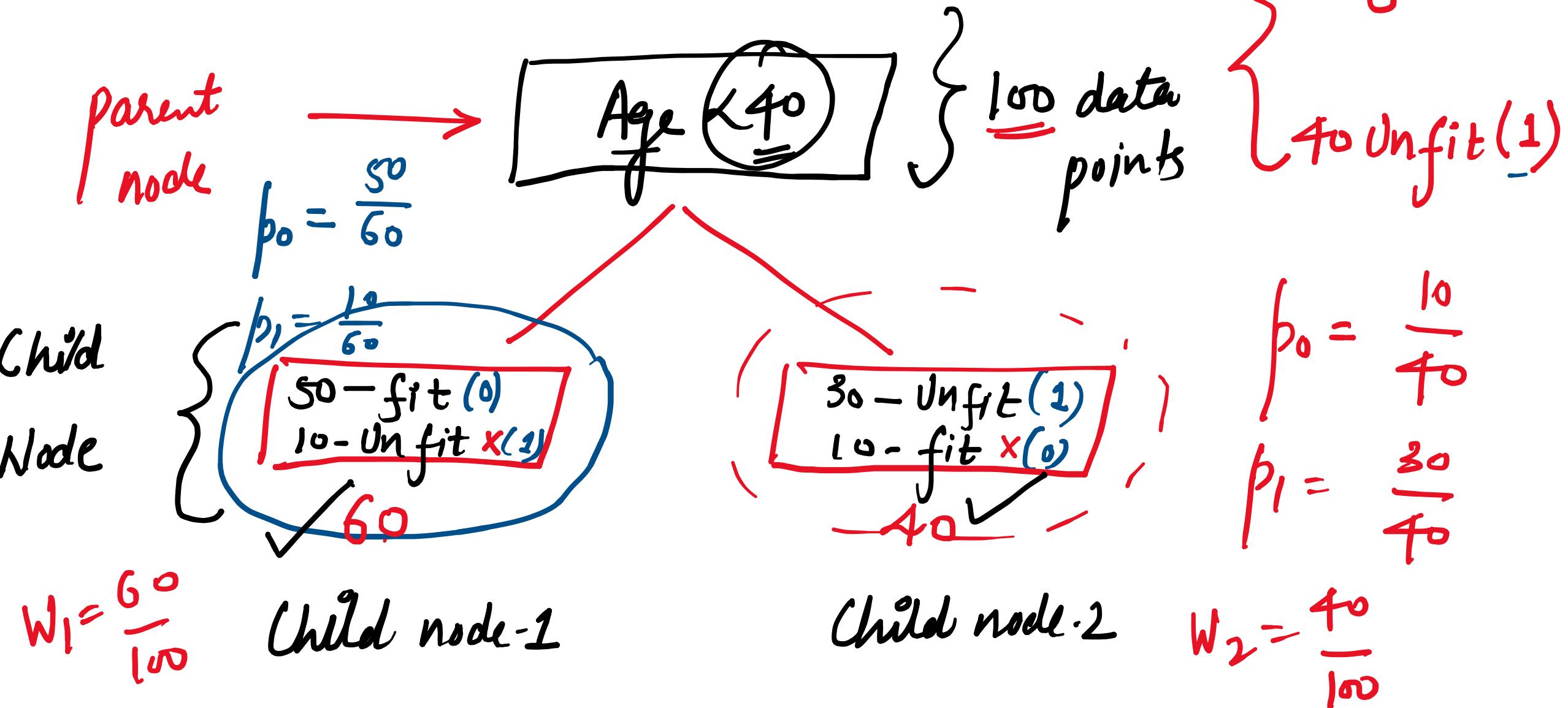
where p_i is the probability of finding a point with the label i , and k is the number of classes.

→ predict whether a person fit or Unfit

Age	Cholesterol	Target
20	70	Fit
55	150	Unfit
..
..



Gini Calculation of Age :-



→ Gini calculation for child node 1 :-

$$G_{C_1} = \sum p_i (1-p_i) = p_0 (1-p_0) + p_1 (1-p_1)$$

$$= \frac{5}{6} \left(1 - \frac{5}{6}\right) + \frac{1}{6} \left(1 - \frac{1}{6}\right)$$

$$= \frac{5}{6} \times \frac{1}{5} \cancel{\times 2} = \frac{5}{18}$$

$G_{C_1} = 0.277$

✓

→ Gini Calculation for child node 2 :-

$$G_{C_2} = p_0(1-p_0) + p_1(1-p_1)$$

$$= \frac{1}{4}\left(1 - \frac{1}{4}\right) + \frac{3}{4}\left(1 - \frac{3}{4}\right)$$

$$= \frac{1}{4} \times \frac{3}{4} \times 2 = \frac{3}{8}$$

$$= \underline{\underline{0.375}} \quad \checkmark$$

→ Gini for split:-

$$G_{age} = w_1 G_{C_1} + w_2 G_{C_2}$$

$$= \frac{60}{100} \times 0.277 + \frac{40}{100} \times 0.375$$

→ $G_{age} = 0.3162$

→ Same process to be repeated for Cholesterol

$$g_{\text{Cholesterol}} = 0.42$$

$$g_{\text{age}} = 0.3162$$

✓ Age (28)

20 → G₁

24 → G₂

28 → G₃

35 → G₄

40 → G₅

 |

 |

✓ Cholesterol (70)

40

45

70 → Min G

75

85

110

125

 |

 |

→ Hyperparameters :-

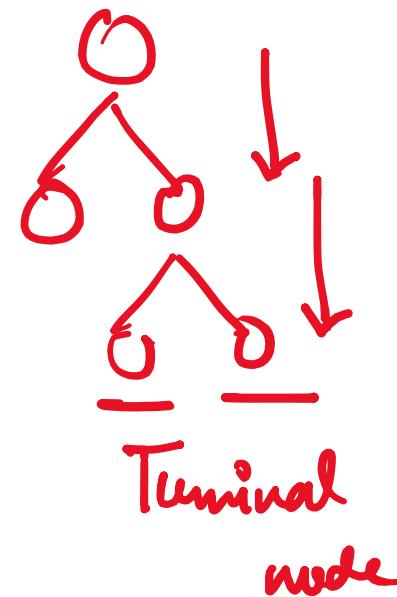
1- Max-depth →

2- Min sample split →

hidden trial

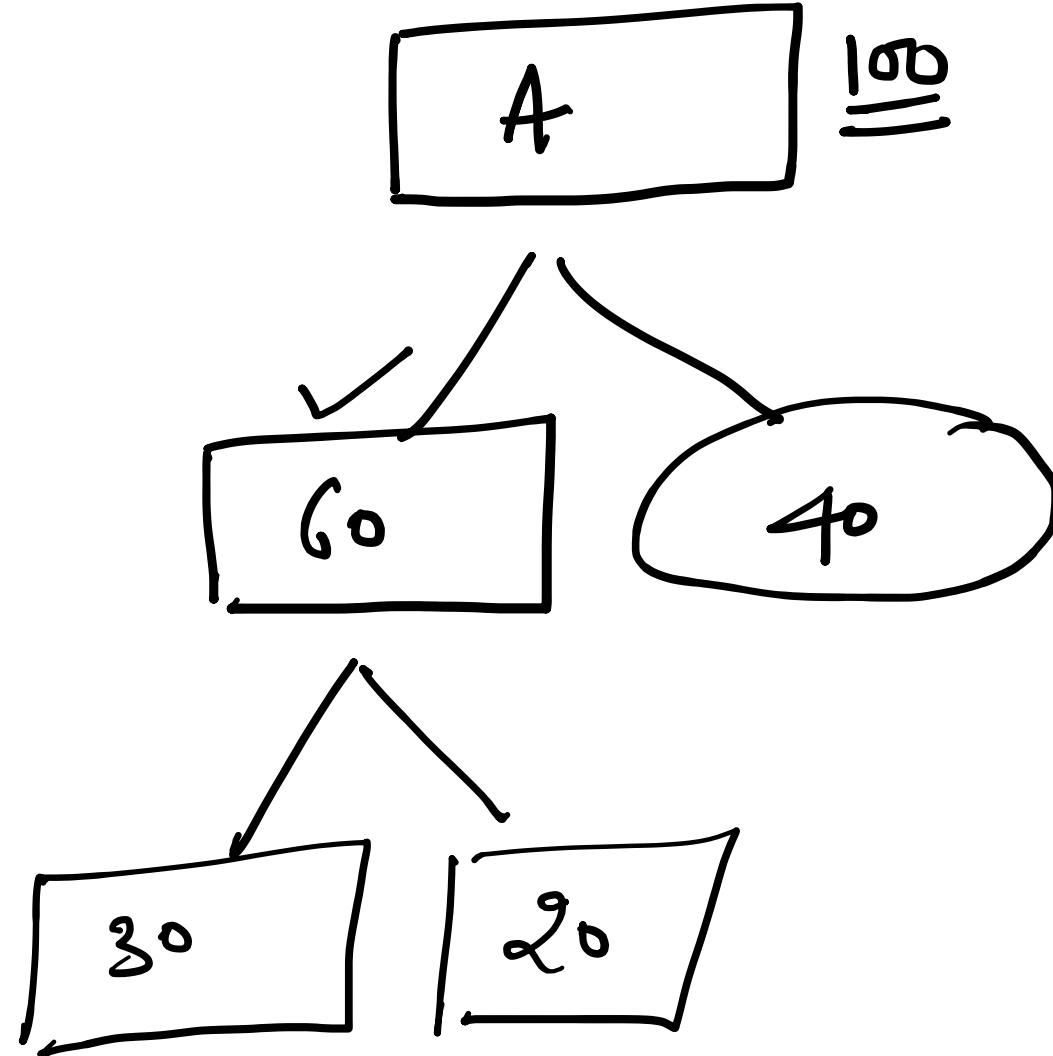
minimum no. of samples

required to make the split



max-depth - 2 ✓

min-sample split → 50 ✓



Doubts !!!

Random Forest Classification:

