



# SNS COLLEGE OF TECHNOLOGY

AN AUTONOMOUS INSTITUTION

COIMBATORE 35

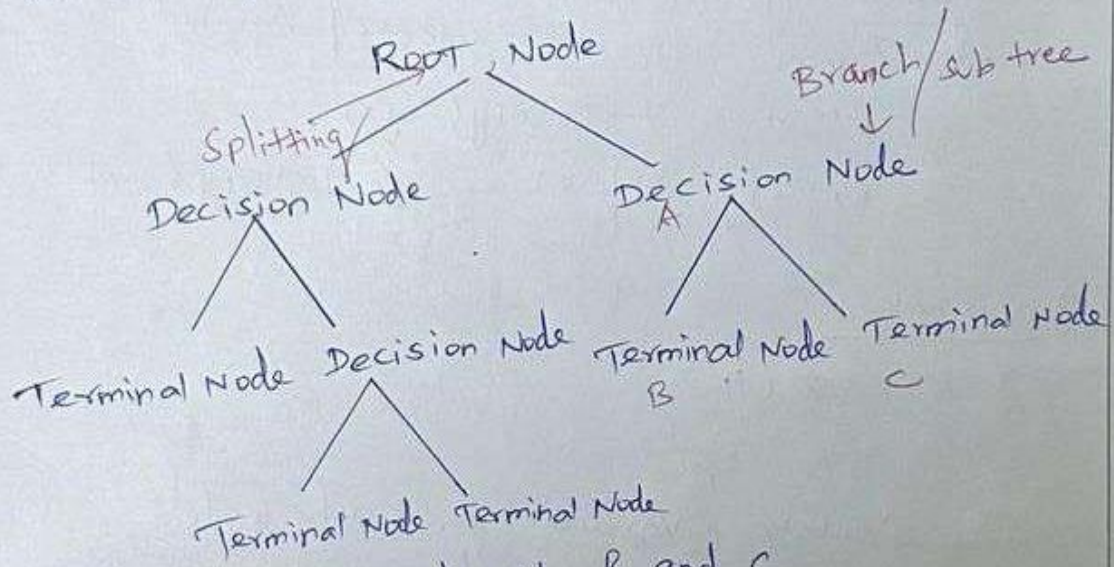


DEPARTMENT OF ELECTRONICS AND COMMUNICATION ENGINEERING

## DECISION TREE BASED METHODS FOR CLASSIFICATION

Terminology Related to Decision Trees

- ① Root Node
- ② Splitting
- ③ Decision Node
- ④ Leaf/Terminal Node
- ⑤ Pruning
- ⑥ Branch
- ⑦ Parent and child Node



A is parent node of B and C

Algorithm Used in Decision Trees

- ↳ ID3 → [Extension of D3]
- ↳ C4.5 → [Successor of ID3]
- ↳ CART → [Classification And Regression Tree]
- ↳ CHAD → [chi-square Automatic Interaction Detection]
- ↳ MARS → Multivariate Adaptive Regression Splines

## IDS Algorithm — Entropy (H) Information gain (IG)

\* Entropy is a measure of randomness in the information being processed. Flipping a coin is an example of an action that provides information that is random.

$$E(S) = \sum_{i=1} -P_i \log_2 P_i$$

play golf

Yes	No	→	= Entropy (S, A)
9	5		= $E(T, X) = \sum_{C \in X} P(C) E(C)$

## \* Information gain

T → Current State X → selected Attribute.

$$= \text{Entropy}(T) - \text{Entropy}(T, X)$$

$$= \text{Entropy}(\text{before}) - \sum_{j=1}^K \text{Entropy}(j, \text{after})$$

## \* Variance

$$\frac{\sum (x - \bar{x})^2}{n}$$

- ① Calculate variance for each node
- ② Calculate variance for each split as the weighted average.

## Pruning

Subdividing the actual training set into 2 sets: training data set, D and validation dataset

