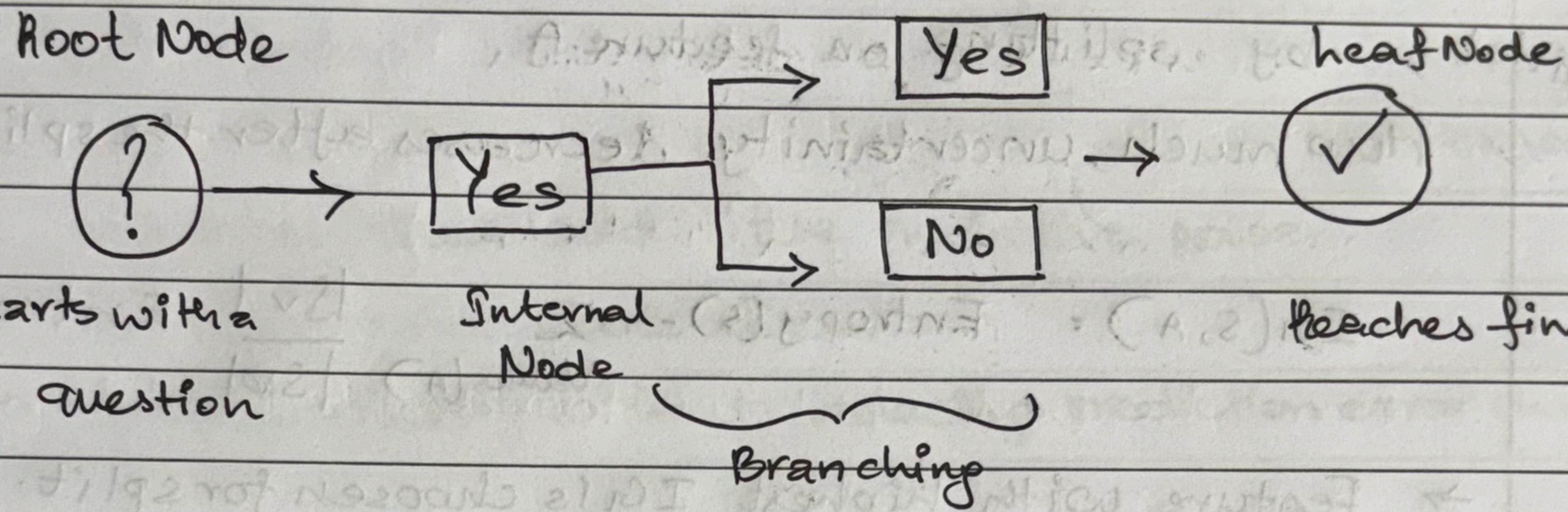


⑨

Decision Trees

A Decision Tree is a flowchart-like model that splits the dataset into branches based on feature values to predict a target variable.



Why?

- ✓ Handles both numerical & categorical data
- ✓ No feature scaling required
- ✓ Non-linear boundaries.

It recursively splits the data using a criterion that measures purity (homogeneity) of that target variable in each node.

→ When building decision tree - which feature and split make the child nodes as pure as possible?

* The splitting criterion is the rule or formula that decides which split to choose - ie, which one reduces the impurity ^{the most}.

Splitting Criterion: Reduction in Impurity

→ Impurity Measures

(a) Entropy, when entropy = 0? perfectly pure node.

↳ disorder / randomness / impurity

$$\text{Entropy}(S) = \sum_{i=1}^k p_i \log_2 p_i [0, 1]$$

ID3
Algos

[Small dataset]

(b) Information Gain: How much impurity is reduced

How good a split is! by splitting on feature A

How much uncertainty decreases after the split

$$IG(S, A) = \text{Entropy}(S) - \sum_{\text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

→ Feature with highest IG is chosen for split.

(c) Gini Impurity: How often a randomly chosen element would be incorrectly identified.

An attribute with lower Gini Index is preferred.

$$G.I. = 1 - \sum_{i=1}^k (P_i)^2 [0, 0.5]$$

CART
Algos

[Large dataset]

- Decision Trees:
1. Start with root node.
 2. Based on impurity find best feature to split
 3. Split the nodes into child nodes
 4. Repeat recursively for each child.

If allowed to grow indefinitely:

A tree will : split till it is pure (impurity = 0)

→ This leads to overfitting

→ We fix this using prepruning and postpruning

- (20) 1. Stopping criteria (Pre-Pruning)
- Prevent the tree from growing too deep or too complex

Some criterias are:

min_samples_split, max_depth, max_leaf_nodes...

- * Balances tradeoff b/w bias & variance
- * Faster, done during training

2. Post Pruning

- CART
- After building a tree fully, we prune it to remove branches that add little predictive power.
 - * Do not improve generalization
 - * Add complexity without reducing prediction error on unseen data

→ Reduced Error Pruning

1. Start with a fully grown tree
2. Remove one node at a time.
3. If removing does not reduce accuracy on validation set
→ keep it pruned.
4. Continue until pruning worsens validation performance.

→ Cost Complexity Pruning (Weakest link pruning)

Used in (CART)

Cost function $R_\alpha(T) = R(T) + \alpha \cdot |T|$

error ↑
 penalty ↑
 no of leaf nodes

- Compute $R_\alpha(T)$ for different α values

- Use cross validation and find α that minimizes val error

(21)

Regression Trees: Predicts a continuous numeric out: y

1. Choose Split based on SSE

$$\text{Best Split} = \arg \min_{j,s} [\text{SSE left} + \text{SSE right}]$$

2. - Weave minimizing the variance within each node.

- After every split child nodes should be more homogeneous in target values than the parent

$$H_{\text{SSE}} = \sum_{m=1}^M \sum_{i \in R_m} (y_i - \bar{y}_{R_m})^2$$

3. Once the best splits are found, the prediction for each terminal node is

$$\hat{y}_{R_m} = \text{mean}(y_i \in R_m)$$

→ Regression minimizes variance

→ Classification minimizes impurity (Maximize IG)

(22)

Categorical Attributes and Splits

↪ split the data based on equality rather than threshold

1. Binary Split (two-way): Divide categories into two groups

- Colour ∈ {Red, Green, Blue}

- Splits: {Red, Green} vs {Blue} or {Blue, Green} vs {Red}

Classification: 1. Compute impurity for left & right groups

2. Compute IG_r

3. choose group with higher / maximum gain

2. Multiway split : we create one branch per category

→ Each branch leads to a subset of data with that value

(ID3,

C4.5)

$$IG_1 = I(\text{Parent}) - \sum_k \frac{|R_k|}{|R|} I(R_k)$$

R_k = subsets of each category

* Overfitting risks - If too many categories

→ Overfitting with high cardinality features (categorical)

- Multiway splits produce thousands of child nodes

- Each node might contain very few samples → overfitting

Sohn:

1. Prefer Binary Splits

2. Use target encoding / one-hot encoding

3. Merge categories with similar class distributions.

→ If you allow multiway splits, nodes with many branches may seem better (higher IG_1)

To fix that C4.5 introduced Gain Ratio

$$\text{Gain Ratio}_k = IG_1 / \text{Split Info}$$

$$\text{where Split Info} = -\sum_i |R_{ik}| / |R| \log_2 |R_{ik}| / |R|$$

This penalizes attributes that create too many splits.

→ Use when categorical variables have many distinct values.

(23)

Missing Values in Decision Trees

Real world datasets often have missing feature values.

- Standard Split cannot be applied if values are missing
- Decision Trees handle missing values naturally

① Imputation (Filling missing values)

Replace missing values with

- Single Imputation:
 - Median
 - Mean (Numerical)
 - Mode (categorical)

* Simple & Fast

* Can bias splits if many values are missing.

② Surrogate Splits (D.T specific)

- Use another feature as backup if/when the primary splitting feature is missing.

1. During tree training:

- For each split, find primary feature and threshold that maximises impurity reduction.
- Look for other features that surrogate that mimic the split
 - ↳ Split the data almost the same way as the primary feature
 - ↳ Compute agreement rate: proportion of samples split the same way.

at a

- 2. If a sample node has missing value in primary feature, follow the surrogate split instead.

- * Even though higher computation \rightarrow no artificial values are generated.