

//_

Decision Trees, $(n \log n) \leftarrow O(N)$

Common metrics:

1. Gini Impurity: $1 - \sum p^2$

2. Information Gain: ~~Entropy of Parent~~

Entropy parent - $\sum (\text{Entropy of split})$

where entropy = $\sum -p \log_2 p$

Stopping Conditions:

1. Max Depth
2. Min Sample per leaf
3. Pruning
4. CV error

~~5.~~

Pseudocode:

func BuildTree(Data, Target)

if stopping-cond: Return leaves

else: if all in same node: return leaves

else: Choose best attribute and split
for each split: BuildTree(split)

Pruning

Pre Pruning: Early stopping conditions like max depth, min split, samples per leaf etc.

Post Pruning: Pruning after fit.

- Cost Complexity Pruning:

Calculate Tree score for each tree and choose the tree with the lowest tree score.

$$\text{Tree Score} = \text{SSR} + \alpha \times \text{Tree Splits}$$

Where α is the tree complexity penalty.

For choosing α : Try various α and choose the one with lowest validation error.