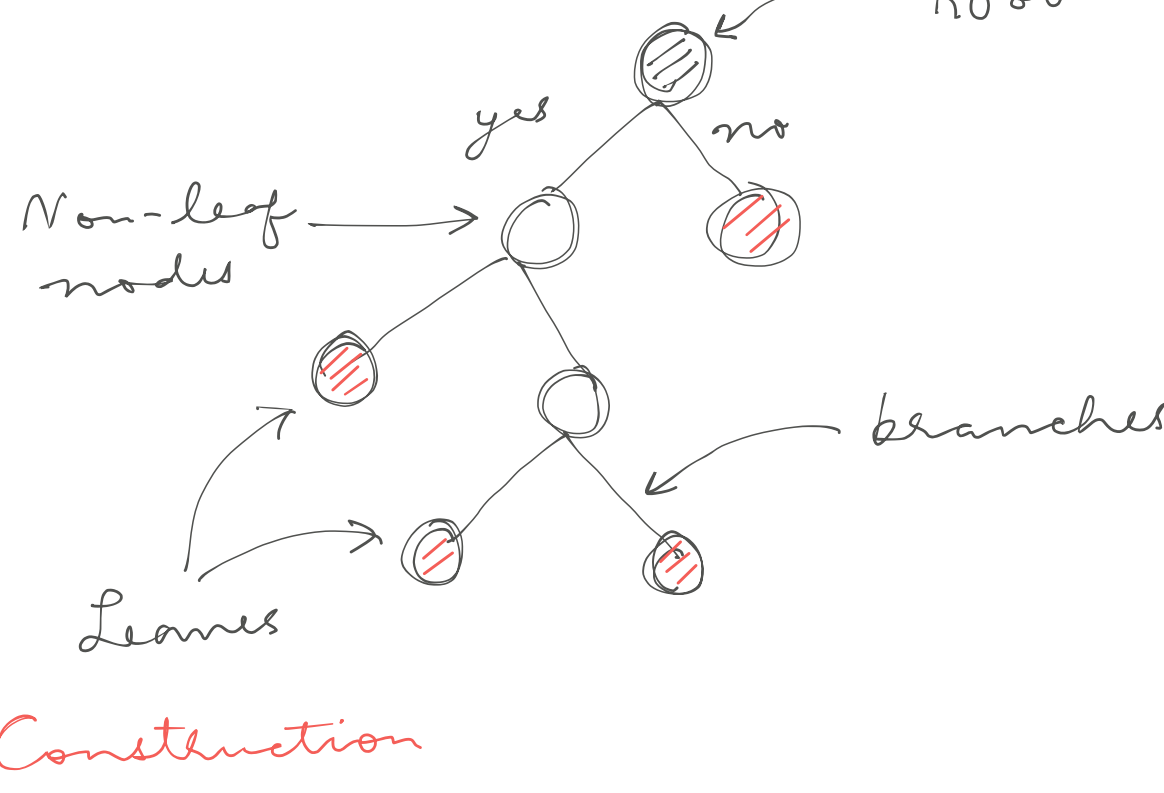


## Decision Tree

Series of yes/no questions  
 ↳ Smartphone

## Terminology



## Construction

Greedy splitting

↳ Classification → Impurity  
 ↳ Regression → MSE

Stopping

100% training / no further improvement

Pre-pruning

↳ max depth

↳ # leaves

↳ # data points in a node

Post pruning

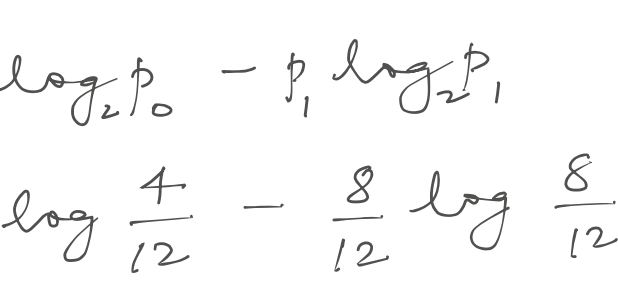
## Impurity functions

1. Entropy → Information gain

2. Gini Index

## ↳ Entropy

$$H(P) = - \sum_c P_c \log P_c \quad c \text{ is a class}$$



$$H(T) = - \frac{4}{12} \log \frac{4}{12} - \frac{8}{12} \log \frac{8}{12}$$

$$= - \frac{4}{12} \log \frac{1}{3} - \frac{8}{12} \log \frac{2}{3}$$

$$= 0.53 + 0.39$$

$$= 0.92$$

$$H(L) = - \frac{4}{6} \log \frac{4}{6} - \frac{2}{6} \log \frac{2}{6}$$

$$= 0.92 \leftarrow$$

$$H(R) = - \frac{0}{6} \log \frac{0}{6} - \frac{6}{6} \log \frac{6}{6}$$

$$= -0 - 0$$

$$= 0$$

Weighted average

$$\rightarrow \frac{6}{12} H(L) + \frac{6}{12} H(R)$$

$$= \frac{1}{2} 0.92$$

$$= 0.46$$

## Gini Index

$$G(P) = 1 - \sum_c P_c^2$$

## Stopping criteria

↳ 100% training acc.

↳ Pre-pruning

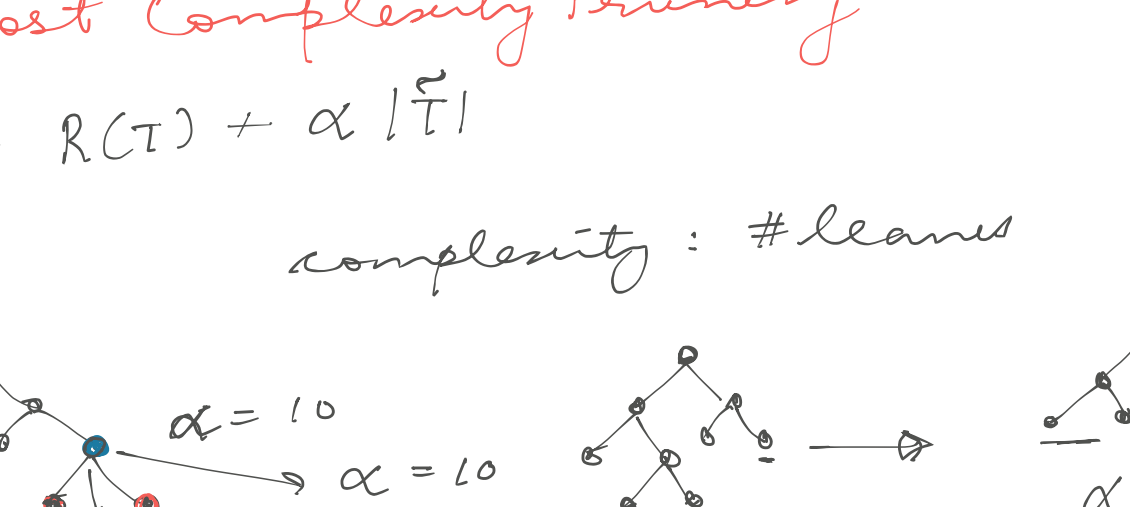
↳ Post-pruning

↳ Validation error

↳ CCP: Cost Complexity Pruning

## Pruning based on val error

Train - Val - Test

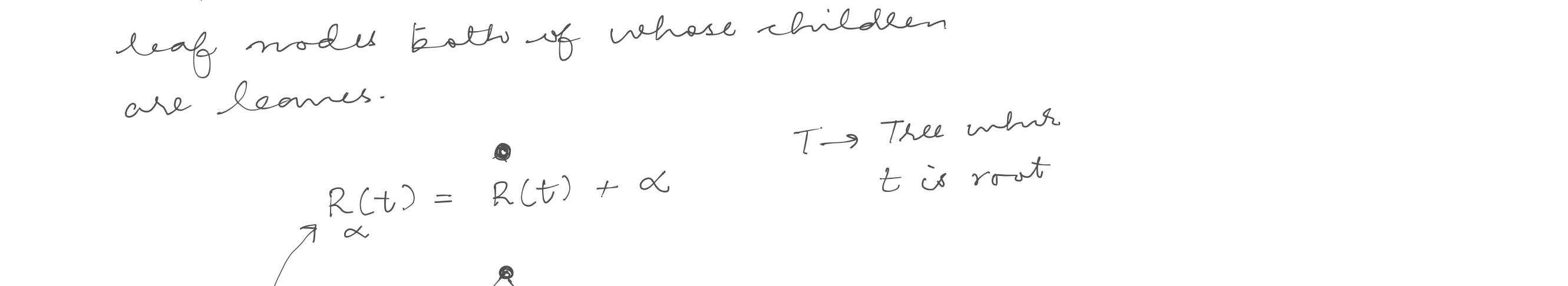


Bottom-up

## CCP → Cost Complexity Pruning

$$R_\alpha(T) = R(T) + \alpha |T|$$

complexity: # leaves



Weighted sum of HCP

↳ RCT

## Total misclassification

1. Train till 100% training acc is achieved.

2. Compute the effective  $\alpha$  ( $\alpha_E$ ) of all non-leaf nodes both of whose children are leaves.

$$R(t) = R(t) + \alpha \quad T \rightarrow \text{Tree where } t \text{ is root}$$

$$R_\alpha(T) = R(T) + \alpha |T|$$

$$R(t) + \alpha = R(T) + \alpha |T|$$

$$\Rightarrow R(t) - R(T) = \alpha (|T| - 1)$$

$$\Rightarrow R(t) - R(T) = \alpha (|T| - 1)$$

$$\Rightarrow \alpha = \frac{R(t) - R(T)}{|T| - 1}$$

if binary & splitting at penultimate level

$$\alpha = R(t) - R(T)$$

3. Prune the leaves of the node that has the minimum  $\alpha_E$ .

4. Repeat this process till only the root node is left. OR Repeat till  $\alpha_E$  is greater than a threshold.

5. Save the pruned tree at each iteration

6. After stopping → Test each pruned tree on the val set.

7. Pick the tree that performs best on the val set.

## Bias / Variance

Bias → Your model is fundamentally incapable of performing well

↑ Var → Overfit

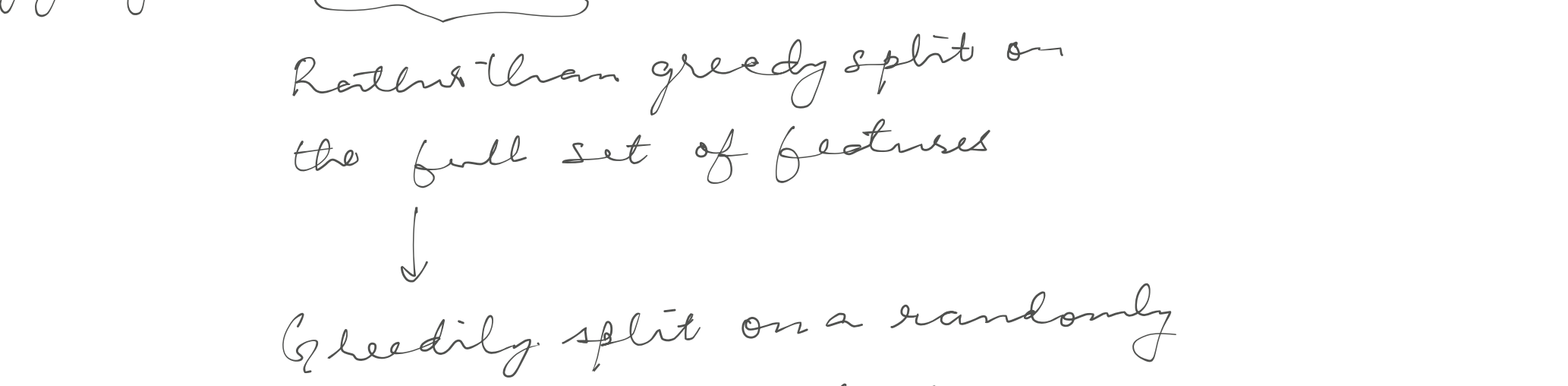
## Bagging → Bootstrapped Aggregation

D → Tr / V / Test

Tr → Tree

↳ Use on test

m = √D



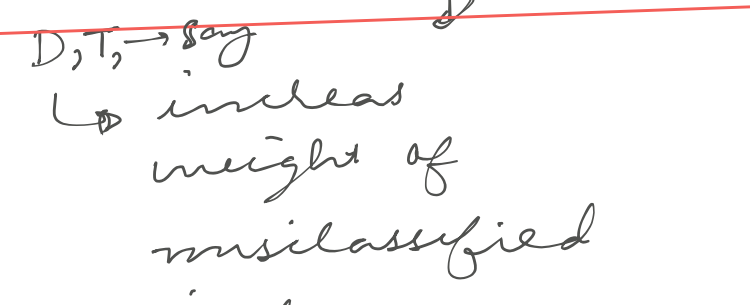
overfit

Aggregating

Voting

ok

## Random Forests



Bagging → T1, T2, ..., Tm

Rather than greedy split on the full set of features

↳ Greedily split on a randomly sampled set of features.

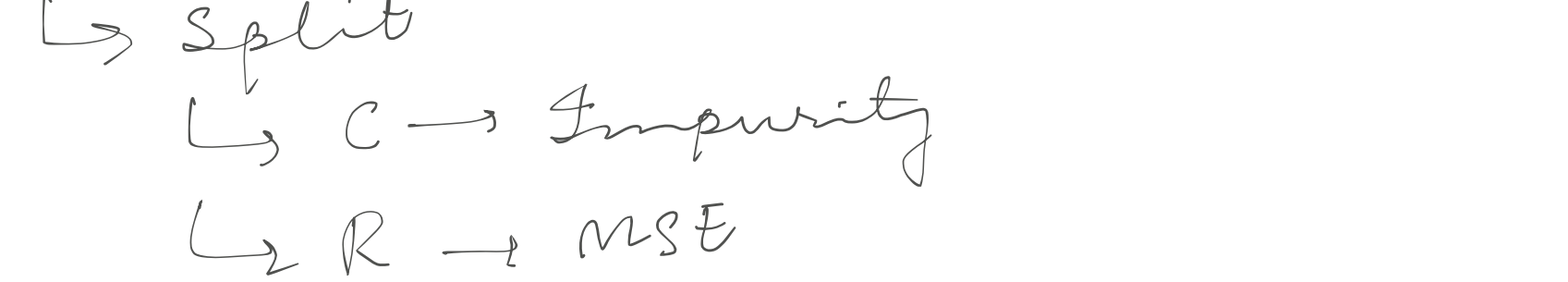
## Bias → Boosting

Combining multiple weak learners to create a strong learner.

↳ Ada Boost

↳ Gradient Boosting

↳ XGBoost



D1, D2, ..., Dm

Tree

80%

increases weights of misclassified instances

↳ reduces the w of others

Each tree gives more weightage to the previous tree's mistakes.

Voting → Soft vote

↳ say

## Summary

↳ Sequence of decisions

↳ Terminology

↳ Construction

↳ Greedy

↳ Split

↳ C → Impurity

↳ R → MSE

↳ stopping

↳ 100% acc train

↳ Overfit

↳ Huge complex trees

↳ Pre-pruning

↳ Post-pruning

↳ Validation error alone

↳ CCP →  $\alpha/|T|$

↳ Bias / Variance

Variance → Bagging → Sampling with replacement

Bias → Boosting