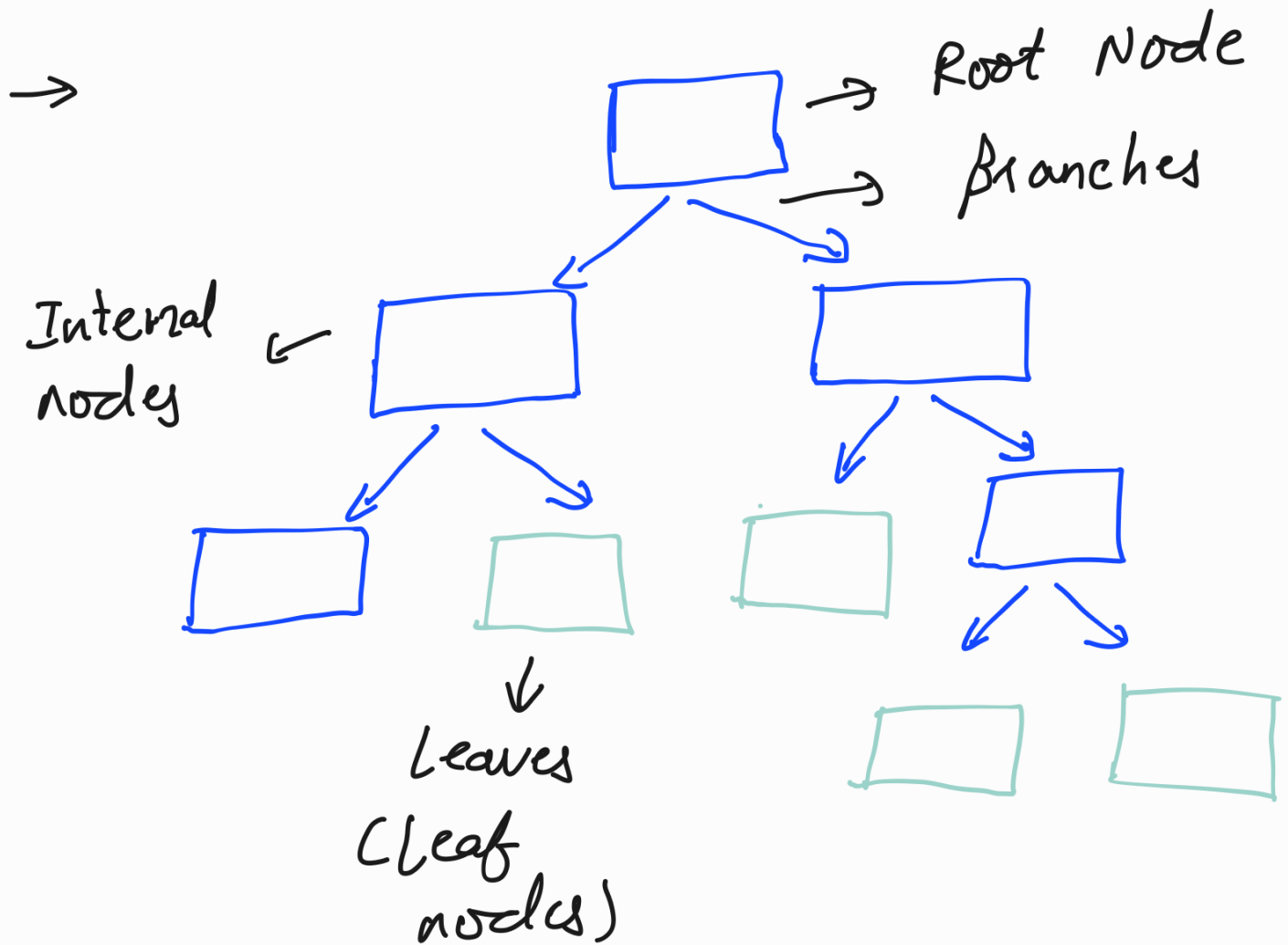→ If the tree is trying to predict a categorical output then it is called classification, if it is predicting a numerical value then it is called regression tree.

→

→ Root Node

→ Branches

Internal nodes ↙

Leaves (leaf nodes)

→ If leaves contain mixture of decision in the leaves then it is called as impure.

→ Ways to quantify impurity

   i) Gini impurity

   ii) Entropy

   iii) Information gain

## Gini impurity of a leaf

$$= 1 - (\text{the probability of yes})^2 - (\text{Probability of No})^2$$

Total Gini impurity = weighted average of gini impurities for the leaves

→ To prevent overfitting we can prune the trees.

→ In regression trees, like how we use gini impurity to decide the root node and subsequent division for internal nodes, in this we use SSR ( Sum of Squared Residuals).

→ As a general rule of thumb min 20 data points should be there to contribute to the split.

→ Tree pruning :

Tree score = $SSR + \alpha T$

$SSR \Rightarrow$ Sum of squared Residuals

$\alpha \Rightarrow$ Found using CV   } Tree complexity penalty

$T \Rightarrow$ Total nr of leaves

→ This process above is known as cost complexity pruning.

# Random Forests

→ create a bootstrapped dataset

→ create a decision tree using the bootstrapped dataset.
(considering a random subset of variables at each step)

→ Bootstrapping the data plus using the aggregate to make a decision is called Bagging.