

- Decision Tree

It is a supervised ML algorithm used for classification and regression.

### → Working Mechanism

Select best feature from dataset.  
Split data based on this feature.  
Repeat the process recursively for each subset.

Stop when
 

- ↳ All samples belong to same class
- ↳ Max depth is reached
- ↳ Min samples per node condition is met.

- DT for classification  
a) Gini Index

$$G_{\text{ini}} = 1 - \sum p_i^2$$

Lower Gini - Better split.

### b) Entropy

$$\text{Entropy} = - \sum p_i \log_2 p_i$$

Measure randomness in data, 0 → pure node

Information Gain

$$\text{Info. Gain} = \text{Entropy}(\text{parent}) - \sum_{n_i} \frac{n_i}{N} \text{Entropy}(\text{child}_i)$$

### • DT for Regression

Criteria → MSE.  $\frac{1}{h} \sum_{j=1}^h (y_j - \hat{y}_j)^2$   
split chosen that minimizes MSE

### • Some challenges of DT

Tend to overfitting. High variance (how much model change with data change).

Small change in data → total different tree

### \* Random forest

It is an ensemble ml alg. that builds multiple decision tree and aggregates their predictions.

It follows Bagging approach → Idea + train models independently on different random samples and avg. their prediction to get final one, it reduce variance & overfitting.

### \* Random forest uses,

• Bootstrapping sampling → technique in which multiple datasets are created from the original dataset with replacement. Each sample is selected randomly, means some data appear multiple times and some not at all purely random.

• Feature Sampling → At each split, consider only subset of feature.

## \* Working of RF

Create multiple bootstrap dataset



Train one decision tree on each dataset



Make prediction

↳ Classification → Majority voting

↳ Regression → Avg of all predictions made by diff DT.

## \* Key Points

- High accuracy as make prediction by ensemble boosting.
- Reduce variance & overfitting.
- Handle large datasets well.
- Higher computation cost.
- Requires more memory.