# Random Forests

A Random forest is a machine learning algorithm that builds many decision trees and then combines their results to make a final decision.

**Advantages :-**
1. High Accuracy : usually better than single trees
2. Handles missing data well
3. Reduces overfitting
4. Works well for both regression & classification

**Disadvantages :-**
1. Slower → many trees → more computation
2. Less interpretable → hard to visualise whole forest
3. Takes more memory

**Usage :-**
1. Finance — fraud detection, loan approval
2. Healthcare — disease prediction from symptom
3. E-commerce — sentiment review
4. Agriculture — crop yield prediction
5. Education — student performance prediction.

**Main Goal :-** Reduce overfitting & improving accuracy

Algorithm with example :-

| ID | Age | Income | Buys-Computer |
|----|-----|--------|---------------|
| 1 | <=30 | H | N |
| 2 | <=30 | H | N |
| 3 | 31-40 | H | Y |
| 4 | >40 | M | Y |
| 5 | >40 | L | Y |
| 6 | >40 | L | N |
| 7 | 31-40 | L | Y |
| 8 | <=30 | M | N |
| 9 | <=30 | L | Y |
| 10 | >40 | M | Y |

Features :- Age, Income
Target :- Buys-Computer

## Step 1

Create Bootstrap Samples
(Random)

Sample 1 → 1, 2, 3, 4, 5, 6, 7, 8, 9, 10   (full dataset)

Sample 2 → 2, 4, 5, 7, 8, 9, 9, 1, 6, 3

Sample 3 → 3, 3, 5, 6, 7, 8, 9, 10, 10, 4

## Step 2

Build Tree 1 (using sample 1), Calculate root Gini for full data

Total Yes → IDs : 3, 4, 5, 7, 9, 10 → 6

Total No → IDs : 1, 2, 6, 8 → 4

$Gini_1 = 1 - \left( \left(\frac{6}{10}\right)^2 + \left(\frac{4}{10}\right)^2 \right) → .48$

$\left\{ Gini = 1 - \left( P_{Yes}^2 + P_{No}^2 \right) \right.$

step 2a :- Try splitting on Age, incomes to decide root node

Calculate ~~root Gini for 31 (Age <=30)~~

~~Age <=30 → 30~~

We'll try splitting on age,

Age divisions → 1. <=30 } → Summarise in 2 categories
                 2. 31-40

                 3. >40

| 1. | <=30 |
|----|------|
| 2. | >30  |

now, Split 1.    Age <=30 vs Age >30

           Left (<=30): IDs = 1,2,8,9    Y=1, N=3

           Right (>30): IDs = 3,4,5,6,7,10    Y=5, N=1

$$\text{Gini}_{left} = 1 - 2\left[(P_{left\ yes})^2 + (P_{left\ no})^2\right] = .375$$

$$\text{Gini}_{Right} = .278$$

$$\text{Gini}_{Age}\ (combined) \rightarrow .375 \times \frac{4}{10} + .278 \times \frac{6}{10} \Rightarrow .316$$

↳ weighted Gini

now, split 2.   a) Income = High vs Income ≠ High

               weighted Gini = .419

       b) Low vs others

               weighted Gini = .45

       c) med vs others

               weighted Gini = .475

Now, what we have done so far is, we calculated root Gini, i.e.
Gini for full data, i.e. .48

& then we have done several splits, based on age & income.

So, the split with the lowest Gini or the one which has the max$^m$ decrease from the root Gini is chosen as the Root Node.

So, best split so far was Age (Gini = ·316)

$$\text{Root} = \text{"Age} <= 30 ?\text{"}$$

ep3 :-

Left Node (Age <=30) : IDs → 1, 2, 8, 9

Y=1, N=3, Gini = ·375

Try splitting on income again for left Node

- High (with Age <=30) → IDs 1, 2 → N & N → Gini = 0

Best split → Income High vs others

Left leafs

High → No

Medium / Low → further split possible

so, Med vs Low,

Med → (ID 8) → No → Gini = 0

Low → (ID 9) → $\overset{\text{Yes}}{\text{No}}$ → Gini = 0

ight Node (Age >30): IDs → 3, 4, 5, 6, 7, 10

Y = 5 N = 1 Gini = ·278

Income: high vs others → G = 0 ✓

med & low can be split again

med → ID (4, 10) → Y & Y, G = 0

low → ID (5, 6, 7) → 2Y & 1N, G → $1 - \left( \left(\tfrac{2}{3}\right)^2 + \left(\tfrac{1}{3}\right)^2 \right)$

→ $1 - \left( \tfrac{4}{9} + \tfrac{1}{9} \right)$ → ·44

<u>Tree 1</u>

Root Gini ⇒ .48        Root Node ⇒ (Age < = 30)?, lowest Gini

[ Root ∴ Age < = 30 ?]

↙ Yes                                    ↘ No

Left Node                              Right Node
(Age < = 30)                          (Age > 30)

[ Income = High? ]                                    [ Income = High? ]
best split.

Yes ↙              ↘ NO                          Yes ↙              ↘ No

Income = High          Income ≠ High, i.e.(med or low)          High          med/L

↓                              ↓                                        ↓                    ↓

NO                    [ Income = Low? ]                    Yes          [Income = m

Yes ↙          ↘ NO                                          yes ↙          ↘N

Income = low          Income ≠ low                          med          lo
                      i.e. med                              ↓            ↓

↓                              ↓

Yes                          NO

The same procedure goes on with other two trees as well
& finally, build a table as

| ID | Tree 1 | Tree 2 | Tree 3 | Majority |
|----|--------|--------|--------|----------|
| 1  | N      | Y      | N      | N        |
| 2  | N      |        | N      | N        |
| 3  | Y      | Y      |        |          |
| .  |        |        |        |          |
| .  |        |        |        |          |
| .  |        |        |        |          |