# Decision Trees → Divide & Conquer → Most widely used
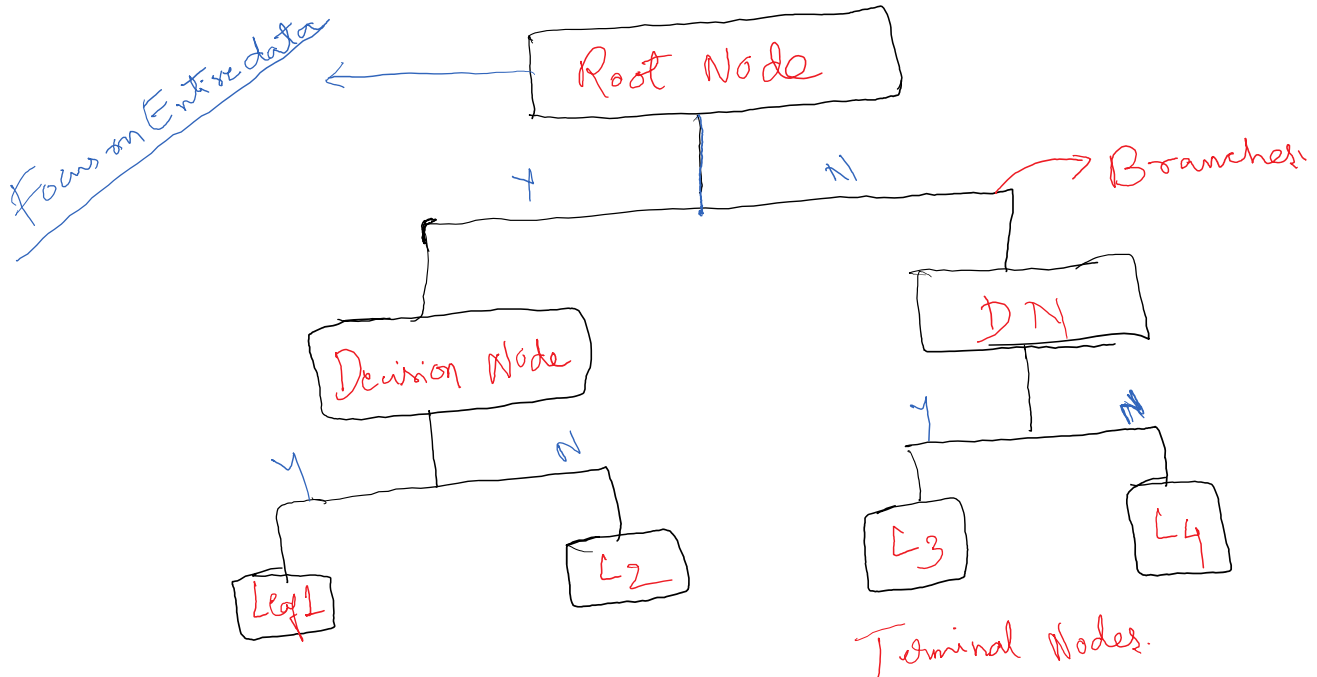
→ Makes Complex choices from simple set of choices & represent their learning in logical form

Focus on Entire data

Root Node

Branches

Y          N

Decision Node

D N

Y          N

Y          N

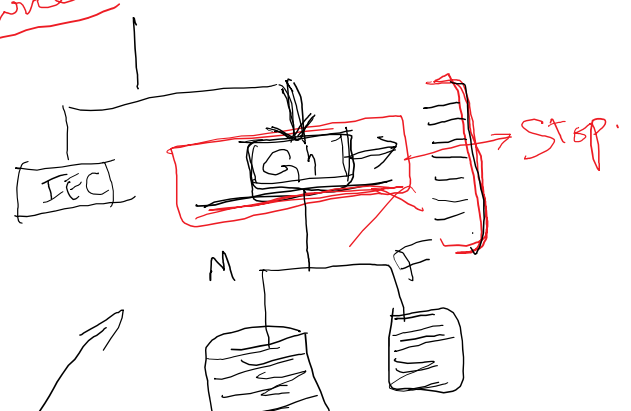Leaf1          L2          L3          L4

Terminal Nodes.

Recursive Partitioning → Enabling a machine to discover/learn Something on its own

→ Splits → Data into Homogenious parts → IN/US

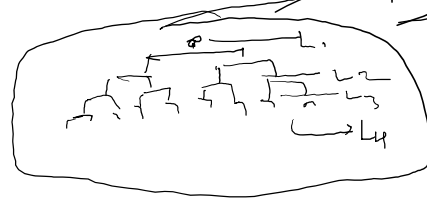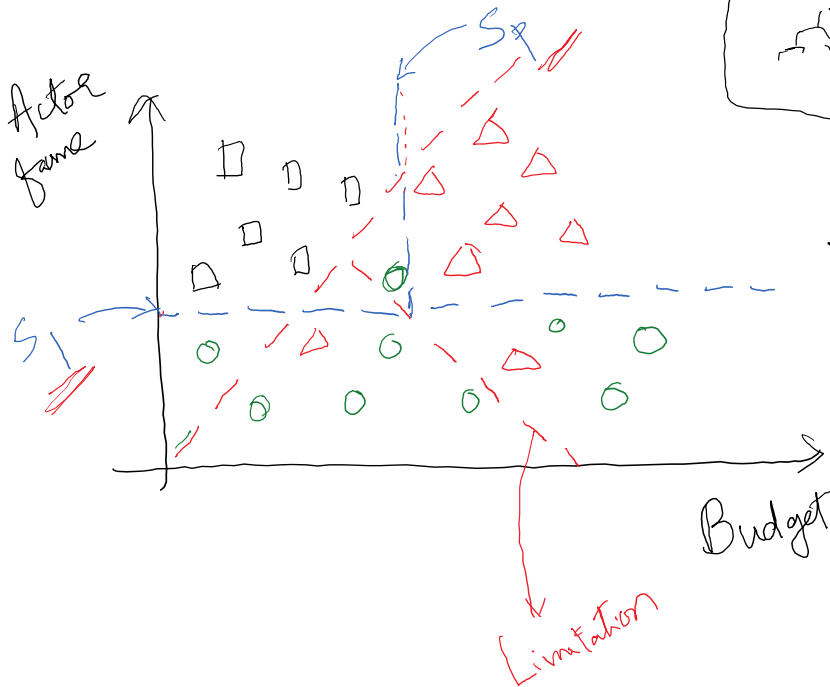→ Terminates → All nodes have same class.

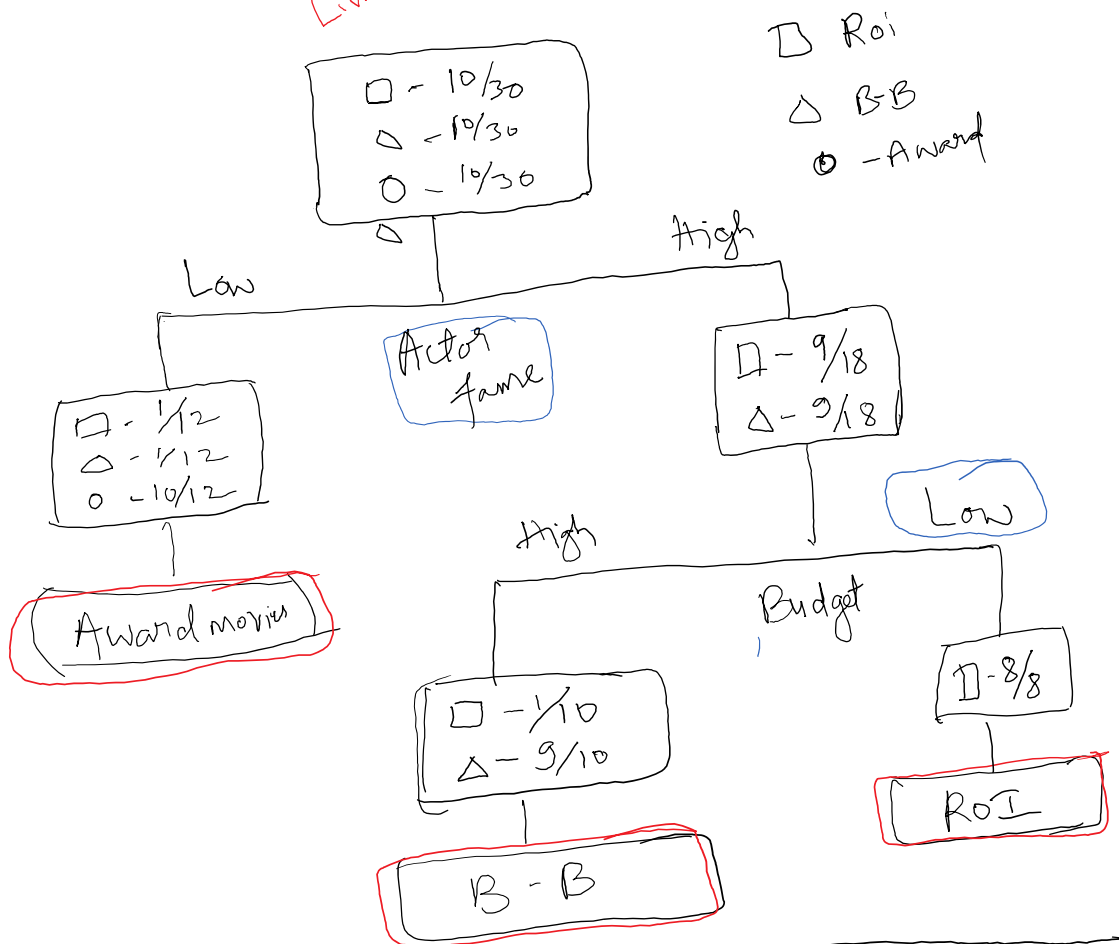C4.5 → Single thread → Open Source
C.5.0
→ Multi thread → Premium

IEC          G1          → STEP.

M          F

→ No features remaining

→ Has grown to a predicted limit → **4 levels**



**Actor fame** (y-axis)

$S_2$

$S_1$

**Budget** (x-axis)

Limitation

□ Roi
△ Blockbuster ⇒
◉ Award. ⇒

⇒ **Axis-Parallel split**

□ Roi
△ B-B
◉ - Award

□ - 10/30
△ - 10/30
○ - 10/30

Low        High

Actor fame

□ - 1/12
△ - 1/12
○ - 10/12

□ - 9/18
△ - 9/18

Award movies

High        Low

High        Budget

□ - 1/10
△ - 9/10

□ - 8/8

B - B        RoI

( C-4.5 → Tree Based ( R-Package RWeka )

C.5.0 Algorithm / C.4.5 → Java Based (R-Package RWeka)

→ Degree to which the subset contains Homogenous elements → "PURITY"

  → All are homogeneous → "PURE" class

→ C.5.0 uses Entropy → Quantifier the Randomness of a set

  Entropy ↑ → Very Diverse
  Entropy ↓ → Very "PURE"

2 possible → Entropy btn 0 & 1
n — h — → —————— 0 to $\log_2(n)$

$$Entropy(S) = \sum_{i=1}^{C} - P_i \log_2(P_i)$$

S → Data Segment
C → no. of class levels.
$P_i$ → Proportion of values falling in a class.

Eg. 2 classes → Red (60%), White (40%)

Entropy = $-0.6 \times \log_2(0.6) - 0.4 \times \log_2(0.4) = 0.97$

NOT PURE

Information gain = $E(S_1) - E(S_2)$

Info gain ↑ & Entropy ↓ →

→ Other forms of splits → Gini index, Chi², gain Ratio.
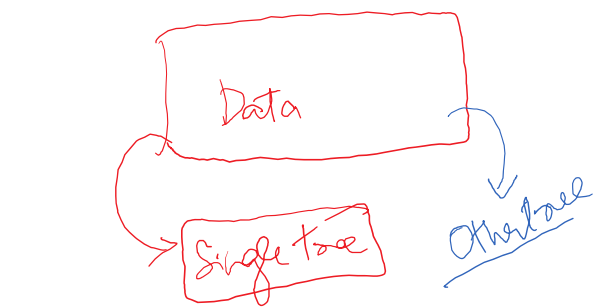
Pruning → Trim → Reduce the size to generalise the data.

    Pre pruning → Early Pruning → get result → stop.
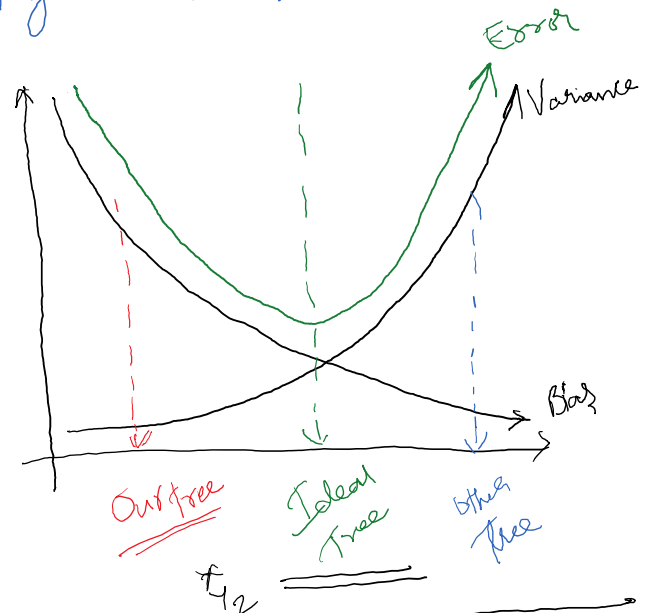
    Post pruning → Allow the tree to grow → Based on need → Trim

---

# Bagging, Boosting & Random Forests

Ensemble method → Group methods.
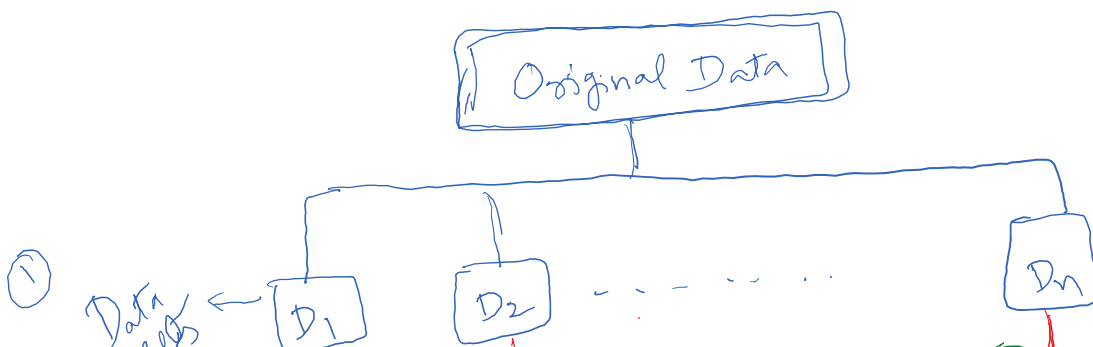
⊛ → Group of predictive models to acheive better Accuracy &
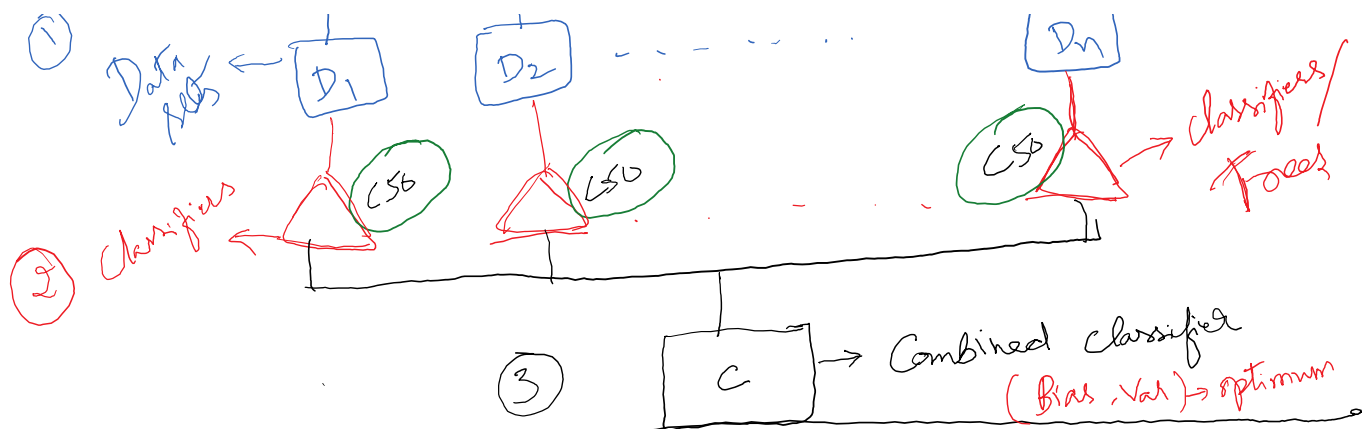    Model stability → It will give a boost to tree models.



$$E = Noise + Bias + Variance$$

---

Bagging    Reduce Variance of the predictors by Combining results
of Multiple classifiers on different Samples of same data.



① Data sets ← $D_1$   $D_2$ - - - - - - -   $D_n$

Classifiers

① Data sets ← [D₁] [D₂] --- [Dₙ]

② Classifiers → (CS0) (CS0) (CS0) → Classifiers/Forest

③ [C] → Combined classifier
(Bias, Var) → optimum

---

① → Sample with Replacement
→ Subset Selection, Shrinkage

② → Use Same classifier.
C.50, RPart, CART

③ → Combine
→ Mean → Mostly
→ Median
→ Modes

---

✳ Random Forest

→ Apply → Regression & also classification

→ Undertakes → Dimensionality Reduction

→ It Treats Missing values, Outliers , Variable transformation

→ Gives the list of Important variables → Contributing O/P.

→ CANNOT Control the mechanism

---

→ Multiple Trees → No pruning

→ Different methods → (CS0, RPART. CART)

⟹ Voting →

Var 1
Vr 2
Vr 3
Vr 1

20        30        45
↓         ↓         ↓
C50       CART      RPART

Nodes ← X-100    X-85     X-96
Variables Y-70   Y-85     Y-80
Partition Z-25   Z-30     Z-92

$V \cdot$

$V \cdot$ ¶

Importre

Partition ← $|2 - 25\rangle$ $|2/ - 30$ $|2| - |$

→ N variables → $\underline{M < N}$ → M → Random Variables

→ Allow the trees to grow fully → No pruning

→ Aggregate data of those N trees

---

Boosting → Refers to a family of Algorithm to make Weak learners To Strong learners.

(Iterative process)

→ classify S/NS

→ Promotional Image → S

→ Only Hyperlink → S

→ You have Won $ ... → S

→ IEC College → NS

→ offer letter Company → NS

→ Bank e-statement → N.S.

→ Weak learning Rule ①

→ Rule ②

→ Rule ③

→ Rule Ⓝ

XGBOOT
GBM

Strong Rule

No of Iterations

Error vs No of Iterations — 50, 100, 1000