

08-11-2025

Agenda:

- Decision tree Regressor
- Ensemble method
 - Bagging & Boosting
- Bootstrap dataset
- Random forest
- code implementation (DT-classifier)

Decision Tree classifier & Regressor

	Classification	Regression
Output	Class Label	Numeric Label
Impurity measure	entropy / gini	Variance / MSE
Leaf prediction	Majority class	Mean of values
Split criteria	Information gain	Variance reduction

Data set

ID	Square feet	Price (₹ in lakhs)
1	800	50
2	850	52
3	900	54
4	1500	90
5	1600	95
6	1700	100

Step 1: Compute parent variance (price)

$$\begin{aligned}
 \mu(\text{price}) &= 73.5 & \text{var} &= \frac{1}{n} \sum (y_i - \mu)^2 \\
 &&&= \frac{1}{6} \sum ((52 - 73.5)^2 + (\dots))^2 \\
 &&&= \frac{2831.5}{6} \\
 &&&= 471.9
 \end{aligned}$$

Step 2: For each feature → Compute below

- we have only one feature:
- we pick the "square feet" column :
- sort values!

800, 850, 900, 1500, 1600, 1700

- split point (midpoint between value):

825, 875, 1200, 1550, 1650

→ we test & compute variance reduction (VR)

split at 1200:

↙ square test ≥ 1200 ↘
left right

800 50

850 52

900 54

1500 90

1600 95

1700 100

$$\mu(\text{left}) \rightarrow 52$$

$$\mu(\text{right}) \rightarrow 95$$

$$\text{var}(\text{left}) \rightarrow 8/3 \rightarrow 2.67$$

$$\begin{aligned}\text{var}(\text{right}) &\rightarrow 50/3 \\ &\rightarrow 16.67\end{aligned}$$

weighted variance:

$$\begin{aligned}\text{var}_{\text{split}} &= \frac{3}{6} \times 2.67 + \frac{3}{6} \times 16.67 \\ &= 1.335 + 8.335 = \underline{\underline{9.67}}\end{aligned}$$

$$471.9 \rightarrow 9.67$$

$$\text{variance Reduction: } 471.9 - 9.67$$

$$= 462.23$$

Sort < 1200 ?
 Yes No

800 50
 850 52
 900 54

1600 90
 1600 95
 1700 100

Step 1 :

midpoint: 825, 875

split 825:

<	>
800, 50	850, 52
↓	↓
var(left)	var(right):
l	r
0	1

midpoint: 1550, 1650

split = 1550

<	>
1500, 90	1600, 95
↓	↓
var(left)	var(right)
↓	↓
0	6.25

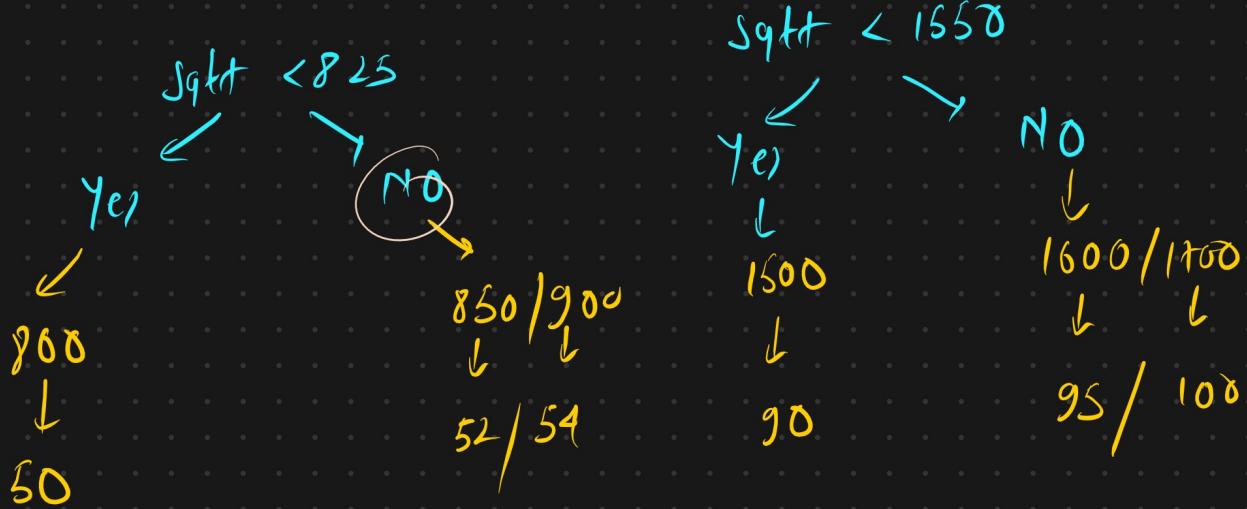
$$\omega \cdot \text{variance} = \frac{1}{3} \times 0 + \frac{2}{3} \times 1 \\ = \frac{2}{3} = \underline{\underline{0.67}}$$

$$\omega \cdot \text{variance:} \\ \frac{1}{3} \times 0 + \frac{2}{3} \times 6.25 \\ = 4.16$$

$$VR = 2.67 - 0.67 \\ = \underline{\underline{2.0}}$$

$$VR: 16.67 - 4.16 \\ : \underline{\underline{12.51}}$$

$Sqft < 1200$?
 Yes No



$$\text{min_samples_leaf} = 1$$

$$\text{min_samples_split} = 2$$

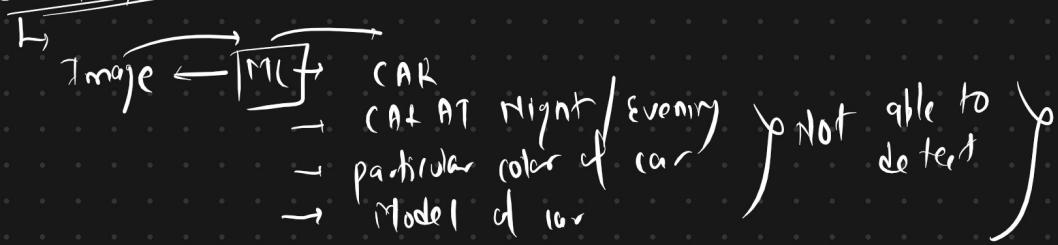
A node must have at least 2 samples
to even consider a split

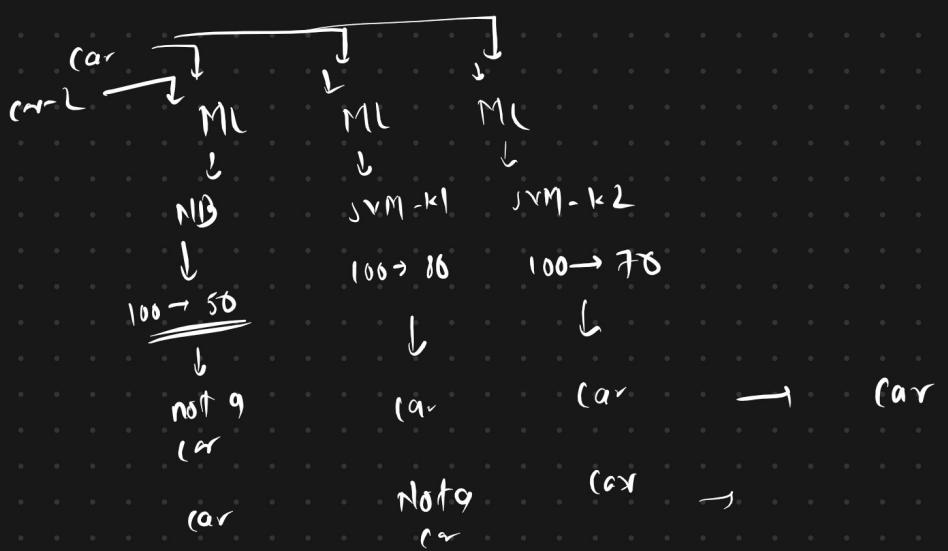
▲ Disadvantages

- ✗ Overfits easily (deep trees)
- ✗ Unstable (small data changes → big tree changes)
- ✗ Not smooth — predictions are piecewise constant

ensemble method | learning | technique:

car detection / classification



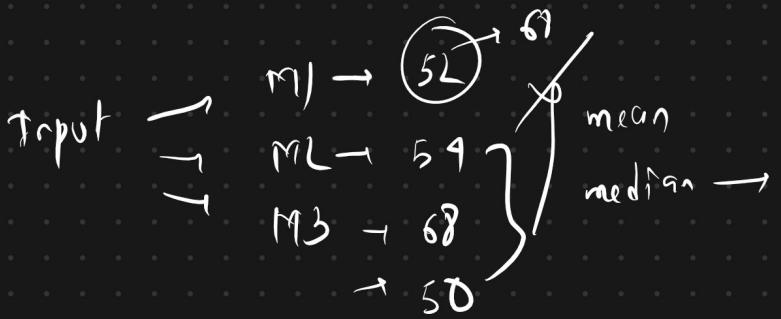


Predict → value of a property → continuous value.

linear regression (5 feature) → ML → 96% → 41.

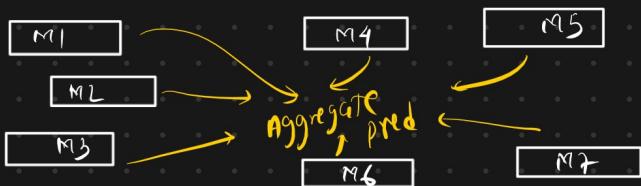
linear regen (8 feature) → ML → 96% → 41.

SVM-k2 (5 fea.) → ML → 97% → 48.



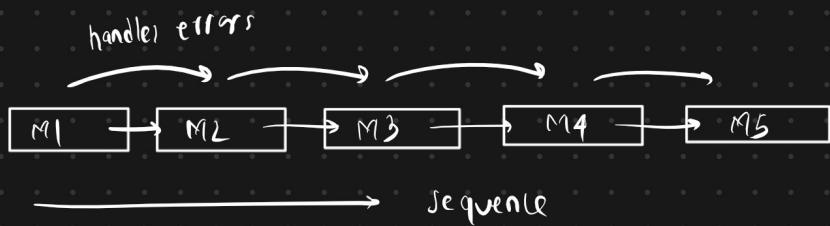
Technique:

Bagging

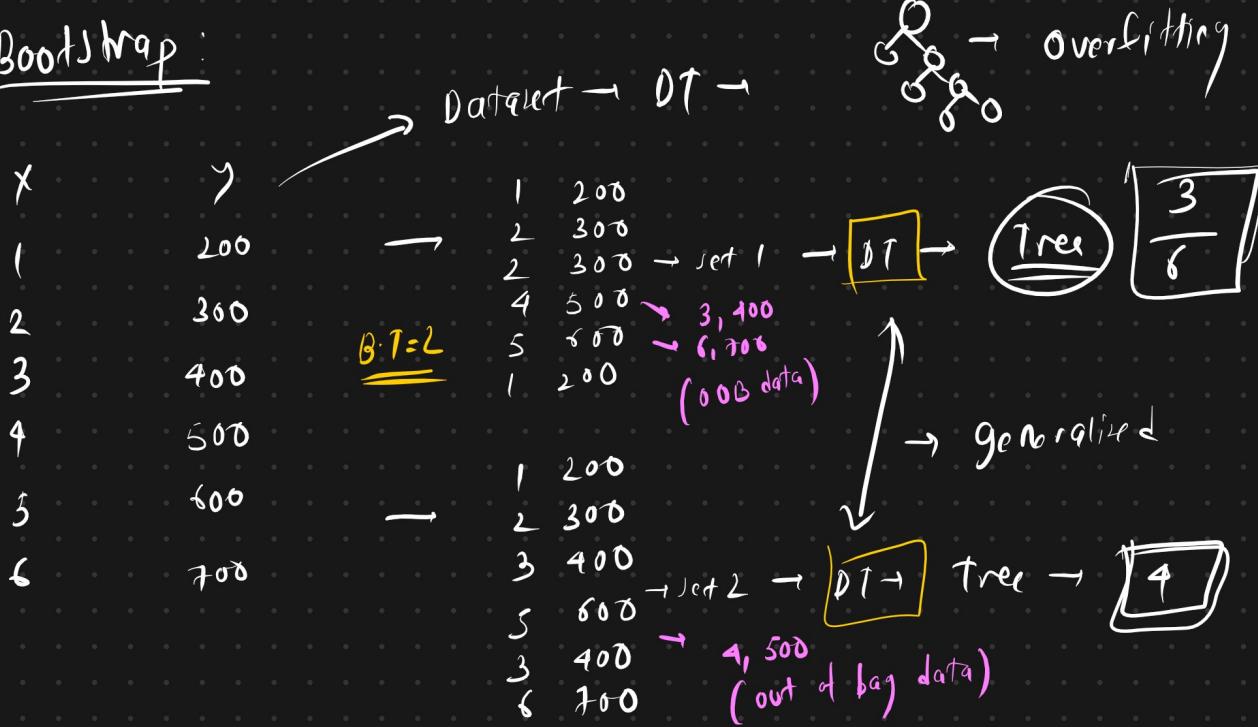


Aggregate → majority
→ average
→ median

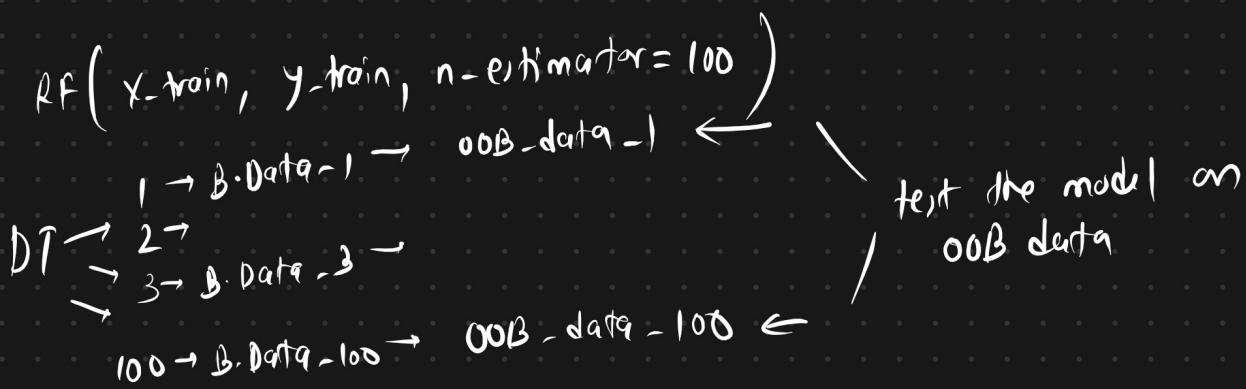
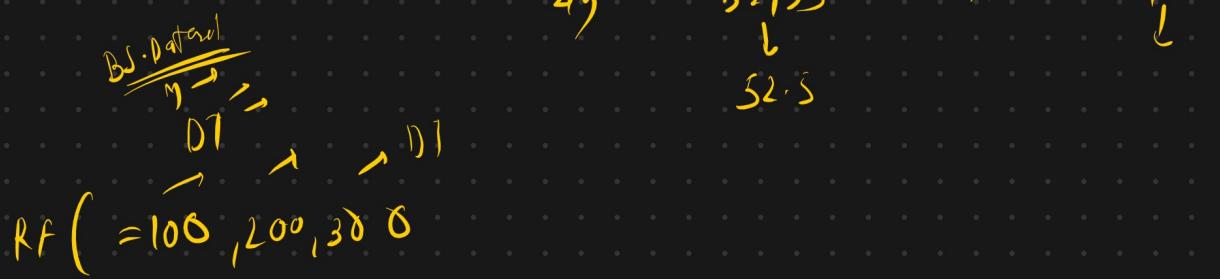
Boosting



Bootstrap



when duplicated are allowed : sampling with replacement



If it's a self-validated model.

Random Forest step by step

max_features = 2

: each split consider a random subset of 2 features.

ID	sqft	Bedroom	Age (year)	Price (lakhs)
A	800	2	10	60
B	900	2	8	70
C	1000	2	7	80
D	1500	3	3	90
E	1600	3	2	95
F	1700	4	1	100

n_estimation = 3, max_features = 2

: each split in a tree will randomly consider 2 of the 3 features.

Step 1: Create Bootstrap samples.

Tree 1 Bootstrap dataset

B(900, 70)

C(1000, 80)

A(800, 68)

B(900, 70) ← repeated

E(1600, 95)

F(1700, 100)

Tree 1. OOB: D(1500, 90)

Tree 2 Bootstrap Data:

- Data
- OOB

Tree 3 Bootstrap Data:

- Data
- OOB

Step 2: for each tree:

- random feature selection [sqft, Bedroom, Age]
- sqft, Age (2 feature)

→ calc parent variance

→ sort data & find mid point.

→ for each mid point:

→ split from mid point

→ calculate variance of right & left side

→ calc weight variance

→ calc variance reduction

→ find the best mid point using the highest variance reduction.

→ final prediction

$$D \mapsto f_i(x) \rightarrow \hat{y}$$

$$\hat{y}_{RF}(x) = \frac{1}{3} \sum_{i=1}^3 f_i(x)$$

$$\hat{y}_{RF}(x) = \frac{1}{3} (97.5 + 93.75 + 90)$$

$$= \underline{\underline{93.75}}$$

$$sqft = 1200$$

$$Bedroom = 3$$

$$Age = 4$$

↓

$$f_1(x) : 97.5$$

$$f_2(x) : 93.75$$

$$f_3(x) : 90$$

