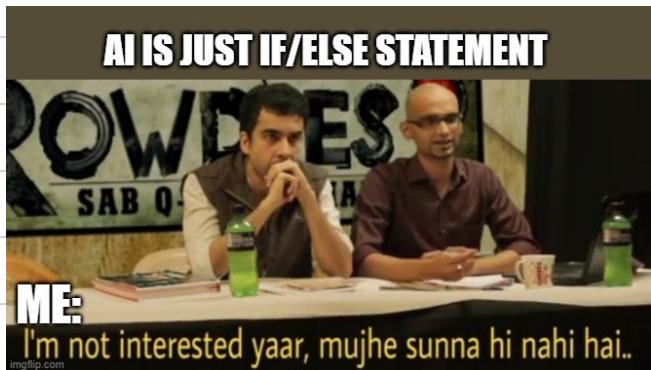
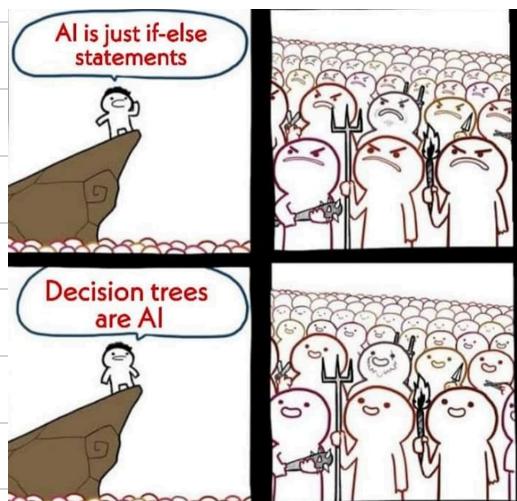


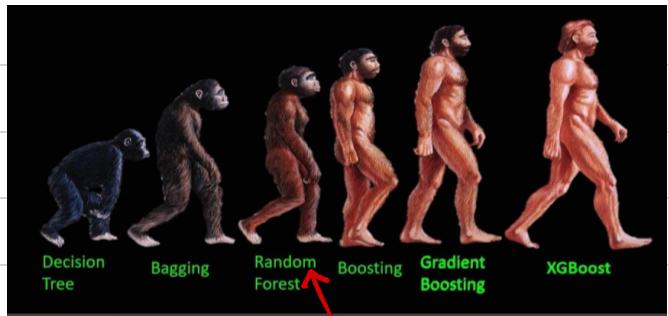
Session-5

Bagging and Random Forest (RF)



AHENDA

- ① ENSEMBLE
- ② BAGGING.



- ③ RF - Random Forest

where we will
reach today

- ④ OOB - Score

- ⑤ Hyper-Parameter Tuning - GRID Search & Random Search

DT - REGRESSION

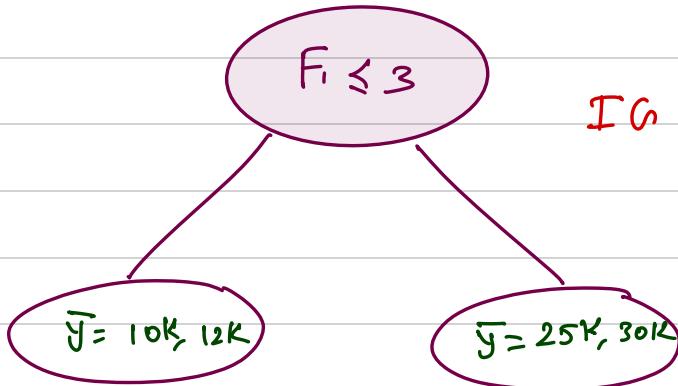
$$G \cdot I = 1 - \sum_{i=1}^k p(y_i)^2$$


$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

f_i	y
3	10K
2	12K
20	25K
23	30K
	$\bar{y} = 19.25$

$$MSE = \frac{1}{4} \left((10 - 19.25)^2 + (12 - 19.25)^2 + (25 - 19.25)^2 + (30 - 19.25)^2 \right)$$

$$= 1/4 * ((10 - 19.25)^2 + (12 - 19.25)^2 + (25 - 19.25)^2 + (30 - 19.25)^2) = 71.68$$



$$= 1/2 * ((10 - 11)^2 + (12 - 11)^2) = 1$$

$$= 1/2 * ((25 - 27.5)^2 + (30 - 27.5)^2) = 6.25$$

$$\frac{2}{4} \times 1 + \frac{2}{4} \times 6.25 = 3.37$$

ENSEMBLE

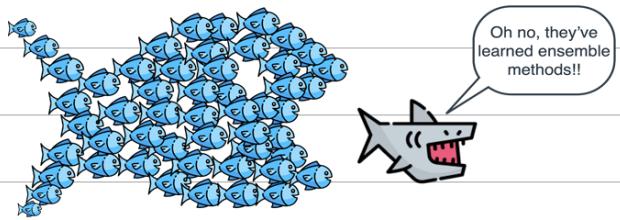
- ① Log-Reg → 1
 - ② Perceptron → 0
 - ③ Decision Trees → 1
- Voting.
1

Models in Ensemble algorithms:



→ TYPES OF ENSEMBLE LEARNING.

① Bagging - R.F



② Boosting (GBDT)

xgboost, lightgbm

③ Stacking

}

Never in companies, but used
in competitions

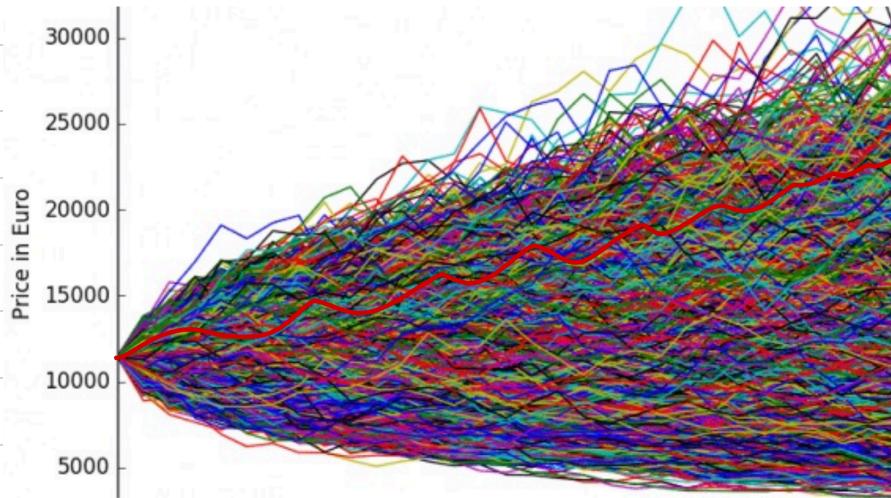
④ Cascading

Which one of these is a type of ensemble learning technique?

0 users have participated

- A Selecting 0%
- B Collecting 0%
- C Bagging 0%
- D Randomizing 0%

[End Quiz Now](#)



WORLDWIDE INDUSTRIAL MARKET REPORT



SM
2



K
1



Deepak
3

Soumyajit Misra
1/1 ⚡ 95.67

Shoreya gupta
4

1/1 ⚡ 95.13

Perisetla Pavan Kalyan
5

1/1 ⚡ 94.53

Gamidi Sri Valli Suprava
6

1/1 ⚡ 94.06

Vishwajit Verma
7

1/1 ⚡ 94.00

OM PRAKASH S
8

1/1 ⚡ 93.73

SHASHANK JHA
9

1/1 ⚡ 93.59

Sumanth Andhavarapu
10

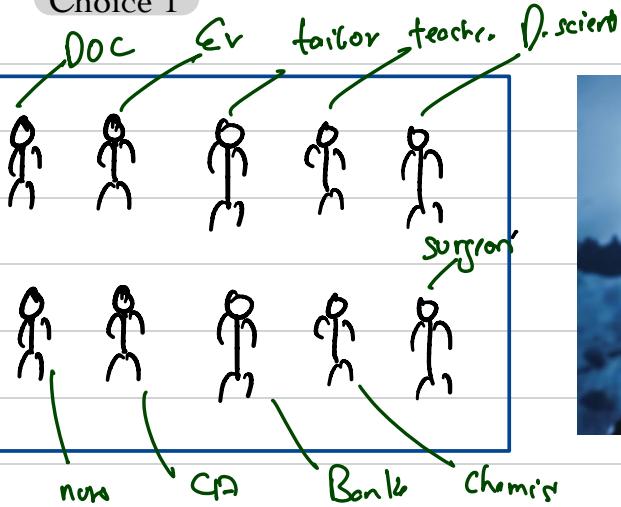
1/1 ⚡ 92.87

Population density

BAGINING

Boot-strapped Aggregation.

Choice 1



Choice 2



Eva 🌸

what movie villain do u secretly agree with?



kira 💜

Thanos, half of y'all need to go

I agree with Thanos everyday a little more



In bagging, each base-learner typically has high variance, also they're called "strong learners".

Dataset
n rows
m Feature

DT₁

Sample 1k rows out of n rows, with replacement.

DT₂

Sample 1k rows out of n rows, with replacement.

DT₃

Sample 1k rows out of n rows, with replacement.

DT₄

Sample 1k rows out of n rows, with replacement.

DT₅

Sample 1k rows out of n rows, with replacement.

DT₆

Sample 1k rows out of n rows, with replacement.

DT₇

Sample 1k rows out of n rows, with replacement.

	f_1	f_2	f_3	f_4	f_5	f_6	y
x_1							
x_2							
x_3							
x_4							
x_5							
x_6							
x_7							
x_8							

$x_1, x_1, x_3, x_5, x_7, x_5$

x_1, x_2

| million row

10¹² rows

RANDOM FOREST

1 mill rows
1000 Feature

DT \rightarrow 10K rows
 \downarrow
20 Feature
 $10K \times 20$



DECISION
TREE



RANDOM
FOREST

Dataset
n rows
m Feature

DT₁

Sample 1k rows out of n rows, with replacement.
Sample x features, out of m features, without replacement.

DT₂

Sample 1k rows out of n rows, with replacement.
Sample x features, out of m features, without replacement.

DT₃

Sample 1k rows out of n rows, with replacement.
Sample x features, out of m features, without replacement.

DT₄

Sample 1k rows out of n rows, with replacement.
Sample x features, out of m features, without replacement.

DT₅

Sample 1k rows out of n rows, with replacement.
Sample x features, out of m features, without replacement.

DT₆

Sample 1k rows out of n rows, with replacement.
Sample x features, out of m features, without replacement.

DT₇

Sample 1k rows out of n rows, with replacement.
Sample x features, out of m features, without replacement.

10K rows

8K - training
1K - testing
1K - validation

Rich man's evaluation



EXAMPLE

Original Set	
Patient A	
Patient B	
Patient C	
Patient D	

Bag 1

Bootstrap Sample	Out-of-Bag-Set
Patient A	
Patient A	
Patient C	
Patient C	

Original Set	
Patient A	
Patient B	
Patient C	
Patient D	

Bag 2

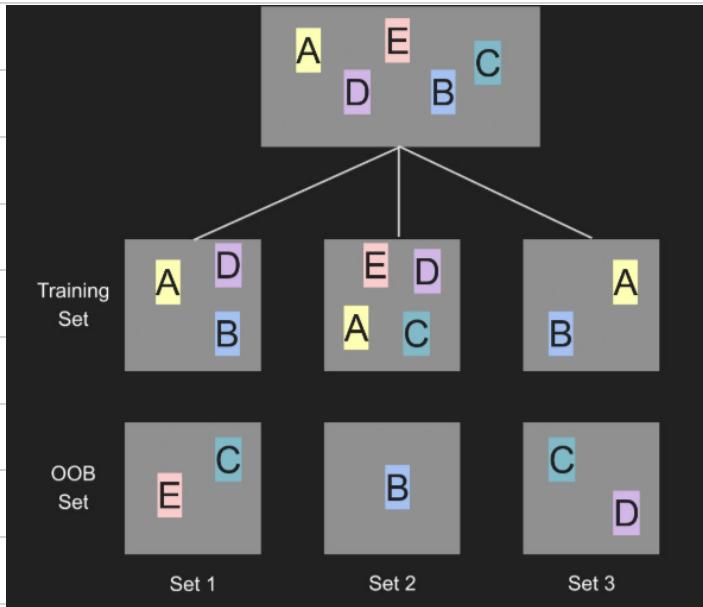
Bootstrap Sample	Out-of-Bag-Set
Patient A	
Patient B	
Patient C	
Patient D	

Original Set	
Patient A	
Patient B	
Patient C	
Patient D	

Bag 3

Bootstrap Sample	Out-of-Bag-Set
Patient A	
Patient D	
Patient D	
Patient D	

OOB-Score



$f_1 \ f_2 \ f_3 \ f_4 \ f_5 \ f_6 \ y$

x_1							
x_2							
x_3							
x_4							
x_5							
x_6							
x_7							
x_8							

$DT_1 \rightarrow x_1, x_2, x_3$
 $(x_4 - x_8) t_{12}$

$DT_2 \rightarrow x_1, x_6, x_7$
 $t_{12} (x_9 - x_5), x_8$

Bias & Variance TRADEOFF

What is the reason for introducing Row/Column Sampling in Random Forests?

0 users have participated

- A For decreasing the training time of model.
- B For tackling the problem of overfitting in base learners.
- C For tackling the problem of underfitting in base learners.
- D None of the above.

[End Quiz Now](#)



Karthik
2/2

184.98



Vishwajit Verma
2/2

186.66



Soumyajit Misra
2/2

179.44

4	SHASHANK JHA	2/2	175.95
5	Kiran Hebasur	2/2	175.71
6	Sumanth Andhavarapu	2/2	172.64
7	Praveen	2/2	171.99
8	Sri Harsha Nanduri	2/2	171.26
9	Deependu Ghosh	2/2	168.35
10	Arpita Saha	2/2	166.23

If a dataset contains "n" rows, and "m" of these rows are sampled to train the base learners in Random Forest, what will be the cross-validation data for each of the models?

1 user has participated

- A Complete dataset with "n" rows
- B A part of "m" sampled rows
- C Remaining "n-m" rows after sampling
- D None of the above.

[End Quiz Now](#)

Leaderboard

Based on all quizzes from the session



Vishwajit Verma
3/3

277.29



Karthik
3/3

279.91



Soumyajit Misra
3/3

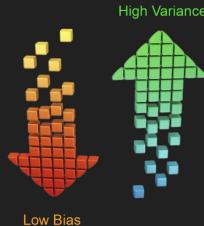
257.89

4	SHASHANK JHA	3/3	256.88
5	Praveen	3/3	255.68
6	Sri Harsha Nanduri	3/3	254.33
7	Deependu Ghosh	3/3	252.87
8	Kiran Hebasur	3/3	250.86
9	OM PRAKASH S	3/3	250.39
10	Arpita Saha	3/3	244.45

Bias - Variance Tradeoff

The base learners in RF are Deep Decision Trees.

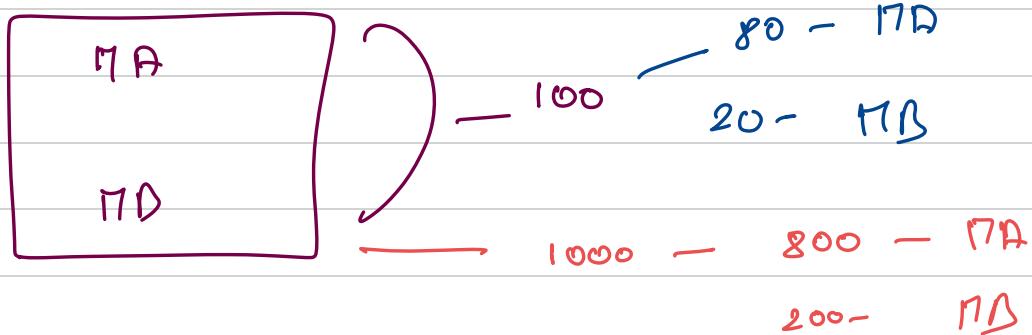
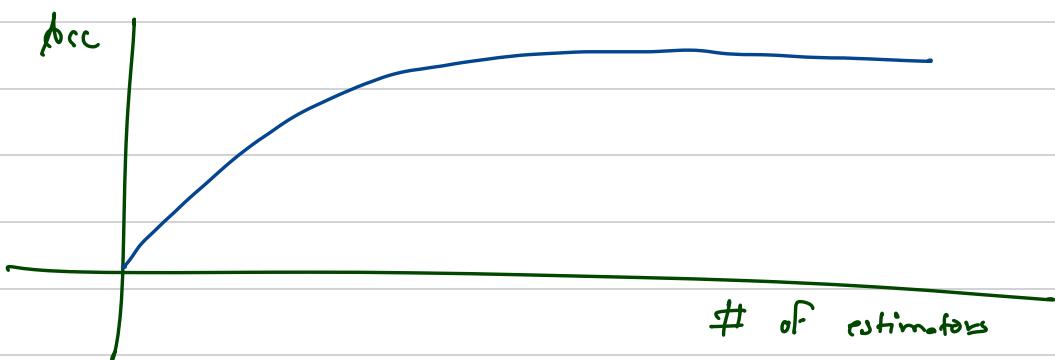
- So they slightly overfit on the sub sample of data
- Which means :



IMP. HYPERPARAMETERS - RF

of trees $\downarrow \rightarrow$ High variance
↓ DT
Overfit on sample data

of trees 11



ccp-alpha

Ridge
/

Log-Reg.

Loss function + $\lambda \times \|w\|^2$

DT →

DT + $ccp\text{-alpha} \times \# \text{ of Leaf nodes.}$

Loss function.

CCP-alpha tries to reduce the # of Leaf nodes

$\uparrow CCP\text{-alpha} \rightarrow \uparrow bias$

$\downarrow CCP\text{-alpha} \rightarrow \downarrow variance$

What is the use of the hyperparameter 'ccp_alpha'?

0 users have participated

- A To set the column sampling ratio 0%
- B To control underfitting or overfitting 0%
- C To optimize the learning rate of RF 0%
- D To set the depth of the base learners 0%

[End Quiz Now](#)

Leaderboard

Based on all quizzes from the session

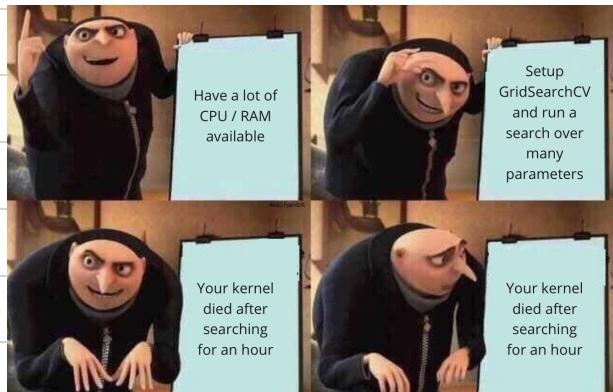
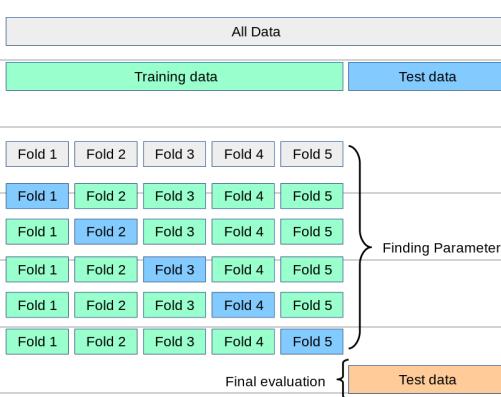
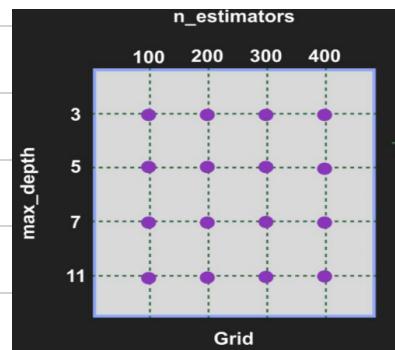
SJ	2	SHASHANK JHA	4/4	351.32	4/4	336.33
SM	1	Soumyajit Misra	4/4	352.42	4/4	330.37
Deependu Ghosh	3	Deependu Ghosh	4/4	345.72		
4	Sri Harsha Nanduri		4/4	336.33		
5	Kiran Hebasur		4/4	330.37		
6	Arpita Saha		4/4	328.72		
7	Praveen		4/4	328.65		
8	N Narayana		4/4	317.29		
9	N Nayana		4/4	314.71		
10	Karthik		3/4	279.91		

FINE-TUNING HYPER-PARAMETERS

```
# Define the hyperparameters grid for dt
param_grid_dt = {
    'criterion': ['gini', 'entropy'],
    'max_depth': [None, 5, 10, 15],
    'min_samples_split': [2, 10, 20],
    'min_samples_leaf': [1, 5, 10],
}
```

For crit in criterion:

```
for depth in max_depth:
    for samp1 in min_samples_split:
        for samp2 in min_samples_leaf:
            train_model (crit, depth, samp1, samp2)
            store-acc
```



If you're performing GridSearchCV on a random forest model for the parameters,

- 'max_depth' having 3 values and
- 'min_samples' having 4 values
- for 10 cross-validations

How many times model.fit() is called?



[End Quiz Now](#)



SHASHANK JHA
5/5 ₹ 432.52



Soumyajit Misra
5/5 ₹ 441.97



Kiran Hebasur
5/5 ₹ 427.93

4	Deependu Ghosh	5/5 ₹ 425.00
5	Sri Harsha Nanduri	5/5 ₹ 416.91
6	Praveen	5/5 ₹ 407.32
7	AS Arpita Saha	5/5 ₹ 403.38
8	N Narayana	5/5 ₹ 396.84
9	V Vishwajit Verma	4/5 ₹ 362.22
10	Shoreya gupta	4/5 ₹ 355.70