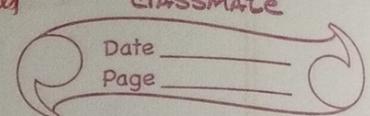


meaning:- taking knowledge from the crowd
— same concept used in Random Forest



Ensemble learning - Wisdom of the crowd

In Ensemble learning, the wisdom of the crowd is the idea that a group of diverse models, when combined, can make accurate predictions than any single model alone — similar to how a group of people can often make a better decision than one person.

In ML:

- Each decision tree in a Random Forest is like one "person" in the crowd.
- By combining their predictions, the final result is more reliable.
- For Classification → Take the majority vote
- For Regression → Take the average prediction

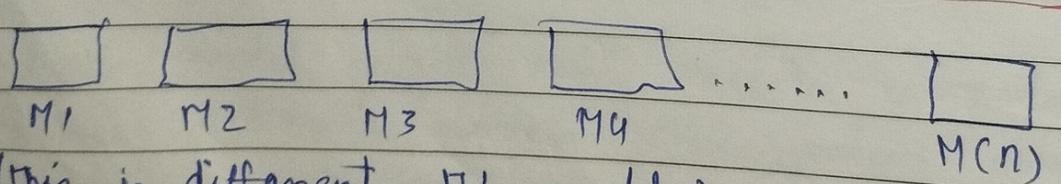
Ensemble Learning - Core Idea

Machine learning algorithm: Training
↓ Testing (Prediction)

What happens in prediction stage:

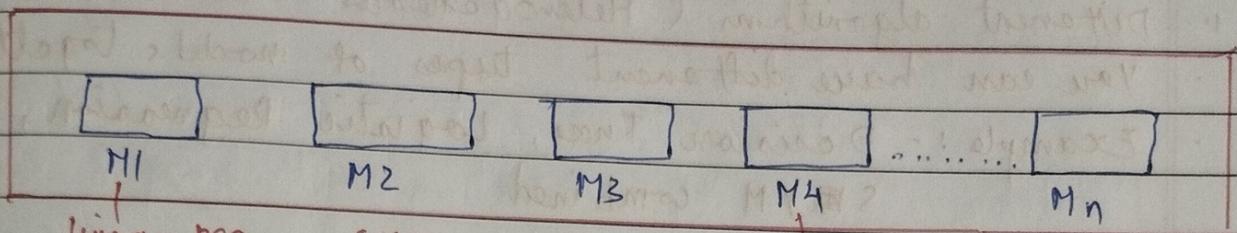
For eg we have Ensemble learning Algorithm :-
collection of smaller basemodels

this is called Ensemble learning model



→ (This is different ML models)

This models are different
& we always need different model in Ensemble learning



→ These are different models

(The model in Ensemble learning needs to be different or at least behave differently because if they are all the same and make the same mistakes, combining them won't help — you'd just be repeating the same error)

Why diversity is important:

- 1. Error cancellation
- 2. Better coverage of patterns
- 3. Reduced Overfitting

→ ① If one model makes wrong prediction, another model might make the correct one
When combined, the mistakes cancel out, improving accuracy

→ ② Different models capture different aspects of the data
Ex: A decision tree might capture non-linear patterns while a logistic regression might capture linear patterns

→ ③ If all models are trained the same way, they might all overfit to the training data
Different models via randomization make them general better

In Ensemble model 3 core Ideas

- 1. Different algorithms (Heterogeneous ensemble)
 - You can have different types of models together
 - Example: - Decision Tree, Logistic Regression, SVM combined
 - Common in stacking

- 2. Same Algorithm, trained differently (Homogeneous Ensemble)
 - on different data
 - You can use the same type of model but train them on different subsets of data or with different settings
 - Example: - Many decision trees trained on different random subsets → Random Forest

we use the same algorithm but we train it on different data



model train differently

variety difference

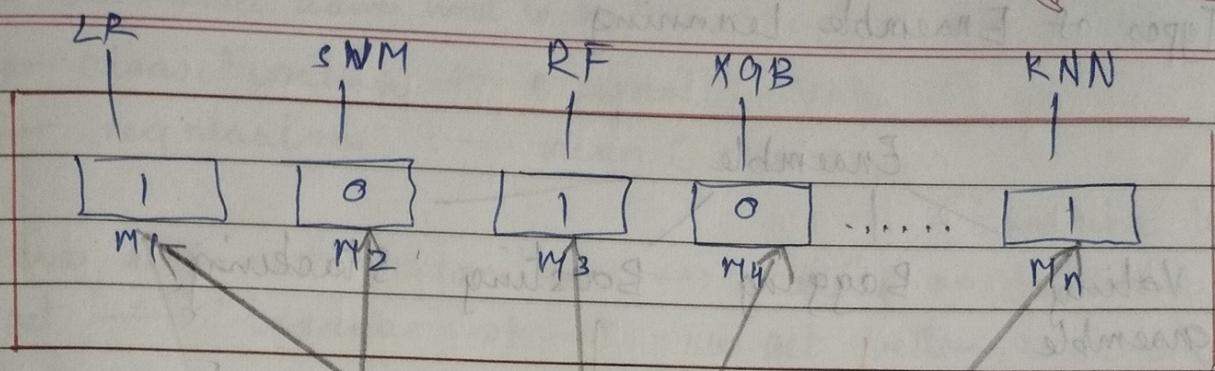
Types

- 1. Heterogeneous Ensemble :- Different Algorithm same data
- 2. Homogeneous Ensemble :- Same Algorithm but different subset of data
 - ↓ → this is mostly used in Bagging & Boosting
- 3. Hybrid Ensemble :- Hybrid Ensemble is a mixture of both where we train on different algorithms (Heterogeneous) and also train on different subsets of data (Homogeneous)

1 - Yes
0 - No

classmate

Date _____
Page _____



A Person with 85% 18 and 7.5 CTPA :- prediction will be placed or not? Classification Problem

We give this data to different algorithms
for eg we have 5

1 :- ~~Linear Reg~~ Logistic Regression (1)

2 :- SVM (0)

3 :- Random Forest (1)

4 :- XGBoost (0)

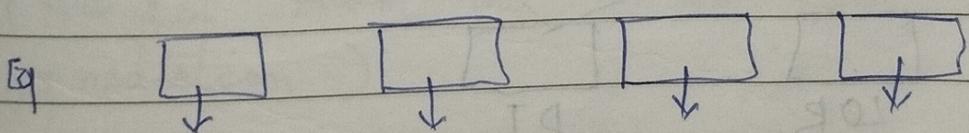
5 :- KNN (1)

out of 5 → 3 said Yes → and 2 said No

∴ final output will be :- Yes because in

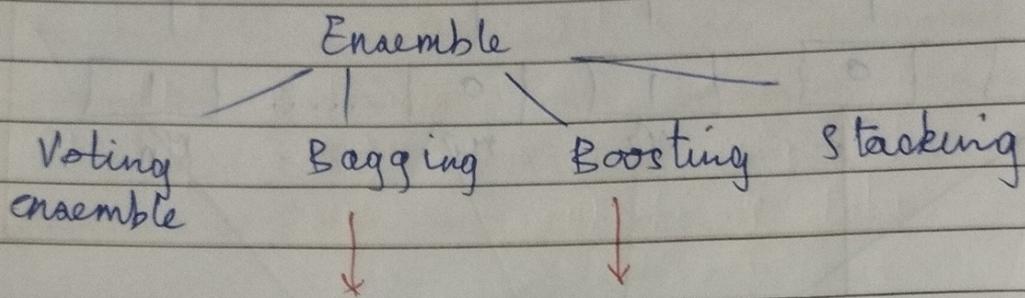
Ensemble learning — majority counts

In Regression :- How much LPA? we will get?



$$8.1 + 3.2 + 4.1 + 7.5 = 22.8 \text{ LPA}$$

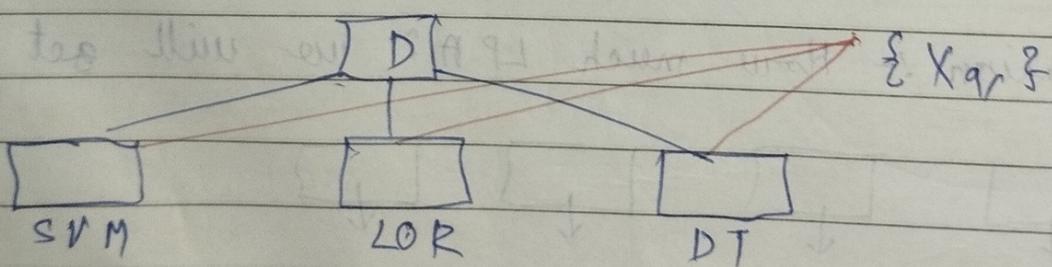
Types of Ensemble Learning



- 1) Random Forest
- 1) AdaBoost
- 2) Gradient Boosting Machines (GBM)
- 3) XGBoost
- 4) Light GBM

- XG - Extreme Gradient Boosting
 - Ada - Adaptive Boosting
 - Light GBM - Light Gradient Boosting Machines
- Why exist different types
 → because of 3 core Ensemble learning model ideas
- 1) Heterogeneous Ensemble
 - 2) Homogeneous Ensemble
 - 3) Hybrid Ensemble

- i) Voting Ensemble :- Whether our base models have (Heterogeneous)
 different Algorithms



Same Dataset D :- Used for different algorithms such as SVM, LOR and DT and now we train it on our data D

- In stacking :- Instead of directly voting, the predictions from the base models becomes the input features for the meta-model
- the meta-model learns how to best combine base models output.
- for classification :- majority vote
 for regression :- mean (average)

classmate

Date _____

Page _____

(Ensemble learning mode)

Since different algorithms → variety or difference we get → because of it we get better, accurate result

(Heterogeneous Ensemble Learning)

2) stacking



$\{X_g\}$

Base Model

MetaModel — (it is trained on validation set) or

KNN

out-of-fold predictions from the base models

(This model does one work:

— assigning weights to the above

algorithm based on the data who has

more weight or correct (accurate)

From where we get data
 (We get from the model itself which makes prediction)

For Eg	SVM	LOR	DT	Result
based on this	1	0	1	1
	0	0	1	1

result it assigns weightable or weight to model (like who's more correct or accurate or not)

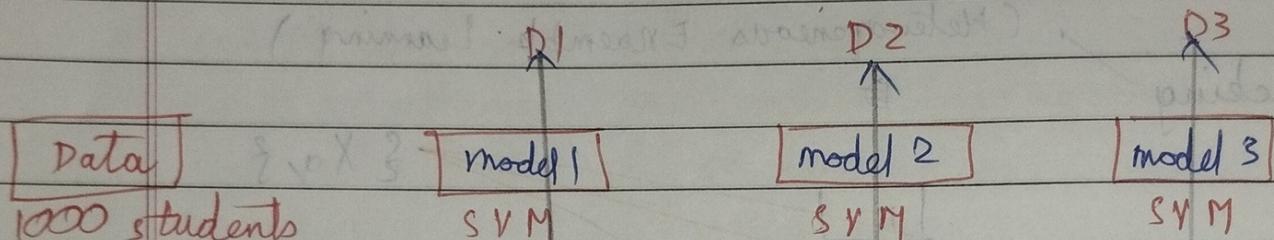
(This could be of any reason)

3) Bagging \rightarrow Bootstrap Aggregation

This follows

(Homogeneous) \rightarrow same algorithm and different subsets of data

$\{ X_i \}$



For eg we gonna feed 500 students data

Random forest



base mode

D trees \rightarrow

In final round

I took

random

500 students

data

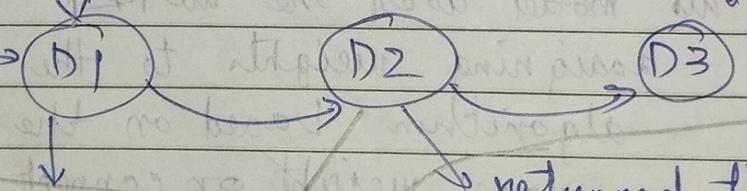
after this

I again returned the 500 students data

and again took 500

random students data

this is called Bootstrapping



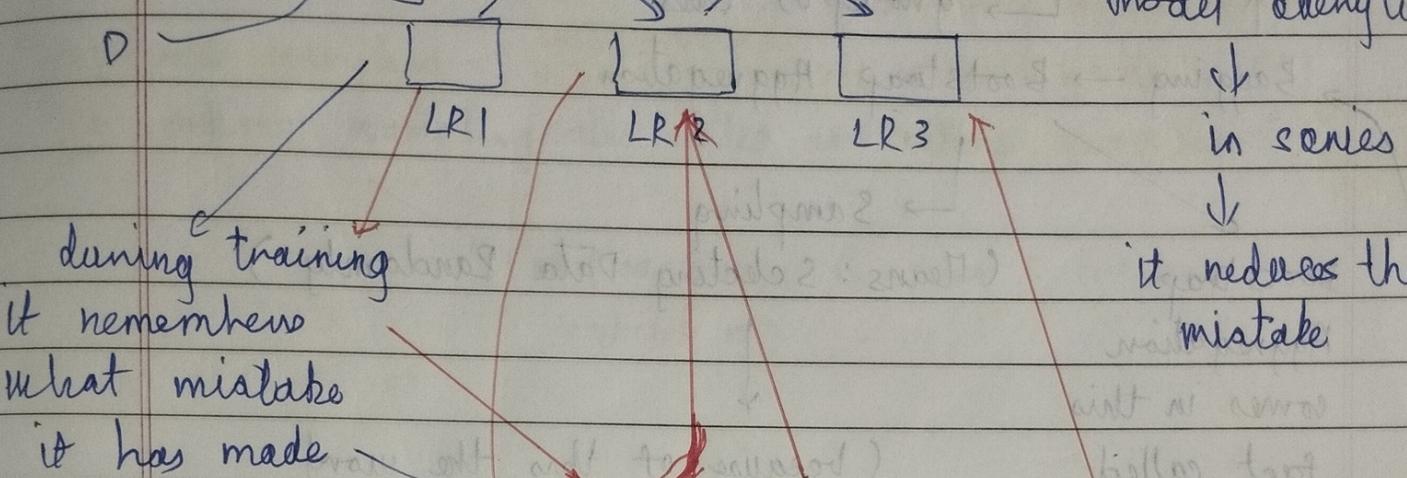
returned the data and again took random 500 students data

Probability of getting the same data is relatively lesser in Percentage as number

most powerful machine learning Algorithm

4) Boosting :-

Boosting your data :- tells the other model everytime



↓
Data
↓
1000 students

It tells this is the student's data I have made mistakes make sure you do correctly in this area

Random Forest

It tries to make sure not do the same mistake and then it does something happens if tells the other model to make sure it doesn't repeat the same mistake and do it correctly

Random Forest (Bootstraping)

Random Forest

↑ ↗ group of trees

→ Bagging → Bootstrap Aggregation

Bootstrap

Aggregation

comes in this

fact called

sampling

Sampling
(Means: Selecting Data Randomly)

↓
(because of this the word
Random came out)

Decision Tree

(any random values chosen randomly)

D = Dataset

D = 1000 rows → (we only give 500 rows out of 1000 rows)

(we give 500 rows out of 1000 rows)

to decision tree but equal to 1000 rows

[DT1] [DT2] [DT3] [DT4] ... [DT100]

not fully) → what we do is: - we create samples out of it:

* 3 Sampling techniques

1 - row sampling

here also 2

types of selection

→ with sampling

→ without

1. with replacement

2. without replacement

we can have duplicate values

because one data comes

and it again goes back to the Database which can repeat again & again :- D = [1, 2, 91, 1001, 2, 1, 91]

with replacement → duplicate values

without replacement → no duplicate values

Data Subset → creation Ideas

without replacement → no duplicate values

increases diversity

1. row sampling < with replacement without replacement - less diversity

2. column sampling also called as feature sampling < with replacement without replacement

eg (randomly choosing 5 columns
out of 10) - D1



for another

(decision tree we again
randomly choose 5 columns)

This is called as

Random subspace method

same we do

for other decision tree

or combining

3. hybrid sampling (row sampling + column sampling)

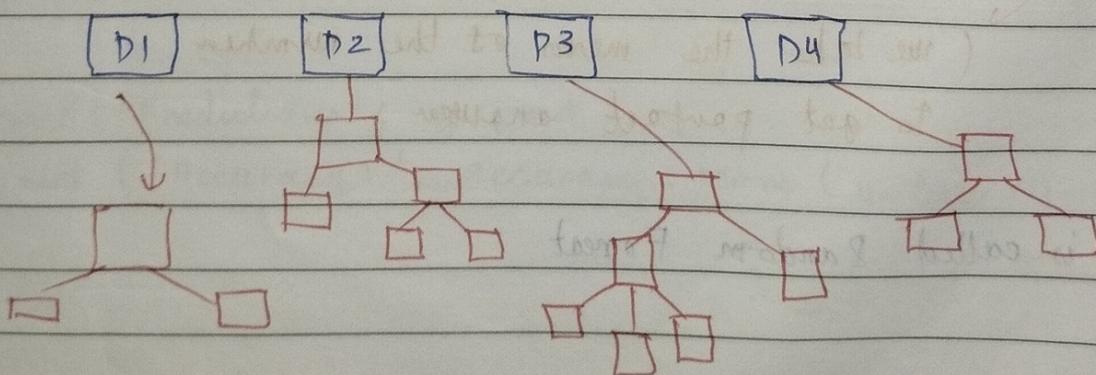
with replacement

(duplicate values)

without replacement

(no duplicate values)

* Since we give different subsets of data to each base model
their structure can also be different

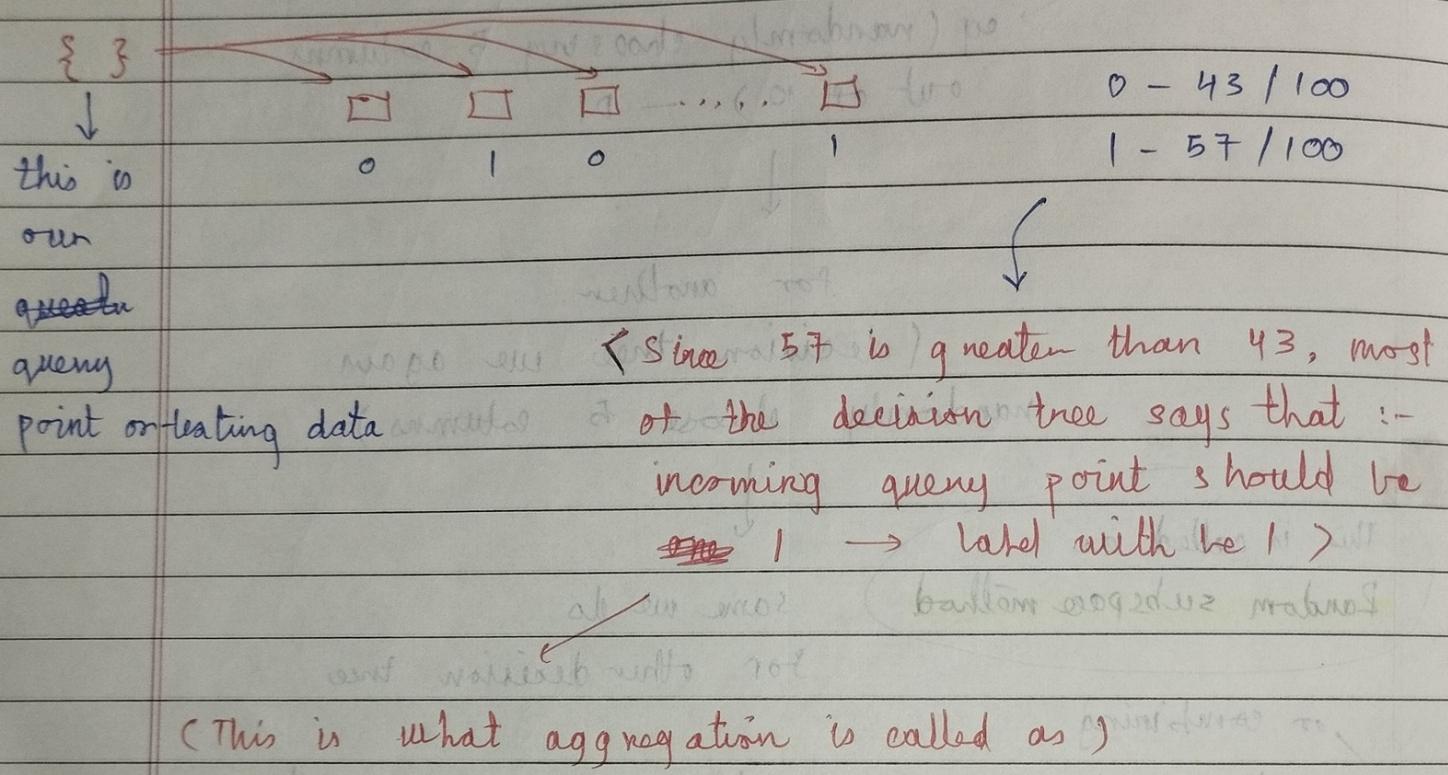


→ This whole thing is called as bootstrapping

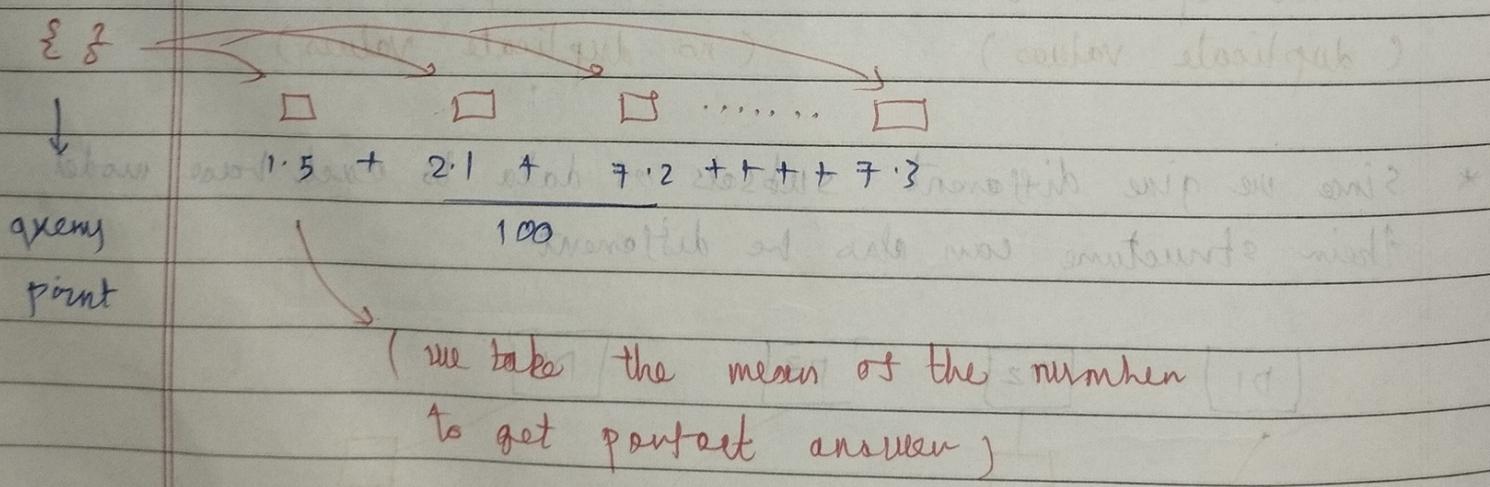
Random Forest - (Aggregation)

Aggregation :- When we want to predict
(classification problem)

For eg 1000 decision trees are trained \rightarrow now they can predict \rightarrow think this as a classification problem



→ If it is a regression point



→ This is called Random Forest