

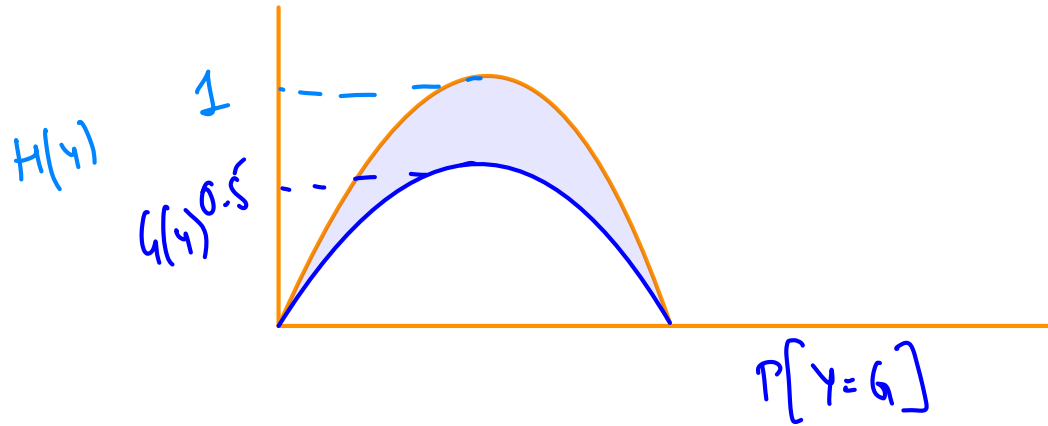
BAGGING
&
RANDOM FOREST

DT \rightarrow Impurity

\rightarrow Entropy $H(y)$

$[0, 1]$
 $P(y_i) \log_2 [P(y_i)]$

Gini Impurity $G(y)$
 $[0, 0.5]$
 $1 - \sum [P(y_i)]^2$



$y = [0, 1, 2]$

$1 - [P(y=0)^2 + P(y=1)^2 + P(y=2)^2]$

Nomrical feature

f_1 f_2 f_3 f_4
C C N N

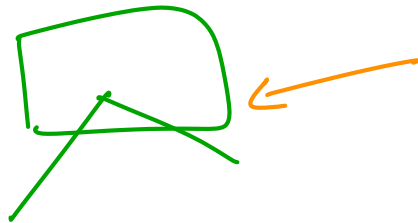
$n = 100$ (unique)

$d = 4$

f_1	f_2	f_3	f_4

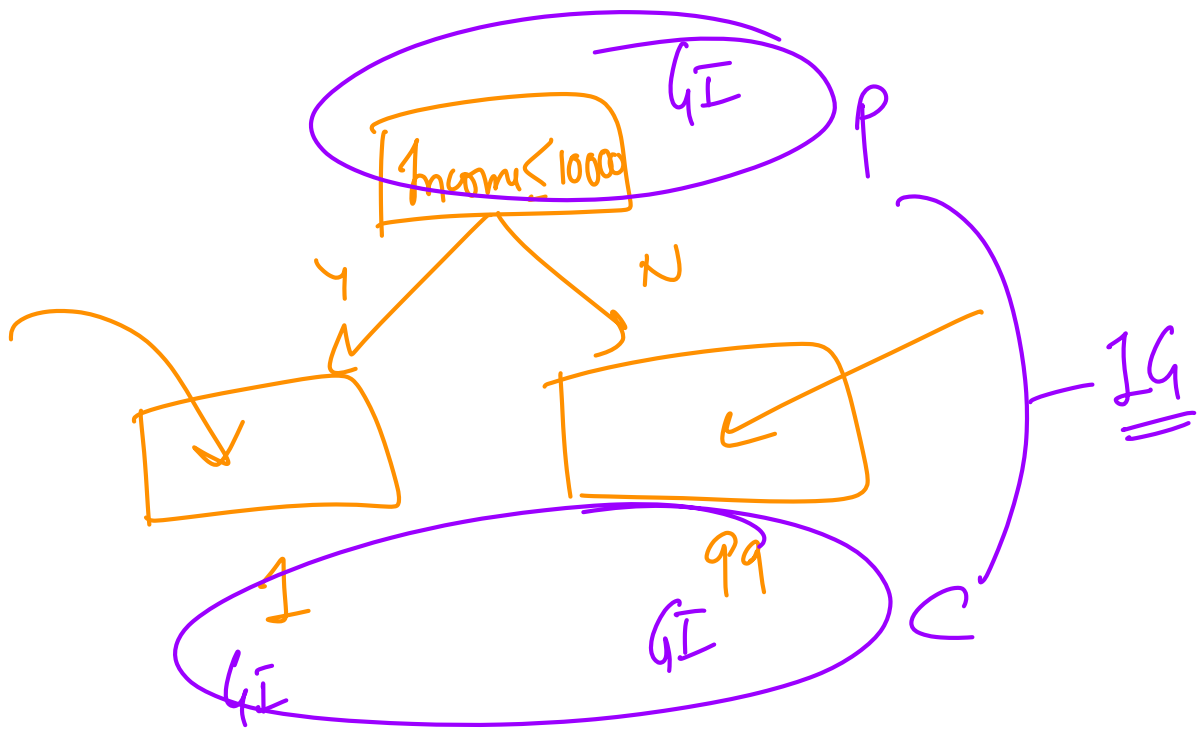
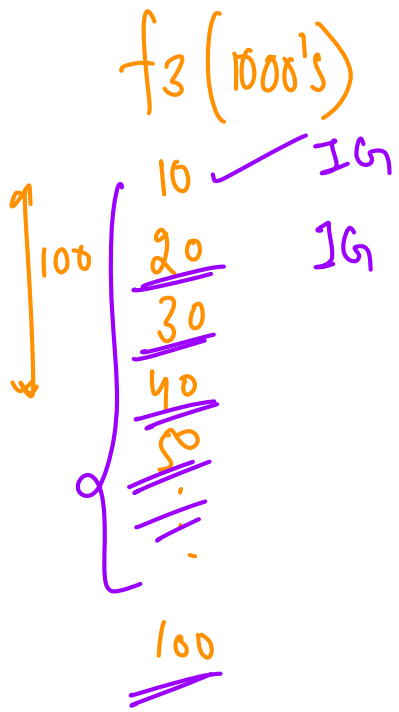
$\Rightarrow \underline{\underline{DT}}$

IG



$f_3 \rightarrow IG \rightarrow 100$
 $f_4 \rightarrow IG \rightarrow 100$
 $f_1 \rightarrow IG - 1$
 $f_2 \rightarrow IG - 1$

202



Q1

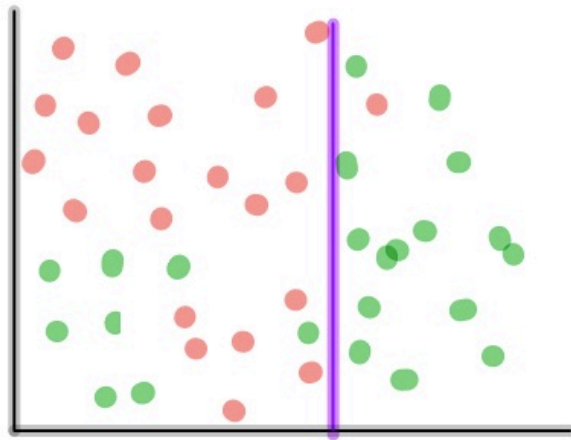
Overfit

Underfit

\Rightarrow Too Complex $\leftarrow d \uparrow$

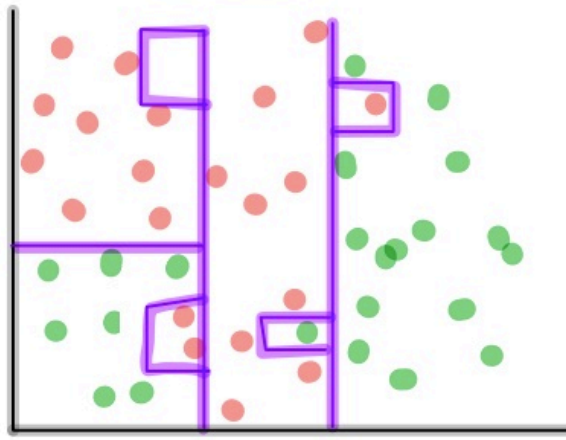
\rightarrow Too simple. $\leftarrow d \downarrow$

Underfitting



Shallow tree

Overfitting



Deep Tree.

Impact of Outlier

d ↑

Impact of outlier is significantly observed.

d ↓

Impact of outlier is either not observed or minimal.

Q

DO WE REQUIRE
feature Scaling

$$\text{Entropy} = - \sum_{i=1}^k P(y_i) \log P[y_i]$$

$$Gini = 1 - \sum_{i=1}^k [P(y_i)]^2$$

Count of Occurrence

calculating based on

frequency

lot of feature \rightarrow high dimensional data

$f \rightarrow 1000's$

$\rightarrow 400 C$
 $\rightarrow 600 N$

} Computational requirement
 \rightarrow

Q Is DT a good idea?

\hookrightarrow Slow

\hookrightarrow Lots of Computational resources required. } \rightarrow

Imbalanced data

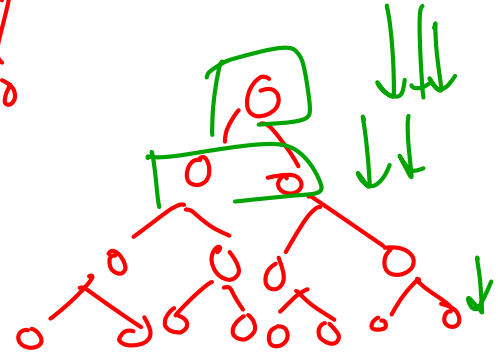
↳ Should I use DT ?

→ Rebalance.

feature importances

How do you find out feature importance?

Total IG → ✓



Regression

2T

MSE

Gini
Entropy

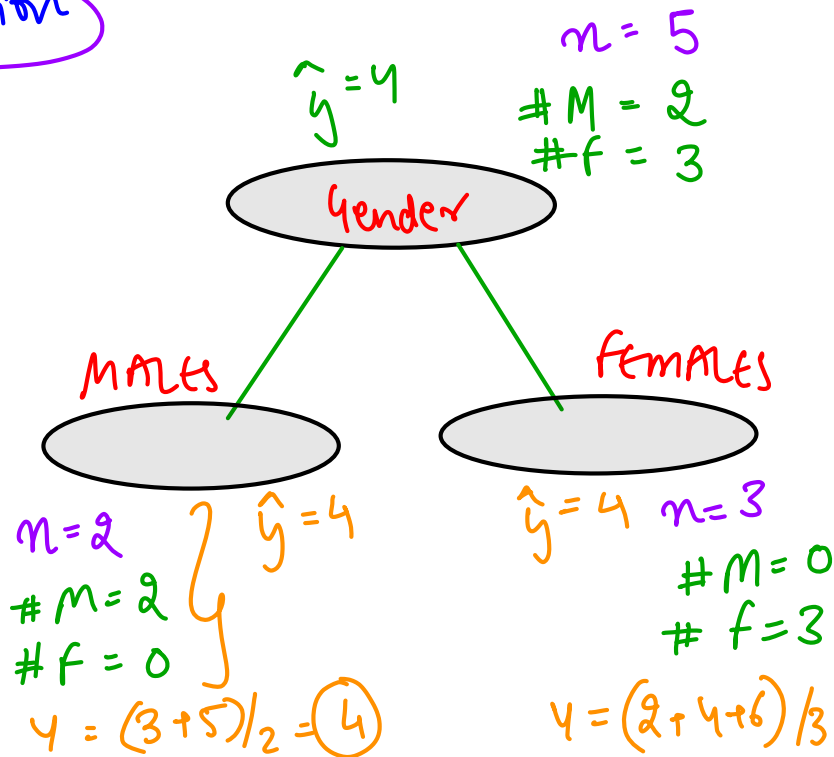
$$-\frac{1}{n} \sum (y - \hat{y})^2$$

Decision Trees for Regression

f_1	f_2	f_3	y	
F	Y		2	D
M	N		3	D-
F	Y		4	D
M	N		5	D-
F	N		6	D

$$\frac{1}{n} \sum (y - \hat{y})^2$$

←



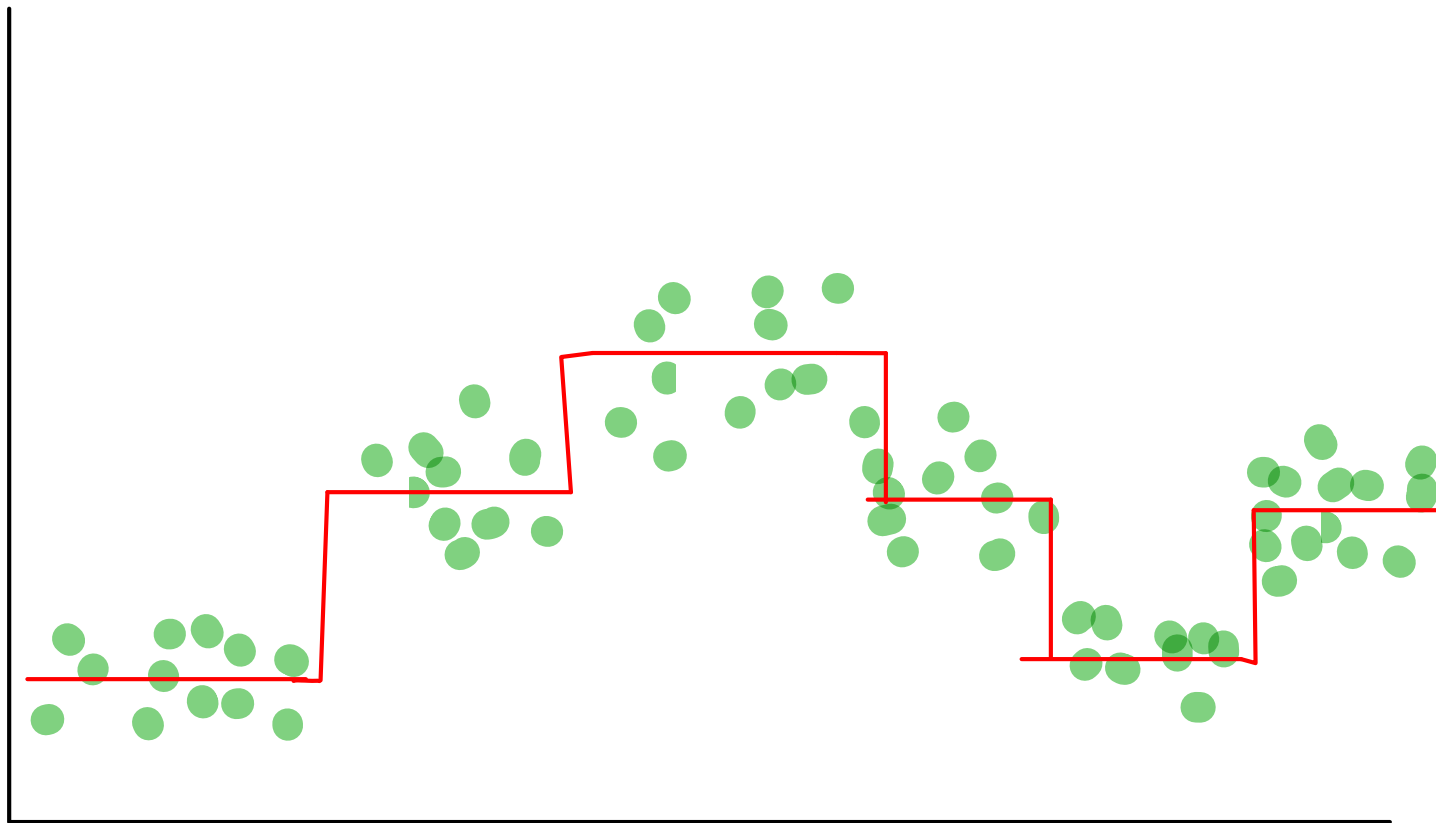
$$C_1 \text{ MSE}_{\text{Males}} = \frac{1}{2} \left[(3-4)^2 + (5-4)^2 \right] = \bigcirc$$

$$C_2 \text{ MSE}_{\text{Females}} = \frac{1}{3} \left[(2-4)^2 + (4-4)^2 + (6-4)^2 \right] = \bigcirc$$

$$\text{MSE}_{\text{Children}} \Rightarrow \frac{2}{5} \text{MSE}_m + \frac{3}{5} \text{MSE}_f$$

A purple arrow points from the first circle to the coefficient $\frac{2}{5}$. A green arrow points from the second circle to the coefficient $\frac{3}{5}$.

$$\text{IG} = \text{MSE}_{\text{parent}} - \text{MSE}_{\text{children (weighted)}}$$



Ensemble

↳ Group of things // → Models

↳ Multiple model → Used together ⇒ Answer.

iphone

Father ⇒ NO

Mother ⇒ YES

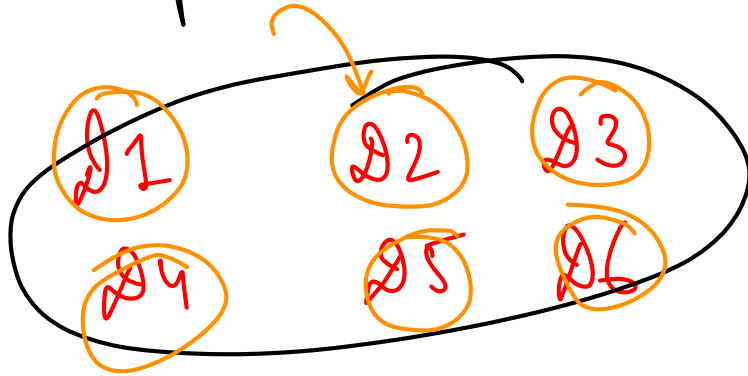
Sister ⇒ YES

F1 ⇒ Y

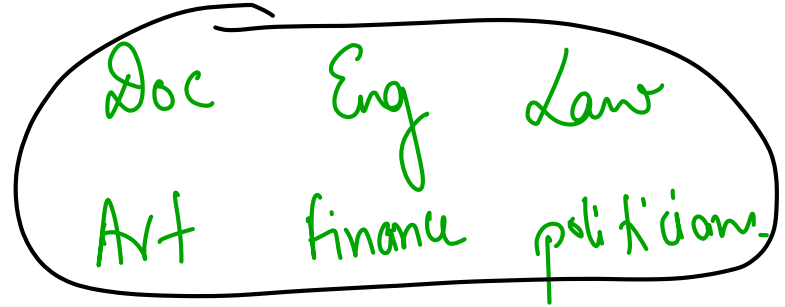
F2 ⇒ Y

F3 ⇒ Y

Implement a new rule related to medicine



I



II

$[m_1 \ m_2 \ m_3 \ \dots \ m_k] \Rightarrow$ Base learner.

Combine \Rightarrow final answer.

① BAGGING

② BOOSTING

③ CASCADING

④ STACKING

Quiz time!

🕒 Quiz Ended!

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

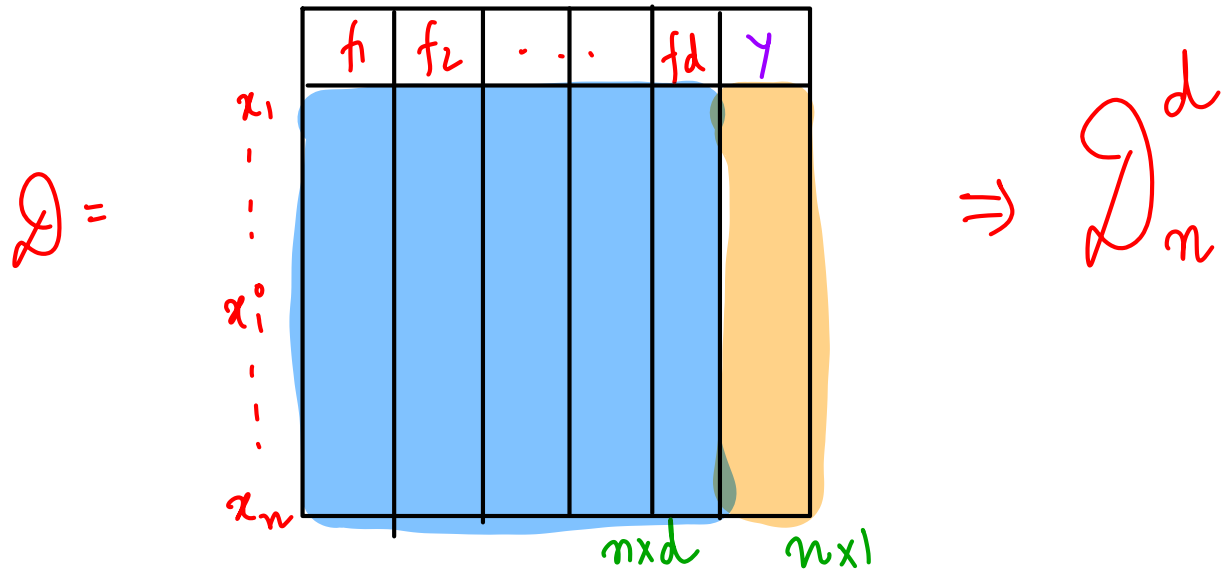
29 users have participated

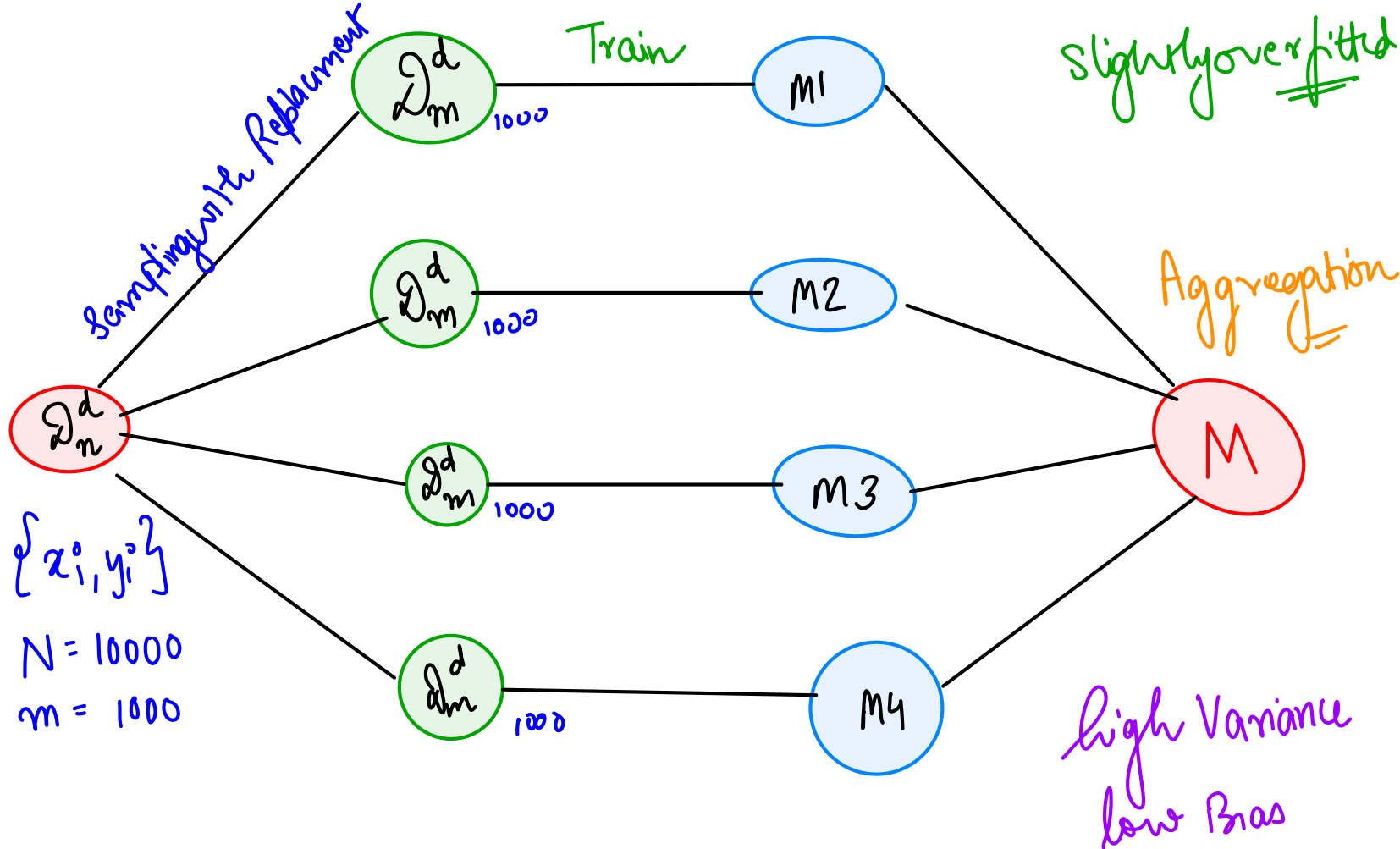


BAGGING

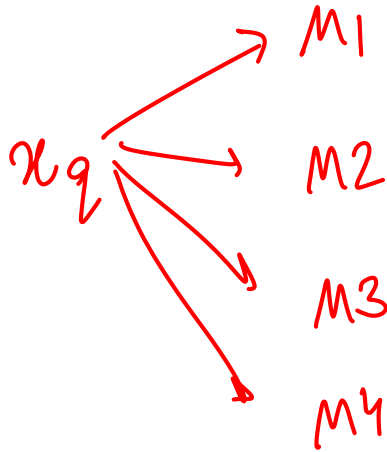
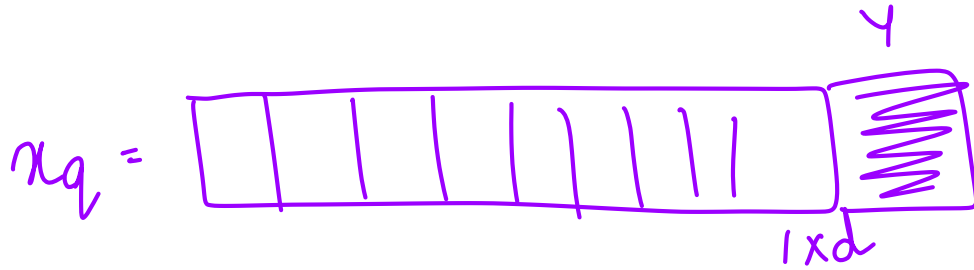
= BOOTSTRAPPED AGGREGATION

Bootstrapped Sampling → Sampling with replacement
Aggregation → Mean, Max, Min.





At Inference



\hat{y}_1

\hat{y}_2

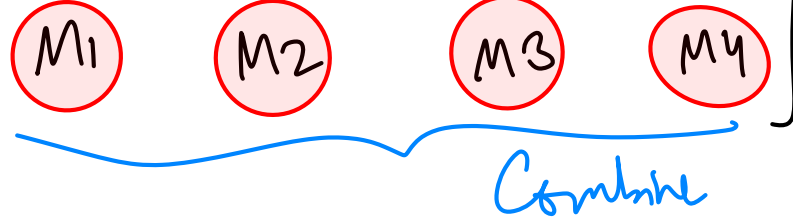
\hat{y}_3

\hat{y}_n

Classification \Rightarrow Majority Voting

Regression \Rightarrow Mean / Median

BAGGING



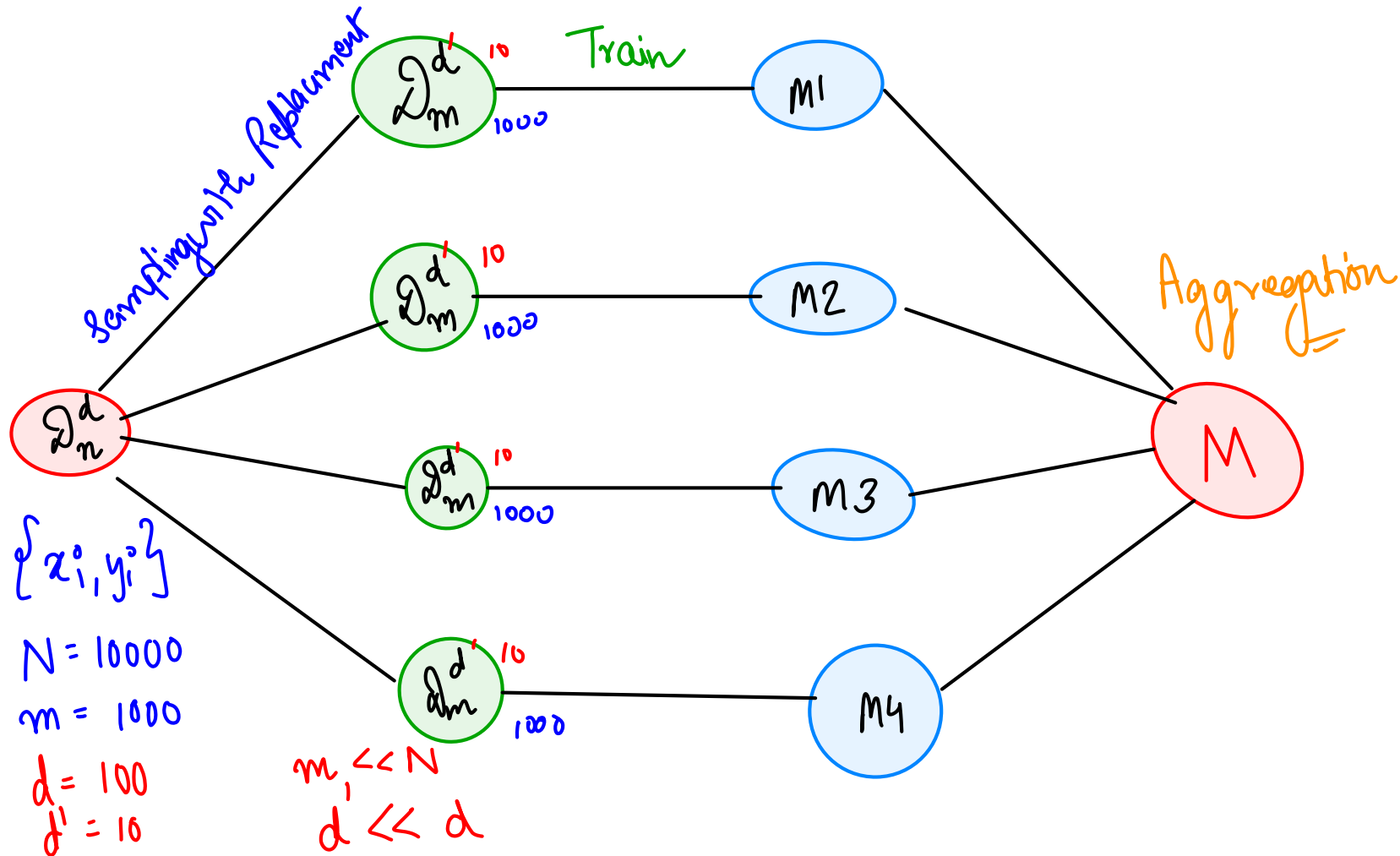
RANDOM FOREST \Rightarrow Collection of trees \Rightarrow decision trees.

① Randomly Select Rows

[Row Sampling]

② Randomly Select Columns

[Column Sampling]
↳



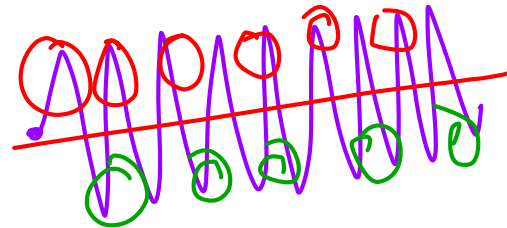
$$RF = DT + \overset{\downarrow \downarrow}{R.S} + \overset{\downarrow \downarrow}{C.S} + \text{Aggregation}$$

Base learner

$$m \ll N$$

$$d' \ll d$$

BIAS \downarrow
VARIANCE \uparrow

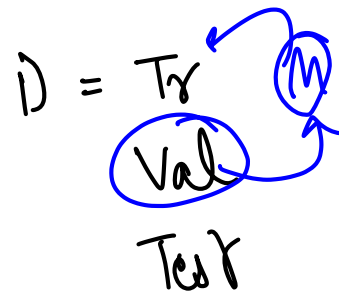


$$'M' \Rightarrow \underline{\text{Low BIAS} + \text{Low VARIANCE}}$$



Aggregation + Randomisation

How do we validate RF?



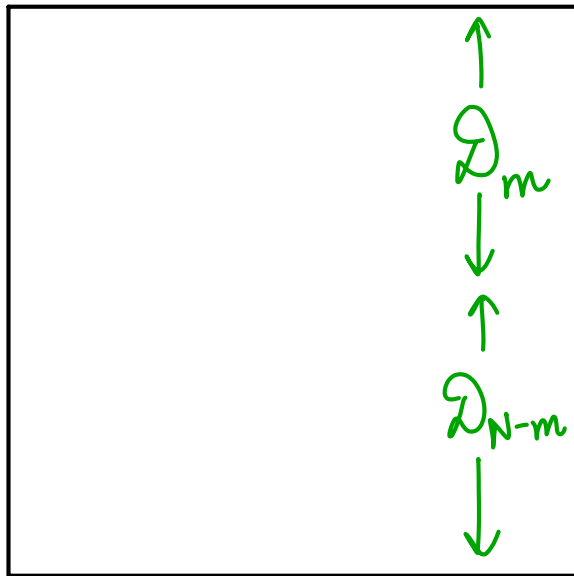
$\Rightarrow k \Rightarrow$ different model

$\Rightarrow 'M'$

$\leftarrow f_d \rightarrow$

D_N^d

\uparrow
 x_N
 \downarrow



D_m

D_{N-m}

Training

Validation

OOB POINTS

Out of Bag

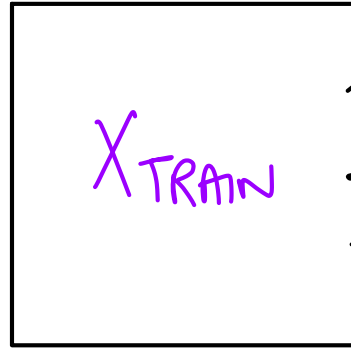
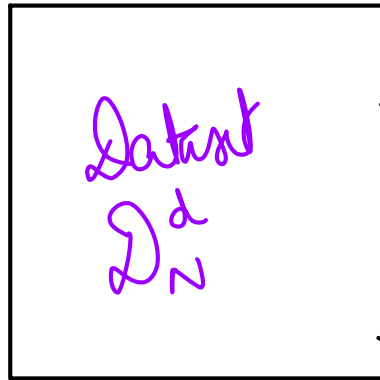
~~Hyper~~

$d =$
⋮

$i =$
⋮

$m =$
⋮

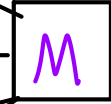
$d' =$
⋮



⋮

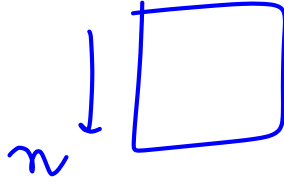


⋮



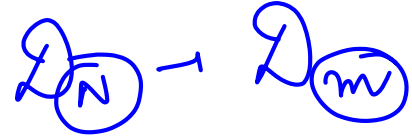
Use this for hyper parameter tuning

Used only once. Used to test the model's performance



Quiz time!

🕒 Quiz Ended!



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?

21 users have participated

- | | | |
|-----|-------------------------------------|-----|
| A | Complete dataset with "n" rows | 5% |
| B | A part of "m" sampled rows | 19% |
| ✓ C | Remaining "n-m" rows after sampling | 76% |
| D | None of the above. | 0% |



