# Ensemble

→ group of musician

↓

In m/c learning

↓

group of models

-

Ensemble

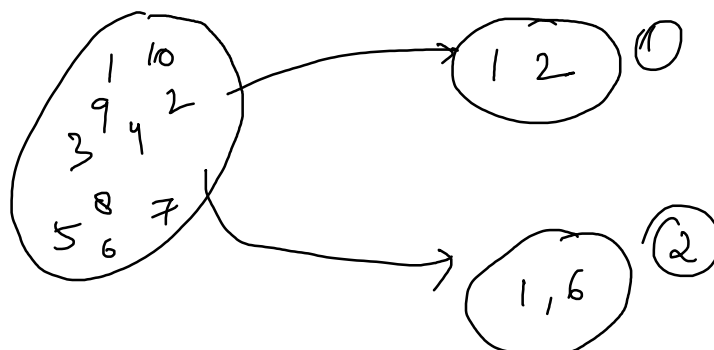↙           ↘

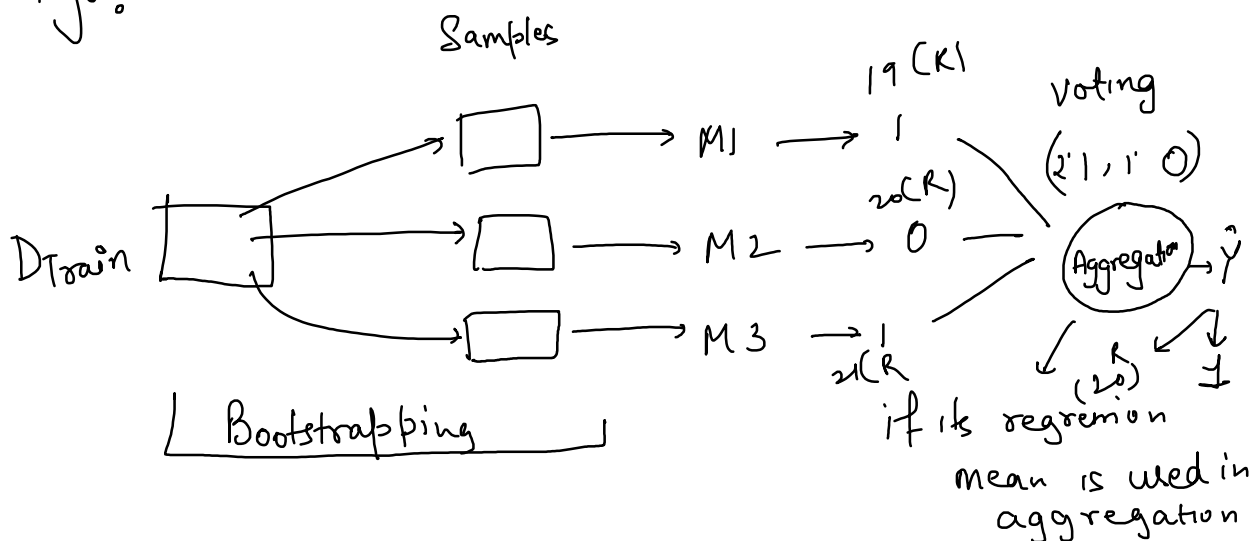Bagging          Boosting

# BAGGING

↙                    ↘

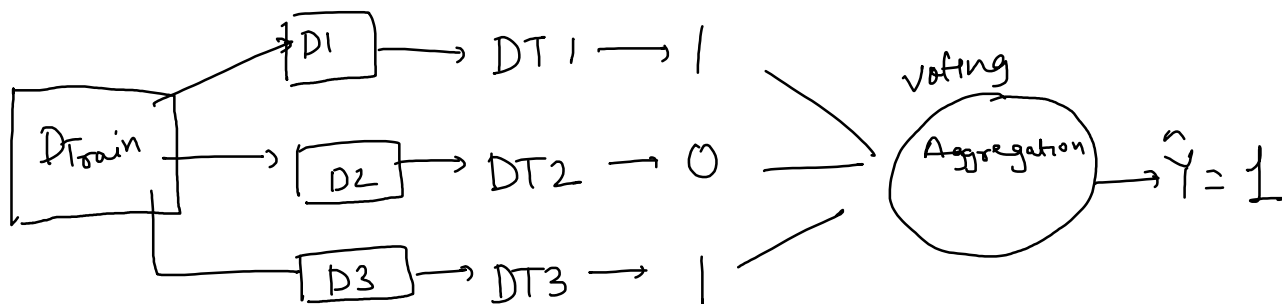Bootstrapping                    Aggregation

↳ (Sampling with replacement)

# Bagging Algo:

Samples

Dtrain → [ ] → M1 → 19 (K) 1

→ [ ] → M2 → 20(R) 0

→ [ ] → M3 → 1, 21(R)

Bootstrapping

voting
(2'1, 1' 0)

Aggregator → $\hat{y}$

(20)

if its regression
mean is used in aggregation

---

## Random Forest:
→ Collection of Decision Trees

→ DT should be of good depth (overfitting)

Dtrain → D1 → DT1 → 1

→ D2 → DT2 → 0

→ D3 → DT3 → 1

voting
Aggregation → $\hat{Y} = 1$

# Models should be different from each other

| | CGPA | IQ | EXTRA | SOCIAL | PLACED |
|---|---|---|---|---|---|
| ① | 7 | 110 | 10 | 9 | 1 |
| ⑪ | 8 | 112 | 9 | 8 | 0 |

| | | | | | |
|---|---|---|---|---|---|
| (II) | 8 | 112 | 9 | 8 | 0 |
| (III) | 9 | 118 | 8 | 7 | 0 |
| (IV) | 10 | 125 | 7 | 6 | 1 |

lets choose I & III : (50% of columns)

$D_1 \rightarrow$ DT1          $D_2 \rightarrow$ DT2

| CGPA | Extra | PLaced | | IQ | SOCIAL | PLACED |
|---|---|---|---|---|---|---|
| 7 | 10 | 1 | | 110 | 9 | 1 |
| 9 | 8 | 0 | | 118 | 7 | 0 |

Different models $\longrightarrow$ different samples $\rightarrow$ column sampling & row sampling

RF $\Rightarrow$ lowbias & high variance DTs + Row sampling + column sampling
$\underbrace{\hspace{4cm}}_{Overfitting}$

Output $\leftarrow$ Aggregation

$d$ (no. of dimensions) 

$d'$ (#dimensions)

$D_{Train}$

$n$ (# rows)

$d > d'$

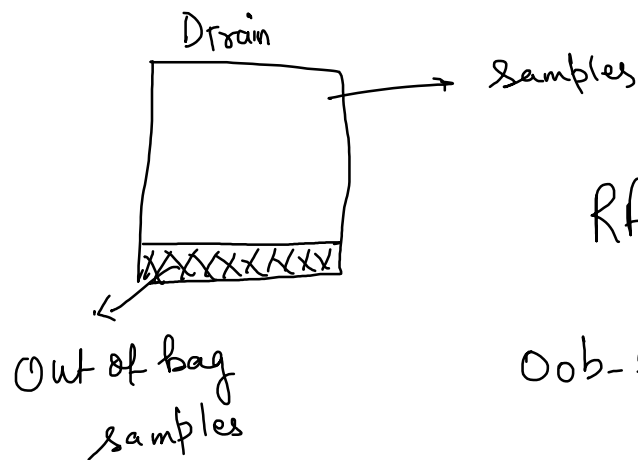$m > n$

m( no. of )

$m(\text{no. of rows})$

$D_{Train}$

$m > n$

## Hyper parameters:

1) No. of models $\uparrow$ (n_estimators) $=$ overfitting $\uparrow$ [100-2000]

2) max_features $\rightarrow$ [ "auto", "sqrt", "lg", 0.7, 0.3]

3) row-sampling $\Rightarrow \dfrac{n}{m}$

4) n_jobs $\Rightarrow -1 \Rightarrow$ Capability of your CPU

5) max_depth

### OOB score
$\downarrow$
out of bag

Dtrain

$\rightarrow$ Samples

Out of bag samples

RF(OOB_score = True)

Oob-score $\propto$ accuracy

Advantages:

Disadvantages:

Advantages:                          Disadvantages:

→ Feature Importances              → Black box


Hyperparameters
Tuning
GSCV ←            → RSCV


Boosting


① Bagging → overfitting DTs

② Boosting → underfitting DTs → ( high bias & low variance)


Flow Chart:

$D_{Train}$ ⟶ $[x_i, y_i]$ ⟶ $M_0$ ⟶ $[\hat{y}_i]$ predictions → $[e_i]$ errors ⟶ $[x_i, e_i]$ $M_1$ ⟶ pred → errors
$(h_0(x))$                    $[h_1(x)]$                          ↓
                                                            repeat
                                                            till errors are
                                                            min

0) $D_{Train} = \left[ x_i^o, y_i^o \right]_{i=1}^n \Rightarrow M_0 \rightarrow$ predictions $\overset{(\hat{y}_i)}{\longrightarrow}$ errors

$[ h_0(x) ]$

$e_1 = y_i^o - \hat{y}_i^o$

$e_i^o = y_i^o - h_0(x)$

2) $M_1 \underset{h_1(x)}{\longrightarrow} \left[ x_1, e_i \right]_{i=1}^n \quad e_i^o = y_i - h_0(x)$

Model at end of stage 1:

$$f_1(x) = \alpha_0 \, h_0(x) + \alpha_1 \, h_1(x)$$

New predictions

old prediction (at stage 0)

$e_i^o = y_i - f_1(x)$

2) $M_2 \underset{h_2(x)}{\longrightarrow} \left\{ x_i^o, e_i^o \right\}_{i=1}^n \quad e_i^o = y_i - f_1(x)$

Model at end of stage 2,

$$f_2(x) = \alpha_0 \, h_0(x) + \alpha_1 \, h_1(x) + \alpha_2 \, h_2(x)$$

new ...

prediction

$$f_2(x) = f_1(x) + \alpha_2 h_2(x) \Rightarrow \text{additive}$$

combination

new prediction  old prediction  current model

n) $$f_n(x) = \alpha_0 h_0(x) + \alpha_1 h_1(x) + \alpha_2 h_2(x) + \cdots + \alpha_n h_n(x)$$

$$\boxed{f(x) = f_{n-1}(x) + \alpha_n h_n(x)} \rightarrow \text{final model}$$

or

$$\boxed{f_n(x) = \sum_{i=1}^{n} \alpha_i h_i(x)}$$

$n \Rightarrow \#$ no of models $\Rightarrow$ hyperparameters

**** **Residual & loss $f^n$:**

$$L(y, f_n(x)) = [y_i - f_n(x)]^2$$

$$\frac{\partial L}{\partial f_n(x)} = \frac{\partial [y_i - f_n(x)]^2}{\partial f_n(x)} = -2[y_i - f_n(x)]$$

$$\partial F_n(x) \qquad\qquad \partial f_n(x)$$

$$-\frac{\partial L}{\partial f_n(x)} = \overbrace{[y_i - F_n(x)]}^{\text{error}}$$

negative gradient
or
pseudo-residual

## Gradient Boosting

$g|p \Rightarrow \langle x_i, y \rangle_{i=1}^{n} + $ differentiable loss $f^n$

0) $F_o = \underset{\gamma}{\text{argmin}} \; \left( \sum_{i=0}^{n} L(y_i, r) \right] \longrightarrow r = \bar{y_i}$

1) for $m=1$ to $M$ $\qquad \Rightarrow m \Rightarrow \#$ models

$$\gamma_m = - \left[ \frac{\partial L (y_i, F_{m-1}(x))}{\partial f_{m-1}(x)} \right]$$

$$L = (y_i - f_1(x)^2$$

for $m=2,$

$$r_m = - \left[ \frac{\partial l(y_i, f_i(x))}{\partial f_i(x)} \right]$$

2) $h_m(x)$ that can fit of $r_m$

     $\hookrightarrow$ train $h_m$ on $[x_i, r_{im}]_{i=1}^n$

         $f_2 = f_1(x) + \alpha_2 h_2(x)$

3) $r_m = \underset{r}{\text{argmin}} \left[ L(y_i, f_{m-1}(x_i) + r_m h_m(x)) \right]$

4)    $f_m(x) = f_{m-1}(x) + r_m h_m(x)$

      $\downarrow$     $\downarrow$     $\hookrightarrow$ Current

     New    old        model

     pred    pred

Hyperparameters $\Rightarrow$ M $\Rightarrow$ # models$\uparrow$ $\rightarrow$ overfitting$\uparrow$

   <u>Shrinkage</u>: $f_m = \underbrace{f_{m-1} + \nu \underbrace{r_m h_m}}$

            $\hookrightarrow$ learning rate

              $\hookrightarrow$ $0 < \nu < 1$

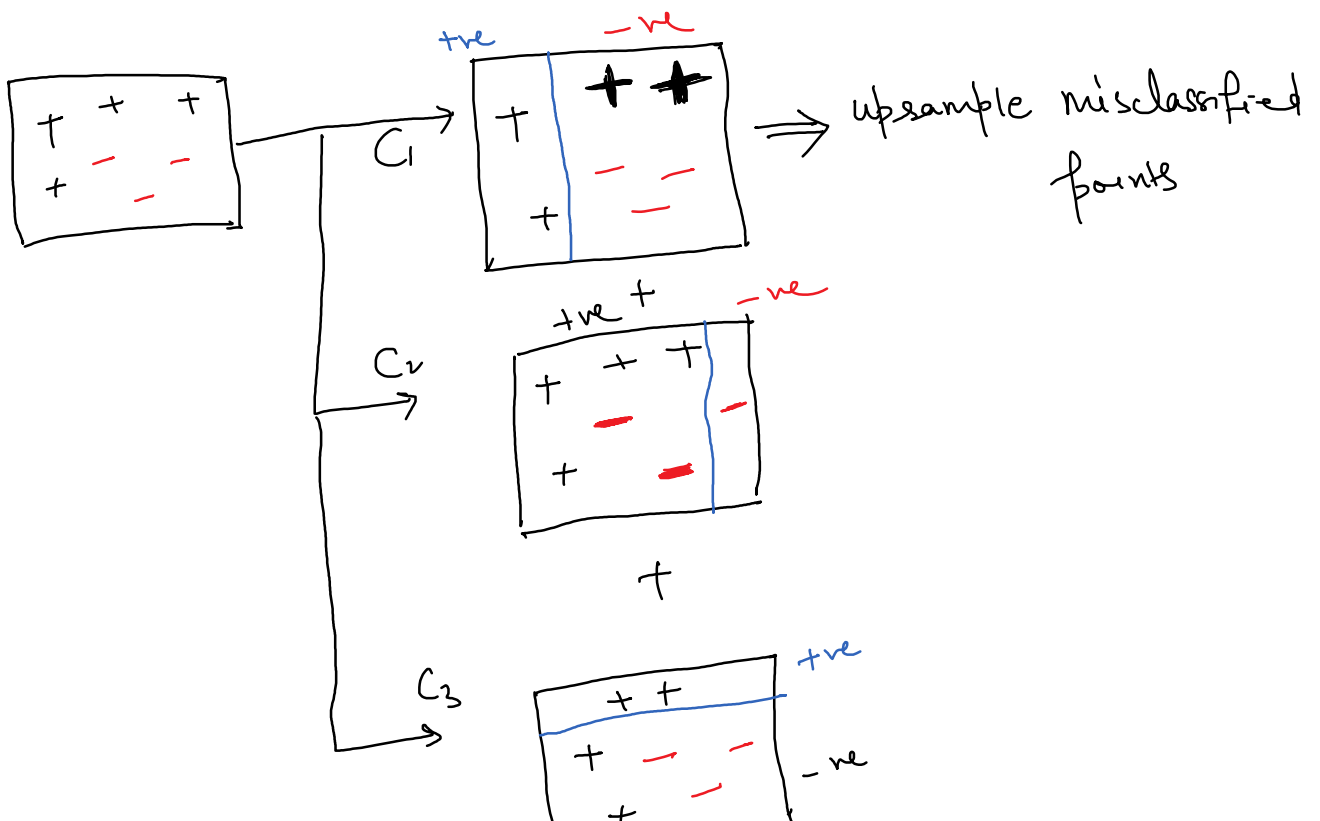$\gamma \rightarrow$ reduces $\gamma_m h_m \rightarrow$ to reduce overfitting

GBDT $\longrightarrow$ models are DTs

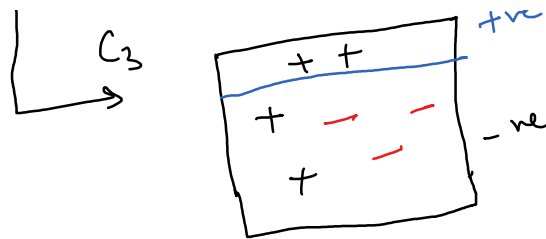$\hookrightarrow$ very slow $\longrightarrow$ optimized $\longrightarrow$ Taylor Series
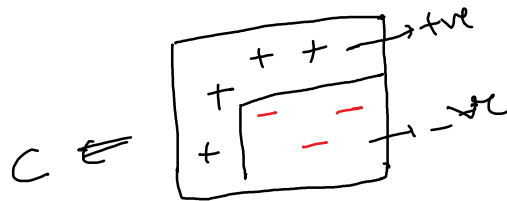
$\downarrow$

XgBoost

$\nwarrow$

¡ pip install XgBoost

Adaboost

Adaptive                    Boosting



$\Rightarrow$ upsample misclassified points

$$C = \gamma_1 C_1 + \gamma_2 C_2 + \gamma_3 C_3 + \, ---$$

| | $X_1$ | $X_2$ | $Y$ | $\hat{Y}$ | weight $= 1/n \Rightarrow \#\text{rows}$ |
|---|---|---|---|---|---|
| 1) | 3 | 9 | 1 | 1 | $0.2 = 1/5$ |
| 2) | 2 | 4 | 0 | 1 ✗ | 0.2 ✓ |
| 3) | 1 | 5 | 1 | 0 ✗ | 0.2 ✓ |
| 4) | 4 | 6 | 0 | 0 | 0.2 |
| 5) | 5 | 7 | 0 | 0 | 0.2 |

$\alpha = $ error rate

error = algebric sum of weights
of misclassified rows/points

error $= 0.2 + 0.2 = 0.4$

$$\alpha = \frac{1}{2} \ln \left( \frac{1 - \text{error}}{\text{error}} \right)$$

$$= \frac{1}{2} \ln \left( \frac{1 - 0.4}{0.4} \right) = 0.2$$

new weight for misclassified pts $= e^{+\alpha} \times$ old weight

$$= e^{0.2} \times 0.2 = 0.24$$

new weight of correctly classified points $= e^{-\alpha} \times$ old weight

$$= e^{-0.2} \times 0.2$$
$$= 0.16$$

| | $X_1$ | $X_2$ | $Y$ | $\hat{Y}$ | Weights | New Weights | Normalized Weights | Range |
|---|---|---|---|---|---|---|---|---|
| ① | 3 | 9 | 1 | 1 | 0.2 | 0.16 | $0.16/0.96 = \frac{1}{6} = 0.167$ | $0 - 0.167$ |
| ② | 2 | 4 | 0 | $1^x$ | 0.2 | 0.24 | $0.24/0.96 = 0.25$ | $0.167 - 0.417$ |
| ③ | 1 | 5 | 1 | $0^x$ | 0.2 | 0.24 | $0.24/0.96 = 0.25$ | $0.417 - 0.667$ |
| ④ | 4 | 6 | 0 | 0 | 0.2 | 0.16 | $0.16/0.96 = \frac{1}{6} = 0.167$ | $0.667 - 0.834$ |
| ⑤ | 5 | 7 | 0 | 0 | 0.2 | 0.16 | $0.16/0.96 = \frac{1}{6} = 0.167$ | $0.834 - 1$ |
| | | | | | $\underline{1}$ | $\underline{0.96}$ | $\underline{1}$ | |

$$0.167 \overset{0.500}{\underline{\qquad}} 0.667$$

Randomly any no b/w 0 & 1

$$0.1 \quad 0.4 \quad 0.5 \quad 0.6 \quad 0.7$$
$$① \quad ② \quad ③ \quad ③ \quad ④$$

New
Dataset

①
upsampling { ②
③ } → misclassified rows in previous model
③
④