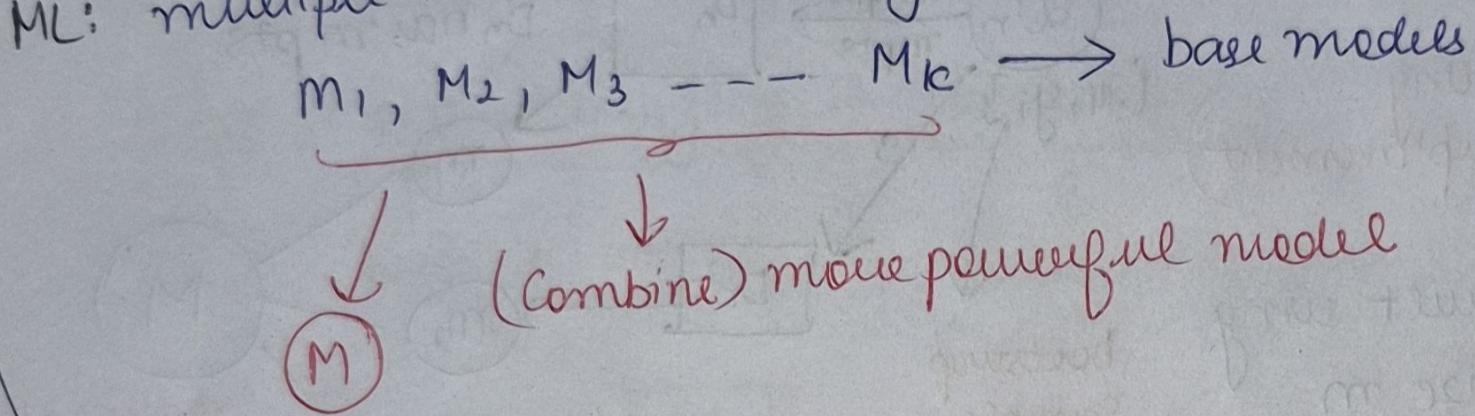


## Ensembles

→ group of musicians

ML: multiple models used together



## 4 Types

- ① Bagging (Bootstrapped Aggregation)
  - very useful in real world
- ② Boosting
  - high performing
- ③ Stacking
  - very powerful
- ④ Cascading
  - kaggle - most

## key aspect

$m_1, m_2, m_3, \dots, m_k$

the more different these models are, the better you can combine them.

{ Problem:  $m_i$ : expert

# Bagging (Bootstrap Aggregation)

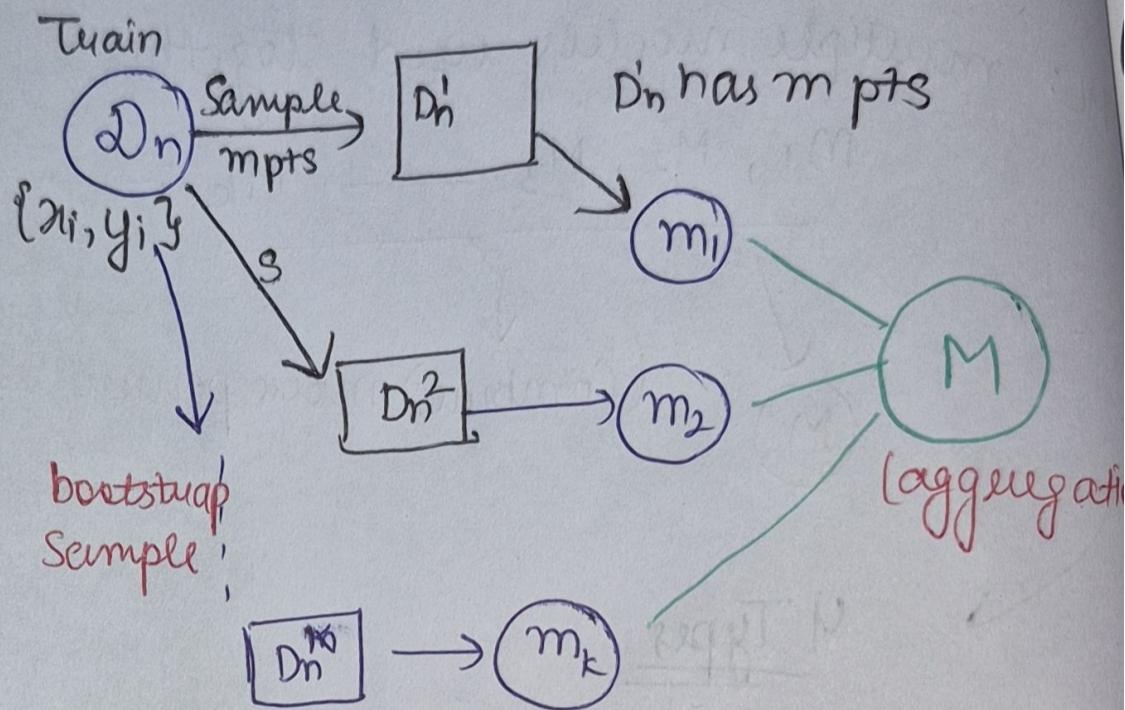
↳ Statistics

(Random Forest)

## Intuition

(Sampling with replacement)

$m_i$  is built using  
 $D_n$  of size  $m$   
( $m \leq n$ )



→ Each model  $m_i$  has seen different subset of data

Aggregation : Classification : Majority vote  
Regression : mean / median

Bagging → can reduce variance in model without impacting the bias

$$\text{model error} = \text{Bias}^2 + \text{Var}$$

base model ( $M_i$ ): low bias, high var model

Bagging ( $M_i$ 's) → low bias; reduced variance

Bagging bunch of low bias, high var models

→ DT of depth

↳ high var

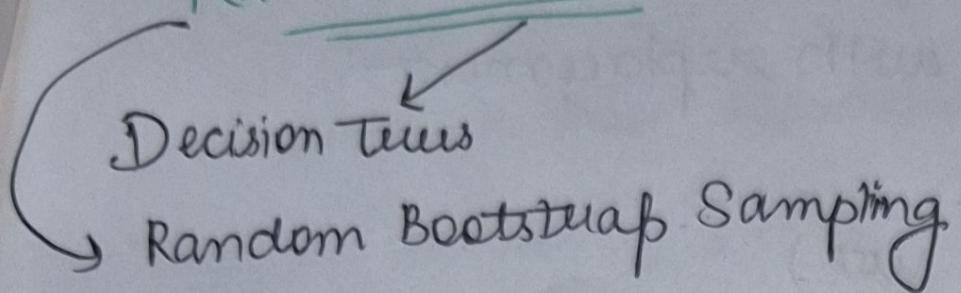
low bias

↳ Random Forest

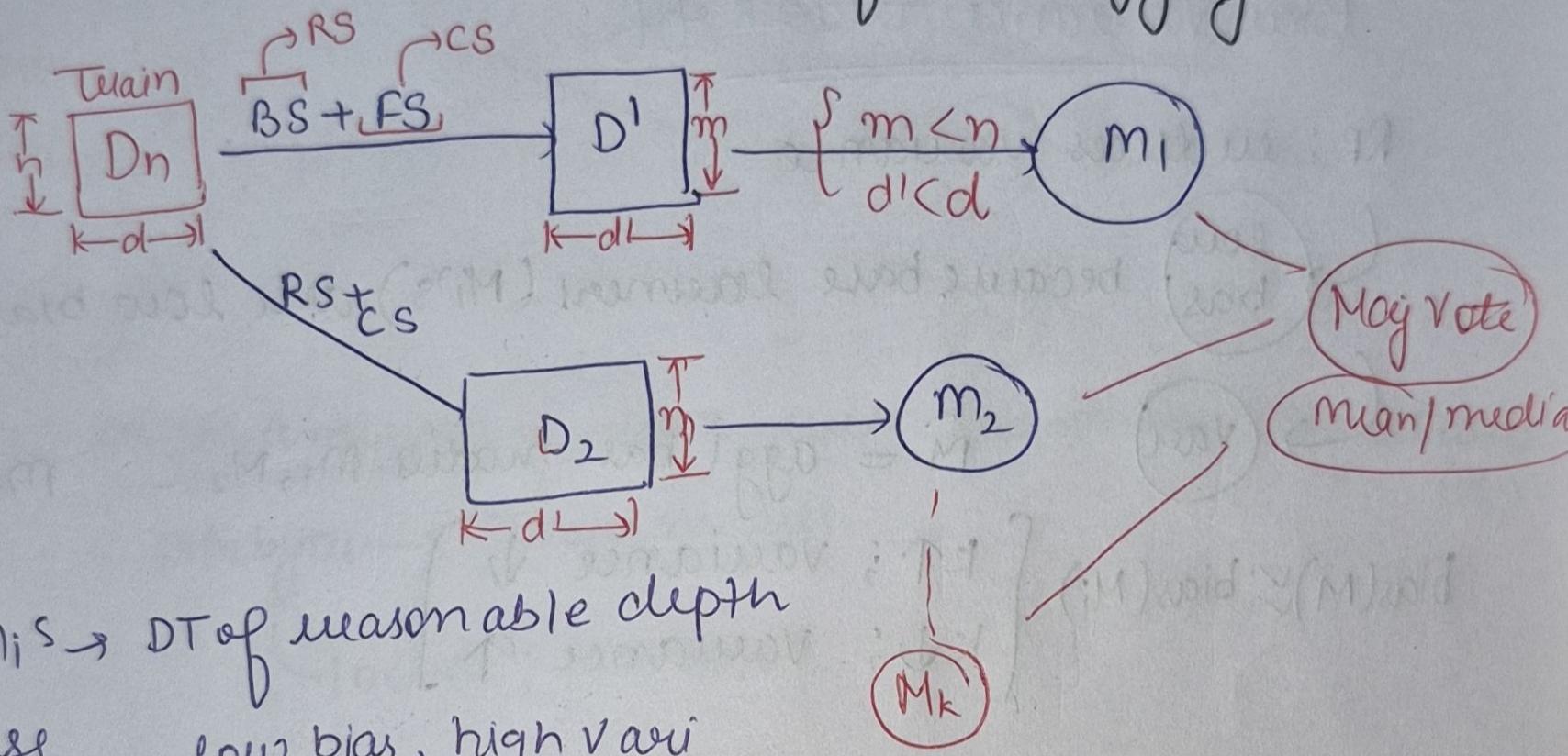
Combine them using bagging

(low bias, reduced var)

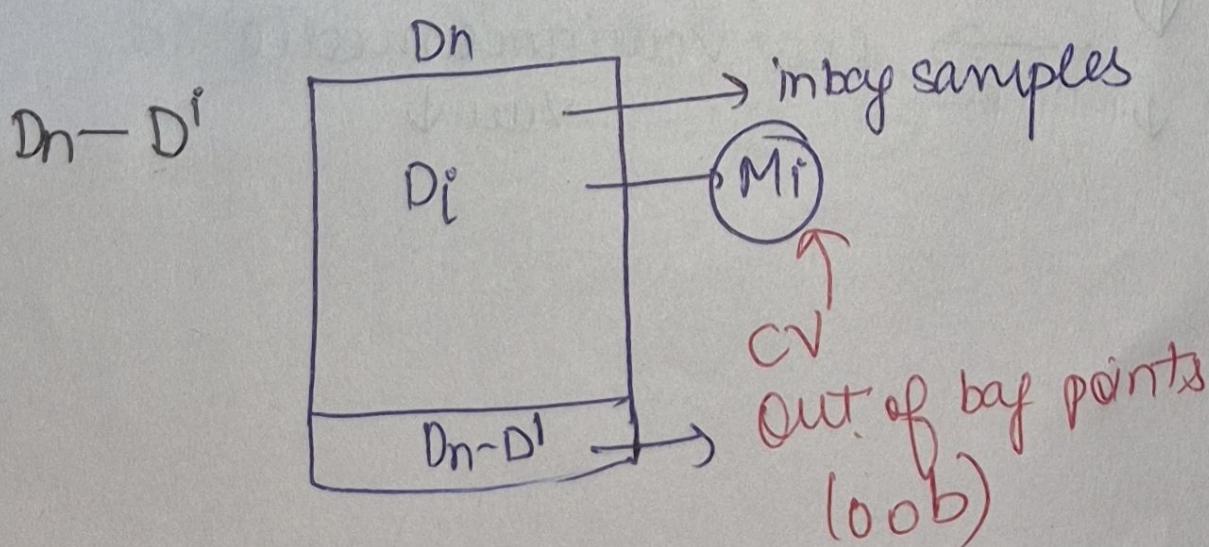
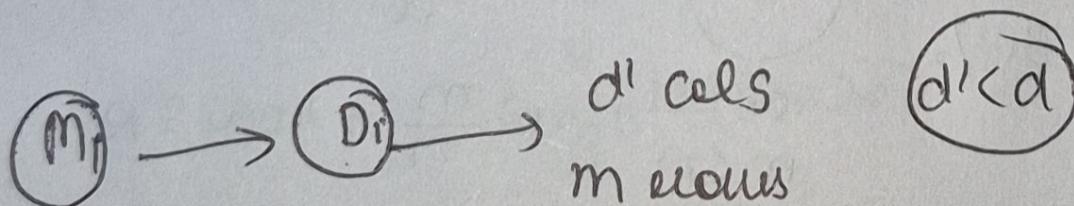
# Random Forest (Bagging)



RF: DT + Bagging + Col. Sampling  
 ↑  
 base RS with replac  
 ↓ feature bagging



$m_i \rightarrow$  DT of reasonable depth  
 base learner  $\rightarrow$  low bias, high vari  
 reduces variance



$m_i \rightarrow D_i$   
 $D_n - D_i \rightarrow$  OOB samples =  $CV - \frac{\text{dataset}}{m_i}$

RF: DT base learners

- + row sampling with replacement
- + col sampling
- + agg (majority vote)  
↳ (mean / median)

### Bias Variance trade off

RF: reduces variance

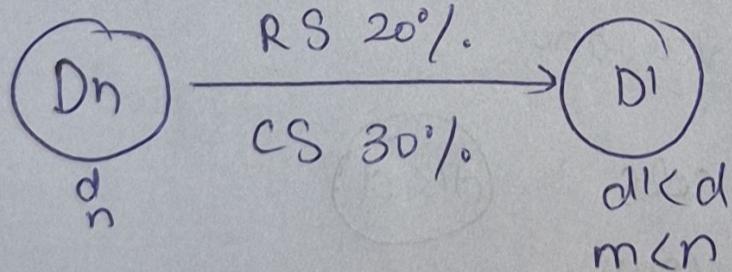
low bias

because base learners ( $M_i$ 's) are low bias

Vari

$$\hat{M} = \text{agg}(\text{base models } M_1, M_2, \dots, M_k)$$

$$\text{bias}(M) \leq \text{bias}(M_i) \begin{cases} K \uparrow; \text{ variance } \downarrow \\ K \downarrow; \text{ variance } \uparrow \end{cases}$$



$$\frac{d1}{d} = \text{colsR}$$

$$\frac{m}{n} = \text{row S.R}$$

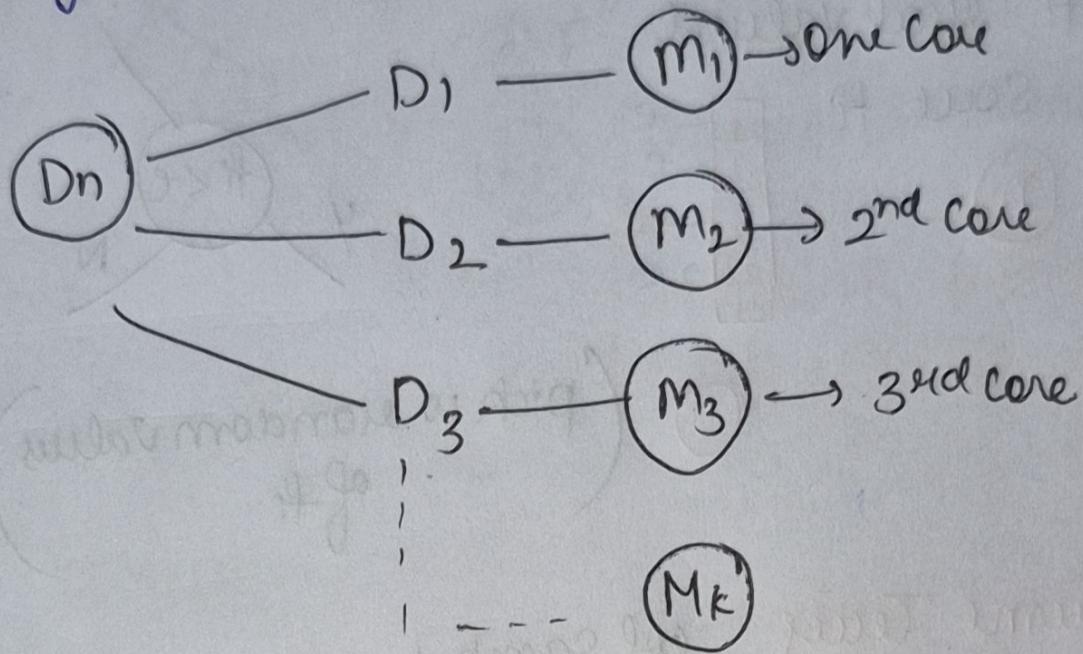
col. SR  $\downarrow$   $\rightarrow$  low variance model  
row SR  $\downarrow$  Vari  $\downarrow$

## Train And Run Time Complexity

RF with K base learners (DT)

Train:  $O(\text{negn } d * K)$  multicore

Trainably  
parallelizable



large amounts of data with reasonable # features (d)

peta bytes → Train modules  
10000; 1000 box

Run Time:  $O(\text{depth} * K)$

large (10 to 20)

Space:  $O(DT * K)$   
If else → takes very large time space

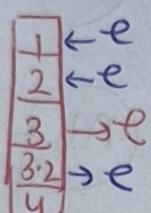
Extremely Randomized Trees → Try out random sample of possible values to determine best split

RF → col-sampling, row sampling, agg determine best DT

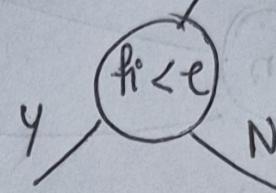
Extreme Trees:

RF/DT {  $f_i$ : real valued  
sort  $f_i$

(n)



$m_i \rightarrow$



(pick 10 random values of  $f_i$ )

DT/ RF  
↓  
try all possible values of  $f_i$  to determine  $e$

Extreme Trees

col samp + row samp + agg + randomization when selecting "e"

Randomization as a way to reduce variance

RF → CS & RS

ET → CS + RS + randomization (e)

(reduce variance better than RF)

Caus:

RF:

$\underbrace{\text{DT} + \text{RS} + \text{CS}}_{\text{reduce variance}} + \text{agg}$

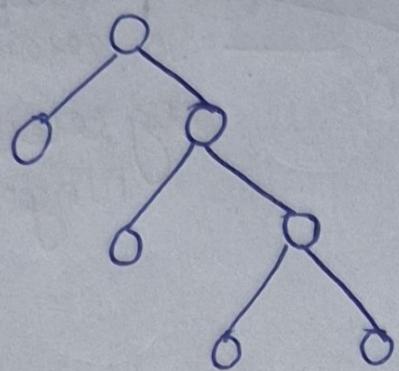
do not handle large dim; cate features with many categories

causes for DT

① Bias Variance of RF  $\propto k \# \text{base learners}$   
L deep trees.

## ② Feature Imp

DT ; fi: - overall reduction in entropy or gini Impurity because of this feature at various levels in DT



(RF) : overall reduction in Hoenig because of fi @ various levels of each of Mi's

## Boosting Intuition

Bagging : high var, low bias base models + randomization + agg (CS, RS)

Boosting low var, high bias + additively combine + reduce bias while keeping var low

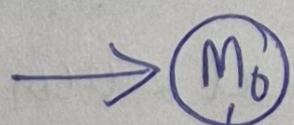
$$\text{error} = \text{bias}^2 + \text{Var} + \varepsilon$$

Core Idea how Boosting reduce bias

$$D_{\text{Train}} = \{(x_i, y_i)\}_{i=1}^n$$

Classif / Regress

①  $D_{\text{Train}}$   
 $\{(x_i, y_i)\}_{i=1}^n$



high train error

high bias, low var  
 {eg: DT which is shallow }  
 low depth

①  $D_{\text{train}} \rightarrow M_0$   $y = h_0(x)$  high bias  $\Rightarrow$  large train error

⑥  $y_i - h_0(x_i) = \text{error}_i$  simple diff error  
 $\{x_i, y_i, \text{error}_i\}_{i=1}^n$  sq. error  
 $\log \text{error}$   
 $\text{hinge error}$

①  $M_1 \leftarrow \{x_i, \text{error}_i\}_{i=1}^n$   
 $h_1(x)$

$f_1(n)$  = model at end of Stage 1.

$f_1(n) = \alpha_0 h_0(n) + \alpha_1 h_1(n)$  weighted sum of 2 base models

②  $\{x_i, \text{error}_i\}_{i=1}^n \rightarrow M_2$   $h_2(n)$   
 $\downarrow$   
 $y_i - f_1(x_i)$

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

end of Stage K:

$f_K(n) = \sum_{i=0}^K \alpha_i h_i(n)$  trained to fit residual error @ end of the previous stage

additive weighted model

$h_i(n) \leftarrow \{x_i, \text{error}_i\}$

Train error  $\downarrow \rightarrow$  bias  $\downarrow$  residual error @ end of stage (i-1)

ends up having a low resid error Gradient Boosted DT

# Residuals, Loss func & Gradients ↗ Gradient Boosting

$$F_k(x) = \sum_{i=0}^k \alpha_i h_i(x)$$

Residual  
@ end of stage k

$$\text{error}_i = y_i - F_k(x)$$

↑ residual

$$M_{k+1} \leftarrow \{x_i, \text{error}_i\}$$

## Loss Minimization

logistic loss ← log reg / clust  
 der Reg ← sq. loss  
 SVM ← hinge loss

regression

$$L(y_i, F_k(x_i)) = (y_i - F_k(x_i))^2$$

↙  
SQ loss

$$\frac{\partial L}{\partial F_k(x_i)} = \frac{\partial L}{\partial z_i}$$

$$\begin{cases} \text{det} \\ F_k(x_i) = z_i \end{cases}$$

$$= \frac{\partial}{\partial z_i} (y_i - z_i)^2$$

$$\frac{\cancel{\partial}}{\partial F_k(x_i)} = (-1) * \cancel{\partial} (y_i - z_i)$$

$$-\underbrace{\frac{\partial L}{\partial F_k(x_i)}}_{\text{negative derivative}} = \cancel{\partial} (y_i - F_k(x_i))$$

↙ residual

negative derivative

\* neg. gradient ↗ residual  
 ↓  
pseudo residual

## gradient boosting

$$h_i(x)$$

$$\leftarrow x_i, e_{ei} = \underbrace{y_i - f_{i-1}(x)}_{\text{err}} \quad \text{err} \rightarrow \text{residual}$$

$e_{ei} \rightarrow \text{residual}$

$e_{ei} \rightarrow \text{pseudo residual}$

$$-\frac{\partial L}{\partial f_{i-1}(x)}$$

Let us have  
any less fn  
which is  
differentiable

$$(RF) \rightarrow$$

$$(GB) \rightarrow \min \text{any less} \rightarrow \text{which is differentiable}$$

Super powerful

Stage i

$$f_{i-1}(x)$$

$$(x_i, e_{ei}) \rightarrow \underbrace{h_i(x)}_{\text{pseudo residual}}$$

$$M_i$$

pseudo residual

$$-\frac{\partial d}{\partial f_{i-1}(x)}$$

Ques Is pseudo-residual  $(-\frac{\partial d}{\partial f_i(x)})$  same as residual  
for non-squared loss func?

ans

Ans

$$\frac{-\partial L}{\partial F_k(x)} \leftarrow \text{residuals for non squared loss func.}$$

$$\frac{-\partial L}{\partial F_k(x)} = \alpha (y_i - F_k(x_i))$$

L :- logg loss

$$\frac{\partial L}{\partial F_k(x)} = F_k(x) = \hat{y}_i$$

$$= p_i - y_i \quad \text{where} \quad p_i = \frac{1}{1 + \exp(-\hat{y}_i)}$$

$p_i = \text{prob of } x_i \in \text{class 1}$

$$= \frac{1}{1 + \exp(-F_k(x))}$$

pseudo residual

$$\frac{-\partial L}{\partial F_k(x)} = p_i - y_i$$

$\downarrow$

$y_i = 0 \text{ or } 1$

$P(x_i \in 1)$

# Gradient Boosting      Regularization & Shrinkage

$$F_m(x) = h_0(x) + \sum_{m=1}^M \gamma_m h_m(x)$$

M: # base models

→ Hyperparameter

C.V. ↘ # base models ↑ ⇒ overfit ↑ ⇒ var ↑

M ↑ ↗ bias ↓ but could var ↑

## Shrinkage

### Shrinkage

$$F_m(x) = F_{m-1}(x) + v \cdot \gamma_m h_m(x), \quad 0 \leq v \leq 1, \quad v \leq 0.1$$

v is small  
↳ 0.0001

v ↓ overfitting ↓ var ↓

(v↑ ⇒ overfitting↑)

## Train & Run Time Complexity

Train:  $O(n \lg nd * M)$       M: # base learners

GBDT take more time to Train than RF

## Run Time

GBDT:  $O(\text{depth} * M)$

↓

Small in GB

$h_m(x)$

$O(\text{depth} * M + M)$

↪  $O(\text{depth} * M)$

Space

$O(1 \text{ store each tree} + \gamma_m)$

if else

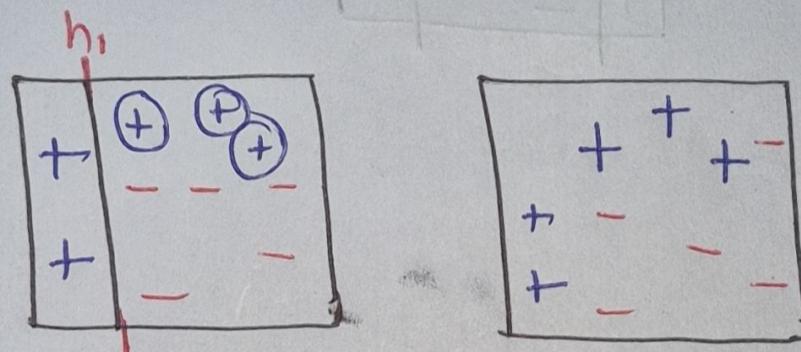
AdaBoosting

→ image processing (face detection)

+	+	+	-
+	-	+	-
+	-	-	-

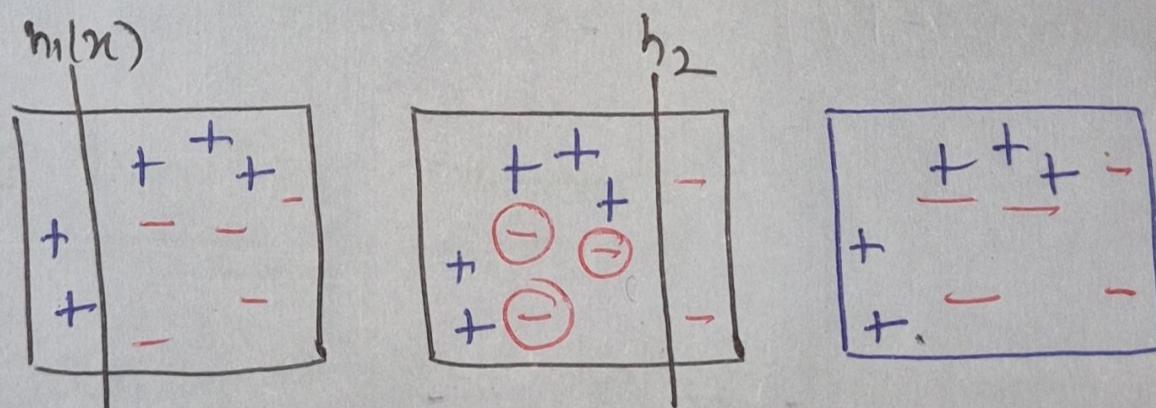
1st Round

$\alpha_1 h_1(x)$

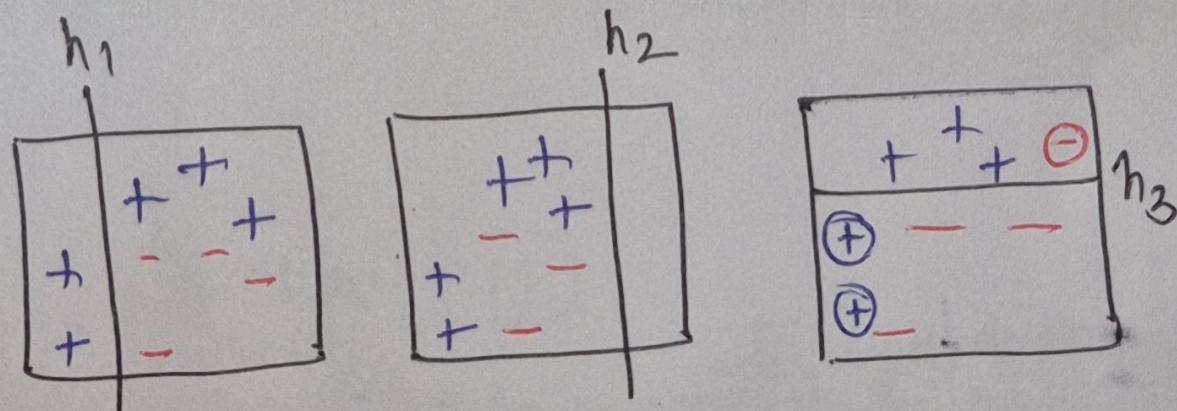


2nd Round

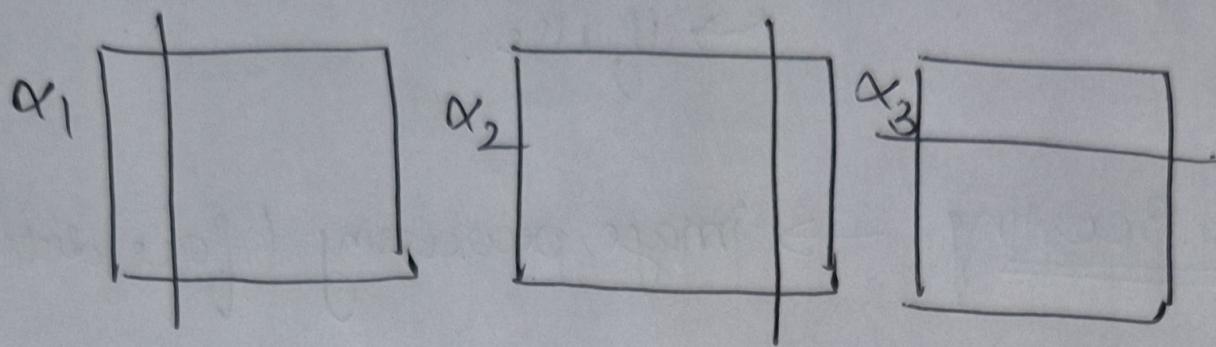
$\alpha_2 h_2$



3rd Round



final Model



$$f_3(n) = \alpha_1 h_1 + \alpha_2 h_2 + \alpha_3 h_3$$

