

## Ensemble Models

- 1 -

✧ group of musicians performing together

In Machine Learning  $\rightarrow$  Multiple models used together

$\{M_1, M_2, M_3, \dots, M_K\} \rightarrow$  base models

$\downarrow$   
(Combine) more "powerful" model

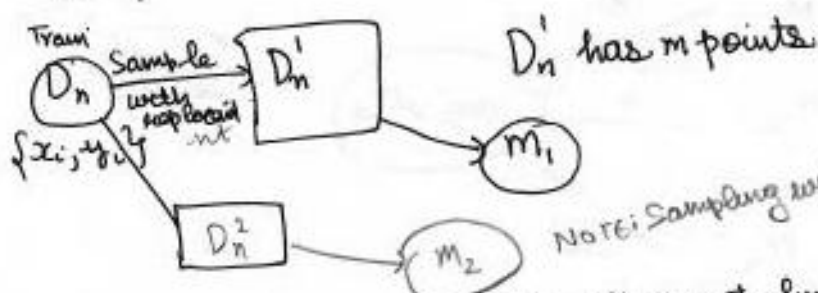
4 Types  $\div$   $\left\{ \begin{array}{l} \text{Bagging (Bootstrapped Aggregation)} \\ \text{Boosting} \\ \text{Stacking} \\ \text{Cascading} \end{array} \right.$

$\rightarrow$  high performing  
 $\rightarrow$  V. powerful  
 $\rightarrow$  Kaggle  $\rightarrow$  most.

1) Bagging (Bootstrap Aggregation)

$\rightarrow$  statistics (Taking new sample from one sample)

Intuition: Random Forest is one of the most popular model used for logging.



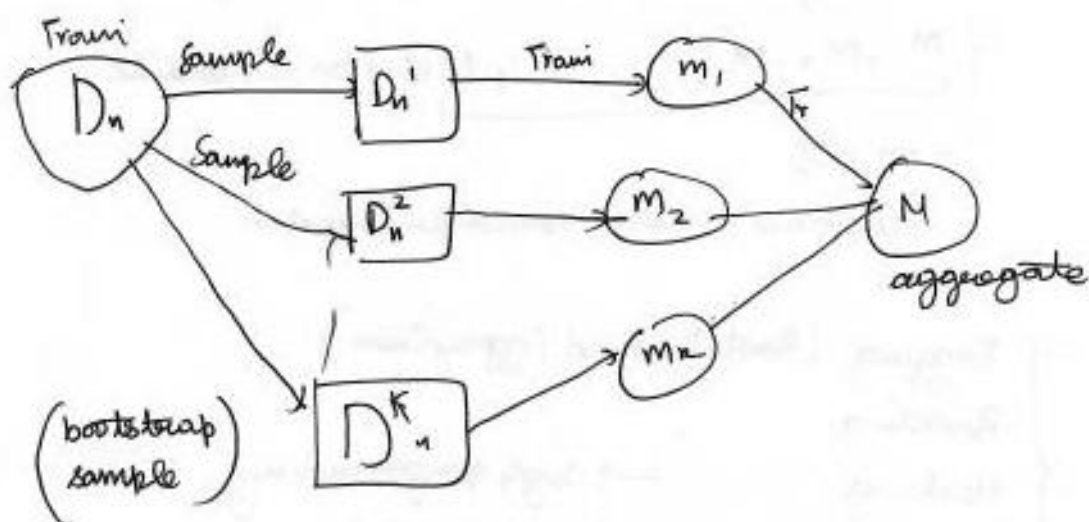
NOTE: Sampling with replacement

NOTE: Sampling with replacement  $\rightarrow$  taken the point from dataset but it remains in dataset from which it has been taken as well as original dataset.



$M_i$  is built using  $D_n^i$  of size  $m$  ( $m \leq n$ )

$\Rightarrow$  Each model  $M_i$  has seen a different subset of data.

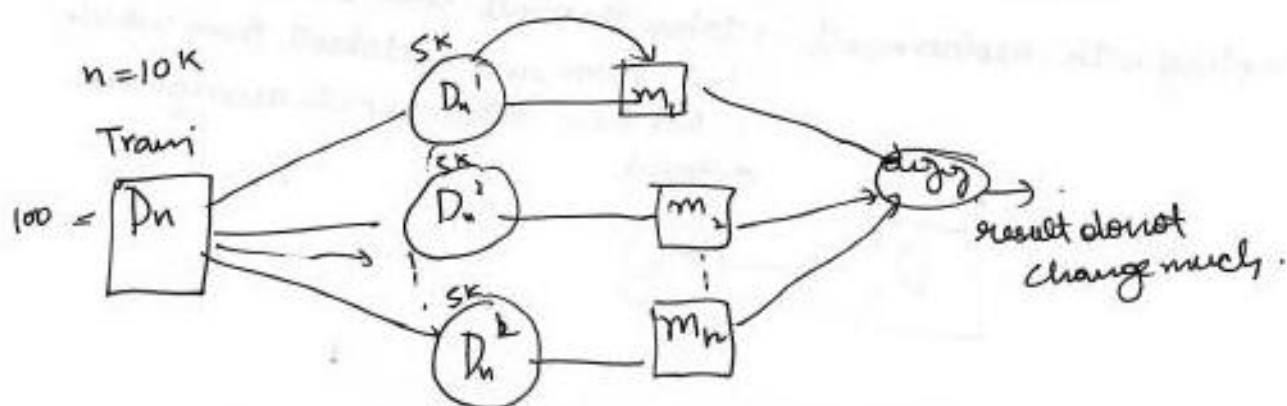
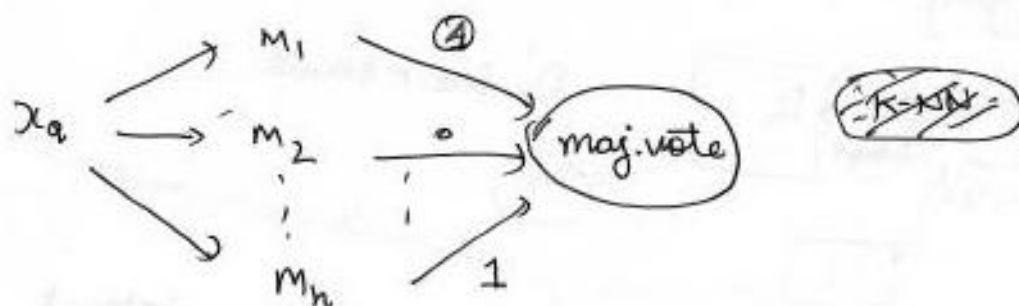


Aggregation : 1) Classification  $\Rightarrow$  Majority vote

2) Regression  $\Rightarrow$  mean / median

$\star$  After training this model  $\uparrow$

I get  $M_1, M_2, \dots, M_k$  we get majority vote.



# Limitations

1) model changes a lot with change in dataset (Variance)

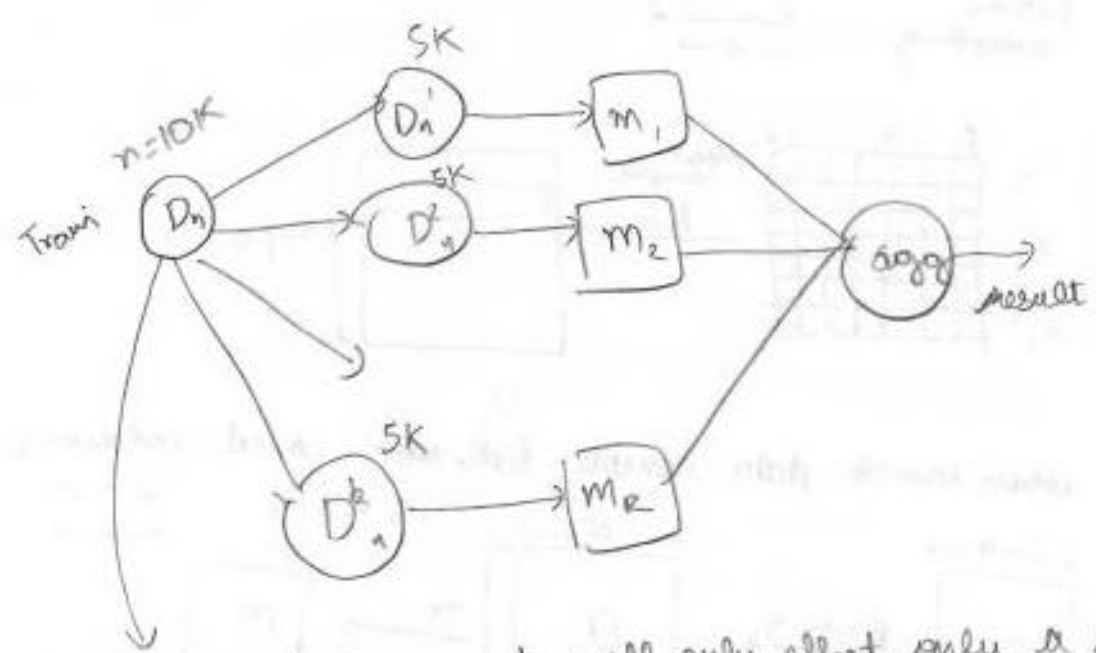
Bagging → can reduce variance in a model without impacting the bias.

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

Base-model ( $M_i$ ) ; - low bias, high Var model

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

Variance → {model changes a lot with change in  $D_n$ }



⇒ If dataset changes it will only affect only a small subset of models which will not affect the aggregate function a lot or by removing a bunch of points from  $D_n$  (train dataset) model don't change much.

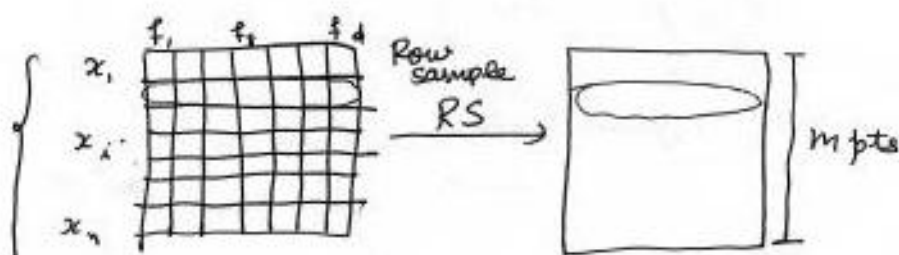
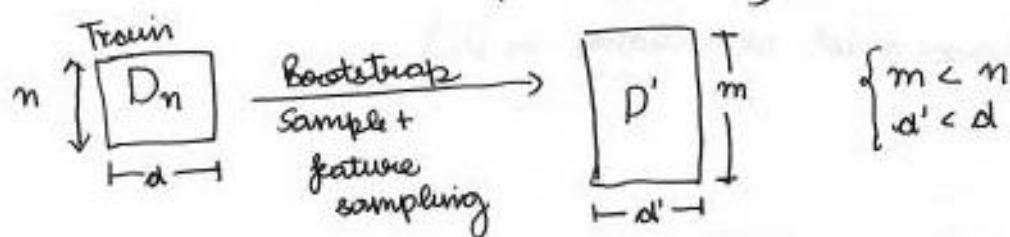
NOTE: Bagging takes bunch bunch of low-bias, high variance models and combine them using bagging and convert them into low bias and low variance models.  
+ reduced.

# Random Forest (Bagging Technique Used)

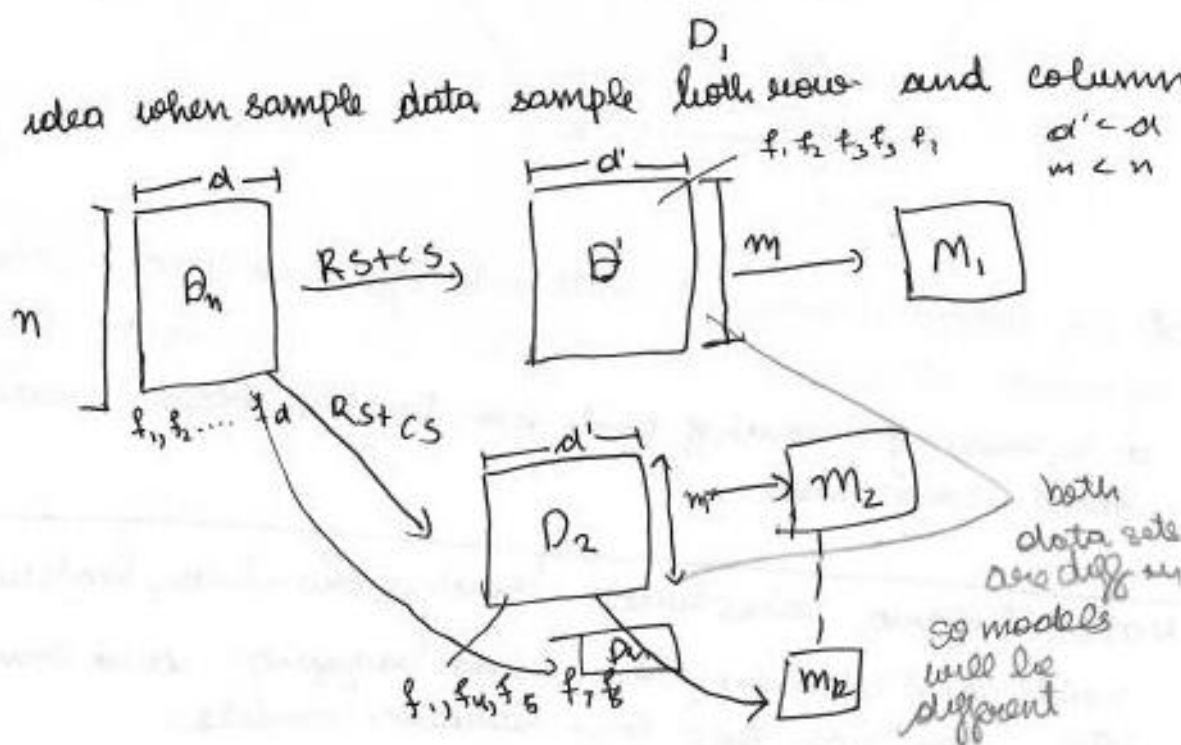
Decision Trees : As decision tree which have large depth has high variance (overfitting) and lower bias it is used to bagging technique to reduce variance and that technique is called Random Bootstrap Sampling

RF : D.T  $\uparrow$  Bagging + column sampling <sup>feature bagging</sup>

$\uparrow$   
 base ( $\uparrow$  Row sampling with replacement)



Core idea when sample data sample both row and column

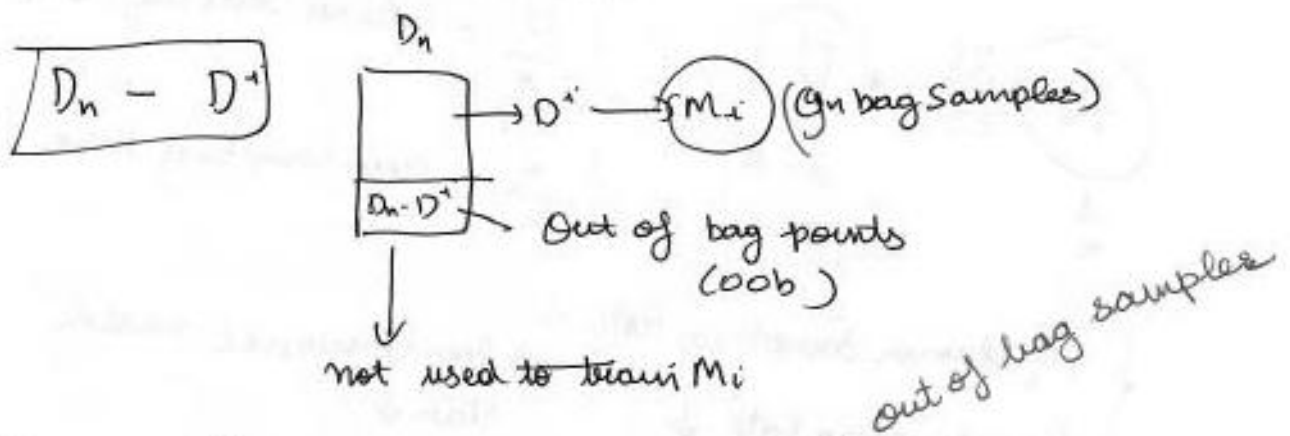


$M_1, M_2, \dots, M_k$  are totally different

so after that we take the majority among  $M_1, M_2, \dots, M_k$

low bias; high variance

$M_i \rightarrow D_i \rightarrow d' \text{ cols } m \text{ rows}$   $d' < d$



Since model  $M_i$  using point in  $D_i$

$D_n - D_i \rightarrow \text{oob samples} = \text{CV-dataset}$

RF : DT base learners of reasonable depth.

- + doing row sampling with replacement
- + column sampling
- + agg. (majority vote in case of classification)  
(mean or median in case of regression)

→ Bias - variance tradeoff

RF → reduce ~~variable~~ variance due to bagging

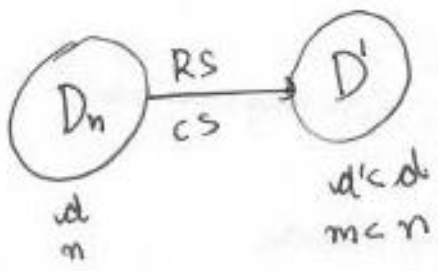
- ↳ low bias because base learners ( $M_i$ 's) are low bias
- ↳ var

$$M = \text{agg}(M_1, M_2, M_3, \dots, M_k)$$

These all are high depth Decision Tree having high variance same bias.

bias (m)  
" " " "  
bias (individual model)

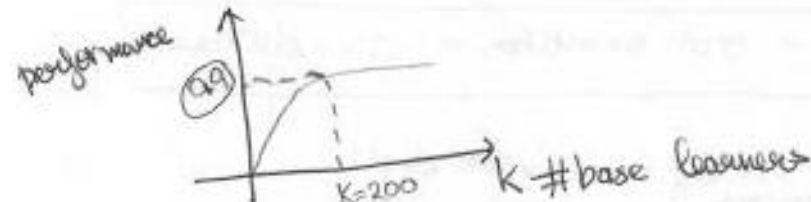
$K \uparrow$ ; variance  $\downarrow$   
 $K \downarrow$ , variance  $\uparrow$



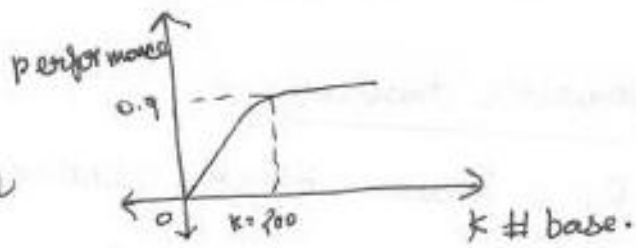
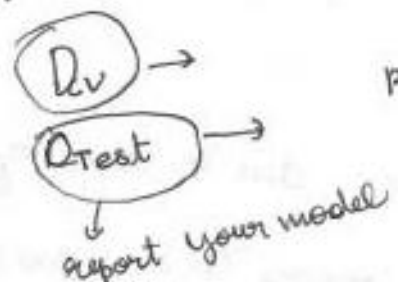
$\left\{ \frac{d'}{d} = \text{Column Sampling Ratio} \right.$

$\left\{ \frac{m}{n} = \text{Row Sampling Ratio} \right.$

$\left\{ \begin{array}{l} \text{If column Sampling Rate } \downarrow \\ \text{Row Sampling Rate } \downarrow \end{array} \right. \rightarrow \text{low variance model} \quad \text{Var } \downarrow$



Only variance reduces but bias remains the same that is bias of trained model  $M$  remains same as the bias of individual model  $M_i$

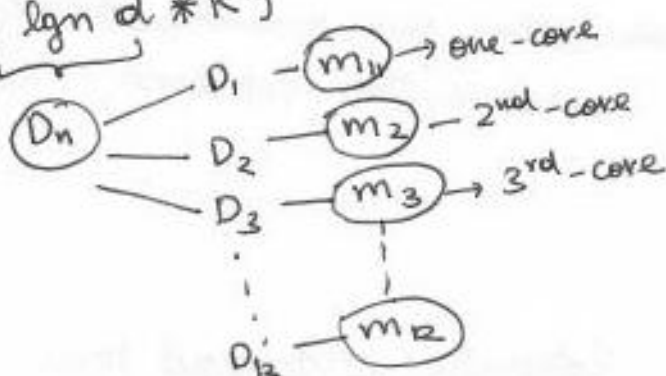


→ most important hyperparameter here is  $k$ . That is number of samples.. (number of base learners)

## Train and Run Time Complexity.

RF with  $k$  base-learners (Decision Tree)

Train:  $O(n \lg n d * k)$



Trivially parallelly

multi-core

System

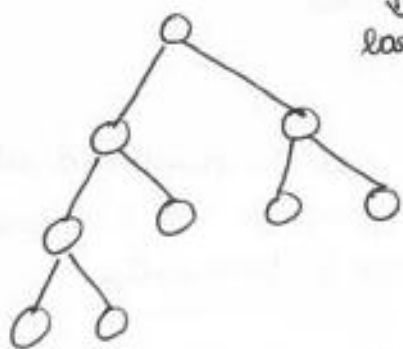
↓  
we can train each model independently parallelly.

When large amount of data with reasonable number of features  
(d) you can train 1000 of trees very fast.

petabytes of data → Train models

Run time:  $O(\text{depth} * k)$

large (10 to 20)



Space:  $O(DT * k)$



Decision Tree: with 200 Trees of depth = 6

Code:

sklearn.ensemble.RandomForestClassifier

Class sklearn.ensemble.RandomForestClassifier (n\_estimators=10, <sup>K</sup>  
 criterion = 'gini', max\_depth = None, min\_samples\_split=2,  
 min\_samples\_leaf=1, min\_weight\_fraction\_leaf=0.0, max\_features='auto',  
 max\_leaf\_nodes = None, min\_impurity\_decrease=0.0, min\_impurity  
 split = None; bootstrap=True, oob\_score=False, ...

Regular decision trees

Extremely Randomized trees

RF → Column sampling ; Row sampling ; aggregation

{ Extreme Trees: →  $M_i$  }

$f_i$ : real valued feature then we can put threshold on feature  $i$ .

RF / Decision Tree {  $f_i$ : real valued values

sort  $f_i$

1	← $\tau$
2	← $\tau$
3	← $\tau$
4	← $\tau$
5	← $\tau$
6	

pick the threshold which  
 give me the max min  
 gini impurity

NOTE: They say instead trying threshold on all values  
 try out on random subset sample of possible values.



Extreme Trees : Col. sampling + Row sampling + aggregation  
+ randomization when selecting

This is the part different from  
Random forest.

Randomization is a way to reduce variance.

R.F  $\rightarrow$  C. Sampling + Row. Sampling

Extremely Randomized trees  $\rightarrow$  Column Samp. + Row Sampling

+ random  $\tau$

Threshold.

They reduce variance better than Random forest but  
bias might increase.

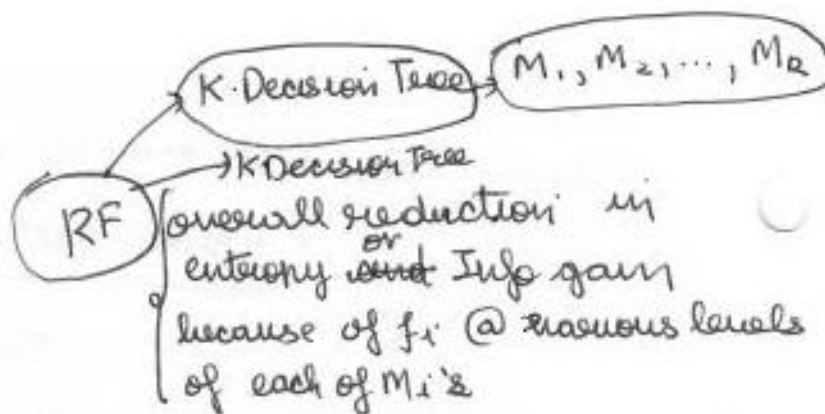
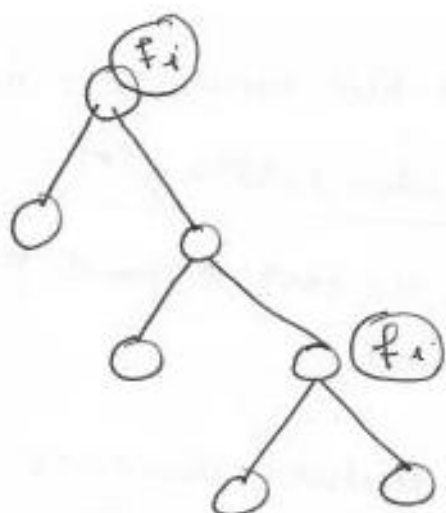
### ◇ Random Forest - Cases

RF : Decision Tree + Row Samp + Column Samp + aggregation reduce variance  
 $\swarrow$   $\searrow$   
 Cases for DT Decision Tree  $\rightarrow$  depth  
 $\searrow$   $\swarrow$   
 $\hookrightarrow$  large dim; Cat. features with many categories

Bias-Variance of Random forest  $\sim K = \#$  of base learners.  
 $\hookrightarrow$  ~~deep trees~~ Train but control variance.

### ② Feature Importance $\rightarrow$

DT  $\div$   $f_i$  : overall reduction in Entropy or  
 Gini Impurity because of this feature  
 @ various levels in the Decision Tree



If feature used in lot of decision tree  
So overall importance of feature increases

⇒ Where decision tree work well Random forest also work well and vice versa. (but not in case of bias variance tradeoff)

## Boosting :

Bagging → high variance & low bias model + randomization (R; C)  
+ aggregation  
to reduce variance.

Boosting ⇒ low variance, <sup>reduce the bias here</sup> high bias model + <sup>additively combining</sup>  
[reduce bias while keep our var. low]

$$\text{Error} = (\text{bias})^2 + \text{Var} + \epsilon$$

→ irreducible error.



$$F_2(x) = \alpha_0 h_0(x) + \alpha_2 h_2(x) + \alpha_1 h_1(x)$$

end of stage  $k$  :

$$F_k(x) = \sum_{i=0}^k \alpha_i \underbrace{h_i(x)}_{\substack{\text{trained to fit the residual error.} \\ \text{additive weighted model}}}$$

$$h_i(x) \leftarrow \{x_i, \text{err}_i\} \quad \begin{matrix} \nearrow y_i - F_{i-1}(x) \\ \uparrow \text{residual error @ end of} \\ \text{stage } (i-1) \end{matrix}$$

◇ Final model

$$K=10 \quad F_K(x) = \sum_{i=0}^K \alpha_i \underbrace{h_i(x)}_{\text{residual error}}$$

ends up having.

a low residual error

Training error reduces

After every iteration residual error reduces.

→ Gradient Boosted Decision Tree (GBDT)

→ Adaptive Boosting → Images (Face detection algorithm)

## BOOSTING NOTES

Reduces bias of the weak learners.

Build learners sequentially.

Sample misclassified by previous learners weighted more in subsequent learners

## Residuals, loss-functions &amp; Gradients

$$F_k(x) = \sum_{i=0}^k \underbrace{\alpha_i}_{\text{weight}} \underbrace{h_i(x)}_{\text{feature}}$$

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

residual  
@ end of  
stage k

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

↑  
Residual

Whenever we solve machine learning problem  
↓

we try to minimize loss function

for Square-loss

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

$$\left\{ \begin{array}{l} \text{Let} \\ F_k(x_i) = z_i \end{array} \right.$$

SQ-loss

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

$$= (-1) * 2 * (y_i - z_i)$$

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

$$- \underbrace{\frac{\partial L}{\partial F_k(x_i)}}_{\text{derivative of loss function w.r.t model at stage k}} = 2 \underbrace{(y_i - F_k(x_i))}_{\text{Residual}}$$

-ve derivative  
of loss function  
w.r.t model at stage k.

\* negative gradient  $\approx$  residual

Gradient boosting is so powerful due to pseudo-residual that is negative gradient

$$h_i(x) \leftarrow (x_i, \text{err}_i) = y_i - F_{i-1}(x)$$

$\text{err}_i \rightarrow$  residual;  
 $\text{err}_i \leftarrow$  pseudo-residual

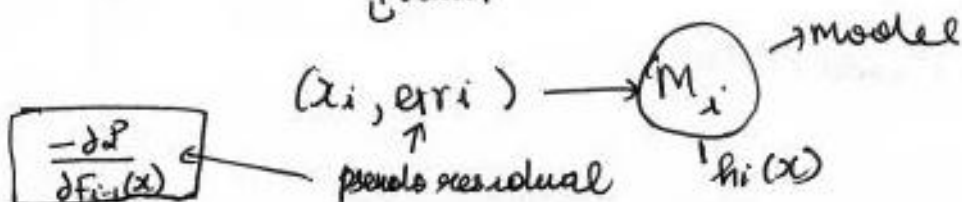
① negative gradient  $\approx$  residual  
 $\uparrow$  pseudo-residual

$\Rightarrow$  why training on pseudo-residual but why not residual?  
 pseudo-residual has an advantage that is ~~any~~ it works on any loss function which is differentiable

for eg RF  $\rightarrow$  min hinge loss  $\rightarrow$  not possible

but with  
 Gradient Boosting  $\rightarrow$  min (any loss)  $\rightarrow$  differentiable  
 $\uparrow$   
 super powerful

Stage: i  $F_{i-1}(x)$   
 $\downarrow$  train



Gradient Boosting  $\rightarrow$  If you have base model having ~~low~~ high bias  $\rightarrow$  combine pseudo-residual  $\hookrightarrow \frac{\partial L}{\partial F_{i-1}(x)}$

Pseudoresidual for log-loss and classification

$\frac{\partial L}{\partial F_k(x)} \approx$  residual/errors for  $n$ m squared-loss functions.

regression  $\frac{\partial L}{\partial F_k(x)} = 2(y_i - F_k(x_i))$

Binary Classification task  $\rightarrow$  Using log-loss

$\frac{\partial L}{\partial F_k(x)}$   $F_x$

$L = y_i \log(p_i) + (1 - y_i) \cdot \log(1 - p_i) \rightarrow$  binary log loss

$p_i = \frac{1}{1 + e^{-\hat{y}_i}} \rightarrow$  logistic function

$\frac{\partial L}{\partial \hat{y}_i} = \frac{y_i - (1 - y_i) \cdot e^{\hat{y}_i}}{1 + e^{\hat{y}_i}}$

After multiplying by  $\frac{e^{-\hat{y}_i}}{e^{-\hat{y}_i}}$



$$L' = \frac{y_i \cdot e^{-\hat{y}_i} + y_i - 1}{1 + e^{-\hat{y}_i}} = \frac{y_i \cdot (1 + e^{-\hat{y}_i})}{(1 + e^{-\hat{y}_i})} - \frac{1}{1 + e^{-\hat{y}_i}}$$

$$= y_i - p_i$$

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

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

$y_i = 0 \text{ or } 1$

$P(x_i \in 1)$

pseudo residual =  $(p_i - y_i)$   $\approx$  error/residual.

In case of SQ loss  $\leftarrow y_i - F_k(x_i)$

$\downarrow$   
 (actual value)  
 $y_i$

$\downarrow$   
 (estimated value)  
 $\hat{y}_i$

$$y_i - \hat{y}_i$$

~~In case of log loss~~

Q pseudo-residual  $\approx$  residual/error always

For SQ-loss  $\rightarrow$  Yes

For log-loss  $\rightarrow$  Yes

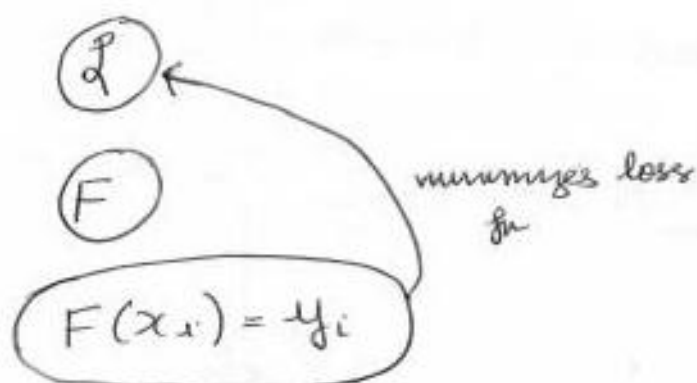
any-loss  $\rightarrow$  No.

It's not always true that pseudo-residual can be interpreted as error/residual.

Q why are we using pseudo-residual if (any loss  $\rightarrow$  No.)  
from previous question but can't we use residual only?

Sol  $\rightarrow$  Proof:

$$D \{x_i, y_i\}$$



In boosting we build a sequence of models

$$F_0(x) = \arg \min \sum_{i=1}^n L(y_i, \delta)$$

$L \leftarrow$  SQ-loss

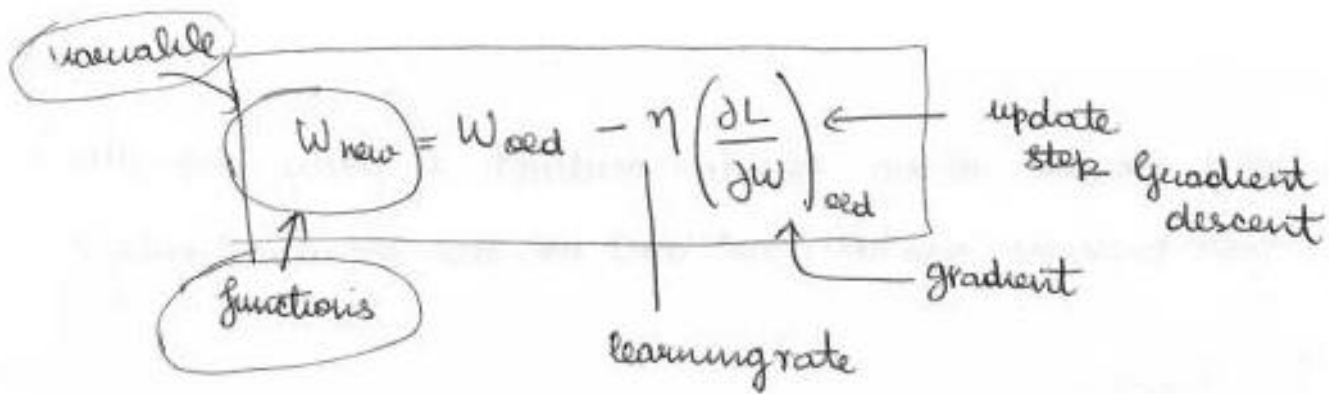
$$= \min_{\delta} \frac{1}{n} \sum_{i=1}^n (y_i - \delta)^2$$

In case regression  $\delta$  would be average of  $y_i$

$H$  = set of all base learners

$$F_m(x) = F_{m-1}(x) + \arg \min_{h_m \in H} \sum_{i=1}^n L(y_i, \underbrace{F_{m-1}(x_i)}_{\text{old}} + h_m(x_i))$$

Loss between  $y_i$  and  $\hat{y}_i$



### Steepest descent

$$W_{\text{new}} = W_{\text{old}} - \text{learnrate} \left( \frac{\partial L}{\partial W} \right)$$

$$F_m(x) = F_{m-1}(x) - \delta_m \sum_{i=1}^n \nabla F_{m-1} L(y_i, F_{m-1}(x_i))_i$$

Annotations for the second equation:

- $F_m(x)$  and  $F_{m-1}(x)$  are indicated by downward arrows.
- $\nabla F_{m-1} L(y_i, F_{m-1}(x_i))_i$  is indicated by a downward arrow to  $\left( \frac{\partial L}{\partial F_{m-1}(x)} \right)_{F_{m-1}(x)}$ .
- $-\nabla F_{m-1}$  is indicated by a downward arrow to  $\parallel$ .
- $\parallel$  is indicated by a downward arrow to **pseudo residual** (boxed).

### Gradient Boosting

$\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss  $J_n L(y, F(x))$ , number of iterations  $M$

1.) Initialize model with a constant value

$$F_0(x) = \arg \min_{\delta} \sum_{i=1}^n L(y_i, \delta)$$

2.) Find  $m=1$  to  $M$ :

1.) compute so-called pseudo-residuals

$$r_{mi} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i=1, \dots, n$$

- 2.) Fit a base learner (eg tree)  $h_m(x)$  to pseudo-residuals, i.e. learn it using the training set  $\{(x_i, r_m)\}_{i=1}^n$

$$\delta_m = \underset{\delta}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \delta h_m(x_i))$$

- 4.) Update the model:

$$F_m(x) = F_{m-1}(x) + \delta_m h_m(x)$$

- 3.) Output  $F_m(x)$ .

In Gradient Boost Algorithm  $\rightarrow$  we need high bias and low variance models that means in Decision Tree we need shallow decision tree as base learners.

Decision tree of depth 1 is called decision stump.

### Regularization By Shrinkage

$$F_m(x) = f_0(x) + \sum \delta_m h_m(x)$$

$m$ : # of base models

# base-models  $\uparrow \Rightarrow$  overfitting  $\uparrow \rightarrow$  Var  $\uparrow$

$m \uparrow \rightarrow$  bias  $\downarrow$

In case of Gradient Boosting  $\rightarrow$  Concept is called Shrinkage:

# base model is hyper-parameter

cv.  $\leftarrow$

$$F_m(x) = \overbrace{F_{m-1}(x)} + \underbrace{v \delta_m}_{\substack{\downarrow \\ \text{learning rate}}} h_m(x), \quad 0 \leq v \leq 1$$

$$V \approx 0.1$$

hyperplane  $\rightarrow CV$

If  $V$  is small less weightage is given to the model

$V \uparrow$  chances of overfitting and variance increases

2 hyperparameters  $\left\{ \begin{array}{l} M \# \text{ base models} \\ V \rightarrow \text{shrinkage coefficient} \end{array} \right.$

$CV \rightarrow \text{gridsearch}(M, V) \rightarrow$  It will reduce the overfitting  $\uparrow$  best value of  $M, V$  should be selected using grid search

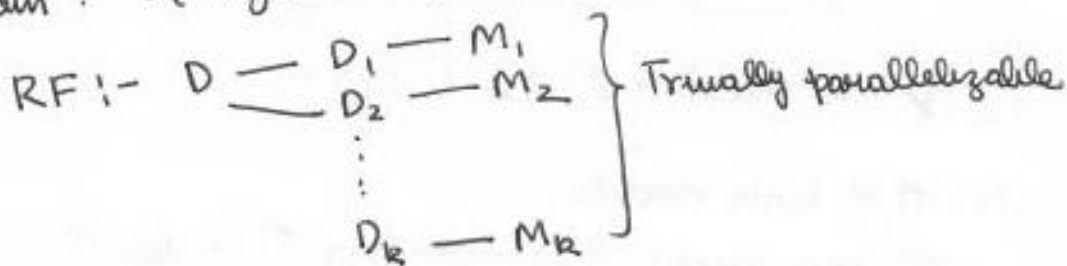
shrinkage

# Base models.

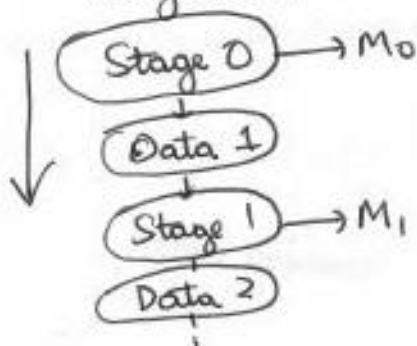
### Train & Run time complexity GB

Train :  $O(n \lg n d * M)$

$M$  : # base learners.



GBDT :- not easy to parallelize  $\rightarrow$  because it is a serial algorithm. (step by step process).



In general GBDT take more time to train than Random forest.

Runtime:

GBDT:

$$O(\underbrace{\text{depth} * M}_{\substack{\text{low latency application} \\ \uparrow \\ \text{small in GB}}})$$

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

$h_m(x)$

$$\text{Space: } O(\underbrace{\text{store each tree} + \delta_m}_{\substack{\text{if...else} \\ \downarrow \\ M}})$$

Suppose GBDT depth = 3  
M = 200 → easy → just if else

GBDT: is used by lot of Internet Companies

⇒ XG Boost: GBDT + row sampling + Column Sampling

~~XXXX data~~

AdaBoost: Geometric Intuition

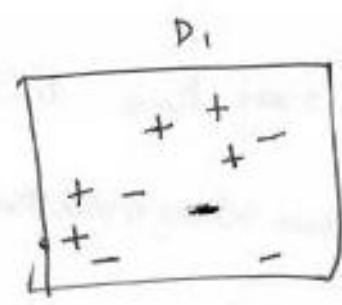
Used in computer vision or Image processing / face detection

Adaptive Boosting → Build learners sequentially

1) samples misclassified by the previous learner

weighted more in subsequent learners.

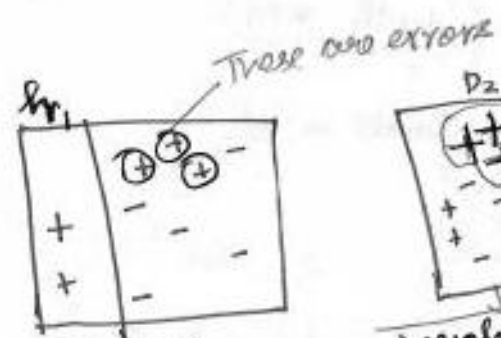
NOTE  
Decision stump is  
either parallel to  
x-axis or y-axis



D.T with depth=1  
⇓  
decision stump

First round:

$\alpha_1 h_1$



$\epsilon_1 = 0.30$   
weight  $\rightarrow \alpha_1 = 0.42$



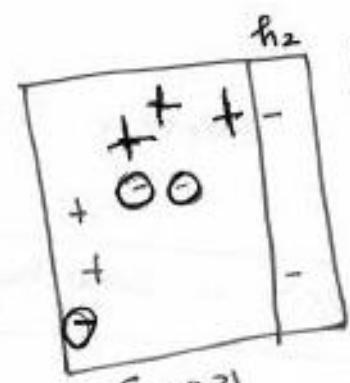
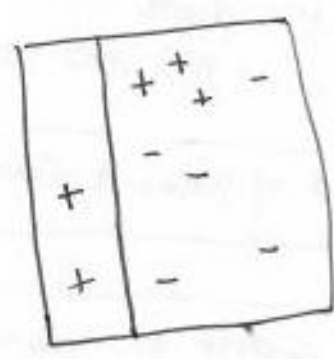
weight on  
wrong point  
increased

weighted models.  
↓  
upsample these  
points

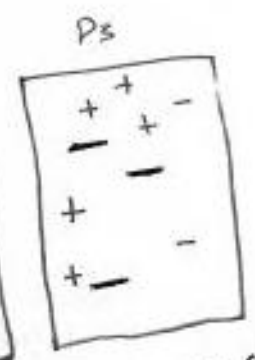
The second round

$\alpha_2 h_2$

$h_1(x)$

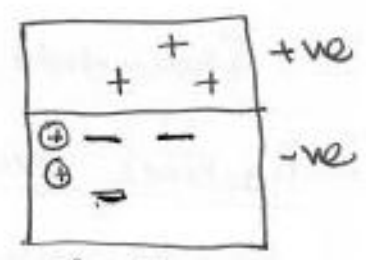
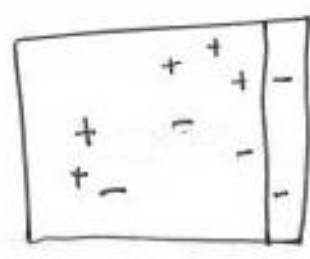
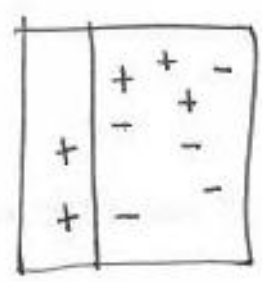


$\epsilon_2 = 0.21$   
 $\alpha_2 = 0.65$



weight of (wrong)  
-ve classified  
points increased

The third round



$\epsilon_3 = 0.14$   
 $\alpha_3 = 0.92$



### The final classifier

$$F_3(x) = \alpha_1 h_1 + \alpha_2 h_2 + \alpha_3 h_3$$

$$H_{\text{final}} = \text{sign} \left( \alpha_1 = 0.42 \begin{array}{|c|} \hline \square \\ \hline \end{array} h_1 + \alpha_2 = 0.65 \begin{array}{|c|} \hline \square \\ \hline \end{array} h_2 + \alpha_3 = 0.92 \begin{array}{|c|} \hline \square \\ \hline \end{array} h_3 \right)$$

		+	+	-
+	-	-	-	-
+	-			

→ In case of GBDT  $\alpha$  are equivalent to  $\gamma$

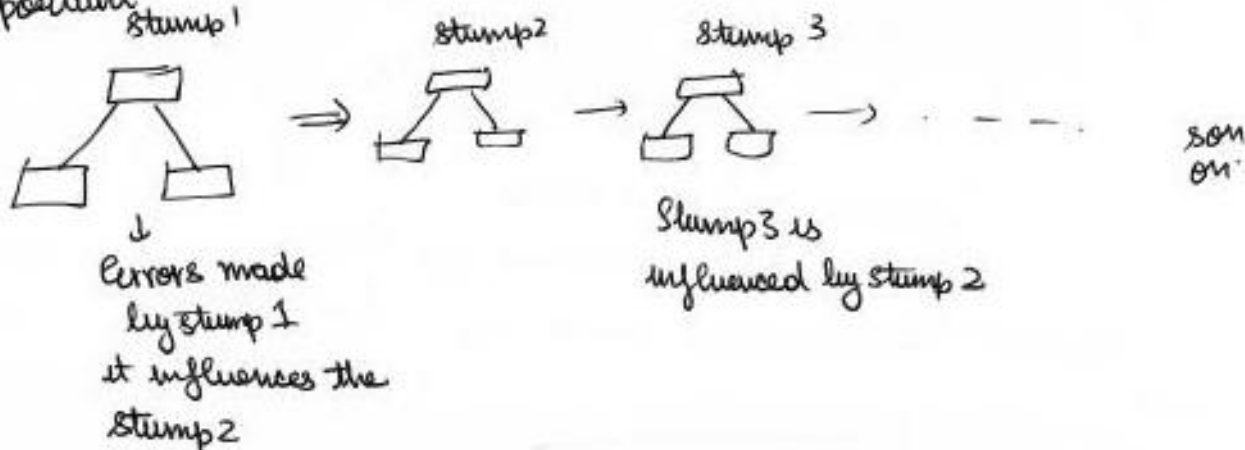
In case of GBDT: we compute pseudo-residuals.

AdaBoost → Forest of Trees made with AdaBoost, the trees are usually just a node and two leaves

Tree with just one node and 2 leaves are called stumps.

Stumps are weak learners that means it has high bias.

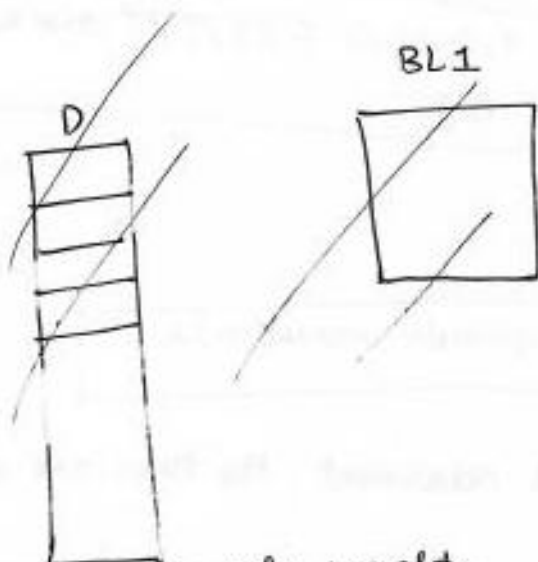
Note → In a Forest of Stumps made with AdaBoost, order is important



### 3 important things about AdaBoost

- 3 important things about AdaBoost:
- 1) It combines a lot of "weak learners" to make classifications.  
The weak learners are almost always stumps.
  - 2) Some stumps get more say in the classification than others.
  - 3) Each stump is made by taking the previous stump's mistakes into account.

Step 1



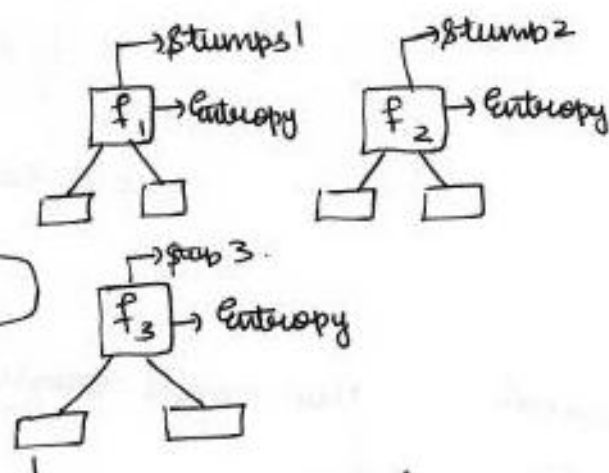
Step 1 → Assign sample weights.

$$w = \frac{1}{n} = \frac{1}{7}$$

[illegible]

Step 1 Sample weight =  $\frac{1}{n} = \frac{1}{7}$

	$f_1$	$f_2$	$f_3$	O/P	Sample weight
				Yes	$\frac{1}{7}$
				No	$\frac{1}{7}$
Incorrect				Yes	$\frac{1}{7}$
				No	$\frac{1}{7}$
					$\frac{1}{7}$
					$\frac{1}{7}$
					$\frac{1}{7}$
					$\frac{1}{7}$



we select the minimum entropy among the 3 features

Step 2 → We calculate the total error of the selected stump  
 Suppose correctly classified = 4  
 incorrect = 1

Step 3 → We calculate total error =  $\frac{1}{7}$

$$\text{Performance of Stump} = \frac{1}{2} \log_e \left( \frac{1 - TE}{TE} \right)$$

$$= \frac{1}{2} \log_e \left[ \frac{1 - \frac{1}{7}}{\frac{1}{7}} \right]$$

$$= \frac{1}{2} \log_e [6]$$

$$= 0.896$$

We just need to increase weight of incorrectly classified stump  
 on the other hand decrease weight of correctly classified stump  
 using performance of stump

or incorrect:

$$\text{New sample weight} = \text{weight} \times e^{\text{performance say}}$$

$$= \frac{1}{7} \times e^{0.895}$$

$$= 0.349$$

For correct:  $\text{New sample weight} = \text{weight} \times e^{-\text{performance say}}$

$$= \frac{1}{7} \times e^{-0.895}$$

$$= 0.05$$

$f_1$	$f_2$	$f_3$	O/P	Sample weight	Updated weight
			Yes	$\frac{1}{7}$	0.05
			No	$\frac{1}{7}$	0.05
			Yes	$\frac{1}{7}$	0.349
			No	$\frac{1}{7}$	0.05
				$\frac{1}{7}$	0.05
				$\frac{1}{7}$	0.05
				$\frac{1}{7}$	0.05

For incorrect updated weight

Step-4

NOTE: Sum of updated weight is not equal to 1 so we need to normalize it.

$f_1$	$f_2$	$f_3$	O/p	Normalized weight
_____			Y	0.07
_____			N	0.07
_____			Y	0.51
_____			N	0.07
_____				0.07
_____				0.07
_____				0.07

Step 5 → We will create a new decision stamp.

$f_1$	$f_2$	$f_3$	O/p	Normalized weight	Cumulative freq
			Yes	0.07	0 - 0.07
			No	0.07	0.07 - 0.14
			Yes	0.51	0.14 - 0.65
			No	0.4907	0.65 - 0.72
				0.07	0.72 - 0.79
				0.07	0.79 - 0.86
				0.07	0.86 - 0.93

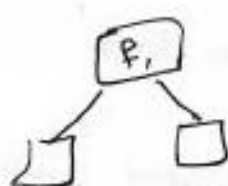
We will use the sample weights to calculate weighted sum indices to determine which variable should split the next stamp.

$f_1$	$f_2$	$f_3$	O/p	Normal	Sum
-------	-------	-------	-----	--------	-----

if number (randomly selected) is between 0.07 - 0.14 we put this sample into new collection of sample

if number is between 0.08 - 0.93 we put corresponding sample into new collection of sample

$f_1$	$f_2$	$f_3$	O/P	cumulative freq	Randomly select
			No	0.07 - 0.14	0.09
			Yes	0.85 - 0.93	0.84



Same process is repeated

### Gradient Boost Part 1:

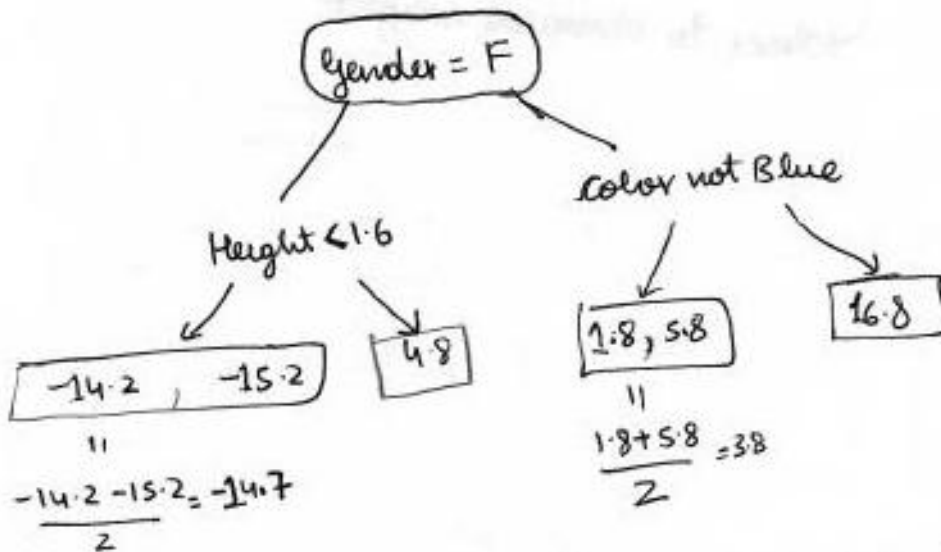
Height	Favorite Color	Gender	Weight
1.6	B	M	88
1.6	G	F	76
1.5	B	F	56
1.8	R	M	73
1.5	G	M	77
1.4	B	F	57

Step I. First Average Height = 71.2

Step 2 → build a tree using error from previous tree

Residual = (Observed - 71.2)

Height	FC	Gender	Weight	Residual
1.6	B	M	88	16.8
1.6	G	F	76	4.8
1.5	B	F	56	-15.2
1.8	R	M	73	1.8
1.5	G	M	77	5.8
1.4	B	F	57	-14.2



so the Predicted weight =  $71.2 + 16.8 = 88$  → It leads to low bias but very high variance

~~Height~~ Gradient Boost deals with this problem by using a learning Rate to scale the contribution from the new tree.

The learning Rate is a value between 0 and 1.

Now the predicted weight =  $71.2 + (0.1 \times 16.8) = 72.9$

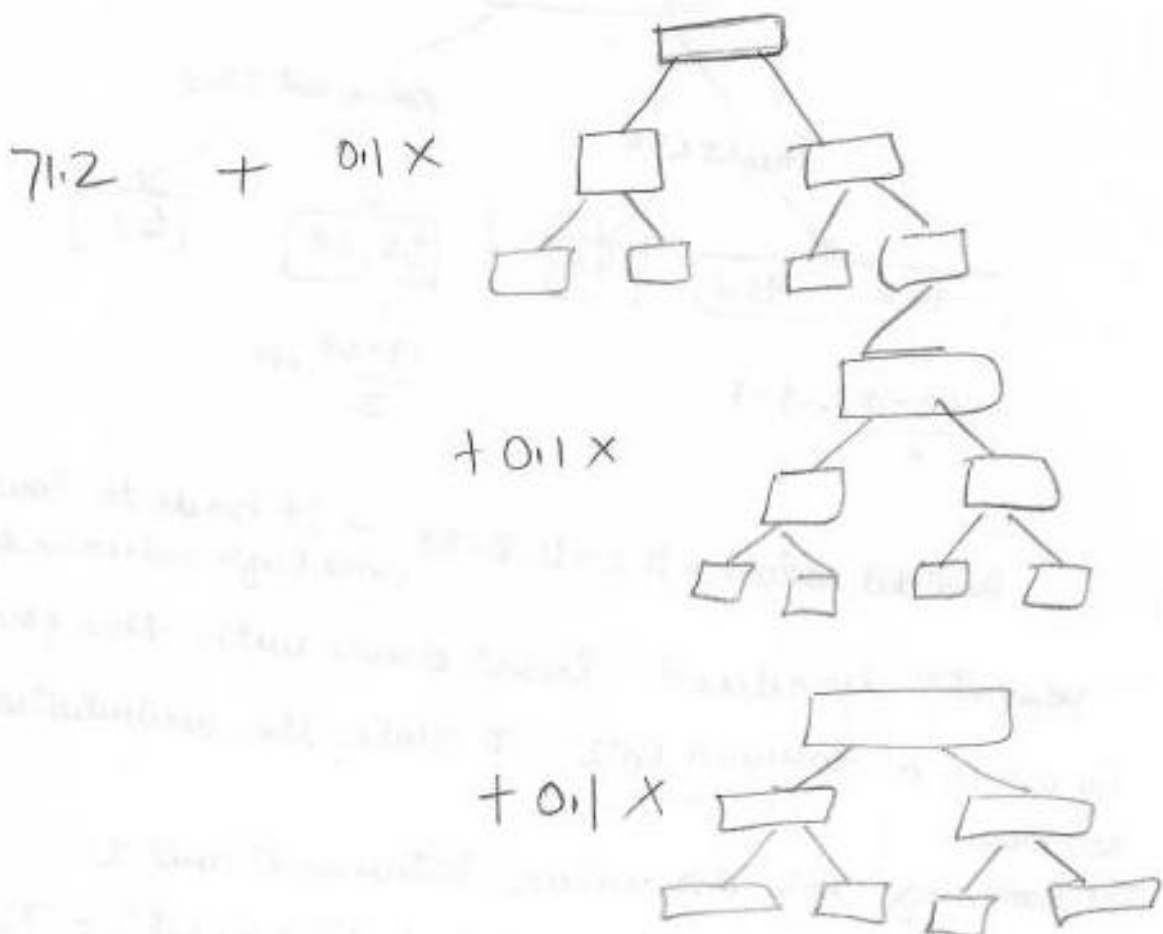


$$\text{Residual} = (\text{Weight} - (w_0 \times \text{Average Residual} + \text{Average}))$$

H	F.C	G	Weight	Residual 1	Residual 2	Residual
			88	16.8	15.1	13.6
				4.8	4.3	3.9
				-15.2	-13.7	-12.4
				1.8	1.4	1.1
				5.8	5.4	5.1
				-14.2	-12.7	-11.4

Residual getting smaller

For first data pt  $\rightarrow 71.2 + (0.1 \times 16.8) + (0.1 \times 15.1)$   
 $= 74.4$   
 $\rightarrow$  closer to observed weight



## Regression details for gradient boost.

Height	F. C	Gender	Weight(kg)
1.6	Blue	M	88
1.6	Green	F	76
1.5	Blue	F	56

←  $x_i$  (Height, F.C., gender)  
 ←  $y_i$  weight

Data  $\{(x_i, y_i)\}_{i=1}^n$  and a differentiable loss fn  $L(y_i, \hat{F}(x))$

loss function for gradient boost  $= \frac{1}{2} (\text{Observed} - \text{Predicted})^2$

$$\frac{d}{d \text{ Predicted}} = \frac{2}{2} (\text{Observed} - \text{Predicted}) \times -1$$

$$= -(\text{Observed} - \text{Predicted}).$$

Step 1 →  $\underline{F_0(x)} = \underset{\gamma}{\text{argmin}} \sum_{i=1}^n L(y_i, \gamma)$

$\downarrow$  loss fn       $\downarrow$  predicted value

$$\frac{1}{2} (88 - \text{Predicted})^2 + \frac{1}{2} (76 - \text{Predicted})^2 + \frac{1}{2} (56 - \text{Predicted})^2 = 0$$

$$3 \text{ Predicted} = 88 + 76 + 56$$

$$F_0(x) \text{ Predicted} = \frac{88 + 76 + 56}{3}$$

$$F_0(x) = 73.3$$

leaf is predicted

Step 2: for  $m=1$  to  $M$ : Calculate residual

A.) Compute  $r_{i,m} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$

$$= (\text{observed} - F_0(x))$$

$$\begin{aligned} r_{1,1} &= (\text{observed} - 73.3) = (88 - 73.3) = 14.7 \\ &= (76 - 73.3) = 2.7 \\ &= (56 - 73.3) = -17.3 \end{aligned}$$

Thus  $r_{i,m}$  are called pseudo residuals

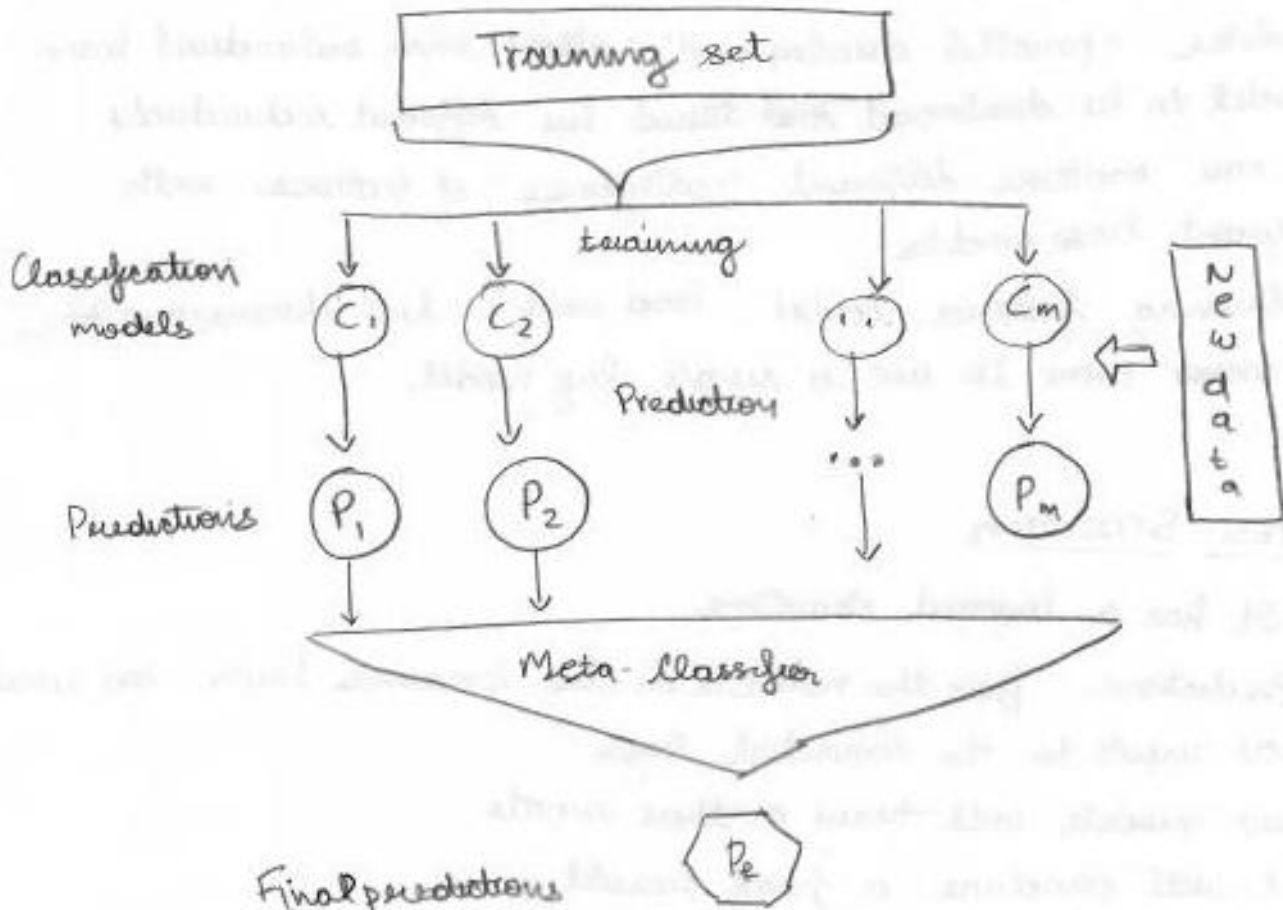
B.) Fit a regression tree to the  $r_{i,m}$  values and create terminal regions  $R_{jm}$  for  $j=1 \dots J_m$ .

C.) For  $j=1 \dots J_m$  compute  $\delta_{jm} = \underset{\delta}{\text{argmin}} \sum_{x_i \in R_{j,m}} L(y_i, F_{m-1}(x_i) + \delta)$

Here we are taking previous prediction.

# Stacking Classifier

## Overview



Meta Classifier  $\rightarrow$  Model can be logistic Regression, Decision Tree

$$y_q = h'(h_1(x_q), h_2(x_q), \dots, h_r(x_q))$$

meta learner

$$\begin{cases} h_1(x) = \text{clf1} \rightarrow \text{NN} \\ h_2(x) = \text{clf2} \rightarrow \text{RF} \\ h_3(x) = \text{clf3} \rightarrow \text{GNB} \end{cases}$$

meta learner  $\leftarrow h'(x) = \text{RF Logistic Regress}$

## Benefits of Stacking

- Stacking increases the diversity of the algorithm and models used
- Stacking can decrease bias - rather than winner takes all, combine datasets to decrease bias.
- enables "parallel development": allow each individual base model to be developed and tuned by different individuals.
- we can capture different "categories" of features with different base models.
- Combining features could lead very high dimensionality, if we were to use a single big model.

## NOTES STACKING

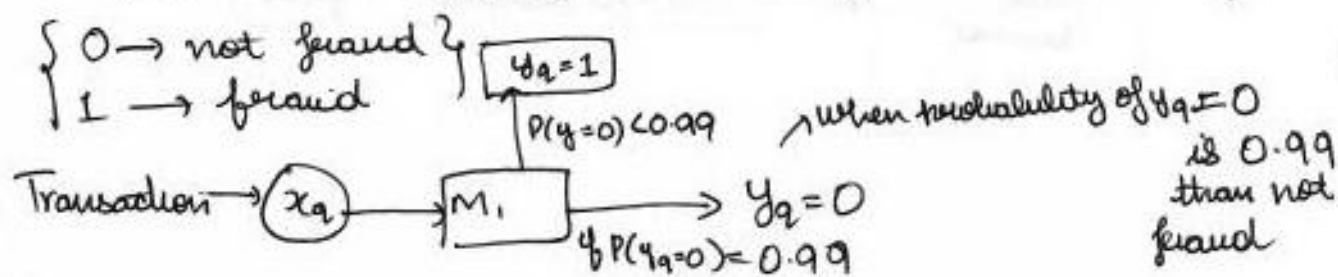
- It has a layered structure
- Predictions from the models in the previous layer are used as inputs to the sequential layer.
- new models will train on these inputs
- It will produce a final result

## Cascading Models:

- eg: predicts credit card transaction is fraudulent or not.

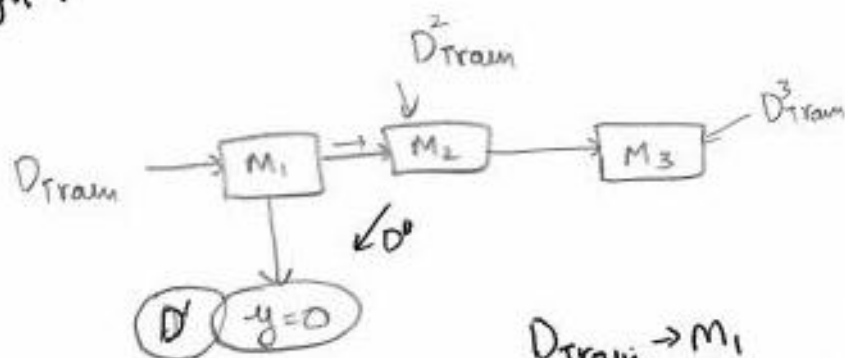
Transaction  $\rightarrow x_q$  (vector) {consisting  $\rightarrow$  location, credit card}  
Credit card company have lot of domain knowledge

$\begin{cases} 0 \rightarrow \text{not fraud} \\ 1 \rightarrow \text{fraud} \end{cases}$



## Cascading model:

Typically used when the cost of making a mistake is high.

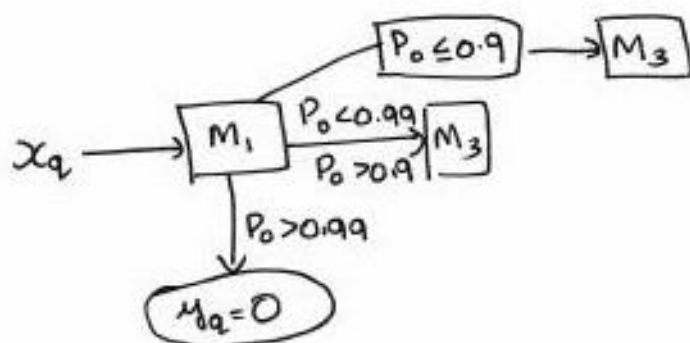


$$D_{\text{Train}} \rightarrow M_1$$

$$D^2_{\text{Train}} \rightarrow M_2$$

$$D^2_{\text{Train}} = D_{\text{Train}} - D'$$

$$D^3_{\text{Train}} = D^2_{\text{Train}} - D''$$



$$P_0 = P(y_q = 0)$$

$$P_1 = P(y_q = 1)$$

$$\text{if } P_0 > 0.99$$

$$y_q = 0$$

$$\text{else } \{ \text{if } P_0 > 0.9$$

$$\text{use } M_2$$

$$\text{else use } M_3.$$

$$y$$

Cascade being built for fraud detection.