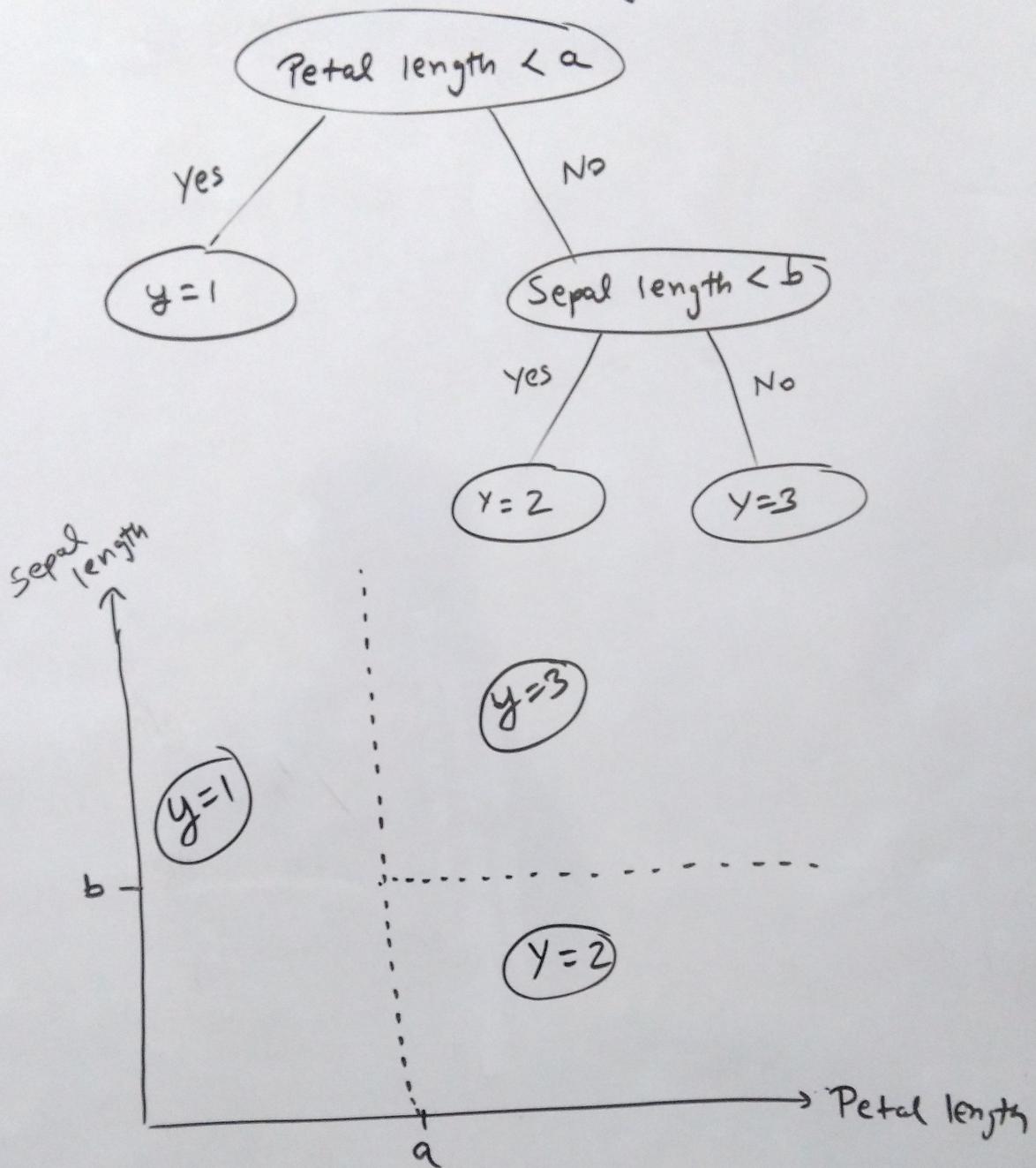


Decision tree

→ g_t is set of axis parallel Hyperplane.



Decision tree

Play Tennis Example

Outlook	Temp.	Humidity	Windy	Play Tennis
Sunny	Hot	High	False	No
Overcast	Mild	Normal	True	Yes
Rainy	Cool	High	False	No

Entropy $H(Y)$: Randomness in target variable.

$$H(Y) = - \sum_{i=1}^k P(y_i) \log(P(y_i))$$

here $k=2$ because classification is Binary

$$P(y_i) = \text{Probability of } y_i = \frac{y_i}{\sum y_i}$$

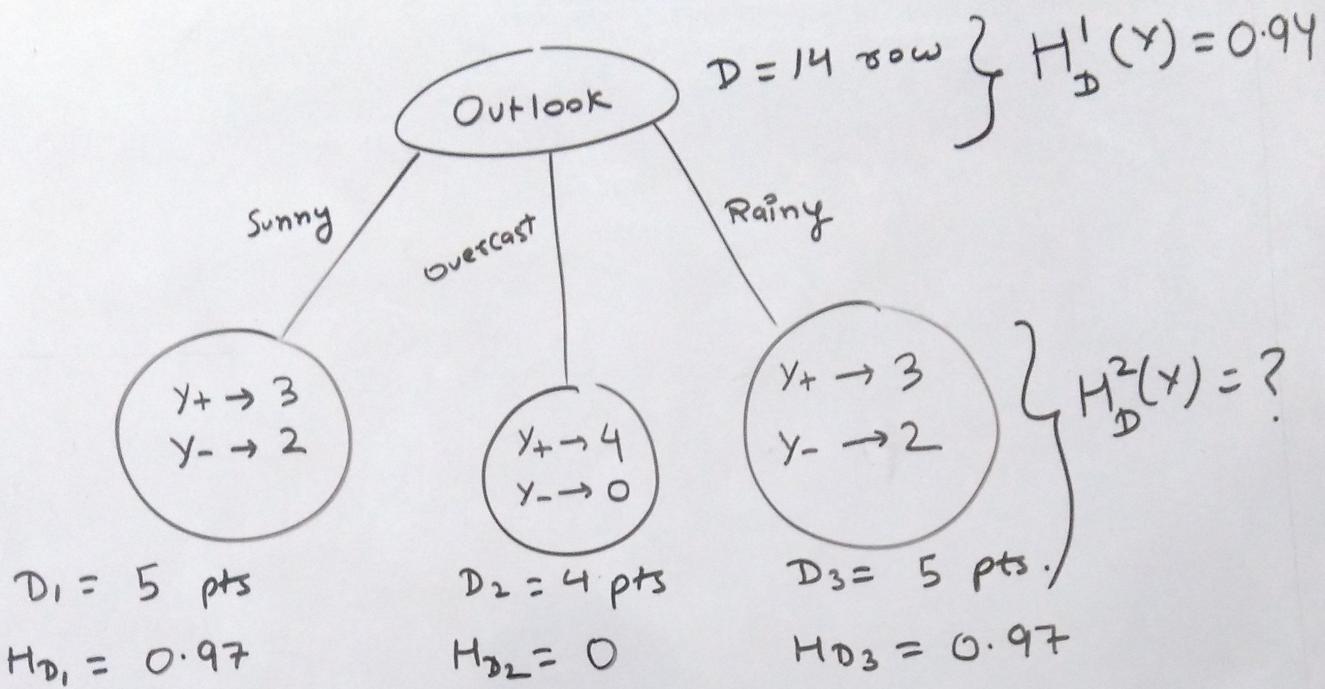
$$\text{if } \begin{cases} Y_+ = 99\% \\ Y_- = 1\% \end{cases} \quad H(Y) = 0.08 \quad \left[\begin{matrix} \text{low} \\ \text{randomness} \end{matrix} \right]$$

$$\text{if } \begin{cases} Y_+ = 50\% \\ Y_- = 50\% \end{cases} \quad H(Y) = 1 \quad \left[\begin{matrix} \text{low} \\ \text{randomness} \end{matrix} \right]$$

$$\text{if } \begin{cases} Y_+ = 0\% \\ Y_- = 100\% \end{cases} \quad H(Y) = 0 \quad \left[\begin{matrix} \text{No randomness} \\ \text{at all} \end{matrix} \right]$$

Binary $H(Y) = - \frac{C_1}{C_1+C_2} \cdot \log \left(\frac{C_1}{C_1+C_2} \right) - \frac{C_2}{C_1+C_2} \log \left(\frac{C_2}{C_1+C_2} \right)$

Decision tree



Entropy at level 2

$$H_D^2(Y) = \frac{D_1}{D} \cdot H_{D_1} + \frac{D_2}{D} \cdot H_{D_2} + \frac{D_3}{D} \cdot H_{D_3}$$

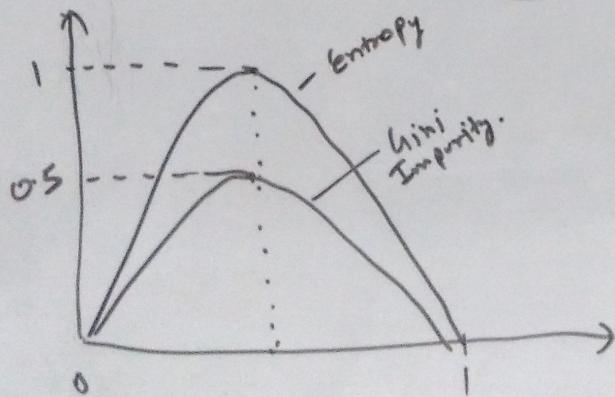
Information Gain

$$IG(Y, \text{Outlook}) = H_D^1(Y) - H_D^2(Y)$$

(Want to maximize this)

Gini Impurity (Same as Entropy) - Easy Computation

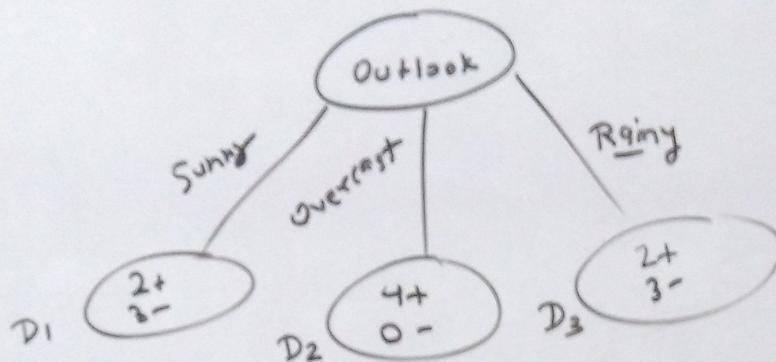
$$I_G(Y) = 1 - \sum_{i=1}^K (P(y_i))^2$$



$$\begin{aligned} & \Rightarrow P(Y_+) = 0.5 \\ & P(Y_-) = 0.5 \\ & I_G(Y) = 0.5 \\ \\ & \Rightarrow P(Y_+) = 1 \\ & P(Y_-) = 0 \\ & I_G = 0 \end{aligned}$$

Decision tree

Construction :



(1) { find $IG(Y, \text{outlook})$
 $IG(Y, \text{Temp.})$
 $IG(Y, \text{Humidity})$
 $IG(Y, \text{Windy})$
 choose Node with Max IG

(2) Break Nodes till:

- we get Pure Node (all +ve or all -ve)
- very few pts in Node (lead to catch Noise and Overfit)
- If we are too deep. (Set max_depth)

Decision tree

Construction :

Handling and Splitting of Numerical feature.

feature (f_i)	y
2.2	1
2.6	1
3.5	0
3.8	0
4.6	1
5.3	0

- Sort Numerical feature in asc. Order.

$$\begin{aligned} f < 2.2 &\rightarrow \text{IG} \\ f < 2.6 &\rightarrow \text{IG} \\ f < 3.5 &\rightarrow \text{IG} \\ f < 3.8 &\rightarrow \text{IG} \\ f < 4.6 &\rightarrow \text{IG} \\ f < 5.3 &\rightarrow \text{IG} \end{aligned}$$

Choose condition which gives Max. Information Gain (I_h)

$x - \quad x - \quad y - \quad x -$

Handling Category feature with many levels
like Zipcode \rightarrow convert into numerical.

Pincode	$y = \{0, 1\}$	$P(Y=1 P_i)$
P_1	0	0.253
P_2	1	0.61
P_3	0	0.001
⋮	⋮	0.65

Decision tree

Complexity: Train time complexity

for breaking one node [det if all feature is numerical.
 f_1, f_2, \dots, f_d

$$f_1 \quad f_2 \quad \dots \quad f_d \\ \downarrow \text{sort} \quad \downarrow \text{sort} \quad \dots \quad \downarrow \\ (n \log n) \quad (n \log n) \quad \dots \quad (n \log n) = O(n \log n d)$$

Test time

$$= O(\text{Max-depth})$$

↓
 Global NP hard
 indeterministic
 problem

Depth \downarrow ; Underfit; interpretability \uparrow

Depth \uparrow ; Overfit; interpretability \downarrow

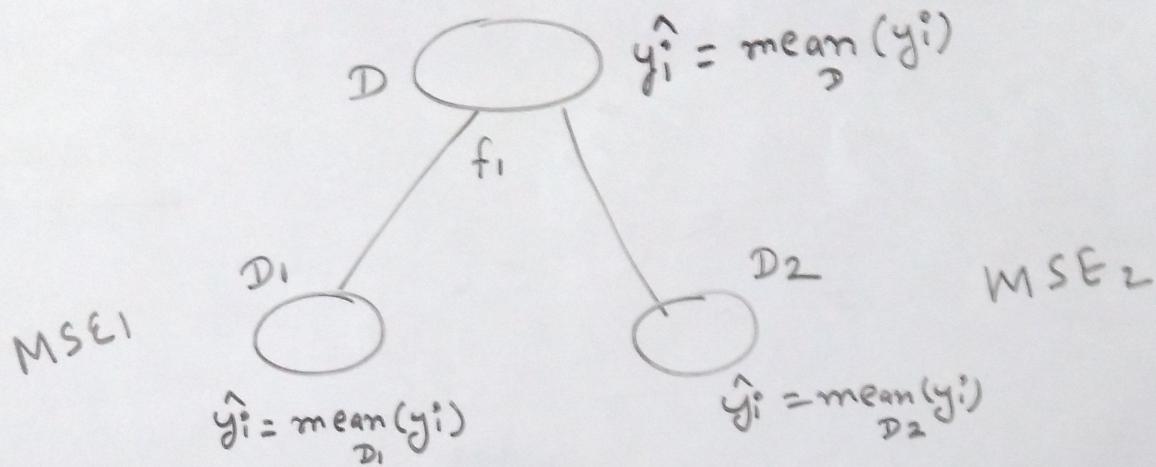
Hyperparameter :- Depth of tree

- Min. point to break Node

Decision tree

Regression

Classification \rightarrow diff. of entropy / IG
 Regression \rightarrow diff. of MSE / MAE



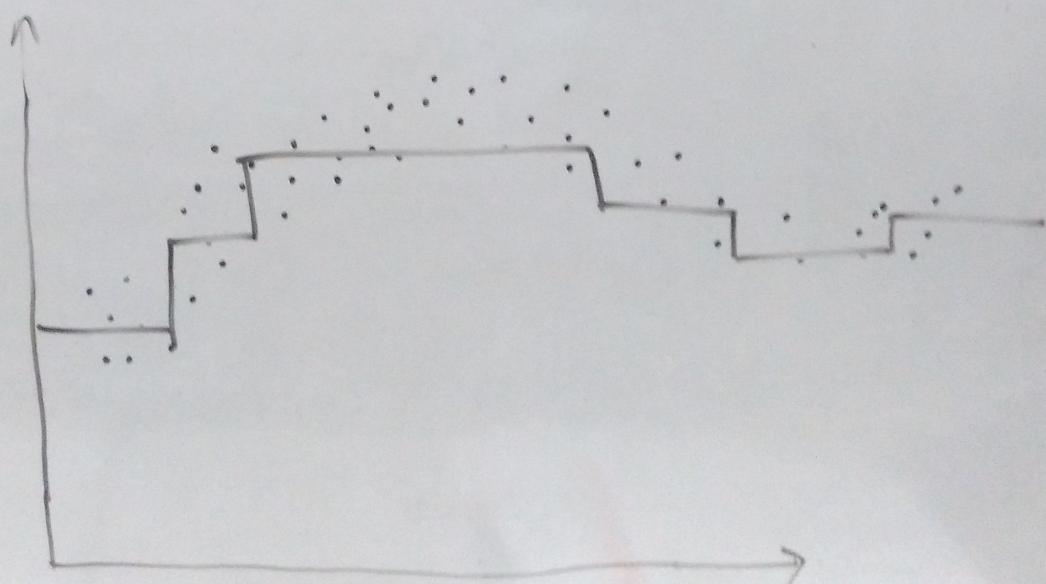
$$\underline{\text{Level 1}} \quad \boxed{MSE^1 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

$$\underline{\text{Level 2}} \quad MSE^2 = \left(\frac{D_1}{D} \cdot MSE_1 + \frac{D_2}{D} \cdot MSE_2 \right)$$

or

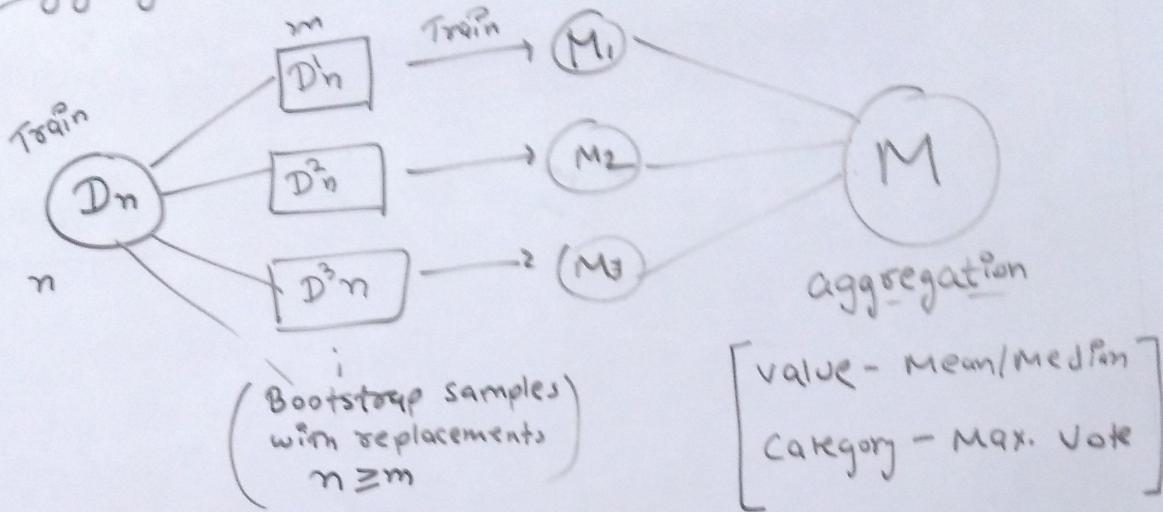
$$\boxed{MSE = w_1 \cdot MSE_{D_1} + w_2 \cdot MSE_{D_2}}$$

Chose feature to Maximize : $(MSE^1 - MSE^2)$



Decision tree - Ensembles

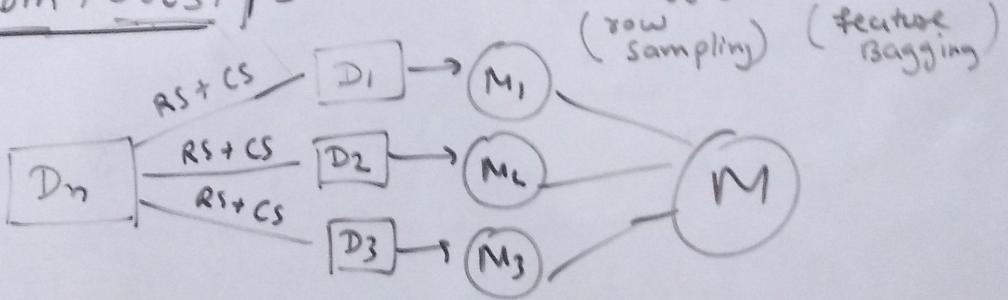
Bagging (Bootstrapped Aggregation)



$$\text{Model Error} = \text{Bias}^2 + \text{Variance} + \text{constant}$$

(Help in Reducing Variance)

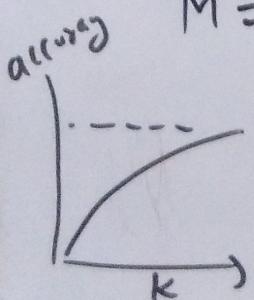
Random Forest = Decision tree + Bagging + Column Sampling



Each Model M_1, M_2, M_k have high variance.

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

$K \uparrow$; Variance \downarrow
 $K \downarrow$; Variance \uparrow



$O(n \log n \cdot K)$

$O(K \cdot \text{depth})$

space:
 $O(\#T \times K)$

$$D_n \xrightarrow[\text{CS (50\%)}]{\text{RS (20\%)}} D'_n$$

CSR \downarrow ; Variance \downarrow
RSR

Boosting

$D_{train} \rightarrow M_0$ (high bias model)

$$\{x_i, y_i\}_{i=1}^n$$

$$y = h_0(x)$$

$$\text{error} = y_i - h_0(x) \quad [\text{Residual error}]$$

Stage 1:

$$\{x_i, \text{error}_i\}_{i=1}^n \rightarrow M_1$$

$$y = h_1(x)$$

Model at the end of stage 1 = $F_1(x)$

$$F_1(x) = \alpha_0 h_0(x) + \alpha_1 h_1(x)$$

$$\text{error} = y_i - F_1(x)$$

Stage 2:

$$\{x_i, \text{error}_i\}_{i=1}^n \rightarrow M_2$$

$$y = h_2(x)$$

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

$$F_2(x) = F_1(x) + \alpha_2 h_2(x)$$

Stage K

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

Boosting

Residual, loss & Gradient

$$F_K(x) = \sum_{i=0}^K \alpha_i h_i(x)$$

at the end of K stage/iteration

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

if Regression Problem:

$$\text{Loss: } L(y_i, F_K(x)) = (y_i - F_K(x))^2$$

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

$$\text{let } F_K(x) = z$$

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

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

negative gradient \approx residual

here $-\frac{\partial L}{\partial z}$ is known as pseudo residual.

[it have ability to penalize any differentiable loss function.]

$$\text{error}_i = y_i - F_K(x) \Rightarrow \underbrace{2(y_i - F_K(x))}_{\substack{\text{residual} \\ \text{replace by}}} \underbrace{2(y_i - F_K(x))}_{\text{pseudo residual.}}$$

{ x_i, error_i }
Pseudo.

Gradient Boosting algorithm

Input: $\{(x_i, y_i)\}_{i=1}^n$;

$L(y, F_k(x))$: A differentiable loss function.

M : No of iteration

1.) Initialize the Model with some constant:

$$F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i - \gamma)$$

here γ is mean(y_i) for regression.

2.) for $m = 1$ to M.

① Compute Pseudo Residuals:

$$g_{im} = \text{error}_i = - \left[\frac{\partial L(y_i, F_{m-1}(x_i))}{\partial F_{m-1}(x_i)} \right]$$

② Fit a Base learner $h_m(x)$ to psuedo residual train using $\{(x_i, g_{im})\}_{i=1}^n$

③ Compute Multiplier γ_m by optimization problem:

$$\gamma_m = \operatorname{argmin}_\gamma \sum L(y_i, F_{m-1}(x_i) + \gamma \cdot h_m(x))$$

④ Update Model

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

3.) Output Final Model.

Gradient Boosting

Regularization :

GBDT are prone to overfitting, so add shrinkage coefficient (γ)

$$F_m(x) = F_{m-1}(x) + \gamma \cdot j_m h_m(x)$$

$0 \leq \gamma \leq 1 \quad \therefore \text{default} = 0.1$

\Rightarrow GBDT Hyperparameter M and γ .

Train : $O(n \log n \cdot d \times M)$

M = No. of Base learner.

Test: $O(\underbrace{\text{depth} \times M}_{\text{Small}} + \underbrace{M}_{\text{add all tree}}) \approx O(\text{depth} \times M)$

Space: $O(\text{each tree} + \underbrace{\text{each } \gamma}_{\text{store M No. of 2}})$.

Adaboost: upsample wrong classified example.
than train again.

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