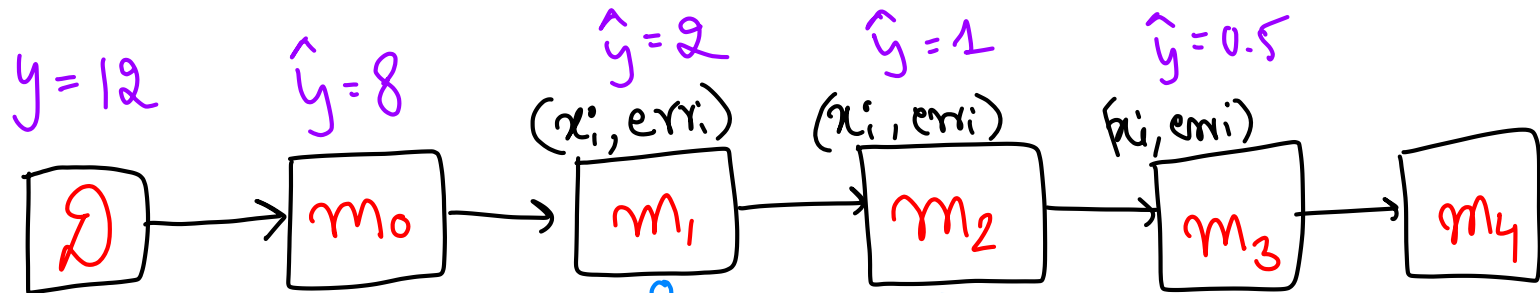


BOOSTING-2

---



$\{x_i, y_i\}$

$$\begin{aligned} \text{error} &= y - \hat{y} \\ &= 12 - 8 \\ &= 4 \end{aligned}$$

$$\begin{aligned} f_0 &= 8 \\ f_1 &= 8 + 2 \end{aligned}$$

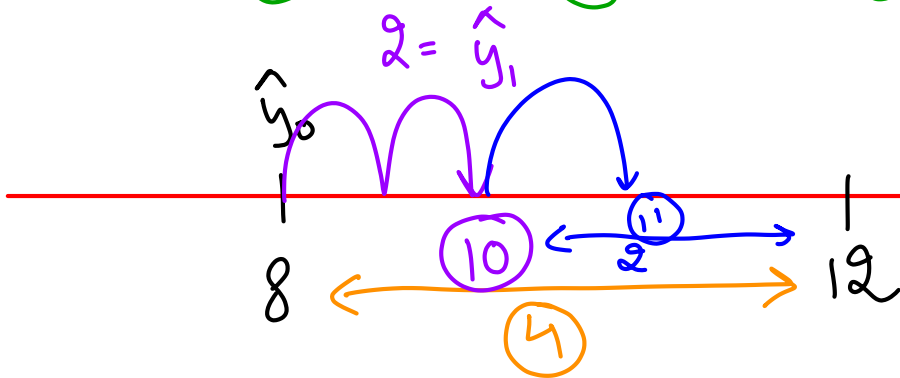
$$\text{error} = 2$$

$$f_2 = f_1 + 1$$




$$\text{error} = 1$$

$$f_3 = f_2 + 0.5$$

$$\text{error} = 0.5$$



$$\begin{aligned} F_0(x) &= 8 \\ f_1(x) &= 8 + 2 \\ f_2(x) &= 8 + 2 + 1 \\ &= 11 \end{aligned}$$

$f_1$	$f_2$	$f_3$	$y$	$m_0$ $\hat{y}_0$	$err_0$	$m_1 \hat{y}_1$ $(f_1, f_2, f_3, err_0)$	$err_1$	$m_2 \hat{y}_2$ $(f_1, f_2, f_3, err_1)$	$err_2$	$m_3 \hat{y}_3$		
			3	5.5	-2.5	-1.5	-1	-0.5	-0.5			
			4	5.5	-1.5	-0.5	-1	-0.3	-0.7			
			6	5.5	0.5	0.5	0	0	0			
			9	5.5	3.5	2	1.5	1	0.5			

mean  
of  $y$ 's

$$y - \hat{y}_0 - \hat{y}_1$$

$$y - \hat{y}_0 - \hat{y}_1 - \hat{y}_2$$

$$12 - 8 - 2 = 2$$

$$\text{err}_m = y - \left[ \hat{y}_0 + \sum_{i=1}^k \hat{y}_i \right]$$

$$\text{err}_0 = y - \hat{y}_0$$

$$\text{err}_1 = y - \hat{y}_0 - \hat{y}_1$$

$$\text{err}_2 = y - \hat{y}_0 - \hat{y}_1 - \hat{y}_2$$

$$y = \text{err}_m + \left[ \hat{y}_0 + \sum_{i=1}^k \hat{y}_i \right]$$

Boosting

$$\mathcal{D}_{\text{Train}} = x_i$$

Step 1

Create  $m_0$

function value  
after stage 0

$$f_0(x) = \hat{y}_0 = \text{Avg}(x) = h_0(x)$$

predicted  
value at stage 0

output of  $m_0$ .

$$\text{err}_0 = y - \hat{y}_0$$

$$y = \underbrace{\hat{y}_0}_{h_0(x)} + \text{err}_0$$

$$\underline{y} = \underbrace{\hat{y}_i}_{h_i(x)} + \text{err}_i$$

for each datapoint  $\{x_i, y_i, \text{err}_i\}_{i=1}^k$

~~Step 1~~ Create M1  $\rightarrow$  Train on  $\{\underbrace{x_i}_\text{input}, \underbrace{\text{err}_i}_\text{label/Target}\}$

$$F_1(x) = h_0(x) + \gamma_1 h_1(x)$$

step 3

Create  $m_2$

$$\{x_i, \underbrace{err_1^{(i)}}_{\rightarrow y - f_1(x)}\}$$

$$\rightarrow y - f_1(x)$$

$$f_2(x) = h_0(x) + \underbrace{\gamma_1 h_1(x) + \gamma_2 h_2(x)}_{\rightarrow \text{weights}}$$

step 6

$$f_G(x) = h_0(x) + \gamma_1 h_1(x) + \gamma_2 h_2(x) + \dots + \gamma_G h_G(x)$$

$$f_G(x) = h_0(x) + \sum_{i=1}^G \gamma_i h_i(x)$$

# Linear Regression

# Gradient Boosting

$$F_3(x) = \underbrace{h_0(x)}_{\text{const}} + \gamma_1 \underbrace{h_1(x)}_{\text{tree}} + \gamma_2 \underbrace{h_2(x)}_{\text{tree}} + \gamma_3 \underbrace{h_3(x)}_{\text{tree}}$$

learning  $\gamma$ 's is a sequential process.

$$\hat{y} = \underbrace{w_0}_{\text{const}} + w_1 \underbrace{x_1}_{\text{feature}} + w_2 \underbrace{x_2}_{\text{feature}} + w_3 \underbrace{x_3}_{\text{feature}}$$

Learning  $w$ 's is a parallel process

## errors / Residuals

$$F_k(x) = h_0(x) + \gamma_1 h_1(x) + \gamma_2 h_2(x) \dots \gamma_k h_k(x)$$

$$\text{Residual} = y_i - F_k^i(x) \quad k=2$$

12

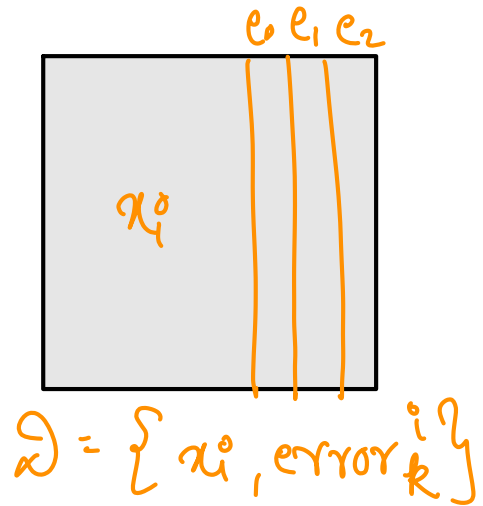
$f_k(x) = 10$

Training for  $(k+1)^{\text{th}}$

$\hookrightarrow m$   
        

$$F_k(x) \Rightarrow \hat{y}_k \Rightarrow \text{Squared loss} = (y - \hat{y}_i)^2$$

Loss function

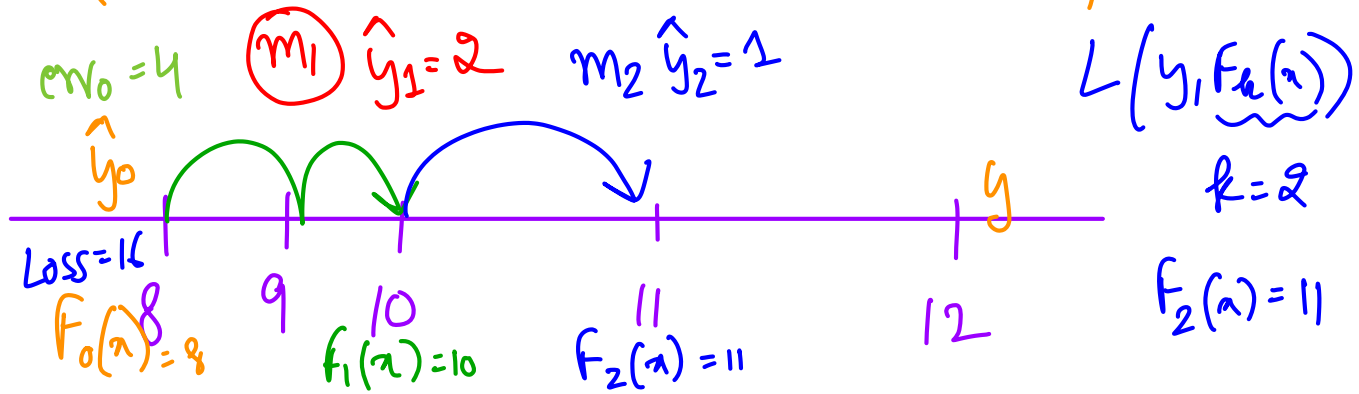




$$L(y, \hat{y}_i) = (y_i - \hat{y}_i)^2$$

$$L(\underset{12}{y}, \underset{10}{F_k(x)}) = (y_i - F_k(x))^2 = 2^2 = 4$$

Let's  $f_k(x) = z_i$  error as function



$$f_k(x) = z_i$$

$$L(y, z_i) = (y - z_i)^2$$

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

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

$$L(y, f_2(x)) = (12 -$$

$$L(y, f_2(x)) = (y - f_2(x))^2$$

$$L(y, f_k(x)) = (y - f_k(x))^2$$

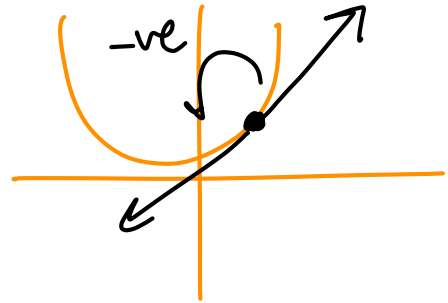
$$\frac{\partial L}{\partial z_i} = -2(y_i - z_i)$$

$$\underbrace{-\frac{\partial L}{\partial z}}_{\text{neg grad.}} = \underbrace{2(y_i - z_i)}_{\text{const}}$$

$$-\frac{\partial L}{\partial z_i} \approx (y_i - z_i)$$

$$-\frac{\partial L}{\partial z_i}$$

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w}$$



$$\underbrace{-\frac{\partial L}{\partial f_k(x_i)}}_{\text{pseudo error}} \approx \underbrace{(y^{(i)} - f_k(x_i))}_{\text{error / Residual}}$$

$L = \underline{\underline{\text{loss}}}$   $\rightarrow$  Boosted trees

$$h_k(x) = \{x^{(i)}, \text{error}_{k-1}^i\}$$

pseudo error

Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:

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

2. For  $m = 1$  to  $M$ :

1. Compute so-called *pseudo-residuals*:

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

2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .

3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

4. Update the model:

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

3. Output  $F_M(x)$ .

Loss function. Sq. Loss.

$$y = 2, 4, 6 \quad (y - \hat{y})^2$$

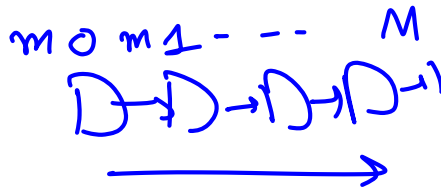
$$y = 2, 4, 6$$

$$L = (2-\gamma)^2 + (4-\gamma)^2 + (6-\gamma)^2$$

$$\frac{\partial L}{\partial \gamma} = 2(2-\gamma) - 2(4-\gamma) - 2(6-\gamma) = 0$$

$$-2(2-\gamma + 4-\gamma + 6-\gamma) = 0$$

$$12 - 3\gamma = 0 \Rightarrow \gamma = 4$$



$$m_1$$

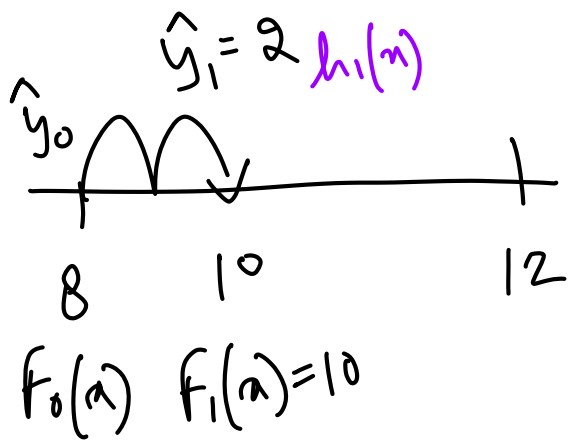
$$\gamma = 1$$

$$L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)) =$$

$$f_1(x) = f_{1-1}(x) + 2 h_1(x)$$

$$(12) = f_1(x) = 8 + 2 \times 2 = 12$$

$$[y_i - F_{m-1}(x_i) - \gamma h_m(x_i)]^2$$



$$\underline{m=1} \quad [y_i - f_{i-1}(x) - \gamma h_1(x)]^2$$

$$f(\gamma) = [y_i - f_0(x) - \gamma h_1(x)]^2$$

$$\frac{\partial}{\partial \gamma} = \underline{2[y_i - f_0(x) - \gamma h_1(x)] h_1(x)}$$

$$\cancel{2[12 - 8 - \gamma 2]} \cancel{2} = 0$$

$$4 = \gamma 2 \Rightarrow \gamma_1 = 2$$

$$y = h_0(x) + \gamma_1 h_1(x)$$

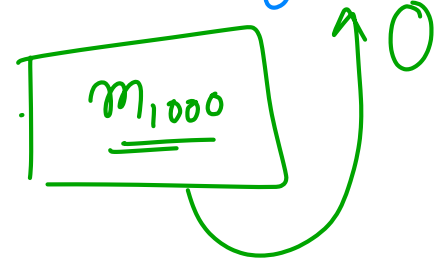
$$y = 8 + (2 \times 2) \Rightarrow 8 + 4$$

$$\Rightarrow \textcircled{12}$$

hyper parameter



$m_0 \quad m_1 \quad m_2 \quad \dots$



# of model  $M \uparrow \rightarrow \text{error tend to } 0 \rightarrow \text{Overfitting}$

$M \downarrow \rightarrow \text{Values will be close to mean value}$   
Underfitting

Value of  $d$   $\rightarrow$  depth of each tree

base learner

$d \uparrow \rightarrow \text{Overfitting}$

$d \downarrow \rightarrow \text{Underfitting}$



# Regularisation by Shrinkage

$$F_M(x) = h_0(x) + \sum_{i=1}^M \gamma_i h_i(x)$$

$M$  → hyperparameter

underfitting

Overfitting

Learning Rate  $(0 < \gamma \leq \infty)$   
recomm

# problem

- \* Slow
- \* Sequential
- \* It overfits very Quickly.

→ Randomisation → RS + CS

GBDT = Pseudo Residual  
+ Additive  
Combining

RS + CS + PR  
+ AC  
Stochastic GBDT

Categorical data

Catboost

Adaboost

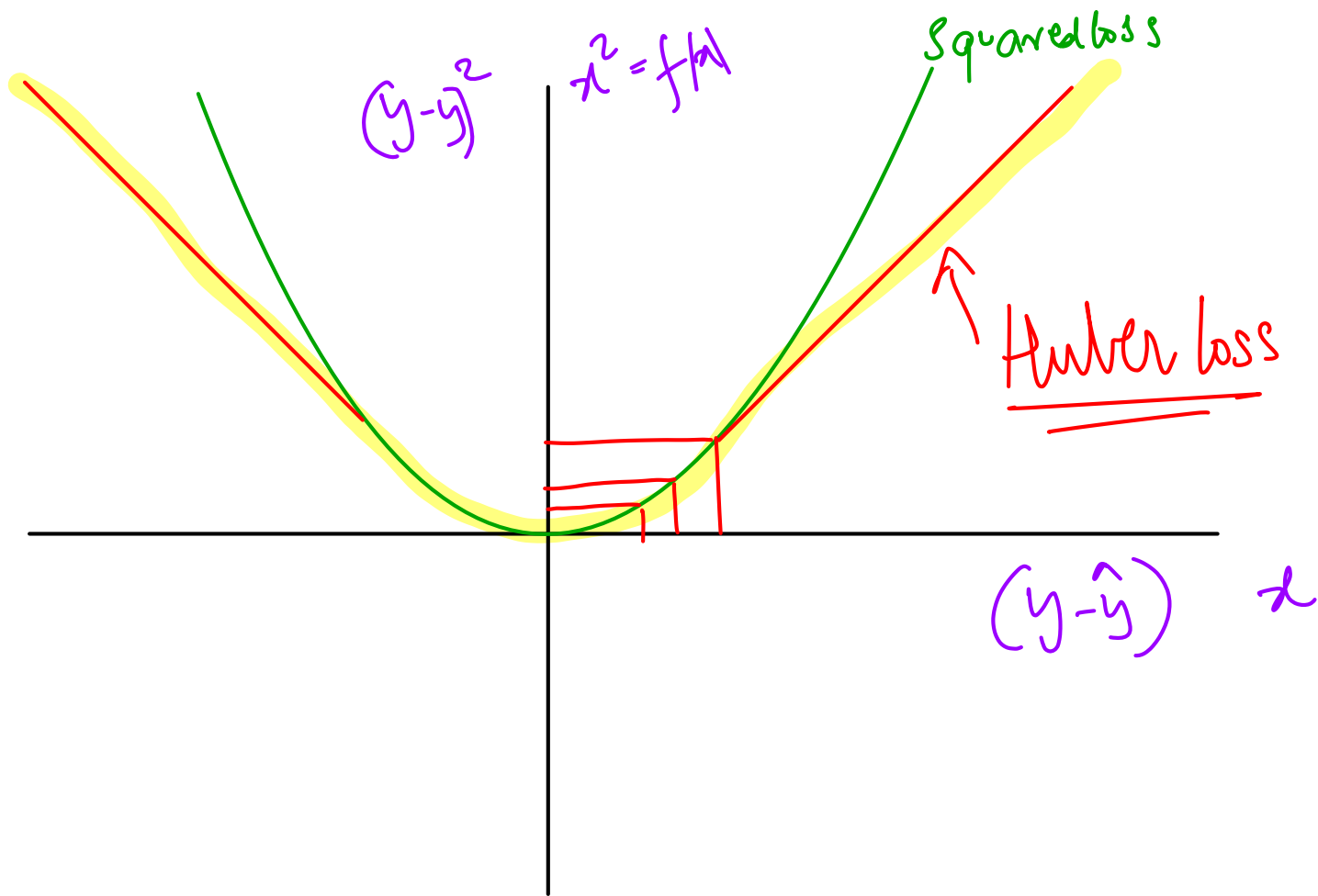
LIGHTGBM

Impact of Outlier  $\rightarrow$  GBDT  
Churnes

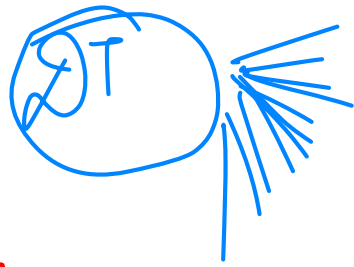
y	$\hat{y}$	err
20	48	-28
4	48	-44
6	48	-42
10	48	-38
200	48	152

Outlier Impact  $\uparrow\uparrow\uparrow\uparrow\uparrow$

model will focus on  
reducing errors associated  
with outliers.



GBDT → slow



X GBOOST → Not slow

LIGHT GBM → Not slow

