

Q7. Quo?

de modellen

⇒ Sigma ?

~ "Sigmoid"

Agenda :

- ① GBDT
- ② Bias Variance
- ③ Sklearn

for train

- function.

$y^{(i)}$

$\hat{f}_k(x^{(i)}) \approx \hat{y}^{(i)}$

$L(y^{(i)}, \hat{f}_k(x^{(i)})$

$z^{(i)} = \hat{f}_k(x^{(i)})$

$\frac{\partial L}{\partial z^{(i)}} \Rightarrow$

$\frac{\partial L}{\partial z^{(i)}} \Rightarrow -$

$\delta F_k(x^i)$

neg. derivative  
of gradient

residual.

The diagram illustrates the relationship between gradients, residuals, and the "Huber loss" function.

- Pseudo Residual:** A box labeled "Pseudo Residual" contains the text "neg. derivative" and "neg. gradients".
- approx. (residuals):** An arrow points from the "Pseudo Residual" box to the text "approx. (residuals)".
- "Huber loss":** A large oval represents the "Huber loss" function. Inside the oval, the formula  $y^{(i)} - f_k(x^i)$  is written. An arrow points from the text "approx. (residuals)" to the oval. The word "mse" is written next to the oval.

At  $j^{th}$  Model iteration

$$h_j(x) \leftarrow \{x^{(i)}, \underline{\text{error}}^{(i)}\} \rightarrow e_m^{(i)} = y^{(i)} - \underline{F_{j-1}(x)}$$

$\text{dy get-gradient } (-)$

$\left[ \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \right]^2$

Automated

$\text{error}^{(i)} \rightarrow \text{residuals}$

$\text{error}^{(i)} \approx \text{pseudo residual}^{(i)}$

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

$Rf \rightarrow m$

$$\underline{\text{Stage } p} \stackrel{\circ}{\circ}$$

$$F_p(x) \Rightarrow F_{p-1}(x) + \underbrace{\alpha_p \cdot h_p(x)}_{\circ}$$

# The Algorithm

Initialize model with a constant value:

$$F_0(x) = \arg \min_x \sum_{i=1}^n L(y_i, x) \rightarrow \text{avg} := h_0(x) = f_0(x)$$

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

- 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
- 3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:
$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)) = 0$$

4. Update the model:  
 $F_m(x) = F_{m-1}(x) + \gamma_m h_m(x)$

3. Output  $F_M(x)$ .

$f_1(x) = F_0(x) + \gamma_1 h_1(x)$

$r_1$   $r_2$   $\dots$   $r_m$   
 $\ll h_i(x)$

$$f_M(x) \Rightarrow h_0(x) + \alpha_1 \cdot h_1(x) + \dots + \alpha_M h_M(x)$$

$\alpha$ 's  $\rightarrow \gamma$ 's

$$f(x) = (x - s)^2$$

At  $x = ?$

$F_1(x) =$

$F_2(x) =$

$\arg \min_x (x - s)^2$

$$f_M(x) = h_0(x) + \sum_{m=1}^M r_m \cdot h_m(x)$$

Overfitting

- [as  $M \uparrow \rightarrow$  overfit  $\uparrow$ ]  
[as  $M \downarrow \rightarrow$  overfit  $\downarrow$ ]
- Regularization by Shin Kage

$$F_m(x) = f_{m-1}(x) + \underbrace{v_m}_{\text{learning rate}} \cdot \dots$$

) if  $v = 0.0001$  (small)  
as  $v \downarrow \rightarrow$  underfit.  
overfit  $\downarrow$

if  $V = 1$  ( $V$ . large)  
as  $V \uparrow \rightarrow$  overfit  $\uparrow$   
(high variance)

GBDT  $\xrightarrow{M}$   $\xrightarrow{\nu}$   $\left\{ \xrightarrow{\text{grid search}}$

## Time & Space Complexity.

If was trivially parallelizable algo,  
GBDT is not easily parallelizable -

Time Complexity:

1 DT  $\rightarrow \mathcal{O}(\text{depth})$   
M DT  $\rightarrow \mathcal{O}(M^{\text{depth}})$

300      2-4-5-6

$$O(M \cdot DT's + M \cdot T's) \\ \hookrightarrow M \text{ gammas.}$$

Diagram illustrating data quality issues:

- A box contains the text Outliers | Missing Values.
- An arrow points from DT (Doubts) to the box.
- An arrow points from PSP ↑ (Probability of Success Probability ↑) to the box.
- Below the box, a bracket spans the text 40 mins.
- Below the bracket, a bracket spans the text 9pm? ghar.

Diagram illustrating the relationship between  $r_{\text{res}}$  and  $\Delta r$ .

The diagram shows a sequence of values:  $8.4 \rightarrow 8.61 \rightarrow 8.62 \rightarrow 8.75 \rightarrow 9.01$ . A bracket below the first four values is labeled  $\Delta r = 0.22$ . Above the last value,  $9.01$ , is the label  $12^-$ .

A large bracket above the entire sequence is labeled  $r_{\text{res}} = 0.12$ .

Diagram illustrating a neural network architecture:

- Input Layer:**  $F_M(x)$  (represented by a blue rounded rectangle containing  $x_1$  and  $x_2$ )
- Hidden Layer:**  $S_{\gamma(i)}$  (represented by an orange rounded rectangle containing  $\gamma(i)$  and  $\exp(i)$ )
- Output Layer:**  $\exp^{(i)}$  (represented by a pink rounded rectangle containing  $\exp(i)$ )

The diagram shows the flow of data from the input layer through the hidden layer to the output layer.

$$\gamma_m = \frac{\partial f_{m-1}(x)}{\partial x} - \alpha L$$