# Recap

## Bagging vs Boosting

Bagging $\Rightarrow$ Bootstrap Aggregation



$L_1$ $\quad$ $L_2$ $\quad$ $L_3$

| Aggregation |
|---|

↓ Final Results

Boosting $\Rightarrow$ Additive Combining



$L_1$ $\quad$ + $\quad$ $L_2$ $\quad \longrightarrow \quad$ $L_3$ $\qquad \Rightarrow$ Final Result

Edo

Selection

MT

Hyp

Sequential Process

# Steps in Boosting

## Step 0 : Mean Model and Error residual

| | Height | Gender | Weight | $h_0(x)$ | Error |
|---|---|---|---|---|---|
| ① | 1·6 | M | 82 | 67 | 15 |
| ② | 1·5 | F | 55 | 67 | −12 |
| ③ | 1·4 | F | 66 | 67 | −1 |
| ④ | 1·4 | M | 65 | 67 | −2 |

$$\frac{82+55+66+65}{4}$$
$$4$$
$$\Rightarrow 67$$

## Step 1 : Model 1 $\Rightarrow$ $D \ni \{ x_i, Error_0 \}$



| | $h_0(x)$ | Error$_0$ | $f_1(x) \Rightarrow h_0x + h_1x$ | Err$_1$ |
|---|---|---|---|---|
| ① | 67 | 15 | 73·5 | |
| ② | 67 | −12 | 60·5 | |
| ③ | 67 | −1 | 60·5 | |
| ④ | 67 | −2 | 73·5 | |

Preds $\Rightarrow$
$$\frac{13}{2} \Rightarrow 6·5$$

Preds :
$$\frac{-12 + -1}{2} \Rightarrow -6·5$$

$h_1 x$

and Combine Stage I with stage 0 Predictions

$$f_1(x) \Rightarrow f_0 x + \gamma f_1(x)$$

| Height | Gender | Weight | Err. | $f_1(x)$ | $f_0(x)$ | $f_1 x = f_0(x) + f_1(x)$ |
|--------|--------|--------|------|----------|----------|---------------------------|
| 1.6 | M | 82 | 15 | 6.5 | 67 | 73.5 |
| 1.5 | F | 55 | -12 | -6.5 | 67 | 60.5 |
| 1.4 | F | 66 | -1 | -6.5 | 67 | 60.5 |
| 1.4 | M | 65 | -2 | 6.5 | 67 | 73.5 |

$$f_2(x) \Rightarrow f_0(x) + f_1(x) + f_2(x)$$

$L_0$

$L_1$

$(x_i, Err)$  $L_3$

$$f_n(x) \Rightarrow f_0 x + \gamma_1 f_1 x + \gamma_2 f_2(x) \cdots f_n(x)$$
$\gamma_n$

⇒ $\gamma_n$ Can't be calculated at Same Time

LR ⇒ $\omega_0 x + \omega_1 x_1 + \omega_2 x_2 \cdots \omega_n x_n$

⇒ all of $\omega$'s can be calculated at Same time

⇒ Boosting is Slow due to Sequential Nature

(Regression

$$L(y_i, \hat{y}_{(i)}) \Rightarrow \text{m.s.e} \Rightarrow \sum_{i=1}^{M} (y_i - \hat{y}_i)^2$$

$$\hat{y} \Rightarrow f_k(x)$$

(Prediction at Stage K)

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

$$\frac{\partial L}{\partial \hat{y}} \Rightarrow -2(y - \hat{y})$$

$$-\frac{\partial L}{\partial \hat{y}} \Rightarrow 2 \cdot (y - \hat{y})$$

negative gradient of Loss w.r.t output

residual

ignore

Pseudo-Residual $\Rightarrow -\frac{\partial L}{\partial \hat{y}}$

Residual is proportional to −ve gradient of Loss w.r.t prediction

How do we use Pseudo-Residual?

⇒ for some model $m_g$

$$M_g \longrightarrow \{x_j, err_{j-1}^i\}$$

$$y_i - F_{j-1}(x)$$

↳ residual

Replace $err_{j-1}$ with it's Pseudo Residual

$$-\frac{\partial L}{\partial F_{j-1}(x.)}$$

Optimizing for Pseudo Residual will also optimize for the Loss function

\* $L$ ⇒ MSE or RMSE for Regression

Log-Loss for Classification

⇒ At any stage $K$, to optimize we need to calculate gradient of Loss function w.r.t output at Stage $K-1$

# GBDT

Pseudo-Residual using Gradient
(optimize for this)

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

$$L(y, \hat{y}(x)) \begin{cases} \text{MSE (Regression)} \\ \text{LogLoss (Classification)} \end{cases}$$

$M$

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:
   $$F_0(x) = \arg\min_{\gamma} \sum_{i=1}^{n} L(y_i, \gamma).$$  ← 6

2. For $m = 1$ to $M$:
   1. Compute so-called *pseudo-residuals*:
      $$r_{im} = -\left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \ldots, n.$$
   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 $\{(x_i, r_{im})\}_{i=1}^n$.
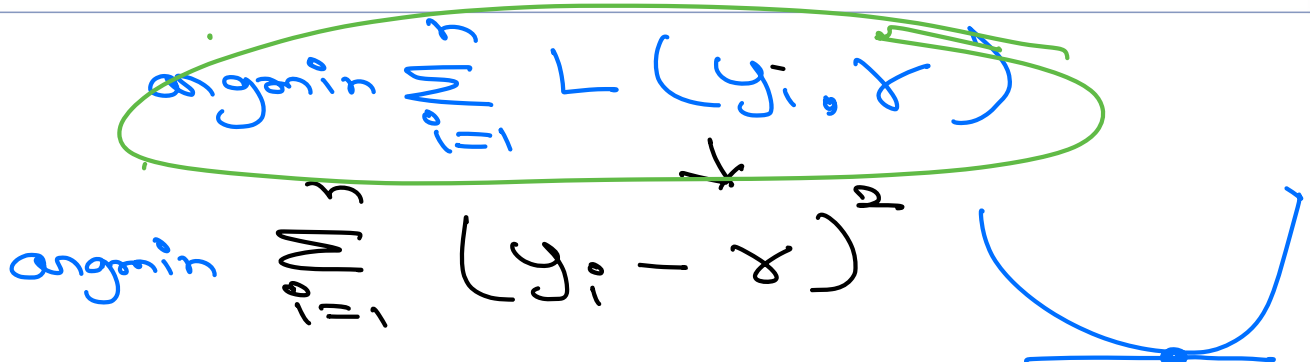   3. Compute multiplier $\gamma_m$ by solving the following one-dimensional optimization problem:
      $$\gamma_m = \arg\min_{\gamma} \sum_{i=1}^{n} L\left(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)\right). \quad \leftarrow$$
      optimal value of $\gamma$
   4. Update the model:
      $$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output $F_M(x).$

$$\arg\min \sum_{i=1}^{M} L(y_i, \gamma)$$

$$\arg\min \sum_{i=1}^{M} (y_i - \gamma)^2$$

To find minima Take diff and set it equal to Zero

$$\frac{\partial L}{\partial \gamma} = \frac{\partial \sum_{i=1}^{M}(y_i - \gamma)^2}{\partial \gamma} \Rightarrow 0$$

$$\Rightarrow -2\left(\sum_{i=1}^{M}\left(y_i - \gamma\right)\right) = 0$$

$$\Rightarrow \sum_{i=1}^{M}(y_i) - (M\gamma) = 0 \times -2$$

$$\Rightarrow M\gamma = \sum_{i=1}^{M}(y_i)$$

$$\Rightarrow \gamma \Rightarrow \frac{1}{M}\sum_{i=1}^{M}(y_i)$$

Step 1 is Nothing but Creation of Mean Model

---

Step 1 $\Rightarrow$ initialize with mean Model $\quad \underset{\gamma}{\text{argmin}}\sum_{i=1}^{M} L(y_i, \gamma)$

For $\Rightarrow$ Avg Model $\quad$ avg

Step 2

2.1 $\Rightarrow$ Calculate pseudo Residual

$$-\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)}$$

2.2 → fit another Low depth
Tree on pseudo-residual

2.3 : Using gradient descent
to calculate the Best
Value of $\gamma_m$
$$\left(\begin{array}{c}\text{the one which}\\ \text{gives min}\\ \text{Loss}\end{array}\right)$$

2.4

$m=1$

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

$$F_1 \rightarrow F_0 \cdot x + \gamma_m f_0(x)$$

Repeat Step 2 for M times

$$\rightarrow \text{model} \atop f_1 x, f_2 x \dots f_m x$$

We are optimizing

$$\rightarrow \gamma_1, \gamma_2 \dots \gamma_m$$
$$\downarrow$$
weighted addition

Base Learner ⎯⎯⎤ High Bias
                    ⎦→ Low Variance

M⊘ Number of Base-learner

M ⟶ ↑ Bias Reduce
                and
        Variance Variance

if m is a very high number, Variance
of Final Boosted Model will also
be Very High

(Overfit)

Base learners ① Depth
                    ↓
            High Depth will
            lead to High
            Variance and
            Overfitting

To Find right Balance we will
have to Tune M
            D

**Q** Is there a regularization Term in GBDT

$$F_m(x) \ni \beta_0 x + \sum_{i=1}^{M} \gamma_m \rho_m(x)$$

$$+$$

Regularization (Shrinkage)

$$\Downarrow$$

$$F_m(x) \ni \beta_0 x + \eta \sum_{i=1}^{M} \gamma_m \rho_m(x)$$

Constant Value (0.1)

(learning Rate)

$\eta \ni 0 \implies$ Underfitted Model

$\eta = \infty \ni$ Overfitting

Issues with GBDT :

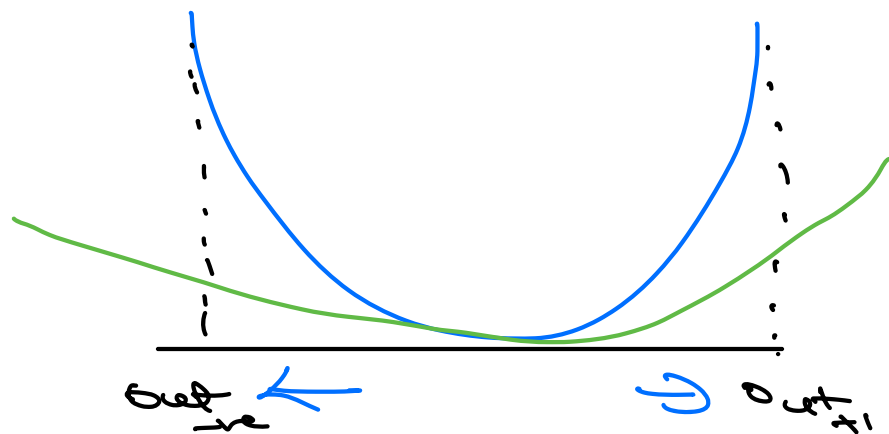1) No Parallelization due to Sequential learning

2) Prone to Overfitting

**Stochastic** ⇒ GBDT + Randomization
**GBDT**                    ( Column Sampling)
                              Row Samplin

$$d = 3$$

|        | $m_1$    | $m_2$    | $m_3$    |
|--------|----------|----------|----------|
| $n_1$  | $e_{11}$ | $e_{12}$ | $e_{13}$ |
| $n_2$  | $e_{21}$ | $e_{22}$ | $e_{23}$ | ⇒ very Right
| $n_3$  | $e_{31}$ | $e_{32}$ | $e_{33}$ |           or very Low

RMSE
  or

Huber
Loss



Out    ←           ⇒ Out
-ve                   +1

\* Replacing Loss function with Something that doesn't Explode ie gives High Value for Outliers

\* GBDT Implemention
\* Variations of GBDT