

## GBDT Framework

→ **Idea**:  $G_m(x) = G_{m-1}(x) + g_m(x)$  Loss function:  $L(G_m, y)$

①  $L(G_m, y)$  on data point  $(x, y)$   
 $= L(G_{m-1}(x) + g_m(x), y)$

**Aim**:  $g_m(x) = \underset{g}{\operatorname{argmin}} L(G_{m-1}(x) + g(x), y) \quad \forall x$

②  $L(G_{m-1}(x) + g(x), y) \approx \left. \frac{\partial L(G, y)}{\partial G} \right|_{G=G_{m-1}(x)} \cdot g(x) + L(G_{m-1}(x), y)$

↓ **first-order**

recommend function  $g_m(x)$  should fit to  $-\left. \frac{\partial L(G, y)}{\partial G} \right|_{G=G_{m-1}(x)}$

**Alternative Interpretation**

this formulation  
gave us the  
choice of LSE

$$\begin{aligned} L(G_{m-1} + g(x), y) &\approx L(G_{m-1}(x), y) + \left. \frac{\partial L(G, y)}{\partial G} \right|_{G=G_{m-1}(x)} g(x) + \frac{1}{2} g(x)^2 \\ &= \text{const} + \frac{1}{2} \left( g(x) + \left. \frac{\partial L(G, y)}{\partial G} \right|_{G=G_{m-1}(x)} \right)^2 \end{aligned}$$

$$\propto \left[ -g(x) - \left. \frac{\partial L(G, y)}{\partial G} \right|_{G=G_{m-1}(x)} \right]^2$$

⇒ as for GBDT, use LSE to achieve  $g_m(x)$  → m-th base classifier

**XGBoost**  $G_m(x) = G_{m-1}(x) + g_m(x)$  Loss function:  $L(G_m, y)$

for simplicity, define  $\begin{cases} g^m(x) = \left. \frac{\partial L(G, y)}{\partial G} \right|_{G=G_m(x)} \\ h^m(x) = \left. \frac{\partial^2 L(G, y)}{\partial G^2} \right|_{G=G_m(x)} \end{cases}$

①  $g_m(x) = \underset{g(x)}{\operatorname{argmin}} L(G_{m-1}(x) + g(x), y)$

Similar Non-parametric Framework

②  $L(G_{m-1}(x) + g(x), y)$

$$\approx L(G_{m-1}(x), y) + g^m(x) \cdot g(x) + \frac{1}{2} h^m(x) \cdot g(x)^2$$

$$\propto (g^m(x) g(x) + \frac{1}{2} h^m(x) g(x)^2)$$

$$\propto h^m(x) \left[ -\frac{g^m(x)}{h^m(x)} + g(x) \right]^2$$

Consider the  $L(G(x), y)$  on total dataset

$$L_{\text{total}} \approx \sum_n h^m(x_n) \left[ -\frac{g^m(x_n)}{h^m(x_n)} - g(x_n) \right]^2$$

$$\propto \sum_n w^m(x_n) \left[ -\frac{g^m(x_n)}{h^m(x_n)} - g(x_n) \right]^2$$

where  $w^m(x_n) = \frac{h^m(x_n)}{\sum_{n'} h^m(x_{n'})}$

Weighted LSE compared

with GBDT (LSE)

AdaBoost

First, Conclusion

AdaBoost

XGBoost with exponential loss  
to achieve  $g_m(\cdot)$   
we primal formulation without  
Taylor approximation on  $\alpha$

$$L_{\text{exp}}(G(x), y) = \exp(-y G(x))$$

$$\textcircled{1} g^m(x) = \frac{\partial L_{\text{exp}}(G, y)}{\partial G} \Big|_{G = G_{m-1}(x)}$$

$$= -y \exp(-y G_{m-1}(x))$$

$$h^m(x) = \frac{\partial^2 L_{\text{exp}}(G, y)}{\partial G^2} \Big|_{G = G_{m-1}(x)}$$

$$= y^2 \exp(-y G_{m-1}(x))$$

$$= 1 \cdot \exp(-y G_{m-1}(x)) = \exp(-y G_{m-1}(x))$$

$\textcircled{2}$  suppose that additive model is  $G_m(x) = G_{m-1}(x) + \alpha_m g_m(x)$

with XGBoost Framework, we have:



$$\begin{aligned}
 (\alpha_m, g_m) &= \operatorname{argmin}_{\alpha, g} \sum_{n=1}^N w^m(x_n) \left[ -\frac{g^m(x_n)}{h^m(x_n)} - \alpha g(x_n) \right]^2 \\
 &= \operatorname{argmin}_{\alpha, g} \sum_{n=1}^N w^m(x_n) [y_n - \alpha g(x_n)]^2 \\
 &= \operatorname{argmin}_{\alpha, g} \sum_{n=1}^N w^m(x_n) [1 + \alpha^2 \underbrace{g^2(x_n)}_1 - 2\alpha y_n g(x_n)] \\
 &= \operatorname{argmin}_{\alpha, g} \sum_{n=1}^N w^m(x_n) [\alpha^2 - 2\alpha y_n g(x_n)]
 \end{aligned}$$

where  $w^m(x_n) = \frac{h^m(x_n)}{\sum_{n'} h^m(x_{n'})} = \frac{\exp(-y_n G_{m-1}(x_n))}{\sum_{n'} \exp(-y_{n'} G_{m-1}(x_{n'}))}$

observe that, we can minimize on  $g(\cdot)$  first for the separability

$$g_m = \operatorname{argmin}_g \sum_{n=1}^N w^m(x_n) \cdot (-y_n g(x_n))$$

$$= \operatorname{argmin}_g \sum_{n=1}^N w^m(x_n) \cdot \frac{1 - y_n g(x_n)}{2}$$

$$= \operatorname{argmin}_g \sum_{n=1}^N w^m(x_n) \mathbb{1}\{y_n \neq g(x_n)\}$$

What we do in  
AdaBoost!

minimize weighted  
training error to  
achieve  $g_m(\cdot)$

then, we use the first objective (without second-order approximation) to  
optimize on  $\alpha$ ! (this is more trivial part)

$$\alpha_m = \operatorname{argmin}_{\alpha} \sum_{n=1}^N \exp(-y_n G_{m-1}(x_n)) \cdot \exp(-y_n \alpha g_m(x_n))$$

we first-order condition, we can achieve

that  $\alpha_m = \log\left(\frac{1 - \epsilon_m}{\epsilon_m}\right) \leftarrow \epsilon_m = \sum_{i=1}^N w^m(x_i) \mathbb{1}\{y_i \neq g_m(x_i)\}$

Rmk: this is what we do in GBDT framework with step  $\alpha$ ,

- {
- ① we generate the  $m$ -th optimal classifier by Taylor Approximation Loss function
  - ② we generate the  $m$ -th optimal step length  $\alpha_m$  by Loss function without Taylor Approximation

★

⇒ To conclude, Adaboost can be viewed as a combination of { XGBoost, GBDT }

in the sense that

- use 2-order Taylor Approximation on exp-loss to achieve the  $m$ -th base classifier
- after achieving classifier, back to exp-loss without Taylor Approx. to directly solve for optimal step length!

XGBoost Part

GBDT Part

Another View Point of Adaboost → with no Taylor Approximation

1) Model:  $G_m(x) = G_{m-1}(x) + \alpha_m g_m(x)$

2) Loss Function:  $L_{\exp}(y_i, G_m(x_i)) = \exp(-y_i G_m(x_i))$

$$R_{\exp}(G_m) = \sum_{n=1}^N L_{\exp}(y_n, G_m(x_n))$$

$$= \sum_{n=1}^N \exp(-y_n G_{m-1}(x_n)) \cdot \exp(-y_n \alpha_m g_m(x_n))$$

$$:= \sum_{n=1}^N \underline{w_n^{(m)}} \exp(-y_n \alpha_m g_m(x_n))$$

$$= \sum_{n=1}^N w_n^{(m)} \exp(-\alpha_m) \mathbb{1}_{\{y_n = g_m(x_n)\}}$$



$$+ \sum_{n=1}^N w_n^{(m)} \exp(\alpha_m) \mathbb{1}\{y_n \neq g_m(x_n)\}$$

$$= \exp(-\alpha_m) \sum_{n=1}^N w_n^{(m)} \mathbb{1}\{y_n = g_m(x_n)\} \\ + \exp(\alpha_m) \sum_{n=1}^N w_n^{(m)} \mathbb{1}\{y_n \neq g_m(x_n)\}$$

$$= (\exp(\alpha_m) - \exp(-\alpha_m)) \cdot \sum_{n=1}^N w_n^{(m)} \mathbb{1}\{y_n \neq g_m(x_n)\} \\ + \exp(-\alpha_m) \sum_{n=1}^N w_n^{(m)}$$

$$\rightarrow (\hat{\alpha}_m, \hat{g}_m) = \underset{\alpha_m, g_m}{\operatorname{argmin}} \operatorname{Remp}(G_m)$$

$$\Rightarrow \textcircled{1} \hat{g}_m = \underset{g_m}{\operatorname{argmin}} \sum_{n=1}^N w_n^{(m)} \mathbb{1}\{y_n \neq g_m(x_n)\}$$

$$\textcircled{2} \hat{\alpha}_m = \log\left(\frac{1 - e_m}{e_m}\right) \quad e_m := \frac{\sum_{n=1}^N w_n^{(m)} \mathbb{1}\{y_n \neq g_m(x_n)\}}{\sum_{n=1}^N w_n^{(m)}}$$

$\textcircled{3}$  update formula for  $w_n^{(m)}$ :

$$w_n^{(m+1)} = \begin{cases} w_n^{(m)} \cdot \exp(-\alpha_m) & \text{if } g_m(x_n) = y_n \\ w_n^{(m)} \exp(\alpha_m) & \text{if } g_m(x_n) \neq y_n \end{cases}$$