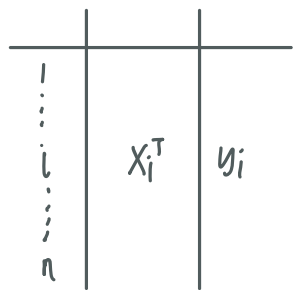
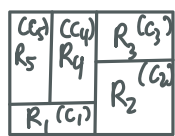


Extreme gradient boosting



$$f(x) = \sum_{k=1}^K h_k(x)$$



$$s_i = f(x_i) = \sum_{k=1}^K h_k(x_i)$$

$$p_i = \frac{e^{s_i}}{1 + e^{s_i}} = \Pr(y_i = 1 | s_i)$$

$$1 - p_i = \frac{1}{1 + e^{s_i}}$$

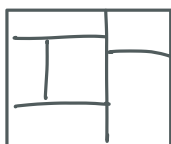
Boosting

At iteration t

$$s_i = \sum_{k=1}^{t-1} h_k(x_i) + h_t(x_i)$$

current committee \hat{s}_i (fixed)

new tree



Δs_i



Recall

$$\text{Likelihood} = \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i} = \prod_{y_i=1} p_i \prod_{y_i=0} (1-p_i)$$

$$\log\text{-lik} = \sum_{i=1}^n [y_i \log p_i + (1-y_i) \log(1-p_i)]$$

$$= \sum_{i=1}^n [y_i (s_i - \log(1+e^{s_i})) + (1-y_i) (-\log(1+e^{s_i}))]$$

$$= \sum_{i=1}^n [y_i s_i - \log(1+e^{s_i})]$$

$$L(s_i) = L(\hat{s}_i + h_t(x_i))$$

$$\approx L(\hat{s}_i) + L'(\hat{s}_i) h_t(x_i) - \frac{1}{2} L''(\hat{s}_i) h_t(x_i)^2$$

$$= \text{const } t$$

$$L'(s_i) = y_i - \frac{e^{s_i}}{1+e^{s_i}} = y_i - p_i = e_i$$

$$L''(s_i) = \frac{d}{ds_i} \left(-\frac{e^{s_i}}{1+e^{s_i}} + 1 \right) = \frac{d}{ds_i} \left(\frac{1}{1+e^{s_i}} \right) = \frac{-e^{s_i}}{(1+e^{s_i})^2} = -\frac{e^{s_i}}{1+e^{s_i}} \frac{1}{1+e^{s_i}} = -p_i(1-p_i) = -w_i$$

$$= -\frac{1}{2} \hat{w}_i \left[h_t(x_i)^2 - 2 \frac{\hat{e}_i}{\hat{w}_i} h_t(x_i) \right] + \text{constant}$$

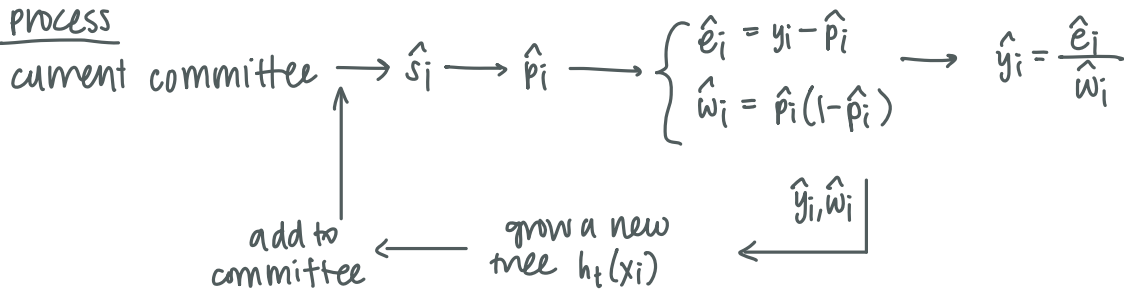
$$= -\frac{1}{2} \hat{w}_i \left[\frac{\hat{e}_i}{\hat{w}_i} - h_t(x_i) \right]^2 + \text{constant}$$

\downarrow
 \hat{y}_i

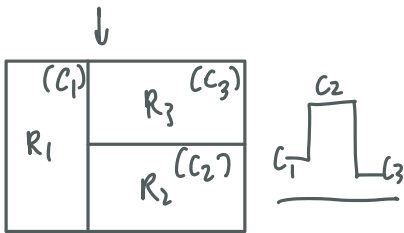
iterated reweighted least squares

each step grows a new tree to expand imperfections

PROCESS



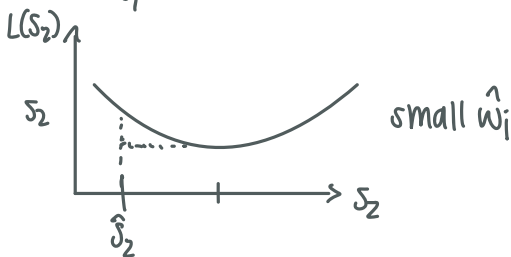
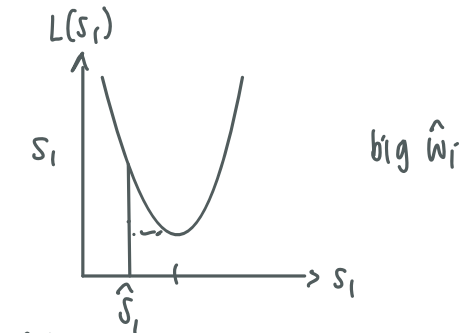
$$\min \sum_{i=1}^n \hat{w}_i (\hat{y}_i - h_t(x_i))^2 + \lambda M + \gamma \sum_{m=1}^M c_m^2$$



we minimize this loss function b/c we want to grow a new tree to fit the data

\hat{y}_i represents the imperfections we want to expand

\hat{e}_i gives us the gradient gradient is scaled by \hat{w}_i



$$h_t(x_i) = \sum_{m=1}^n c_m \mathbb{1}(x_i \in R_m)$$

$$\hat{c}_m = \frac{\sum_{i=1}^n y_i \mathbb{1}(x_i \in R_m)}{\sum_{i=1}^n \mathbb{1}(x_i \in R_m) + \gamma}$$

shrinkage parameter

$$\sum_{i=1}^n \hat{w}_i (\hat{y}_i - \hat{c}_m) \mathbb{1}(x_i \in R_m) + \lambda M$$

weighted loss function

grow tree by recursive partitioning, only depends on regions/partitioning

Adaboost (origin)

in XGboost, $h_k(x) \in \mathbb{R}$, regression tree as opposed to binary classifier

$$f(x) = \sum_{k=1}^K \beta_k h_k(x) \rightarrow \text{classifier} \rightarrow \{+, -\}$$

$$\hat{y} = \text{sign}(f(x)) = \begin{cases} + & f(x) > 0 \\ - & f(x) < 0 \end{cases}$$

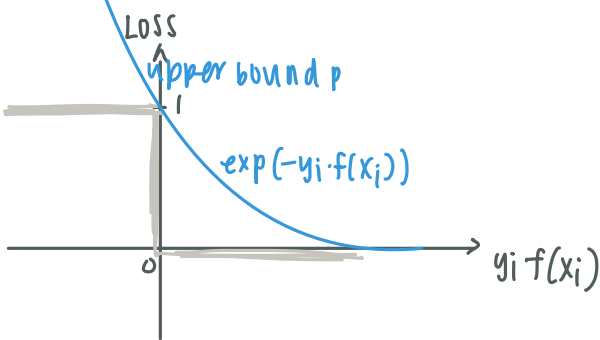
final result is a classifier; each tree is a classifier
withholds information as it does not include info on magnitude

purpose is to answer a theoretical question

weak learn = strong learn??

can the combination of weak classifiers yield a strong classifier?

non continuous loss funct



$$\text{Loss} = \sum_{i=1}^n \exp(-y_i s_i) = \sum_{i=1}^n \exp(-y_i (\hat{s}_i + \beta_t h_t(x_i)))$$

$$s_i = \sum_{k=1}^{t-1} \beta_k h_k(x_i) + \beta_t h_t(x_i)$$

current committee \downarrow \hat{s}_i fixed

$\beta_t h_t(x_i)$ \downarrow to be learned

$$= \sum_{i=1}^n \exp(-y_i \hat{s}_i) \exp(\beta_t y_i h_t(x_i))$$

D_i (distribution)

$$D_i \xrightarrow{\text{normalize}} \frac{D_i}{\sum_{i=1}^n D_i}$$

$\sum_{i=1}^n D_i = 1$ difficult examples

$D = (D_i, i=1, \dots, n)$ pays attention to examples that current committee did not do well

$$\text{Loss} = \sum_{i=1}^n D_i \exp(\beta_t y_i h_t(x_i))$$

$+1 : y_i = h_t(x_i)$
 $-1 : y_i \neq h_t(x_i)$ did not classify correctly

$$= \sum_{y_i = h_t(x_i)} D_i e^{-\beta_t} + \sum_{y_i \neq h_t(x_i)} D_i e^{\beta_t}$$

a b

where $\sum_{y_i = h_t(x_i)} D_i = 1 - \epsilon$, $\sum_{y_i \neq h_t(x_i)} D_i = \epsilon$

$$a+b \geq 2\sqrt{ab}$$

$$(\sqrt{a}-\sqrt{b})^2 \geq 0$$

"=" $a=b$

$\min \beta_t \geq (1-\epsilon)(\epsilon)$ lower bound

$$(1-\epsilon)e^{-\beta t} = \epsilon e^{\beta t}$$

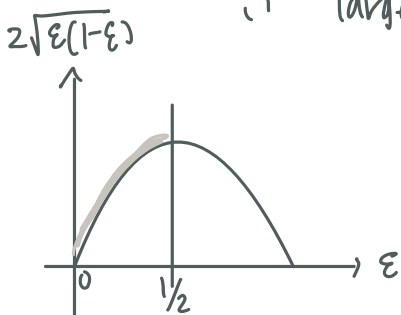
$$\frac{1-\epsilon}{\epsilon} = e^{2\beta t}$$

$$\hat{\beta}_t = \frac{1}{2} \log \frac{1-\epsilon}{\epsilon}$$

voting right of $h_t(\cdot)$

If ϵ is small, then ratio of $\frac{1-\epsilon}{\epsilon}$ is very large, so $\hat{\beta}_t$ is large

" large, " " small " small



from range $(0, \frac{1}{2})$ function is increasing

recut $h_t(\cdot)$ with minimal ϵ

assign $\hat{\beta}_t = \frac{1}{2} \log \frac{1-\epsilon}{\epsilon}$

PROCESS:

current committee $\rightarrow D \rightarrow$ grow $h_t(\cdot)$
 $\hat{\beta}_t$

