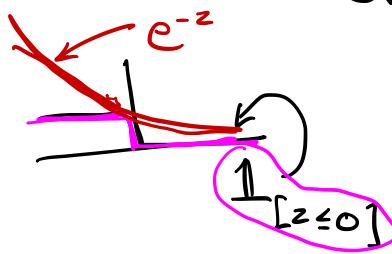


Boosting

Boosting - "the statistical view"

misclassification error = $\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[y_i f(x_i) \leq 0]}$



$$\leq \frac{1}{n} \sum_{i=1}^n e^{-y_i f(x_i)} \quad \text{exp loss}$$

choose f to be a linear model, a linear combination of "weak" classifiers.

$$f(x) = \sum_{j=1}^p \lambda_j h_j(x)$$

if I'm boring, $h_j(x) = x_{\cdot j}$ so $f(x) = \sum_{j=1}^p \lambda_j x_{\cdot j}$

$$R^{\text{train}}(\bar{\lambda}) = \frac{1}{n} \sum_i e^{-y_i \sum_j \lambda_j h_j(x_i)} = \frac{1}{n} \sum_i e^{-\sum_j y_i h_j(x_i) \lambda_j} = \frac{1}{n} \sum_i e^{-(\bar{M}\bar{\lambda})_i}$$

where $M = \begin{bmatrix} & & & \\ i & & & j \\ & & & \\ n & & & \end{bmatrix}^{(n \times p \cdot p \times 1)}$

"matrix of margins" $i \quad j \quad n$

margin of i^{th} point
for j^{th} weak classifier.

Assume all weak classifiers are binary so $h_j(x_i) = \pm 1$
Then $M_{ij} = \pm 1$.

$$R^{\text{train}}(\bar{\lambda}) = \frac{1}{n} \sum_i e^{-(\bar{M}\bar{\lambda})_i}$$

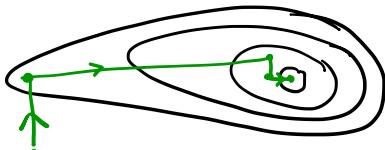
Do coordinate descent for λ

Step 1 - "coordinates" are j's.

Find the steepest coordinate

Step 2 - do a line search in that direction

Step 1: Say we are at λ_t and we want steepest direction.



$$\begin{aligned} \lambda_t \in \operatorname{argmax}_j & \left[-\frac{d R^{\text{train}}(\lambda_t + \alpha e_j)}{d \alpha} \Big|_{\alpha=0} \right] \\ & - \frac{d}{d \alpha} \left(\sum_i e^{-[\bar{M}(\bar{\lambda}_t + \alpha e_j)]_i} \right) \Big|_{\alpha=0} \\ & - \frac{d}{d \alpha} \left(\sum_i e^{-(\bar{M}\bar{\lambda}_t)_i - \alpha M_{ij}} \right) \Big|_{\alpha=0} \\ & \frac{1}{n} \sum_i \left[-\frac{d}{d \alpha} \left(e^{-(\bar{M}\bar{\lambda}_t)_i - \alpha M_{ij}} \right) \right] \Big|_{\alpha=0} \\ & \frac{1}{n} \sum_i M_{ij} e^{-(\bar{M}\bar{\lambda}_t)_i} + 0 \\ \lambda_t \in \operatorname{argmax}_j & \frac{1}{n} \sum_i M_{ij} e^{-(\bar{M}\bar{\lambda}_t)_i} \quad \text{"steepest direction"} \end{aligned}$$

"Frechet derivative"
find steepest direction j
 $e_j = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}_j$

Step 2: Line search along direction j_t :

$$O = \frac{\partial R^{\text{train}}(\bar{x}_t + \alpha \bar{e}_{j_t})}{\partial \alpha} \Big|_{\alpha_t} \quad \begin{matrix} \text{how far} \\ \alpha_t \leftarrow \text{to go} \end{matrix}$$

$$= \frac{d}{d\alpha} \frac{1}{n} \sum_i e^{-[M(\bar{x}_t + \alpha \bar{e}_{j_t})]_i} \Big|_{\alpha_t}$$

$$= \frac{1}{n} \sum_i \frac{d}{d\alpha} [e^{-(\bar{M}\bar{x}_t)_i - \alpha M_{ij_t}}] \Big|_{\alpha_t}$$

$$= \frac{1}{n} \sum_i (-M_{ij_t}) e^{-(\bar{M}\bar{x}_t)_i - \alpha M_{ij_t}} \Big|_{\alpha_t}$$

$$= -\frac{1}{n} \sum_{i: M_{ij_t}=1} M_{ij_t} e^{-(\bar{M}\bar{x}_t)_i - \alpha M_{ij_t}} - \frac{1}{n} \sum_{i: M_{ij_t}=-1} M_{ij_t} e^{-(\bar{M}\bar{x}_t)_i - \alpha M_{ij_t}} \Big|_{\alpha_t}$$

$$= -\cancel{\sum_{i: M_{ij_t}=1} e^{-(\bar{M}\bar{x}_t)_i} e^{-\alpha}} - \cancel{\sum_{i: M_{ij_t}=-1} -e^{-(\bar{M}\bar{x}_t)_i} e^{-\alpha}} \Big|_{\alpha_t}$$

$$= -e^{-\alpha} \sum_{i: M_{ij_t}=1} e^{-(\bar{M}\bar{x}_t)_i} + e^{\alpha} \sum_{i: M_{ij_t}=-1} e^{-(\bar{M}\bar{x}_t)_i} \Big|_{\alpha_t}$$

$$= -e^{-\alpha_t} \underbrace{\left[\sum_{i: M_{ij_t}=1} e^{-(\bar{M}\bar{x}_t)_i} \right]}_{Z_t} + e^{\alpha_t} \underbrace{\left[\sum_{i: M_{ij_t}=-1} e^{-(\bar{M}\bar{x}_t)_i} \right]}_{Z_t} \quad \leftarrow \text{normalization}$$

$$\text{define } d_{t,i} = \frac{e^{-(\bar{M}\bar{x}_t)_i}}{Z_t}$$

$$\frac{d_{t,i}}{\prod_{i=1}^n \prod_{i=2}^n} \quad \begin{matrix} \text{---} \\ i=1 \end{matrix} \quad \begin{matrix} \text{---} \\ i=2 \end{matrix} \quad \begin{matrix} \text{---} \\ i=n \end{matrix}$$

$$O = -e^{-\alpha_t} \underbrace{\sum_{i: M_{ij_t}=1} d_{t,i}}_{d_+} + e^{\alpha_t} \underbrace{\sum_{i: M_{ij_t}=-1} d_{t,i}}_{d_-}$$

total weight for correctly classified points

$$O = -e^{-\alpha_t} \underbrace{\sum_{i: M_{ij_t}=1} d_{t,i}}_{d_+} + e^{\alpha_t} \underbrace{\sum_{i: M_{ij_t}=-1} d_{t,i}}_{d_-}$$

total weight for correctly classified points

$$= -e^{-\alpha_t} d_+ + e^{\alpha_t} d_-$$

$$\bar{e}^{\alpha_t} d_+ = e^{\alpha_t} d_-$$

$$\frac{d_+}{d_-} = e^{2\alpha_t}$$

$$\frac{1}{2} \ln \left(\frac{d_+}{d_-} \right) = \alpha_t \quad \text{Since } Z_t \text{ was normalization}$$

$$d_+ = 1 - d_-$$

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1-d_-}{d_-} \right) \quad \therefore$$

Simplify one last thing:

$$\text{Step 1: } j_t \in \operatorname{argmax}_j \frac{1}{n} \sum_i M_{ij} e^{-(\bar{M} \lambda_t)_i};$$

$$j_t \in \operatorname{argmax}_j \sum_i M_{ij} d_{t,i}$$

$$j_t \in \operatorname{argmax}_j (\bar{d}_t^\top \bar{M})_j$$

Finally the algorithm:

$$d_{t,i} = \frac{1}{n} \quad i=1 \dots n$$

$$\bar{\lambda}_t = \bar{0}$$

for $t=1 \dots T$

$$\text{Step 1: } j_t \in \operatorname{argmax}_j (\bar{d}_t^\top \bar{M})_j$$

$$\text{Notation: } d_- = \sum_{M_{ij}=-1} d_{t,i}$$

$$\text{Step 2: } \alpha_t = \frac{1}{2} \ln \left(\frac{1-d_-}{d_-} \right)$$

$$\text{Take the step: } \bar{\lambda}_{t+1} = \bar{\lambda}_t + \alpha_t e_{j_t}$$

$$\begin{bmatrix} 3 \\ 2 \\ 1 \\ -1 \\ -2 \\ -3 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \\ 1 \\ \vdots \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

$$Z_{t+1} = \sum_i e^{-(M\bar{\lambda}_{t+1})_i}$$

$$\text{Notation: } d_{t+1,i} = e^{-(M\bar{\lambda}_{t+1})_i} / Z_{t+1}$$

This \nearrow is AdaBoost. Except for one thing.

$j_t \in \operatorname{argmax}_j (\bar{d}_t^\top \bar{M})_j$ is replaced by a "weak learning algorithm"

$$= \operatorname{argmax}_j \left[\underbrace{\sum_{i: M_{ij}=1} d_{t,i}} + \sum_{i: M_{ij}=-1} -d_{t,i} \right]$$

$$= \operatorname{argmax}_j \left[1 - \sum_{i: M_{ij}=-1} d_{t,i} - \sum_{i: M_{ij}=1} d_{t,i} \right]$$

$$= \operatorname{argmin}_j \left[\sum_{i: M_{ij}=-1} d_{t,i} \right] \quad \begin{array}{l} \text{pick the weak classifier} \\ \text{that minimizes the "weight"} \\ \text{of misclassified points} \end{array}$$

One last bit of notation:

$$\text{weight update: } d_{t+1} = e^{-(M\lambda_{t+1})_i} / Z_{t+1}$$

$$\propto [e^{-(M\lambda_t)_i} e^{-M_{ij_t} \alpha_t}]$$

$$d_{t+1,i} \propto \begin{cases} d_{t,i} e^{-\alpha_t} & \text{if } M_{ij_t} = 1 \\ d_{t,i} e^{\alpha_t} & \text{if } M_{ij_t} = -1 \end{cases}$$

multiplicative weight update

