

Boosting "the statistical view"

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

$$\leq \frac{1}{n} \sum_i e^{-y_i f(x_i)} \quad \text{exponential loss}$$

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 $f(x) = \sum_{j=1}^p \lambda_j x \cdot j$ as before.

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

where $M = \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix}$ $\begin{matrix} p \\ \left[\begin{matrix} y_i h_j(x_i) \end{matrix} \right] \end{matrix}$ "matrix of margins"
margin of i^{th} point for j^{th} weak classifier

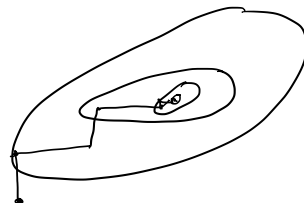
Assume all weak classifiers are binary, so $h_j(x_i) = \pm 1$
Then $M_{ij} = \pm 1$

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

Do coordinate descent for $\bar{\lambda}$

Step 1: "coordinates" are j 's. Find steepest coordinate.

Step 2: do a linesearch in that direction



Step 1. Say we are at $\bar{\lambda}_t$ and want steepest direction.

$$j_t \in \operatorname{argmax}_j \left[\underbrace{- \frac{dR^{\text{train}}(\bar{\lambda} + \alpha \bar{e}_j)}{d\alpha}}_{\alpha=0} \right]$$

"Frechet derivative"
find steepest direction j_t

$$\bar{e}_j = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}_j$$

$$- \frac{d}{d\alpha} \left(\frac{1}{n} \sum_i e^{-[\bar{M}(\bar{\lambda}_t + \alpha \bar{e}_j)]_i} \right) \bigg|_{\alpha=0}$$

$$- \frac{d}{d\alpha} \left(\frac{1}{n} \sum_i e^{-(\bar{M}\bar{\lambda}_t)_i - \alpha M_{ij}} \right) \bigg|_{\alpha=0}$$

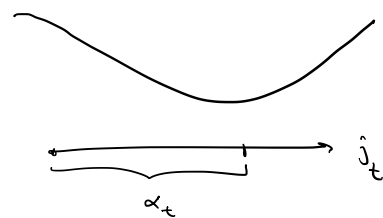
$$\frac{1}{n} \sum_i \left[- \frac{d}{d\alpha} \left(e^{-(\bar{M}\bar{\lambda}_t)_i - \alpha M_{ij}} \right) \bigg|_{\alpha=0} \right]$$

$$\frac{1}{n} \sum_i + M_{ij} e^{-(\bar{M}\bar{\lambda}_t)_i} + 0$$

$$j \in \operatorname{argmax}_j \frac{1}{n} \sum_i M_{ij} e^{-(\bar{M}\bar{\lambda}_t)_i} \quad \text{"steepest direction"}$$

Step 2: Linesearch along direction j

$$0 = \frac{dR^{\text{train}}(\bar{\lambda}_t + \alpha \bar{e}_{j_t})}{d\alpha} \Big|_{\alpha_t} \quad \leftarrow \text{how far to go}$$



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

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

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

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

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

$$\text{define } d_{t,i} = \frac{e^{-(M\lambda)_i}}{Z_t} \quad \leftarrow \text{normalization}$$

$$0 = -\sum_{i: M_{ij_t}=1} d_{t,i} e^{-\alpha} - \sum_{i: M_{ij_t}=-1} -d_{t,i} e^{\alpha} \Big|_{\alpha_t}$$

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

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

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

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

$$\frac{1}{2} \ln\left(\frac{d_+}{d_-}\right) = \alpha_t \quad \leftarrow \text{Since } d\text{'s are normalized } d_+ = 1 - d_-$$

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

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

$$j_t \in \operatorname{argmax}_j \underbrace{\sum_i M_{ij} d_{t,i}}_{(\bar{d}_t^T M)_j}$$

Finally the algorithm:

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

$$\vec{\lambda}_1 = \vec{0}$$

for $t=1 \dots T$

$$\text{Step 1: } j_t \in \arg\max_j (\vec{d}_t^T \vec{M})_j$$

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

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

$$\text{Take the step: } \vec{\lambda}_{t+1} = \vec{\lambda}_t + \underbrace{\alpha_t}_{\text{size}} \underbrace{\vec{e}_{j_t}}_{\text{direction}}$$

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

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

end

This \uparrow is AdaBoost. Except for one thing

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

$$\propto e^{-(\vec{M} \vec{\lambda}_t)_i} e^{-M_{ij_t} \alpha_t}$$

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

Replaced in practice by weak learning algorithm.

$$j_t \in \arg\max_j (\vec{d}_t^T \vec{M})_j$$

$$\arg\max_j \left[\sum_{i: M_{ij}=1} d_{t,i} + \sum_{i: M_{ij}=-1} -d_{t,i} \right]$$

$$\arg\max_j \left[1 - \sum_{i: M_{ij}=1} d_{t,i} - \sum_{i: M_{ij}=-1} d_{t,i} \right]$$

$$\arg\min_j \left[\sum_{i: M_{ij}=-1} d_{t,i} \right] \leftarrow \text{weak classifier minimizes "weight" of misclassified points}$$