

$$\begin{aligned} \max_{\gamma} \quad & \gamma \quad \text{s.t.} \quad y_i \frac{\boldsymbol{\lambda}^T \mathbf{x}_i + \lambda_0}{\|\boldsymbol{\lambda}\|_2} \geq \gamma \quad i = 1 \dots n \\ \min_{\boldsymbol{\lambda}, \lambda_0} \quad & \frac{1}{2} \|\boldsymbol{\lambda}\|_2^2 \quad \text{s.t.} \quad y_i (\boldsymbol{\lambda}^T \mathbf{x}_i + \lambda_0) - 1 \geq 0 \quad i = 1 \dots n \end{aligned}$$

$$\mathcal{L}([\boldsymbol{\lambda}, \lambda_0], \boldsymbol{\alpha}) = \frac{1}{2} \sum_{j=1}^n \lambda_j^2 + \sum_{i=1}^n \alpha_i [-y_i (\boldsymbol{\lambda}^T \mathbf{x}_i + \lambda_0) + 1]$$

$$\begin{aligned} \nabla_{\boldsymbol{\lambda}} \mathcal{L}([\boldsymbol{\lambda}, \lambda_0], \boldsymbol{\alpha}) = \boldsymbol{\lambda} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = \mathbf{0} &\implies \boldsymbol{\lambda} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i. \\ \frac{\partial}{\partial \lambda_0} \mathcal{L}([\boldsymbol{\lambda}, \lambda_0], \boldsymbol{\alpha}) = - \sum_{i=1}^n \alpha_i y_i = 0 &\implies \sum_{i=1}^n \alpha_i y_i = 0. \\ \alpha_i \geq 0 \quad \forall i &\quad \text{(dual feasibility)} \\ \alpha_i [-y_i (\boldsymbol{\lambda}^T \mathbf{x}_i + \lambda_0) + 1] = 0 \quad \forall i &\quad \text{(complementary slackness)} \\ -y_i (\boldsymbol{\lambda}^T \mathbf{x}_i + \lambda_0) + 1 \leq 0. &\quad \text{(primal feasibility)} \end{aligned}$$

$$\mathcal{L}(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,k} \alpha_i \alpha_k y_i y_k \mathbf{x}_i^T \mathbf{x}_k \quad \text{s.t.} \quad \begin{cases} \alpha_i \geq 0 & i = 1 \dots n \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$$

$$\frac{\gamma^0 \gamma^1 \xi}{\prod_{\mathbf{I}} \mathbf{I}} \|\mathbf{y}\|_{\mathbf{I}}^5 + \mathbf{C} \sum_{\mathbf{I}} \xi_{\mathbf{I}}^{\dagger} \quad \mathbf{I}^{\dagger} \mathbf{I}^{\dagger} \quad \begin{cases} \xi_{\mathbf{I}}^{\dagger} \geq 0 \\ \mathbf{I}^{\dagger} (\mathbf{y}_{\mathbf{I}} \mathbf{x}^{\dagger} + \mathbf{y}^0) \geq \mathbf{I} - \xi_{\mathbf{I}}^{\dagger} \end{cases}$$

$$\mathcal{L}(\boldsymbol{\lambda}, \lambda_0, \xi, \boldsymbol{\alpha}, r) = \frac{1}{2} \|\boldsymbol{\lambda}\|_2^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [y_i (\boldsymbol{\lambda}^T \mathbf{x}_i + \lambda_0) - 1 + \xi_i] - \sum_{i=1}^n r_i \xi_i$$

where α_i 's and r_i 's are Lagrange multipliers (constrained to be ≥ 0). The dual turns out to be (after some work)

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,k=1}^n \alpha_i \alpha_k y_i y_k \mathbf{x}_i^T \mathbf{x}_k \quad \text{s.t.} \quad \begin{cases} 0 \leq \alpha_i \leq C & i = 1 \dots n \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases} \quad (6)$$

$$b^* = - \frac{\max_{i: y_i = -1} w^{*T} x_i + \min_{i: y_i = 1} w^{*T} x_i}{2} \quad \text{多项式核 } K(x_1, x_2) = (\langle x_1, x_2 \rangle + R)^d,$$

$$\begin{aligned} \text{高斯核 } K(x_1, x_2) &= \exp(-\|x_1 - x_2\|^2 / 2\sigma^2), \\ \alpha_i = 0 &\Leftrightarrow y_i u_i \geq 1, \quad \begin{cases} L = \max(0, \alpha_i^{old} - \alpha_i^{old}), H = \min(C, C + \alpha_i^{old} - \alpha_i^{old}) & \text{if } y_i \neq y_2 \\ L = \max(0, \alpha_i^{old} + \alpha_i^{old} - C), H = \min(C, \alpha_i^{old} + \alpha_i^{old}) & \text{if } y_i = y_2 \end{cases} \\ 0 < \alpha_i < C &\Leftrightarrow y_i u_i = 1, \\ \alpha_i = C &\Leftrightarrow y_i u_i \leq 1. \end{aligned}$$

令 $E_i = u_i - y_i$ (表示预测值与真实值之差), $\eta = K(x_1, x_1) + K(x_2, x_2) - 2K(x_1, x_2)$, 然后上式两边同时除以 η , 得到一个关于单变量 α_2 的解:

$$\alpha_2^{new, naive} = \alpha_2^{old} + \frac{y_2(E_1 - E_2)}{\eta} \quad (3.44)$$

这个解没有考虑其约束条件 $0 \leq \alpha_2 \leq C$, 即是未经剪辑时的解。然后考虑约束 $0 \leq \alpha_2 \leq C$ 可得到经过剪辑后的 α_2^{new} 的解析解为:

$$\alpha_2^{new, naive} = \begin{cases} H_i & \alpha_2^{new, naive} > H \\ \alpha_2^{new, naive} & L \leq \alpha_2^{new, naive} \leq H \\ L & \alpha_2^{new, naive} < L \end{cases} \quad (3.45)$$

求出了后 α_2^{new} , 便可以求出 α_1^{new} , 得 $\alpha_1^{new} = \alpha_1^{old} + y_1 y_2 (\alpha_2^{old} - \alpha_2^{new})$.

- 对于 α_1 , 即第一个乘子, 可以通过刚刚说的那 3 种不满足 KKT 的条件来找;
- 而对于第二个乘子 α_2 可以寻找满足条件: $\max\{E_1 - E_2\}$ 的乘子。

而 b 在满足下述条件:

$$b = \begin{cases} b_1, & 0 < \alpha_1^{new} < C \\ b_2 & 0 < \alpha_2^{new} < C \\ (b_1 + b_2)/2 & \text{otherwise} \end{cases} \quad (3.46)$$

下更新 b :

$$b_1^{new} = b^{old} - E_1 - y_1 (\alpha_1^{new} - \alpha_1^{old}) K(x_1, x_1) - y_2 (\alpha_2^{new} - \alpha_2^{old}) K(x_1, x_2) \quad (3.47)$$

$$b_2^{new} = b^{old} - E_2 - y_1 (\alpha_1^{new} - \alpha_1^{old}) K(x_1, x_2) - y_2 (\alpha_2^{new} - \alpha_2^{old}) K(x_2, x_2) \quad (3.48)$$

且每次更新完两个乘子的优化后, 都需要再重新计算 b , 及对应的 E_i 值。

$$\begin{aligned} j_t &\in \operatorname{argmax}_j \left[\left. \frac{\partial R^{\text{train}}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_j)}{\partial \alpha} \right|_{\alpha=0} \right] & d_{t,i} &= e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i} / Z_t \text{ where } Z_t = \sum_{i=1}^n e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i} \\ &= \operatorname{argmax}_j \left[\left. \frac{\partial}{\partial \alpha} \left[\frac{1}{n} \sum_{i=1}^n e^{-(\mathbf{M}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_j))_i} \right] \right|_{\alpha=0} \right] & 0 &= \left. \frac{\partial R^{\text{train}}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_{j_t})}{\partial \alpha} \right|_{\alpha_t} \\ &= \operatorname{argmax}_j \left[\left. \frac{\partial}{\partial \alpha} \left[\frac{1}{n} \sum_{i=1}^n e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i - \alpha (\mathbf{M}\mathbf{e}_j)_i} \right] \right|_{\alpha=0} \right] & &= -\frac{1}{n} \sum_{i=1}^n M_{ij_t} e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i - \alpha_t M_{ij_t}} \\ &= \operatorname{argmax}_j \left[\left. \frac{\partial}{\partial \alpha} \left[\frac{1}{n} \sum_{i=1}^n e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i - \alpha M_{ij}} \right] \right|_{\alpha=0} \right] & &= -\frac{1}{n} \sum_{i: M_{ij_t}=1} e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i} e^{-\alpha_t} - \frac{1}{n} \sum_{i: M_{ij_t}=1} -e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i} e^{\alpha_t}. \\ &= \operatorname{argmax}_j \left[\frac{1}{n} \sum_{i=1}^n M_{ij} e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i} \right]. \end{aligned}$$

- (Primal feasibility) $g_i(\mathbf{x}^*) \leq 0, i = 1, ..., m$ and $h_i(\mathbf{x}^*) = 0, i = 1, ..., p$.
- (Dual feasibility) $\alpha_i^* \geq 0, i = 1, \dots, m$.
- (Complementary Slackness) $\alpha_i^* g_i(\mathbf{x}^*) = 0, i = 1, \dots, m$.
- (Lagrangian stationary) $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) = \mathbf{0}$.

$$P(Y=y_i|\boldsymbol{\lambda},\mathbf{x}_i)=\frac{1}{1+e^{-y_i\boldsymbol{\lambda}^T\mathbf{x}_i}}.$$

$$\begin{aligned} \boldsymbol{\lambda}^* &\in \operatorname{argmax}_{\boldsymbol{\lambda}} \log L(\boldsymbol{\lambda}) \\ &= \operatorname{argmax}_{\boldsymbol{\lambda}} \sum_{i=1}^n \log \frac{1}{1+e^{-y_i\boldsymbol{\lambda}^T\mathbf{x}_i}} \\ &= \operatorname{argmin}_{\boldsymbol{\lambda}} \sum_{i=1}^n \log(1+e^{-y_i\boldsymbol{\lambda}^T\mathbf{x}_i}). \end{aligned}$$

$$\begin{aligned} d_{1,i} &= \frac{1}{n} \text{ for all } i \\ d_{t+1,i} &= \frac{d_{t,i}}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } y_i = h_{(t)}(x_i) \text{ (smaller weights for easy examples)} \\ e^{\alpha_t} & \text{if } y_i \neq h_{(t)}(x_i) \text{ (larger weights for hard examples)} \end{cases} \\ &\text{where } Z_t \text{ is a normalization constant for the discrete distribution} \\ &\text{that ensures } \sum_i d_{t+1,i} = 1 \\ &= \frac{d_{t,i}}{Z_t} e^{-y_i \alpha_t h_{(t)}(x_i)} \end{aligned}$$

$$H(x)=\text{sign}\left(\sum_{t=1}^T\alpha_th_{(t)}(x)\right)\qquad\qquad\alpha_t=\frac{1}{2}\ln\left(\frac{1-\epsilon_t}{\epsilon_t}\right).$$

$$\epsilon_t = P_{t \sim \mathbf{d}_t}[h_{(t)}(x_t) \neq y_t] = \sum_i d_{t,i} \mathbf{1}_{[h_{(t)}(x_i) \neq y_i]}$$

Consider the misclassification error:

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

which is upper bounded by the exponential loss:

$$\frac{1}{n} \sum_{i=1}^n e^{-y_i f(x_i)}.$$

Choose f to be some linear combination of weak classifiers,

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

Define an $m \times n$ matrix \mathbf{M} so that $M_{ij} = y_i h_j(x_i)$.

$$\mathbf{M} = \begin{matrix} & \begin{matrix} \text{weak classifiers} \\ j \end{matrix} \\ \begin{matrix} \text{examples} \\ i \end{matrix} & \begin{bmatrix} & \\ & \pm 1 \\ & \end{bmatrix} \end{matrix}$$

So matrix \mathbf{M} encodes all of the training examples and the whole weak learning algorithm. In other words, \mathbf{M} contains all the inputs to AdaBoost. The i^{th} entry in the matrix is 1 whenever weak classifier j correctly classifies example i (Note: we might never write out the whole matrix \mathbf{M} in practice!)

Then

$$y_i f(x_i) = \sum_j \lambda_j y_i h_j(x_i) = \sum_j \lambda_j M_{ij} = (\mathbf{M}\boldsymbol{\lambda})_i.$$

Then the exponential loss is:

$$R^{\text{train}}(\boldsymbol{\lambda}) = \frac{1}{n} \sum_i e^{-y_i f(x_i)} = \frac{1}{n} \sum_i e^{-(\mathbf{M}\boldsymbol{\lambda})_i}. \tag{3}$$

$$\begin{aligned} d_{1,i} &= 1/n \text{ for } i = 1 \dots n \\ \boldsymbol{\lambda}_1 &= \mathbf{0} \\ \text{loop } t &= 1 \dots T \\ j_t &\in \operatorname{argmax}_j (d_t^T \mathbf{M})_j \\ d_- &= \sum_{M_{ij_t}=-1} d_{t,i} \\ \alpha_t &= \frac{1}{2} \ln \left(\frac{1-d_-}{d_-} \right) \\ \boldsymbol{\lambda}_{t+1} &= \boldsymbol{\lambda}_t + \alpha_t \mathbf{e}_{j_t} \\ d_{t+1,i} &= e^{-(\mathbf{M}\boldsymbol{\lambda}_{t+1})_i} / Z_{t+1} \text{ for each } i, \text{ where } Z_{t+1} = \sum_{i=1}^n e^{-(\mathbf{M}\boldsymbol{\lambda}_{t+1})_i} \\ \text{end} \end{aligned}$$

$$\begin{aligned} \mathbf{E}e^{-Yf(x)} &= P(Y=1|x)e^{-f(x)} + P(Y=-1|x)e^{f(x)} \\ 0 = \frac{d\mathbf{E}(e^{-Yf(x)}|x)}{df(x)} &= -P(Y=1|x)e^{-f(x)} + P(Y=-1|x)e^{f(x)} \\ P(Y=1|x)e^{-f(x)} &= P(Y=-1|x)e^{f(x)} \\ \frac{P(Y=1|x)}{P(Y=-1|x)} &= e^{2f(x)} \Rightarrow f(x) = \frac{1}{2} \ln \frac{P(Y=1|x)}{P(Y=-1|x)}. \end{aligned}$$

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{[y_i \neq H(x_i)]} \leq R^{\text{train}}(\boldsymbol{\lambda}_T) \leq e^{-2 \sum_{i=1}^T \gamma_i^2} \leq e^{-2 \gamma_{WLA}^2 T}.$$