

مسئله‌ی نامشاهی بزرگی داده شده را به صورت مسئله با بُند محدود می‌توان بازنویسی کرد

تابع محدب $f(x)$ را می‌توان با piecewise linear convex Function به صورت زیر تقریب زد:

$$f(x) = \max_{j=1 \dots m} \{a_j x + b_j\}$$

m تعداد piece است

$$a_1 \leq \dots \leq a_m$$

جایگزینی $\rightarrow \min_{a_j, b_j} \sum_{i=1}^n \left(\max_{j=1 \dots m} \{a_j x_i + b_j\} - y_i \right)^2$

$$\text{s.t. } a_1 \leq \dots \leq a_m$$

مسئله بهینه‌سازی با بُند محدود به دست آمده را می‌توان توسط QP و یا gradient descent حل کرد.
زیرا تابع هزینه quadratic است و قیدها نیز خطی هستند.
بنابراین به دلیل محدب بودن توابع و محدود بودن n مسئله قابل یادگیری است

می‌دانیم که ریسک هر weak learner به تنهایی باید کمتر از $\frac{1}{2}$ باشد. ریسک h_t را در مرحله $t+1$ محاسبه می‌کنیم:

$$R_{D_{t+1}}(h_t) = \sum_{i=1}^n \underbrace{D^{t+1}(i)}_{\frac{D^t(i) e^{-\alpha_t y_i h_t(x_i)}}{Z_t}} \mathbb{1}[y_i \neq h_t(x_i)] = \sum_{y_i h_t(x_i) < 0} \frac{D^t(i) e^{\alpha_t}}{Z_t}$$

$$= \frac{e^{\alpha_t}}{Z_t} \underbrace{\sum_{i=1}^n D^t(i) \mathbb{1}[y_i \neq h_t(x_i)]}_{\varepsilon_t} \stackrel{\textcircled{I}}{=} \frac{\sqrt{\frac{1}{\varepsilon_t} - 1}}{2\sqrt{\varepsilon_t(1-\varepsilon_t)}} \varepsilon_t = \frac{\sqrt{\varepsilon_t - \varepsilon_t^2}}{2\sqrt{\varepsilon_t - \varepsilon_t^2}} = \frac{1}{2}$$

بنابراین h_t برای مرحله $t+1$ یادگیرنده باید مناسبی نیست و h_{t+1} باید متفاوت از h_t باشد

$$\begin{aligned} Z_t &= \sum_{i=1}^n D^t(i) e^{-\alpha_t y_i h_t(x_i)} \\ &= e^{\alpha_t} \underbrace{\sum_{i \in \text{corr}} D^t(i)}_{\varepsilon_t} + e^{-\alpha_t} \underbrace{\sum_{i \in \text{err}} D^t(i)}_{1-\varepsilon_t} \quad \textcircled{I} \\ &= 2\sqrt{\varepsilon_t(1-\varepsilon_t)} \end{aligned}$$

4

فرض می‌کنیم n یک بردار با d المان غیر صفر است و w یک بردار با k المان غیر صفر است

الف بدون فرض

$$\|x\|_\infty \|w\|_1 \approx \|w\|_1 \approx k$$

if $k < d$ (sparse w) $\rightarrow L1$ is better

$$\|x\|_2 \|w\|_2 \approx \sqrt{d} \|w\|_2 \approx \sqrt{dk}$$

if $k > d$ (dense w) $\rightarrow L2$ is better

در حالت کلی مدنی برای کران نمی‌توان زد ولی $L1$ به طور کلی جواب را sparse تر می‌کند

ب sparse w

$$\|x\|_\infty \|w\|_1 \approx k$$

$k < d \rightarrow L1$ is better

$$\|x\|_2 \|w\|_2 \approx \sqrt{dk}$$

ج $(\|w\|_1 \approx \sqrt{d} \|w\|_2)$ dense w

$$\|x\|_\infty \|w\|_1 \approx \sqrt{d} \|w\|_2$$

$$\|x\|_2 \|w\|_2 \approx \sqrt{d} \|w\|_2 \rightarrow \text{comparable}$$

در این حالت نیز مدنی برای کران نمی‌توان زد و تقریباً برابر هستند

5.a

$$R_\ell(h) = \mathbb{E}_{(n,y) \sim \mathcal{D}} [e^{-y h(n)}] = \mathbb{E}_n \mathbb{E}_{y|n} [e^{-y h(n)}]$$

$$\underline{\underline{y \in \{-1, 1\}}} \mathbb{E}_n [n(n) e^{-h(n)} + (1-n(n)) e^{h(n)}] \geq \mathbb{E}_n [2 \sqrt{n(n)(1-n(n))}]$$

حال اگر نشان دهیم کران پایین $R_\ell(h)$ به ازای h می‌تواند achieve می‌شود، آن h همان h_{exp} خواهد بود

$$\text{if } h(n) = \frac{1}{2} \log \frac{n(n)}{1-n(n)} \rightarrow R_\ell(h) = \mathbb{E}_n \left[\underbrace{n(n) \frac{1}{e^{\frac{1}{2} \log \frac{n(n)}{1-n(n)}}}}_{\sqrt{(1-n(n))n(n)}} + \underbrace{(1-n(n)) e^{\frac{1}{2} \log \frac{n(n)}{1-n(n)}}}_{\sqrt{n(n)(1-n(n))}} \right]$$

$$= \mathbb{E}_n [2 \sqrt{n(n)(1-n(n))}]$$

بنابراین $h_{exp}(n) = \frac{1}{2} \log \frac{n(n)}{1-n(n)}$ پاسخ بهینه می‌باشد

5.6

$$R(h) = \mathbb{E}_{(x,y) \sim D} [\mathbb{1}[\text{sign}(h(x)) \neq y]] = \mathbb{E}_x \mathbb{E}_{y|x} [\mathbb{1}[\text{sign}(h(x)) \neq y]]$$

$$\stackrel{y \in \{-1,1\}}{\geq} \mathbb{E}_x [\eta(x) \mathbb{1}[h(x) < 0] + (1-\eta(x)) \mathbb{1}[h(x) \geq 0]]$$

$$\geq \mathbb{E}_x [\min\{\eta(x), 1-\eta(x)\}]$$

چون $\mathbb{1}$ ها اشتراکی ندارند و در هر حالت فقط یکی از آنها انتخاب می شود
پس اگر \min را برداریم صوابه کوچکتر مساوی می شود

اگر h_{exp} بجای h بگذاریم
تبدیل را در نظر بگیریم $\rightarrow R(h_{exp}) \stackrel{\textcircled{I}}{=} \mathbb{E}_x [\eta(x) \mathbb{1}_{\eta(x) < \frac{1}{2}} + (1-\eta(x)) \mathbb{1}_{\eta(x) \geq \frac{1}{2}}]$

$$\left. \begin{aligned} h_{exp}(x) = \frac{1}{2} \lg \frac{\eta(x)}{1-\eta(x)} < 0 &\rightarrow \eta(x) < \frac{1}{2} \\ \sim \sim \sim \geq 0 &\rightarrow \eta(x) \geq \frac{1}{2} \end{aligned} \right\} \textcircled{I}$$

عبارت داخل \mathbb{E}_x برابر $\min\{\eta(x), 1-\eta(x)\}$ است زیرا
اگر $\eta(x) < \frac{1}{2}$ باشد خودی انتخاب می شود، اگر
 $\eta(x) \geq \frac{1}{2}$ باشد مکمل این انتخاب می شود
بنابراین کمترین مقدار $R(h)$ که آنرا R^* می نامیم
همان $R(h_{exp})$ است

6

می‌دانیم که ریسک هر weak learner به تنهایی باید کمتر از $\frac{1}{2}$ باشد. ریسک h_t را در مرحله $t+1$ می‌سازیم:

$$R_{D_{t+1}}(h_t) = \frac{\sum_{i=1}^n \underbrace{D^{t+1}(i)}_{\frac{D^t(i) e^{-\alpha_t y_i h_t(x_i)}}{Z_t}} \mathbb{1}[y_i \neq h_t(x_i)]}{Z_t} = \sum_{y_i h_t(x_i) < 0} \frac{D^t(i) e^{\alpha_t}}{Z_t}$$

$$= \frac{e^{\alpha_t}}{Z_t} \underbrace{\sum_{i=1}^n D^t(i) \mathbb{1}[y_i \neq h_t(x_i)]}_{\varepsilon_t} \stackrel{\textcircled{I}}{=} \frac{\sqrt{\frac{1}{\varepsilon_t} - 1}}{2\sqrt{\varepsilon_t(1-\varepsilon_t)}} \varepsilon_t = \frac{\sqrt{\varepsilon_t - \varepsilon_t^2}}{2\sqrt{\varepsilon_t - \varepsilon_t^2}} = \frac{1}{2}$$

بنابراین h_t برای مرحله $t+1$ یادگیرنده پایدار مناسبی نیست و باید متفاوت از h_t باشد.

$$\begin{aligned} Z_t &= \sum_{i=1}^n D^t(i) e^{-\alpha_t y_i h_t(x_i)} \\ &= e^{\alpha_t} \underbrace{\sum_{i \in \text{incorr}} D^t(i)}_{\varepsilon_t} + e^{-\alpha_t} \underbrace{\sum_{i \in \text{corr}} D^t(i)}_{1-\varepsilon_t} \quad \textcircled{I} \\ &= 2\sqrt{\varepsilon_t(1-\varepsilon_t)} \end{aligned}$$

در برخی سدها با علی بابایک همفکری داریم