



تمرین چهارم

مسئله ۱. رگرسیون محدب (۱۰ نمره)

فرض کنید می‌خواهیم مسئله رگرسیون به شکل زیر را حل کنیم:

$$\min_{f \in \mathcal{F}} \sum_{i=1}^n (f(x_i) - y_i)^2$$

که در اینجا \mathcal{F} مجموعه همه توابع محدب است. روشی برای حل این مسئله بهینه‌سازی نامتناهی بعدی ارائه دهید. (نکته جالب: برخلاف حالت classification این مسئله قابل یادگیری است)

حل. ابتدا توجه می‌کنیم که برای هر تابع محدب f داریم:

$$f(x) \geq f(x_i) + (x - x_i)^T z_i$$

که در اینجا z_i زیرگرادیان تابع محدب در نقطه‌ی z_i می‌باشد. و حال مسئله متناهی البعد زیر را حل می‌کنیم:

$$\min_{\{\tilde{y}_i, z_i\}_{i=1}^n} (y_i - \tilde{y}_i)^2$$

$$\text{s.t. } \tilde{y}_j \geq \tilde{y}_i + \tilde{z}_i^T (x_j - x_i)$$

و از جواب نهایی برای تعریف تابع نهایی به شکل زیر استفاده می‌کنیم:

$$\hat{f}(x) = \max_{i=1,2,\dots,n} \{\tilde{y}_i + \tilde{z}_i^T (x - x_i)\}$$

که با توجه به این که max تعدادی تابع آفین (محدب) است خود محدب خواهد بود. \triangleright

مسئله ۲. به‌روزرسانی وزن‌ها (۱۵ نمره)

فرض کنید که فرض main weak learner در الگوریتم AdaBoost برقرار باشد. فرض کنید h_t یادگیرنده پایه‌ای باشد که در مرحله t انتخاب شده است. نشان دهید که یادگیرنده پایه h_{t+1} که در مرحله $t+1$ انتخاب شده است باید متفاوت از h_t باشد.

حل. با توجه به فرض یادگیری ضعیف، یک فرضیه $h \in \mathcal{H}$ وجود دارد که خطای D_{t+1} آن کمتر از نصف است. خطای تجربی h_t را برای توزیع D_{t+1} بررسی کنید. از آنجا که $Z_t = 2\sqrt{\epsilon_t(1-\epsilon_t)}$ و $\alpha_t = \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$

$$\begin{aligned}\hat{R}_{\mathcal{D}_{t+1}}(h_t) &= \sum_{i=1}^m \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t} \cdot \mathbb{1}_{[h_t(x_i) y_i < 0]} \\ &= \sum_{h_t(x_i) y_i < 0}^m \frac{D_t(i) \exp(\alpha_t)}{Z_t} \\ &= \frac{e^{\alpha_t}}{Z_t} \sum_{h_t(x_i) y_i < 0}^m D_t(i) \\ &= \frac{\sqrt{\frac{1-\epsilon_t}{\epsilon_t}}}{2\sqrt{\epsilon_t(1-\epsilon_t)}} \cdot \epsilon_t = \frac{1}{2}\end{aligned}$$

این نشان می دهد که h_t نمی تواند در دور $t+1$ انتخاب شود. \triangleright

مسئله ۳. لاسو گروهی با گروه های همپوشان (۱۵ نمره)

نشان دهید اگر از لاسو گروهی برای گروه های همپوشان استفاده کنیم در این که پرامترهای غیر صفر ما در اجتماعی از این گروه ها قرار بگیرند کمکی نمی کند. به طور دقیق تر استفاده از این هموارساز باعث می شود که پرامترها در مکمل اجتماعی از این گروه ها قرار بگیرند. (نیازی به اثبات این مورد نیست، برای حل این مشکل Overlapping group (lasso) طراحی شده است که می توانید در مورد آن بیشتر بخوانید.

حل. فرض کنید که $\theta \in \mathbb{R}^3$ فضای پرامترهای ما باشد. و دو گروه مدنظر ما برابر $\{1, 2\}$ و $\{1, 3\}$ باشند. فرض کنیم که $\theta_1 = \theta_2 = 1$ در این صورت:

$$f(\theta_3) = \sqrt{1 + \theta_3^2} - 1$$

این تابع در نقطه ی صفر مشتق پذیر است و مشتقی برابر صفر دارد، این موضوع باعث می گردد که θ_3 به تنگ بودن تشویق نشود. با عوض کردن جای θ_2, θ_3 می توان حکم مشابه را برای θ_2 نتیجه گرفت و لذا اگر متغیر اول فعال باشد معمولا دو متغیر دیگر نیز فعال خواهند بود. \triangleright

مسئله ۴. در جستجوی هموارساز خوب (۱۰ نمره)

در دو کرانی که برای پیچیدگی راداماخر مدل های خطی به دست آوردیم دو عبارت ظاهر می شد:

(الف) برای نرم ۱ عبارتی که ظاهر می شد برابر $\|w\|_1 \|x\|_\infty$ بود (در واقع کران هایی روی این ها ولی برای سادگی فعلا این عبارت ها را در نظر می گیریم.)

(ب) برای نرم ۲ عبارت $\|w\|_2 \|x\|_2$ ظاهر می شد.

فرض کنید w, x متغیرهای تصادفی نزدیک مقادیر $\{-1, 1\}$ باشند. در حالت های زیر این دو کران را بررسی کنید و نتیجه بگیرید هموارسازی با نرم ۱ و ۲ چه زمان هایی مناسب تر هستند.

الف) بدون هیچ فرض اضافه ای.

ب) w تنک با حداکثر k مولفه ناصفر باشد.

ج) w چگال باشد مثلاً با فرض $\|w\|_1 \approx \sqrt{d}\|w\|_2$

حل.

الف) بدون هیچ فرض اضافه ای عبارت شامل نرم ۲ تقریباً برابر $\sqrt{d}\sqrt{d}$ خواهد بود عبارت شامل نرم ۱ برابر d و لذا تفاوتی بین دو عبارت نخواهد بود.

ب) در این حالت عبارت شامل نرم ۱ برابر k خواهد بود و عبارت شامل نرم ۲ برابر $\sqrt{k}\sqrt{d}$ و لذا استفاده از نرم ۱ در این گونه مسائل مفید خواهد بود.

ج) در این حالت خواهیم داشت:

$$\|w\|_2\|x\|_2 \leq \frac{1}{\sqrt{d}}\|w\|_1\sqrt{d}\|x\|_\infty = \|w\|_1\|x\|_\infty$$

و لذا در این حالت بهتر است از نرم ۲ به عنوان هموارساز استفاده کنیم.

▷

مسئله ۵. خطای نمایی Bayes (۱۵ نمره)

فرض کنید مجموعه ورودی \mathcal{X} و فضای برچسب $Y = \{-1, +1\}$ باشد. در الگوریتم AdaBoost از تابع خطای نمایی زیر استفاده می شود:

$$\ell(h(x), y) = \exp(-yh(x)).$$

برای یک توزیع \mathcal{D} روی $\mathcal{X} \times Y$ ، خطای نمایی به صورت زیر تعریف می شود:

$$R_l(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(h(x), y)].$$

خطای نمایی بیز برای یک توزیع \mathcal{D} به صورت خطای کمینه روی توابع قابل اندازه گیری تعریف می شود:

$$R_l^* = \inf_{h: \mathcal{X} \rightarrow \mathbb{R} \text{ measurable}} R_l(h).$$

فرضیه ای h_{exp} با $R_l(h_{\text{exp}}) = R_l^*$ به عنوان راه حل بهینه بیز نامیده می شود. $\eta(x) = P[y = +1|x]$ را به صورت $\eta(x) = P[y = +1|x]$ تعریف کنید.

(a) عبارت رابطه راه حل بهینه بیز h_{exp} را برای خطای نمایی بر حسب $\eta(x)$ بدهید.

(b) خطای generalization و خطای بیز برای طبقه بندی باینری را به صورت زیر تعریف کنید:

$$R(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathbb{I}[\text{sign}(h(x)) \neq y]], \quad R^* = \inf_{h: \mathcal{X} \rightarrow \mathcal{Y} \text{ measurable}} R(h).$$

$$\text{sign}(t) = \mathbb{1}_{t \geq 0} - \mathbb{1}_{t < 0} \text{ که}$$

$$R(h_{\text{exp}}) = R^* \text{ نشان دهید}$$

حل. الف)

بر اساس تعریف، $R_\ell(h)$ به صورت زیر بیان می شود:

$$\begin{aligned} R_\ell(h) &= \mathbb{E}_{(x,y) \sim \mathcal{D}} [\exp(-yh(x))] \\ &= \mathbb{E}_x \mathbb{E}_{y|x} [\exp(-yh(x))] \\ &= \mathbb{E}_x [\eta(x) \exp(-h(x)) + (1 - \eta(x)) \exp(h(x))] \\ &\geq \mathbb{E}_x \left[\sqrt{\eta(x)(1 - \eta(x))} \right], \end{aligned}$$

که برابری برقرار است اگر و تنها اگر برای هر $x \in X$ ، داشته باشیم:

$$h(x) = \frac{1}{2} \log \left(\frac{\eta(x)}{1 - \eta(x)} \right).$$

بنابراین،

$$R_\ell^* = \mathbb{E}_x \left[\sqrt{\eta(x)(1 - \eta(x))} \right]$$

خطای بیز برای خطای نمایی است و

$$h_{\text{exp}} : x \mapsto \frac{1}{2} \log \left(\frac{\eta(x)}{1 - \eta(x)} \right)$$

جواب بهینه بیز است.

ب)

بر اساس تعریف، $R(h)$ به صورت زیر بیان می شود:

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

$$R^* = \inf_{h: X \rightarrow \mathbb{R}} R(h),$$

که در آن $\text{sign}(t) = 1_{t \geq 0} - 1_{t < 0}$. اثبات کنید که $R(h_{\text{exp}}) = R^*$

بر اساس تعریف، $R(h)$ به صورت زیر بیان می‌شود:

$$\begin{aligned} R(h) &= \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathbb{1}_{\text{sign}(h(x)) \neq y}] \\ &= \mathbb{E}_x \mathbb{E}_{y|x} [\mathbb{1}_{\text{sign}(h(x)) \neq y}] \\ &= \mathbb{E}_x [\eta(x) \mathbb{1}_{h(x) < \cdot} + (1 - \eta(x)) \mathbb{1}_{h(x) \geq \cdot}] \\ &\geq \mathbb{E}_x [\min\{\eta(x), 1 - \eta(x)\}], \end{aligned}$$

که برابری برقرار است اگر و تنها اگر برای هر $x \in X$ ، $\text{sign}(h(x)) = \text{sign}(\eta(x) - 1/2)$ باشد. از آنجا که برای هر $x \in X$ ، $\text{sign}(h_{\text{exp}}(x)) = \text{sign}(\eta(x) - 1/2)$ ، ثابت شد که $R(h_{\text{exp}}) = R^*$.

▷

مسئله ۶. خطای طبقه‌بند ضعیف (۱۰ نمره)

نشان دهید خطای h_t با توجه به توزیع $D^{(t+1)}$ دقیقاً برابر با $1/2$ است. به عبارت دیگر، نشان دهید که برای هر $t \in [T]$ داریم:

$$\sum_{i=1}^m D_i^{(t+1)} \mathbb{I}[y_i \neq h_t(x_i)] = \frac{1}{2}.$$

حل. طبق تعریف ϵ_t :

$$\sum_{i=1}^m D_i^{(t)} \exp(-w_t y_i h_t(x_i)) \cdot \mathbb{1}_{[h_t(x_i) \neq y_i]} = \epsilon_t \cdot \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} = \sqrt{\epsilon_t(1 - \epsilon_t)}. \quad (1)$$

به طور مشابه،

$$\sum_{j=1}^m D_j^{(t)} \cdot \exp(-w_t y_j h_t(x_j)) = \epsilon_t \cdot \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} + (1 - \epsilon_t) \cdot \sqrt{\frac{\epsilon_t}{1 - \epsilon_t}} = 2\sqrt{\epsilon_t(1 - \epsilon_t)}. \quad (2)$$

برابری مورد نظر با مشاهده این موضوع به دست می‌آید که خطای h_t نسبت به توزیع D_{t+1} با تقسیم معادله (۱) بر معادله (۲) قابل محاسبه است.

▷