

Math 281C Homework 5 Solutions

1. Consider the pair $\mathbf{z} = (\mathbf{x}, y) \in \mathbb{R}^d \times \{-1, 1\}$. Recall from Math 281A that the logistic loss is

$$m_{\boldsymbol{\theta}}(\mathbf{z}) = \log(1 + \exp(-y \cdot \langle \mathbf{x}, \boldsymbol{\theta} \rangle)),$$

and the population expectation is $M(\boldsymbol{\theta}) = \mathbb{E}[m_{\boldsymbol{\theta}}(\mathbf{X}, Y)]$, for $(\mathbf{X}, Y) \sim P$.

(a) Show that if $\Theta \in \mathbb{R}^d$ is a compact set and $\mathbb{E}[\|\mathbf{X}\|] < \infty$ for some norm $\|\cdot\|$ on \mathbb{R}^d , then

$$\sup_{\boldsymbol{\theta} \in \Theta} |P_n m_{\boldsymbol{\theta}}(\mathbf{X}, Y) - M(\boldsymbol{\theta})| \xrightarrow{P} 0.$$

Solution: First we show that $m_{\boldsymbol{\theta}}(\mathbf{z})$ is $\|\mathbf{x}\|_*$ -Lipschitz in $\boldsymbol{\theta}$. For any $\boldsymbol{\theta}$ and $\boldsymbol{\theta}'$ in Θ ,

$$|m_{\boldsymbol{\theta}}(\mathbf{x}, y) - m_{\boldsymbol{\theta}'}(\mathbf{x}, y)| \leq |(\boldsymbol{\theta} - \boldsymbol{\theta}') \cdot \mathbf{x}| \leq \|\boldsymbol{\theta} - \boldsymbol{\theta}'\| \cdot \|\mathbf{x}\|_*.$$

Then consider a ϵ -net $\{\boldsymbol{\theta}^i\}_{i=1}^N$ for Θ , for each $\boldsymbol{\theta}^i$ in the ϵ -net, construct function pairs $\{\ell_i, u_i\}$ with

$$\ell_i = m_{\boldsymbol{\theta}^i}(\mathbf{x}, y) - \epsilon \|\mathbf{x}\|_* \quad \text{and} \quad u_i = m_{\boldsymbol{\theta}^i}(\mathbf{x}, y) + \epsilon \|\mathbf{x}\|_*,$$

such that for any $m_{\boldsymbol{\theta}}(\mathbf{x}, y)$, we can find a pair $\{\ell_i, u_i\}$ satisfying $\ell_i(\mathbf{x}) \leq m_{\boldsymbol{\theta}}(\mathbf{x}, y) \leq u_i(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}^d$, and

$$\mathbb{E}[u_i(\mathbf{x}) - \ell_i(\mathbf{x})] \leq 2\epsilon \mathbb{E}\|\mathbf{x}\|_*.$$

This means $\{\ell_i, u_i\}_{i=1}^N$ form a $2\epsilon \mathbb{E}\|\mathbf{x}\|_*$ -bracketing of $\{m_{\boldsymbol{\theta}}(\cdot) | \boldsymbol{\theta} \in \Theta\}$ with respect to ℓ_1 -norm, and by construction, we have

$$N_{[]}(\{m_{\boldsymbol{\theta}}\}, \ell_1, 2\epsilon \mathbb{E}\|\mathbf{x}\|_*) \leq N(\Theta, \|\cdot\|, \epsilon) < \infty,$$

which implies the uniform consistency.

(b) Assume that Θ is contained in the norm ball $\{\boldsymbol{\theta} \in \mathbb{R}^d : \|\boldsymbol{\theta}\| \leq r\}$ and that \mathbf{X} is supported on the dual norm ball $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_* \leq M\}$. Show that there is a constant $C < \infty$ such that for all $0 < \delta < 1$,

$$\mathbb{P}\left(\sup_{\boldsymbol{\theta} \in \Theta} |P_n m_{\boldsymbol{\theta}}(\mathbf{X}, Y) - M(\boldsymbol{\theta})| \geq \epsilon_n(\delta)\right) \leq \delta,$$

where

$$\epsilon_n(\delta) = C \sqrt{\frac{r^2 M^2}{n} \left(d \log n + \log \frac{1}{\delta}\right)}.$$

Solution: First we have

$$\log(1 + \exp(-Mr)) \leq m_{\boldsymbol{\theta}}(\mathbf{x}, y) \leq \log(1 + \exp(Mr)),$$

which means $m_{\boldsymbol{\theta}}(\mathbf{x}, y) - \log(2) \in [-Mr, Mr]$. Thus, for any fixed $\boldsymbol{\theta} \in \Theta$ and $t > 0$, applying Hoeffding's inequality gives

$$\mathbb{P}(|P_n m_{\boldsymbol{\theta}}(\mathbf{X}, Y) - M(\boldsymbol{\theta})| \geq t) \leq 2 \exp\left(-\frac{nt^2}{2M^2 r^2}\right). \quad (1)$$

Then consider a minimal ϵ -net $\{\boldsymbol{\theta}^i\}_{i=1}^N$ for Θ satisfying $N \leq (1 + 2r/\epsilon)^d$. For any $\boldsymbol{\theta} \in \Theta$, there is $\boldsymbol{\theta}^i$ such that $\|\boldsymbol{\theta} - \boldsymbol{\theta}^i\| \leq \epsilon$, and

$$\begin{aligned} \sup_{\boldsymbol{\theta} \in \Theta} |P_n m_{\boldsymbol{\theta}} - M(\boldsymbol{\theta})| &\leq \sup_{\boldsymbol{\theta} \in \Theta} |P_n m_{\boldsymbol{\theta}} - P_n m_{\boldsymbol{\theta}^i}| + \max_{i=1, \dots, N} |P_n m_{\boldsymbol{\theta}^i} - P m_{\boldsymbol{\theta}^i}| + \sup_{\boldsymbol{\theta} \in \Theta} |P m_{\boldsymbol{\theta}^i} - P m_{\boldsymbol{\theta}}| \\ &\leq 2M\epsilon + \max_{i=1, \dots, N} |P_n m_{\boldsymbol{\theta}^i} - P m_{\boldsymbol{\theta}^i}|. \end{aligned}$$

Therefore, combining (1) with union bound, we have

$$\begin{aligned} \mathbb{P}\left(\sup_{\theta \in \Theta} |P_n m_\theta - M(\theta)| \geq 2M\epsilon + t\right) &\leq \mathbb{P}\left(\max_{i=1, \dots, N} |P_n m_{\theta^i} - P m_{\theta^i}| \geq t\right) \\ &\leq \left(1 + \frac{2r}{\epsilon}\right)^d 2 \exp\left(-\frac{nt^2}{2M^2 r^2}\right). \end{aligned}$$

Finally, choosing

$$t = \sqrt{\frac{2M^2 r^2 (d \log(1 + 2r/\epsilon) + \log(2/\delta))}{n}}$$

and $\epsilon = r/n$ gives the desired bound.

2. Consider a binary classification problem with data in pair $(\mathbf{x}, y) \in \mathbb{R}^d \times \{-1, 1\}$, and let $\phi: \mathbb{R} \rightarrow \mathbb{R}_+$ be a 1-Lipschitz non-increasing convex function, for example, $\phi(t) = \log(1 + e^{-t})$ or $\phi(t) = [1 - t]_+$. Define $m_\theta(\mathbf{x}, y) = \phi(y \cdot \langle \mathbf{x}, \theta \rangle)$. Given an i.i.d. sample $\{\mathbf{X}_i, Y_i\}_{i=1}^n$ and consider the empirical risk minimization procedure

$$\widehat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} P_n m_\theta = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n m_\theta(\mathbf{X}_i, Y_i). \quad (2)$$

Ledoux-Talagrand contraction inequality may be useful. Let $\phi \circ \mathcal{F} = \{h: h(x) = \phi(f(x)), f \in \mathcal{F}\}$ denote the composition of $\phi(\cdot)$ with functions in \mathcal{F} . If $\phi(\cdot)$ is L -Lipschitz, then $\mathcal{R}_n(\phi \circ \mathcal{F}) \leq L \mathcal{R}_n(\mathcal{F})$.

- (a) In one word, is the procedure (2) likely to give a reasonably good classifier?

Solution: Yes. Intuitively, we will pick $\widehat{\theta}_n$ such that $Y_i \cdot \langle \mathbf{X}_i, \widehat{\theta}_n \rangle > 0$ for most cases.

Before we proceed to parts (b) and (c), let us prove some general results. Let $\|\cdot\|$ be an arbitrary norm and define $\mathcal{F} = \{f(x) = \langle \mathbf{x}, \theta \rangle \mid \|\theta\| \leq r\}$, then

$$\mathcal{R}_n(\mathcal{F}) = \frac{1}{n} \mathbb{E} \left[\sup_{\|\theta\| \leq r} \sum_{i=1}^n \epsilon_i \langle \mathbf{X}_i, \theta \rangle \right] \leq \frac{r}{n} \mathbb{E} \left[\left\| \sum_{i=1}^n \epsilon_i \mathbf{X}_i \right\|_* \right]. \quad (3)$$

Moreover, suppose a function class \mathcal{F} is b -uniformly bounded, then for any $t > 0$, we have

$$\mathbb{P}(\|P_n - P\|_{\mathcal{F}} \geq 2\mathcal{R}_n(\mathcal{F}) + t) \leq \exp\left(-\frac{nt^2}{2b^2}\right). \quad (4)$$

To prove (4), first we show concentration around mean. If we define

$$G(x_1, \dots, x_n) := \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(x_i) - \mathbb{E}f(X)) \right|,$$

then it can be checked that

$$|G(x_1, \dots, x_i, \dots, x_n) - G(x_1, \dots, x'_i, \dots, x_n)| \leq \frac{2b}{n}.$$

Therefore, by bounded difference inequality, we have

$$\mathbb{P}(\|P_n - P\|_{\mathcal{F}} \geq \mathbb{E}\|P_n - P\|_{\mathcal{F}} + t) \leq \exp\left(-\frac{nt^2}{2b^2}\right).$$

Furthermore, by Theorem 1.1 of Lecture 5, $\mathbb{E}\|P_n - P\|_{\mathcal{F}} \leq 2\mathcal{R}_n(\mathcal{F})$. Combining this with the above display completes the proof of (4).

- (b) Let $\Theta \subset \{\theta \in \mathbb{R}^d: \|\theta\|_2 \leq r\}$ and let $\{\mathbf{X}_i\}_{i=1}^n$ be supported on the ℓ_2 -ball $\{\mathbf{x} \in \mathbb{R}^d: \|\mathbf{x}\|_2 \leq M\}$. Give the smallest $\epsilon_n(\delta, d, r, M)$ you can (ignoring the constants) such that

$$\mathbb{P}\left(\sup_{\theta \in \Theta} |P_n m_\theta - P m_\theta| \geq \epsilon_n(\delta, d, r, M)\right) \leq \delta.$$

How does your ϵ_n compare with Question 1?

Solution: The dual norm of ℓ_2 -norm is ℓ_2 -norm. Consider $\mathcal{F} = \{f(x) = \langle \mathbf{x}, \boldsymbol{\theta} \rangle | \boldsymbol{\theta} \in \Theta\}$. By the independence of Rademacher variables and Jensen's inequality,

$$\mathbb{E} \left[\left\| \sum_{i=1}^n \epsilon_i \mathbf{X}_i \right\|_2 \right] \leq \sqrt{\mathbb{E} \left[\left\| \sum_{i=1}^n \epsilon_i \mathbf{X}_i \right\|_2^2 \right]} = \sqrt{\mathbb{E} \left[\sum_{i=1}^n \|\mathbf{X}_i\|_2^2 \right]} \leq \sqrt{n}M,$$

which, together with (3), implies

$$\mathcal{R}_n(\mathcal{F}) \leq \frac{Mr}{\sqrt{n}}.$$

Applying Ledoux-Talagrand contraction inequality yields

$$\mathcal{R}_n(\{m_{\boldsymbol{\theta}}\}_{\boldsymbol{\theta} \in \Theta}) \leq \frac{Mr}{\sqrt{n}}.$$

Now, by the Lipschitz continuity of $\phi(\cdot)$,

$$\sup_{\boldsymbol{\theta} \in \Theta, \|\mathbf{x}\|_2 \leq M, y \in \{-1, 1\}} |\phi(y \cdot \langle \mathbf{x}, \boldsymbol{\theta} \rangle) - \phi(0)| \leq Mr,$$

and it follows by (4) that for any $t > 0$,

$$\mathbb{P} \left(\sup_{\boldsymbol{\theta} \in \Theta} |P_n m_{\boldsymbol{\theta}} - P m_{\boldsymbol{\theta}}| \geq \frac{2Mr}{\sqrt{n}} + t \right) \leq \exp \left(-\frac{nt^2}{2M^2 r^2} \right).$$

Finally, setting

$$t = \sqrt{\frac{2M^2 r^2 \log(1/\delta)}{n}}$$

gives

$$\epsilon_n = \frac{Mr}{\sqrt{n}} (2 + \sqrt{2 \log(1/\delta)}).$$

This bound is apparently sharper than the bound in Question 1, since it's independent of d .

- (c) Let $\Theta \subset \{\boldsymbol{\theta} \in \mathbb{R}^d : \|\boldsymbol{\theta}\|_1 \leq r\}$ and let $\{\mathbf{X}_i\}_{i=1}^n$ be supported on the ℓ_∞ -ball $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_\infty \leq M\}$. Give the smallest $\epsilon_n(\delta, d, r, M)$ you can (ignoring the constants) such that

$$\mathbb{P} \left(\sup_{\boldsymbol{\theta} \in \Theta} |P_n m_{\boldsymbol{\theta}} - P m_{\boldsymbol{\theta}}| \geq \epsilon_n(\delta, d, r, M) \right) \leq \delta.$$

How does your ϵ_n compare with Question 1?

Solution: Let $\{X_i\}_{i=1}^n$ be a sequence of centered sub-Gaussian random variables with parameter σ , then

$$\mathbb{E} \left[\max_{i=1, \dots, n} |X_i| \right] \leq \sigma \sqrt{2 \log(2n)} \quad (5)$$

To prove (5), by the definition of sub-Gaussian and Jensen's inequality, for any $\lambda > 0$,

$$\begin{aligned} \mathbb{E} \left[\max_{i=1, \dots, n} |X_i| \right] &= \frac{1}{\lambda} \mathbb{E} \log \exp \left(\lambda \max_{i=1, \dots, n} |X_i| \right) \\ &\leq \frac{1}{\lambda} \log \mathbb{E} \exp \left(\lambda \max_{i=1, \dots, n} |X_i| \right) \\ &\leq \frac{1}{\lambda} \log \left[2n \exp \left(\frac{\lambda^2 \sigma^2}{2} \right) \right] \\ &= \frac{\log 2n}{\lambda} + \frac{\lambda \sigma^2}{2}. \end{aligned}$$

Taking $\lambda = \sqrt{2 \log(2n)}/\sigma$ completes the proof of (5).

Now notice that the dual norm of ℓ_1 -norm is ℓ_∞ -norm, and the j -th coordinate of $\sum_{i=1}^n \epsilon_i \mathbf{X}_i$ is sub-Gaussian with parameter $\sqrt{\sum_{i=1}^n \mathbf{X}_{ij}^2} \leq M\sqrt{n}$. By (5), we have

$$\mathbb{E} \left[\left\| \sum_{i=1}^n \epsilon_i \mathbf{X}_i \right\|_\infty \right] \leq M \sqrt{2n \log(2d)}.$$

Applying this to (3) gives

$$R_n(\mathcal{F}) \leq \frac{Mr\sqrt{2\log(2d)}}{\sqrt{n}}.$$

The remaining arguments are the same as part (b), and we can achieve

$$\epsilon_n = \frac{Mr}{\sqrt{n}}(2\sqrt{2\log(2d)} + \sqrt{2\log(1/\delta)}).$$

This bound is also sharper than the one in Question 1.