# 1 The minimum $\ell_1$ - norm interpolator

$\{(x_i, y_i)\}_{i=1}^n$  $x_i \in \mathbb{R}^d$, $y_i = w_*^T x_i \in \mathbb{R}$

Consider  $\min \|w\|_1$  s.t. $w^T x_i = y_i$, $i = 1, \dots n$   solution $\hat{w}_n$

$$\hat{R}_n(w) = \frac{1}{n} \sum_{i=1}^n |w^T x_i - y_i|^2 , \quad R(w) = \mathbb{E}_{x,y}\left[(w^T x - y)^2\right]$$

$\|x\|_\infty \leq 1$

(a) $L_B = \{|w^T x - w_*^T x|^2 : \|w\|_1 \leq B\}$, $H_B = \{w^T x : \|w\|_1 \leq B\}$

Prove $\widehat{Radn}(L_B) \lesssim (B + \|w_*\|_1) \widehat{Radn}(H_B)$

**Proof.** $\widehat{Radn}(L_B) = \mathbb{E}_\xi \sup_{\|w\|_1 \leq B} \frac{1}{n} \sum_{i=1}^n \xi_i \left[(w^T - w_*^T)x_i\right]^2$

$$\leq 2(B + \|w_*\|_1) \widehat{Radn}(H_B)$$

As $x^2$ is $2\beta$-Lipschitz continuous in $[0, \beta]$

and $|(w^T - w_*^T)x_i| \leq B + \|w_*\|_1$  #

(b) $\forall \delta \in (0,1)$, with probability $\geq 1 - \delta$ over sampling of $x_1, \dots x_n$,

$$R(\hat{w}_n) \lesssim \|w_*\|_1^2 \left( \sqrt{\frac{\log(2d)}{n}} + \sqrt{\frac{\log(4/\delta)}{n}} \right)$$

**Proof.** $\widehat{Radn}(L_B) \lesssim (B + \|w_*\|_1) \widehat{Radn}(H_B)$

$$\lesssim (B + \|w_*\|_1) B \sqrt{\frac{\log(2d)}{n}}$$

According to Generalization error based on Rademacher complexity,

$$\sup_w |R(w) - \hat{R}(w)| \leq 2 \widehat{Radn}(L_B) + 4B \sqrt{\frac{2\log(4/\delta)}{n}}$$

$$\lesssim \|w_*\|_1^2 \left( \sqrt{\frac{\log(2d)}{n}} + \sqrt{\frac{\log(4/\delta)}{n}} \right)$$

Notice that $\hat{R}(\hat{w}_n) = 0$, we are done. #

# 2 The reproducing kernel property

$k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ be a SPD kernel. $\mathcal{H}_k$: associated RKHS.

(a) if $k(x,x) \leq C$ for all $x \in \mathcal{X}$, then $|f(x)| \leq \sqrt{C}$ for all $f$ in unit ball of $\mathcal{H}_k$.

**Proof.** We know $k(x,x') = \langle K_x, K_{x'} \rangle_{\mathcal{H}_k}$. For any $\|f\|_{\mathcal{H}_k} \leq 1$,

$$|f(x)| = |\langle f, K_x \rangle|$$

$$\leq \|f\|_{\mathcal{H}_k} \|K_x\|_{\mathcal{H}_k} \quad (\text{Cauchy-Schwarz})$$

But $\|K_x\|_{\mathcal{H}_k}^2 = \langle K_x, K_x \rangle_{\mathcal{H}_k} = k(x,x) \leq C$

So we have $|f(x)| \leq \sqrt{C}$. #


(b) MMD  P, Q pro. distribution over $\mathcal{X}$

$$MMD(P, Q) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim P}[f(x)] - \mathbb{E}_{x \sim Q}[f(x)]$$

Show $MMD^2(P,Q) = \mathbb{E}_{x,x' \sim P}[k(x,x')] + \mathbb{E}_{z,z' \sim Q}[k(z,z')] - 2\mathbb{E}_{x \sim P, z \sim Q}[k(x,z)]$

**Proof.** Let $k(x,y) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}_k}$. (e.g. $\varphi = K_x$)

$$MMD(P, Q) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim P}[f(x)] - \mathbb{E}_{y \sim Q}[f(y)]$$

$$= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim P} \langle f, \varphi(x) \rangle_{\mathcal{H}_k} - \mathbb{E}_{y \sim Q} \langle f, \varphi(y) \rangle_{\mathcal{H}_k}$$

$$= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \langle f, \mathbb{E}_{x \sim P}[\varphi(x)] \rangle_{\mathcal{H}_k} - \langle f, \mathbb{E}_{y \sim Q}[\varphi(y)] \rangle_{\mathcal{H}_k}$$

$$= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \langle f, \mathbb{E}_{x \sim P}\varphi(x) - \mathbb{E}_{y \sim Q}\varphi(y) \rangle$$

$$= \|\mathbb{E}_{x \sim P}\varphi(x) - \mathbb{E}_{y \sim Q}\varphi(y)\|_{\mathcal{H}_k}$$

$$MMD^2(P,Q) = \langle \mathbb{E}_{x \sim P}\varphi(x), \mathbb{E}_{x' \sim P}\varphi(x') \rangle_{\mathcal{H}_k} + \langle \mathbb{E}_{y \sim Q}\varphi(y), \mathbb{E}_{y' \sim Q}\varphi(y') \rangle_{\mathcal{H}_k}$$
$$- 2 \langle \mathbb{E}_{x \sim P}\varphi(x), \mathbb{E}_{y \sim Q}\varphi(y) \rangle_{\mathcal{H}_k}$$

$$= \mathbb{E}_{x,x' \sim P}[k(x,x')] + \mathbb{E}_{z,z' \sim Q}[k(z,z')] - 2\mathbb{E}_{x \sim P, z \sim Q}[k(x,z)]$$ #

# 3. $L^\infty$ approximation of two-layer NN

$S^{d-1} = \{x \in \mathbb{R}^d : \|x\|_2 = 1\}$  $f \in C(S^{d-1})$ target function

$\exists$ prob. distribution $\pi$ s.t. $f(x) = \mathbb{E}_{(a,b) \sim \pi}[a \sigma(b^T x)]$, $b \in S^{d-1}$, $|a| \leq 1$ a.s.

$\sigma = \text{ReLu}$.

(a) $h_x : [-1,1] \times S^{d-1} \mapsto \mathbb{R}$, $h_x(a,b) = a\sigma(b^T x)$  $\mathcal{H} = \{h_x : \|x\|_2 \leq 1\}$

Prove  $\widehat{Rad}_m(\mathcal{H}) \leq \frac{2}{\sqrt{m}}$

**Proof.** $\widehat{Rad}_m(\mathcal{H}) = \mathbb{E}_\xi \sup_{h_x \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \xi_i a_i \sigma(b_i^T x_i) = \frac{1}{m} \mathbb{E}_\xi \sup_{h_x \in \mathcal{H}} \sum_{i=1}^m \text{sign}(a) (|a| b_i^T x)$

$= \frac{1}{m} \mathbb{E}_\xi \sup_{h_x \in \mathcal{H}} \sum_{i=1}^m \xi_i' \sigma(|a_i| b_i^T x) \leq \beta \widehat{Rad}_m(\mathcal{H})$  by concentration lemma.

where  $\mathcal{H} = \{w^T x : \|w\|_2 \leq 1\}$  and $\beta = 2$  (Lip Const for ReLu for $-1 \leq x \leq 1$)

We also have  $\widehat{Rad}_m(\mathcal{H}) \leq \sqrt{\frac{1}{m}}$  (Linear class)

So  $\widehat{Rad}_m(\mathcal{H}) \leq \frac{2}{\sqrt{m}}$. #

(b) Let $(a_i, b_i) \overset{iid}{\sim} \pi$  $\forall \delta \in (0,1)$, with probability $1-\delta$,

$\sup_{x \in S^{d-1}} \left| \frac{1}{m} \sum_{i=1}^m a_i \sigma(b_i^T x) - f(x) \right| \lesssim \frac{1}{\sqrt{m}} + \sqrt{\frac{\log(4/\delta)}{m}}$

**Proof.**  $0 \leq f \leq 1$ a.s.

From Generalization error based on Rademacher complexity,

$\sup_{x \in S^{d-1}} \left| \frac{1}{m} \sum_{i=1}^m a_i \sigma(b_i^T x) - f(x) \right|$

$\leq 2 \widehat{Rad}_m(\mathcal{H}) + 4\sqrt{\frac{2\log(4/\delta)}{m}}$

$\lesssim \frac{1}{\sqrt{m}} + \sqrt{\frac{\log(4/\delta)}{m}}$     #

# 4 Margin-based bounds for classification

$S = \{(x_i, y_i)\}, \quad y_i^2 = 1 \qquad \hat{n}_\gamma(f) = |\{i \in [n]: f(x_i)y_i < \gamma\}|$

$$\ell_{0-1}(t) = \begin{cases} 1 & \text{if } t < 0 \\ 0 & \text{else} \end{cases} \qquad \ell_\gamma(t) = \begin{cases} 1 & \text{if } t < 0 \\ 1 - t/\gamma & \text{if } 0 \leq t < \gamma \\ 0 & \text{if } t \geq \gamma \end{cases}$$



$$R_{0-1}(f) = \mathbb{E}_{x,y}[\ell_{0-1}(f(x)y)]$$

(a) $R_\gamma(f) = \mathbb{E}_{x,y}[\ell_\gamma(f(x)y)] \quad \hat{R}_\gamma(f) = \frac{1}{n}\sum_{i=1}^{n}\ell_\gamma(f(x_i)y_i)$

$\forall f \in \mathcal{F}, \; \gamma > 0, \quad R_{0-1}(f) \leq R_\gamma(f), \quad \hat{R}_\gamma(f) \leq \frac{\hat{n}_\gamma(f)}{n}$

**Proof.** As $\ell_{0-1}(t) \leq \ell_\gamma(t)$ for all $t \in \mathbb{R}$

$\mathbb{E}_{x,y}\,\ell_\gamma(f(x)y) \geq \mathbb{E}_{x,y}\,\ell_{0-1}(f(x)y), \quad R_{0-1}(f) \leq R_\gamma(f)$

$\ell_\gamma(t) = 0$ for all $t \geq \gamma$ and $\ell_\gamma(t) \leq 1$ for all $t \in \mathbb{R}$.

So $\hat{R}_\gamma(f) = \frac{1}{n}\sum_{i=1}^{n}\ell_\gamma(f(x_i)y_i) \leq \frac{1}{n}\hat{n}_\gamma(f)$. #

(b) $G = \{(x,y) \mapsto f(x)y : f \in \mathcal{F}\} \quad \mathcal{L}_\gamma = \{(x,y) \mapsto \ell_\gamma(f(x)y) : f \in \mathcal{F}\}$

Show that $\widehat{Rad}_n(\mathcal{L}_\gamma) \leq \frac{1}{\gamma}\widehat{Rad}_n(\mathcal{F})$

**Proof.** $\mathcal{L}_\gamma = \ell_\gamma \circ G$

By contraction lemma,

$$\widehat{Rad}_n(\mathcal{L}_\gamma) \leq Lip(\ell_\gamma)\,\widehat{Rad}_n(G)$$

$$= \frac{1}{\gamma}\mathbb{E}_\xi \sup_{f \in \mathcal{F}} \frac{1}{n}\sum_{i=1}^{n}\xi_i f(x_i)y_i$$

$$= \frac{1}{\gamma}\mathbb{E}_{\xi'} \sup_{f \in \mathcal{F}} \frac{1}{n}\sum_{i=1}^{n}\xi_i' f(x_i)$$

$$= \frac{1}{\gamma}\widehat{Rad}_n(\mathcal{F}) \qquad \#$$

where $\xi, \xi' \sim B(\frac{1}{2}, n)$.

(c) Fix $\gamma > 0$, $\delta \in (0,1)$. With prob. $\geq 1-\delta$ over sampling of $S$, $\forall f \in \mathcal{F}$

we have $R_{0-1}(f) \lesssim \frac{\hat{n}_\gamma(f)}{n} + \frac{1}{\gamma}\widehat{Rad}_n(\mathcal{F}) + \sqrt{\frac{\log(2/\delta)}{n}}$

**Proof.** We have

$$\begin{cases} R_{0-1}(f) \leq R_\gamma(f) \\[4pt] \sup_{f \in \mathcal{F}} |R_\gamma(f) - \hat{R}_\gamma(f)| \leq 2\,Rad_n(L_\gamma) + \sqrt{\frac{2\log(2/\delta)}{n}} \\[4pt] 0 \leq \hat{R}_\gamma(f) \leq \frac{\hat{n}_\gamma(f)}{n} \\[4pt] Rad_n(L_\gamma) = \mathbb{E}_{x,y}\,\widehat{Rad}_n(L_\gamma) \leq \frac{1}{\gamma}\widehat{Rad}_n(\mathcal{F}) \end{cases}$$

Thus $R_{0-1}(f) \leq R_\gamma(f) \leq \hat{R}_\gamma(f) + 2\,Rad_n(L_\gamma) + \sqrt{\frac{2\log(2/\delta)}{n}}$

$$\leq \frac{\hat{n}_\gamma(f)}{n} + \frac{1}{\gamma}\widehat{Rad}_n(\mathcal{F}) + \sqrt{\frac{2\log(2/\delta)}{n}} \quad . \#$$

(d) $\exists\, f^* \in \mathcal{F}$ s.t. $\mathbb{P}_{x,y}\{f^*(x)y \geq \gamma^*\} = 1$

$$\hat{f} = \arg\max_{f \in \mathcal{F}} \min_{i \in [n]} f(x_i)y_i$$

Show that $R_{0-1}(\hat{f}) \leq \frac{1}{\gamma^*}\widehat{Rad}_n(\mathcal{F}) + \sqrt{\frac{\log(2/\delta)}{n}}$

**Proof.**

Because we have $\mathbb{P}_{x,y}(f^*(x)y \geq \gamma^*) = 1$

we have $\min_{i \in [n]} \hat{f}(x_i)y_i \geq \min_{i \in [n]} f^*(x_i)y_i \geq \gamma^*$ a.s.

This implies $n_\gamma(\hat{f}) = 0$ a.s.

Using results in (c), we have

$$R_{0-1}(\hat{f}) \lesssim \frac{1}{\gamma^*}\widehat{Rad}_n(\mathcal{F}) + \sqrt{\frac{\log(2/\delta)}{n}} \quad \#$$