

Lecture 25: concentration / Rademacher

Announcement: HW3 out today

* Concentration

- * "Modern concerns"
- * Hoeffding / sub-Gaussian
- * Rademacher complexity
- * Rademacher for DL
- * Rademacher basics
- * Interpolation / double descent

Theorem (Hoeffding). Suppose (z_1, \dots, z_n) given, independent, $z_i \in [a_i, b_i]$ a.s.

$$\forall \epsilon > 0, \quad \mathbb{P} \left[\frac{1}{n} \sum z_i \leq \frac{1}{n} \sum \mathbb{E} z_i - \epsilon \right] \leq \sqrt{\frac{\sum (b_i - a_i)^2}{2n^2}} \ln \frac{1}{\delta}$$

Rem. write $F(z_1, \dots, z_n) = \frac{1}{n} \sum z_i$
note this convex & Lipschitz in $\vec{z} = (z_1, \dots, z_n)$;
a special case of convex + Lip Theorem.

Defn. X is σ^2 -sub-Gaussian $\iff \forall t \quad \ln \mathbb{E} \exp(tX) \leq \frac{t^2 \sigma^2}{2}$.
3 4 equiv defns of sub-G

Examples. (last time, but (X_1, \dots, X_n) , $(\sigma_1^2, \dots, \sigma_n^2)$ -sub-G
 $\Rightarrow \frac{1}{n} \sum X_i$ is $\frac{\sum \sigma_i^2}{n^2}$
"concentration" $\rightarrow n^2$

Chernoff bounding technique:

Lemma (Markov). $\forall X \geq 0$ a.s., $a > 0$,

$$\Pr[X \geq a] \leq \frac{\mathbb{E} X}{a}$$

Proof. Note $a \mathbb{1}_{\{X \geq a\}} \leq X$, apply $\mathbb{E}(\cdot)$ to both sides //

Chernoff bounding technique (X sub-G v.v., possibly negative)
 $\forall t > 0$

$$\Pr[X \geq \epsilon] \leq \Pr[\exp(tX) \geq \exp(t\epsilon)]$$

$$\Rightarrow \Pr[X \geq \epsilon] \leq \inf_{t > 0} \Pr[\exp(tX) \geq \exp(t\epsilon)]$$

$$\leq \inf_{t > 0} \frac{\mathbb{E} \exp(tX)}{\exp(t\epsilon)} \leq \inf_{t > 0} \exp\left(\frac{t^2 \sigma^2}{2} - t\epsilon\right)$$

$$= \exp\left(\inf_{t > 0} \left[\frac{t^2 \sigma^2}{2} - t\epsilon\right]\right) = \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right)$$

Rem

$$\frac{\sum X_i}{n} \xrightarrow{a.s.} \mathbb{E} X, \quad \text{SLLN}$$

$$\frac{\sum X_i}{\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, 1) \quad \text{CLT}$$

$$\limsup_{n \rightarrow \infty} \frac{\left| \frac{\sum X_i}{\sqrt{2n \ln n}} \right|}{\sqrt{2n \ln n}} = 1 \quad \text{a.s. LIL}$$

One more tail inequality

Theorem (McDiarmid)

Suppose (z_1, \dots, z_n) independent, $z_i \in \mathcal{S}_i$ a.s.
let $F: \mathcal{H} \rightarrow \mathbb{R}$ be given,
suppose $\exists (c_1, \dots, c_n)$ s.t. bounded differences hold: $\forall i$

$$\sup_{\substack{z_1, \dots, z_n \in \mathcal{S}_i \\ z_i'}} \left| F(z_1, \dots, z_n) - F(z_1, \dots, z_i', \dots, z_n) \right| \leq c_i$$

Then $F(z_1, \dots, z_n)$ is $\frac{\sum c_i^2}{4}$ -sub-Gaussian

and w/ $\forall \delta \geq 1 - \delta$

$$\mathbb{P}[F(z_1, \dots, z_n) \leq \mathbb{E} F(z_1, \dots, z_n) - \epsilon] \leq \sqrt{\frac{\sum c_i^2}{2n^2}} \ln \frac{1}{\delta}$$

Rademacher complexity

Abstract goal: w/ priors, $\forall f \in \mathcal{F}$ $R(f) \leq \hat{R}(f) + \sqrt{\frac{\ln 1/\delta}{n}} + \sqrt{\frac{\text{compl}(f)}{n}}$ brute force way to handle the output

We'll get this shortly

Define. Given $V \subseteq \mathbb{R}^n$, $URad(V) = \mathbb{E}_{\epsilon} \sup_{u \in V} \langle u, \epsilon \rangle$
 $\epsilon = \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i e_i$; $\epsilon_i \sim \text{Unif}(\{-1, +1\})$

- Rem: ① Others use $Rad(V) = \frac{1}{\sqrt{n}} URad(V)$
 ② Typical usage $V = \mathcal{G}_{1z} = \{ (g(z_1), \dots, g(z_n)) : g \in \mathcal{G} \}$
 or $(\mathcal{L} \circ \mathcal{F})_{1s} = \{ (L(y_1, \theta_{k1}), \dots, L(y_n, \theta_{kn})) : f \in \mathcal{F} \} \subseteq \mathbb{R}^n$
"loss class"

Theorem (core Rademacher bound): Let z_1, \dots, z_n iid Let \mathcal{G} be given with $\mathcal{G}(z_i) \in [a, b]$ a.s. Then w/ priors, $\sup_{g \in \mathcal{G}} (\mathbb{E} g(z_i) - \frac{1}{n} \sum g(z_i)) \leq \frac{2}{\sqrt{n}} URad(\mathcal{G}_{1z}) + (b-a) \sqrt{\frac{\ln 2/\delta}{2n}}$

Remarks:

- ① (History) term coined by Bartlett-Mendelson (2002), but is much older (fied "symmetrization")
- ② Many weaknesses due to where that proof was cut, e.g., $1/n$ rates
- ③ Not oldest technique (covering Kolmogorov-Tikhonov)
- ④ Non-weakness: note $URad(\sum_{i=1}^n \epsilon_i e_i) = \mathbb{E}_{\epsilon} \sup_{u \in \sum_{i=1}^n \epsilon_i e_i} \langle \epsilon, u \rangle = n$

Recall DL is a universal approximator (cf. "rethinking generalization")
easy non-fix: norm bound!
 (spirit of criticism: DL theorists need to think harder.)

Rademacher calculus

Sanity checks

$$URad(\sum \epsilon u_0) = \mathbb{E}_{\epsilon} \langle u_0, \epsilon \rangle = 0$$

$$URad(V + \sum \epsilon u_0) = \mathbb{E}_{\epsilon} \sup_{v \in V} \langle \epsilon, v + u_0 \rangle = \mathbb{E}_{\epsilon} \langle u_0, \epsilon \rangle + \mathbb{E}_{\epsilon} \sup_{v \in V} \langle \epsilon, v \rangle = URad(V)$$

$$V \subseteq V' \Rightarrow URad(V) \leq URad(V')$$

$$URad(\sum_{i=1}^n \epsilon_i e_i) = n, \quad URad(\sum_{i=1}^n \epsilon_i (1, \dots, 1), (+1, \dots, +1)) = \mathbb{E}_{\epsilon} \sqrt{|\sum \epsilon_i|^2} \in [\sqrt{\frac{n}{2}}, \sqrt{n}]$$

Khintchine inequality

$$URad(V) \geq 0.$$

Lemma.

- ① $URad(cV + \sum \epsilon u_0) = |c| \cdot URad(V)$
- ② $URad(\text{conv}(V)) = URad(V)$
- ③ Let $(V_i)_{i=1}^n$ given w/ $\sup_{u \in V_i} \langle \epsilon, u \rangle \geq 0 \quad \forall \epsilon \in \sum_{i=1}^n \epsilon_i e_i$
 Then $URad(\bigcup_{i=1}^n V_i) \leq \sum_{i=1}^n URad(V_i)$

- Rem. ① union rule often loose (even exponentially)
 ② Another definition: $\mathbb{E}_{\epsilon} \sup_{u \in V} |\langle \epsilon, u \rangle| = URad(V \cup -V)$

Lemma. $URad(\{x \mapsto w^T x : \|w\| \leq R\}) \leq \|X\|_F \cdot R$

Proof. for any $\epsilon \in \sum_{i=1}^n \epsilon_i e_i$, $\sup_{\|w\| \leq R} \langle \epsilon, w \rangle = \langle w_0, \epsilon \rangle + \sup_{\|w-w_0\| \leq R} \langle \epsilon, w-w_0 \rangle$
ops, forgot dotz.

So $\sup_{\|w\| \leq R} \langle \epsilon, w \rangle = \sup_{\|w\| \leq R} \sum \epsilon_i v^T x_i = \sup_{\|w\| \leq R} \langle \sum \epsilon_i x_i, w \rangle = R \|\sum \epsilon_i x_i\|$

$\mathbb{E}_{\epsilon} \sup_{\|w\| \leq R} \sum \epsilon_i v^T x_i = \mathbb{E}_{\epsilon} R \|\sum \epsilon_i x_i\|$
 $\leq R \sqrt{\mathbb{E}_{\epsilon} \langle \sum \epsilon_i x_i, \sum \epsilon_i x_i \rangle}$
 $= R \sqrt{\mathbb{E}_{\epsilon} \sum \epsilon_i \epsilon_i \|x_i\|^2} = R \sqrt{\sum \|x_i\|^2} = R \|X\|_F$