

Generalization & uniform convergence

Def A (centered) RV X is σ -subGaussian (or σ -SG; variance proxy σ^2) if $\mathbb{E} \exp(\lambda X) \leq \exp(\lambda^2 \sigma^2 / 2)$, $\forall \lambda \in \mathbb{R}$.

Lemma If X is σ -subGaussian, then for any $\varepsilon > 0$,

$$\mathbb{P}(X \geq \varepsilon) \leq \exp(-\varepsilon^2 / 2\sigma^2).$$

Pf: Exercise. Hint: note that $\mathbb{P}(X \geq \varepsilon) = \sup_{t \geq 0} \mathbb{P}(\exp(tX) \geq \exp(t\varepsilon))$.

- Bounded RV's are SG. Exercise: If $X \in [a, b]$ a.s., X is SG w/ variance proxy $(b-a)^2/4$.
- Gaussians are SG.
- Sums of SG are SG.

Homework 0: If X_1, \dots, X_n are indep. σ_i -SG RV's, then

$Z := \sum_{i=1}^n X_i$ is SG with variance proxy $\sum_{i=1}^n \sigma_i^2$; and if $\alpha > 0$ then αX_i is $\alpha \sigma_i$ -SG (variance proxy: $\alpha^2 \sigma_i^2$).

Thus, if

$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ | then \bar{X} is SG w/ variance proxy $\frac{1}{n^2} \sum_{i=1}^n \sigma_i^2$.

→ Lemma says $\mathbb{P}(\bar{X}_n \geq \varepsilon) \leq \exp\left(\frac{-\varepsilon}{\frac{1}{n^2} \sum_{i=1}^n \sigma_i^2}\right) = \exp\left(\frac{-n^2 \varepsilon}{2 \sum_{i=1}^n \sigma_i^2}\right)$
for indep 0-SG X_i , Similarly, $\mathbb{P}(\bar{X}_n \leq -\varepsilon) \leq \exp\left(\frac{-n^2 \varepsilon}{2 \sum_{i=1}^n \sigma_i^2}\right)$

Thus, for indep 0-SG X_i , (each have $\mathbb{E} X_i = 0$)

$$\mathbb{P}(|\bar{X}_n| > \varepsilon) \leq 2 \exp\left(-\frac{n^2 \varepsilon}{2 \sum_{i=1}^n \sigma_i^2}\right).$$

Generalizing to non-mean zero, we get

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n (X_i - \mu_i)\right| > \varepsilon\right) \leq 2 \exp\left(-\frac{n^2 \varepsilon}{2 \sum_{i=1}^n \sigma_i^2}\right)$$

Letting $\sigma^2 := \frac{1}{n} \sum_{i=1}^n \sigma_i^2$, we get $\leq 2 \exp\left(\frac{-n \varepsilon}{2 \sigma^2}\right)$.

For $\varepsilon := \sigma \cdot \sqrt{\frac{2 \log(2/\delta)}{n}}$ we get $2 \exp(-n \varepsilon / 2 \sigma^2) = \delta$, so

$$\text{w.p. } > 1 - \delta, \quad \left|\frac{1}{n} \sum_{i=1}^n (X_i - \mu_i)\right| \leq \sqrt{\frac{1}{n} \sum_{i=1}^n \sigma_i^2} \cdot \sqrt{\frac{2 \log(2/\delta)}{n}}.$$

If X_i are iid, says w.p. $> 1 - \delta$, $\left|\mu - \frac{1}{n} \sum_{i=1}^n X_i\right| \leq \sqrt{\frac{\sigma^2 \log(2/\delta)}{n}}$.
each 0-SG,

→ as n gets larger, sample mean closer to pop mean.

EX. Let $(x_i, y_i) \stackrel{iid}{\sim} P$, $x_i \in \mathbb{R}^d$, $y_i \in \{\pm 1\}$, $f: \mathbb{R}^d \rightarrow \{\pm 1\}$,
and let $z_i := 1(f(x_i) \neq y_i)$. Then each z_i is iid,
bounded (hence SG: with variance proxy $\frac{1}{4}$) in $[0, 1]$. So by above,

$$\text{w.p.} > 1 - \delta, \quad \mathbb{P}(y \neq f(x)) \leq \frac{1}{n} \sum_{i=1}^n 1(y_i \neq f(x_i)) + \sqrt{\frac{\log(2/\delta)}{2n}}.$$

→ test error is bounded by train error $+ \tilde{O}(\sqrt{\frac{1}{n}})$.

Example. Suppose $(x_i, y_i)_i^n$ are iid. For any $n \in \mathbb{N}$, define:

$$f_n(x) := \begin{cases} y_i: x \in \{x_1, \dots, x_n\}, \\ -10: \text{otherwise} \end{cases}$$

Consider two situations:

① X has finite support. Then $\frac{1}{n} \sum_{i=1}^n 1(y_i \neq f_n(x_i)) = 0$ for all n by defn.

and $\mathbb{P}(y \neq f_n(x)) \rightarrow 0$ as well, since we recover all pts.

② X has continuous distri Then $\frac{1}{n} \sum_{i=1}^n 1(y \neq f(x_i)) = 0$ by construction,
but $\mathbb{P}(y \neq f_n(x)) = 1 \quad \forall n$.

What broke subG concentration?

f_n is a random variable. Although (x_i, y_i) are iid,

$Z_i := 1(y_i \neq f_n(x_i))$ are not independent.

② is overfitting: $\hat{L}(f) = 0$ but $L(f) = 1$.

How can we guarantee test error is small when looking at training error?

We'll see how via uniform convergence:

For iid Z_i , loss $\ell(Z_i)$,

$$L(f) := \mathbb{E} \ell(f), \quad \hat{L}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(Z_i)),$$

Goal: bound $L(f)$. Suppose $f \in \mathcal{F}$, some function class \mathcal{F} .
And suppose we use $S = \{Z_i\}_{i=1}^n$ to fit $\hat{f} = \hat{f}(S)$. Then we typically lose indep. of $f(Z_i; S)$.

Approach is then:

$$L(f) = L(f) - \hat{L}(f) + \hat{L}(\hat{f})$$

$$\leq \hat{L}(\hat{f}) + \sup_{f \in \mathcal{F}} \{ L(f) - \hat{L}(f) \}.$$

Seems very silly, but we will see very fruitful to do so.

We'll prove deviation bounds that hold uniformly over $f \in \mathcal{F}$.

Example. Let $\mathcal{F} = \{f_1, \dots, f_k\}$, $|\mathcal{F}| = k$. If $(x_i, y_i)_{i=1}^n$ are iid, $f_i: \mathcal{X} \rightarrow \{\pm 1\}$, then SG concentration as before gives for fixed f_l ,

$$\mathbb{P}_{(x_i, y_i)_{i=1}^n} \left(\left| \mathbb{P}(f_l(X) \neq y) - \frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i \neq f_l(x_i)) \right| > \sqrt{\frac{\log 2/\delta}{2n}} \right) \leq \delta.$$

for fixed l , w.p. $> 1 - \delta$, $\left| \mathbb{P}(f_l(X) \neq y) - \hat{\mathbb{P}}(f_l(X) \neq y) \right| \leq \sqrt{\frac{\log 2/\delta}{n}}.$

Union bound:

$$\mathbb{P} \left(\exists l \in [k] : \left| \mathbb{P}(f_l(X) \neq y) - \hat{\mathbb{P}}(f_l(X) \neq y) \right| \geq \sqrt{\frac{\log 2k/\delta}{2n}} \right) \leq k \cdot \frac{\delta}{k} = \delta.$$

i.e. w.p. $> 1 - \delta$, for all $l \in [k]$, $\left| \mathbb{P}(y \neq f_l(X)) - \hat{\mathbb{P}}(f_l(X) \neq y) \right| \leq \sqrt{\frac{\log(2k/\delta)}{2n}}.$

$$\leq \sqrt{\frac{\log |\mathcal{F}|}{2n}} + \sqrt{\frac{\log(4/\delta)}{2n}}.$$

For finite classes, get $\sqrt{\frac{\log |\mathcal{F}|}{2n}}$ extra term. We'll see next that Rademacher complexity allows for dealing w/ $|\mathcal{F}| = \infty$.

Def. For $V \subset \mathbb{R}^n$, the unnormalized/normalized Rademacher complexity is

$$\text{URad}(V) := \mathbb{E}_{\varepsilon} \sup_{u \in V} \langle \varepsilon, u \rangle, \quad \text{Rad}(V) = \frac{1}{n} \text{URad}(V),$$

where $\varepsilon \in \mathbb{R}^n$ is iid Rademacher: $\varepsilon_i \sim \text{Unif}(\{\pm 1\})$.

We will typically apply this to outputs of a function class over training data.
e.g. for $z_i = (x_i, y_i)$, $S = \{z_i\}_1^n$, for a class \mathcal{F} ,

$$\mathcal{F}_S := \{ (f(z_1), \dots, f(z_n)) : f \in \mathcal{F} \}.$$

$$\rightarrow \text{URad}(\mathcal{F}_S) = \mathbb{E}_{\varepsilon} \sup_{f \in \mathcal{F}} \langle \varepsilon, u \rangle = \mathbb{E}_{\varepsilon} \sup_{f \in \mathcal{F}} \sum_1^n \varepsilon_i f(z_i).$$

- $\text{URad}(\mathcal{F}_S)$ is large if, for any $\varepsilon_i \in \{\pm 1\}$, there is some $f \in \mathcal{F}$ st $f(z_i) \geq \varepsilon_i$.

- If we think of $f(z_i) \in \{\pm 1\}$, then this corresponds to \mathcal{F} fitting "random labels" ε_i .

- We'll often look at URad for losses, i.e. for ℓ ,

$$\text{URad}((\ell \circ \mathcal{F})_S) = \text{URad}(\ell(y_1, f(x_1)), \dots, \ell(y_n, f(x_n)) : f \in \mathcal{F}).$$

- $\text{Rad}(V)$ vaguely measures how large/complicated V is.

Properties : (1) $\text{URad}(\{u\}) = \mathbb{E} \langle \varepsilon, u \rangle = 0$.

(2) $\text{URad}(V + \{u\}) = \text{URad}(\{v + u : v \in V\}) = \text{URad}(V)$.

(3) If $V \subset V'$, $\text{URad}(V) \subset \text{URad}(V')$.

(4) $\text{URad}(\{\pm 1\}^n) = \mathbb{E}_{\varepsilon} \sup_{x \in \{\pm 1\}^n} \langle \varepsilon, x \rangle = \mathbb{E}_{\varepsilon} \|\varepsilon\|^2 = n$.

$\rightarrow \{\pm 1\}^n$ is as large as possible among vectors taking vals in ± 1 .

(5) $\text{URad}(\{(-1, -1, \dots, -1), (1, 1, \dots, 1)\}) = \mathbb{E}_{\varepsilon} \max\{\sum \varepsilon_i, -\sum \varepsilon_i\} = \mathbb{E}_{\varepsilon} \left| \sum_1^n \varepsilon_i \right|$.

$$\left| \sum_1^n \varepsilon_i \right| = \left| \sum_1^n (2 \cdot \text{Ber}(\tfrac{1}{2}) - 1) \right| = \left| 2 \cdot \text{Bin}(n, \tfrac{1}{2}) - n \right|.$$

Anti-concentration of Binomial shows $|2 \text{Bin}(n, \tfrac{1}{2}) - n| = \Theta(\sqrt{n})$.

You will also sometimes see an absolute value version of Rad. complexity,

$$\widetilde{\text{URad}}(V) := \mathbb{E}_{\varepsilon} \sup_{u \in V} |\langle \varepsilon, u \rangle|.$$

Similar idea, but a bit less nice for reasons we won't get into.

Theorem. Let \mathcal{F} be a fun class w $f(z) \in [a, b] \forall z, \forall f \in \mathcal{F}$, let \mathbb{P} :
 distr over \mathcal{Z} .

① For any $\delta \in (0, 1)$, w.p. $> 1 - \delta$,

$$\sup_{f \in \mathcal{F}} \left\{ \mathbb{E} f(z) - \frac{1}{n} \sum_{i=1}^n f(z_i) \right\} \leq \mathbb{E}_{z_1} \left(\sup_{f \in \mathcal{F}} \left\{ \mathbb{E} f(z) - \frac{1}{n} \sum_{i=1}^n f(z_i) \right\} \right) + (b-a) \sqrt{\frac{\log \frac{1}{\delta}}{2n}}.$$

② w.p. $> 1 - \delta$,

$$\mathbb{E}_{z_i} \text{Rad}(\mathcal{F}|_S) \leq \text{Rad}(\mathcal{F}|_S) + (b-a) \sqrt{\frac{\log \frac{1}{\delta}}{2n}}$$

③ w.p. $> 1 - \delta$,

$$\sup_{f \in \mathcal{F}} \left\{ \mathbb{E} f(z) - \frac{1}{n} \sum_{i=1}^n f(z_i) \right\} \leq 2 \text{Rad}(\mathcal{F}|_S) + 3(b-a) \sqrt{\frac{\log \frac{2}{\delta}}{n}}$$

To prove this, we'll use MacDiarmid's ineq:

Thm (MacDiarmid). Suppose $F: \mathbb{R}^n \rightarrow \mathbb{R}$ satisfies bounded differences:

$\forall i \in [n], \exists c_i$ st $\sup_{z_1, \dots, z_n, z_i'} |F(z_1, \dots, z_i, z_{i+1}, \dots, z_n) - F(z_1, \dots, z_i', z_{i+1}, \dots, z_n)| \leq c_i$. Then,

$$\text{w.p. } > 1 - \delta, \quad \mathbb{E}_{z_i} F(z_1, \dots, z_n) \leq F(z_1, \dots, z_n) + \sqrt{\frac{\sum_i c_i^2 \log \frac{1}{\delta}}{2}}.$$

Lemma. Let $(z_1, \dots, z_n), (z_1', \dots, z_n')$ be iid from \mathbb{P} .

Let $\hat{\mathbb{P}}_n$: uniform on (z_1, \dots, z_n) ; $\hat{\mathbb{P}}_n'$: uniform on (z_1', \dots, z_n') . Same for $\mathbb{P}_n, \mathbb{P}_n'$.

$$\text{Then } \mathbb{E}_n \left[\sup_{f \in \mathcal{F}} \left\{ \mathbb{E} f - \hat{\mathbb{E}}_n f \right\} \right] \leq \mathbb{E}_n \left[\mathbb{E}_n' \left(\sup_{f \in \mathcal{F}} \left\{ \hat{\mathbb{E}}_n' f - \hat{\mathbb{E}}_n f \right\} \right) \right].$$

Pf.

First note that since $z_i' \stackrel{d}{=} z_i$,

$$\mathbb{E} f_\varepsilon = \mathbb{E}_{z \sim \mathbb{P}} f_\varepsilon(z) = \mathbb{E}_n' \hat{\mathbb{E}}_n' f_\varepsilon, \text{ since } z_i' \sim \mathbb{P} \text{ so } \mathbb{E}_{z_i'} f_\varepsilon(z_i') = \mathbb{E} f_\varepsilon.$$

Let $\varepsilon > 0$. Then $\exists f_\varepsilon \in \mathcal{F}$ s.t. $\sup_{f \in \mathcal{F}} \left\{ \mathbb{E} f - \hat{\mathbb{E}}_n f \right\} \leq \mathbb{E} f_\varepsilon - \hat{\mathbb{E}}_n f_\varepsilon + \varepsilon$.

$$\Rightarrow \mathbb{E}_n \left[\sup_{f \in \mathcal{F}} \left\{ \mathbb{E} f - \hat{\mathbb{E}}_n f \right\} \right] \leq \mathbb{E}_n \left[\mathbb{E} f_\varepsilon - \hat{\mathbb{E}}_n f_\varepsilon + \varepsilon \right].$$

$$= \mathbb{E}_n \left[\mathbb{E}_n' \hat{\mathbb{E}}_n' f_\varepsilon - \hat{\mathbb{E}}_n f_\varepsilon + \varepsilon \right]$$

$$= \mathbb{E}_n' \mathbb{E}_n \left[\hat{\mathbb{E}}_n' f_\varepsilon - \hat{\mathbb{E}}_n f_\varepsilon \right] + \varepsilon$$

$$\leq \mathbb{E}_n' \mathbb{E}_n \left[\sup_{f \in \mathcal{F}} \left\{ \hat{\mathbb{E}}_n' f_\varepsilon - \hat{\mathbb{E}}_n f_\varepsilon \right\} \right] + \varepsilon.$$

Since $\varepsilon > 0$ is arbitrary, this completes the proof. \square

Lemma. $\mathbb{E}_n [\mathbb{E}_n' \sup_{f \in \mathcal{F}} \{ \hat{\mathbb{E}}_n' f - \hat{\mathbb{E}}_n f \}] \leq 2 \mathbb{E}_n \text{Rad}(\mathcal{F}|_S)$.

If For fixed $\varepsilon \in \{\pm 1\}^n$, let RV $\xi_i := (u_i, u_i') := \begin{cases} (z_i, z_i'), & \varepsilon = 1; \\ (z_i', z_i), & \varepsilon = -1. \end{cases}$

By defⁿ,

$$\begin{aligned} \mathbb{E}_n \mathbb{E}_n' \left[\sup_{f \in \mathcal{F}} \{ \hat{\mathbb{E}}_n' f - \hat{\mathbb{E}}_n f \} \right] &= \mathbb{E}_n \mathbb{E}_n' \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_i \varepsilon_i (f(z_i') - f(z_i)) \right\} \right] \\ &= \mathbb{E}_n \mathbb{E}_n' \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_i \varepsilon_i \cdot (f(u_i') - f(u_i)) \right\} \right] \end{aligned}$$

Since $\{z_i, z_i'\}$ are iid, if ε_i are iid Rademacher, so are $\{u_i, u_i'\}$, and in particular $(z_1, \dots, z_n, z_1', \dots, z_n') = (u_1, \dots, u_n, u_1', \dots, u_n')$. Thus,

$$\begin{aligned} \mathbb{E}_\varepsilon \mathbb{E}_n \mathbb{E}_n' \left[\sup_{f \in \mathcal{F}} \{ \hat{\mathbb{E}}_n' f - \hat{\mathbb{E}}_n f \} \right] &= \mathbb{E}_\varepsilon \mathbb{E}_n \mathbb{E}_n' \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_i \varepsilon_i (f(u_i') - f(u_i)) \right\} \right] \\ &= \mathbb{E}_\varepsilon \mathbb{E}_n \mathbb{E}_n' \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_i \varepsilon_i (f(z_i') - f(z_i)) \right\} \right] \\ &\leq \mathbb{E}_\varepsilon \mathbb{E}_n \mathbb{E}_n' \left[\sup_{f, f' \in \mathcal{F}} \left\{ \frac{1}{n} \sum_i \varepsilon_i (f(z_i') - f'(z_i)) \right\} \right] \\ &= \mathbb{E}_\varepsilon \mathbb{E}_n' \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_i \varepsilon_i f(z_i') \right] + \mathbb{E}_\varepsilon \mathbb{E}_n \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_i \varepsilon_i (-\varepsilon_i) f'(z_i) \right] \\ &= 2 \mathbb{E}_\varepsilon \mathbb{E}_n \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_i \varepsilon_i f(z_i) \right] \quad \text{since } z_i \stackrel{d}{=} z_i', \varepsilon_i \stackrel{d}{=} -\varepsilon_i \\ &= 2 \mathbb{E}_n \left[\frac{1}{n} \mathbb{E}_\varepsilon \sup_{f \in \mathcal{F}} \sum_i \varepsilon_i f(z_i) \right] = 2 \mathbb{E}_n \text{Rad}(\mathcal{F}|_S). \quad \square \end{aligned}$$

This shows that

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \left\{ \mathbb{E} f - \hat{\mathbb{E}}_n f \right\} \right] \leq 2 \mathbb{E}_n \text{Rad}(\mathcal{F}|_S).$$

We'll now work on making a high-probability version of this

Theorem (McDiarmid bounded differences):

Suppose $g: \mathbb{R}^n \rightarrow \mathbb{R}$ is st. $\forall i \in \{1, \dots, n\}$, $\exists c_i$ s.t.

$$\sup_{z_1, \dots, z_n, z_i'} |g(z_1, \dots, z_i, \dots, z_n) - g(z_1, \dots, z_i', \dots, z_n)| \leq c_i.$$

Then, w.p. $> 1 - \delta$,

$$\mathbb{E}_n g(z_1, \dots, z_n) \leq g(z_1, \dots, z_n) + \sqrt{\frac{\sum_i c_i^2}{2} \log(1/\delta)}$$

Pf omitted, see linked notes from Daniel Hsu

We'll now prove Thm xx.

① We will verify that $\sup_{f \in \mathcal{F}} \{ \mathbb{E} f(Z) - \hat{\mathbb{E}}_n f \}$ satisfies bounded differences with constant $\frac{b-a}{n}$.

Consider z_1, \dots, z_n, z'_i . For $j \neq i$, call $z'_j = z_j$. Then,

$$\begin{aligned} & \left| \sup_{f \in \mathcal{F}} \{ \mathbb{E} f - \hat{\mathbb{E}}_n f \} - \sup_{g \in \mathcal{F}} \{ \mathbb{E} g - \hat{\mathbb{E}}'_n g \} \right| \\ &= \left| \sup_{f \in \mathcal{F}} \{ \mathbb{E} f - \hat{\mathbb{E}}_n f \} - \sup_{g \in \mathcal{F}} \left\{ \mathbb{E} g - \frac{1}{n} \sum_{i=1}^n g(z_i) + \frac{1}{n} g(z_i) - \frac{1}{n} g(z'_i) \right\} \right| \\ &= \left| \sup_{f \in \mathcal{F}} \{ \mathbb{E} f - \hat{\mathbb{E}}_n f \} - \sup_{g \in \mathcal{F}} \left\{ \mathbb{E} g - \frac{1}{n} \sum_{i=1}^n g(z_i) + \frac{g(z_i)}{n} - \frac{g(z'_i)}{n} \right\} \right| \\ &\leq \sup_{h \in \mathcal{F}} \left\{ \left| \sup_{f \in \mathcal{F}} \{ \mathbb{E} f - \hat{\mathbb{E}}_n f \} - \sup_{g \in \mathcal{F}} \left\{ \mathbb{E} g - \hat{\mathbb{E}}_n g + \frac{h(z_i)}{n} - \frac{h(z'_i)}{n} \right\} \right| \right\} \\ &= \sup_{h \in \mathcal{F}} \left| \frac{h(z_i) - h(z'_i)}{n} \right| \leq \frac{b-a}{n}. \end{aligned}$$

\Rightarrow satisfies bounded diff. w/ $c_i = \frac{b-a}{n} \forall i$. $\sum c_i^2 = \frac{n(b-a)^2}{n^2}$ so
Thus, w.p. $> 1-\delta$,

$$\sup_{f \in \mathcal{F}} \{ \mathbb{E} f - \hat{\mathbb{E}}_n f \} \leq \mathbb{E}_n \left[\sup_{f \in \mathcal{F}} \{ \mathbb{E} f - \hat{\mathbb{E}}_n f \} \right] + \sqrt{\frac{(b-a)^2 \log(2/\delta)}{n}}.$$

Let $S = \{z_i\}$, $S' = \{z'_i\}$.

$$\begin{aligned}
(2) \quad & |URad(\mathcal{F}|_S) - URad(\mathcal{F}|_{S'})| \\
&= |URad(\mathcal{F}|_S) - \mathbb{E}_\varepsilon \sup_{f \in \mathcal{F}} \sum_{i=1}^n \varepsilon_i f(z'_i)| \\
&\leq |URad(\mathcal{F}|_S) - \mathbb{E}_\varepsilon \sup_{f \in \mathcal{F}} \{ \sum_{i=1}^n \varepsilon_i f(z_i) - \varepsilon_i f(z_i) + \varepsilon_i f(z_i) \}| \\
&= |URad(\mathcal{F}|_S) - \mathbb{E}_\varepsilon \sup_{f \in \mathcal{F}} \{ \sum_{i=1}^n \varepsilon_i f(z_i) - \varepsilon_i f(z_i) + \varepsilon_i f(z'_i) \}| \\
&\leq \sup_{h \in \mathcal{F}} |URad(\mathcal{F}|_S) - \mathbb{E}_\varepsilon \sup_{f \in \mathcal{F}} \{ \sum_{i=1}^n \varepsilon_i f(z_i) - \varepsilon_i h(z_i) + \varepsilon_i h(z'_i) \}| \\
&= \sup_{h \in \mathcal{F}} | \mathbb{E}_\varepsilon [\varepsilon_i h(z_i) - \varepsilon_i h(z'_i)] | \\
&\leq \sup_{h \in \mathcal{F}} \mathbb{E}_\varepsilon | \varepsilon_i h(z_i) - \varepsilon_i h(z'_i) | \leq (b-a).
\end{aligned}$$

Satisfies bounded differences w/ $c_i \equiv b-a$, so $\sum c_i^2 = (b-a)^2 n$;
 dividing by n to get normalized Rad complexity we get

w.p. $> 1-\delta$,

$$\mathbb{E}_n \text{Rad}(\mathcal{F}|_S) \leq \text{Rad}(\mathcal{F}|_S) + \sqrt{\frac{(b-a)^2 \log 2/\delta}{n}}.$$

Putting everything together,

Thus, w.p. $> 1-\delta$,

$$\sup_{f \in \mathcal{F}} \{ \mathbb{E}f - \hat{\mathbb{E}}_n f \} \leq \mathbb{E}_n \left[\sup_{f \in \mathcal{F}} \{ \mathbb{E}f - \hat{\mathbb{E}}_n f \} \right] + \sqrt{\frac{(b-a)^2 \log(2/\delta)}{n}}.$$

$$\leq \mathbb{E}_n \left[\mathbb{E}_n' \left(\sup_{f \in \mathcal{F}} \{ \mathbb{E}_n' f - \hat{\mathbb{E}}_n f \} \right) \right] + (b-a) \sqrt{\frac{\log 2/\delta}{n}}$$

$$\leq 2 \mathbb{E}_n \text{Rad}(\mathcal{F}|_S) + (b-a) \sqrt{\frac{\log 2/\delta}{n}}$$

$$\leq 2 \text{Rad}(\mathcal{F}|_S) + 3(b-a) \sqrt{\frac{\log 2/\delta}{n}}. \quad \square$$

Thus Rademacher complexity provides a distribution-dependent (via $\mathcal{F}|_S$; S depends on \mathcal{Q}) way to guarantee uniform convergence.

We'll now instantiate for particular function classes.

Example logistic regression with bounded weights.

$$l(y, f(x)) := \log(1 + \exp(-y f(x)));$$

$$\mathcal{F} = \{ w \in \mathbb{R}^d : \|w\| \leq B \};$$

$$(\mathcal{l} \circ \mathcal{F})|_S := \{ (l(y_1, w^T x_1), \dots, l(y_n, w^T x_n)) : \|w\| \leq B \},$$

$$R(w) := \mathbb{E} l(y \langle w, x \rangle), \quad \hat{R}(w) = \frac{1}{n} \sum_{i=1}^n l(y_i \langle w, x_i \rangle).$$

Via prev theorem, suffices to bound $\text{Rad}((\mathcal{l} \circ \mathcal{F})|_S)$.

Lemma Let $l: \mathbb{R}^n \rightarrow \mathbb{R}^n$ have components l_i which are univariate & L -Lip.

Then $\text{Rad}(l \circ V) \leq L \cdot \text{Rad}(V)$.

$$\begin{aligned}
 \text{pf } \mathbb{U}\text{Rad}(l \circ V) &= \mathbb{E}_{\varepsilon} \sup_{u \in V} \sum_{i=1}^n \varepsilon_i l_i(u_i) \\
 &= \mathbb{E}_{\varepsilon} \left[\sup_{u \in V} \left\{ \varepsilon_1 l_1(u_1) + \sum_{i=2}^n \varepsilon_i l_i(u_i) \right\} \right] \\
 &= \frac{1}{2} \mathbb{E}_{\varepsilon_{2:n}} \left[\sup_{u \in V} \left\{ l_1(u_1) + \sum_{i=2}^n \varepsilon_i l_i(u_i) \right\} \right. \\
 &\quad \left. + \sup_{u \in V} \left\{ -l_1(u_1) + \sum_{i=2}^n \varepsilon_i l_i(u_i) \right\} \right] \\
 &= \frac{1}{2} \mathbb{E}_{\varepsilon_{2:n}} \left[\sup_{u_1, w_1 \in V} \left\{ l_1(u_1) - l_1(w_1) + \sum_{i=2}^n \varepsilon_i (l_i(u_i) + l_i(w_i)) \right\} \right] \\
 &\leq \frac{1}{2} \mathbb{E}_{\varepsilon_{2:n}} \left[\sup_{u_1, w_1 \in V} \left\{ L|u_1 - w_1| + \sum_{i=2}^n \varepsilon_i (l_i(u_i) + l_i(w_i)) \right\} \right] \\
 &= \frac{1}{2} \mathbb{E}_{\varepsilon_{2:n}} \left[\sup_{u_1, w_1 \in V} \left\{ L(u_1 - w_1) + \sum_{i=2}^n \varepsilon_i (l_i(u_i) + l_i(w_i)) \right\} \right] \\
 &= \frac{1}{2} \mathbb{E}_{\varepsilon_{2:n}} \left[\sup_{u \in V} \left\{ Lu_1 + \sum_{i=2}^n \varepsilon_i l_i(u_i) \right\} + \sup_{w \in V} \left\{ -Lw_1 + \sum_{i=2}^n \varepsilon_i l_i(w_i) \right\} \right] \\
 &= \mathbb{E} \sup_{u \in V} \left\{ L\varepsilon_1 u_1 + \sum_{i=2}^n \varepsilon_i l_i(u_i) \right\} \\
 &= \dots = \mathbb{E} \sup_{u \in V} L \langle u, \varepsilon \rangle = \mathbb{U}\text{Rad}(L \cdot V) = L \cdot \mathbb{U}\text{Rad}(V).
 \end{aligned}$$

Corollary If ℓ is L -Lip. $\exists \ell \circ f \in [a, b]$ a.s., then

$$\text{w.p. } > 1 - \delta, \quad \forall f \in \mathcal{F}, \quad |R_\ell(f)| \leq \hat{R}_\ell(f) + 2L \text{Rad}(\mathcal{F}|_S) + 3(b-a) \sqrt{\frac{\log \frac{2}{\delta}}{n}}.$$

Pf: $|\ell(-y_i f(x_i)) - \ell(-y_i f'(x_i))| \leq L |-y_i f(x_i) - (-y_i f'(x_i))|$
 $\leq L |f(x_i) - f'(x_i)|. \quad \square$

Theorem Given $S = (x_1, \dots, x_n)$, $X \in \mathbb{R}^{n \times d}$ with rows x_i^T ,
 $\text{Rad}(\{x \mapsto \langle w, x \rangle : \|w\|_2 \leq B\} |_S) \leq B \|X\|_F.$

Pf Let $\varepsilon \in \{\pm 1\}^n$. Then,

$$\begin{aligned} \sup_{\|w\| \leq B} \sum_i \varepsilon_i \langle w, x_i \rangle &= \sup_{\|w\| \leq B} \langle w, \sum_i \varepsilon_i x_i \rangle \\ &= \sup_{\|w\| \leq B} \langle w, \sum_i \varepsilon_i x_i \rangle \\ &= \left\| \sum_i \varepsilon_i x_i \right\|_2 \end{aligned}$$

By Jensen's inequality (for convex ϕ , $\phi(\mathbb{E} X) \leq \mathbb{E} \phi(X)$; reversed for concave)

$$\mathbb{E} \|\sum_i \varepsilon_i x_i\|_2 = \mathbb{E} \sqrt{\|\sum_i \varepsilon_i x_i\|_2^2} \leq \sqrt{\mathbb{E} \|\sum_i \varepsilon_i x_i\|_2^2}.$$

$$\begin{aligned} \mathbb{E} \|\sum_i \varepsilon_i x_i\|_2^2 &= \mathbb{E} \left[\sum_i \varepsilon_i^2 \|x_i\|^2 + \sum_{i \neq j} \varepsilon_i \varepsilon_j \langle x_i, x_j \rangle \right] \\ &= \mathbb{E} \sum_i \|x_i\|^2 + 0 \\ &= \sum_i \|x_i\|^2 = \|x\|_F^2. \end{aligned}$$

$$\Rightarrow \text{Rad}(\{x \mapsto \langle w, x \rangle : \|w\|_2 \leq b\} | S) \leq \frac{1}{n} \mathbb{E} \|\sum_i \varepsilon_i x_i\| \leq \frac{\|x\|_F}{n}.$$