# Generalization & uniform convergence

**Def** A (centered) RV $X$ is $\sigma$-subGaussian (or $\sigma$-SG; variance proxy $\sigma^2$) if

$$\mathbb{E}\, e^{\lambda X} \leq e^{\lambda^2 \sigma^2 / 2}, \quad \forall \lambda \in \mathbb{R}.$$

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

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

**Pf**: Exercise. **Hint**: note that $\mathbb{P}(X \geq \varepsilon) = \inf_{t \geq 0} \mathbb{P}\left(e^{tX} \geq e^{t\varepsilon}\right)$.

— 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 $X_1, \ldots, X_n$ are indep. $\sigma_i$-SG RV's, then $Z := \sum_1^n X_i$ is SG with variance proxy $\sum_1^n \sigma_i^2$; and if $\alpha > 0$ then $\alpha X_i$ is $\alpha \sigma_i$-SG (variance proxy: $\alpha^2 \sigma_i^2$).

Thus, if

$\overline{X}_n := \frac{1}{n}\sum_i^n X_i$ / then $\overline{X}$ is SG w/ variance proxy $\frac{1}{n^2}\sum_i^n \sigma_i^2$.

$\Rightarrow$ Lemma says $\mathbb{P}(\overline{X}_n \geq \varepsilon) \leq \exp\left(\frac{-\varepsilon}{\frac{2}{n^2}\sum_i^n \sigma_i^2}\right) = \exp\left(\frac{-n^2\varepsilon}{2\sum_i^n \sigma_i^2}\right)$.

for indep $\sigma$-SG $X_i$,

Similarly, $\mathbb{P}(\overline{X}_n \leq -\varepsilon) \leq \exp\left(-n^2\varepsilon / 2\sum_i^n \sigma_i^2\right)$

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

$$\mathbb{P}(|\overline{X}_n| > \varepsilon) \leq 2\exp\left(-n^2\varepsilon / 2\sum_i^n \sigma_i^2\right).$$

Generalizing to non-mean zero, we get

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

Letting $\sigma^2 := \frac{1}{n}\sum_i^n \sigma_i^2$, 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 $2\exp\left(-n\varepsilon / 2\sigma^2\right) = \delta$, so

up $> 1-\delta$, $\left|\frac{1}{n}\sum_i^n (X_i - \mu_i)\right| \leq \sqrt{\frac{1}{n}\sum_i^n \sigma_i^2} \cdot \sqrt{\frac{2\log(2/\delta)}{n}}$.

If $X_i$ are iid, says up $> 1-\delta$, $\left|\mu - \frac{1}{n}\sum_i^n X_i\right| \leq \sqrt{\frac{\sigma^2 \log(2/\delta)}{n}}$.

each $\sigma$-SG,

$\longrightarrow$ as $n$ gets larger, sample mean closer to pop mean.

EX. Let $(X_i, Y_i) \overset{iid}{\sim} P$, $X \in \mathbb{R}^d$, $Y \in \{\pm 1\}$, $f: \mathbb{R}^d \to \{\pm 1\}$,

and let $Z_i := \mathbb{1}(f(X_i) \neq Y_i)$. Then each $Z_i$ is iid, bounded (hence SG: with variance proxy $\frac{1}{4}$). So by above,
in $[0,1]$

w.p. $> 1 - \delta$, $\quad \mathbb{P}(y \neq f(x)) \leq \frac{1}{n} \sum_1^n \mathbb{1}(y_i \neq f(x_i)) + \sqrt{\frac{\log(2/\delta)}{2n}}$.

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

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 $\frac{1}{n} \sum_1^n \mathbb{1}(y_i \neq f_n(x_i)) = 0$ for all $n$ by def$^n$.

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

② $X$ has continuous distn. Then $\frac{1}{n} \sum_1^n \mathbb{1}(y \neq f(x_i)) = 0$ by construction,

but $\mathbb{P}(y \neq f_n(x)) = 1$ th.

What broke sub-G concentration?

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

$$z_i := \mathbb{1}(y_i \neq f_n(x_i))$$ 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 $\ell(z_i)$,

$$L(f) := \mathbb{E} f(z), \quad \hat{L}(f) = \frac{1}{n} \sum_i^n f(z_i),$$

Goal: bound $L(f)$. Suppose $f \in \mathcal{F}$, some fcn class $\mathcal{F}$.
And suppose we use $S = \{z_i\}_i^n$ to fit $f = f(S)$. Then we typically lose indep. of $f(z_i; s)$.

Approach is then:

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

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

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 $\mathcal{F} = \{f_1, \ldots, f_k\}$, $|\mathcal{F}| = k$. If $(x_i, y_i)_i^n$ are iid, $f_i : \mathbb{R}^d \to \{\pm 1\}$, then SG concentration as before gives for fixed $f_\ell$,

$$\mathbb{P}\left( \left| \mathbb{P}(f_\ell(x) \neq y) - \frac{1}{n}\sum_1^n \mathbb{1}(y_i \neq f_\ell(x_i)) \right| > \sqrt{\frac{\log \frac{2}{\delta}}{2n}} \right) \leq \delta.$$

$(x_i, y_i)_i^n$

for fixed $\ell$, wp $> 1 - \delta$, $\left| \mathbb{P}(f_\ell(x) \neq y) - \hat{\mathbb{P}}(f_\ell(x) \neq y) \right| \leq \sqrt{\frac{\log \frac{2}{\delta}}{n}}$.

**Union bound:**

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

i.e. wp $> 1 - \delta$, for <u>all</u> $\ell \in [k]$, $\left| \mathbb{P}(y \neq f_\ell(x)) - \hat{\mathbb{P}}(f_\ell(x) \neq y) \right| \leq \sqrt{\frac{\log(2k/\delta)}{2n}}$.

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

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

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

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 $z_i = (x_i, y_i)$, $S = \{z_i\}_1^n$, fcn class $\mathcal{F}$,

$$\mathcal{F}_{|S} := \{ (f(z_1), \ldots, f(z_n)) : f \in \mathcal{F} \}.$$

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

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

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

- We'll often look at URad for _losses_, ie. for $\ell$,
  $$\text{URad}\big( (\ell \cdot f)_{|S} \big) = \text{URad}\big( (\ell(y_1, f(x_1)), \ldots, \ell(y_n, f(x_n))) : f \in \mathcal{F} \big).$$

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

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

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

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

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

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

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

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

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

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

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

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

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

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

$$\sup_{f \in \mathcal{F}} \left\{ \mathbb{E} f(z) - \frac{1}{n} \sum_1^n f(z_i) \right\} \leq \mathbb{E}_{z_i}\left( \sup_{f \in \mathcal{F}} \left\{ \mathbb{E} f(z) - \frac{1}{n} \sum_1^n f(z_i) \right\} \right) + (b-a)\sqrt{\frac{\log 1/\delta}{2n}}.$$

② w.p. $> 1-\delta$,

$$\mathbb{E}_{z_i} \mathrm{URad}(\mathcal{F}_{|s}) \leq \mathrm{URad}(\mathcal{F}_{|s}) + (b-a)\sqrt{\frac{n \log(1/\delta)}{2}}$$

③ w.p. $> 1-\delta$,

$$\sup_{f \in \mathcal{F}} \left\{ \mathbb{E} f(z) - \frac{1}{n} \sum_1^n f(z_i) \right\} \leq \frac{2}{n} \mathrm{URad}(\mathcal{F}_{|s}) + 3(b-a)\sqrt{\frac{\log \frac{2}{\delta}}{n}}$$

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

**Thm (MacDiarmid).** Suppose $F: \mathbb{R}^n \to \mathbb{R}$ satisfies <u>bounded differences</u>:
$\forall i \in [n]$, $\exists c_i$ $\mathfrak{st}$ $\displaystyle\sup_{z_1,..,z_n, z_i'} |F(z_1,..,z_i,z_{i+1},..,z_n) - F(z_1,..,z_i',..,z_n)| \leq c_i$. Then,

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

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

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

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

**Pf.**

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

$$\mathbb{E}f_\varepsilon = \mathbb{E}_{z\sim\mathbb{P}} f_\varepsilon(z) = \mathbb{E}_n' \hat{\mathbb{E}}_n' f_\varepsilon, \text{ since } z_0' \sim \mathbb{P} \text{ so } \mathbb{E}_{z_i'} f_\varepsilon(z_i') = \mathbb{E}f_\varepsilon.$$

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

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

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

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

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

Since $\varepsilon > 0$ is arbitrary, this completes the proof. $\boxed{}$