# 10. Rademacher, part 2

Sunday, January 23, 2022    7:49 PM

Interlude : McDiarmid's Inequality

- Thm (McDiarmid)

    Let $f$ be a function s.t. $\exists \; c_i < \infty$ s.t.

    $\forall i \in [m],$ $\quad | f(x_1, \ldots, x_i, \ldots, x_m) - f(x_1, \ldots, x_i', \ldots, x_m) | \leq c_i$

    "bounded difference"

    for arbitrary $(x_1, \ldots, x_m)$ and $x_i'$

    Let $S = (X_1, \ldots, X_m)$        prob. notation

    independent r.v.

    Then $\quad \mathbb{P}\left[ f(S) - \mathbb{E} f(S) \geq \varepsilon \right] \leq \exp\left( \frac{-2\varepsilon^2}{\Sigma_i c_i^2} \right)$

    $\qquad \mathbb{P}\left[ f(S) - \mathbb{E} f(S) \leq -\varepsilon \right] \leq \text{``} \underline{\hspace{2cm}} \text{''}.$

    (in Hoeffding, $f(S) = \frac{1}{m} \Sigma_i' X_i$ ... $f$ is linear in $S$ )


Thm ( main R.C. generalization bound )        Thm 3.3 in Mohri

   Let $\mathcal{F}$ be a family of functions from $X$ to $[0,1]$ (or really $[a,b]$ works)

   then $\forall \; \delta > 0$, w.p. $\geq 1 - \delta$ ( w.r.t. $S \sim D^m$ iid )        [any D]

   $(\forall f \in \mathcal{F})$

   $\underset{z \sim D}{\mathbb{E}} f(z) - \frac{1}{m} \sum_{i=1}^{m} f(z_i) \leq \begin{cases} \text{Ⓐ} \quad 2\,\mathcal{R}_m(\mathcal{F}) + \sqrt{\dfrac{\log(1/\delta)}{2m}} \\[4mm] \text{Ⓑ} \quad 2\,\hat{\mathcal{R}}_S(\mathcal{F}) + 3 \cdot \sqrt{\dfrac{\log(2/\delta)}{2m}} \end{cases}$

   risk if $f = \ell \circ h$     empirical risk

   D shows up here    Right Hand Side  RHS

   data dependent (sometimes nice, sometimes not)

   ( see also Thm. 26.5 in [SS]...
        [SS] usually sets RHS $= \varepsilon(m)$, solves for $m(\varepsilon)$
        which is nice to have but requires many approximations/bounds
        to make it tractable )

   proof   Recall $\text{Rep}_D(\mathcal{F}, S) = \underset{f \in \mathcal{F}}{\sup} \; \underset{z \sim D}{\mathbb{E}} f(z) - \frac{1}{m} \sum_{i=1}^{m} f(z_i)$

   $\mathbb{E}_S f$  shorthand

   so by observation, we just need $\text{Rep}_D(\mathcal{F}, S) \leq \text{RHS}$.

   Idea: use McDiarmid's, applied to $\text{Rep}_D(\mathcal{F}, S) = \text{Rep}(S)$

   drop D, $\mathcal{F}$ notation since those are fixed

Check the assumptions:

Let $S = S'$ except some $z_i \neq z_i'$

we'll use 
$$\sup_a g(a) = \sup_a ( g(a) - h(a) + h(a) )$$
$$\leq \sup_a ( g(a) - h(a) ) + \sup_a h(a) \qquad \text{like earlier}$$

$$\text{Rep}(S) - \text{Rep}(S') := \left( \sup_{f \in \mathcal{F}} \mathbb{E} f - \hat{\mathbb{E}}_S f \right) - \left( \sup_{f \in \mathcal{F}} \mathbb{E} f - \hat{\mathbb{E}}_{S'} f \right)$$

$$\leq \sup_{f \in \mathcal{F}} \left( ( \mathbb{E} f - \hat{\mathbb{E}}_S f ) - ( \mathbb{E} f - \hat{\mathbb{E}}_{S'} f ) \right) \qquad \text{via}$$

$$= \sup_{f \in \mathcal{F}} \left( \frac{1}{m} ( f(z_i') - f(z_i) ) \right)$$
$$\underbrace{\qquad\qquad}_{\leq 1} \quad \text{if } f : X \to [0,1]$$

$$\leq 1/m .$$

So $\forall i \in [m]$, the "$c_i$" in McDiarmid's is $1/m$

Applying McDiarmid's

$$\mathbb{P}\left[ \text{Rep}(S) - \mathbb{E} \text{Rep}(S) \geq \varepsilon \right] \leq \exp\left( \frac{-2\varepsilon^2}{\sum_i (1/m)^2} \right) = \exp( -2 \varepsilon^2 m )$$
$$=: \delta$$

Now fix $\delta$, solve for $\varepsilon$
$$-2\varepsilon^2 m = \log(\delta) , \qquad \varepsilon = \sqrt{\frac{-\log(\delta)}{2m}} = \sqrt{\frac{\log(1/\delta)}{2m}}$$

So w.p. $\geq 1-\delta$,
$$\text{Rep}(S) \leq \mathbb{E}_{S \sim \mathcal{D}^m} \text{Rep}(S) + \varepsilon$$

$$\leq 2 \mathcal{R}_m(\mathcal{F}) + \varepsilon \qquad \text{via previous lemma.}$$
$$\text{which proves } \boxed{A}$$

To prove $\boxed{B}$, restate $\boxed{A}$ using $\delta/2$ :

w.p. $\geq 1 - \delta/2$, $\quad \text{Rep}(S) \leq 2 \mathcal{R}_m(\mathcal{F}) + \sqrt{\frac{\log(2/\delta)}{2m}}$ $\qquad (**)$

Recall $\mathcal{R}_m(\mathcal{F}) := \mathbb{E}_S \hat{\mathcal{R}}_S(\mathcal{F})$
$\underbrace{\qquad}$ write as $\hat{\mathcal{R}}(S)$ since $\mathcal{F}$ fixed for now

$$\hat{\mathcal{R}}(S) := \mathbb{E}_\sigma \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i f(z_i)$$

Let's see if we can't apply McDiarmid's again, but this time the $\mathbb{P}[ \dots \leq -\varepsilon ] \leq \dots$ side instead of $\mathbb{P}[ \dots \geq \varepsilon ] \leq \dots$

Turns out $\hat{\mathcal{R}}(S) - \hat{\mathcal{R}}(S') \leq 1/m$ $\qquad S = (z_1, \dots, z_m)$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad S' = (z_1, \dots, z_i', \dots, z_m)$

just like for Rep(S) ( only difference is now $\mathbb{E}_\sigma$ (...) which

is linear and doesn't affect anything )

So

McDiarmid: $\mathbb{P}\left[ \hat{R}(S) - \underbrace{\mathbb{E}\hat{R}(S)}_{R_m} \leq -\varepsilon \right] \leq \exp(-2\varepsilon^2 m) =: \delta/2$

i.e. $\mathbb{P}\left[ \hat{R}(S) - R_m(\mathcal{F}) > -\varepsilon \right] \geq 1 - \delta/2$

So $\varepsilon = \sqrt{\dfrac{\log(2/\delta)}{2m}}$

So (**)

w.p. $\geq 1 - \delta/2$    $\text{Rep}(S) \leq 2 \cdot R_m(\mathcal{F}) + \sqrt{\dfrac{\log(2/\delta)}{2m}}$

and

w.p. $\geq 1 - \delta/2$    $R_m(\mathcal{F}) < \hat{R}_S(\mathcal{F}) + \sqrt{\dfrac{\log(2/\delta)}{2m}}$

So combining:

w.p. $\geq 1 - \delta,$    $\text{Rep}(S) \leq 2 \hat{R}_S(\mathcal{F}) + 3 \cdot \sqrt{\dfrac{\log(2/\delta)}{2m}}$   Ⓑ  □

**Specific Application:** Binary classification, 0-1 loss   ( Lemma 3.4 Mohri )

$\mathcal{H}$ a set of functions (hypotheses) of the form $h: X \longrightarrow \mathcal{Y} := \{\pm 1\}$

and use 0-1 loss as usual, $\ell(h, (x,y)) := \mathbb{1}_{h(x) \neq y}$

$S = ( (x_1, y_1), \dots, (x_m, y_m) )$

$S_x = ( x_1, \dots, x_m )$

Observe $\underbrace{\hat{R}_S(\ell \circ \mathcal{H})}_{\substack{\text{what we care} \\ \text{about}}} := \mathbb{E}_\sigma \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i \cdot \mathbb{1}_{h(x_i) \neq y_i}$

check:
$y_i h(x_i) = \pm 1$
$+1$ if agree
$-1$ else

$= \underbrace{\cancel{1} - y_i h(x_i)}_{} \over 2$

$= \frac{1}{2} \mathbb{E}_\sigma \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m h(x_i)$

$\mathbb{E}\, \sigma_i \cdot 1 = 0$
$\mathbb{E} - \sigma_i y_i = \mathbb{E} \sigma_i y_i$
$\mathbb{E}\, \sigma_i y_i h(x_i) = \mathbb{E} \sigma_i h(x_i)$
since $y_i \in \{\pm 1\} \Rightarrow \sigma_i y_i \sim \text{Rad}.$

$= \frac{1}{2} \hat{R}_{S_x}(\mathcal{H})$

No reference to labels or loss  (0-1 is implicit)
only $\mathcal{H}$, $\mathcal{D}_x$

So $\hat{R}_S(\ell \circ \mathcal{H}) = \frac{1}{2} \hat{R}(\mathcal{H})$

and similarly

$R_m(\ell \circ \mathcal{H}) = \frac{1}{2} R_m(\mathcal{H})$

[ implicit:
labels
$y = f(x_i)$ ]
Not agnostic

and now apply our Thm to this case

**Thm 3.5 [Mohri]** Rademacher Complexity thm for binary classif. w/ 0-1 loss

Let $\mathcal{H} \subseteq \{\pm 1\}^X$, $\mathcal{D}$ any distribution $X$,                 Then $\forall \delta \in (0,1),$

w.p. $\geq 1 - \delta$ (wrt $S_x \sim \mathcal{D}^m$ iid), $\forall h \in \mathcal{H}$

$$L_D(h) \leq \hat{L}_S(h) + \begin{cases} R_m(\mathcal{H}) + \sqrt{\dfrac{\log(1/\delta)}{m}} \\[2ex] \hat{R}_{S_x}(\mathcal{H}) + 3 \cdot \sqrt{\dfrac{\log(2/\delta)}{2m}} \end{cases}$$

proof via previous lemma.

Typically $R_m(\mathcal{H}) = O(1/m)$ or $O(1/\sqrt{m})$. Either way,

  we're bounding our estimation error by $O(1/\sqrt{m})$

  so need $m = \Omega\left(1/\varepsilon_{est}^2\right)$

$R_m(\mathcal{H})$ is one measure of complexity of $\mathcal{H}$

Also, we could possibly compute it via optimization:

$$\hat{R}_S(\mathcal{H}) = \mathbb{E}_\sigma \sup_{f \in \mathcal{H}} \frac{1}{m} \sum' \sigma_i f(z_i)$$

$$= - \mathbb{E}_\sigma \inf_{f \in \mathcal{H}} \frac{1}{m} \sum' \tilde{\sigma}_i f(z_i) \qquad \tilde{\sigma}_i = -\sigma_i \text{, both} \sim \text{Rademach}$$

  but at least as hard as solving

$$\inf_{f \in \mathcal{H}} \frac{1}{m} \sum' f(z_i) = ERM_{\mathcal{H}}$$

  and as we'll see, that's often intractable.

  But for simple $\mathcal{H}$, we can often compute (or at least
    upper bound) $\hat{R}_S$ or $R_m$ by hand