

Rademacher Complexity based learning bounds

$$\text{Let } T = \sup_{f \in \mathcal{F}} \left\{ R[f] - R_{\text{emp}}[f] \right\}$$

$$\hookrightarrow \frac{1}{m} \sum_{i=1}^m \ell(f(x_i), y_i)$$

Let us denote $z_i = (x_i, y_i)$ & $h_f(z_i) = \ell(f(x_i), y_i)$
 \hookrightarrow training examples (together with labels)

It is easy to see that T is a function of z_1, z_2, \dots, z_m , which are

iid! Let $T = g(z_1, z_2, \dots, z_m) = \sup_{f \in \mathcal{F}} \left\{ R[f] - R_{\text{emp}}[f] \right\}$

Now let us also assume that $a \leq \ell(f(x), y) \leq b$

i.e. assume loss function is bounded. $\forall f \in \mathcal{F}$ & $\forall (x, y)$

Then it is easy to see that T (which is a function of iid r.v.s) satisfies the following property: (known as Bounded difference property)

$$\max_{z_i, z'_i} \left| g(z_1, z_2, \dots, z_i, z_{i+1}, \dots, z_m) - g(z_1, z_2, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_m) \right| \leq \frac{b-a}{m} \quad \forall i$$

Here's why: $\left| g(z_1, \dots, z_i, \dots, z_m) - g(z_1, \dots, z'_i, \dots, z_m) \right|$

$$\leq \left\{ R[f^*] - R_{\text{emp}}[f^*] \right\} - \left(R[f^*] - \frac{1}{m} \sum_{j \neq i} h_{f^*}(z_j) - \frac{1}{m} h_{f^*}(z'_i) \right)$$

note this \leftarrow

where $f^* = \arg \max_{f \in \mathcal{F}} \left\{ R[f] - R_{\text{emp}}[f] \right\}$

$$= \frac{1}{m} (h_{f^*}(z_i) - h_{f^*}(z'_i)) \leq \frac{b-a}{m} \quad (\text{since } h_{f^*} \text{ is bounded between } a \text{ \& } b)$$

Similarly

$$\begin{aligned} & g(z_1, z_2, \dots, z_i, \dots, z_m) - g(z_1, \dots, z'_i, \dots, z_m) \\ & \leq \left(R[f^*] - \frac{1}{m} \sum_{j \neq i} h_{f^*}(z_j) - \frac{1}{m} h_{f^*}(z_i) \right) - \left(R[f^*] - \frac{1}{m} \sum_{j \neq i} h_{f^*}(z_j) - \frac{1}{m} h_{f^*}(z'_i) \right) \\ & = \frac{1}{m} (h_{f^*}(z_i) - h_{f^*}(z'_i)) \leq \frac{b-a}{m} \quad \text{here, } f^* = \arg \max_{f \in \mathcal{F}} \end{aligned}$$

$$\therefore |g(z_1, \dots, z_i, z_{i+1}, \dots, z_m) - g(z_1, \dots, z'_i, \dots, z_m)| \leq \frac{b-a}{m}$$

Hence T is a function of m iid rvs satisfies the bounded difference property. $\forall z_i, z'_i$

→ We saw that Chernoff bounding gives efficient bounds in case of rvs which are ^{in turn} sum of independent rvs.

Using similar ideas one can handle generic functions of independent rvs (which satisfy properties like bounded difference) using the McDiarmid's

inequality:

McDiarmid's Inequality: Let g be a function of m independent rvs z_1, \dots, z_m satisfying Bounded difference property with bound $c (= \frac{b-a}{m})$

$$P[g - E[g] > \epsilon] \leq e^{-\frac{2\epsilon^2}{mc^2}} \quad \text{need not be iid}$$

$$P[E[g] - g > \epsilon] \leq e^{-\frac{2\epsilon^2}{mc^2}} \quad \text{in our case } c = \frac{b-a}{m}$$

(!!!) we can also show that $P[E[g] - g > \epsilon] \leq e^{-\frac{2\epsilon^2}{mc^2}}$

Note that McDiarmid's inequality is an extension of Chernoff's bounding technique from the case of bounding sum of independent RVs to handling generic functions of independent RVs satisfying mild conditions like the bounded difference property.

Proof: $T = g(z_1, \dots, z_m)$ & g satisfies bounded diff. property.

Recall that Chernoff bounds apply to case of sum of independent RVs. So the idea here is to rewrite T as sum of independent RVs so that we can hopefully apply Chernoff techniques.

In view of this, define $V_i = E[T/z_1, \dots, z_i] - E[T/z_1, \dots, z_{i-1}]$.

$$V_1 = E[T/z_1] - E[T] \xrightarrow{\text{func. of } z_1} z_1$$

↪ also known as
Martingale
difference

$$V_2 = E[T/z_1, z_2] - E[T/z_1] \xrightarrow{\text{func. of } z_1, z_2} z_1, z_2$$

⋮

$$V_i = E[T/z_1, \dots, z_i] - E[T/z_1, \dots, z_{i-1}] \xrightarrow{\text{func. of } z_1, z_2, \dots, z_i} z_1, z_2, \dots, z_i$$

⋮

$$V_m = E[T/z_1, \dots, z_m] - E[T/z_1, \dots, z_{m-1}] \xrightarrow{\text{func. of } z_1, z_2, \dots, z_m} z_1, z_2, \dots, z_m$$

$$\sum_{i=1}^m V_i = E[T/z_1, \dots, z_m] - E[T]$$

$$= T - E[T]$$

Moreover, V_1, V_2, \dots, V_m are conditionally independent:

$$\begin{array}{l}
 V_1 \rightarrow f.c. \text{ of } Z_1 \\
 V_2/Z_1 \rightarrow f.c. \text{ of } Z_2 \\
 \vdots \\
 V_i/Z_1, Z_2, \dots, Z_{i-1} \rightarrow f.c. \text{ of } Z_i \\
 \vdots \\
 V_m/Z_1, \dots, Z_{m-1} \rightarrow f.c. \text{ of } Z_m
 \end{array}
 \left. \vphantom{\begin{array}{l} V_1 \\ V_2 \\ \vdots \\ V_i \\ \vdots \\ V_m \end{array}} \right\} \begin{array}{l} \text{all are} \\ \text{independent!} \\ \text{(since } Z_i \text{ are ind.)} \end{array}$$

Now the idea is to apply Chernoff bounding technique to $\sum V_i$ noting the conditional independence.

$$\begin{aligned}
 P[g - E[g] > \varepsilon] &= P\left[\sum_i V_i > \varepsilon\right] = P\left[e^{\sum V_i} > e^{\varepsilon}\right] \quad (\Delta > 0) \\
 &\leq \frac{E\left[e^{\sum V_i}\right]}{e^{\varepsilon}} \quad (\because \text{Markov inequality}) \\
 &= \frac{E\left[\prod_{i=1}^m e^{\Delta V_i}\right]}{e^{\varepsilon}}
 \end{aligned}$$

~~Note~~ Now we exploit cond. independence to bound mgf:

$$\begin{aligned}
 E\left[\prod_{i=1}^m e^{\Delta V_i}\right] &= E\left[E\left[\prod_{i=1}^m e^{\Delta V_i} / Z_1, \dots, Z_{m-1}\right]\right] \\
 &\quad \underbrace{\hspace{10em}}_{\text{total expectation rule}} \\
 &= E\left[E\left[\left(\prod_{i=1}^{m-1} e^{\Delta V_i}\right) e^{\Delta V_m} / Z_1, \dots, Z_{m-1}\right]\right] \\
 &= E\left[\prod_{i=1}^{m-1} e^{\Delta V_i} E\left[e^{\Delta V_m} / Z_1, \dots, Z_{m-1}\right]\right] \quad (\because \text{conditional independence})
 \end{aligned}$$

Now, we will upper bound
 $E[e^{SV_m} / z_1, \dots, z_{m-1}]$

9. In general, we ~~will~~ ^{can} upper bound:

$$E[e^{SV_i} / z_1, \dots, z_{i-1}]$$

$\forall i=2 \text{ to } m$

The idea is to note that $V_i / z_1, \dots, z_{i-1}$ is a mean zero RV
 \mathcal{B} is bounded (with interval length c again) and then apply

Hoeffding inequality:

Now,

$$\begin{aligned} E[V_i / z_1, \dots, z_{i-1}] &= E\left[\left(E[T / z_1, \dots, z_i] - E[T / z_1, \dots, z_{i-1}]\right) / z_1, \dots, z_{i-1}\right] \\ &= E\left[E[T / z_1, \dots, z_i] / z_1, \dots, z_{i-1} - E\left[E[T / z_1, \dots, z_{i-1}] / z_1, \dots, z_{i-1}\right]\right] \\ &\quad \downarrow \text{total prob. rule} \qquad \downarrow \\ &= E[T / z_1, \dots, z_{i-1}] - E[T / z_1, \dots, z_{i-1}] \\ &= 0 \end{aligned}$$

Also, we know T satisfies bounded difference

$$\max_{z_i, z_i'} |g(z_1, \dots, z_i, \dots, z_m) - g(z_1, \dots, z_i', \dots, z_m)| \leq c$$

$$\Rightarrow \max_{z_i, z_i'} |E[T / z_1, \dots, z_i] - E[T / z_1, \dots, z_i']| \leq c$$

$$\Rightarrow \max_{z_i, z_i'} \left| \underbrace{E[T / z_1, \dots, z_i] - E[T / z_1, \dots, z_{i-1}]}_{V_i / z_1, \dots, z_{i-1}} - \underbrace{(E[T / z_1, \dots, z_i] - E[T / z_1, \dots, z_{i-1}])}_{V_i' / z_1, \dots, z_{i-1}} \right| \leq c$$

\nearrow add-subtract the same term

$$\Rightarrow \max_{z_i, z'_i} \left| V_i / z_{1, \dots, z_{i-1}} - V'_i / z_{1, \dots, z_{i-1}} \right| \leq c$$

$$\Rightarrow \max_{z_i} V_i / z_{1, \dots, z_{i-1}} - \min_{z'_i} V'_i / z_{1, \dots, z_{i-1}} \leq c$$

Here $V_i / z_{1, \dots, z_{i-1}}$ is has mean zero randoms in interval of length $\leq c$.

$$\therefore E[e^{\Delta V_i} / z_{1, \dots, z_{i-1}}] \leq e^{\frac{c^2}{8}}$$

$$\Rightarrow E\left[\prod_{i=1}^m e^{\Delta V_i}\right] = E\left[\prod_{i=1}^{m-1} e^{\Delta V_i} E\left[e^{\Delta V_m} / z_{1, \dots, z_{m-1}}\right]\right]$$

$$\leq e^{\frac{c^2}{8}} E\left[\prod_{i=1}^{m-1} e^{\Delta V_i}\right]$$

$$\leq e^{\frac{c^2 m}{8}}$$

(repeated application of conditional independence)

$$\Rightarrow P[g - E[g] > \epsilon] \leq \frac{e^{\frac{\epsilon^2 m}{8}}}{e^{\epsilon}} \quad (\ast) > 0$$

$$\therefore P[g - E[g] > \epsilon] \leq e^{\min_{\epsilon} \left\{ \frac{\epsilon^2 m}{8} - \epsilon \right\}} = e^{-\frac{2\epsilon^2}{mc^2}}$$

Here McDiarmid's inequality is proved.

Proof for $P\{E[g] - g > \epsilon\} \leq e^{-\frac{2\epsilon^2}{mc^2}}$ is also similar and is skipped here.

Getting back to our problem we have $T = \sup_{f \in \mathcal{F}} \{R(f) - R_p^m(f)\}$ and it satisfies bounded diff. prob with $c = \frac{b-a}{m}$. \therefore By McDiarmid's inequality we have:

$$P\left[\sup_{f \in \mathcal{F}} \{R(f) - R_p^m(f)\} - E\left[\sup_{f \in \mathcal{F}} \{R(f) - R_p^m(f)\}\right] > \epsilon\right] \leq e^{-\frac{2m\epsilon^2}{(b-a)^2}} = \delta(\epsilon, m)$$

Now re-write ϵ in terms of δ : $\epsilon = (b-a) \sqrt{\frac{\ln 1/\delta}{2m}}$

$\Rightarrow \forall f \in \mathcal{F}$, with at least prob. $1-\delta$, we have:

$$R[f] \leq R_{\text{emp}}^m[f] + E \left[\sup_{f \in \mathcal{F}} \{ R[f] - R_{\text{emp}}^m[f] \} \right] + (b-a) \sqrt{\frac{\ln 1/\delta}{2m}}$$

(New leading bound)

~~key differences with VC-type bounds:~~

~~confidence term~~

difficult to compute
 here we will further
 upper bound this term
 by a term which reflects
 complexity of \mathcal{F} (resembles VC-dim bound)

confidence term
 (free of complexity
 terms like
 VC dim. etc.)

So the task now is to bound the $E[T]$.

$$E[T] = E \left[\sup_{f \in \mathcal{F}} \{ R[f] - R_{\text{emp}}^m[f] \} \right]$$

$$= E \left[\sup_{f \in \mathcal{F}} \left\{ E \left[\bar{R}_{\text{emp}}^m[f] \right] - R_{\text{emp}}^m[f] \right\} \right]$$

$$= E \left[\sup_{f \in \mathcal{F}} \left\{ E \left[\bar{R}_{\text{emp}}^m[f] - R_{\text{emp}}^m[f] \right] \right\} \right]$$

$$\leq E \left[E \left[\sup_{f \in \mathcal{F}} \left\{ \bar{R}_{\text{emp}}^m[f] - R_{\text{emp}}^m[f] \right\} \right] \right] \quad (\because \text{Jensen's inequality})$$

with Z_i with Z_i'

(Symmetrization trick)

$$\text{define } \bar{R}_{\text{emp}}^m[f] = \frac{1}{m} \sum_{i=1}^m h_f(Z_i')$$

Z_i' are also iid
 copies of Z_i .

Now we will again resort to ~~the~~ Rademacher variables
 ↓ introducing

$$E\{T\} \leq E \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{m} \sum_i (h_f(z'_i) - h_f(z_i)) \right\} \right]$$

wrt z_i, z'_i

$$= E \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{m} \sum_i \sigma_i (h_f(z'_i) - h_f(z_i)) \right\} \right]$$

wrt $z_i, z'_i, \sigma_i \rightarrow$ Rademacher r.v.s.

$$\leq E \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{m} \sum_i \sigma_i h_f(z'_i) \right\} \right] + E \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{m} \sum_i \sigma_i h_f(z_i) \right\} \right]$$

wrt z'_i, σ_i

wrt z_i, σ_i

$$= 2 E \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{m} \sum_i \sigma_i h_f(z_i) \right\} \right]$$

Now given a loss function, & set of learning functions \mathcal{F} ,
 we get an induced set of loss functions $\rightarrow \left\{ h_f / \ell(f(x), y) = h_f(z) \right\}$
 $\forall f \in \mathcal{F}$
 let us call this set as \mathcal{L}

$$\Rightarrow E\{T\} \leq 2 E \left[\sup_{h \in \mathcal{L}} \left\{ \frac{1}{m} \sum_{i=1}^m \sigma_i h(z_i) \right\} \right]$$

wrt z_i, σ_i

known as Rademacher Complexity
 or
 Rademacher average of function class \mathcal{L}
 denoted by $\mathcal{R}(\mathcal{L})$.

$$\therefore \text{We have: } R[f] \leq R_{\text{emp}}[f] + 2 \mathcal{R}(\mathcal{L}) + \frac{(\epsilon - \alpha) \sqrt{\ln 1/\delta}}{2m}$$

↓
 purely a property of \mathcal{L} & underlying distribution.

Now $R(Z)$ is also difficult to compute as underlying distribution of Z_1, \dots, Z_m is unknown.

So we further upper bound it by quantity which is easily computable:

$$R_m(Z) = E \left[\sup_{h \in \mathcal{H}} \left\langle \frac{1}{m} \sum_{i=1}^m \sigma_i h(Z_i) \right\rangle \middle| Z_1, \dots, Z_m \right]$$

nothing but conditional Rademacher Average
 wrt σ_i alone

$R_m(Z)$ is easy to compute (at least numerically) since the distribution of σ_i is known.
 average computed for given training dataset Z_1, \dots, Z_m

Now we need to bound $R(Z)$ by terms involving $R_m(Z)$

For this we note that

- (i) $E[R_m(Z)] = R(Z)$ (\because total prob. rule)
- (ii) $R_m(Z)$ is a function of Z_1, \dots, Z_m (which are iid) and satisfies bounded difference property.

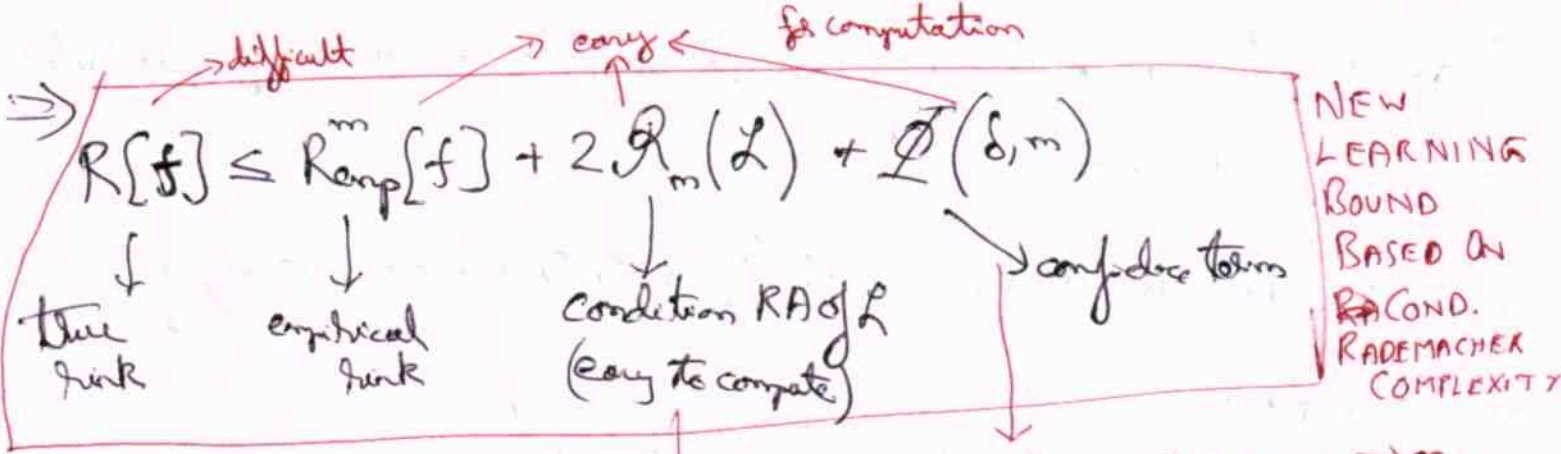
Student must show this

\therefore We can use McDiarmid's inequality again:

$$P[R(Z) - R_m(Z) > \epsilon] \leq e^{-\frac{2\epsilon^2 m}{(b-a)^2}} = \delta \text{ (say)}$$

$$\downarrow$$

$$E[R_m(Z)] \Rightarrow R(Z) \leq R_m(Z) + (b-a) \sqrt{\frac{\log 1/\delta}{2m}}$$



we know $\rightarrow 0$ as $m \rightarrow \infty$
 does this go to 0 as $m \rightarrow \infty$?
 If no, then ERM is constant!

We will now take a specific case of binary classification. i.e. \mathcal{F} as set of all functions taking values $\{-1, 1\}$ and let loss function be 0-1 loss.

$$l(f(x), y) = \frac{1}{2} (1 - y f(x)) = \begin{cases} 1 & \text{if } y \neq f(x) \\ 0 & \text{if } y = f(x) \end{cases}$$

$$\begin{aligned}
 \mathcal{R}_m(\mathcal{L}) &= E \left[\sup_{h \in \mathcal{L}} \frac{1}{m} \sum_i \sigma_i h(z_i) \right] \\
 &= E \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{m} \sum_i \sigma_i \frac{1}{2} (1 - y_i f(x_i)) \right\} \right] \\
 &= \frac{1}{2} E \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{m} \sum_i -\sigma_i y_i f(x_i) \right\} \right] \quad \text{II} \\
 &= \frac{1}{2} E \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{m} \sum_i \sigma_i f(x_i) \right\} \right] \quad \text{absgion} \\
 &= \frac{1}{2} \mathcal{R}_m(\mathcal{F}) \rightarrow \text{easily computed given training data points alone.}
 \end{aligned}$$

∴ For binary classification with 0-1 loss we have learning bound:

$$R[f] \leq R_{\text{exp}}[f] + \mathcal{R}_m(\mathcal{F}) + \mathcal{I}(\delta, m)$$

↓
Cod. RA of net classifiers.

In fact one can show that:

$$\mathcal{R}_m(\mathcal{F}) \leq \sqrt{\frac{h}{m}} \xrightarrow{\text{non constant}} \text{VC dim. of } \mathcal{F}!$$

↳ This recovers our good old bounds dependent on VC dim. releases extra term of $(\log \frac{m}{h} + 1)$, is termed. (which is very good)

→ Intuitively also **(II)** says why $\mathcal{R}_m(\mathcal{F})$ is a measure of "complexity" of net of classifiers.

If \mathcal{F} were rich to match $\sigma_i \gamma_i$ then we would have $\mathcal{R}_m(\mathcal{F}) = \frac{1}{2}$ ^{always $\leq m$} $\forall m$ which $\not\rightarrow 0$ as $m \rightarrow \infty$ ∴ ERM is not consistent in this case & is expected. (||h to VC dim = ∞)

→ Now ~~lets compute~~ SRM principle will suggest minimizing Empirical risk as well as $\mathcal{R}_m(\mathcal{F})$ simultaneously.

Now lets get ~~a~~ bound on $\mathcal{R}_m(\mathcal{F})$ for hyperplane classifiers in feature space induced by a kernel k .

we want something (previously we got $h \leq \|w\|^2 k^2$)
numbers

Cond. R.A. for Canonical hyperplane classifiers in feature space (RKHS)

Suppose we have $\mathcal{F} = \{f \in \mathcal{H} \mid \|f\|_{\mathcal{H}} \leq W, \}$ RKHS induced by
 $f(x) = \sum_{i=1}^m \alpha_i k(\cdot, x_i)$ given kernel 'k'
 we wish to find $R_m(\mathcal{F})$. & at least upper bound it.

$$R_m(\mathcal{F}) = E_{\sigma_i} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{m} \sum_i \sigma_i f(x_i) \right\} \right]$$

$$f(x_i) = \langle f, k(x_i, \cdot) \rangle_{\mathcal{H}}$$

$$= \frac{1}{m} E_{\sigma_i} \left[\sup_{f \in \mathcal{F}} \left\{ \langle f, \sum_i \sigma_i k(x_i, \cdot) \rangle_{\mathcal{H}} \right\} \right] \quad (\because \text{linearity of } \langle \cdot, \cdot \rangle_{\mathcal{H}})$$

$$\leq \|f\|_{\mathcal{H}} \left\| \sum_i \sigma_i k(x_i, \cdot) \right\|_{\mathcal{H}} \quad \text{by Cauchy-Schwarz inequality.}$$

$$\leq \frac{1}{m} E_{\sigma_i} \left[\sup_{f \in \mathcal{F}} \left\{ \|f\|_{\mathcal{H}} \sqrt{\sum_i \sum_j \sigma_i \sigma_j k(x_i, x_j)} \right\} \right]$$

$$\leq \frac{W}{m} E_{\sigma_i} \left[\sqrt{\sum_i \sum_j \sigma_i \sigma_j k(x_i, x_j)} \right]$$

$$\leq \frac{W}{m} \sqrt{\sum_i \sum_j E[\sigma_i \sigma_j k(x_i, x_j)]}$$

(\because Jensen's inequality)

$$= \frac{W}{m} \sqrt{\text{trace}(K)}$$

\rightarrow Gram matrix with kernel 'k' of training data points.