

3.8 Rademacher Complexity

Rademacher complexity helps to measure the complexity of a hypothesis class.

① Set up:

Define $G_n := \sup_{h \in H} L(h) - \hat{L}(h)$, then G_n is r.v. depending on Z_1, \dots, Z_n (data).

Define $G'_n := \sup_{h \in H} \hat{L}(h) - L(h)$, so that

$$\begin{aligned} P\{L(\hat{h}) - L(h^*) \geq \epsilon\} &\leq P\left\{\sup_{h \in H} |L(h) - \hat{L}(h)| \geq \frac{\epsilon}{2}\right\} \\ &\leq P\{G_n \geq \frac{\epsilon}{2}\} + P\{G'_n \geq \frac{\epsilon}{2}\} \end{aligned}$$

② Concentration:

Let g be the deterministic function s.t. $G_n = g(Z_1, \dots, Z_n)$,

Then g satisfies the bounded differences condition:

$$|g(Z_1, \dots, Z_i, \dots, Z_n) - g(Z_1, \dots, Z'_i, \dots, Z_n)| \leq \frac{1}{n}$$

Pf: $\hat{L}(h) := \frac{1}{n} \sum_{i=1}^n \ell(Z_i, h)$, we have

$$\begin{aligned} &\left| \sup_{h \in H} \{L(h) - \hat{L}(h)\} - \sup_{h \in H} \{L(h) - \hat{L}(h) + \frac{1}{n} [\ell(Z_i, h) - \ell(Z'_i, h)]\} \right| \leq \frac{1}{n} \\ \Rightarrow & |g(Z_1, \dots, Z_i, \dots, Z_n) - g(Z_1, \dots, Z'_i, \dots, Z_n)| \leq \frac{1}{n} \quad \square \end{aligned}$$

Now apply McDiarmid's inequality:

$$P[G_n \geq E[G_n] + \epsilon] \leq \exp\{-2n\epsilon^2\}.$$

③ Symmetrization. Bound $E[G_n]$.

Introduce ghost dataset Z'_1, \dots, Z'_n drawn i.i.d. from p^* .

Let $\hat{L}'(h) = \frac{1}{n} \sum_{i=1}^n \ell(Z'_i, h)$

Rewriting $L(h)$ in terms of the ghost dataset.

$$E[G_n] = E\left[\sup_{h \in H} L(h) - \hat{L}(h)\right] = E\left[\sup_{h \in H} \hat{L}'(h) - \hat{L}(h)\right]$$

$$\begin{aligned}
\mathbb{E}[G_n] &= \mathbb{E} \left[\sup_{h \in H} \mathbb{E}[\hat{L}'(h) - \hat{L}(h) | Z_{1:n}, Z'_{1:n}] \right] \\
&= \mathbb{E} \left[\sup_{h \in H} \mathbb{E}[\hat{L}'(h) - \hat{L}(h) | Z_{1:n}] \right] \\
&\leq \mathbb{E} \left[\mathbb{E} \left[\sup_{h \in H} \hat{L}'(h) - \hat{L}(h) | Z_{1:n} \right] \right] \\
&= \mathbb{E} \left[\sup_{h \in H} \hat{L}'(h) - \hat{L}(h) \right]
\end{aligned}$$

$$\Rightarrow \mathbb{E}[G_n] \leq \mathbb{E} \left[\sup_{h \in H} \frac{1}{n} \sum_{i=1}^n [\ell(Z'_i, h) - \ell(Z_i, h)] \right]$$

To remove the dependence on the ghost dataset $Z'_{1:n}$, we introduce i.i.d. Rademacher variables $\sigma_1, \dots, \sigma_n$ independent of $Z_{1:n}, Z'_{1:n}$, where σ_i is uniform over $\{-1, +1\}$. Since $\ell(Z'_i, h) - \ell(Z_i, h)$ is symmetric around 0, multiplying by σ_i doesn't change its distribution. Then without dependence of σ_i 's distribution:

$$\begin{aligned}
\mathbb{E}[G_n] &\leq \mathbb{E} \left[\sup_{h \in H} \frac{1}{n} \sum_{i=1}^n \sigma_i [\ell(Z'_i, h) - \ell(Z_i, h)] \right] \\
&\leq \mathbb{E} \left[\sup_{h \in H} \frac{1}{n} \sum_{i=1}^n \sigma_i \ell(Z'_i, h) + \sup_{h \in H} \frac{1}{n} \sum_{i=1}^n (-\sigma_i) \ell(Z_i, h) \right] \\
&= 2 \mathbb{E} \left[\sup_{h \in H} \frac{1}{n} \sum_{i=1}^n \sigma_i \ell(Z_i, h) \right] \quad \text{here we consider } \sigma_i \stackrel{d}{=} -\sigma_i
\end{aligned}$$

RHS is defined as Rademacher complexity.

[Definition 9] Rademacher complexity.

Let F be a class of real-valued functions $f: \mathcal{Z} \rightarrow \mathbb{R}$.

① Define the Rademacher complexity of F to be

$$R_n(F) := \mathbb{E} \left[\sup_{f \in F} \frac{1}{n} \sum_{i=1}^n \sigma_i f(Z_i) \right]$$

where (i) $Z_1, \dots, Z_n \sim p^*$, i.i.d.

(ii) $\sigma_1, \dots, \sigma_n \sim \mathcal{U}\{-1, +1\}$ **uniform**.

② Define the empirical Rademacher complexity of F to be

$$\hat{R}_n(F) := \mathbb{E} \left[\sup_{f \in F} \frac{1}{n} \sum_{i=1}^n \sigma_i f(Z_i) | Z_{1:n} \right]$$

[Theorem 9] generalization bounds based on Rademacher complexity.

Define $A := \{Z \mapsto \ell(Z, h) : h \in H\}$ to be the loss class. With

$$\text{prob.} \geq 1 - \delta, \quad L(\hat{h}) - L(h^*) \leq 4R_n(A) + \sqrt{\frac{2 \log(2/\delta)}{n}}$$

Pf: Note that $E[G_n] \leq 2R_n(A)$, and $R_n(A) = R_n(-A)$

$$E[G'_n] \leq 2R_n(-A)$$

$$\begin{aligned} P\{G_n \geq \frac{\varepsilon}{2}\} &\leq \exp(-2n(\frac{\varepsilon}{2} - E[G_n])^2) \\ &\leq \exp(-2n(\frac{\varepsilon}{2} - 2R_n(A))^2) \quad \text{for } \varepsilon \geq 4R_n(A) \\ &:= \frac{\delta}{2} \end{aligned}$$

Similarly, $P\{G'_n \geq \frac{\varepsilon}{2}\} \leq \frac{\delta}{2}$. And we have

$$\varepsilon = 2\sqrt{\log(2/\delta)/2n} + 4R_n(A) \quad (\varepsilon \geq 4R_n(A)). \quad \square$$

Then we just need to study $R_n(A)$ over different A . now we first discuss some basic properties:

[Basic Properties of Rademacher complexity]

- ✓ (i) Boundedness: $R_n(F) \leq \max_{f \in F} \max_z f(z)$
- ✓ (ii) Singleton: $R_n(\{f\}) = 0$
- ✓ (iii) Monotonicity: $R_n(F_1) \leq R_n(F_2)$ if $F_1 \subseteq F_2$.
- ✓ (iv) Linear combination: $R_n(F_1 + F_2) = R_n(F_1) + R_n(F_2)$
- ✓ (v) Scaling: $R_n(cF) = |c| R_n(F)$
- * (vi) Lipschitz composition: $R_n(\phi \circ F) \leq C_\phi R_n(F)$, C_ϕ is the L-constant.
- * (vii) Convex hull: $R_n(\text{convex-hull}(F)) = R_n(F)$ for finite F .