

3.14 Algorithm stability

Motivation:

Using uniform convergence: [for $\forall h \in H$]

$$P[L(\hat{h}) - \hat{L}(\hat{h}) \geq \epsilon] \leq P\left[\sup_{h \in H} L(h) - \hat{L}(h) \geq \epsilon\right]$$

Actually, LHS holds for any estimator.

Question: what if an algorithm A do not use all of H ?

Let's define the regularized empirical risk minimizer as:

$$\hat{h} = \operatorname{argmin}_{h \in H} \hat{L}(h) + \frac{\lambda}{2} \|h\|_H^2$$

where $\|h\|_H^2$ is $\|w\|_2^2$ of the associated weight vector.

and constraint is $\|h\|_H \leq B$. [B depends on data].

Define the training set S and the perturbed version S^i as

(i) $S = (z_1, \dots, z_n)$: drawn i.i.d. from p^*

(ii) $S^i = (z_1, \dots, z_i', \dots, z_n)$: i.i.d. copy of the i -th example.

(iii) z_0 is a new test sample.

[Definition 15] Uniform stability

(i) an algorithm $A: Z^n \rightarrow H$ has uniform stability β w.r.t. a loss function l if for all $S \in Z^n$, $S^i \in Z^n$ and $z_0 \in Z$,

$$|l(z_0, A(S)) - l(z_0, A(S^i))| \leq \beta$$

Note, this is a strong condition and not reliant on distrib.

[Example 11] stability of mean estimation

Assume $\forall z \in Z \subseteq \mathbb{R}^d$, $\|z\|_2^2 \leq B$.

Define loss: $l(z, h) = \frac{1}{2} \|z - h\|_2^2$

Define algorithm: $A(S) := \operatorname{argmin}_{h \in \mathbb{R}^d} \hat{L}(h) + \frac{\lambda}{2} \|h\|_2^2$

$$\Rightarrow A(S) = \frac{1}{(1+\lambda)n} \sum_{i=1}^n z_i$$

$$\text{Th. ... } R = \frac{6B^2}{\lambda}$$

then $P = \frac{1}{(1+\lambda)n}$.

To derive this, define $V_1 = A(S) - z_0$, $V_2 = \frac{1}{(1+\lambda)n} [z'_i - z_i]$

$$\begin{aligned} |\ell(z_0, A(S)) - \ell(z_0, A(S^i))| &= \left| \frac{1}{2} \|V_1\|_2^2 - \frac{1}{2} \|V_1 + V_2\|_2^2 \right| \\ &\leq \|V_2\|_2 (\|V_1\|_2 + \frac{1}{2} \|V_2\|_2) \quad \text{Minkowski ineq.} \\ &\leq \frac{2B}{(1+\lambda)n} (2B + B) \quad \square \end{aligned}$$

[Theorem 16] generalization under uniform stability.

Let A be an algorithm with uniform stability β .

Assume the loss bounded: $\sup_{z, h} |\ell(z, h)| \leq M$

Then with prob $\geq 1 - \delta$, we have

$$L(A(S)) \leq \hat{L}(A(S)) + \beta + (\beta n + M) \sqrt{\frac{2 \log(1/\delta)}{n}}$$

Note: we must have that $\beta \sim o(\frac{1}{\sqrt{n}})$.

Pf of theorem 16:

Define $D(S) := L(A(S)) - \hat{L}(A(S))$

① Bound $E[D(S)]$

$$\begin{aligned} E[D(S)] &= E\left[\frac{1}{n} \sum_{i=1}^n [\ell(z_0, A(S)) - \ell(z_i, A(S))]\right] \\ &= E\left[\frac{1}{n} \sum_{i=1}^n [\ell(z'_i, A(S)) - \ell(z_i, A(S))]\right] \quad \text{rename } z_0 \\ &\leq \beta \end{aligned}$$

② show $D(S)$ satisfies the bounded differences property.

$$\begin{aligned} |D(S) - D(S^i)| &= |L(A(S)) - \hat{L}(A(S)) - L(A(S^i)) + \hat{L}^i(A(S^i))| \\ &\leq |L(A(S)) - L(A(S^i))| + |\hat{L}(A(S)) - \hat{L}(A(S^i))| \\ &\quad + |\hat{L}(A(S)) - \hat{L}^i(A(S^i))| \\ &\leq \beta + \beta + \frac{2M}{n} \end{aligned}$$

where $\hat{L}^i(A(S^i)) = \frac{1}{m} \sum_{k=1}^m L(A(S_k^i))$ $S_k^i = \{(z_1, \dots, z_i^{(k)}, \dots, z_n) : k\}$

③ McDiarmid's inequality:

$$-2\epsilon^2$$

$$\begin{aligned}
 P\{D(S) - E[D(S)] \geq \varepsilon\} &\leq \exp\left\{-\frac{n(2\beta + \frac{2M}{n})\varepsilon^2}{2}\right\} \\
 &= \exp\left\{-\frac{n\varepsilon^2}{2(\beta n + M)^2}\right\} =: \delta
 \end{aligned}
 \quad \square$$