

Lecture 18: November 2

Lecturer: Csaba Szepesvári

Scribes: Vedant Vyas

Note: *LaTeX template courtesy of UC Berkeley EECS dept. ([link](#) to directory)*

Disclaimer: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

18.1 Leave One Out (L.O.O) Stability

Definition 18.1. $k : \overset{\text{target}}{Z^n} \cup \overset{\text{auxiliary}}{Z^{n+1}} \rightarrow \mathcal{M}_1(\mathbb{R}^{\mathbb{Z}})$ is ε -LOO (Leave One Out) stable with $\varepsilon : \mathbb{Z}^{n+1} \cup \mathbb{Z} \rightarrow [0, \infty)$
if $\forall z_{1:n+1} \in Z^{n+1} \mu_k(z_{1:n}, z_{n+1}) \leq \mu_k(z_{1:n+1}, z_{n+1}) + \varepsilon(z_{1:n+1}, z_{n+1})$

Proposition 18.2. Assume $\mu_k(z_{1:n+1}^{i \leftrightarrow n+1}, z_{n+1}) = \mu_k(z_{1:n+1}, z_{n+1})$ and $\varepsilon(z_{1:n+1}^{i \leftrightarrow n+1}, z_{n+1}) = \varepsilon(z_{1:n+1}, z_{n+1})$
& K is ε -LOO stable then

$$\mathbb{E}[\mu_k(z_{1:n}, z_{n+1})] \leq \mathbb{E}P_{n+1}\mu_k(z_{1:n+1}) + \mathbb{E}P_{n+1}\varepsilon(z_{1:n+1})$$

Note: From our assumptions on μ_k we can get complete symmetry as to swap $i \leftrightarrow j$, we can do $i \leftrightarrow n+1$, and then $n+1 \leftrightarrow j$

Proof.

$$\begin{aligned} z_{1:n+1}^{i \leftrightarrow n+1} &= (z_1, \dots, z_{i-1}, z_{n+1}, \dots, z_i) \\ (z_{1:n+1}^{i \leftrightarrow n+1}, z_i) &\stackrel{P}{=} (z_{1:n+1}, z_{n+1}) \quad (\text{Same joint distribution}) \\ (\mu_k + \varepsilon)(z_{1:n+1}^{i \leftrightarrow n+1}, z_i) &\sim (\mu_k + \varepsilon)(z_{1:n+1}, z_{n+1}) \\ \text{So, } \mu_k(z_{1:n}, z_{n+1}) &\leq (\mu_k + \varepsilon)(z_{1:n+1}, z_{n+1}) \\ &\stackrel{P}{=} (\mu_k + \varepsilon)(z_{1:n+1}^{i \leftrightarrow n+1}, z_i) \end{aligned}$$

Taking expectation on both sides, and averaging over i would lead to the desired result □

18.2 First Order Optimality

Lemma 18.3. $f : C \rightarrow \mathbb{R}, C \subseteq \mathbb{R}^d$ closed, convex, $C \neq \emptyset$

$$x^* \in \operatorname{argmin}_{x \in C} f(x).$$

- (1.) $\exists \theta \in \partial f(x^*)$ s.t. $\theta^T(x - x^*) \geq 0$.
- (2.) Assume f is λ -SOC, $\tilde{\tau} \in \partial f(x)$, $\tilde{\theta}^\top(x^* - x) \geq -g\|x - x^*\|$ for some $g > 0 \Rightarrow \|x - x^*\| \leq \frac{g}{\lambda}$.

Proof.

$$f'(x^*; x - x^*) = \theta^\top(x - x^*) \text{ for some } \theta \in \partial f(x)$$

$$\begin{aligned} f(x^*) &\geq f(x) + \tilde{\theta}^\top(x^* - x) + \frac{\lambda}{2}\|x^* - x\|^2 \\ &\geq f(x) - g\|x - x^*\| + \frac{\lambda}{2}\|x^* - x\|^2 \end{aligned}$$

$$\begin{aligned} f(x) &\geq f(x^*) + \theta^\top(x - x^*) + \frac{\lambda}{2}\|x - x^*\|^2 \\ &\geq f(x^*) + \frac{\lambda}{2}\|x - x^*\|^2 \quad (\text{Using 1}) \end{aligned}$$

$$g\|x - x_*\| \geq \lambda\|x - x^*\|^2$$

Q.E.D.

□

Theorem 18.4. $G = \{g_w : w \in C\}$, $C \subseteq \mathbb{R}^d$ closed and convex. $w \mapsto g_w(z)$ $(W_1, \|\cdot\|_2) \rightarrow (\mathbb{R}, |\cdot|)$ $G(z)$ -Lipschitz
 $h : \mathbb{R}^d \rightarrow \mathbb{R}$; $\bar{g}_w(z) = g_w(z) + h(w)$.
Assume $\omega \mapsto P_{z_{1:n}}\bar{g}_w$ is λ -SOC (λ Strongly Convex); $\forall z_{1:n} \in Z^n$

$$\begin{aligned} \mathcal{A}(z_{1:n}) &:= \operatorname{argmin}_{w \in C} P_{z_{1:n}}\bar{g}_w = \operatorname{argmin}_{w \in C} P_{z_{1:n}}g_w + h(w) \\ \mathcal{A}(z_{1:n+1}) &= \operatorname{argmin}_{w \in C} P_{z_{1:n+1}}g_w + \frac{n}{n+1}h(w) \end{aligned}$$

Then: 1.) \mathcal{A} is $\varepsilon(z_{1:n+1}, z'_{n+1}) = \frac{G(z'_n+1)^2}{\lambda(n+1)}$ - LOO Stable

2.) \mathcal{A} is $\frac{2\|G\|_\infty^2}{\lambda n}$ uniformly stable

3.) $z_{1:n} \sim p^{\otimes n}$. Then, for $w_n = \mathcal{A}(z_{1:n})$,

$$\mathbb{E}P_{z_{1:n}}g_{w_n} \leq \inf_{\omega \in C} (P_{z_{1:n}}g_{w_n} + h(\omega)) + \frac{\mathbb{E}G^2(z_1)}{\lambda(n+1)}$$

Proof.

Part 1: L.O.O Stability

$$\begin{aligned} \text{Let } Z_{1:n+1} &\sim P^{\otimes n+1} \\ L_n(w) &= P_{z_{1:n}}\bar{g}_w \\ w_n &= \mathcal{A}(z_{1:n}) = \operatorname{argmin}_w L_n(w) \\ L_{n+1}(w) &= P_{z_{1:n+1}}g_w + \frac{n}{n+1}h(w) \\ w_{n+1} &= \mathcal{A}(z_{1:n+1}) = \operatorname{argmin}_w L_{n+1}(w) \end{aligned}$$

Assume L_n, L_{n+1} are differentiable. By F.O.O. Lemma (18.4): $\theta_n := \nabla L_n(w_n)$ s.t.

$$\theta_n^\top (w_{n+1} - w_n) \geq 0. \quad (*)$$

$$\begin{aligned}
\text{Now, } L_{n+1}(w) &= \frac{n}{n+1} L_n(w) + \frac{1}{n+1} g_w(z_{n+1}), \\
\nabla L_{n+1}(w_n) &= \frac{n}{n+1} \nabla L_n(w_n) + \frac{1}{n+1} \nabla g_{w_n}(z_{n+1}), \\
\Rightarrow \nabla L_{n+1}(w_n)^T (w_{n+1} - w_n) &\stackrel{(*)}{\geq} \frac{1}{n+1} \nabla g_{w_n}(z_{n+1})^\top (w_{n+1} - w_n), \\
\geq -\frac{G(z_{n+1})}{n+1} \|w_{n+1} - w_n\|, \\
\downarrow \\
\|\nabla g_w\| &\leq G, \quad (\text{Cauchy Schwartz}).
\end{aligned}$$

Now by using F.O.O Lemma (Part 2), we get

$$\begin{aligned}
\Rightarrow \|W_{n+1} - W_n\| &\leq \frac{G(z_{n+1})}{\lambda(n+1)} \\
\Rightarrow g_{w_n}(z_{n+1}) - g_{w_{n+1}}(z_{n+1}) &\leq G(z_{n+1}) \|w_n - w_{n+1}\| \\
\leq \frac{G^2(z_{n+1})}{\lambda(n+1)} \\
\Rightarrow \mathcal{A} \text{ is } \frac{Q^2}{\lambda(t+1)} - L \cdot O.O \text{ stable}
\end{aligned}$$

Part 2: Uniform Stability Proof for Uniform stability follows similarly but with the use of $\|\cdot\|_\infty$

Part 3: Symmetry

$$\begin{aligned}
\mathcal{A}(z_{1:n+1}^{i \leftrightarrow n+1}) &= \mathcal{A}(z_{1:n+1}), \\
\varepsilon(z_{1:n+1}^{i \oplus n+1}) &= \varepsilon(z_1 : n-1).
\end{aligned}$$

LOO Theorem: $Z_{1:n} \sim P^{\otimes n}$; $W_n : A(Z_{1:n})$, $W_{n+1} = A(Z_{1:n+1})$.

$$\begin{aligned}
\mathbb{E} P g_{w_n} &\leq \mathbb{E} P_{n+1} g_{w_{n+1}} + \mathbb{E} P_{n+1} \varepsilon(z_{n:n+1}) \\
&\leq \mathbb{E} L_{n+1}(w_{n+1}) + \frac{\mathbb{E} G^2(z_1)}{\lambda(n+1)} \\
&\leq \mathbb{E} L_{n+1}(\omega) + \frac{1}{n+1} h(\omega) + \frac{\mathbb{E} G^2(z_1)}{\lambda(n+1)} \\
&= P g_w + h(\omega) + \frac{\mathbb{E} G^2(z_1)}{\lambda(n+1)}
\end{aligned}$$

□

Example)

$$\begin{aligned}
 l(f_1(x, y)) &= \max(1 - f(\tilde{x})y, 0), y \in \pm 1 \\
 f_w(\alpha) &= w^\top \psi(x), w \in \mathbb{R}^d \\
 g(w) &= \frac{\lambda}{2} \|w\|^2; \quad h = 0. \\
 \text{Then } \mathbb{E}[P\ell \circ f_{w_n}] &\leq \inf_w P\ell \circ f_w + \frac{\lambda}{2} \|w\|^2 + \frac{\mathbb{E}(\|\psi(X)\| + \sqrt{2\lambda})^2}{\lambda(n+1)}
 \end{aligned}$$

Proof. Consider the function defined by

$$g_\omega(x, y) = \ell(f_{w_1}(x, y)) + \frac{\lambda}{2} \|\omega\|^2. \quad (18.1)$$

The mapping $\omega \mapsto g_w(z)$ induces a gradient

$$\nabla g_{w(1)} = \underbrace{l'(f_w, z)}_{\in \{0, 1\}} \psi + \underbrace{\lambda \omega}_{\text{unbounded}}. \quad (18.2)$$

Given $z_{1:n} \in Z^n$ and $w_n \in \operatorname{argmin} L_n(w)$, we have

$$\frac{\lambda}{2} \|w_n\|^2 \leq L_n(w_n) \leq L_n(0) \leq 1 \Rightarrow \|w_n\| \leq \sqrt{\frac{2}{\lambda}}, \quad (18.3)$$

where

$$\mathcal{A}(z_{1:n}) = \operatorname{argmin}_{w \in C} L_n(w) \quad \text{and} \quad C = \left\{ w : \|w\| \leq \sqrt{\frac{2}{\lambda}} \right\}. \quad (18.4)$$

□

Theorem 18.5. Assume $\forall z \in Z$, the mapping $w \mapsto g_w(z)$ is λ -strongly convex and L -smooth. Then the following properties hold:

$$\mathcal{A}(z_{1:n}) = \operatorname{argmin}_w P_{z_{1:n}} g_w, \quad (18.5)$$

$$\mathcal{A}(z_{1:n+1}) = \operatorname{argmin}_w P_{z_{1:n+1}} g_w, \quad (18.6)$$

$$\varepsilon(z_{1:n+1}, z'_{n+1}) = \left(1 + \frac{L}{2\lambda n}\right) \frac{\|\nabla g_{\mathcal{A}(z_{1:n})}(z'_{n+1})\|_2^2}{\lambda n}. \quad (18.7)$$

1. \mathcal{A} is ε -L.O.O stable.

2. If $L \leq 0.2\lambda n$, then

$$\mathbb{E} P g_{\omega_n} \leq \inf_w \left(P g_\omega + \frac{2.2}{\lambda n} P \|\nabla g_\omega\|_2^2 \right). \quad (18.8)$$

Example: let the loss function l and regularization term g be defined as:

$$l(f_1(x, y)) = (f(x) - y)^2,$$

$$g(w) = \frac{\lambda}{2} \|w\|_2^2.$$

Then, the model f parameterized by weights w , and the composite objective function g_ω can be expressed as:

$$\begin{aligned} f &= f_w = w^\top \psi, \\ g_\omega(z) &= l(f_w, z) + g(\omega), \text{ implying strong convexity } (\lambda - \text{SOC}). \\ \nabla g_w(z) &= 2(f_w(x) - y)\psi(x) + \lambda\omega, \\ \nabla^2 g_w(z) &= 2\psi(x)^\top \psi(x) + \lambda I, \\ \lambda_{\max}(\nabla^2 g_\omega(z)) &\leq 2\|\psi(x)\|^2 + \lambda. \end{aligned}$$

We define the Lipschitz constant L as:

$$L := \sup_x (2\|\psi(x)\|^2 + \lambda).$$

If $L \leq 0.2\lambda n$, it follows that:

$$\mathbb{E}Pg_{w_n} \leq \inf_\omega \left(Pg_\omega + \frac{8.8}{\lambda n} \times \mathbb{E}\|\psi(X)\|^2 \|(f_w(x) - y)^2\| \right).$$

Assuming there exists a w_* such that:

$$\mathbb{E}[(f_{w_*}(X) - Y)^2 | X] \leq 0^2 \text{ almost surely (a.s.)},$$

We can deduce that:

$$\mathbb{E}Pg_{w_n} \leq \sigma^2 + \frac{\lambda}{2}\|w_*\|_2^2 + \frac{8.8\sigma^2}{\lambda n} \mathbb{E}\|\psi(x)\|^2.$$

H.W: Compare to result that does not use smoothness!

18.3 Bibliography