

Stability and Generalization (cont.)

Recap:

data z_1, z_2, \dots $\stackrel{iid}{\sim} P$ $P \in \mathcal{P}$, on \mathcal{Z}

$A: \mathcal{Z}^n \mapsto A(z^n) \in \mathcal{F}$

so:

$A: \mathcal{Z}^* \rightarrow \mathcal{F}$ $\mathcal{Z}^* := \bigcup_{n \geq 1} \mathcal{Z}^n$

loss: $l: \mathcal{F} \times \mathcal{Z} \rightarrow \mathbb{R}$

\mathcal{F} : closed, convex subset of H.S. \mathcal{H}



$$L(A(z^n)) = \int_{\mathcal{Z}} l(A(z^n), z) P(dz)$$

↙
r.v.!

$$= \mathbb{E}[l(A(z^n), z_{n+1}) | \mathcal{Z}^n]$$

$$\mathbb{E}[L(A(z^n))] = \mathbb{E}[l(A(z^n), z_{n+1})]$$

$$L_n(A(z^n)) = \frac{1}{n} \sum_{i=1}^n l(A(z^n), z_i) \quad (\text{empirical risk})$$

Note: $\mathbb{E}[L_n(A(z^n))] \neq \mathbb{E}[L(A(z^n))]$

Generalization gap of A :

$$g_n(A) := \sup_{P \in \mathcal{P}} \mathbb{E}_P [L_n(A(z^n)) - L(A(z^n))]$$

A generalizes if $\lim_{n \rightarrow \infty} g_n(A) = 0$

(— as sample gets large enough, $L_n(A(z^n))$ is a good proxy for $L(A(z^n))$).

$$z^n \rightarrow \boxed{A} \rightarrow A(z^n)$$

$$(z_{n+1}, \dots, z_{2n}) : E\left[\frac{1}{n} \sum_{i=1}^n l(A(z^n), z_{n+i}) \mid z^n\right] = L(A(z^n))$$

Consistency!

$$\begin{aligned} L^* &:= \inf_{f \in \mathcal{F}} L(f) \\ &= \inf_{f \in \mathcal{F}} E_p[l(f, z)] = L(f^*) \end{aligned}$$

A is consistent if

$$c_n(A) := \sup_{P \in \mathcal{P}} E[\underbrace{L(A(z^n)) - L^*}_{\geq 0}] \rightarrow 0$$

as $n \rightarrow \infty$

Asymptotically ERM algorithms!

$$L_n^* := \inf_{f \in \mathcal{F}} L_n(f) = \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n l(f, z_i)$$

$$\text{AERM: } c_n(A) := \sup_{P \in \mathcal{P}} E[L_n(A(z^n)) - L_n^*] \rightarrow 0 \text{ as } n \rightarrow \infty.$$

ERM \Rightarrow AERM (trivially)

Generalization on average:

$$\bar{g}_n(A) := \sup_{P \in \mathcal{P}} |\mathbb{E}[L_n(A(z^n)) - L(A(z^n))]|$$

A generalizes on avg if $\bar{g}_n(A) \rightarrow 0$ as $n \rightarrow \infty$

$$\bar{g}_n(A) \leq g_n(A) \quad \forall n$$

Thm For any A,

$$e_n(A) \leq \underbrace{e_n(A)}_{\substack{\rightarrow 0 \\ \text{(AERM)}}} + \underbrace{\bar{g}_n(A)}_{\substack{\rightarrow 0 \\ \text{(gen on avg)}}$$

\Rightarrow AERM + (gen on avg) \Rightarrow consistent.

Proof

$$\mathbb{E} L(A(z^n))$$

$$= \mathbb{E}[L(A(z^n)) - L_n(A(z^n))] + \mathbb{E}[L_n(A(z^n))]$$

$$\leq \bar{g}_n(A) + \mathbb{E}[L_n(A(z^n))]$$

$$= \bar{g}_n(A) + \mathbb{E}[L_n(A(z^n)) - L_n^*] + \mathbb{E}[L_n^*]$$

$$\leq \bar{g}_n(A) + e_n(A) + \mathbb{E}[L_n^*]$$

$$\Rightarrow \mathbb{E}[L(A(z^n))] \leq \mathbb{E}[L_n^*] + \bar{g}_n(A) + e_n(A)$$

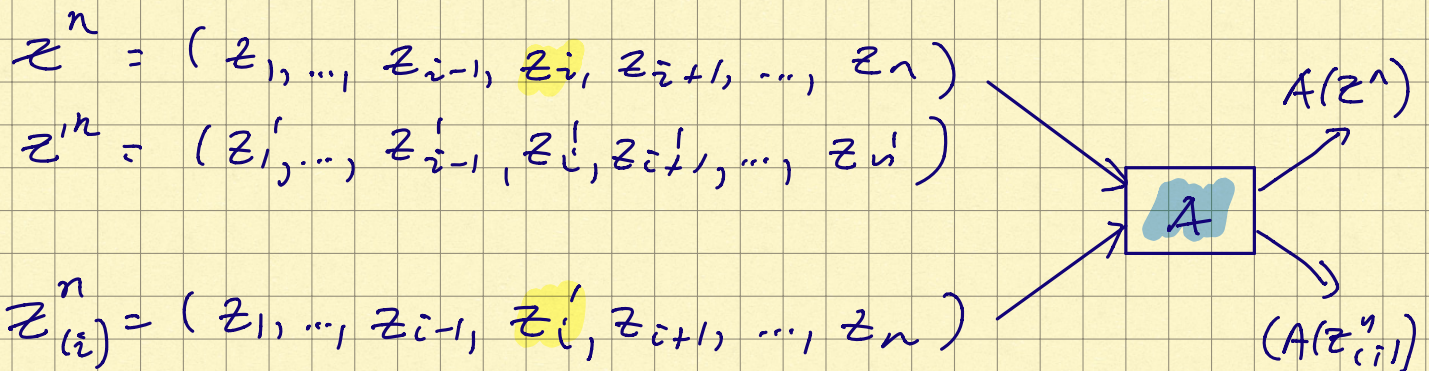
where $\mathbb{E}[L_n^*] = \mathbb{E}[L_n^* - L_n(f^*)] + \mathbb{E}[L_n(f^*)]$
 where f^* achieves L^* $= L(f^*)$

$$\therefore \mathbb{E}[L_n(A(z^n)) - L^*] \leq e_n(A) + \bar{g}_n(A)$$

AERM + (gen on avg) \Rightarrow consistency

gen on avg.: $\bar{g}_n(A) = \sup_{P \in \mathcal{P}} |\mathbb{E}[L_n(A(z^n)) - L(A(z^n))]|$

Stability



Stability: $A(z^n) \approx A(z_{(i)}^n)$ for each i

$$\mathbb{E}[L_n(A(z^n)) - L(A(z^n))]$$

$$= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \ell(A(z^n), z_i) - \frac{1}{n} \sum_{i=1}^n \ell(A(z^n), z'_i)\right]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\ell(A(z^n), z_i) - \ell(A(z^n), z_i')]]$$

• for each $i \in [n]$, $(A(z^n), z_i) \stackrel{d}{=} (A(z_{(i)}^n), z_i')$

$$\Rightarrow \mathbb{E} [\ell(A(z^n), z_i)] = \mathbb{E} [\ell(A(z_{(i)}^n), z_i')] \quad !$$

so

$$\mathbb{E} [L_n(A(z^n)) - L(A(z^n))]]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\ell(A(z_{(i)}^n), z_i') - \ell(A(z^n), z_i')]]$$

• if $A(z_{(i)}^n) \approx A(z^n) \quad \forall i$, then

$\mathbb{E} [\ell(A(z_{(i)}^n), z_i') - \ell(A(z^n), z_i')]]$ is "small".

e.g. $O\left(\frac{1}{\sqrt{n}}\right)$

$$\Rightarrow \mathbb{E} [L_n(A(z^n)) - L(A(z^n))] = O\left(\frac{1}{\sqrt{n}}\right)$$

• Stability \Leftrightarrow also not too "sensitive" to modifications of individual samples

$$I \sim \text{Unif}([n]), \quad I \perp (z^n, z'^n)$$

$$\frac{1}{n} \sum_{i=1}^n [\ell(A(z_{(i)}^n), z_i') - \ell(A(z^n), z_i')]]$$

$$= \mathbb{E} [\ell(A(z_{(I)}^n), z_I') - \ell(A(z^n), z_I') | z^n, z'^n]$$

Overall,

$$\mathbb{E}[L_n(A(z^n)) - L(A(z^n))]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\ell(A(z_{i-1}^n), z_{i-1}') - \ell(A(z^n), z_{i-1}')]]$$

Stability on Avg.:

A is stable on avg if

$$S_n(A)$$

$$:= \sup_{\mathcal{Z}^n} \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\ell(A(z_{i-1}^n), z_{i-1}') - \ell(A(z^n), z_{i-1}')]]$$

$$\rightarrow 0 \quad \text{as } n \rightarrow \infty$$

Key points:

1) $\overline{g}_n(A) = S_n(A)$ (gen on avg \Leftrightarrow stab. on avg)

2) in many cases, $S_n(A)$ can be easily estimated

- sufficient condition: uniform stability

$$\sup_{z \in \mathcal{Z}} |\ell(A(z^n), z) - \ell(A(z_{i-1}^n), z)| = o(n)$$

3) For any A , $c_n(A) \leq \underbrace{e_n(A)}_{\rightarrow 0 \text{ AERM}} + \underbrace{S_n(A)}_{\rightarrow 0 \text{ stable on avg}}$

Thm Suppose $0 \leq l(f, z) \leq 1$ for all $f \in \mathcal{F}, z \in \mathcal{Z}$.

Then, for any A ,

$$g_n(A) \leq \bar{g}_n(A) + 2(\epsilon_n(A) + \frac{1}{\sqrt{n}}).$$

Remark: $\bar{g}_n(A) = \sup_P |E L_n(A(z^n)) - E L(A(z^n))|$

$$g_n(A) = \sup_P E |L_n(A(z^n)) - L(A(z^n))|$$

$$\bar{g}_n(A) \leq g_n(A)$$

— AERM + (gen on avg) \Rightarrow gen.

Stochastic Gradient Descent (SGD)

$$L_n(f) := \frac{1}{n} \sum_{i=1}^n l(f, z_i)$$

• $f \in \mathcal{F}$, closed convex subset of \mathcal{H}
Hilbert space \mathcal{H}

• $f \mapsto l(f, z)$ is differentiable for each $z \in \mathcal{Z}$

$\nabla l(f, z)$: gradient (Fréchet derivative)

$$l(f', z) = l(f, z) + \underbrace{\nabla l(f, z)(z' - z)}_{\text{linear map } \mathcal{H} \rightarrow \mathbb{R}} + \underbrace{o(\|z' - z\|)}_{\text{error} \rightarrow 0 \text{ as } z' \rightarrow z}$$

$$\text{ERM: } f_n^* := \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \ell(f, z_i)$$

alternative (approx. ERM): \hookrightarrow

$$f_0 \in \mathcal{F} \quad (\text{init})$$

at $t = 1, 2, \dots$

$$f_t = \Pi \left(f_{t-1} - \alpha_t \nabla L_n(f_{t-1}) \right)$$

where Π is the projection onto \mathcal{F} :

$$\Pi(g) := \underset{f \in \mathcal{F}}{\operatorname{argmin}} \|g - f\|$$

$\alpha_t > 0$ - stepsize

$$t = 1, \dots, T : \quad \hat{f}_n = f_T$$

bound $L_n(\hat{f}_T) - L_n^*$ (AERM)

However: computation of $\nabla L_n(f)$ has $O(n)$ complexity

$$\nabla L_n(f) = \frac{1}{n} \sum_{i=1}^n \nabla \ell(f, z_i)$$

$$\text{SGD: } L_n(f) = \frac{1}{n} \sum_{i=1}^n \ell_i(f) \quad \ell_i(f) := \ell(f, z_i)$$

$I_1, I_2, I_3, \dots, I_t, \dots$: stochastic process in $[n]$
indep. of z^n

$$f_t = \Pi \left(f_{t-1} - \alpha_t \nabla \ell_{I_t}(f_{t-1}) \right)$$

each step has now $O(1)$ complexity

Examples of (I_t) :

- sampling w/ replacement: $I_t \stackrel{iid}{\sim} \text{unif}([n])$
- cyclic coordinate descent:
 $\pi(1), \dots, \pi(n)$: random perm of $[n]$

$(I_1, I_2, \dots, I_n, I_{n+1}, \dots)$

$= (\pi(1), \pi(2), \dots, \pi(n), \pi(1), \pi(2), \dots, \pi(n), \dots)$

- new perm. every n rounds

Next lecture: SGD is stable (under regularity assumptions)