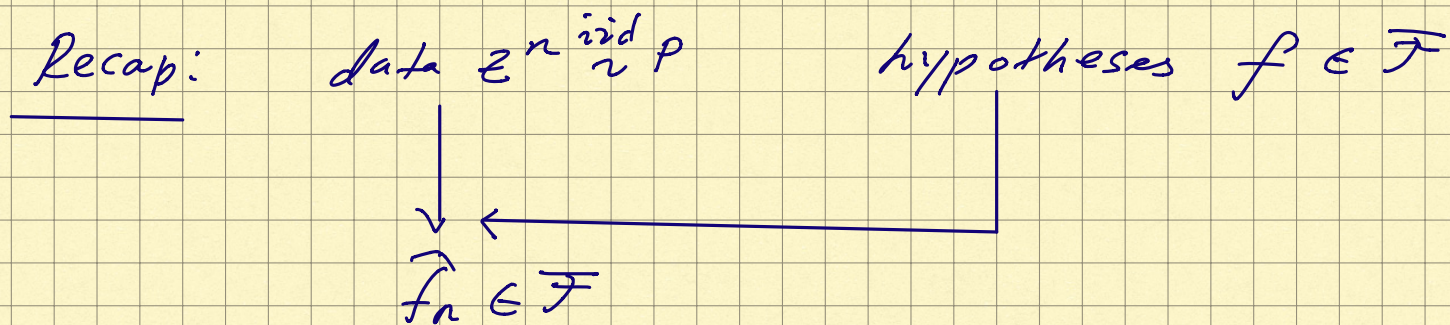


Learning Algorithms in Depth: Stability

Ch. 13 (read Ch. 3 on optimization)



e.g. ERM $\hat{f}_n = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n f(z_i)$

$$L(\hat{f}_n) = \int_{\mathcal{Z}} \hat{f}_n(z) P(dz)$$

(universal)
"consistency":

$$L^* = \min_{f \in \mathcal{F}} L(f)$$
$$\sup_P \{L(\hat{f}_n) - L^*\} \xrightarrow{\text{as } n \rightarrow \infty} 0$$

Sufficient condition for consistency: UCEM

$$P_n(f) := \frac{1}{n} \sum_{i=1}^n f(z_i) = \mathbb{E}_{P_n}(f) \quad (\text{empirical risk})$$

$$P(f) = \mathbb{E}_P(f)$$

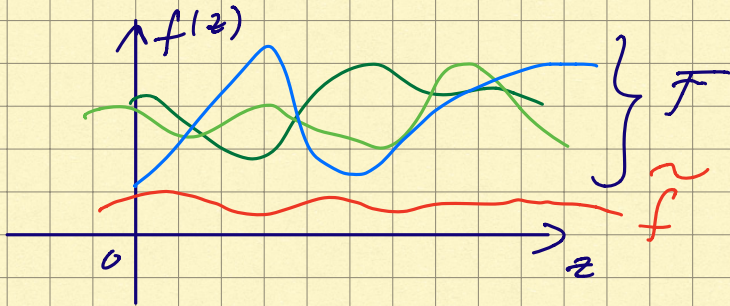
$$\sup_P \mathbb{E}_P \sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \xrightarrow{n \rightarrow \infty} 0$$

UCEM (property of \mathcal{F}) \Rightarrow ERM is consistent

But: learnability is possible w/o UCEM!

• \mathcal{F} given no assumptions

• assume $\exists \tilde{f}$ s.t. $\tilde{f}(z) < \inf_{f \in \mathcal{F}} f(z) \quad \forall z$



$$\tilde{\mathcal{F}} := \mathcal{F} \cup \{\tilde{f}\}$$

ERM on $\tilde{\mathcal{F}}$ will always return \tilde{f} !

closer look at learning algo

Framework (Vapnik, 1995)

- \mathcal{Z} : instance space
- \mathcal{P} : class of prob. dist. on \mathcal{Z}
- \mathcal{F} : closed, convex subset of a Hilbert space \mathcal{H}
- $l: \mathcal{F} \times \mathcal{Z} \rightarrow \mathbb{R}$: loss fcn

$$L_P(f) := \mathbb{E}_P[l(f, z)]$$

$$= \int_{\mathcal{Z}} l(f, z) P(dz)$$

Examples: 1) \mathcal{H} : H.S. of fcn $\mathcal{Z} \rightarrow \mathbb{R}$
 $l(f, z) := f(z)$

$$2) \mathcal{H} = \mathcal{H}_K \quad (\text{RKHS for some } K)$$

$$z = (x, y) \in \mathcal{X} \times \{\pm 1\}$$

$$l(f, z) = l(f, (x, y)) = \varphi(-yf(x))$$

↳ pen. fun.

$$3) l(f, (x, y)) = (y - f(x))^2$$

...

Main idea: want to bring out the dependence of $l(\cdot, \cdot)$ on both f and z

• Learning algos:

$$A_n: \mathcal{Z}^n \rightarrow \mathcal{F} \quad (\text{one for each } n)$$

$$\mathcal{Z}^* := \bigcup_{n \geq 1} \mathcal{Z}^n \quad : \text{ all tuples over } \mathcal{Z}$$

$$A: \mathcal{Z}^* \rightarrow \mathcal{F}$$

$A(z^n)$: random element of \mathcal{F}

• Risks:

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

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

Motivating Example: ERM w/ Strongly Convex Losses (Ch. 3)

$$(f, z) \mapsto l(f, z)$$

$f \in \mathcal{F}$: closed, convex subset of $(\mathcal{F}, \|\cdot\|)$

Assume:

1) $l(f, z)$ is L -Lipschitz in $f \in \mathcal{F}$, unif. in z :

$$\sup_{z \in \mathcal{Z}} |l(f, z) - l(f', z)| \leq L \|f - f'\|$$

2) $l(f, z)$ is m -strongly convex in $f \in \mathcal{F}$, unif. in z ($m > 0$)

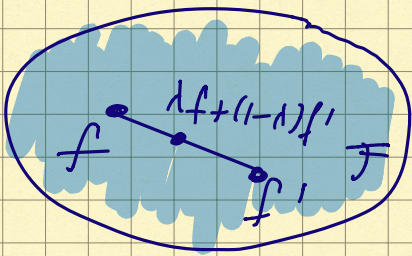
Review:

$$\varphi: \mathcal{F} \rightarrow \mathbb{R}$$

\mathcal{F} - convex set

$$f, f' \in \mathcal{F}$$

$$\Rightarrow \lambda f + (1-\lambda)f' \in \mathcal{F} \quad \forall \lambda \in [0, 1]$$



\mathcal{F} : closed

$$f_n \in \mathcal{F}$$

$$f_n \rightarrow f \Rightarrow f \in \mathcal{F}$$

$\varphi: \mathcal{F} \rightarrow \mathbb{R}$ is convex if

$$\varphi(\lambda f + (1-\lambda)f') \leq \lambda \varphi(f) + (1-\lambda) \varphi(f')$$

for all $f, f' \in \mathcal{F}$, $\lambda \in [0, 1]$

$\varphi: \mathcal{F} \rightarrow \mathbb{R}$ is m -strongly convex ($m \geq 0$)

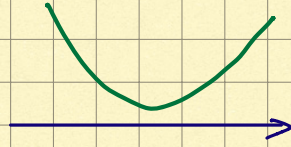
if $\tilde{\varphi}(f) := \varphi(f) - \frac{m}{2} \|f\|^2$ is convex

[note: not requiring φ to be differentiable]

Consequences of strong convexity:

$$\bullet \varphi(\lambda f + (1-\lambda)f') \leq \lambda \varphi(f) + (1-\lambda)\varphi(f') - \frac{\mu}{2} \lambda(1-\lambda) \|f - f'\|^2 \quad (\text{Conv.}) \quad (A)$$

• unique minimizers



• if $\varphi(f^*) = \min_{f \in \mathcal{F}} \varphi(f)$, then

$$\varphi(f) \geq \varphi(f^*) + \frac{\mu}{2} \|f - f^*\|^2$$

(can be derived from (A))

• let $B(f)$ be an L -Lip. fun, let

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \varphi(f)$$

$$\bar{f} = \operatorname{argmin}_{f \in \mathcal{F}} \left\{ \varphi(f) + \underbrace{B(f)}_{\text{Lip. perturbation}} \right\}$$

$$\text{Then } \|\bar{f} - f^*\| \leq \frac{L}{\mu}$$

(stability of minimizers under Lip. perturbations)

• if $\varphi(\cdot)$ is convex, then $\varphi(\cdot) + \frac{\mu}{2} \|\cdot\|^2$ is μ -strongly convex

Back to learning:

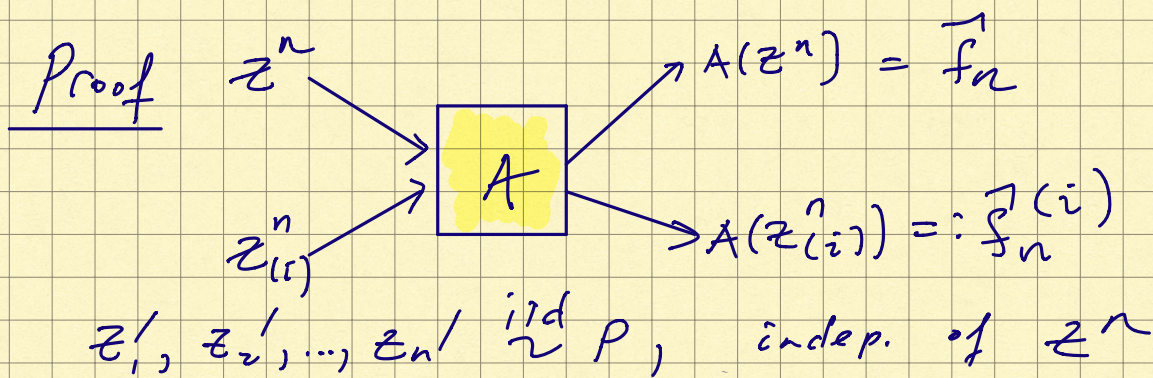
$\ell(f, z)$: L -Lip., μ -str. conv in $f \in \mathcal{F}$
(uniformly in z)

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

Thm With $pr. \geq 1 - \delta,$

$$L(\hat{f}_n) - \inf_{f \in \mathcal{F}} L(f) \leq \frac{2L^2}{5mn}$$

Note: $\frac{1}{\delta}$ dependence instead of $\log(\frac{1}{\delta})$



$z_1', z_2', \dots, z_n' \stackrel{iid}{\sim} P,$ indep. of z^n

$$z_{(i)}^n := (z_1, \dots, z_{i-1}, z_i', z_{i+1}, \dots, z_n)$$

$$z^n = (z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_n)$$

• Fix $f \in \mathcal{F}$

$$L_n(f) := \frac{1}{n} \sum_{i=1}^n \ell(f, z_i) \quad [\text{loss on } z^n]$$

$$L_n^{(i)}(f) := \frac{1}{n} \ell(f, z_i') + \frac{1}{n} \sum_{j: j \neq i} \ell(f, z_j) \quad [\text{loss on } z_{(i)}^n]$$

$$= \frac{1}{n} \ell(f, z_i') + L_n(f) - \frac{1}{n} \ell(f, z_i)$$

$$= L_n(f) + \frac{1}{n} [\ell(f, z_i') - \ell(f, z_i)]$$

\therefore for each $i \in [n]:$

$$L_n^{(i)}(f) - L_n(f) = \frac{1}{n} [\ell(f, z_i') - \ell(f, z_i)]$$

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} L_n(f)$$

$$\hat{f}_n^{(i)} = \arg \min_{f \in \mathcal{F}} L_n^{(i)}(f)$$

Claim: $\|\hat{f}_n - \hat{f}_n^{(i)}\| \leq \frac{2L}{mn}$

Proof (of claim)

1) $f \mapsto L_n(f)$ is m -str. conv. (b/c l is)

2) $f \mapsto \frac{1}{n} (l(f, z_i') - l(f, z_i))$ is $\frac{2L}{n}$ -Lip.

$$B(f) := \frac{1}{n} (l(f, z_i') - l(f, z_i))$$

$$\begin{aligned} |B(f) - B(f')| &\leq \frac{1}{n} |l(f, z_i') - l(f', z_i')| \\ &\quad + \frac{1}{n} |l(f, z_i) - l(f', z_i)| \end{aligned}$$

$$\leq \frac{2L}{n} \|f - f'\| \quad (\text{since } l \text{ is } L\text{-Lip. in } f)$$

$$\Rightarrow L_n^{(i)}(f) = \underbrace{L_n(f)}_{m\text{-s.c.}} + \underbrace{B(f)}_{\frac{2L}{n}\text{-Lip.}}$$

By stab. of minimizers, $\|\hat{f}_n - \hat{f}_n^{(i)}\| \leq \frac{2L}{nm}$. □

• Claim: $\mathbb{E} [L(\hat{f}_n) - L_n(\hat{f}_n)] \leq \frac{2L^2}{mn}$.

Proof (of claim)

$$\hat{f}_n = A(z^n)$$

$$\mathbb{E}[L(\hat{f}_n)] = \mathbb{E}[L(A(z^n))]$$

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

$\forall i$

$$z^n \perp z_{i'}$$

$$\Rightarrow \mathbb{E}[L(\hat{f}_n)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\ell(A(z^n), z_{i'})]$$

$$\mathbb{E}[L_n(\hat{f}_n)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\ell(A(z^n), z_i)]$$

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

$$\Rightarrow \mathbb{E}[L_n(\hat{f}_n)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\ell(A(z_{i'}^n), z_{i'})]$$

$$\therefore \mathbb{E}[L(\hat{f}_n) - L_n(\hat{f}_n)]$$

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

$$\leq \frac{1}{n} \sum_{i=1}^n L \cdot \mathbb{E} \|\hat{f}_n - \hat{f}_n^{(i)}\|$$

$$\leq \frac{2L^2}{nm}$$

□

$$\bullet \mathbb{E}[L(\hat{f}_n) - L(f^*)]$$

$$L(f^*) = \min_{f \in \mathcal{F}} L(f)$$

$$= \mathbb{E}[L(\hat{f}_n) - L_n(\hat{f}_n)]$$

$$+ \underbrace{\mathbb{E}[L_n(\hat{f}_n) - L_n(f^*)]}_{\leq 0} + \underbrace{\mathbb{E}[L_n(f^*) - L(f^*)]}_{= 0}$$

$$\leq \frac{2L^2}{nm}$$

• By Markov's inequality,

$$P\{L(\hat{f}_n) - L(f^*) \geq t\} \leq \frac{E[L(\hat{f}_n) - L(f^*)]}{t}$$
$$\leq \frac{2L^2}{nm t}$$

let $\frac{2L^2}{nm t} = \delta \quad \Rightarrow \quad t = \frac{2L^2}{\delta nm}$ □

Key points:

• $\|A(z^n) - A(z_{(i)}^n)\| \leq \frac{2L}{nm} \quad \forall i$
- hypothesis stability

• $E\left[\frac{1}{n} \sum_{i=1}^n (\ell(A(z^n), z_{i'}^n) - \ell(A(z_{(i)}^n), z_{i'}^n))\right]$
 $\leq \frac{2L^2}{nm}$
- replace-one stability