

Chapter 5

Uniform Laws of Large Numbers

- Rademacher complexity
- Glivenko-Cantelli Thm
- Upper bounds on Rademacher complexity
- VC dimension
- Learning linear classifiers
- Learning without uniform laws with convex losses.

Goal: Bound

$$\|P_n - P\|_F \triangleq \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}[f(x)] \right|$$

Defn:

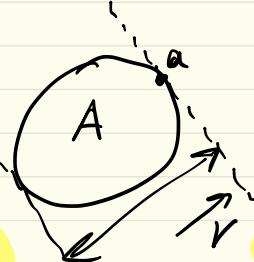
The Rademacher Complexity of a set $A \subseteq \mathbb{R}^n$, is the quantity

$$C(A) \triangleq \mathbb{E}_{\epsilon} \sup_{a \in A} \langle a, \epsilon \rangle$$

where $\epsilon = (\epsilon_1, \dots, \epsilon_n)$ are i.i.d. Rademacher random variables

Recall that support function of a set A is

$$\tilde{\sigma}_A(v) := \sup_{a \in A} |\langle v, a \rangle|$$



So

$$R(A) = \mathbb{E}_{\epsilon} \tilde{\sigma}_A(\epsilon)$$

Defn: Consider a sequence of random ^{independent} variables $(x_1, \dots, x_n) \subseteq \mathcal{X}^n$ and a class of functions \mathcal{F} on \mathcal{X} . The Rademacher complexity of \mathcal{F} is

$$\begin{aligned} C_n(\mathcal{F}) &\triangleq \mathbb{E}_{\substack{x, \epsilon \\ x \in \mathcal{X}, \epsilon \in \{-1, 1\}^n}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right| \\ &= \mathbb{E}_x C(\mathcal{F}^{\perp}(x), \cdot) \\ &\leq 2 \mathbb{E}_x C(\mathcal{F}(x)/n) \end{aligned}$$

where $\mathcal{F}(x) = \{(f(x_1), \dots, f(x_n)) : f \in \mathcal{F}\}$

Thm: $\mathbb{E} \|P_n - P\|_{\mathcal{F}} \leq 2 C_n(\mathcal{F})$

Pf: Let $y_1, \dots, y_n \stackrel{\text{iid}}{\sim} P$ and independent of x_1, \dots, x_n . Then

$$\begin{aligned} \mathbb{E} \|P_n - P\|_{\mathcal{F}} &= \mathbb{E}_x \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}_y \left[f(y_i) \right] \right| \right] \\ &= \mathbb{E}_x \left[\sup_{f \in \mathcal{F}} \left| \mathbb{E}_y \left[\frac{1}{n} \sum (f(x_i) - f(y_i)) \right] \right| \right] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_X \left[\sup_{f \in \mathcal{F}} \left| \mathbb{E}_y \left[\frac{1}{n} \sum (f(x_i) - f(y_i)) \right] \right| \right] \\
&\leq \mathbb{E}_X \mathbb{E}_y \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum (f(x_i) - f(y_i)) \right| \\
&= \mathbb{E}_{X, Y, \epsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(x_i) - f(y_i)) \right| \\
&\leq 2 \mathbb{E}_{X, \epsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right| = 2 C_n(\mathcal{F})
\end{aligned}$$

Thm: Suppose \mathcal{F} is b -uniformly bounded, meaning

$$\|f\|_\infty = \sup_{x \in X} |f(x)| \leq b \quad \forall f \in \mathcal{F}.$$

Then for any $n \in \mathbb{N}$, $t \geq 0$, it holds

$$\Pr \left(\|P_n - P\|_{\mathcal{F}} \leq 2 C_n(\mathcal{F}) + t \right) \geq 1 - e^{-\frac{n t^2}{2b^2}}$$

pf: All we have to do is show
 the bounded difference property
 for the function $\|P_n - P\|_F$
 with $L_i = \frac{2b}{n}$.

Define $\bar{f}(x) = f(x) - E f(x)$.

$$\text{Then } \|P_n - P\|_F = \sup_{f \in F} \left| \frac{1}{n} \sum_{i=1}^n \bar{f}(x_i) \right|$$

As before, let x' differ from x only in i 'th entry. Then

$$\left| \frac{1}{n} \sum_{i=1}^n \bar{f}(x_i) \right| - \sup_{g \in F} \left| \frac{1}{n} \sum_{i=1}^n \bar{g}(x'_i) \right| \\ \leq \left| \frac{1}{n} \sum_{i=1}^n \bar{f}(x_i) \right| - \left| \frac{1}{n} \sum_{i=1}^n \bar{f}(x'_i) \right| \leq \frac{2b}{n}$$

Take $\sup_{f \in F}$ and exchange x, x' 12

Rademacher complexity characterizes the asymptotics of $\|P_n - P\|_{\mathcal{F}}$.

$$\text{Prop } \frac{1}{2} C_n(\bar{\mathcal{F}}) \leq \mathbb{E}_X [\|P_n - P\|_{\mathcal{F}}] \leq 2C_n(\mathcal{F})$$

where $\bar{\mathcal{F}} = \{f - \mathbb{E}f : f \in \mathcal{F}\}$.

[See Prop 4.1.1 in Wainwright]

Exercise: Suppose \mathcal{F} is b -bounded
Then $\forall n \in \mathbb{N}, t > 0$, if holds

$$\begin{aligned} \text{pr} \left(\|P_n - P\|_{\mathcal{F}} \geq \frac{1}{2} C_n(\bar{\mathcal{F}}) - \frac{\sup_{f \in \mathcal{F}} |\mathbb{E}f(X)|}{2\sqrt{n}} - t \right) \\ \geq 1 - e^{-\frac{nt^2}{2b^2}} \end{aligned}$$

Goal: Bound the Rademacher Complexity of interesting sets.

Lemma (Basic Calculus)

For any $A, B \subseteq \mathbb{R}^n$, $c \in \mathbb{R}$, $a_0 \in \mathbb{R}^n$ it holds

- $C(cA) = |c| C(A)$
- $C(A+B) = C(A) + C(B)$
- $C(A+a_0) = C(A)$
- $C(A) = C(\text{conv}(A))$

Lemma: (Massart) Let $A = \{a_1, \dots, a_N\} \subseteq \mathbb{R}^n$.

Define $\bar{a} = \frac{1}{N} \sum_{i=1}^N a_i$. Then

$$C(A) \leq \max_{a \in A} \|a - \bar{a}\| \cdot \sqrt{2 \log(N)}$$

pf: WLOG assume $\bar{a} = 0$. Let $\lambda > 0$ and
 $A' = \{\lambda a_1, \dots, \lambda a_N\}$. Then

$$C(A') = \mathbb{E}_\varepsilon \max_{a \in A'} \langle \varepsilon, a \rangle$$

$$= \mathbb{E}_\varepsilon \log \left(\max_{a \in A'} e^{\langle \varepsilon, a \rangle} \right)$$

$$\leq \mathbb{E}_\varepsilon \log \left(\sum_{a \in A'} e^{\langle \varepsilon, a \rangle} \right)$$

Jensen

$$\leq \log \left(\sum_{a \in A'} \mathbb{E}_\varepsilon e^{\langle \varepsilon, a \rangle} \right)$$

Observe

$$\mathbb{E}_\varepsilon e^{\langle \varepsilon, a \rangle} = \prod_{i=1}^d \mathbb{E}_\varepsilon e^{a_i \varepsilon_i} \leq \prod_{i=1}^d e^{\frac{a_i^2}{2}} = e^{\|a\|_2^2}$$

$$\leq \log \left(\sum_{a \in A'} e^{\|a\|_2^2} \right)$$

$$\leq \log (|A'| \max_{a \in A'} (e^{\|a\|_2^2}))$$

$$= \log (|A'|) + \max_{a \in A'} \frac{\|a\|_2^2}$$

$$\text{Thus } C(A) = \frac{1}{\lambda} C(A') \leq \frac{\log(|A|) + \max_{a \in A} \frac{\lambda^2 \|a\|^2}{2}}{\lambda}$$

Optimize over $\lambda > 0$. □

Defn: F has polynomial discrimination of order $V \geq 1$, if $\forall n \in \mathbb{N}$ and

$x_1, \dots, x_n \in \mathcal{X}$, we have

$$\text{card}(F(x_1, \dots, x_n)) \leq (n+1)^V.$$

Cor: Suppose F has polynomial discrimination of order V . Then

$$C_n(F) \leq 4 \mathbb{E} \left[\sup_{f \in F} \sqrt{\frac{1}{n} \sum_{i=1}^n f^2(x_i)} \right] \cdot \sqrt{V \log(n+1)}$$

pf: Fix $x = (x_1, \dots, x_n)$. Then Massart implies

$$C(F(x)/n) \leq 4 \cdot \frac{\max_{f \in F} \sqrt{\sum_{i=1}^n f_i^2(x)}}{n} \cdot \sqrt{V \log(n+1)}$$

□

In particular, if F is b -bounded, then

$$C_n(F) \leq 4b \sqrt{\frac{v \log(n+1)}{n}}$$

Cor: (Glivenko-Cantelli)

Let $g(t) = P[X \leq t]$ be the CDF of X , and let $g_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}[X_i \leq t]$ where $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} P$. Then

$$P\left[\|g_n - g\|_\infty \geq 8 \sqrt{\frac{\log(n+1)}{n}} + \delta\right] \leq \exp\left(-\frac{n\delta^2}{2}\right)$$

Pf: Let $\mathcal{F} = \left\{ \mathbb{1}_{(-\infty, t]} : t \in \mathbb{R} \right\}$

$$\begin{aligned} \text{Then } \|g_n - g\|_\infty &= \sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}[X_i \leq t] - P[X \leq t] \right| \\ &= \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - E[f(X)] \right| \end{aligned}$$

Since \mathcal{F} is 1-bounded

$$P[\|g_n - g\|_\infty \leq 2C_n(\mathcal{F}) + \delta] \geq 1 - \exp\left(-\frac{n\delta^2}{2}\right)$$

Observe \mathcal{F} has poly discrimination with $V=1$ \square

Vapnik-Chervonenkis (VC) Theory

Method for bounding polynomial /
discrimination of $\{0,1\}$ -valued F .

Consider a class F of binary valued
functions on X .

Def: We say that $x = (x_1, \dots, x_n)$ is
shattered by F if $\text{card}(F(x)) = 2^n$.

The VC dimension is

$$\text{VC}(F) = \sup \{ n \in \mathbb{N} : \exists x \in X^n \text{ shattered by } F \}$$

Notation: If $F = \{ \bigcup_S : \text{some sets } S \}$

set $S(x) := F(x)$ and $\text{VC}(S) := \text{VC}(F)$

Ex: $S_{\text{left}} = \{ (-\infty, a] : a \in \mathbb{R} \} \Rightarrow \text{VC}(S_{\text{left}}) = 1$

$S_{\text{two}} = \{ (a, b] : a, b \in \mathbb{R}, a < b \} \Rightarrow \text{VC}(S_{\text{two}}) = 2$



Thm (Sauer and Shelah)

For any $x = (x_1, \dots, x_n)$ with $n \geq VC(S)$, we have

$$\text{card}(S(x)) \leq \sum_{i=0}^{VC(S)} \binom{n}{i} \leq (n+1)^{VC(S)}$$

Therefore F has polynomial discriminability of order $VC(S)$ and

$$C_n(F) \leq 2 \sqrt{\frac{VC(S) \cdot \log(n+1)}{n}} \leftarrow \begin{matrix} \log(n+1) \\ \text{can be removed.} \end{matrix}$$

[See Prop 4.18 in Wainwright for a proof.]

Examples: Let G be a class of functions. For any $g: X \rightarrow \mathbb{R}$ define

$$S_g = \{x \in X : g(x) \leq 0\}$$

Prop: Let G be a vector space of functions $g: \mathbb{R}^d \rightarrow \mathbb{R}$ with $\dim(G) < \infty$.

Then

$$VC(S(G)) \leq \dim(G)$$

Pf: Set $n = \dim(G) + 1$ and fix $x = \{x_1, \dots, x_n\}$ with $x_i \in \mathcal{X}$. Define $L: G \rightarrow \mathbb{R}^n$ by

$$L(g) = (g(x_1), \dots, g(x_n))$$

Since $n > \dim(G)$, there exists $\alpha \neq \gamma \in \mathbb{R}^n$ s.t. $\langle \alpha, L(g) \rangle = 0 \forall g \in G$.

$$\Rightarrow \sum_{\{i : \gamma_i \leq 0\}} (-\gamma_i) g(x_i) = \sum_{\{i : \gamma_i > 0\}} \gamma_i g(x_i) \quad \forall g \in G$$

WLOG suppose $\gamma_i > 0$ for some i .

Suppose there were $g \in G$ such that S_g includes only $\{x_i : \gamma_i \leq 0\}$. Then

$$0 \geq \{LHS\} = \{RHS\} > 0$$

Contradiction \square

Ex: (Halfspaces)
Define $S_{ab} = \{x \in \mathbb{R}^d : \langle a, x \rangle + b \leq 0\}$

$$S = \{S_{a,b} : a, b \in \mathbb{R}\}$$

$$G = \{x \mapsto \langle a, x \rangle + b : a, b \in \mathbb{R}\}$$

Then

$$VC(S) \leq \dim(G) = d+1.$$

actually equality

Ex: (Balls)

$$\text{Define } S_{a,b} = \{x \in \mathbb{R}^d : \|x - a\|_2 \leq b\}$$

$$S = \{S_{a,b} : a \in \mathbb{R}^d, b \geq 0\}$$

$$\text{Define } g_{a,b}(x) = \|x\|_2^2 - 2\langle a, x \rangle + \|a\|_2^2 - b^2$$

Trick: Define $\Theta : \mathbb{R}^d \rightarrow \mathbb{R}^{d+2}$ by

$$\Theta(x) = (1, x_1, \dots, x_d, \|x\|_2^2)$$

$$g_c(x) = \langle c, \Theta(x) \rangle \text{ where } c \in \mathbb{R}^{d+2}$$

Then

$$\left\{ g_{a,b} : \begin{array}{l} a \in \mathbb{R} \\ b \geq 0 \end{array} \right\} \subseteq \left\{ g_c : c \in \mathbb{R}^d \right\}$$

vector space of dimension $d+2$.

$$\rightarrow VC(\mathcal{S}) \leq d+2$$

[Exact bound is $d+1$: harder to prove]

Rademacher Complexity and VC-dim often scale with the dimension of the ambient space $x \in \mathcal{X}$.

Thm: Consider $\min_{w \in W \subseteq \mathbb{R}^d} f(w) = \mathbb{E}_z f(w, z)$ where

$\max_{w \in W} \|w\| \leq B$ and $f(\cdot, z)$ is L -Lipschitz

Then
$$\mathbb{E} \left[\sup_{w \in W} \left(\frac{1}{n} \sum_{i=1}^n f(w, z_i) - \mathbb{E} f(w, z) \right) \right] \stackrel{\approx f(w)}{\leftarrow} \text{tight} \leq O \left(\sqrt{\frac{L^2 B^2 d \log(n)}{n}} \right)$$

[We'll prove this later!].

Dimension Independent Bound for Generalization

- linear models
- convexity

Linear Models:

Consider the problem

$$\min_{w \in W} \mathbb{E}_{(a,b) \sim P} l(\langle w, a \rangle, b)$$

To get generalization bounds we need to compute $C_n(\mathcal{F})$ where

$$\mathcal{F} = \{(a, b) \mapsto l(\langle w, a \rangle, b) : w \in W\}$$

Lemma: (Contraction)

Consider a set $A \subseteq \mathbb{R}^n$ and let

$\phi_i : \mathbb{R} \rightarrow \mathbb{R}$ be a β -Lipschitz function.

Let $A' = \{\phi_1(a_1), \dots, \phi_n(a_n) : a \in A\}$

Then $C(A') \leq C(A)$

pf: WLOG, assume $\beta = 1$. It suffices
to assume $A' = \{(\emptyset(a_i), a_2, \dots, a_n) : a \in A\}$

Then

$$C(A') = E_{\varepsilon} \left[\sup_{a \in A'} \sum_{i=1}^n \varepsilon_i a_i \right]$$

$$= \frac{1}{2} E_{\varepsilon_2, \dots, \varepsilon_n} \left[\sup_{a \in A} \left\{ \emptyset(a_i) + \sum_{i=2}^n \varepsilon_i a_i \right\} \right. \\ \left. + \sup_{\hat{a} \in A} \left\{ -\emptyset(\hat{a}_i) + \sum_{i=2}^n \varepsilon_i \hat{a}_i \right\} \right]$$

$$= \frac{1}{2} E_{\varepsilon_2, \dots, \varepsilon_n} \left[\sup_{\substack{a, \hat{a} \in A}} (\emptyset(a_i) - \emptyset(\hat{a}_i) + \sum_{i=2}^n \varepsilon_i (a_i + \hat{a}_i)) \right]$$

$$\leq \frac{1}{2} E_{\varepsilon_2, \dots, \varepsilon_n} \left[\sup_{\substack{a, \hat{a} \in A}} |a_i - \hat{a}_i| + \sum_{i=2}^n \varepsilon_i (a_i + \hat{a}_i) \right]$$

$$= \frac{1}{2} E_{\varepsilon_2, \dots, \varepsilon_n} \left[\sup_{\substack{a, \hat{a} \in A}} (a_i - \hat{a}_i) + \sum_{i=2}^n \varepsilon_i (a_i + \hat{a}_i) \right]$$

$$= C(A)$$

End

Rademacher Complexity of linear classes

Lemma: Consider a set of vectors $x_1, \dots, x_n \in \mathbb{R}^d$ and define

$$A = \left\{ (\langle w, x_1 \rangle, \dots, \langle w, x_n \rangle) : \|w\|_2 \leq 1 \right\}$$

Then $C(A) \leq \sqrt{\sum_{i=1}^n \|x_i\|_2^2}$

and therefore

$$C_n(\{x \mapsto (\langle w, x \rangle : \|w\|_2 \leq 1)\}) \leq 2 \frac{\max_{i=1, \dots, n} \|x_i\|_2}{\sqrt{n}}$$

Pf:

$$\begin{aligned} C(A) &= \mathbb{E}_\epsilon \sup_{a \in A} \sum_{i=1}^m \epsilon_i \cdot a_i \\ &= \mathbb{E}_\epsilon \sup_{\|w\|_2 \leq 1} \sum_{i=1}^m \epsilon_i \cdot \langle w, x_i \rangle \\ &= \mathbb{E}_\epsilon \sup_{\|w\|_2 \leq 1} \left\langle w, \sum_{i=1}^m \epsilon_i \cdot x_i \right\rangle \\ &= \mathbb{E}_\epsilon \left\| \sum_{i=1}^m \epsilon_i \cdot x_i \right\|_2 \\ &\leq \sqrt{\mathbb{E}_\epsilon \left\| \sum_{i=1}^m \epsilon_i \cdot x_i \right\|_2^2} \leq \sqrt{\sum_{i=1}^m \|x_i\|_2^2} \end{aligned}$$

□

Lemma: Consider a set of vectors $x_1, \dots, x_m \in \mathbb{R}^d$ and define

$$A = \left\{ (\langle w, x_i \rangle, \dots, \langle w, x_m \rangle) : \|w\| \leq 1 \right\}$$

Then

$$C(A) \leq \sqrt{2n \log(2d)} \cdot \max_{i=1, \dots, n} \|x_i\|_\infty$$

and therefore

$$C_n(\{x \mapsto (\langle w, x \rangle : \|w\| \leq 1)\}) \leq \frac{\sqrt{8 \log(6)}}{\sqrt{n}} \cdot \max_{i=1, \dots, n} \|x_i\|_\infty$$

Pf:

$$\begin{aligned} C(A) &= \mathbb{E}_\epsilon \sup_{a \in A} \sum_{i=1}^m \epsilon_i \cdot a_i \\ &= \mathbb{E}_\epsilon \sup_{\|w\| \leq 1} \sum_{i=1}^m \epsilon_i \cdot \langle w, x_i \rangle \\ &= \mathbb{E}_\epsilon \sup_{\|w\| \leq 1} \left\langle w, \sum_{i=1}^m \epsilon_i \cdot x_i \right\rangle \\ &= \mathbb{E}_\epsilon \left\| \sum_{i=1}^m \epsilon_i \cdot x_i \right\|_\infty \\ &= \mathbb{E}_\epsilon \sup_{V \in \{ \pm e_j \}_{j=1}^d} \sum_{i=1}^m \epsilon_i \langle x_i, V \rangle \end{aligned}$$

$$\mathbb{E}_\varepsilon \sup_{V \in \{\pm e_i\}_{i=1}^d} \sum_{i=1}^n \varepsilon_i \langle x_i, v \rangle$$

$$= C \left((\langle x_1, v \rangle, \dots, \langle x_n, v \rangle) : v \in \{\pm e_i\} \right)$$

$$\leq \sqrt{2 \log(2d)} \cdot \max_{V \in \{\pm e_i\}} \|(\langle x_1, v \rangle, \dots, \langle x_n, v \rangle)\|_2$$

$$\leq \sqrt{2n \log(2d)} \cdot \max_{i=1, \dots, n} \|x_i\|_\infty. \quad \square$$

Back to $\min_{w \in W} \mathbb{E}_{(a,b) \sim P} l(\langle w, a \rangle, b)$ where $W = B_2$ or B_1

So if $l(\cdot, b)$ is β -Lipschitz $\forall b$,
then the Rademacher bounds are

$$\underline{l_2\text{-case}}: P \cdot \frac{\mathbb{E} \max_{i=1, \dots, n} \|a_i\|_2}{\sqrt{n}}$$

$$\underline{l_1\text{-case}}: P \cdot \sqrt{\log(d)} \cdot \frac{\mathbb{E} \max_{i=1, \dots, n} \|a_i\|_\infty}{\sqrt{n}}$$

Convexity: Generalization without uniform laws

Suppose we want to solve

$$\textcircled{X} \quad \min_{w \in W} f(w) = \mathbb{E}_{x \sim P} f(w, x)$$

Let $S = \{x_1, \dots, x_n\}$ be iid from P , and let $t(S)$ be an output of an algorithm which aims to approximately solve \textcircled{X} . Let $S^i = (x_1, \dots, x_{i-1}, x', x_{i+1}, \dots, x_n)$ where $x' \sim P$ independent of x_1, \dots, x_n . Which $t(S)$ generalizes?

Intuition: $f(A(S^i), x_i) - f(A(S), x_i)$ should not be big, otherwise overfitting

Thm: $E_S [f(A(S)) - \frac{1}{n} \sum_{i=1}^n f(A(S), x_i)]$

$$= E_{\substack{(S, x') \sim P \\ i \sim U(n)}} [f(A(S^i), x_i) - f(A(S), x_i)]$$

Pf: for every i ,

$$\mathbb{E}_S f(A(S)) = \mathbb{E}_{S^{(2)}} f(A(S), x') = \mathbb{E}_{S^{(2)}} f(A(S^i), x_i)$$

Observe

$$\mathbb{E}_S \frac{1}{n} \sum_{i=1}^n f(A(S), x_i) = \mathbb{E}_{S, i} [f(A(S), x_i)]$$

Defn: $A(\cdot)$ is leave-one-out stable with rate $\epsilon(n)$ if

$$\mathbb{E}_{\substack{(S' x') \sim P^{n+1} \\ i \sim l(n)}} [f(A(S^i), x_i) - f(A(S), x_i)] \leq \epsilon(n).$$

Henceforth, fix $\lambda > 0$ and we'll analyze

$$A(S) := \arg \min_{w \in W} \frac{1}{n} \sum_{i=1}^n f(w, x_i) + \frac{\lambda}{2} \|w\|^2.$$

Also suppose W and $f(\cdot; x_i)$ are convex.

Defn: A function $g: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ is λ -strongly convex if $g - \frac{\lambda}{2} \|\cdot\|^2$ is convex.

Lemma: If g is λ -strongly convex, then it has a unique minimizer \bar{w} and

$$g(w) - g(\bar{w}) \geq \frac{\lambda}{2} \|w - \bar{w}\|^2 \quad \forall w$$

Thm: Suppose f is convex, and $f(x)$ is convex and ρ -Lipschitz. Then the rule

$$A(S) = \arg \min_{w \in W} \frac{1}{n} \sum_{i=1}^n f(w, x_i) + \frac{\lambda}{2} \|w\|^2$$

is leave-one-out-stable with rate $\frac{2\rho}{\lambda n}$

$$\text{Therefore } E_S [f(A(S)) - \frac{1}{n} \sum_{i=1}^n f(A(S), x_i)] \leq \frac{2\rho}{\lambda n}^2$$

Pf: Define $f_S(w) := \frac{1}{n} \sum_{i=1}^n f(w, x_i) + \frac{\lambda}{2} \|w\|^2$

$$\Rightarrow f_S(w) - f_S(A(S)) \geq \frac{\lambda}{2} \|w - A(S)\|^2 \quad \forall w$$

For all $w, v \in \mathbb{R}^n$ it holds

$$\begin{aligned} f_S(w) - f_S(v) &= \frac{1}{n} \sum_{x \in S} f(w, x) + \frac{\lambda}{2} \|w\|^2 \\ &\quad - \frac{1}{n} \sum_{x \in S} f(v, x) - \frac{\lambda}{2} \|v\|^2 \\ &= \frac{1}{n} \sum_{x \in S^c} f(w, x) + \frac{\lambda}{2} \|w\|^2 \\ &\quad - \frac{1}{n} \sum_{x \in S^c} f(v, x) - \frac{\lambda}{2} \|v\|^2 \\ &\quad + \frac{f(w, x_i) - f(v, x_i)}{n} + \frac{f(v, x^*) - f(w, x^*)}{n} \end{aligned}$$

Setting $w = A(S^c)$, $v = A(S)$ get

$$f_S(A(S^c)) - f_S(A(S)) \leq f_{S^c}(A(S^c)) - f_{S^c}(A(S))$$

$$\begin{aligned} &\leq -\frac{\lambda}{2} \|A(S^c) - A(S)\|^2 \\ &\quad + \frac{2f\|A(S^c) - A(S)\|}{n} \end{aligned}$$

$$\Rightarrow \|A(S^c) - A(S)\| \leq \frac{2f}{\lambda n} \quad \square$$

$$\text{Cor: } \mathbb{E}_S f(A(S)) \leq \min_{\bar{w}} f + \frac{\lambda}{2} \|\bar{w}\|^2 + \frac{2P^2}{n}$$

where \bar{w} is any minimizer of f on W .
 Therefore under optimal choice $\lambda = \sqrt{\frac{4P^2}{n\|\bar{w}\|^2}}$
 get

$$\mathbb{E}_S (f(A(S))) \leq \min f + 2 \sqrt{\frac{P^2 \|\bar{w}\|^2}{n}}$$

$$\begin{aligned} \text{pf: } \mathbb{E}_S f(A(S)) &= \mathbb{E}_S \left[\frac{1}{n} \sum_{x \in S} f(A(S), x) \right] \\ &\quad + \mathbb{E}_S [f(A(S)) - \frac{1}{n} \sum_{x \in S} f(A(S), x)] \\ &\leq \mathbb{E}_S \left[\frac{1}{n} \sum_{x \in S} f(A(S), x) \right] + \frac{2P^2}{n} \end{aligned}$$

For any w , we have

$$\begin{aligned} \mathbb{E}_S \left[\underbrace{\frac{1}{n} \sum_{x \in S} f(A(S), x)}_{f_S(A(S))} + \frac{\lambda}{2} \|A(S)\|^2 \right] &\leq \mathbb{E}_S f_S(w) \\ &\leq f(w) + \frac{\lambda}{2} \|w\|^2 \\ \Rightarrow \mathbb{E}_S \left[\frac{1}{n} \sum_{z \in S} f(A(S), z) \right] &\leq \min f + \frac{\lambda}{2} \|\bar{w}\|^2 \end{aligned}$$