

Chapter 3

Uniform Laws:

- Glivenko-Cantelli:
- sample complexity of learning
- uniform law w/ Rademacher Complexity
- Upper Bounds on Rademacher Complexity:
 - Massart Lemma
 - Finite sample Glivenko-Cantelli:
 - VC dimension
 - Dimension Independent generalization bounds
 - linear models
 - convex losses

Central Set-up:

Let \mathcal{F} be a family of integrable functions on some probability space (X, Σ, P) .

Goal: Estimate

$$\|P_n - P\|_{\mathcal{F}} := \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \int_{X \times P} f(x) \right|$$

where $x_1, \dots, x_n \stackrel{i.i.d.}{\sim} P$.

Ex: (Empirical CDF)

Suppose we want to estimate the CDF

$$g(t) := P[X \leq t]$$

We can use the empirical CDF

$$g_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{(-\infty, t]}(x_i)$$

Want to control

$$\|g_n - g\|_{\infty} := \sup_t |g_n(t) - g(t)|$$

This is exactly

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}f(X) \right|$$

where $\mathcal{F} = \left\{ \mathbb{1}_{(-\infty, t]}(\cdot) : t \in \mathbb{R} \right\}$.

Thm: (Glivenko-Cantelli.) For any distribution

$$\|g_n - g\|_{\infty} \xrightarrow{\text{a.s.}} 0$$

[We'll prove a finite sample version soon]

Cor (Plug in Estimator) [HW]

Let $\gamma(\cdot)$ be a $\|\cdot\|_{\infty}$ -continuous functional of CDF's. Then

$$\gamma(g_n) \rightarrow \gamma(g) \text{ in probability}$$

Ex: γ quantiles, median, mean
uniform distance to hypothesis' CDF.

Ex: (Statistical Learning)

Consider a family of probability measures

$$\{P_\theta : \theta \in \Omega\}$$

where Ω is some set.

Suppose we are given

$$x_1, \dots, x_n \stackrel{\text{iid}}{\sim} P_{\theta^*}$$

for some θ^* [may not be in Ω]

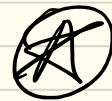
Fix a "goodness of fit" function
 $\theta \mapsto L_\theta(x)$

Goal is to minimize population risk

$$\min_{\theta \in \Omega} R(\theta) := \mathbb{E}_{x \sim P_{\theta^*}} [L_\theta(x)]$$

One approach is to instead minimize
the Empirical Risk :

$$\min_{\theta \in \Omega} R_n(\theta) := \frac{1}{n} \sum_{i=1}^n L_\theta(x_i)$$



Goal: Bound the excess risk

$$E(\hat{\theta}_n) := R(\hat{\theta}_n) - \inf_{\theta \in \Omega} R(\theta)$$

where $\hat{\theta}_n$ is the minimizer of R

Ex: (Maximum Likelihood)

Let p_θ be the density of P_θ

Define

$$L_\theta(x) = \log\left(\frac{p_{\theta^*}(x)}{p_\theta(x)}\right)$$

Then

$$\hat{\theta}_n \in \operatorname{argmax}_{\theta \in \Omega} \left\{ \frac{1}{n} \sum_{i=1}^n \log(p_\theta(x_i)) \right\}$$

Ex: (Binary Classification)

Data $(x_i, y_i) \in \mathbb{R}^d \times \{-1, +1\}$ according to P .

We want to estimate $f: \mathbb{R}^d \rightarrow \{-1, +1\}$ solving

$$\min_{f \in \mathcal{F}} P[f(x) \neq y] = \min_{f \in \mathcal{F}} \underbrace{E_{(x,y) \sim P} \mathbb{1}\{f(x) \neq y\}}_{R(f)}$$

↑
Some function class

Try instead to solve

$$\min_{f \in T} \frac{1}{n} \sum_{i=1}^n \mathbb{1}[f(x_i) \neq y_i]$$

Decomposition of Access Risk

$$E(\hat{\theta}_n, \theta^*) = R(\hat{\theta}_n) - \inf_{\theta \in \Omega} R(\theta)$$

$$= R(\hat{\theta}_n) - R_n(\hat{\theta}_n)$$

$$+ R_n(\hat{\theta}_n) - R_n(\theta_0)$$

$$+ R_n(\theta_0) - R(\theta_0)$$

where $\theta_0 = \operatorname{argmin}_{\theta \in \Omega} R(\theta)$

is approximation error How to control?

is optimization error. If $\hat{\theta}_n$ is true minimizer of $R_n(\theta)$, then

$$\leq 0$$

easy to control approximation error:

$$E R_n(\theta_0) = R(\theta_0)$$

$$\text{Bias} = R(\hat{\theta}_n) - R_n(\hat{\theta}_n)$$

Main Difficulty: $\hat{\theta}_n$ is not fixed!

But we can estimate

$$\begin{aligned} \text{Bias} &\leq \sup_{\theta \in \Omega} |R_n(\theta) - R(\theta)| \\ &= \sup_{\theta \in \Omega} \left| \frac{1}{n} \sum_{i=1}^n L_{\theta}(x_i) - \mathbb{E}_{X \sim P_{\theta}} L_{\theta}(X) \right| \end{aligned}$$

So

$$\mathbb{E}(\hat{\theta}_n, \theta^*) \leq 2 \|P_n - P\|_F$$

where $F := \{x \mapsto L_{\theta}(x) : \theta \in \Omega\}$

Remark: Function class F s.t.

$\|P_n - P\|_F \rightarrow 0$ in probability is called Glivenko-Cantelli.

- If F is too big, it may not be Glivenko-Cantelli (see Ex 4.7 in W)

Our first approach will be based on Rademacher Complexity.

Defn:

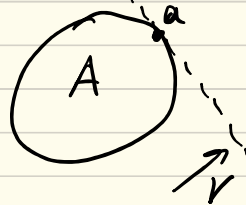
The Rademacher Complexity of a set $A \subseteq \mathbb{R}^n$, is the quantity

$$R(A) := \mathbb{E}_{\epsilon} \sup_{a \in A} \langle a, \epsilon \rangle$$

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

Recall that support function of a set A is

$$\sigma_A(v) := \sup_{a \in A} \langle v, a \rangle$$



So

$$R(A) = \mathbb{E}_{\epsilon} \sigma_A(\epsilon)$$

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

$$\begin{aligned} \mathcal{R}_n(\mathcal{F}) &:= \mathbb{E}_{x \in \mathcal{X}^n} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right| \\ &= \mathbb{E}_x \mathcal{R}(\mathcal{F}^\pm(x), \cdot) \\ &\leq 2 \mathbb{E}_x \mathcal{R}(\mathcal{F}(x)/n) \end{aligned}$$

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

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

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

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

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

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

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

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

$$\Pr \left(\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq 2 \mathcal{R}_n(\mathcal{F}) + t \right) \geq 1 - e^{-\frac{nt^2}{2b^2}}$$

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

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

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

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

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

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

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

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

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

[See Prop 4.1.1 in Wainwright]

Prop (HW) Suppose \mathcal{F} is b -bounded
Then $\forall n \in \mathbb{N}, t > 0$, it holds
$$\Pr(\|P_n - P\|_{\mathcal{F}} \geq \frac{1}{2} \mathcal{R}_n(\mathcal{F}) + \frac{\sup_{f \in \mathcal{F}} |\mathbb{E}f(x)|}{2\sqrt{n}} + t)$$
$$\geq 1 - e^{-\frac{nt^2}{2b^2}}$$

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

- $\mathcal{R}(cA) = |c| \mathcal{R}(A)$
- $\mathcal{R}(A+B) \leq \mathcal{R}(A) + \mathcal{R}(B)$
with equality if A and B are convex
- $\mathcal{R}(A+a_0) = \mathcal{R}(A)$
- $\mathcal{R}(A) = \mathcal{R}(\text{conv}(A))$

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

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

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

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

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

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

Jensen

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

Observe

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

$$\begin{aligned} &\leq \log \left(\sum_{a \in A'} e^{\frac{\|a\|^2}{2}} \right) \\ &\leq \log \left(|A'| \max_{a \in A'} \left(e^{\frac{\|a\|^2}{2}} \right) \right) \\ &= \log(|A'|) + \max_{a \in A'} \frac{\|a\|^2}{2} \end{aligned}$$

Thus

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

Optimize over $\lambda > 0$. \square

Defn: \mathcal{F} has polynomial discrimination of order $\nu \geq 1$ if $\forall n \in \mathbb{N}$ and $x_1, \dots, x_n \in \mathcal{X}$, we have $\text{card}(\mathcal{F}(x_1, \dots, x_n)) \leq (n+1)^\nu$.

Cor: Suppose \mathcal{F} has polynomial discrimination of order ν . Then

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

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

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

\square

In particular, if \mathcal{F} is b -bounded, then

$$\mathcal{R}_n(\mathcal{F}) \leq 4b \sqrt{\frac{5 \log(n+1)}{n}}$$

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

$$P\left[\|g_n - g\|_\infty \geq 8 \sqrt{\frac{\log(n+1)}{n}} + \delta\right] \leq \exp\left(-\frac{n\delta^3}{2}\right)$$

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

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

Since \mathcal{F} is 1-bounded

$$P\left[\|g_n - g\|_\infty \leq 2\mathcal{R}_n(\mathcal{F}) + \delta\right] \geq 1 - \exp\left(-\frac{n\delta^3}{2}\right)$$

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

Vapnik-Chervonenkis (VC) Theory

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

Consider a class \mathcal{F} of binary valued functions on X .

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

The VC dimension is

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

Notation: If $\mathcal{F} = \{ \mathbb{1}_S : \text{some sets } S \}$

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

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

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



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 \mathcal{F} has polynomial discrimination of order $VC(S)$ and

$$\mathcal{R}_n(\mathcal{F}) \leq 2 \sqrt{\frac{VC(S) \cdot \log(n+1)}{n}} \leftarrow \begin{array}{l} \log(n+1) \\ \text{can be} \\ \text{removed.} \end{array}$$

[See Prop 4.18 in Wainwright for a proof.]

Examples: Let G be a class of functions.

For any $g: \mathcal{X} \rightarrow \mathbb{R}$ define

$$\mathcal{S}_g = \{x \in \mathcal{X} : g(x) \leq 0\}$$

$$\mathcal{S}(G) = \{\mathcal{S}_g : g \in G\}$$

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

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

pt: 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 $0 \neq \gamma \in \mathbb{R}^n$ s.t. $\langle \gamma, L(g) \rangle = 0 \quad \forall g \in G$.

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

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

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

Contradiction \square

Ex: (Halfspaces)

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

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

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

Then

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

actually equality

Ex: (Balls)

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

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

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

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

$g_c(x) = \langle c, \Theta(x) \rangle$ 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 $\min_{w \in \mathcal{W}} f(w) = \mathbb{E}_z f(w, z)$ where

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

Then

$$\mathbb{E} \left[\sup_{w \in \mathcal{W}} \left| \frac{1}{n} \sum_{i=1}^n f(w, z_i) - \mathbb{E} f(w, z) \right| \right] \leq O \left(\sqrt{\frac{L^2 B^2 d \log(4)}{n}} \right)$$

[We'll prove this later!]

Dimension Independent Bound for Generalization

- linear models
- convexity

Linear Models:

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

To get generalization bounds we need to compute $R_n(\mathcal{F})$ where
$$\mathcal{F} = \{ (a,b) \mapsto \ell(\langle w, a \rangle, b) : w \in W \}$$

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 $R(A') \leq \beta R(A)$

pf: WLOG, assume $p=1$. It suffices to assume $A' = \{(\phi_1(a_1), a_2, \dots, a_n) : a \in A\}$

Then

$$\mathcal{R}(A') = \mathbb{E}_{\epsilon} \left[\sup_{a \in A'} \sum_{i=1}^n \epsilon_i a_i \right]$$

$$= \frac{1}{2} \mathbb{E}_{\epsilon_2, \dots, \epsilon_n} \left[\sup_{a \in A} \left\{ \phi_1(a) + \sum_{i=2}^n \epsilon_i a_i \right\} + \sup_{\hat{a} \in A} \left\{ -\phi_1(\hat{a}_1) + \sum_{i=2}^n \epsilon_i \hat{a}_i \right\} \right]$$

$$= \frac{1}{2} \mathbb{E}_{\epsilon_2, \dots, \epsilon_n} \left[\sup_{a, \hat{a} \in A} \left\{ \phi_1(a) - \phi_1(\hat{a}_1) + \sum_{i=2}^n \epsilon_i (a_i + \hat{a}_i) \right\} \right]$$

$$\leq \frac{1}{2} \mathbb{E}_{\epsilon_2, \dots, \epsilon_n} \left[\sup_{a, \hat{a} \in A} \left\{ a_1 - \hat{a}_1 + \sum_{i=2}^n \epsilon_i (a_i + \hat{a}_i) \right\} \right]$$

$$= \mathcal{R}(A)$$

□

Rademacher Complexity of linear classes

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

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

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

and therefore

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

Prf:

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

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

$$A = \{ (\langle w, x_1 \rangle, \dots, \langle w, x_m \rangle) : \|w\| \leq 1 \}$$

Then

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

and therefore

$$R_n(\{x_i \mapsto \langle w, x_i \rangle : \|w\| \leq 1\}) \leq \frac{\sqrt{2n \log(2d)} \cdot \max_{i=1, \dots, n} \|x_i\|_\infty}{\sqrt{n}}$$

Prf:

$$\begin{aligned} R(A) &= \mathbb{E}_\varepsilon \sup_{a \in A} \sum_{i=1}^m \varepsilon_i a_i \\ &= \mathbb{E}_\varepsilon \sup_{\|w\| \leq 1} \sum_{i=1}^m \varepsilon_i \langle w, x_i \rangle \\ &= \mathbb{E}_\varepsilon \sup_{\|w\| \leq 1} \left\langle w, \sum_{i=1}^m \varepsilon_i x_i \right\rangle \end{aligned}$$

$$= \mathbb{E}_\varepsilon \left\| \sum_{i=1}^m \varepsilon_i x_i \right\|_\infty$$

$$= \mathbb{E}_\varepsilon \sup_{v \in \{ \pm e_j \}_{j=1}^d} \sum_{i=1}^m \varepsilon_i \langle x_i, v \rangle$$

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

$$= \mathbb{R} \left((\langle x_1, v \rangle, \dots, \langle x_n, v \rangle) : v \in \{\pm e_j\} \right)$$

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

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

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

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

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

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

Convexity: Regret without uniform laws

Suppose we want to solve

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

Let $S = \{x_1, \dots, x_n\}$ be iid from P ,
and let $A(S)$ be an output of
an algorithm which aims to approximately
solve \star . 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 $A(S)$ generalizes?

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

Thm:

$$\mathbb{E}_S [f(A(S)) - \frac{1}{n} \sum_{i=1}^n f(A(S), x_i)] \\ = \mathbb{E}_{\substack{(S, x') \sim P^{n+1} \\ i \sim U(n)}} [f(A(S^i), x_i) - f(A(S), x_i)]$$

Pf: For every i ,

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

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

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

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

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

Also suppose \mathcal{W} and $f(\cdot, x_i)$ are convex.

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

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

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

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

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

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

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

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 $w, v \in \mathcal{W}$ it holds

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

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

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

$$\begin{aligned}
 &+ \frac{\lambda}{2} \|A(S^i) - A(S)\|^2 \\
 &+ \frac{2\rho \|A(S^i) - A(S)\|}{n}
 \end{aligned}$$

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

Cor: $\mathbb{E}_S f(A(S)) \leq \min_W f + \frac{\lambda}{2} \|\bar{w}\|^2 + \frac{2\rho^3}{\lambda n}$
 where \bar{w} is any minimizer of f on W .
 Therefore under optimal choice $\lambda = \sqrt{\frac{4\rho^2}{n\|\bar{w}\|^2}}$

get

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

pf: $\mathbb{E}_S f(A(S)) = \mathbb{E}_S \left[\frac{1}{n} \sum_{x \in S} f(A(S), x) \right]$
 $+ \mathbb{E}_S \left[f(A(S)) - \frac{1}{n} \sum_{x \in S} f(A(S), x) \right]$
 $\leq \mathbb{E}_S \left[\frac{1}{n} \sum_{x \in S} f(A(S), x) \right] + \frac{2\rho^3}{\lambda n}$

For any w , we have

$$\mathbb{E}_S \left[\frac{1}{n} \sum_{x \in S} f(A(S), x) + \frac{\lambda}{2} \|A(S)\|^2 \right] \leq \mathbb{E}_S f_S(w)$$

$$\leq f(w) + \frac{\lambda}{2} \|w\|^2$$

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