

STAT576
Empirical Process Theory

Pingbang Hu

October 24, 2023

Abstract

This is a graduate-level theoretical statistics course taught by [Sabyasachi Chatterjee](#) at University of Illinois Urbana-Champaign, aiming to provide an introduction to empirical process theory with applications to statistical M -estimation, non-parametric regression, classification and high dimensional statistics.

While there are no required textbooks, some books do cover (almost all) part of the material in the class, e.g., Van Der Vaart and Wellner's *Weak Convergence and Empirical Processes* [[VW96](#)].



This course is taken in Fall 2023, and the date on the covering page is the last updated time.

Contents

1	Introduction	2
1.1	What is Empirical Process Theory?	2
1.2	Applications of Uniform Law of Large Numbers	3
1.3	Bounding Supremum of Empirical Process	5
2	Concentration Inequalities	6
2.1	Gaussian Distribution	6
2.2	MGF Trick	7
2.3	Hoeffding's Inequality	8
2.4	Bernstein's Inequality	11
2.5	Bounded Difference Concentration Inequality	14
3	Expected Supremum of Empirical Process	20
3.1	Statistical Learning	20
3.2	Vapnik-Chervonenkis Dimension	26
3.3	Metric Entropy Methods	29
3.4	Bracketing Bound	50
4	Applications to M-Estimation	54
4.1	The M -Estimation Problem	54
4.2	Consistency	55
4.3	Rate of Convergence	56
4.4	Non-Parametric Regression	65

Chapter 1

Introduction

Lecture 1: Introduction to Mathematical Statistics

1.1 What is Empirical Process Theory?

21 Aug. 9:00

This subject started in the 1930s with the study of the [empirical CDF](#).

Definition 1.1.1 (Empirical CDF). Given inputs i.i.d. data points $X_1, \dots, X_n \sim \mathbb{P}$, the *empirical CDF* is

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_i \leq t}.$$

The classical result is that, fixing t , $F_n(t) \rightarrow F(t)$ almost surely.

Note. At the same time, $\sqrt{n}(F_n(t) - F(t)) \rightarrow \mathcal{N}(0, F(t)(1 - F(t)))$ in distribution.

On the other hand, we can also ask does this convergence happen if we jointly consider all possible $t \in \mathbb{R}$. By the [Glivenko-Cantelli theorem](#), $\sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \xrightarrow{n \rightarrow \infty} 0$ almost surely, so the answer is again yes.

Now, we're ready to see a "canonical" example of an [empirical process](#).

Example (Canonical empirical process). The *canonical empirical process* is the family of random variables $\{F_n(t)\}_{t \in \mathbb{R}}$, i.e., a stochastic process.

By considering a general class of functions, we have the following.

Definition 1.1.2 (Empirical process). Let χ be the domain, \mathbb{P} be a distribution on χ , and \mathcal{F} be the class of function such that $\chi \rightarrow \mathbb{R}$. The *empirical process* is the stochastic process indexed by functions in \mathcal{F} , $\{G_n(f) : f \in \mathcal{F}\}$ where

$$G_n(f) = \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}[f(X)]$$

and $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$.

Remark. The [empirical process](#) is a family of mutually dependent random variables, all of them being functions of the same inherent randomness in the i.i.d. data X_1, \dots, X_n .

Now, two questions arises.

1.1.1 Uniform Law of Large Numbers

As $n \rightarrow \infty$, whether

$$S_n(\mathcal{F}) := \sup_{f \in \mathcal{F}} |G_n(f)| \rightarrow 0,$$

and if, at what rate?

Remark. The rate of convergence of law of large numbers uniformly over a class of functions \mathcal{F} determines the performance of many types of statistical estimators as we will see.

We will spend most of this course just on this topic with applications. We will show that $S(\mathcal{F})$ concentrates around its expectation and will bound $\mathbb{E}[S(\mathcal{F})]$.

1.1.2 Uniform Central Limit Theorem

The most general probabilistic question one can ask is the following.

Problem. What is the joint distribution of the [empirical process](#)?

Answer. For a given sample size, it's most often intractable to be able to calculate the joint distribution exactly. One can then use asymptotics when the sample size n is very large to derive limiting distributions. By the regular central limit theorem, $\sqrt{n}G_n(f) \xrightarrow{d} \mathcal{N}(0, \text{Var}[f(X)])$ for any f . We want to understand if this holds uniformly (jointly) over $f \in \mathcal{F}$ in some sense. \circledast

We first motivate this through an example.

Example (Uniform empirical process). Consider

- X_1, \dots, X_n i.i.d. from $\mathcal{U}(0, 1)$.^a
- $\mathcal{F} = \{\mathbb{1}_{[-\infty, t]} : t \in \mathbb{R}\}$
- $U_n(t) = \sqrt{n}(F_n(t) - t)$ where F_n is the [empirical CDF](#).

We can view $U_n(t)$ as collection of random variables one for each $t \in (0, 1)$, or just as a random function. Then this stochastic process $\{U_n(t) : t \in (0, 1)\}$ is called the *uniform empirical process*.

Then, the CLT states that for each $t \in [0, 1]$, $U_n(t) \rightarrow \mathcal{N}(0, t - t^2)$ as $n \rightarrow \infty$. Moreover, for fixed t_1, \dots, t_k , the multivariate CLT implies that $(U_n(t_1), \dots, U_n(t_k)) \xrightarrow{d} \mathcal{N}(0, \Sigma)$ where $\Sigma_{ij} = \min(t_i, t_j) - t_i t_j$.

^a \mathcal{U} denotes the uniform distribution.

From this example, one can ask question like the following.

Problem. Does the entire process $\{U_n(t) : t \in [0, 1]\}$ converge in some sense? If so, what is the limiting process?

Answer. The limiting process is an object called the *Brownian Bridge*. This was conjectured by Doob and proved by Donsker. \circledast

Other than that, how do we characterize convergence of stochastic processes in distribution to another stochastic process? How do we generalize this result for a general function class \mathcal{F} defined on a probability space χ ? What are some statistical applications of such process convergence results? This is a classical topic and in the last few weeks of this course, we will touch upon some of these questions.

1.2 Applications of Uniform Law of Large Numbers

Next, we see one major example where uniform law of large numbers can be applied.

1.2.1 M -Estimators

Consider the class of estimators called “ M -estimator”, which is of the form

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n M_{\theta}(X_i),$$

where X_1, \dots, X_n taking values in χ , Θ is the parameter space, and $M_{\theta}: \chi \rightarrow \mathbb{R}$ for each $\theta \in \Theta$. Let's see some examples.

Example (Maximum log-likelihood). $M_{\theta}(X) = -\log p_{\theta}(X)$ for a class of densities $\{p_{\theta}: \theta \in \Theta\}$, then $\hat{\theta}$ is the *Maximum log-likelihood* of θ .

There are lots of examples on “local estimators” as well.

Example (Mean). $M_{\theta}(x) = (x - \theta)^2$.

Example (Median). $M_{\theta}(x) = |x - \theta|$.

Example (τ quantile). $M_{\theta}(x) = Q_{\tau}(x - \theta)$ where $Q_{\tau}(x) = (1 - \tau)x\mathbb{1}_{x < 0} + \tau x\mathbb{1}_{x \geq 0}$.

Example (Mode). $M_{\theta}(x) = -\mathbb{1}_{|x - \theta| \leq 1}$.

Now, the target quantity for the estimator $\hat{\theta}$ is

$$\theta_0 = \arg \max_{\theta \in \Theta} \mathbb{E} [M_{\theta}(X_1)]$$

where $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$. In the asymptotic framework, the two key questions are the following.

Problem. Is $\hat{\theta}$ consistent for θ_0 ? Does $\hat{\theta}$ converge to θ_0 almost surely or in probability as $n \rightarrow \infty$? I.e., is $d(\hat{\theta}, \theta_0) \rightarrow 0$ for some metric d ?

Problem. What is the rate of convergence of $d(\hat{\theta}, \theta_0)$? For example is it $O(n^{-1/2})$ or $O(n^{-1/3})$?

To answer these questions, one is led to investigate the closeness of the empirical objective function to the population objective function in some uniform sense. Consider $M_n(\theta) = \frac{1}{n} \sum_{i=1}^n M_{\theta}(X_i)$ and $M(\theta) = \mathbb{E} [M_{\theta}(X_1)]$, then

$$\begin{aligned} \mathbb{P}(d(\hat{\theta}, \theta_0) > \epsilon) &\leq \mathbb{P} \left(\sup_{\theta: d(\theta, \theta_0) > \epsilon} M_n(\theta_0) - M_n(\theta) \geq 0 \right) \\ &= \mathbb{P} \left(\sup_{\theta: d(\theta, \theta_0) > \epsilon} (M_n(\theta_0) - M(\theta_0) - [M_n(\theta) - M(\theta)]) \geq \inf_{\theta: d(\theta, \theta_0) > \epsilon} (M(\theta) - M(\theta_0)) \right) \\ &\leq \mathbb{P} \left(2 \sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \geq \inf_{\theta: d(\theta, \theta_0) > \epsilon} (M(\theta) - M(\theta_0)) \right). \end{aligned}$$

We see that the left-hand side $2 \sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)|$ is just $S(\mathcal{F})$ for $\mathcal{F} = \{f_{\theta}: \theta \in \Theta, f_{\theta} = M_{\theta}(\cdot)\}$, while the right-hand side $\inf_{\theta: d(\theta, \theta_0) > \epsilon} M(\theta) - M(\theta_0)$ is larger than 0.

Remark. The last step could be too loose in some problems.

Lecture 2: Sub-Gaussian Random Variables and the MGF Trick

1.3 Bounding Supremum of Empirical Process

Most of this course will focus on bounding suprema of the [empirical process](#). Let's define it rigorously.

Problem 1.3.1 (Bounding supremum of empirical process). Given a domain χ , a probability measure \mathbb{P} on χ , data $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$, and a function class $\mathcal{F} \ni f: \chi \rightarrow \mathbb{R}$. We want to find an (non-asymptotically) bound on

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

Answer. To do this, broadly speaking, we will go through a route with three basic steps:

- (a) $S_n(\mathcal{F})$ “concentrates” around its expectation $\mathbb{E}[S_n(\mathcal{F})]$.
- (b) $\mathbb{E}[S_n(\mathcal{F})] \leq$ the [Rademacher complexity](#) of \mathcal{F} via “[symmetrization](#)”.
- (c) Bounding the [Rademacher complexity](#)’s expected supremum of a “sub-Gaussian process” by a technique called *chaining*.

*

Toward this end, we need some basic and fundamental concentration inequalities which are of wide interest and use.

Chapter 2

Concentration Inequalities

As we just saw, to solve [Problem 1.3.1](#), we need some basic tools on concentration inequalities. The most celebrated concentration inequality might be the Gaussian tail, which achieves a quadratic exponential decay. Combine this with the classical central limit theorem, we can expect that as $n \rightarrow \infty$, approximately the Gaussian tail bound kicks in.

However, to get a concrete, non-asymptotic bound for $S_n(\mathcal{F})$, we would need more sophisticated tools. Let's start with the basics, i.e., the Gaussian distribution.

2.1 Gaussian Distribution

For us, the gold standard for concentration would be the Gaussian distribution. The property of the Gaussian distribution we are interested in is its rapid tail decay as we mentioned:

Lemma 2.1.1. For $Z \sim \mathcal{N}(0, 1)$,

$$\left(\frac{1}{t} - \frac{1}{t^3}\right) \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \leq \mathbb{P}(Z \geq t) \leq \frac{1}{t} \cdot \frac{1}{\sqrt{2\pi}} e^{-t^2/2}.$$

Proof. We want to show

$$\begin{aligned} \left(\frac{1}{t} - \frac{1}{t^3}\right) \frac{1}{\sqrt{2\pi}} e^{-t^2/2} &\leq \int_t^\infty \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \leq \frac{1}{t} \cdot \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \\ \Leftrightarrow \left(\frac{1}{t} - \frac{1}{t^3}\right) e^{-t^2/2} &\leq \int_t^\infty e^{-x^2/2} dx \leq \frac{1}{t} \cdot e^{-t^2/2}. \end{aligned}$$

Observe that from integration by part (with x/x introduced),

$$\int_t^\infty \frac{x}{x} \cdot e^{-x^2/2} dx = -\frac{e^{-x^2/2}}{x} \Big|_t^\infty - \int_t^\infty \frac{e^{-x^2/2}}{x^2} dx = \frac{e^{-t^2/2}}{t} - \int_t^\infty \frac{e^{-x^2/2}}{x^2} dx \leq \frac{1}{t} \cdot e^{-t^2/2}$$

since the integrand $e^{-x^2/2}/x^2$ is non-negative, which is the desired upper-bound. For the lower bound, if we again apply integration by part (with x/x introduced again), then

$$\begin{aligned} \int_t^\infty e^{-x^2/2} dx &= \frac{e^{-t^2/2}}{t} - \int_t^\infty \frac{x}{x} \cdot \frac{e^{-x^2/2}}{x^2} dx \\ &= \frac{e^{-t^2/2}}{t} - \left(-\frac{e^{-x^2/2}}{x^3} \Big|_t^\infty - \int_t^\infty 3 \frac{e^{-x^2/2}}{x^4} dx \right) \\ &= \frac{e^{-t^2/2}}{t} - \frac{e^{-t^2/2}}{t^3} + \int_t^\infty 3 \frac{e^{-x^2/2}}{x^4} dx \\ &\geq \left(\frac{1}{t} - \frac{1}{t^3}\right) e^{-t^2/2}, \end{aligned}$$

since, again, the integrand $3e^{-x^2/2}/x^4$ is non-negative, so the last term is positive, hence we get the desired lower-bound. ■

Corollary 2.1.1. For all $t \geq 1$, we have

$$\mathbb{P}(\mathcal{N}(0, \sigma^2) \geq t) \leq e^{-t^2/2\sigma^2}.$$

Now, as is suggested by CLT, the following question arises.

Problem. Does [Corollary 2.1.1](#) hold for sums of independent random variables? That is, given i.i.d. X_1, \dots, X_n with mean μ and variance σ^2 , whether for all $t \geq 0$,

$$\mathbb{P}(\sqrt{n}(\bar{X} - \mu) \geq t) \leq e^{-t^2/2\sigma^2}?$$

Answer. Just invoking CLT is not enough, we need to handle the error term in the normal approximation. We can show this directly for a class of distributions with fast tail decay. ⊛

To go beyond Gaussian tail bound, let start with the [moment generating function \(MGF\) trick](#).

2.2 MGF Trick

The [MGF trick](#) is easy to develop, but it gives a foundation of all the concentration inequalities we're going to develop. Hence, although it's short, it's worth to make it a separate section.

2.2.1 Markov's Inequality

To start with, the most basic tool to bound tail probabilities is the [Markov's inequality](#).

Lemma 2.2.1 (Markov's inequality). For a non-negative random variable $X \geq 0$,

$$\mathbb{P}(X \geq t) \leq \frac{\mathbb{E}[X]}{t}.$$

Note. [Markov's inequality](#) is valid as soon as $\mathbb{E}[X] < \infty$. That is, it holds even when the second moment does not exist.

Remark. The rate of tail decay is slow ($O(1/t)$). For the Gaussian, by [Lemma 2.1.1](#), it's $O(e^{-t^2/2})$.

By the above remark, one might ask the following.

Problem. Can we derive faster tail decay bounds in general?

Answer. Yes, if we assume more moments exist. If all moments exist and in particular the MGF exists, like for the Gaussian, then we can expect faster tail decay. ⊛

2.2.2 Chebyshev Inequality

Continuing the discussion on the previous problem, for example, if we assume second moment exists, then we can get an $O(1/t^2)$ tail decay by [Chebyshev inequality](#).

Lemma 2.2.2 (Generalized Chebyshev inequality). Given a random variable X ,

$$\mathbb{P}(|X - \mu| \geq t) = \mathbb{P}(|X - \mu|^p \geq t^p) \leq \min_{p \geq 1} \frac{\mathbb{E}[|X - \mu|^p]}{t^p}.$$

Proof. This is directly implied by the [Markov's inequality](#). ■

Remark (Chebyshev Inequality). For $p = 2$, we have the usual form

$$\mathbb{P}(|X - \mu| \geq t) \leq \frac{\text{Var}[X]}{t^2}$$

Remark. All tail bounds are derived using [Markov's inequality](#); the clever part is to apply it to the right random variable. In this sense, every tail bound is just [Markov's inequality](#).

2.2.3 Cramer-Chernoff Method

In the same vein, developed by Cramer and Chernoff, if we now assume the MGF exists and apply [Markov's inequality](#), we get the [MGF trick](#).

Lemma 2.2.3 (MGF trick (Cramer-Chernoff method)). Given a random variable X ,

$$\mathbb{P}(X - \mu \geq t) = \mathbb{P}(e^{\lambda(X-\mu)} \geq e^{\lambda t}) \leq \inf_{\lambda > 0} \frac{\mathbb{E}[e^{\lambda(X-\mu)}]}{e^{\lambda t}}.$$

We will use the [MGF trick](#) rather than the [generalized Chebyshev's inequality](#) to derive tail bounds because MGF of a sum of independent random variables decomposes as the product of the MGF's. It is messier to work with the p^{th} moment of a sum of independent random variables.

2.3 Hoeffding's Inequality

2.3.1 Sub-Gaussian Random Variables

We will now consider a class of distributions whose MGF is dominated by the MGF of a Gaussian. Then, in a very clean way, the [MGF trick](#) will give us Gaussian tail bounds for these distributions.

Definition 2.3.1 (Sub-Gaussian). Given a random variable X with $\mathbb{E}[X] = 0$, we say X is *sub-Gaussian* with variance factor^a σ^2 if for all $\lambda \in \mathbb{R}$,

$$\mathbb{E}[e^{\lambda X}] \leq e^{\frac{\sigma^2 \lambda^2}{2}}.$$

^aAlso called proxy, sub-Gaussian norm, etc.

Notation. We write $\text{Subg}(\sigma^2)$ for a compact representation of the class of [sub-Gaussian](#) random variables with variance factor σ^2 .

Remark. Observe that if $X \in \text{Subg}(\sigma^2)$:

- $-X \in \text{Subg}(\sigma^2)$;
- $X \in \text{Subg}(t^2)$ if $t^2 > \sigma^2$;
- $cX \in \text{Subg}(c\sigma^2)$.

Lemma 2.3.1 (Equivalent conditions). Given a random variable X with $\mathbb{E}[X] = 0$, the following are equivalent for absolute constants $c_1, \dots, c_5 > 0$.

- (a) $\mathbb{E}[e^{\lambda X}] \leq e^{c_1^2 \lambda^2}$ for all $\lambda \in \mathbb{R}$.
- (b) $\mathbb{P}(|X| \geq t) \leq 2e^{-t^2/c_2^2}$.
- (c) $(\mathbb{E}[|X|^p])^{1/p} \leq c_3 \sqrt{p}$.

Add proof

(d) For all λ such that $|\lambda| \leq 1/c_4$, $\mathbb{E} [e^{\lambda^2 X^2}] \leq e^{c_4^2 \lambda^2}$.

(e) For some $c_5 < \infty$, $\mathbb{E} [e^{X^2/c_5^2}] \leq 2$.

Proof. Let's just see the first implication from (a) to (b). Given $X \in \text{Subg}(\sigma)$,

$$\mathbb{P}(X \geq t) \leq \inf_{\lambda > 0} e^{\lambda^2 \sigma^2 / 2 - \lambda t} \leq e^{-\frac{t^2}{2\sigma^2}}$$

where the last inequality follows from minimizing the quadratic function $\lambda^2 \sigma^2 / 2 - \lambda t$ whose minimizer is $\lambda^* = t/\sigma^2$. The same bound holds for the left tail and a union bound gives the two-sided version. ■

Let's see some examples of the **sub-Gaussian** random variables.

Example (Rademacher random variable). $\epsilon = \pm 1$ with probability $1/2$ is a $\text{Subg}(1)$ random variable.

Proof. We see that

$$\mathbb{E} [e^{\lambda \epsilon}] = \frac{1}{2} e^{\lambda} + \frac{1}{2} e^{-\lambda} = \frac{1}{2} \sum_{k=1}^{\infty} \left(\frac{\lambda^k}{k!} + \frac{(-\lambda)^k}{k!} \right) = \sum_{k=1}^{\infty} \frac{\lambda^{2k}}{(2k)!} \leq 1 + \sum_{k=1}^{\infty} \frac{(\lambda^2)^k}{2^k k!} = e^{\lambda^2/2}$$

since $(2k)! \geq 2^k \cdot k!$. *

In fact, the above can be generalized for any bounded random variable.

Lemma 2.3.2. Given $X \in [a, b]$ such that $\mathbb{E} [X] = 0$. Then

$$\mathbb{E} [e^{\lambda X}] \leq \exp \left(\lambda^2 \frac{(b-a)^2}{8} \right)$$

for all $\lambda \in \mathbb{R}$, i.e., $X \in \text{Subg}((b-a)^2/4)$.

Proof. We will prove this with a worse constant. Let $X' \stackrel{\text{i.i.d.}}{\sim} X$ be an i.i.d. copy, then

$$\mathbb{E} [e^{\lambda X}] = \mathbb{E} [e^{\lambda(X - \mathbb{E}[X'])}] = \mathbb{E} [e^{\lambda X} \cdot e^{-\lambda \mathbb{E}[X']}] \leq \mathbb{E} [e^{\lambda X}] \cdot \mathbb{E} [e^{-\lambda X'}] = \mathbb{E} [e^{\lambda(X - X')}] ,$$

where we have used the **Jensen's inequality** for $e^{-\lambda \mathbb{E}[X']} \leq \mathbb{E} [e^{-\lambda X'}]$.^a Now we introduce a **Rademacher random variable** $\epsilon = \pm 1$, to further write

$$\mathbb{E} [e^{\lambda X}] \leq \mathbb{E}_{X, X'} [e^{\lambda(X - X')}] = \mathbb{E}_{X, X', \epsilon} [e^{\lambda \epsilon (X - X')}] = \mathbb{E}_{X, X'} [\mathbb{E}_{\epsilon} [e^{\lambda \epsilon (X - X')}]] ,$$

and $\mathbb{E}_{\epsilon} [e^{\lambda \epsilon (X - X')}] \leq \mathbb{E} [e^{\frac{\lambda^2 (X - X')^2}{2}}] \leq e^{\frac{\lambda^2 (b-a)^2}{2}}$, where we used the known bound on MGF of a **Rademacher random variable**, hence overall, we get

$$\mathbb{E} [e^{\lambda X}] \leq \mathbb{E}_{X, X'} \left[e^{\frac{\lambda^2 (b-a)^2}{2}} \right] = e^{\frac{\lambda^2 (b-a)^2}{2}} .$$

■

^aThis is a trick called symmetrization. A basic example is $\text{Var} [X] = \frac{1}{2} \mathbb{E} [(X - X')^2]$.

Note. If $a = -1$ and $b = 1$, we get back to the earlier example.

Just like independent Gaussians, sums of independent **sub-Gaussians** remain **sub-Gaussian**.

Lemma 2.3.3 (Closed under convolution). Let X_i be independent random variables with $\mathbb{E} [X_i] = \mu_i$,

and $X_i - \mu_i \in \text{Subg}(\sigma_i^2)$. Then

$$\sum_{i=1}^n X_i - \sum_{i=1}^n \mu_i \in \text{Subg}\left(\sum_{i=1}^n \sigma_i^2\right).$$

Proof. We simply observe that

$$\mathbb{E}\left[e^{\lambda \sum_{i=1}^n (X_i - \mu_i)}\right] = \prod_{i=1}^n \mathbb{E}\left[e^{\lambda (X_i - \mu_i)}\right] \leq e^{\frac{\lambda^2 (\sum_{i=1}^n \sigma_i^2)}{2}}.$$

■

2.3.2 Hoeffding's Inequality

We can now immediately prove the famous [Hoeffding's inequality](#), which is the main tool in our interest.

Theorem 2.3.1 (Hoeffding's inequality for sub-Gaussian random variables). Let X_i be independent random variables with $\mathbb{E}[X_i] = \mu_i$, and $X_i - \mu_i \in \text{Subg}(\sigma_i^2)$. Then for all $t \geq 0$,^a

$$\mathbb{P}\left(\left|\sum_{i=1}^n (X_i - \mu_i)\right| \geq t\right) \leq 2 \exp\left(\frac{-t^2}{2 \sum_{i=1}^n \sigma_i^2}\right).$$

^aOne-sided version holds without the factor 2.

Proof. It's immediate from [Lemma 2.3.3](#) and the equivalent condition (b) in [Lemma 2.3.1](#). ■

Lecture 3: Sub-Exponential Random Variables

For bounded random variables, we can apply [Hoeffding's inequality](#) to obtain the following.

25 Aug. 9:00

Corollary 2.3.1. Let $X_i \in [a, b]$ be random variables with mean μ_i ,

$$\mathbb{P}\left(\sum_i (X_i - \mu_i) \geq t\right) \leq \exp\left(-\frac{2t^2}{n(b-a)^2}\right).$$

As a consequence, if X_i are i.i.d., then

$$\mathbb{P}(\sqrt{n}(\bar{X} - \mu) \geq t) \leq \exp\left(-\frac{2t^2}{(b-a)^2}\right).$$

Compare this with Gaussian approximation, we then have

$$\mathbb{P}(\sqrt{n}(\bar{X} - \mu) \geq t) \approx \mathbb{P}(\mathcal{N}(0, \sigma^2) \geq t) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right),$$

i.e., $\sigma^2 \sim (b-a)^2/4$.¹

Remark (Comparison between Hoeffding's bound and Gaussian tail bound). We see that

- (a) [Hoeffding's inequality](#) can be used for any sample size, but Gaussian approximation can only be used when n is large.
- (b) As $\sigma^2 \leq (b-a)^2/4$, we see that Gaussian approximation gives a tighter tail bound.
- (c) Another way to state this is that from CLT we get the asymptotically valid confidence interval

¹Actually, $\sigma^2 \leq (b-a)^2/4$ always holds.

for μ as

$$\left[\bar{X} \pm \frac{\sigma}{\sqrt{n}} Z_{\alpha/2} \right],$$

while from the [Hoeffding's inequality](#), we have (finite sample valid) confidence interval

$$\left[\bar{X} \pm \frac{b-a}{2\sqrt{n}} \sqrt{\log \frac{2}{\alpha}} \right],$$

which is much larger.

The above discussion suggests that if the range is very large compared to the variance, then [Hoeffding's inequality](#) may not perform very well. Clearly, such random variables exist. Here are some examples.

Example. Suppose

$$\mathbb{P}(X = 0) = 1 - 1/k^2$$

$$\mathbb{P}(X = \pm K) = 1/2k^2$$

with $\mathbb{E}[X] = 0$ and $\text{Var}[X] \leq 1$. The range is $2K$, which is very large compared to the variance. This is a case where [Hoeffding's inequality](#) would not perform very well, in the sense that the confidence interval based on it would be too wide.

Another example is the following.

Example. Let X_1, \dots, X_n be i.i.d. Bernoulli(λ/n), where each one of them has range 1, but its variance is at most $\frac{\lambda}{n} \ll 1$. Then a direct application of [Hoeffding's inequality](#) gives

$$\mathbb{P}\left(\sum_i X_i - \lambda \geq t\right) \leq \exp\left(\frac{-2t^2}{n}\right).$$

This suggests that $\sum_i X_i = O(\sqrt{n})$ whereas we know that in this case that the distribution of $\sum_i X_i$ is close to the Poisson(λ) and thus should be $O(1)$.

On the other hand, the CLT inspired bound would give the right order. This points out that we would like to be able to replace the range term by the variance in [Hoeffding's inequality](#). This is what is done in [Bernstein's inequality](#) which we will discuss next.

Let's see some non-examples.

Example (Not sub-Gaussian). Some examples of random variables which are not [sub-Gaussians](#) random variables are Cauchy, exponential, and Poisson random variables.

What about mixture?

Problem. Suppose $Z_1, Z_2 \in \text{Subg}(\sigma^2)$ with mean 0, and consider

$$X = \begin{cases} Z_1, & \text{w.p. } p; \\ Z_2, & \text{w.p. } 1 - p. \end{cases}$$

Is this a [sub-Gaussian](#) random variable?

2.4 Bernstein's Inequality

2.4.1 Sub-Exponential Random Variables

The main reason for considering the class of [sub-Gaussian](#) random variables is that the MGF is finite and thus the [MGF trick](#) works. So if we want to extend the [MGF trick](#), we would like to ask the following:

Problem. How fat could the tails of a distribution be so that the MGF is finite?

Answer. It turns out that we can allow fatter tails than [sub-Gaussian](#), essentially the PDF can decay no slower than an exponential with a proper exponent. \circledast

Consider the following example.

Example. Let $Z^2 \sim \chi^2$, then for all $t \geq 1$, $\mathbb{P}(Z^2 > t) = 2\mathbb{P}(Z \geq \sqrt{t}) \leq 2e^{-t/2}$. It is seen that the rate of decrease of the χ^2 tail probability is slower than that of normal. In fact, the MGF of χ^2 is

$$\mathbb{E} \left[e^{\lambda(Z^2-1)} \right] = \begin{cases} \frac{e^{-\lambda}}{\sqrt{1-2\lambda}}, & \text{if } 0 \leq \lambda < 1/2; \\ \infty, & \text{if } \lambda \geq 1/2, \end{cases}$$

where we see that the MGF exists in a neighborhood around 0, but not everywhere.

This motivates the following definition.

Definition 2.4.1 (Sub-exponential). A random variable X is *sub-exponential* with parameters (σ^2, α) with mean λ if for all $|\lambda| < 1/\alpha$

$$\mathbb{E} \left[e^{\lambda(X-\mu)} \right] \leq e^{\frac{\lambda^2 \sigma^2}{2}}.$$

It's then immediate to see that $\text{SubExp}(\sigma^2, \alpha)$ random variables have the same bound on their MGF as a $\text{SubG}(\sigma^2)$ but only for λ in the interval $(-\frac{1}{\alpha}, \frac{1}{\alpha})$.

Example. For the χ^2 random variable Z^2 , we have $Z^2 \in \text{SubExp}(2, 4)$.

Proof. This is immediate from [Definition 2.4.1](#) since For all $|\lambda| < 1/4$, we have

$$\frac{e^{-\lambda}}{\sqrt{1-2\lambda}} \leq e^{2\lambda^2}.$$

\circledast

With [Definition 2.4.1](#), we can extend the [MGF trick](#) naturally.

Lemma 2.4.1 (Tail decay for sub-exponential random variable). Let $X \in \text{SubExp}(\sigma^2, \alpha)$ with mean μ . Then

$$\mathbb{P}(X - \mu \geq t) \leq \begin{cases} e^{-\frac{t^2}{2\sigma^2}}, & \text{if } 0 \leq t \leq \frac{\sigma^2}{\alpha}; \\ e^{-\frac{t}{2\alpha}}, & \text{if } t > \frac{\sigma^2}{\alpha}. \end{cases}$$

Proof. We see that

$$\mathbb{P}(X - \mu \geq t) \leq \inf_{0 \leq \lambda < 1/\alpha} \frac{\mathbb{E} \left[e^{\lambda(X-\mu)} \right]}{e^{\lambda t}} \leq \inf_{0 \leq \lambda < 1/\alpha} e^{\frac{\lambda^2 \sigma^2}{2} - \lambda t}.$$

Now, we just need to minimize the exponent, which is a convex quadratic function, in the range $(0, \frac{1}{\alpha})$. The infimum depends on the value of α :

- $\frac{t}{\sigma^2} < \frac{1}{\alpha}$: we get the Gaussian bound.
- $\frac{t}{\sigma^2} \geq \frac{1}{\alpha}$: the minimizer is $1/\alpha$, and we get the exponential bound.

■

Corollary 2.4.1. Let $X \in \text{SubExp}(\sigma^2, \alpha)$ with mean μ . Then

$$\mathbb{P}(|X - \mu| \geq t) \leq 2 \exp\left(-\frac{t^2}{2(\sigma^2 + t\alpha)}\right)$$

for all $t \geq 0$.

Proof. We see that

$$\mathbb{P}(|X - \mu| \geq t) \leq 2 \exp\left(-\min\left\{\frac{t^2}{2\sigma^2}, \frac{t}{2\alpha}\right\}\right) \leq 2 \exp\left(-\frac{t^2}{2(\sigma^2 + t\alpha)}\right)$$

by observing $\min(1/u, 1/v) \geq 1/(u+v)$. ■

Just like [Lemma 2.3.3](#) for [sub-Gaussian](#) random variables, [sub-exponential](#) random variables are also closed under convolution.

Lemma 2.4.2 (Closed under convolution). Let $X_i \in \text{SubExp}(\sigma_i^2, \alpha_i)$ be all independent with mean μ_i , then

$$\sum_i (X_i - \mu_i) \in \text{SubExp}\left(\sum_i \sigma_i^2, \|\alpha\|_\infty\right).$$

Proof. Since

$$\mathbb{E}\left[e^{\lambda \sum_i (X_i - \mu_i)}\right] = \prod_{i=1}^n \mathbb{E}\left[e^{\lambda (X_i - \mu_i)}\right] \leq \prod_{i=1}^n e^{\lambda^2 \sigma_i^2 / 2} = e^{\lambda^2 \sum_i \sigma_i^2 / 2}$$

where the inequality holds if $|\lambda| < 1/\alpha_i$ for all i , i.e., $|\lambda| < 1/\|\alpha\|_\infty$. ■

2.4.2 Bernstein's Inequality

We are now ready to state the generalization of [Hoeffding's inequality](#) to sums of independent [sub-exponential](#) random variables.

Theorem 2.4.1 (Bernstein's inequality for sub-exponential random variables). Let $X_i \sim \text{SubExp}(\sigma_i^2, \alpha_i)$ be all independent with mean μ_i , then

$$\mathbb{P}\left(\left|\sum_{i=1}^n (X_i - \mu_i)\right| \geq t\right) \leq 2 \exp\left(-\min\left\{\frac{t^2}{2 \sum_i \sigma_i^2}, \frac{t}{2 \|\alpha\|_\infty}\right\}\right).$$

Proof. This is immediate from [Lemma 2.4.1](#) and [Lemma 2.4.2](#). ■

We can restate [Bernstein's inequality](#) in a convenient way.

Corollary 2.4.2. Let $X_i \sim \text{SubExp}(\sigma_i^2, \alpha_i)$ be all independent with mean μ_i , and let $k \geq \sigma_i, \alpha_i$ for all i . Then for all $a_i \in \mathbb{R}$, we have

$$\mathbb{P}\left(\left|\sum_{i=1}^n a_i (X_i - \mu_i)\right| \geq t\right) \leq 2 \exp\left(-\min\left\{\frac{t^2}{k^2 \|a\|^2}, \frac{t}{k \|a\|_\infty}\right\}\right).$$

Note. If we let $a_i = 1/\sqrt{n}$, we obtain an absolute constant c (depending on k only)

$$\mathbb{P}\left(\left|\frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu_i)\right| \geq t\right) \leq \begin{cases} 2e^{-ct^2}, & \text{if } 0 < t < c\sqrt{n}; \\ 2e^{-t\sqrt{n}}, & \text{if } t > c\sqrt{n}. \end{cases}$$

Remark. Bernstein's inequality gives the [sub-Gaussian](#) tail decay expected from CLT for most t . Only in the very rare event regime, does the slower exponential tail decay come in.

Lecture 4: McDiarmid's Inequality

2.5 Bounded Difference Concentration Inequality

28 Aug. 9:00

2.5.1 Applications of Bernstein's Inequality to Bounded Random Variables

Now we see some applications of [Bernstein's inequality](#), addressing weaknesses of [Hoeffding's inequality](#).

Lemma 2.5.1. Let $|X - \mu| \leq b$ and $X - \mu$ is Subg(b^2). It's also true that $X - \mu \in \text{SubExp}(2\sigma^2, 2b)$ where $\text{Var}[X] = \sigma^2$.

Proof. From $(X - \mu)^k \leq (X - \mu)^2 |X - \mu|^{k-2} \leq (X - \mu)^2 b^{k-2}$, we have

$$\mathbb{E} \left[e^{\lambda(X-\mu)} \right] = 1 + \frac{\lambda^2}{2} \sigma^2 + \sum_{k=3}^{\infty} \lambda^k \frac{\mathbb{E}[X - \mu]^k}{k!} \leq 1 + \frac{\lambda^2 \sigma^2}{2} + \frac{\lambda \sigma^2}{2} \sum_{k=3}^{\infty} (|\lambda|b)^{k-2}.$$

The last sum is a geometric series, which converges if $|\lambda| < 1/b$ to

$$1 + \frac{\lambda^2 \sigma^2}{2} \left(\frac{1}{1 - b|\lambda|} \right).$$

Then from $1 + x \leq e^x$, we see that for $|\lambda| < 1/2b$,

$$\mathbb{E} \left[e^{\lambda(X-\mu)} \right] \leq e^{\frac{\lambda^2 \sigma^2}{2(1-b|\lambda|)}} \leq e^{\lambda^2 \sigma^2}.$$

■

From this, by directly apply [Bernstein's inequality](#), we have the following.

Corollary 2.5.1. Let X be a random variable such that $|X - \mu| \leq b$. For any $t > 0$,

$$\mathbb{P}(|X - \mu| \geq t) \leq 2 \exp \left(\frac{-t^2}{2(2\sigma^2 + t \cdot 2b)} \right).$$

Furthermore, let X_1, \dots, X_n be independent random variables with $\mathbb{E}[X_i] = \mu_i$ and $\text{Var}[X_i] = \sigma_i^2$ such that $|X_i - \mu_i| \leq b$ for all i . Then for any $t > 0$,

$$\mathbb{P} \left(\left| \sum_{i=1}^n (X_i - \mu_i) \right| \geq t \right) \leq 2 \exp \left(\frac{-t^2}{4(\sum_i \sigma_i^2 + tb)} \right).$$

In particular, if $\mu_i = \mu$ for all i , then

$$\Pr \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| \geq t \right) \leq 2 \exp \left(-\frac{nt^2}{4(\sigma^2 + tb)} \right).$$

Remark. Observe that in the last line of the proof of [Lemma 2.5.1](#), the inequality is quite loose. This means that we can explicitly maximize the quantity in the exponent over $|\lambda| \in (0, 1/2b)$ to get a higher bound and hence, a better variance factor. This leads to a tighter version of [Corollary 2.5.1](#).

Corollary 2.5.2. Let X_1, \dots, X_n be independent random variables with $\mathbb{E}[X_i] = \mu$ and $\text{Var}[X_i] = \sigma^2$ such that $|X_i - \mu| \leq b$ for all i . Then for any $t > 0$,

$$\mathbb{P} \left(\left| \sum_{i=1}^n X_i - \mu \right| \geq t \right) \leq 2 \exp \left(\frac{-t^2/2}{n\sigma^2 + bt/3} \right).$$

In particular,

$$\mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| \geq t \right) \leq 2 \exp \left(\frac{-nt^2/2}{\sigma^2 + bt/3} \right).$$

From [Corollary 2.5.2](#):

- if $t \leq 3\sigma^2/b$, the tail of the sample mean behaves like a [sub-Gaussian](#) tail;
- if $t > 3\sigma^2/b$, the tail of the sample mean behaves like a [sub-exponential](#) tail.

Remark. In practice, since we know that sample mean is \sqrt{n} -consistent, we generally look at a sequence of quantiles of the sample mean that is of $O(n^{-1/2})$. Therefore, the tail behavior when t gets large, is practically irrelevant.

By choosing the appropriate t in the above tail bound, we can get the following confidence interval for μ .

Corollary 2.5.3. Under the assumption of [Corollary 2.5.2](#),

$$\mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| \leq \frac{\sigma}{\sqrt{n}} \sqrt{2 \log \frac{2}{\alpha}} + \frac{3b}{3n} \log \frac{2}{\alpha} \right) \geq 1 - \alpha$$

Proof. Let

$$\alpha = 2 \exp \left(\frac{-t^2}{2(V + bt/3)} \right),$$

then

$$t^2 - \frac{2tb}{3} \log \frac{2}{\alpha} - 2V \log \frac{2}{\alpha} = 0.$$

■

In [Corollary 2.5.3](#), we have an $O(1/\sqrt{n})$ term, which is similar to the [one](#) derived from [Hoeffding's inequality](#) for bounded random variables. In contrary to the Hoeffding's bound, we have an additional lower order term here.

Remark. Observe that the higher order term in [Corollary 2.5.3](#) involves the variance, whereas in the case of [Hoeffding](#), it involves the range. Therefore, for random variables with large range but highly concentrated around its mean, the [Hoeffding confidence interval](#) would be much wider.

The above remark is demonstrated by the following example.

Example. Let $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(p)$. Suppose we observe $X_i = 0$ for all i , then $\hat{p} = \bar{X} = 0$ and the estimate of $\text{Var}[X_1]$ would be $\hat{p}(1 - \hat{p}) = 0$.

Hence, if we plug this estimate of variance into the [confidence bound from Bernstein](#), the length of which would be $O(1/n)$. However, in the case of [Hoeffding](#) (which works with the range, in this case, 1), the length would be $O(1/\sqrt{n})$.

2.5.2 McDiarmid's Inequality

Now we go back to the discussion about [empirical process](#). We do the first step, i.e., we want to show

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

“concentrates” when \mathcal{F} is bounded provided that

$$\sup_{x \in \mathcal{X}, f \in \mathcal{F}} |f(x)| \leq B.$$

One simple example of bounded function class arises in the task of classification.

Example (Classification). Consider $f(x)$ corresponds to the class label of an observation with feature value x , then the class is bounded.

However, since S_n falls neither into the category of [Hoeffding](#) nor [Bernstein](#), we would need a more general concentration inequalities: the [McDiarmid's inequality](#).²

Theorem 2.5.1 (McDiarmid's inequality). Let X_1, \dots, X_n be i.i.d. random variables on χ , and let $f: \chi^n \rightarrow \mathbb{R}$ satisfying the *bounded difference property*, i.e.,

$$\sup_{x_1, \dots, x_n, x'_i} |f(x_1, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| \leq c_i$$

for all i . Then for any $t > 0$,

$$\mathbb{P}(f(X_1, \dots, X_n) - \mathbb{E}[f(X_1, \dots, X_n)] \geq t) \leq \exp\left(\frac{-2t^2}{\sum_i c_i^2}\right).$$

The same bound holds for the left tail.

Remark. The qualitative statement for [McDiarmid's inequality](#) is that “a random variable that depends on the influence of many independent random variables but not too many on any one of them concentrates”.

Proof. Typically, $\sum_i c_i = O(1)$ concentration will happen if $\sum_i c_i^2 = o(1)$. For example, if each $c_i = O(1/n)$, then concentration happens but not when all $c_i = 0$ except one of them is 1. \ast

Remark. [McDiarmid's inequality](#) is a generalization of [Hoeffding's inequality](#).

Proof. Let

$$f(x_1, \dots, x_n) = \frac{1}{n}(x_1 + \dots + x_n).$$

When X_i 's are independent and $X_i \in [a_i, b_i]$ for all i , it's easy to observe that when we change the i^{th} argument of f , the value of f can change at most by $(b_i - a_i)/n$, i.e., [McDiarmid's inequality](#) is satisfied with $c_i := (b_i - a_i)/n$, plugging in, we get back [Hoeffding's inequality](#). \ast

With [McDiarmid's inequality](#), we can check that the following holds for bounded function classes \mathcal{F} :

$$|S_n(x_1, \dots, x_n) - S_n(x_1, \dots, x'_i, \dots, x_n)| \leq \frac{2B}{n} =: c_i.$$

Then from [McDiarmid's inequality](#), for any $t > 0$,

$$\mathbb{P}(S_n \geq \mathbb{E}[S_n] + t) \leq \exp\left(\frac{-nt^2}{2B^2}\right) =: \delta,$$

or equivalently, $S_n \leq \mathbb{E}[S_n] + B\sqrt{\frac{2}{n} \log \frac{1}{\delta}}$ with probability at least $1 - \delta$.

Note. $B\sqrt{\frac{2}{n} \log \frac{1}{\delta}}$ is a lower order term, i.e., $\mathbb{E}[S_n]$ dominates it.

Proof. Since^a

$$O(B) \geq \mathbb{E}[S_n] \geq \mathbb{E}\left[\left|\frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}[f(X)]\right|\right] = O\left(\sqrt{\frac{\text{Var}[f(X_1)]}{n}}\right) \approx O\left(\frac{1}{\sqrt{n}}\right).$$

\ast

^aThis upper bound is pretty weak, and we will eventually work on getting better bounds.

All these imply that *it's enough to bound* $\mathbb{E}[S_n]$.

²It's also known as the *bounded difference inequality*.

Lecture 5: Proof of McDiarmid's Inequality

We should note that the usual proof of [McDiarmid inequality](#) involves [martingale decomposition](#) and [Azuma-Hoeffding inequality](#), a generalization of [Hoeffding's inequality](#) for [martingale difference sequence](#).

1 Sep. 9:00

Definition 2.5.1 (Martingale difference sequence). A *martingale difference sequence* is a sequence of random variables Δ_1, \dots such that $\mathbb{E}[\Delta_i \mid \Delta_{1:i-1}] = 0$ for all i .

However, we will not go with this route; instead, we prove something weaker but trickier.³

Note. The condition $\sup_{x_1, \dots, x_n, x'_i} |f(x_1, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| \leq c_i$ is equivalent to

$$|f(x_1, \dots, x_n) - f(z_1, \dots, z_n)| \leq \sum_{i=1}^n c_i \mathbb{1}_{x_i \neq z_i}.$$

Now, we need one last lemma to prove [McDiarmid inequality](#).

Lemma 2.5.2. For all $x \neq y \in \mathbb{R}$,

$$\frac{e^x - e^y}{x - y} \leq \frac{e^x + e^y}{2} \Rightarrow |e^x - e^y| \leq |x - y| \left(\frac{e^x + e^y}{2} \right).$$

Proof. Since

$$\frac{e^x - e^y}{x - y} = \int_0^1 e^{sx + (1-s)y} ds = \frac{1}{x - y} \int_x^y e^t dt$$

where we let $t = sx + (1-s)y$. On the other hand, due to convexity, we also have

$$\frac{e^x - e^y}{x - y} = \int_0^1 e^{sx + (1-s)y} ds \leq \int_0^1 s \cdot e^x + (1-s)e^y ds = \frac{e^x + e^y}{2}.$$

■

We're now ready.

Proof of Theorem 2.5.1. Firstly, we note that it's equivalent to show that $f(X_1, \dots, X_n) - \mathbb{E}[f] \in \text{Subg}(\sum_i c_i^2/4)$. Without loss of generality, let $\mathbb{E}[f] = 0$, and we want to show that

$$\mathbb{E} \left[e^{\lambda(f(X) - \mathbb{E}[f])} \right] \leq e^{\frac{\lambda^2 \sum_i c_i^2}{8}} \Leftrightarrow M(\lambda) = \mathbb{E} \left[e^{\lambda f(X)} \right] \leq \exp \left(\frac{\lambda^2 (\sum_i c_i^2)}{8} \right) \Leftrightarrow \log M(\lambda) \leq \frac{\lambda^2 \sum_i c_i^2}{8}.$$

Observe that since both sides of the inequality is 0 at $\lambda = 0$, it's enough to show

$$\frac{d \log M(\lambda)}{d\lambda} = \frac{M'(\lambda)}{M(\lambda)} \leq \lambda \cdot \frac{\sum_i c_i^2}{4}$$

Let $\mathbb{X} = (X_1, \dots, X_n)$, and $\mathbb{X}' \stackrel{\text{i.i.d.}}{\sim} \mathbb{X}$ be the i.i.d. copy of \mathbb{X} . Then define the following.

Notation. $\mathbb{X}^{(i)} := (X'_1, \dots, X'_i, X_{i+1}, \dots, X_n)$ and $\mathbb{X}^{[i]} := (X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n)$.

³In fact, what we're going to prove is not even a weaker version: we prove something weaker while we really need the original (stronger) statement to hold.

Note that this implies $\mathbb{X}^{(0)} = \mathbb{X}$ and $\mathbb{X}^{(n)} = \mathbb{X}'$. Then, we can show that

$$\begin{aligned} M'(\lambda) &= \mathbb{E} \left[f(\mathbb{X}) e^{\lambda f(\mathbb{X})} \right] && \text{As } \mathbb{E}[f] = 0 \text{ and } \mathbb{X}, \mathbb{X}' \text{ are independent} \\ &= \mathbb{E} \left[(f(\mathbb{X}) - f(\mathbb{X}')) e^{\lambda f(\mathbb{X})} \right] \\ &= \mathbb{E} \left[\sum_{i=1}^n (f(\mathbb{X}^{(i-1)}) - f(\mathbb{X}^{(i)})) \cdot e^{\lambda f(\mathbb{X})} \right] \end{aligned}$$

if i^{th} position of \mathbb{X} and \mathbb{X}' are swapped, then for the new data $\mathbb{X}^{(i-1)}$ and $\mathbb{X}^{(i)}$ will also be swapped,

$$\begin{aligned} &= \mathbb{E} \left[\frac{1}{2} \sum_{i=1}^n \left(f(\mathbb{X}^{(i-1)}) - f(\mathbb{X}^{(i)}) \right) \cdot \left(e^{\lambda f(\mathbb{X})} - e^{\lambda f(\mathbb{X}^{[i]})} \right) \right] \\ &\leq \mathbb{E} \left[\frac{\lambda}{2} \sum_{i=1}^n \left| f(\mathbb{X}^{(i-1)}) - f(\mathbb{X}^{(i)}) \right| \cdot \left| f(\mathbb{X}) - f(\mathbb{X}^{[i]}) \right| \cdot \left(\frac{e^{\lambda f(\mathbb{X})} + e^{\lambda f(\mathbb{X}^{[i]})}}{2} \right) \right] \\ &\quad \text{from Lemma 2.5.2} \\ &\leq \frac{\lambda}{2} \left(\sum_{i=1}^n c_i^2 \right) \cdot M(\lambda). \end{aligned}$$

■

We note the following.

Note. The above proof doesn't even show a weaker version of [McDiarmid's inequality](#).

Proof. While in the proof, we need to show

$$\frac{d \log M(\lambda)}{d\lambda} = \frac{M'(\lambda)}{M(\lambda)} \leq \lambda \cdot \frac{\sum_i c_i^2}{4},$$

we only show

$$\frac{d \log M(\lambda)}{d\lambda} = \frac{M'(\lambda)}{M(\lambda)} \leq \lambda \cdot \frac{\sum_i c_i^2}{2}.$$

⊛

2.5.3 Applications of McDiarmid's Inequality

U-Statistics

Let $g: \mathbb{R}^2 \rightarrow \mathbb{R}$ be a symmetric function, and let $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$. Consider

$$U(X) = \frac{1}{\binom{n}{2}} \sum_{j < k} g(X_j, X_k).$$

Here're some examples of g .

Example. $g(x, y) = (x - y)^2$.

Example. $g(x, y) = |x - y|$.

Example (Wilcoxon's ranksum test). $g(x, y) = \mathbb{1}_{x_1 + x_2 > 0}$.

We're interested to know about $\mathbb{E}[g(X_1, X_2)]$. Assume g is bounded by B , then

$$U(\mathbb{X}) - U(\mathbb{X}^{[k]}) \leq \frac{1}{\binom{n}{2}} (n-1) 2B \leq \frac{4B}{n},$$

implying

$$\mathbb{P}(U - \mathbb{E}[U] \geq t) \leq e^{-\frac{nt^2}{8b^2}}$$

from [McDiarmid's inequality](#) with $c_i := 2B$.

Beyond McDiarmid's Inequality

Let's see some more advanced inequalities. In many cases, we want variance to be small. While

$$\text{Var}[X_1 + \dots + X_n] \leq \sum_{i=1}^n \text{Var}[X_i],$$

to have an inequality for a non-linear function, we have the following.

Theorem 2.5.2 (Efron-Stein inequality). Let X_1, \dots, X_n be independent random variables, and X'_1, \dots, X'_n be i.i.d. copies of X_i 's. Then

$$\text{Var}[f(\mathbb{X})] \leq \frac{1}{2} \sum_{i=1}^n \mathbb{E}[(f(\mathbb{X}) - f(\mathbb{X}^{[i]}))^2].$$

Note. We see that since $\text{Var}[X] = \frac{1}{2} \mathbb{E}[(X - X')^2]$, by letting $f(X_1, \dots, X_n) = \sum_i X_i$, if f satisfies bounded condition, then $\text{Var}[f] \leq \frac{1}{2} \sum_i c_i^2$.

Now, recall that by using [McDiarmid's inequality](#), we can show that for $\mathcal{F} \ni f$ being B -bounded,

$$S_n \leq \mathbb{E}[S_n] + B \sqrt{\frac{2}{n} \log \frac{1}{\delta}}$$

with probability at least $1 - \delta$. However, what if the variance $\text{Var}[f(X)]$ is small, but the maximum spread (B) is very large? In this case, we would want to replace B in the inequality by $\text{Var}[f(X)]$.

Notation (Empirical process notation). Let $\mathbb{P}f = \mathbb{E}[f]$ and $\mathbb{P}_n f = \sum_i f(X_i)/n$.

This is achieved by the following, although it's much harder to prove [[BLM13](#), §12].

Theorem 2.5.3 (Talagrand's concentration inequality). Let \mathcal{F} is B -bounded, and $S_n = \sup_{f \in \mathcal{F}} |\mathbb{P}_n f - \mathbb{P}f|$. Then

$$S_n \leq c \cdot \mathbb{E}[S_n] + c \sqrt{\frac{\sup_{f \in \mathcal{F}} \text{Var}[f(X_1)]}{n} \log \frac{1}{\alpha}} + c \cdot \frac{B}{n} \log \frac{1}{\alpha}$$

with probability at least $1 - \alpha$.

Remark. We might encounter an explicit situation where [Talagrand's concentration](#) is more profitable to use than [bounded differences inequality](#) later in the course.

Chapter 3

Expected Supremum of Empirical Process

Lecture 6: A Glance at Statistical Learning Theory

3.1 Statistical Learning

6 Sep. 9:00

3.1.1 Goodness of Fit Testing

Let's first see another motivation on studying uniform law of large numbers, i.e., the *goodness of fit testing*. Given $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$, we want to distinguish between $H_0: \mathbb{P} = \mathbb{P}_0$ and $H_1: \mathbb{P} \neq \mathbb{P}_0$.

Many tests are possible. One approach could be the **Kolmogorov-Smirnov test**: assume F is the CDF of \mathbb{P}_0 , then consider the **Kolmogorov-Smirnov statistics**:

Definition 3.1.1 (Kolmogorov-Smirnov statistics). The *Kolmogorov-Smirnov statistics* for a distribution \mathbb{P} is defined as

$$D_n = \sup_{t \in \mathbb{R}} |F_n(t) - F(t)|$$

where $F_n(t)$ and F is the **empirical CDF** and the CDF of \mathbb{P} , respectively.

From **Glivenko-Cantelli theorem**, $D_n \rightarrow 0$ under H_0 , and D_n should not converge to 0, under some alternative. Assuming continuity of F , Kolmogorov showed that

- (a) the distribution D_n does not depend on F ;
- (b) $D_n = O_p(1/\sqrt{n})$;
- (c) $\sqrt{n}D_n \rightarrow \sup_{t \in [0,1]} |B(t)|$ where $B(t)$ is the **Broweian bridge** on $[0, 1]$.
- (d) $\mathbb{P}(\sqrt{n}D_n \leq 2.4) \approx 0.999973$.

We'll take a non-asymptotic approach to this problem, i.e., we may not get such sharp constants.

3.1.2 Empirical Risk Minimization

Consider the following problem.

Problem 3.1.1 (Empirical risk minimization). Let $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ be n i.i.d. copies of $(X, Y) \in \mathcal{X} \times \mathcal{Y} \subseteq \mathbb{R}^d \times \mathbb{R}$ with distribution $\mathbb{P} = \mathbb{P}_X \times \mathbb{P}_{Y|X}$. Given a loss function $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ and a function class $\mathcal{F} = \{f: \mathcal{X} \rightarrow \mathcal{Y}\}$, the *empirical risk minimization* is

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i).$$

Example. \mathcal{F} can be the set of neural networks, decision trees, linear functions.

Example (Linear regression). Consider $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \mathbb{R}$, with $\mathcal{F} = \{x \rightarrow w^\top x : w \in \mathbb{R}^d\}$ and $\ell(a, b) = (a - b)^2$.

Example (Linear classification). Consider $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \{0, 1\}$, with $\mathcal{F} = \{x \rightarrow (\text{sgn}(w^\top x) + 1)/2 : w \in B_2^d\}$ where B_2^d is the unit ball in d -dimension, and $\ell(a, b) = \mathbb{1}_{a \neq b}$.

We also define the following.

Definition. Consider the set-up of [empirical risk minimization](#).

Definition 3.1.2 (Expected loss). The *expected loss*^a of $f \in \mathcal{F}$ is defined as

$$L(f) = \mathbb{E}_{(X,Y) \sim \mathbb{P}} [\ell(f(X), Y)].$$

^aAlso called *population loss* and *test error*.

Definition 3.1.3 (Empirical loss). The *empirical loss* is defined as

$$\hat{L}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i).$$

The main question in statistical learning is that, what is an upper-bound on the [expected loss](#) of [ERM](#)? If we plug in \hat{f} instead of f , this is asking the [test error](#) of \hat{f} .

To be specific, \hat{f} is basically a function of training data S , but when we look at

$$L(\hat{f}) = \mathbb{E}_{(X,Y)} [\ell(\hat{f}(x), Y)],$$

it is the expectation of future data points, i.e., it becomes a random variable, which is a function of S .

Lemma 3.1.1. For any \mathcal{F} , the [ERM](#) \hat{f} satisfies

$$\mathbb{E}[L(\hat{f})] - \inf_{f \in \mathcal{F}} L(f) \leq \mathbb{E} \left[\sup_{f \in \mathcal{F}} (L(f) - \hat{L}(f)) \right].$$

Proof. Let $f^* = \inf_{f \in \mathcal{F}} L(f)$. Then

$$L(\hat{f}) - L(f^*) = [L(\hat{f}) - \hat{L}(\hat{f})] + [\hat{L}(\hat{f}) - \hat{L}(f^*)] + [\hat{L}(f^*) - L(f^*)].$$

We see that

- $\hat{L}(\hat{f}) - \hat{L}(f^*) \leq 0$ by [definition](#);
- $\hat{L}(f^*) - L(f^*) = 0$ in expectation since f^* is fixed,
- We can't say $\mathbb{E}[L(\hat{f}) - \hat{L}(\hat{f})] = 0$ since \hat{f} is also random.

Combine all these, we have

$$\mathbb{E}[L(\hat{f})] - \inf_{f \in \mathcal{F}} L(f) = \mathbb{E}[L(\hat{f}) - L(f^*)] \leq \mathbb{E}[L(\hat{f}) - \hat{L}(\hat{f})] \leq \mathbb{E} \left[\sup_{f \in \mathcal{F}} (L(f) - \hat{L}(f)) \right].$$

■

Note. Let us decode what [Lemma 3.1.1](#) is claiming.

- Since $L(f)$ is the [population error](#) of f and $\hat{L}(f)$ is the [empirical loss](#) of f , $\sup_{f \in \mathcal{F}} (L(f) - \hat{L}(f))$ is the supremum of an [empirical process](#).
- For the left-hand side, it represents the [expected loss](#) of \hat{f} and the best possible out-of-sample error.^a This is often called the [excess risk](#).

^aOr the best possible prediction error of \mathcal{F} .

Notation (Excess risk). $\mathbb{E}[L(\hat{f})] - \inf_{f \in \mathcal{F}} L(f)$ is often called the *excess risk* of an [ERM](#).

Remark. For “curved” loss function like square loss, supremum can be further “localized”.

Remark. The bound in [Lemma 3.1.1](#) can be vacuumed for now, e.g., for linear regression.

Example (1-D classification with thresholds). Let $\ell(a, b) = \mathbb{1}_{a \neq b} = a + (1 - 2a)b$ for $a, b \in \{0, 1\}$. Then consider $a = y$ and $b = f(x)$,

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} (L(f) - \hat{L}(f)) \right] = \mathbb{E} \left[\sup_{f \in \mathcal{F}} \left(\mathbb{E} [Y + (1 - 2Y)f(X)] - \frac{1}{n} \sum_{i=1}^n (y_i + (1 - 2y_i)f(x_i)) \right) \right],$$

which can be viewed essentially as^a the [empirical process](#) on the function f instead of ℓ ,

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

For 1-D case, assume that $\mathcal{F} = \{x \mapsto \mathbb{1}_{x \leq \theta} : \theta \in \mathbb{R}\}$, then

$$\mathbb{E} \left[\sup_{\theta \in \mathbb{R}} \left(\mathbb{P}(X \leq \theta) - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \theta} \right) \right] = \mathbb{E} \left[\sup_{\theta \in \mathbb{R}} (F(\theta) - F_n(\theta)) \right],$$

i.e., $P(X \leq \theta)$ is the CDF of the marginal distribution of X , $F(\theta)$, and $\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \theta}$ is the [empirical CDF](#) $F_n(\theta)$. Therefore, we go back to the same problem we introduced in the beginning of the chapter, i.e., the [Kolmogorov-Smirnov statistics](#).

Let the term $\mathbb{P}(X \leq \theta) - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \theta}$ to be a random variable U_θ . One problem here is, we have infinitely many random variables, and they are also correlated with each other quite a lot. So how does this supremum behave?

Since each U_θ is at most 1, for any θ , i.e., $\sup U_\theta \leq 1$. So the worst case here is 1, and probably the best case is $O(1/\sqrt{n})$.

^aSince $Y - \sum_i y_i/n$ is independent of f , so let's drop it; and $1 - 2Y$ is the sign, so can be dropped essentially.

Lecture 7: Bracketing and Symmetrization

Our main [empirical process](#) is so far $\mathbb{E} [\sup_{f \in \mathcal{F}} \mathbb{P}_n f - \mathbb{P} f]$. Let's first focus on the [1-D thresholds classification](#), i.e., we want to bound the supremum

$$\mathbb{E} \left[\sup_{\theta \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \theta} - \mathbb{P}(X \leq \theta) \right| \right].$$

There are 2 approaches to bound this supremum: bracketing and symmetrization.

3.1.3 Bracketing

The main idea of bracketing is the following.

Intuition. Reduce an infinite number of random variables to finite, which will be more manageable.

Assume that \mathbb{P} is continuous, and consider a finite set $\{\theta_i\}_{i=0}^{N+1}$ with $\theta_0 = -\infty$, $\theta_{N+1} = \infty$, such that they correspond to quantile of \mathbb{P} , i.e.,

$$\mathbb{P}(\theta_i \leq X \leq \theta_{i+1}) = \frac{1}{N+1}.$$

Given a θ , X will lie in between two adjacent θ_i 's in the sequence. Denote the upper-bound as $u(\theta)$ and the lower-bound as $\ell(\theta)$ for this θ , then

$$\begin{aligned} \mathbb{P}(X \leq \theta) - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \theta} &\leq \mathbb{P}(X \leq u(\theta)) - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \ell(\theta)} \\ &\leq \mathbb{E} [\mathbb{1}_{X \leq u(\theta)}] - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \ell(\theta)} \\ &\leq \mathbb{E} [\mathbb{1}_{X \leq \ell(\theta)}] - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \ell(\theta)} + \mathbb{P}(\ell(\theta) \leq X \leq u(\theta)) \\ &\leq \mathbb{E} [\mathbb{1}_{X \leq \ell(\theta)}] - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \ell(\theta)} + \frac{1}{N+1} \end{aligned}$$

if we take the supremum over $\ell(\theta) \in \mathbb{R}$ instead of θ ,

$$\leq \frac{1}{N+1} + \mathbb{E} \left[\max_{0 \leq j \leq N} \mathbb{E} [\mathbb{1}_{X \leq \theta_j}] - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \theta_j} \right]. \quad (3.1)$$

To further bound Equation 3.1, recall the following.

As previously seen. If $X_i \sim \text{Subg}(\sigma^2)$ independent, $\sum_i a_i X_i \sim \text{Subg}((\sum_i a_i^2) \sigma^2)$ from Lemma 2.3.3.

Remark. Let $a_i = 1/n$, we see that $\mathbb{E} [\mathbb{1}_{X \leq \theta_j}] - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \theta_j} \in \text{Subg}(1/n)$.^a

^aSince it's bounded between 0 and 1.

Finally, recall what we have proved in the homework.

Lemma 3.1.2. Let $X_1, \dots, X_n \sim \text{Subg}(\sigma^2)$,^a then $\mathbb{E} [\max_i X_i] \leq \sqrt{2\sigma^2 \log n}$.

^aNot necessary independent.

Then, we can show the final bound.

Proposition 3.1.1 (Bracketing). Let $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$, and $\mathcal{F} = \{\mathbb{1}_{X \leq \theta} : \theta \in \mathbb{R}\}$. Then

$$\mathbb{E}_X \left[\sup_{f \in \mathcal{F}} \left(\mathbb{P}(X \leq \theta) - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \theta} \right) \right] = O \left(\sqrt{\frac{\log n}{n}} \right).$$

Proof. From Lemma 3.1.2, since we have $(N+1)$ random variables with variance factor $1/n$, by choosing $N+1 := n$,^a Equation 3.1 can be further bounded by

$$\sqrt{\frac{2 \log(N+1)}{n}} + \frac{1}{N+1} = O \left(\sqrt{\frac{\log n}{n}} \right).$$

■

^aRecall that n is the sample size, so we can choose the corresponding n to meet the requirement.

3.1.4 Symmetrization

Another technique called symmetrization, which is essentially stated in the following lemma.

Lemma 3.1.3 (Symmetrization). Given a function class $\mathcal{F} = \{f: \mathcal{X} \rightarrow \mathcal{Y}\}$ and $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$, and $\epsilon_1, \dots, \epsilon_n$ be i.i.d. [Rademacher random variables](#). Then

$$\max \left(\mathbb{E} \left[\sup_{f \in \mathcal{F}} \mathbb{P}_n f - \mathbb{P} f \right], \mathbb{E} \left[\sup_{f \in \mathcal{F}} \mathbb{P} f - \mathbb{P}_n f \right] \right) \leq 2 \mathbb{E}_{\epsilon_i, X_i} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right].$$

In particular,

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{P}_n f - \mathbb{P} f| \right] \leq 2 \mathbb{E}_{\epsilon_i, X_i} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right].$$

Proof. Let X'_i 's be i.i.d. copies of X_i 's for all i . Since adding a sign ϵ_i won't change the expectation,^a

$$\begin{aligned} \mathbb{E} \left[\sup_{f \in \mathcal{F}} \mathbb{E} [f(X)] - \frac{1}{n} \sum_{i=1}^n f(X_i) \right] &= \mathbb{E} \left[\sup_{f \in \mathcal{F}} \mathbb{E}_{X'_i} \left[\frac{1}{n} \sum_{i=1}^n f(X'_i) - \frac{1}{n} \sum_{i=1}^n f(X_i) \right] \right] \\ &\leq \mathbb{E}_{X_i} \left[\mathbb{E}_{X'_i} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(X'_i) - f(X_i)) \right] \right] \\ &= \mathbb{E}_{X_i, X'_i, \epsilon_i} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(X'_i) - f(X_i)) \epsilon_i \right] \\ &\leq \mathbb{E}_{X'_i, \epsilon_i} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f(X'_i) \epsilon_i \right] + \mathbb{E}_{X_i, \epsilon_i} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f(X_i) \epsilon_i \right] \\ &= 2 \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right]. \end{aligned}$$

■

^aSince the distributions of $f(X'_i) - \sum_i f(X_i)$ and $f(X_i) - \sum_i f(X'_i)$ are the same.

Intuition. If we condition on X_i 's, the bound can be seen as linear combination of [Rademacher random variables](#). Thus, we can refer to properties of [sub-Gaussian](#) random variables.

The upper-bound deserves a special name.

Definition 3.1.4 (Rademacher complexity). Let $X_i \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$ be independent and ϵ_i be i.i.d. [Rademacher random variables](#). The *Rademacher complexity* of a function class \mathcal{F} w.r.t. \mathbb{P} is

$$R_n(\mathcal{F}) := \mathbb{E}_{\epsilon_i, X_i} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right].$$

On the other hand, the opposite direction of [symmetrization lemma](#) also holds.

Lemma 3.1.4. Given a function class $\mathcal{F} = \{f: \mathcal{X} \rightarrow \mathcal{Y}\}$ and $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$, and $\epsilon_1, \dots, \epsilon_n$ be i.i.d. [Rademacher random variables](#). Then

$$\mathbb{E}_{X_i, \epsilon_i} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right] \leq 2 \mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{P}_n f - \mathbb{P} f| \right] + \frac{1}{\sqrt{n}} \sup_{f \in \mathcal{F}} |\mathbb{P} f|.$$

Proof. This technique is so-called *desymmetrization*: Consider

$$\begin{aligned} & \mathbb{E}_{\epsilon_i, X_i} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right] \\ & \leq \mathbb{E}_{\epsilon_i, X_i} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(X_i) - \mathbb{E}[f(X)]) \right| \right] + \mathbb{E}_{\epsilon_i} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbb{E}[f(X)] \right| \right] \\ & = \mathbb{E}_{\epsilon_i, X_i, X'_i} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(X_i) - \mathbb{E}[f(X'_i)]) \right| \right] + \mathbb{E} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbb{E}_{\epsilon_i}[f(X_i)] \right| \right]. \end{aligned}$$

The first term can be further bounded by

$$\begin{aligned} \mathbb{E}_{\epsilon_i, X_i, X'_i} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(X_i) - \mathbb{E}[f(X'_i)]) \right| \right] & \leq \mathbb{E}_{\epsilon_i, X_i, X'_i} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(X_i) - f(X'_i)) \right| \right] \\ & = \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(X_i) - f(X'_i)) \right] \\ & = \mathbb{E} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - f(X'_i) + (\mathbb{E}[f] - \mathbb{E}[f])) \right| \right] \\ & = 2 \mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{P}_n f - \mathbb{P} f| \right], \end{aligned}$$

and the second term can be bounded by

$$\mathbb{E}_{\epsilon_i} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbb{E}[f(X)] \right| \right] \leq \sup_{f \in \mathcal{F}} |\mathbb{E}[f(X)]| \cdot \mathbb{E} \left[\left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \right| \right] \leq \frac{1}{\sqrt{n}} \sup_{f \in \mathcal{F}} |\mathbb{P} f|$$

where $\mathbb{E} \left[\left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \right| \right] \leq \frac{c}{\sqrt{n}}$ with $c = 1$. Combine them together, we have the final result. \blacksquare

Lecture 8: Symmetrization on 1-D Threshold Classification

Analogous to the [Rademacher complexity](#) defined for a function class w.r.t. \mathbb{P} , we can define it on a set. 11 Sep. 9:00

Definition 3.1.5 (Rademacher width). Let ϵ_i be i.i.d. [Rademacher random variables](#). Then the *Rademacher width*^a of a set $A \subseteq \mathbb{R}^n$ is defined as

$$R_n(A) = \mathbb{E}_{\epsilon_i} \left[\sup_{a \in A} \frac{1}{n} \sum_{i=1}^n \epsilon_i a_i \right].$$

^aAlso called *Rademacher average*.

Notation. People sometimes just say “Rademacher complexity” for [Rademacher width](#).

Now, applying the [symmetrization lemma](#) to $\mathcal{F} = \{\mathbb{1}_{X \leq \theta} : \theta \in \mathbb{R}\}$, we have the following result that is comparable to [Proposition 3.1.1](#).

Proposition 3.1.2. Let $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$, and $\mathcal{F} = \{\mathbb{1}_{x \leq \theta} : \theta \in \mathbb{R}\}$. Then

$$\mathbb{E}_X \left[\sup_{f \in \mathcal{F}} \left(\mathbb{P}(X \leq \theta) - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \theta} \right) \right] = O \left(\sqrt{\frac{\log n}{n}} \right).$$

Proof. From the [symmetrization lemma](#),

$$\begin{aligned} \mathbb{E}_{X, x_i} \left[\sup_{\theta \in \mathbb{R}} \left(\mathbb{P}(X \leq \theta) - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq \theta} \right) \right] &\leq 2 \mathbb{E}_{\epsilon_i, x_i} \left[\sup_{\theta \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbb{1}_{x_i \leq \theta} \right] \quad \text{condition on } x_1, \dots, x_n \\ &= 2 \mathbb{E}_{x_i} \left[\mathbb{E}_{\epsilon_i | x_i} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbb{1}_{x_i \leq \theta} \middle| x_1, \dots, x_n \right] \right]. \end{aligned}$$

Let $V_\theta := \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbb{1}_{x_i \leq \theta}$, we see that there are only $n+1$ distinct V_θ 's, and it's constant in the intervals $\theta \in [X_{(k)}, X_{(k+1)})$ for $k = 0, \dots, n-1$ where $X_{(k)}$ are the order statistics with $X_{(0)} := -\infty$. Now, define $\theta_k := X_{(k)}$, we can then write

$$\sup_{\theta \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbb{1}_{x_i \leq \theta} = \max_{k=0, \dots, n} \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbb{1}_{x_i \leq \theta_k},$$

hence,

$$2 \mathbb{E}_{x_i} \left[\mathbb{E}_{\epsilon_i | x_i} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbb{1}_{x_i \leq \theta} \middle| x_1, \dots, x_n \right] \right] = 2 \mathbb{E}_{x_i} \left[\mathbb{E}_{\epsilon_i | x_i} \left[\max_{k=0, \dots, n} V_{\theta_k} \middle| x_1, \dots, x_n \right] \right]$$

with $V_{\theta_k} \sim \text{Subg}(1/n)$ and [Lemma 3.1.2](#),

$$\leq 2 \mathbb{E}_{x_i} \left[\sqrt{\frac{2}{n} \log(n+1)} \right] = O \left(\sqrt{\frac{\log n}{n}} \right).$$

■

Remark. Looking back to the [example of 1-D thresholds classification](#), we see that the [excess risk](#) of [ERM](#) is $O(\sqrt{\log n/n})$.

3.2 Vapnik-Chervonenkis Dimension

3.2.1 Glivenko-Cantelli Class

From [bracketing](#) and [symmetrization](#), we see that there are classes of functions such that

$$\sup_{f \in \mathcal{F}} |\mathbb{P}f - \mathbb{P}_n f| \rightarrow 0$$

as $n \rightarrow \infty$. They deserve their own name.

Definition 3.2.1 (Glivenko-Cantelli). A function class $\mathcal{F} = \{f: \chi \rightarrow \mathbb{R}\}$ is *Glivenko-Cantelli* w.r.t. \mathbb{P} if as $n \rightarrow \infty$,

$$\sup_{f \in \mathcal{F}} |\mathbb{P}f - \mathbb{P}_n f| \rightarrow 0.$$

From [bracketing](#) and [symmetrization](#), we know that $\mathcal{F} = \{\mathbb{1}_{X \leq \theta} : \theta \in \mathbb{R}\}$ is [Glivenko-Cantelli](#). Let's see some counterexamples.

Example. Let $\chi = \mathbb{R}$, $\mathcal{F} = \{\mathbb{1}_A : A \subseteq \chi, |A| < \infty\}$, and \mathbb{P} be any continuous measure on χ . Then \mathcal{F} is not [Glivenko-Cantelli](#) w.r.t. \mathbb{P} .

Proof. For $f = \mathbb{1}_A$, $\mathbb{P}f = \mathbb{P}(X \in A) = 0$ since $|A| < \infty$. On the other hand, let $A_0 = \{X_1, \dots, X_n\}$ be the observed empirical data, $\mathbb{P}_n f = 1$, i.e., $\sup_{f \in \mathcal{F}} |\mathbb{P}f - \mathbb{P}_n f| = 1$ for all $n \in \mathbb{N}$. ⊗

Example. Let $\chi = \mathbb{R}$, $\mathcal{F} = \{f: \chi \rightarrow \mathbb{R} \text{ bounded and continuous}\}$, and $\mathbb{P} = \mathcal{U}[0, 1]$. Then \mathcal{F} is not [Glivenko-Cantelli](#).

Proof. Consider $f(X_i) = 1$ for $i = 1, \dots, n$ and $f = 0$ elsewhere (continuously),^a then we can make $\int_0^1 f(t) dt < \delta$ for some $\delta \in (0, 1)$. This implies $\sup_{f \in \mathcal{F}} |\mathbb{P}f - \mathbb{P}_n f| \geq 1 - \delta$ for all $n \in \mathbb{N}$. \circledast

^aE.g., sharp peak at X_i 's.

3.2.2 Vapnik-Chervonenkis Dimension

Notation. Let $\mathcal{F}(x_1, \dots, x_n) := \{(f(x_1), \dots, f(x_n))\}_{f \in \mathcal{F}} \subseteq \mathbb{R}^n$.

We can relate the [Rademacher width](#) of $\mathcal{F}(X_1, \dots, X_n)$ to the [Rademacher complexity](#) of \mathcal{F} since¹

$$\mathbb{E}_{X_i} [R_n(\mathcal{F}(X_1, \dots, X_n))] = \mathbb{E}_{X_i, \epsilon_i} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right] = R_n(\mathcal{F}).$$

Moreover, we see that if $\mathcal{F}(X_1, \dots, X_n)$ is finite, by the same proof as in [Proposition 3.1.2](#),

$$\mathbb{E}_{X_i} [R_n(\mathcal{F}(X_1, \dots, X_n))] \leq 2 \sqrt{\frac{2 \log |\mathcal{F}(X_1, \dots, X_n)|}{n}}.$$

The up-shot is the following.

Remark. If $|\mathcal{F}(X_1, \dots, X_n)| \leq n^d$ for some $d \in \mathbb{N}^+$, then we again get an $O(\sqrt{\log n/n})$ bound.

This is captured by the [polynomial discrimination](#), where we're going to focus on boolean functions.

Definition 3.2.2 (Polynomial discrimination). We say that a boolean function class \mathcal{F} on χ has a *polynomial discrimination* if for all $x_1, \dots, x_n \in \chi$, $|\mathcal{F}(x_1, \dots, x_n)| \leq \text{poly}(n)$.

To characterize $|\mathcal{F}(x_1, \dots, x_n)|$, we will look at the [VC dimension](#) of \mathcal{F} , which is related to the size of the discrimination of \mathcal{F} in a non-trivial way.

Definition. Let \mathcal{F} be a boolean function class on χ .

Definition 3.2.3 (Shatter). A finite set $\{x_1, \dots, x_D\} \subseteq \chi$ is *shattered* by \mathcal{F} if $\mathcal{F}(x_1, \dots, x_D) = \{0, 1\}^D$.^a

^aWe take the convention that \emptyset is always *shattered*.

Definition 3.2.4 (Vapnik-Chervonenkis dimension). The *VC dimension* of \mathcal{F} on χ is the maximum integer D such that there exists a size D finite set $A \subseteq \chi$ *shattered* by \mathcal{F} .

Let's consider some examples on $\chi = \mathbb{R}$.

Example. The [VC dimension](#) of $\mathcal{F} = \{\mathbb{1}_{X \leq \theta} : \theta \in \mathbb{R}\}$ is 1.

Example. The [VC dimension](#) of $\mathcal{F} = \{\mathbb{1}_{[a,b]} : a, b \in \mathbb{R}\}$ is 2.

Let's look at one example with $\chi = \mathbb{R}^2$.

Example. The [VC dimension](#) of $\mathcal{F} = \{\mathbb{1}_{[a,b] \times [c,d]} : a, b, c, d \in \mathbb{R}\}$ is 4.

Lecture 9: VC Dimension

Firstly, given [VC dimension](#), we can upper-bound the size of the discrimination.

¹This is why people overload R_n for both [Rademacher width](#) and [Rademacher complexity](#).

Lemma 3.2.1 (Sauer-Shelah lemma). Let \mathcal{F} be a boolean function class such that $\text{VC}(\mathcal{F}) = d$, then for every $\{x_1, \dots, x_n\} \subseteq \chi$ such that $n \geq d$,

$$|\mathcal{F}(x_1, \dots, x_n)| \leq \binom{n}{0} + \binom{n}{1} + \dots + \binom{n}{d} \leq \left(\frac{en}{d}\right)^d.$$

To prove [Sauer-Shelah lemma](#), we need [Pajor's lemma](#).

Lemma 3.2.2 (Pajor's lemma). Given a boolean function class \mathcal{F} on a finite set Ω , then

$$|\mathcal{F}| \leq |\{S \subseteq \Omega: S \text{ shattered by } \mathcal{F}\}|.$$

Proof. We prove this by induction on n . For $n = 1$ (base case), it holds trivially since

$$|\mathcal{F}| = 2 \leq |\{S \subseteq \Omega: S \text{ shattered by } \mathcal{F}\}|.$$

Assume the statement holds for all Ω such that $|\Omega| = n$. For $|\Omega| = n + 1$, we write $\Omega = (\Omega \setminus \{x_0\}) \cup \{x_0\} =: \Omega_0 \cup \{x_0\}$ and let \mathcal{F}_0 and \mathcal{F}_1 be two boolean function classes defined on Ω_0 as

$$\mathcal{F}_0 = \{f \in \mathcal{F}: f(x_0) = 0\}, \quad \mathcal{F}_1 = \{f \in \mathcal{F}: f(x_0) = 1\}.$$

We further define $S_{\mathcal{F}'}$ as $S_{\mathcal{F}'} = \{S \subseteq \Omega': S \text{ shattered by } \mathcal{F}'\}$ for any function class \mathcal{F}' defined on Ω' . Then, by induction hypothesis, $|\mathcal{F}_i| \leq |S_{\mathcal{F}_i}|$, hence $|\mathcal{F}| = |\mathcal{F}_0| + |\mathcal{F}_1| \leq |S_{\mathcal{F}_0}| + |S_{\mathcal{F}_1}|$. Finally, we claim the following.

Claim. $|S_{\mathcal{F}_0}| + |S_{\mathcal{F}_1}| \leq |S_{\mathcal{F}}|$.

Proof. Let $S \subseteq \Omega_0$ shattered by both \mathcal{F}_0 and \mathcal{F}_1 , then S is shattered by \mathcal{F} too. Moreover, Observe that $S \cup \{x_0\}$ is shattered by \mathcal{F} but not \mathcal{F}_i ($f(x_0)$ is fixed for $f \in \mathcal{F}_i$). Now, when

- S is shattered by only one of the \mathcal{F}_i 's: S contributes one unit both to $|S_{\mathcal{F}}|$ and $|S_{\mathcal{F}_i}|$;
- S is shattered by both \mathcal{F}_i 's, S and $S \cup \{x_0\}$ are shattered by \mathcal{F} : S contributes two units to $|S_{\mathcal{F}}|$ and one unit to both $|S_{\mathcal{F}_i}|$'s.

By counting, we're done (it's possible that S is shattered by \mathcal{F} but not \mathcal{F}_i 's, so \leq). ⊗

This implies $|\mathcal{F}| \leq |S_{\mathcal{F}}|$ for $|\Omega| = n + 1$, i.e., the induction is done. ■

We can then prove the [Sauer-Shelah lemma](#).

Proof of Lemma 3.2.1. Let Ω be a set of size n , then the number of subsets with size $\leq d$ is $\binom{n}{0} + \binom{n}{1} + \dots + \binom{n}{d}$, hence by the definition of [VC dimension](#),

$$|\{S \subseteq \Omega: S \text{ shattered by } \mathcal{F}\}| \leq \binom{n}{0} + \binom{n}{1} + \dots + \binom{n}{d}.$$

■

Then, as our motivation suggests, the same proof of [Proposition 3.1.2](#) applies, giving the following.

Proposition 3.2.1. For any function class \mathcal{F} , if $n \geq \text{VC}(\mathcal{F})$, for some constant c ,

$$R_n(\mathcal{F}) \leq c \sqrt{\frac{\text{VC}(\mathcal{F})}{n} \log \left(\frac{en}{\text{VC}(\mathcal{F})} \right)}.$$

Remark. We see that [Proposition 3.2.1](#) is independent of \mathbb{P} , i.e., the bounds still holds after taking $\sup_{\mathbb{P}}$ on the left-hand side. However, if $\text{VC}(\mathcal{F}) = \infty$, then this “distribution-free” uniform convergence fails.

However, if we don't care about distribution-free property, we do have examples that the uniform convergence holds for a particular \mathbb{P} when $\text{VC}(\mathcal{F}) = \infty$.

Example. For $\mathcal{F} = \{\mathbb{1}_A : \text{compact convex } A \subseteq [0, 1]^d\}$, $\text{VC}(\mathcal{F}) = \infty$. If \mathbb{P} is continuous w.r.t. Lebesgue's measure, then the uniform law of large number still holds.

Remark. The $\sqrt{\log n}$ factors in Proposition 3.2.1 is superfluous (Corollary 3.3.5).

Example. Let V be a vector space of real function on χ with $\dim(V) = D$, and $\mathcal{F} = \{\mathbb{1}_{f \geq 0} : f \in V\}$. Then $\text{VC}(\mathcal{F}) \leq D$.

Proof. We want to show that for any $\{x_1, \dots, x_{D+1}\}$ can't be shattered. Let

$$T = \{(f(x_1), \dots, f(x_{D+1})) : f \in V\},$$

which is a linear subspace of \mathbb{R}^{D+1} such that $\dim(T) \leq D$. Hence, there exists a non-zero $y \in \mathbb{R}^{D+1}$ such that $\sum_{i=1}^{D+1} y_i f(x_i) = 0$ for all $f \in V$. Now, without loss of generality, there exists an index k such that $y_k > 0$. If \mathcal{F} shatters $\{x_1, \dots, x_{D+1}\}$, then there exists $f \in V$ such that

$$\begin{cases} f(x_i) < 0, & \forall i: y_i > 0; \\ f(x_i) \geq 0, & \forall i: y_i \leq 0. \end{cases}$$

But then $\sum_i y_i f(x_i) < 0$, which is a contradiction. \circledast

Example (Half-space). For $\mathcal{F} = \{\mathbb{1}_H : \text{half space } H \subseteq \mathbb{R}^d\}$, $\text{VC}(\mathcal{F}) = d + 1$.

Although it seems like $\text{VC}(\mathcal{F}) \approx \#\text{parameters of } \mathcal{F}$; however, it's not true in general.

Example. Consider $\mathcal{F} = \{x \mapsto \mathbb{1}_{\sin tx \geq 0} : t \in \mathbb{R}^+\}$, then $\text{VC}(\mathcal{F}) = \infty$.

Lecture 10: Discretization of a Space

3.3 Metric Entropy Methods

15 Sep. 9:00

We have been focusing on boolean function class with finite VC dimension, and our goal now is to generalize beyond the boolean case. This can be done by discretizing of a space.

Intuition (Informal principle). We want to bound $\mathbb{E}[\sup_{t \in T} X_t]$. If $\{X_t\}_{t \in T}$ is “sufficiently continuous”, then $\mathbb{E}[\sup_{t \in T} X_t]$ is governed by metric properties of T (metric entropy!).

Definition 3.3.1 (Pseudo-metric). Given a space T , a function $d: T \times T \rightarrow \mathbb{R}^+$ is a *pseudo-metric* if

- (a) $d(x, x) = 0$ for all $x \in T$;^a
- (b) $d(x, y) = d(y, x)$ for all $x, y \in T$;
- (c) $d(x, y) \leq d(x, z) + d(y, z)$ for all $x, y, z \in T$.

^aIf d further satisfies that $d(x, y) > 0$ for all $x \neq y$, then it becomes a *metric*.

Note. The motivation of looking at pseudo-metric instead of the usual metric is because, consider observed data x_1, \dots, x_n at hands, the most natural distance might be

$$(f, g) \mapsto \sqrt{\frac{1}{n} \sum_{i=1}^n (f(x_i) - g(x_i))^2},$$

which is a **pseudo-metric** since f and g can agree only on x_i 's and vary elsewhere.

3.3.1 Covering Number and Packing Number

Now, let (T, d) denote a **pseudo-metric** space in the remaining of this section, unless specified.

Definition 3.3.2 (ϵ -net). A set N is an ϵ -net of (T, d) if for all $t \in T$, there exists $\pi(t) \in N$ such that $d(t, \pi(t)) \leq \epsilon$.

Definition 3.3.3 (Covering number). The ϵ -covering number $N(T, d, \epsilon)$ of (T, d) is defined as

$$N(T, d, \epsilon) := \inf\{|N| : N \text{ is an } \epsilon\text{-net for } (T, d)\}.$$

Remark. N is not necessary a subset of T for convenience. Furthermore, if $N \not\subseteq T$, one can construct another **net** $N' \subseteq T$ and N' is a **2 ϵ -net**.

Definition 3.3.4 (Totally bounded). (T, d) is *totally bounded* if for all $\epsilon > 0$, $N(T, d, \epsilon) < \infty$.

Definition 3.3.5 (ϵ -packing). A set $N \subseteq T$ is an ϵ -packing of (T, d) if for all $t \neq t'$ in N , $d(t, t') > \epsilon$.

Definition 3.3.6 (Packing number). The ϵ -packing number $M(T, d, \epsilon)$ of (T, d) is defined as

$$M(T, d, \epsilon) = \sup\{|N| : N \text{ is an } \epsilon\text{-packing of } (T, d)\}.$$

As the title suggests, we define the following **metric** properties, which is an essential notion helps us bound the expected **empirical process** supremum.

Definition 3.3.7 (Metric entropy). The *metric entropy* of (T, d) is defined as $\log M(T, d, \epsilon)$.

The fact that we're using **packing number** $M(T, d, \epsilon)$ when defining **metric entropy** is not relevant here due to the following.

Lemma 3.3.1. For any $\epsilon > 0$,

$$M(T, d, 2\epsilon) \leq N(T, d, \epsilon) \leq M(T, d, \epsilon).$$

Proof. We show them one by one.

Claim. $M(T, d, 2\epsilon) \leq N(T, d, \epsilon)$.

Proof. Take \mathcal{M} to be a **2 ϵ -packing** and \mathcal{N} to be an **ϵ -net**. Then for any $t \in \mathcal{N}$, consider $B(t, \epsilon)$. We see that there is at most one $x \in \mathcal{M}$ such that $d(t, x) \leq \epsilon$ since otherwise, if $x, x' \in \mathcal{M}$ such that $x \neq x'$ and $d(t, x), d(t, x') \leq \epsilon$, then $d(x, x') \leq 2\epsilon$, a contradiction to \mathcal{M} . \otimes

Claim. $N(T, d, \epsilon) \leq M(T, d, \epsilon)$.

Proof. Take \mathcal{M} to be a maximum **ϵ -packing**, it suffices to show that \mathcal{M} is also an **ϵ -net**, i.e., for all $t \in T$, there exists $x \in \mathcal{M}$ such that $d(x, t) \leq \epsilon$. Suppose not, then $d(t, x) > \epsilon$ for all $x \in \mathcal{M}$, i.e., we can add x to \mathcal{M} , contradiction. \otimes

■

For simplicity, we will use the following notations.

Notation. If (T, d) and ϵ are clear from the context, we write $N := N(T, d, \epsilon)$ and $M := M(T, d, \epsilon)$.

Turns out that there's a characterization of the [packing number](#) of the unit ball in euclidean space.

Proposition 3.3.1. Consider $(\mathbb{R}^d, \|\cdot\|)$ where $\|\cdot\|$ is any norm. Denote $B = \{x: \|x\| \leq 1\}$, then for all $\epsilon > 0$,

$$(1/\epsilon)^d \leq M(B, \|\cdot\|, \epsilon) \leq (1 + 2/\epsilon)^d.$$

Proof. For the lower-bound, we see that

$$N \text{Vol}(\epsilon B) \geq \text{Vol}(B) \Rightarrow N\epsilon^d \geq 1.$$

With $N \leq M$ from [Lemma 3.3.1](#), we get the lower-bound.

For the upper-bound, since $\epsilon/2$ balls around points in M are disjoint, union of these $\epsilon/2$ balls will lie in $(1 + \epsilon/2)B$. This implies

$$M \times \left(\frac{\epsilon}{2}\right)^d \times \text{Vol}(B) \leq \left(1 + \frac{\epsilon}{2}\right)^d \times \text{Vol}(B) \Rightarrow M \leq \left(1 + \frac{2}{\epsilon}\right)^d.$$

■

Note. From [Proposition 3.3.1](#), $\log M(\mathbb{R}^d, \|\cdot\|, \epsilon) \approx d \log 1/\epsilon$.

3.3.2 Hölder Smooth Functions

We are interested in looking at function spaces, and the following are the canonical smooth function classes studied in *nonparametric regression*.

Definition 3.3.8 (Hölder smooth function class). Let $\alpha > 0$ and $\beta = \lfloor \alpha \rfloor$. Then the *Hölder smooth function class* \mathcal{S}_α is defined to be the class of functions on $[0, 1]$ such that

- (a) f continuous on $[0, 1]$;
- (b) f is β -times differentiable;
- (c) $|f^{(k)}| \leq 1$ for all $k = 0, \dots, \beta$;
- (d) $|f^{(\beta)}(x) - f^{(\beta)}(y)| \leq |x - y|^{\alpha - \beta}$ for all $x, y \in [0, 1]$.

Note. When $\alpha = 1$, \mathcal{S}_α is a class of 1-Lipschitz functions.

Remark. The [Hölder smooth function classes](#) are nested, so it's not surprising that the [metric entropies](#) decrease as α increases.

Now, let $d(f, g) = \sup_{x \in [0, 1]} |f(x) - g(x)|$, then (\mathcal{S}_α, d) is a [pseudo-metric](#) space.

Theorem 3.3.1. There exists c_1, c_2 such that for all $\epsilon > 0$,

$$\exp\left(c_2 \epsilon^{-1/\alpha}\right) \leq M(\mathcal{S}_\alpha, d, \epsilon) \leq \exp\left(c_1 \epsilon^{-1/\alpha}\right).$$

Proof sketch. Here we illustrate the basic idea when $\alpha = 1$, i.e., the set of $[0, 1]$ valued 1-Lipschitz functions on $[0, 1]$. We only sketch the proof of the upper-bound, since the lower-bound is similar.

Firstly, we partition both the domain and the range of f with small intervals with width ϵ , resulting in $1/\epsilon$ small intervals on both the x -axis and the y -axis.

Take any function $f \in \mathcal{F}$. We construct a piece-wise constant function \tilde{f} which approximates f . On each small interval in the x -axis, we can define \tilde{f} to be constant, taking value equal to the

midpoint of the interval in the y -axis where the value of f at the left endpoint of this interval (in the x -axis) lies. Then, we have the following.

Claim. $\sup_{x \in [0,1]} |f(x) - \tilde{f}(x)| \leq C\epsilon$.

Proof. Since f cannot vary by more than ϵ in any interval of length ϵ . ⊗

Now, as we vary $f \in \mathcal{F}$, consider the following.

Problem. What is the number of distinct \tilde{f} we can get?

A trivial bound is that, in each small interval on the x -axis, it takes one of the midpoints of the intervals on the y -axis and hence, the number of such functions is bounded by $(\frac{1}{\epsilon})^{\frac{1}{\epsilon}}$.

We can do slightly better. Note that, for the first interval, the number of possible values of \tilde{f} is $\frac{1}{\epsilon}$. However, after that, in the next interval, the value of \tilde{f} can only go up one interval, down one interval, or stay the same (due to 1-Lipschitzness of f), i.e., there are only 3 choices afterward for every interval, going from left to right, resulting an upper bound on the number of distinct \tilde{f} as

$$\frac{1}{\epsilon} 3^{\frac{1}{\epsilon}-1} \leq \exp\left(\frac{C}{\epsilon}\right).$$

■

Remark. Comparing [Proposition 3.3.1](#) and [Theorem 3.3.1](#), we see that the [metric entropy](#) is logarithmic in $1/\epsilon$ versus some exponent of $1/\epsilon$. This is typically the hallmark of a parametric versus a nonparametric function class.

Lecture 11: Gaussian and Sub-Gaussian Process

3.3.3 Sub-Gaussian Process

18 Sep. 9:00

As previously seen. Given a stochastic process $\{X_t\}_{t \in T}$ with (T, d) , we want to bound $\mathbb{E}[\sup_{t \in T} X_t]$.

Recall our [informal principle](#), i.e., if $\{X_t\}_{t \in T}$ is “sufficiently continuous” w.r.t. d , then $\mathbb{E}[\sup_{t \in T} X_t]$ is governed by metric properties (e.g., [metric entropy](#)) of T . We start by considering the [Gaussian process](#).

Definition 3.3.9 (Gaussian process). A stochastic process $\{X_t\}_{t \in T}$ is a *Gaussian process* if for any finite set of indices t_1, \dots, t_k , $(X_{t_1}, \dots, X_{t_k}) \sim \mathcal{N}(0, \Sigma)$.

Clearly, this is a very strong notion due to the following.

Note. For $d(t, t') = \sqrt{\mathbb{E}[(X_t - X_{t'})^2]}$, we have

$$\mathbb{E}\left[e^{\lambda(X_t - X_{t'})}\right] = e^{\lambda^2/2 \mathbb{E}[X_t - X_{t'}]} = \exp\left(\frac{\lambda^2}{2} d^2(t, t')\right).$$

The following generalized process characterizes the concept of “sufficiently continuous”.

Definition 3.3.10 (Sub-Gaussian process). A stochastic process $\{X_t\}_{t \in T}$ is a *sub-Gaussian process* w.r.t. d if $X_t - X_s \sim \text{Subg}(d^2(t, s))$. Assume $\mathbb{E}[X_t] = 0$ for all $t \in T$, then equivalently, for all $t \neq s \in T$ and $\lambda \in \mathbb{R}$,

$$\mathbb{E}\left[e^{\lambda(X_t - X_s)}\right] \leq \exp\left(\frac{\lambda^2}{2} d^2(t, s)\right).$$

It's clear that the [sub-Gaussian](#) condition encodes a strong notion of continuity (in probability) of the stochastic process $\{X_t\}_{t \in T}$ w.r.t. d .

Example (Gaussian process). We see that $d(t, t') = \sqrt{\mathbb{E}[(X_t - X_{t'})^2]}$ is the naturally induced **pseudo-metric** such that a **Gaussian process** is **sub-Gaussian**.

Another interesting example is the following.

Example (Rademacher process). Consider the unnormalized **Rademacher width** of a set $T \subseteq \mathbb{R}^n$,

$$R_n(T) = \mathbb{E} \left[\sup_{t \in T} \sum_{i=1}^n \epsilon_i t_i \right].$$

Let $X_t = \langle \epsilon, t \rangle$, then from **Lemma 2.3.3**, $X_t - X_{t'} = \langle \epsilon, t - t' \rangle \sim \text{Subg}(\|t - t'\|_2^2)$, i.e., $X_t \sim \text{Subg}$ w.r.t. $\|\cdot\|_2$. This is the so-called **Rademacher process**.

Inspired by the above example, one can also define the **Gaussian width**.

Definition 3.3.11 (Gaussian width). Let $g_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$. Then the **Gaussian width** of a set $A \subseteq \mathbb{R}^n$ is defined as

$$\text{GW}_n(A) = \mathbb{E} \left[\sup_{a \in A} \sum_{i=1}^n \frac{1}{n} g_i a_i \right].$$

This means that the **Rademacher process** can be slightly modified as follows.

Example. Consider $X_t = \langle g, t \rangle$ where g is a random Gaussian vector. We then have $X_t \sim \text{Subg}$ w.r.t. $\|\cdot\|_2$.

Theorem 3.3.2 (Gaussian width and Rademacher width are similar). For any $n \geq 1$ and any set $T \subseteq \mathbb{R}^n$,

$$R_n(T) \leq \text{GW}_n(T) \leq \sqrt{\log n} R_n(T).$$

Let's look at some examples of (unnormalized) **Rademacher width**.

Example. $R(B_\infty^n) = n$, $R(B_2^n) = \sqrt{n}$, and $R(B_1^n) = 1$.

Proof. We see that

- for ℓ_∞ , the supremum is achieved by matching signs of ϵ , which gives $R_n(B_\infty^n) = n$;
- for ℓ_2 the supremum is achieved by choosing $t = \epsilon / \|\epsilon\|_2$, then we get $R(B_2^n) = \mathbb{E}[\|\epsilon\|_2] = \sqrt{n}$;
- for ℓ_1 , from Hölder's inequality, $\langle \epsilon, t \rangle \leq \|\epsilon\|_\infty \|t\|_1 = 1$.

⊛

Example (Supremum of empirical process). Let \mathcal{F} be a class of functions bounded by 1. Let $X_f = \sqrt{n}(\mathbb{P}_n f - \mathbb{P} f)$, and consider $\{X_f\}_{f \in \mathcal{F}}$. Then,

$$X_f - X_g = \sqrt{n} \cdot \frac{1}{n} \sum_{i=1}^n \underbrace{(f(x_i) - g(x_i) - \mathbb{P} f + \mathbb{P} g)}_{\leq 2\|f - g\|_\infty} \sim \text{Subg}(4\|f - g\|_\infty^2),$$

hence $\{X_f\}_{f \in \mathcal{F}} \sim \text{Subg}$ w.r.t. $d(f, g) = 2\|f - g\|_\infty$.

These are all simple sets. For an arbitrary set however, we need more general tools in order to compute the **Rademacher width**. Firstly, recall the following.

Definition 3.3.12 (Diameter). The **diameter** of (T, d) is defined as $\text{Diam}(T) = \sup_{t, t' \in T} d(t, t')$.

3.3.4 Single Scale Bound for Expected Supremum of Sub-Gaussian Process

We're going to see the most sophisticated tools in this course. We first see a preliminary version of which and generalize it later.

Lemma 3.3.2 (Single scale bound). Let $\{X_t\}_{t \in T}$ be a centered [sub-Gaussian process](#) on (T, d) w.r.t. d . Then

$$\mathbb{E} \left[\sup_{t \in T} X_t \right] \leq \inf_{\epsilon > 0} \left(\mathbb{E} \left[\sup_{\substack{t, t' \in T: \\ d(t, t') \leq \epsilon}} X_t - X_{t'} \right] + \text{Diam}(T) \sqrt{2 \log N(T, d, \epsilon)} \right).$$

Proof. We first note that $\mathbb{E} [\sup_{t \in T} X_t] = \mathbb{E} [\sup_{t \in T} X_t - X_{t_0}]$ for some fixed $t_0 \in T$. Now, take an [\$\epsilon\$ -net](#) N with $\pi(t) \in N$ denotes the point such that $d(t, \pi(t)) \leq \epsilon$, then

$$\mathbb{E} \left[\sup_{t \in T} X_t - X_{t_0} \right] \leq \mathbb{E} \left[\sup_{t \in T} X_t - X_{\pi(t)} \right] + \mathbb{E} \left[\sup_{t \in T} X_{\pi(t)} - X_{t_0} \right]$$

Observe that $X_{\pi(t)} - X_{t_0} \sim \text{Subg}(\text{Diam}^2(T))$, then the second term is a finite maximum such that

$$\mathbb{E} \left[\sup_{t \in T} X_{\pi(t)} - X_{t_0} \right] \leq \sqrt{2 \text{Diam}^2(T) \log N(T, d, \epsilon)} = \text{Diam}(T) \sqrt{2 \log N(T, d, \epsilon)}$$

from [Lemma 3.1.2](#). By rewriting the first term, we have

$$\mathbb{E} \left[\sup_{t \in T} X_t \right] \leq \inf_{\epsilon > 0} \left(\mathbb{E} \left[\sup_{\substack{t, t' \in T: \\ d(t, t') \leq \epsilon}} X_t - X_{t'} \right] + \text{Diam}(T) \sqrt{2 \log N(T, d, \epsilon)} \right).$$

■

Notation (Approximation error). The first term in the [single scale bound](#) is the *approximation error*.

We see that the first term in the [single scale bound](#) is still an infinite maximum, so it is not clear how to bound it. Typically, we have to do something crude here. There are some exceptions, though.

Example. For [Rademacher processes](#), we have $\mathbb{E} \left[\sup_{t, t' \in T: \|t - t'\| \leq \delta} \langle \epsilon, t - t' \rangle \right] \leq \|\epsilon\| \delta \leq \sqrt{n} \delta$.

Remark. As ϵ decreases, the approximation error should get smaller and the finite maximum increases. Therefore, when we use the [single scale bound](#) we can then choose an optimum ϵ to minimize the sum of these two.

Let's see some applications of [single scale bound](#) which show that the [single scale bound](#) may not get the optimal rate.

Example. Consider a finite set $T = \{(0, 0, \dots, 0), (1, 0, \dots, 0), \dots, (1, 1, \dots, 1)\} \subseteq \mathbb{R}^n$, i.e., the footprint of the boolean function class on \mathbb{R} given by $\{\mathbb{1}_{x \leq \theta}\}_{\theta \in \mathbb{R}}$. By [Lemma 3.1.2](#), $R_n(T) \leq \sqrt{n \log n}$.

As previously seen. We [claimed](#) that $\log n$ is superfluous.

We still can't remove the $\sqrt{\log n}$: from the [single scale bound](#), with $\text{Diam}(T) = \sqrt{n}$,

$$R_n(T) \leq \sqrt{n} \epsilon + \sqrt{n} \sqrt{\log N(T, \|\cdot\|_2, \epsilon)}.$$

To remove $\log n$, ϵ needs to be $O(1)$ for the first term. But then $\log N(T, \|\cdot\|_2, \epsilon) \rightarrow \infty$, and we fail.

Now, let's revisit the [previous example](#), and recall the following.

As previously seen. For a class of functions bounded by 1, $X_f \sim \text{Subg}(2^2\|f - g\|_\infty^2)$, i.e., $X_f - X_g \leq c\sqrt{n}\|f - g\|_\infty$ almost surely.

Example (Empirical process supremum of S_1). Consider $X_f = \sqrt{n}(\mathbb{P}_n f - \mathbb{P}f)$ on $\mathcal{F} = S_1$, i.e., functions bounded by 1 and are 1-Lipschitz on $[0, 1]$. So in particular, $X_f - X_g \leq c\sqrt{n}\|f - g\|_\infty$. From the [single scale bound](#) and [Theorem 3.3.1](#),

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} X_f \right] \leq c \left(\sqrt{n\epsilon} + \sqrt{1/\epsilon} \right) = c(\sqrt{n} \cdot n^{-1/3})$$

by letting $\epsilon = n^{-1/3}$ (where this bound is minimized), giving us

$$\mathbb{E} \left[\sup_{f \in S_1} \mathbb{P}_n f - \mathbb{P}f \right] \leq \frac{c}{n^{1/3}}.$$

This is the first non-trivial bound we have shown besides boolean function classes.

However, observe that $X_f - X_g \leq C\sqrt{n}\|f - g\|_\infty$ implies $X_f - X_g \leq \|f - g\|_\infty$ in probability. The fact that we are stuck with the above almost surely bound and don't know how to incorporate this additional information, suggests that this bound is still not optimal.

Remark. The optimal bound for S_1 is c/\sqrt{n} , i.e., the CLT rate.

It's perhaps surprising that for the class of functions S_1 , we get the $O(n^{-1/2})$ rate for the supremum of the [empirical process](#), because even for a single function $f \in S_1$, we would still have got the $O(n^{-1/2})$ rate. This is not always the case, though. For Lipschitz function defined on $[0, 1]^d$, the rates are slower. We state this result without proof for now.

Lemma 3.3.3. Let $S_{1,d}$ to be the set of 1-bounded 1-Lipschitz functions w.r.t. the Euclidean norm defined on $[0, 1]^d$. Then there exists a universal constant $C > 0$ such that

$$\mathbb{E} \left[\sup_{f \in S_{1,d}} \mathbb{P}_n f - \mathbb{P}f \right] \leq \begin{cases} Cn^{-1/2}, & \text{if } d = 1; \\ Cn^{-1/2} \log n, & \text{if } d = 2; \\ Cn^{-1/d} \log n, & \text{if } d > 2. \end{cases}$$

These rates are tight and corresponding lower bounds are also known [[Han16](#), Problem 5.11 (d)].

Lecture 12: Chaining Method and Dudley's Entropy Bound

3.3.5 Dudley's Entropy Bound

20 Sep. 9:00

To overcome the limitation of the [single scale bound](#), we can repeatedly taking [\$\epsilon\$ -net](#), which is considered as a *multi-scale bound*. The theorem requires one technical assumption of the stochastic process.

Definition 3.3.13 (Separable). We say that $\{X_t\}_{t \in T}$ is a *separable* process if there exists a countable $T_0 \subseteq T$ such that (outside a null set) for all $t \in T$, there exists $\{t_n \in T_0\}_n$ such that $d(t_n, t) \rightarrow 0$ satisfying $\lim_{n \rightarrow \infty} X_{t_n} = X_t$.

It's clear that $\sup_{t \in T_0} X_t = \sup_{t \in T} X_t$. Moreover, this notion is consistent with the separability of a topological space² we saw in real analysis.

Example (Separable metric space). If (T, d) is separable (as a topological space), $\{X_t\}$ has countable sample path almost surely, then $\{X_t\}$ is [separable](#).

²A topological space is *separable* if it contains a countable dense subset.

Now, we can state the bound we want.

Theorem 3.3.3 (Dudley's entropy bound). Let $\{X_t\}_{t \in T}$ be a centered and separable sub-Gaussian process on (T, d) w.r.t. d . Then

$$\mathbb{E} \left[\sup_{t \in T} X_t \right] \leq 6 \sum_{k \in \mathbb{Z}} 2^{-k} \sqrt{\log N(T, d, 2^{-k})}.$$

Proof. Consider the case that $|T| < \infty$ and $|T| = \infty$.

Claim. The result holds for $|T| < \infty$.

Proof. Let K_0 be the largest integer such that $2^{-K_0} \geq \text{Diam}(T)$, and let K_1 be the smallest integer such that $0 < 2^{-K_1} < \min_{s \neq t \in T} d(s, t)$. Then we let N_k be a 2^{-k} -net of T such that

- $k = K_0$: $N_{K_0} = \{t_0\}$ is a 2^{-K_0} -net of T for a fixed $t_0 \in T$.
- $k = K_1$: $N_{K_1} = T$ is a 2^{-K_1} -net of T .

Write $\pi_k(t)$ for the closest element in N_k to t , in particular, $d(t, \pi_k(t)) \leq 2^{-k}$. By writing

$$X_t = X_{\pi_{K_1}(t)} - X_{\pi_{K_0}(t)} = X_{\pi_{K_1}(t)} - X_{\pi_{K_1-1}(t)} + X_{\pi_{K_1-1}(t)} - \cdots + X_{\pi_{K_0+1}(t)} - X_{\pi_{K_0}(t)},$$

we have

$$\begin{aligned} \mathbb{E} \left[\sup_{t \in T} X_t \right] &= \mathbb{E} \left[\sup_{t \in T} X_t - X_{t_0} \right] \\ &= \mathbb{E} \left[\sup_{t \in T} \sum_{k=K_0+1}^{K_1} (X_{\pi_k(t)} - X_{\pi_{k-1}(t)}) \right] \leq \sum_{k=K_0+1}^{K_1} \mathbb{E} \left[\sup_{t \in T} (X_{\pi_k(t)} - X_{\pi_{k-1}(t)}) \right]. \end{aligned}$$

Since the cardinality of $\{X_{\pi_k(t)} - X_{\pi_{k-1}(t)}\}_{t \in T}$ is $|N_k| |N_{k-1}| \leq |N_k|^2$, with

$$X_{\pi_k(t)} - X_{\pi_{k-1}(t)} \sim \text{Subg}(d(\pi_k(t), \pi_{k-1}(t)))$$

where $d(\pi_k(t), t) + d(t, \pi_{k-1}(t)) \leq 2^{-k} + 2^{-(k+1)} \leq 3 \cdot 2^{-k}$. From Lemma 3.1.2, for each k ,

$$\mathbb{E} \left[\sup_{t \in T} (X_{\pi_k(t)} - X_{\pi_{k-1}(t)}) \right] \leq 3 \times 2^{-k} \sqrt{2 \log |N_k|^2} = 6 \times 2^{-k} \sqrt{\log |N_k|}.$$

⊗

Claim. The result holds for $|T| = \infty$.

Proof. From separability, there exists a countable T_0 such that $\mathbb{E} [\sup_{t \in T_0} X_t] = \mathbb{E} [\sup_{t \in T} X_t]$. Let T_k be a countable approximation of T_0 , then $\sup_{t \in T_k} X_t \rightarrow \sup_{t \in T_0} X_t$ as $k \rightarrow \infty$, so

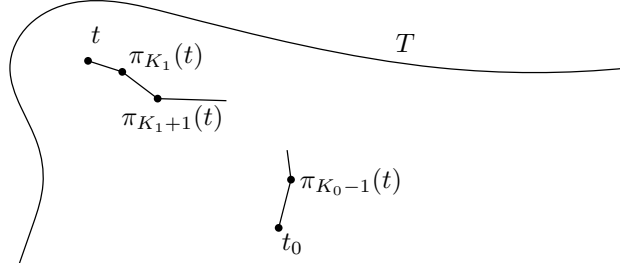
$$\mathbb{E} \left[\sup_{t \in T_k} X_t \right] \rightarrow \mathbb{E} \left[\sup_{t \in T_0} X_t \right] = \mathbb{E} \left[\sup_{t \in T} X_t \right] \text{ as } k \rightarrow \infty$$

from monotone convergence theorem. Hence, it suffices to bound $\mathbb{E} [\sup_{t \in T_k} X_t]$ instead of $\mathbb{E} [\sup_{t \in T} X_t]$ for each k . As $|T_k| < \infty$ and $N(T_k, d, 2^{-k}) \leq N(T_0, d, 2^{-k})$ for all k ,

$$6 \sum_{k \in \mathbb{Z}} 2^{-k} \sqrt{\log N(T_k, d, 2^{-k})} \leq 6 \sum_{k \in \mathbb{Z}} 2^{-k} \sqrt{\log N(T_0, d, 2^{-k})}.$$

⊗

Note (Chaining method). This method is called *chaining* since we're constructing a chain of $X_{\pi_k(t)}$, with smaller and smaller distances.



An alternative integral form of [Dudley's entropy bound](#) is given by the following.

Corollary 3.3.1 (Dudley integral entropy bound). Let $\{X_t\}_{t \in T}$ be a centered and [separable sub-Gaussian process](#) on (T, d) w.r.t. d . Then

$$\mathbb{E} \left[\sup_{t \in T} X_t \right] \leq 12 \int_0^{\text{Diam}(T)} \sqrt{\log N(T, d, \epsilon)} d\epsilon.$$

Proof. Observe that

$$\begin{aligned} \sum_{k \in \mathbb{Z}} 2^{-k} \sqrt{\log N(T, d, 2^{-k})} &= 2 \sum_{k \in \mathbb{Z}} \int_{2^{-k-1}}^{2^{-k}} \sqrt{\log N(T, d, 2^{-k})} d\epsilon \\ &\leq 2 \sum_{k \in \mathbb{Z}} \int_{2^{-k-1}}^{2^{-k}} \sqrt{\log N(T, d, \epsilon)} d\epsilon && N(T, d, \epsilon) \nearrow \text{ as } \epsilon \searrow \\ &= 2 \int_0^\infty \sqrt{\log N(T, d, \epsilon)} d\epsilon \\ &= 2 \int_0^{\text{Diam}(T)} \sqrt{\log N(T, d, \epsilon)} d\epsilon. && \epsilon > \text{Diam}(T), N(T, d, \epsilon) = 1 \end{aligned}$$

Now, we note that we have finally reached the optimal bound for S_1 , solving the problems we saw in the [previous example](#).

Example (Supremum of empirical process for S_1). Consider the [separable sub-Gaussian process](#) $X_f = \sqrt{n}(\mathbb{P}_n f - \mathbb{P} f)$ for $\mathcal{F} = S_1$. In particular, $f, g \in \mathcal{F}$ are 1-Lipschitz on $[0, 1]$ satisfying $|f|, |g| \leq 1$ and $X_f - X_g \in \text{Subg}(2^2 \|f - g\|_\infty^2)$. Since $\text{Diam}(\mathcal{F}) = 2$ and for all $\epsilon < 1/2$,

$$N(S_1, \|\cdot\|_\infty, \epsilon) = \exp(c/\epsilon)$$

from [Theorem 3.3.1](#). Then by the [Dudley's integral entropy bound](#),

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} X_f \right] \leq 12 \int_0^2 \sqrt{\log N(S_1, \|\cdot\|_\infty, \epsilon)} d\epsilon = 12 \int_0^2 \sqrt{\frac{c}{\epsilon}} d\epsilon = 24\sqrt{2c} < O_n(1).$$

Dividing both sides by \sqrt{n} , we achieve the optimal rate $\mathbb{E} [\sup_{f \in \mathcal{F}} (\mathbb{P}_n f - \mathbb{P} f)] = O(1/\sqrt{n})$.

Remark. The [Dudley's integral entropy bound](#) for \mathcal{S}_α is also finite; while for function classes with [covering number](#) as $\exp(c/\epsilon^2)$ is divergent.

Lecture 13: More on Chaining

Let's see some alternate forms of [Dudley's integral entropy bound](#). In the following, assume that $\{X_t\}_{t \in T}$ is a centered and [separable sub-Gaussian process](#) on (T, d) w.r.t. d . 22 Sep. 9:00

Corollary 3.3.2 (Difference form). The same bound as the [Dudley's integral entropy bound](#) holds for $\mathbb{E} [\sup_{t \in T} |X_t - X_{t_0}|]$ and $\mathbb{E} [\sup_{s, t \in T} |X_s - X_t|]$.

Proof. This can be proved by the same [chaining argument](#) with triangle inequality. ■

Corollary 3.3.3 (High probability form). The high probability bound version holds:

$$\mathbb{P} \left(\sup_{s, t \in T} |X_s - X_t| \leq C \left(\int_0^\infty \sqrt{\log N(T, d, \epsilon)} d\epsilon + u \text{Diam}(T) \right) \right) \geq 1 - 2e^{-u^2}.$$

Proof. The proof is the same, where we first show the high probability bound for finite case. ■

Corollary 3.3.4 (Finite resolution form). The following generalizes the [Dudley's integral entropy bound](#) in the sense that $\delta > 0$:

$$\mathbb{E} \left[\sup_{t \in T} X_t \right] \leq C \left(\left[\sup_{\substack{t, t' \in T \\ d(t, t') \leq \delta}} X_t - X_{t'} \right] + \int_\delta^\infty \sqrt{\log N(T, d, \epsilon)} d\epsilon \right).$$

Proof. The proof is still the same, but instead we start with finite resolution. ■

The [finite resolution version](#) is useful since the [entropy](#), integral can diverge, e.g., if $\log N(t, d, \epsilon) = \Omega(1/\epsilon^2)$. Moreover, this can be used to show [Lemma 3.3.3](#).

Remark. We can moreover write

$$\begin{aligned} \mathbb{E} \left[\sup_{t \in T} X_t \right] &\leq C \int_0^{\text{Diam}(T)} \sqrt{\log N(T, d, \epsilon)} d\epsilon \\ &\leq C \left(\int_0^{\text{Diam}(T)/2} \sqrt{\log N(T, d, \epsilon)} d\epsilon + \int_{\text{Diam}(T)/2}^{\text{Diam}(T)} \sqrt{\log N(T, d, \epsilon)} d\epsilon \right) \\ &\leq 2C \int_0^{\text{Diam}(T)/2} \sqrt{\log N(T, d, \epsilon)} d\epsilon. \end{aligned}$$

3.3.6 Uniform Entropy Integral Bound

Let's discuss some limitation of the [Dudley's integral entropy bound](#). First, recall the following.

As previously seen. In the [example of the optimal rate for \$S_1\$](#) , whenever

$$\int_0^\infty \sqrt{\log N(T, d, \epsilon)} d\epsilon < \infty \Rightarrow \mathbb{E} \left[\sup_f \mathbb{P}_n f - \mathbb{P} f \right] \leq c/\sqrt{n}.$$

Note that we're doing [chaining](#) w.r.t. $\|\cdot\|_\infty$ on \mathcal{F} so far. To see its limitation, consider again the boolean function classes \mathcal{F} and let $X_f = \sqrt{n}(\mathbb{P}_n f - \mathbb{P} f) \sim \text{Subg}(c^2 \|\cdot\|_\infty^2)$. From [Proposition 3.2.1](#),

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \mathbb{P}_n f - \mathbb{P} f \right] \leq \sqrt{\frac{\text{VC}(\mathcal{F}) \log n}{n}}.$$

Now, observe that for any $f \neq g$ in \mathcal{F} , $\|f - g\|_\infty = 1$. This implies that by taking $\epsilon \in (0, 1)$,

$$N(\mathcal{F}, \|\cdot\|_\infty, \epsilon) = |\mathcal{F}| = \infty$$

for any interesting case, e.g., $\mathcal{F} = \{\mathbb{1}_{x \leq \theta}\}_{\theta \in \mathbb{R}}$, i.e., [chaining](#) w.r.t. $\|\cdot\|_\infty$ only gives a vacuous bound.

Intuition. To fix this, we can use the idea of the [symmetrization](#).

Firstly, given some observed i.i.d. data x_1, \dots, x_n , recall the following.

As previously seen. By conditioning on the data x_1, \dots, x_n , [symmetrization](#) shows that

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \sqrt{n}(\mathbb{P}_n f - \mathbb{P} f) \right] \leq \frac{2}{\sqrt{n}} R_n(\mathcal{F}) = 2 \mathbb{E}_{x, \epsilon} \left[\sup_{f \in \mathcal{F}} \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(x_i) \right].$$

Specifically, we want to look at $\mathbb{E}_x [\mathbb{E}_\epsilon [\sup_{f \in \mathcal{F}} X_f]]$ and compute the [Rademacher width](#). Let $X_f = \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(x_i)$, we have

$$X_f - X_g = \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i (f(x_i) - g(x_i)) \sim \text{Subg}(\|(f(x_i))_i - (g(x_i))_i\|_2^2) = \text{Subg}\left(\frac{1}{n} \sum_{i=1}^n (f(x_i) - g(x_i))^2\right),$$

where $(f(x_i))_i = (f(x_1), \dots, f(x_n))$. Hence, $X_f \sim \text{Subg}\left(\frac{1}{n} \sum_{i=1}^n (f(x_i) - g(x_i))^2\right)$.

Note. We're already doing better since $\sqrt{\frac{1}{n} \sum_{i=1}^n (f(x_i) - g(x_i))^2} \leq \|f - g\|_\infty$.

We see that $\sqrt{\frac{1}{n} \sum_{i=1}^n (f(x_i) - g(x_i))^2}$ is similar to $\|f - g\|_2$, but just on the empirical measure (with i.i.d. data x_i 's). Hence, consider the following notation.

Notation. Let $L_2(\mathbb{P}_n)$ denote the [metric](#) w.r.t. \mathbb{P}_n ^a such that

$$L_2^2(\mathbb{P}_n)(f, g) := \frac{1}{n} \sum_{i=1}^n (f(x_i) - g(x_i))^2.$$

^aFormally, \mathbb{P}_n is the empirical measure uniform on $\{x_i\}_{i=1}^n$.

In our new notation, $X_f \sim \text{Subg}(L_2(\mathbb{P}_n))$. Now, we can do the [chaining argument](#) on $L_2(\mathbb{P}_n)$ and get

$$\mathbb{E} [\sup \sqrt{n}(\mathbb{P}_n f - \mathbb{P} f)] \leq C \int_0^{\text{Diam}(\mathcal{F})} \sqrt{\log N(\mathcal{F}, L_2(\mathbb{P}_n), \epsilon)} d\epsilon,$$

where

$$\text{Diam}(\mathcal{F}) = \sup_{f, g} L_2(\mathbb{P}_n)(f, g) = \frac{1}{n} \sum_{i=1}^n (f(x_i) - g(x_i))^2 \leq \sup_{f \in \mathcal{F}} \sqrt{\mathbb{P}_n f^2},$$

hence we have

$$\mathbb{E} [\sup \sqrt{n}(\mathbb{P}_n f - \mathbb{P} f)] \leq C \cdot \mathbb{E}_x \left[\int_0^{\sup_{f \in \mathcal{F}} \sqrt{\mathbb{P}_n f^2}} \sqrt{\log N(\mathcal{F}, L_2(\mathbb{P}_n), \epsilon)} d\epsilon \right].$$

However, there's a problem.

Problem. $L_2(\mathbb{P}_n)$ is a “random” [metric](#), so $N(\mathcal{F}, L_2(\mathbb{P}_n), \epsilon)$ is hard to compute.

Answer. To resolve this, we take the supremum over all measures μ supported on χ , i.e.,

$$C \mathbb{E}_x \left[\int_0^{\sup_{f \in \mathcal{F}} \sqrt{\mathbb{P}_n f^2}} \sqrt{\log N(\mathcal{F}, L_2(\mathbb{P}_n), \epsilon)} d\epsilon \right] \leq C \mathbb{E}_x \left[\int_0^{\sup_{f \in \mathcal{F}} \sqrt{\mathbb{P}_n f^2}} \sqrt{\sup_{\mu} \log N(\mathcal{F}, L_2(\mu), \epsilon)} d\epsilon \right].$$

⊛

This might seem very bad, but actually it's not since $L_2(\mu) < L_\infty$. Specifically, to bound this supremum over all measures, consider the following.

Definition 3.3.14 (Koltchinskii-Pollard entropy). The *Koltchinskii-Pollard entropy* of \mathcal{F} is defined as

$$\sup_{\mu} \log N(\mathcal{F}, L_2(\mu), \epsilon).$$

Example. For boolean function classes, $\sup_f \sqrt{\mathbb{P}_n f^2} \leq 1$.

We then have the following for the boolean function classes.

Intuition (Main bound). Let \mathcal{F} be a boolean function class, then since $\sup_{f \in \mathcal{F}} \sqrt{\mathbb{P}_n f^2} \leq 1$,

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \sqrt{n} |\mathbb{P}_n f - \mathbb{P} f| \right] \leq C \mathbb{E}_x \left[\int_0^1 \sqrt{\sup_{\mu} \log N(\mathcal{F}, L_2(\mu), \epsilon)} d\epsilon \right].$$

More generally, if we have $F \geq f$ (called *envelope*) for all $f \in \mathcal{F}$ such that $\mathbb{P} F^2 < \infty$, this holds.

Problem. How can we compute the *Koltchinskii-Pollard entropy*?

Answer. We can use some notions of combinatorial dimension (e.g., *VC dimension*) upper-bounds the *Koltchinskii-Pollard entropy* such that

$$\sup_{\mu} \log N(\mathcal{F}, L_2(\mu), \epsilon) \leq (c_1/\epsilon)^{c_2 \times \text{VC}(\mathcal{F})} \approx \epsilon^{-d}$$

for d being “dimension” (parametric). ⊛

Remark. This implies a $\sqrt{\text{VC}(\mathcal{F})/n}$ rate (without a log term!) for $\mathbb{E} [\sup (\mathbb{P}_n f - \mathbb{P} f)]$.

Lecture 14: Uniform Entropy Integral Bound

As previously seen. Motivated by the fact that boolean function classes are not *totally bounded* w.r.t. ℓ_{∞} , $\|f - g\|_{\infty} = 1$, we’re trying to establish the *bound* on

25 Sep. 9:00

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \sqrt{n} (\mathbb{P}_n f - \mathbb{P} f) \right] \leq 2 \mathbb{E}_x \left[\mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(x_i) \right] \right] = \frac{1}{\sqrt{n}} R_n(\mathcal{F})$$

where the inner expectation is just the *Rademacher width* $R_n(\{f(x_1), \dots, f(x_n)\}_{f \in \mathcal{F}})$.^a Let $X_f = \langle \epsilon, f \rangle / \sqrt{n}$, then $\{X_f\}_{f \in \mathcal{F}}$ is *sub-Gaussian* w.r.t. $L_2(\mathbb{P}_n)$.^b

^aWe can also abuse the notation $R_n(\mathcal{F})$ by $R_n(\{f(x_1)/\sqrt{n}, \dots, f(x_n)/\sqrt{n}\}_{f \in \mathcal{F}})$, but let’s not do this.

^bRecall that $L_2^2(\mathbb{P}_n)(f, g) = \frac{1}{n} \sum_{i=1}^n (f(X_i) - g(X_i))^2$, compared to $L_2(\mathbb{P})(f, g) = \int (f(x) - g(x))^2 d\mathbb{P}$.

The obstacle we’re facing is the lack of control of $\sup_{f \in \mathcal{F}} \sqrt{\mathbb{P}_n f^2}$. Consider the following notion.

Definition 3.3.15 (Envelope). A non-negative valued function $F: \chi \rightarrow [0, \infty]$ is an *envelope* for \mathcal{F} if $\sup_{f \in \mathcal{F}} |f(x)| \leq F(x)$ for all $x \in \chi$.

Example. For boolean function classes, $F(x) = 1$ is an *envelope*.

Remark. Let F be an *envelope* of \mathcal{F} , then $\sup_{f \in \mathcal{F}} \sqrt{\mathbb{P}_n f^2} \leq \sqrt{\mathbb{P}_n F^2}$, as we want.

With this new notion, we can state the main bound we want, i.e., the *uniform entropy integral bound*.

Theorem 3.3.4 (Uniform entropy integral bound). Given a function class \mathcal{F} and an [envelope](#) F of \mathcal{F} such that $\mathbb{P}F^2 < \infty$, then for $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$,

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \sqrt{n} |\mathbb{P}_n f - \mathbb{P} f| \right] \leq \frac{2}{\sqrt{n}} R_n(\mathcal{F}) \leq C \|F\|_{L_2(\mathbb{P})} \int_0^1 \sqrt{\log \sup_{\mu} N(\mathcal{F}, L_2(\mu), \epsilon \sqrt{\mu F^2})} d\epsilon.$$

Proof. Summarizing what we have established, we have

$$\begin{aligned} \mathbb{E} \left[\sup_{f \in \mathcal{F}} \sqrt{n} (\mathbb{P}_n f - \mathbb{P} f) \right] &\leq 2 \mathbb{E}_x \left[\mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(x_i) \right] \right] && \text{symmetrization} \\ &= \frac{2}{\sqrt{n}} R_n(\mathcal{F}) \\ &\leq C \mathbb{E}_x \left[\int_0^{\sup_{f, g \in \mathcal{F}} \frac{L_2(\mathbb{P}_n)(f, g)}{2}} \sqrt{\log N(\mathcal{F}, L_2(\mathbb{P}_n), \epsilon)} d\epsilon \right] && \text{Dudley's bound} \\ &\leq C \mathbb{E}_x \left[\int_0^{\sup_{f \in \mathcal{F}} \sqrt{\mathbb{P}_n f^2}} \sqrt{\log N(\mathcal{F}, L_2(\mathbb{P}_n), \epsilon)} d\epsilon \right] \\ &\leq C \mathbb{E}_x \left[\int_0^{\sqrt{\mathbb{P}_n F^2}} \sqrt{\log N(\mathcal{F}, L_2(\mathbb{P}_n), \epsilon)} d\epsilon \right] \\ &= C \mathbb{E}_x \left[\sqrt{\mathbb{P}_n F^2} \int_0^1 \sqrt{\log N(\mathcal{F}, L_2(\mathbb{P}_n), \epsilon \sqrt{\mathbb{P}_n F^2})} d\epsilon \right] && \epsilon \leftarrow \sqrt{\mathbb{P}_n F^2} \epsilon \\ &\leq C \mathbb{E}_x \left[\sqrt{\mathbb{P}_n F^2} \int_0^1 \sqrt{\sup_{\mu} \log N(\mathcal{F}, L_2(\mu), \epsilon \sqrt{\mu F^2})} d\epsilon \right] \\ &\leq C \left[\int_0^1 \sqrt{\sup_{\mu} \log N(\mathcal{F}, L_2(\mu), \epsilon \sqrt{\mu F^2})} d\epsilon \right] \mathbb{E}_x \left[\sqrt{\mathbb{P}_n F^2} \right] \\ &\text{from Jensen's inequality, } \mathbb{E} \left[\sqrt{\mathbb{P}_n F^2} \right] \leq \sqrt{\mathbb{E} [\mathbb{P}_n F^2]} = \sqrt{\mathbb{P} F^2} = \|F\|_{L_2(\mathbb{P})}, \\ &\leq C \|F\|_{L_2(\mathbb{P})} \left[\int_0^1 \sqrt{\sup_{\mu} \log N(\mathcal{F}, L_2(\mu), \epsilon \sqrt{\mu F^2})} d\epsilon \right]. \end{aligned}$$

■

Notation. We sometimes denote $\int_0^1 \sqrt{\log \sup_{\mu} N(\mathcal{F}, L_2(\mu), \epsilon \sqrt{\mu F^2})} d\epsilon$ by $\mathbf{J}(F, \mathcal{F})$.

\mathcal{F} needs not to be bounded, instead, what we really need is an [envelope](#):

Remark. If we apply the above bound to a bounded function class then we do not really need the notion of an envelope. The assumption of an [envelope](#) is slightly more general than assuming boundedness.

Remark. The [Koltchinskii-Pollard entropy](#) integral is free from \mathbb{P} , while $\|F\|_{L_2(\mathbb{P})}$ depends on \mathbb{P} . Thus, the overall bound is distribution free, as has been all of our bounds so far, if the function class is bounded.

Example. For boolean function classes, $\|F\|_{L_2(\mathbb{P})} \leq 1$, so the [bound](#) is uniform over \mathbb{P} .

3.3.7 Uniform L_2 Entropy is Bounded by Combinatorial Dimension

We're now ready to revisit the [problem](#) we asked in the previous lecture:

Problem. How can we bound the uniform L_2 -entropy, i.e., [Koltchinskii-Pollard entropy](#)?

Answer. For boolean function classes, [Koltchinskii-Pollard entropy](#) can be bounded in terms of [VC dimension](#); for non-boolean function classes, an extended notion of [VC dimension](#) is needed. \circledast

Theorem 3.3.5 (Dudley). Let \mathcal{F} be a boolean function class, then there exist absolute constants c_1, c_2 such that for all $0 < \epsilon < 1$,

$$\sup_{\mu} N(\mathcal{F}, L_2(\mu), \epsilon) \leq \left(\frac{c_1}{\epsilon}\right)^{c_2 \text{VC}(\mathcal{F})}$$

Proof. It suffices to bound the [packing number](#) instead of the [covering number](#) from [Lemma 3.3.1](#). To upper-bound the [packing number](#) via $d := \text{VC}(\mathcal{F})$, first fix a probability measure μ on χ , consider a maximum ϵ -packing $M = \{f_1, \dots, f_N\}$ of \mathcal{F} w.r.t. $L_2(\mu)$. Then for all $i \neq j$,

$$\int (f_i - f_j)^2 d\mu = \mu(f_i \neq f_j) > \epsilon^2.$$

Then, sample K points W_1, \dots, W_K i.i.d. from μ , we want that all $\{f_i\}_{i=1}^N$ to have different values on (w_1, \dots, w_K) . Note that for $i \neq j$,

$$\mu(f_i = f_j \text{ on } w_1, \dots, w_K) \leq (1 - \epsilon^2)^K \leq e^{-K\epsilon^2}$$

from $(1 - x)^k \leq e^{-kx}$. This implies

$$\mu(\exists \text{ at least one pair } i \neq j \text{ such that } f_i = f_j \text{ on } w_1, \dots, w_K) \leq \binom{N}{2} e^{-K\epsilon^2},$$

hence

$$\mu(\text{all the } f_i \text{'s are distinct on } w_1, \dots, w_K) \geq 1 - \binom{N}{2} e^{-K\epsilon^2} \geq \frac{1}{2}$$

by choosing $K\epsilon^2 \approx 2 \log N$. We conclude that there exists K points w_1, \dots, w_K such that all the f_i 's are distinct on $\{w_1, \dots, w_K\}$. From [Sauer-Shelah lemma](#),^a

$$N = |\mathcal{F}(w_1, \dots, w_K)| \leq \left(\frac{eK}{d}\right)^d = \left(\frac{2e \log N}{\epsilon^2 d}\right)^d.$$

We see that $N \leq (\log N)^d$. To further bound N , consider^b

$$N^{1/d} \leq \frac{4e \log N}{2d\epsilon^2} = \frac{4e}{\epsilon^2} \log N^{1/2d} \leq \frac{4e}{\epsilon^2} N^{1/2d}$$

from $\log x \leq x$, hence $N^{1/2d} \leq 4e/\epsilon^2$, or equivalently,

$$N \leq (4e)^{2d} \left(\frac{1}{\epsilon}\right)^{4d} = \left(\frac{2\sqrt{e}}{\epsilon}\right)^{4d} =: \left(\frac{c_1}{\epsilon}\right)^{c_2 d}.$$

■

^aNote that we have only shown the case for $K \geq \text{VC}(\mathcal{F})$. However, for $K < \text{VC}(\mathcal{F})$, it's also easy to show.

^bWe want to make the exponent a of $\log N^a$ to be less than $1/d$.

Remark. For boolean function classes, while $N(\mathcal{F}, L_\infty, \epsilon) = \infty$, the above shows that

$$\sup_{\mu} N(\mathcal{F}, L_2(\mu), \epsilon) < \infty.$$

We can now finally generalize [Proposition 3.2.1](#), where for boolean function classes, we have

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \sqrt{n} |\mathbb{P}_n f - \mathbb{P} f| \right] \leq c \sqrt{\text{VC}(\mathcal{F}) \log \frac{en}{\text{VC}(\mathcal{F})}}.$$

Corollary 3.3.5. Let \mathcal{F} be a boolean function class, for some constant C , we have

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \sqrt{n} |\mathbb{P}_n f - \mathbb{P} f| \right] \leq C \sqrt{\text{VC}(\mathcal{F})}.$$

Proof. Applying [uniform entropy integral bound](#), with $\|F\|_{L_2(\mathbb{P})} = 1$ and [Theorem 3.3.5](#),

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \sqrt{n} |\mathbb{P}_n f - \mathbb{P} f| \right] \leq C \int_0^1 \sqrt{c_2 \text{VC}(\mathcal{F}) \log \frac{c_1}{\epsilon}} d\epsilon \leq C' \sqrt{\text{VC}(\mathcal{F})} \int_0^1 \log \frac{1}{\epsilon} d\epsilon \leq C' \sqrt{\text{VC}(\mathcal{F})}.$$

■

Remark. Compare [Proposition 3.2.1](#) to [Corollary 3.3.5](#), the extra $\sqrt{\log n}$ factor in the bound which we have now got rid of thanks to [chaining](#). The bound really holds for [Rademacher complexity](#) and as a consequence for suprema of [empirical process](#).

Note. Consider the [classification problem](#) in the statistical learning setting, where \hat{f} is the [ERM](#) over a given boolean function class. Then the [excess risk](#) is also bounded by $C \sqrt{\text{VC}(\mathcal{F})/n}$.

Proof. From [symmetrization](#), the [excess risk](#) can be bounded as

$$\mathbb{E} [L(\hat{f})] - \inf_{f \in \mathcal{F}} \mathbb{E} [L(f)] \leq \mathbb{E} \left[\sup_{f \in \mathcal{F}} \left(\mathbb{E} [f(X)] - \frac{1}{n} \sum_{i=1}^n f(x_i) \right) \right] \leq 2R_n(\mathcal{F}).$$

From [Corollary 3.3.5](#), we finally show that in this example, the $\sqrt{\log n}$ factor is superfluous as well.

⊛

Remark. [Corollary 3.3.5](#) is a distribution-free result, i.e.,

$$\sup_{\mathbb{P}} \mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{P}_n f - \mathbb{P} f| \right] \leq C \sqrt{\frac{\text{VC}(\mathcal{F})}{n}}.$$

This means that boolean function classes with finite [VC dimension](#) are uniform [Glivenko-Cantelli](#).

We note that the “machinery” we have done is the following:

1. Bound [uniform entropy w.r.t. \$L_2\$](#) ([uniform entropy integral bound](#)).
2. Uniform L_2 entropy \leq [VC dimension](#) ([Theorem 3.3.5](#)).

Problem. How to extend this “machinery” to non-boolean function classes?

Answer. Define [VC dimension](#) for non-boolean function classes.

⊛

Lecture 15: Parametric v.s. Non-Parametric

3.3.8 Parametric versus Non-Parametric Function Classes

27 Sep. 9:00

We start by asking the following question.

Problem. What makes a function class “parametric” or “non-parametric”?

Answer. If the function class is a vector space, we can usually use the linear algebra notion of dimension. \ast

Consider the following (not very precise) definition in terms of the [uniform \$L_2\$ entropy](#).

Definition 3.3.16 (Parametric). A function class \mathcal{F} is *parametric* if there exists a notion of dimension $\dim \mathcal{F}$ and a constant C such that

$$\sup_{\mu} N(\mathcal{F}, L_2(\mu), \epsilon) \leq \left(\frac{C}{\epsilon}\right)^{\dim(\mathcal{F})}.$$

Example. Boolean function class on χ with finite [VC dimension](#) is [parametric](#).

Proof. [Dudley’s result](#) directly applies. \ast

Definition 3.3.17 (Non-parametric). A function class \mathcal{F} is *non-parametric* if there is a $p > 0$ and a constant C such that

$$\sup_{\mu} \log N(\mathcal{F}, L_2(\mu), \epsilon) \leq \left(\frac{C}{\epsilon}\right)^p.$$

Let’s consider any [parametric](#) class \mathcal{F} uniformly bounded by 1 with $\dim \mathcal{F} = d$. Then from [Dudley’s theorem](#), the [Rademacher complexity](#) can be bounded as

$$\mathbb{E}_x \left[\mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right] \right] \leq \frac{12}{\sqrt{n}} \int_0^1 \sqrt{\sup_{\mu} \log N(\mathcal{F}, L_2(\mu), \epsilon)} d\epsilon \leq \frac{12}{\sqrt{n}} \int_0^1 \sqrt{d \log \frac{C}{\epsilon}} d\epsilon \leq C' \sqrt{\frac{d}{n}}.$$

Hence, we get the [parametric](#) rate $O(\sqrt{d/n})$, and this result is distribution-free.

Analogously, we want to know what we will get for a [non-parametric](#) function class (uniform bounded by 1)? Now, since for a [non-parametric](#) class, the [uniform \$L_2\$ entropy](#) is $\leq (C/\epsilon)^p$,

$$\begin{aligned} & \mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(x_i) \right] \quad \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(x_i) = X_f \sim \text{Subg}(L_2(\mathbb{P}_n)) \\ & \leq \mathbb{E} \left[\sup_{\substack{f, g \in \mathcal{F}: \\ L_2(\mathbb{P}_n)(f, g) \leq \delta}} X_f - X_g \right] + \int_{\delta}^1 \sqrt{\sup_{\mu} \log N(\mathcal{F}, L_2(\mu), \epsilon)} d\epsilon \quad \text{modified Corollary 3.3.4} \\ & \leq \sqrt{n} \cdot \delta + \int_{\delta}^1 \left(\frac{C}{\epsilon}\right)^{p/2} d\epsilon \end{aligned}$$

Since the choice of δ is arbitrary, we can optimize in terms of δ . There are three cases:

- $p < 2$: Take $\delta = 0$ because the integral converges, and we get a [parametric](#) rate bound for R_n with some constant c :

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right] \leq \frac{c}{\sqrt{n}}.$$

Remark. $O(1/\sqrt{n})$ is a [parametric](#) rate.

Example. Consider [Hölder smooth classes](#). Even though these function classes are [non-parametric](#) according to [Definition 3.3.17](#), in terms of R_n or supremum of [empirical process](#), the rate is still [parametric](#).

- $p > 2$: We see that for all $\delta \in (0, 1)$,

$$\begin{aligned}
 \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right] &\leq \delta + \frac{1}{\sqrt{n}} \int_{\delta}^1 \left(\frac{C}{\epsilon} \right)^{p/2} d\epsilon \\
 &= \delta + \frac{C^{p/2}}{\sqrt{n}} \cdot \frac{\epsilon^{-p/2+1}}{1-p/2} \Big|_{\delta}^1 \\
 &\approx \delta + \frac{1}{\sqrt{n}} (-\epsilon^{-p/2+1}) \Big|_{\delta}^1 \quad \text{dropping constant } (1-p/2 < 0) \\
 &\approx \delta + \frac{1}{\sqrt{n}} \delta^{1-p/2}.
 \end{aligned}$$

Now by optimizing over δ , setting $\delta = \delta^{1-p/2}/\sqrt{n}$, we get the bound $O(n^{-1/p})$.

Remark. $O(n^{-1/p})$ is a **non-parametric** rate, strictly slower than the **parametric** rate $O(n^{-1/2})$.

This upper bound is also tight for certain function classes.

Example. For 1-bounded and 1-Lipschitz functions on $[0, 1]^d$, the **uniform L_2 entropy** (in fact the L_{∞} **entropy**) grows like $(1/\epsilon)^d$.

Proof. Since $|f(x) - f(y)| \leq \|x - y\|_2$, $O(n^{-1/d})$ rate here is tight for $d > 2$. ⊗

- $p = 2$: From the same calculation, we have

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right] \leq \delta + \frac{1}{\sqrt{n}} \int_{\delta}^1 \frac{C}{\epsilon} d\epsilon = \delta + \frac{C}{\sqrt{n}} \ln \frac{1}{\delta} = O \left(\frac{1}{\sqrt{n}} \log n \right)$$

by setting $\delta = O(1/\sqrt{n})$. Hence, the **Rademacher complexity** is bounded by $\log n/\sqrt{n}$, which is “almost” the **parametric** rate up to the extra log factor. This might not be tight.

Remark. To summarize, we have the following bounds on the **Rademacher complexity**:

- **Parametric class:** $C\sqrt{d/n}$.
- **Non-parametric class:**
 - $p < 2$: $C\sqrt{1/n}$;
 - $p = 2$: $C \log n/\sqrt{n}$;
 - $p > 2$: $C \cdot n^{-1/p}$.

Example (Linear function class). Let $\chi = B_2^d$, and $\mathcal{F} = \{x \mapsto w^\top x : w \in B_2^d\}$. For a given data $x_1, \dots, x_n \in \mathbb{R}^d$,

$$\mathcal{F}|_{x_1, \dots, x_n} = \left\{ Xw : w \in B_2^d, X_{n \times d} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix} \right\}.$$

To determine whether \mathcal{F} is **parametric** or **non-parametric**, we need to bound $N(\mathcal{F}, L_2(\mathbb{P}_n), \epsilon)$:

$$\begin{aligned}
& \sqrt{\frac{1}{n} \sum_{i=1}^n (\langle w, x_i \rangle - \langle w', x_i \rangle)^2} && N(\mathcal{F}, L_2(\mathbb{P}_n), \epsilon) \\
& \leq \max_{i \in [n]} |\langle w - w', x_i \rangle| && \leq N(\mathcal{F}, L_\infty(\mathbb{P}_n), \epsilon) \\
& \leq \max_{x \in B_2^d} |\langle w - w', x \rangle| && \Leftrightarrow \leq N(\mathcal{F}, \|\cdot\|_\infty, \epsilon) \\
& \leq \|w - w'\|_2 && \leq N(B_2^d, \|\cdot\|_2, \epsilon) \\
& \leq \epsilon,
\end{aligned}$$

where $\max_{x \in B_2^d} |\langle w - w', x \rangle| \leq \|w - w'\|_2$ since $\|x\|_2 \leq 1$. Then, from [Proposition 3.3.1](#),

$$N(B_2^d, \|\cdot\|_2, \epsilon) \leq \left(1 + \frac{2}{\epsilon}\right)^d,$$

so we get a $\sqrt{d/n}$ rate since this satisfies **parametric** condition.^a

^aIn high dimension situation, this bound can be loose.

It turns out that one can use the norm constraints to also show that

$$\sup_{\mu} \log N(\mathcal{F}, L_2(\mathbb{P}_n), \epsilon) \leq \frac{C}{\epsilon^2},$$

i.e., a dimension-free bound, hence \mathcal{F} also behaves like a **non-parametric** class! This bound will be useful in high dimensional setting when d is not small compared to n . Thus, we make this an important remark.

Remark. A function class can be viewed as **parametric** and **non-parametric** at the same time.

There are other examples.

Example. Neural networks are like this: we can either measure its complexity by the *number of parameters*, or get dimension independent bounds on its **Rademacher complexity** by using *norm constraints*.

Problem. For what function classes can be bound in terms of uniform L_2 **entropy**?

Answer. We have seen boolean function classes with finite **VC dimension**, and **Hölder smooth function classes**. *

Lecture 16: Beyond VC Dimension: Fat-Shattering Dimension

3.3.9 Fat-Shattering Dimension

4 Oct. 9:00

To generalize **VC dimension**, consider the following.

Definition. Let \mathcal{F} be a real-valued function class on χ .

Definition 3.3.18 (ϵ -shattered). A set $\{x_1, \dots, x_n\}$ of χ is ϵ -shattered by \mathcal{F} if there exists t_1, \dots, t_n such that for all $S \subseteq [n]$, there exists $f \in \mathcal{F}$ such that

$$\begin{cases} f(x_s) \leq t_s, & \text{if } s \in S; \\ f(x_s) \geq t_s + \epsilon, & \text{if } s \notin S. \end{cases}$$

Definition 3.3.19 (Fat-shattering dimension). The *fat-shattering dimension* $VC(\mathcal{F}, \epsilon)$ of \mathcal{F} on χ is the maximum integer D such that there exists a size D finite set $A \subseteq \chi$ ϵ -shattered by \mathcal{F} .

Remark. $VC(\mathcal{F}, \epsilon)$ is a non-increasing function of ϵ .

Proof. If a set is ϵ -shattered then for any $\delta \leq \epsilon$, it's also δ -shattered. *

Note. If \mathcal{F} is boolean, then $VC(\mathcal{F}, \epsilon) = VC(\mathcal{F})$ for all $\epsilon \in (0, 1)$ by setting $t_s = 0$ for all s .

Example. Consider $\mathcal{M} = \{f: I \rightarrow [-1, 1] \text{ non-decreasing}\}$. Then for all $\epsilon > 0$, $VC(\mathcal{M}, \epsilon) \leq 1 + 2/\epsilon$.

Definition 3.3.20 (Total variation). The *total variation* of a function $f: \mathbb{R} \supseteq I \rightarrow \mathbb{R}$ on an interval I is defined as

$$TV(f) := \sup_{n \geq 1} \sup_{x_1 < \dots < x_n \in I} \sum_{i=1}^n |f(x_i) - f(x_{i-1})|.$$

Intuition. The *total variation* is some measure of smoothness of functions. It's more general than differentiability since $TV(f)$ can also be defined for discontinuous functions.

Remark. *Total variation* is actually a norm, i.e., triangle inequality holds.

Let's first see one example.

Example. Consider $BV(2) := \{f: TV(f) \leq 2\}$. If $f'(x)$ exists, then $TV(f) = \int |f'(x)| dx$. We see that $BV(2) \supseteq \mathcal{M}$ since for any non-decreasing function f ranging in $[a, b]$, $TV(f) = b - a$.

In general, we have the following.

Lemma 3.3.4. Let $\mathcal{F} \ni f: I \rightarrow \mathbb{R}$ with $TV(f) \leq v$ for all $f \in \mathcal{F}$, i.e., $\mathcal{F} = BV(v)$. Then

$$VC(BV(v), \epsilon) = 1 + \left\lfloor \frac{v}{\epsilon} \right\rfloor.$$

Proof. We prove this by proving two directions.

Claim. $VC(BV(v), \epsilon) \leq 1 + \lfloor v/\epsilon \rfloor$.

Proof. Let $\{x_1, \dots, x_n\}$ be ϵ -shattered by \mathcal{F} , then there exists t_1, \dots, t_n and f_1, f_2 such that

$$\begin{cases} f_1(x_i) \leq t_i, & \text{if } i \text{ is odd;} \\ f_1(x_i) \geq t_i + \epsilon, & \text{if } i \text{ is even;} \end{cases}, \quad \begin{cases} f_2(x_i) \leq t_i, & \text{if } i \text{ is even;} \\ f_2(x_i) \geq t_i + \epsilon, & \text{if } i \text{ is odd.} \end{cases}$$

Consider $f = (f_1 - f_2)/2$, then

$$\begin{cases} f(x_i) \leq -\epsilon/2, & \text{if } i \text{ is odd;} \\ f(x_i) \geq \epsilon/2, & \text{if } i \text{ is even;} \end{cases} \Rightarrow TV(f) \geq (n-1)\epsilon$$

by considering this particular partition $\{x_i\}$ of I . Furthermore, since TV is a norm,

$$TV(f) \leq \frac{TV(f_1) + TV(f_2)}{2} \leq v$$

from triangle inequality, hence $(n-1)\epsilon \leq v$, i.e., $n \leq 1 + \lfloor v/\epsilon \rfloor$. *

Claim. $\text{VC}(\text{BV}(v), \epsilon) \geq 1 + \lfloor v/\epsilon \rfloor$.

Proof. Let $d = \lfloor v/\epsilon \rfloor$, and consider $y_1 < y_2 < \dots < y_d$, which induces $d + 1$ intervals

$$I_0 = (-\infty, y_1), \quad I_j = [y_j, y_{j+1}), \quad I_d = [y_d, \infty)$$

Let \mathcal{G} be the set of piece-wise continuous functions on I_0, \dots, I_d taking values between $\{0, \epsilon\}$, so $|\mathcal{G}| = 2^{d+1}$. Then, the set

$$\{x_1, \dots, x_{d+1} : x_j \in I_{j-1}\}$$

is ϵ -shattered by \mathcal{G} . Finally, since $\text{TV}(g) \leq d\epsilon \leq v$ for all $g \in \mathcal{G}$, we're done. \otimes

■

Example (Linear function class). Let $\chi = B_2^d$, and $\mathcal{F} = \{x \mapsto w^\top x : w \in B_2^d\}$. Then $\text{VC}(\mathcal{F}, \epsilon) \leq d$; and if we consider $\text{sgn}(w^\top x)$, we get $\text{VC}(\mathcal{F}, \epsilon) = d + 1$.

Consider the following result (which we will not prove).

Theorem 3.3.6 (Mendelson-Vershynin). Let \mathcal{F} be a class of functions that is uniformly bounded by 1. Then there exists $c > 1$ such that for every $0 < \epsilon \leq 1$,

$$\sup_{\mu} M(\mathcal{F}, L_2(\mu), \epsilon) \leq \left(\frac{2}{\epsilon}\right)^{c \text{VC}(\mathcal{F}, \frac{\epsilon}{c})}.$$

Remark. For $\text{BV}(2)$, [Mendelson-Vershynin theorem](#) gives $\sup_{\mu} M(\text{BV}(2), \epsilon, L_2(\mu)) \leq \exp\left(\frac{c}{\epsilon} \log 2\epsilon\right)$.^a

^aNote that $\log 2\epsilon$ is superfluous again.

Lecture 17: Perceptron Algorithm

As previously seen. In the [previous example](#), we state that the [fat-shattering dimension](#) for a linear function class is $\text{VC}(\mathcal{F}, \epsilon) \leq d$, which is not dimension free. 1

6 Oct. 9:00

If we further impose a norm constraint $\|w\|_2$, then a dimension-free bound can be obtained; specifically, for all $\epsilon > 0$, we can show that

$$\text{VC}(\mathcal{F}, \epsilon) \leq \frac{C}{\epsilon^2}$$

where C is some constant. To prove the above, we can use the [perceptron algorithm](#).

Algorithm 3.1: Perceptron Algorithm

Data: A data sequence $\{(x_i, y_i)\}_{i=1}^T$, observed one-by-one

Result: A linear function with weight w

```

1  $\hat{w}_1 \leftarrow 0$ 
2 for  $t = 1, \dots, T$  do
3   observe  $x_t \in \chi$                                      // data
4    $\hat{y}_t \leftarrow \text{sgn}(\hat{w}_t^\top x_t)$                        // predict  $\hat{y}_t$ 
5   observe  $y_t \in \{\pm 1\}$                                 // true label
6    $\hat{w}_{t+1} \leftarrow \hat{w}_t + \mathbb{1}_{\hat{y}_t \neq y_t} y_t x_t$ 
7 return  $\hat{w}_{T+1}$ 
```

Remark. Suppose there exists w such that $y_t = 1$ whenever $w^\top x_t \geq 0$, and $y_t = -1$ whenever $w^\top x_t < 0$, then $y_t w^\top x_t > 0$.

The following lemma (see, e.g., [\[Nov62\]](#)) provides an error bound for the [perceptron algorithm](#).

Lemma 3.3.5 (Perceptron Mistake Bound). For any sequence $(x_1, y_1), \dots, (x_T, y_T) \in B_2^d \times \{\pm 1\}$, the [perceptron algorithm](#) makes at most $1/\gamma^2$ mistakes, where

$$\gamma = \max_{w \in B_2^d} \min_{1 \leq t \leq T} y_t w^\top x_t$$

is the margin of the data sequence $\{(x_i, y_i)\}_{i=1}^T$.

Proof. Let M be the total number of mistakes made when running the [perceptron algorithm](#). Suppose a mistake is made in round t , then

$$\|\hat{w}_{t+1}\|^2 = \|\hat{w}_t + y_t x_t\|^2 = \|\hat{w}_t\|^2 + 2 \underbrace{\langle \hat{w}_t, y_t x_t \rangle}_{\leq 0} + \|x_t\|^2 \leq \|\hat{w}_t\|^2 + 1,$$

implying that $\|\hat{w}_{t+1}\| \leq \sqrt{M}$.

Now, consider the margin γ : when γ is achieved at $w = w^*$, $\gamma \leq y_t (w^*)^\top x_t$; moreover, if there is a mistake at round t ,

$$\gamma \leq y_t (w^*)^\top x_t = (w^*)^\top (\hat{w}_{t+1} - \hat{w}_t)$$

since $y_t x_t = \hat{w}_{t+1} - \hat{w}_t$ in this case. Summing the above over t results in a telescoping sum

$$M\gamma \leq (w^*)^\top \hat{w}_{T+1} \leq \|w^*\| \|\hat{w}_{T+1}\| \leq \|\hat{w}_{T+1}\| \leq \sqrt{M},$$

hence $M \leq 1/\gamma^2$. ■

We are now ready to prove the following.

Theorem 3.3.7. Let $\chi = B_2^d$, and $\mathcal{F} = \{x \mapsto w^\top x : w \in B_2^d\}$. Then for all $\epsilon > 0$,

$$\text{VC}(\mathcal{F}, \epsilon) \leq \frac{4}{\epsilon^2}.$$

Proof. Suppose $\{x_1, \dots, x_T\}$ is [\$\epsilon\$ -shattered](#) by \mathcal{F} . Then, there exists $t_1, \dots, t_n \in [-1, 1]$ such that for all $S \subset \{1, \dots, n\}$, there exists $w_S \in B_2^d$ such that

$$\begin{cases} w_S^\top x_i \geq t_i + \frac{\epsilon}{2}, & \text{if } i \in S; \\ w_S^\top x_i \leq t_i - \frac{\epsilon}{2}, & \text{if } i \notin S. \end{cases}$$

Let $\tilde{x}_i = (x_i, t_i) \in \mathbb{R}^{d+1}$ and $\tilde{w}_S = (w_S, -1) \in \mathbb{R}^{d+1}$, we can rewrite the above as

$$\begin{cases} \tilde{w}_S^\top \tilde{x}_i \geq \frac{\epsilon}{2}, & \text{if } i \in S; \\ \tilde{w}_S^\top \tilde{x}_i \leq -\frac{\epsilon}{2}, & \text{if } i \notin S. \end{cases}$$

Equivalently, for any sign vector $y \in \{\pm 1\}^T$, there exists \tilde{w}_y such that for all i , $y_i \tilde{w}_y^\top \tilde{x}_i \geq \epsilon/2$, i.e., the margin $\gamma \geq \epsilon/2$. This means that if we run the [perceptron algorithm](#) with $\{\tilde{x}_i, y_i\}_{i=1}^T$ such that

$$y_i = -\hat{y}_i = \text{sgn}(\hat{w}_t^\top x_t),$$

i.e., the [perceptron algorithm](#) makes mistake every round, i.e., $M = T$. But since $\gamma \geq \epsilon/2$, $M \leq 4/\epsilon^2$ from the [perceptron mistake bound](#). Combining these two, we have $T = M \leq 4/\epsilon^2$ if a size T subset of χ is [\$\epsilon\$ -shattered](#) by \mathcal{F} , i.e., $\text{VC}(\mathcal{F}, \epsilon) \leq 4/\epsilon^2$. ■

Remark. We can try to use the above together with the [Mendelson-Vershynin theorem](#).

Remark. Suppose f_1, \dots, f_d are linearly independent. If $\mathcal{F} = \{\sum_{i=1}^d c_i f_i\}$, $\text{VC}(\mathcal{F}, \epsilon) = d$ for all $\epsilon > 0$.

Lecture 18: Beyond Uniform Entropy Bound: Bracketing Bound

3.4 Bracketing Bound

11 Oct. 9:00

As previously seen. So far, we have the [uniform entropy bound](#)

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{P}_n f - \mathbb{P} f| \right] \leq \frac{c}{\sqrt{n}} \|F\|_{L_2(\mathbb{P})} \int_0^1 \sqrt{\log \sup_{\mu} N(\mathcal{F}, L_2(\mu), \epsilon \|F\|_{L_2(\mathbb{P})})} dx,$$

where F is an [envelope](#) of \mathcal{F} .

In this lecture, we will see another bound using bracketing (recall [Proposition 3.1.1](#)).

3.4.1 Bracketing Number

Consider the following.

Definition 3.4.1 (ϵ -bracket). Given a probability measure \mathbb{P} on χ and two functions $\ell, u: \chi \rightarrow \mathbb{R}$, an ϵ -bracket, denoted as $[\ell, u]$, is defined as

$$[\ell, u] := \{f: \chi \rightarrow \mathbb{R}: \ell(x) \leq f(x) \leq u(x) \text{ for all } x \in \chi\}$$

such that $\|u - \ell\|_{L_2(\mathbb{P})} \leq \epsilon$.^a

^aThis is the $L_2(\mathbb{P})$ size of $[\ell, u]$, i.e., $\left(\int_{\chi} (u(x) - \ell(x))^2 \mathbb{P}(dx) \right)^{1/2}$.

Definition 3.4.2 (Bracketing number). For every $\epsilon > 0$, the ϵ -bracketing number $N_{[\cdot]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon)$ of a function class \mathcal{F} from χ to \mathbb{R} is defined as the smallest number of ϵ -brackets such that every $f \in \mathcal{F}$ belongs to only one of the brackets.

Lemma 3.4.1. For every $\epsilon > 0$, $N(\mathcal{F}, L_2(\mathbb{P}), \epsilon/2) \leq N_{[\cdot]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon)$.

Proof. Consider ϵ -brackets $[\ell_i, u_i]$ for $i = 1, \dots, N_{[\cdot]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon)$, then $\{(\ell_i + u_i)/2\}_i$ forms an $\epsilon/2$ -net since for any $f \in \mathcal{F}$ and any $x \in \chi$,

$$\left\| f - \frac{u_i + \ell_i}{2} \right\|_{L_2(\mathbb{P})} \leq \left\| \frac{u_i - \ell_i}{2} \right\|_{L_2(\mathbb{P})} \leq \frac{\epsilon}{2}$$

from the fact that $u_i \geq f \geq \ell_i$ and $\|u_i - \ell_i\|_{L_2(\mathbb{P})} \leq \epsilon$. ■

Let's see one simple example of computing [bracketing functions](#).

Example. Let $\mathcal{F} = \{\mathbb{1}_{[-\infty, t]}: t \in \mathbb{R}\}$ and \mathbb{P} be a probability measure on \mathbb{R} . Then for all $\epsilon > 0$,

$$N_{[\cdot]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon) \leq 1 + 1/\epsilon^2.$$

Proof. We show this by finding a collection of $\sqrt{\epsilon}$ -brackets with at most $1 + 1/\epsilon$ many of the brackets. Let $t_0 = -\infty$, and recursively define

$$t_i = \sup\{x: x > t_{i-1}: \mathbb{P}((t_{i-1}, x]) \leq \epsilon\}.$$

Finally, let $k \geq 1$ be the smaller integer such that $t_k = \infty$. We then have

- $\mathbb{P}((t_{i-1}, t_i)) \leq \epsilon$: for every $\delta > 0$, $\mathbb{P}((t_{i-1}, t_i - \delta]) \leq \epsilon$, as $\delta \rightarrow 0$, $\mathbb{P}((t_{i-1}, t_i)) \leq \epsilon$.
- $\mathbb{P}((t_{i-1}, t_i]) \geq \epsilon$: for every $\delta > 0$, $\mathbb{P}((t_{i-1}, t_i + \delta]) > \epsilon$, as $\delta \rightarrow 0$, $\mathbb{P}((t_{i-1}, t_i]) \geq \epsilon$.

Then,

$$1 = \mathbb{P}((-\infty, \infty)) \geq \sum_{i=1}^k \mathbb{P}((t_{i-1}, t_i]) \geq (k-1)\epsilon,$$

implying $k \leq 1 + \frac{1}{\epsilon}$. Now, consider **brackets** $[\mathbb{1}_{(-\infty, t_{i-1})}, \mathbb{1}_{(-\infty, t_i)}]$ which cover \mathcal{F} with $L_2(\mathbb{P})$ size equal to $\sqrt{\mathbb{P}((t_{i-1}, t_i))} \leq \sqrt{\epsilon}$. Hence, this is a collection of valid **$\sqrt{\epsilon}$ -brackets** of size $\leq 1 + 1/\epsilon$, i.e., by replacing $\epsilon \leftarrow \sqrt{\epsilon}$, we have $N_{[\cdot]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon) \leq 1 + 1/\epsilon^2$. \circledast

Proposition 3.4.1. Let \mathcal{F} to be a function class such that $N_{[\cdot]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon) < \infty$ for all $\epsilon > 0$. Then as $n \rightarrow \infty$,

$$\sup_{f \in \mathcal{F}} |\mathbb{P}_n f - \mathbb{P} f| \xrightarrow{\text{a.s.}} 0.$$

Proof. Fix $\epsilon > 0$, let $[\ell_i, u_i]$ for $i = 1, \dots, N$ to be a set of **ϵ -brackets**. Then, it suffices to show^a

$$\sup_{f \in \mathcal{F}} |\mathbb{P}_n f - \mathbb{P} f| \leq \left(\max_{1 \leq i \leq N} \max(|\mathbb{P}_n u_i - \mathbb{P} u_i|, |\mathbb{P}_n \ell_i - \mathbb{P} \ell_i|) \right) + \epsilon.$$

To show this, let $f \in [\ell_i, u_i]$ for some i , then

$$\mathbb{P}_n f - \mathbb{P} f \leq (\mathbb{P}_n u_i - \mathbb{P} u_i) + (\mathbb{P} u_i - \mathbb{P} f) \leq (\mathbb{P}_n u_i - \mathbb{P} u_i) + \mathbb{P}(u_i - \ell_i) \leq \mathbb{P}_n u_i - \mathbb{P} u_i + \epsilon$$

since $\mathbb{P}(u_i - \ell_i) \leq \|u_i - \ell_i\|_{L_2(\mathbb{P})} \leq \epsilon$. On the other hand, we also have

$$\mathbb{P} f - \mathbb{P}_n f \leq (\mathbb{P} f - \mathbb{P} \ell_i) + (\mathbb{P} \ell_i - \mathbb{P}_n \ell_i) \leq (\mathbb{P} u_i - \mathbb{P} \ell_i) + (\mathbb{P} \ell_i - \mathbb{P}_n \ell_i) \leq |\mathbb{P}_n \ell_i - \mathbb{P} \ell_i| + \epsilon,$$

hence we're done. \blacksquare

^aIt then implies $\limsup_{n \rightarrow \infty} |\mathbb{P}_n f - \mathbb{P} f| \leq \epsilon$ almost surely just by the law of large number. By taking $\epsilon = 1/m$ to 0, we can say that $|\mathbb{P}_n f - \mathbb{P} f| \rightarrow 0$ almost surely.

3.4.2 Bracketing Bound

The main theorem of this section is the following.

Theorem 3.4.1 (Bracketing bound). Let F be an **envelope** of \mathcal{F} such that $\mathbb{P} F^2 < \infty$. Then for some constant $C > 0$,

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \sqrt{n} (\mathbb{P}_n f - \mathbb{P} f) \right] \leq C \|F\|_{L_2(\mathbb{P})} \int_0^1 \sqrt{\log N_{[\cdot]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon \|F\|_{L_2(\mathbb{P})})} d\epsilon.$$

Remark. The main differences between the **bracketing bound** and the **uniform entropy bound** are

- **covering number** is replaced by **bracketing number**;
- We do not have the \sup_{μ} , hence the **bracketing bound** is only w.r.t. \mathbb{P} .

Lecture 19: Applications to M -Estimators

The following shows that using the **bracketing number** is more tractable than using the uniform L_2 **covering number**. 13 Oct. 9:00

Lemma 3.4.2 (Parametric Lipschitz function class). Let $\Theta \subseteq \mathbb{R}^d$ be constrained in an L_2 ball of radius R , and let $\mathcal{F} = \{m_{\theta} : \chi \rightarrow \mathbb{R} : \theta \in \Theta\}$ be a function class indexed by θ . Suppose there exists a

function $M(x)$ with $\|M\|_{L_2(\mathbb{P})} < \infty$ such that

$$|m_{\theta_1}(x) - m_{\theta_2}(x)| \leq M(x)\|\theta_1 - \theta_2\|_2$$

for all $x \in \mathcal{X}$ and $\theta_1, \theta_2 \in \Theta$. Then, for all $\epsilon > 0$,

$$N_{[]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon \|M\|_{L_2(\mathbb{P})}) \leq \left(1 + \frac{4R}{\epsilon}\right)^d.$$

Proof. Let $\{\theta_i\}_{i=1}^N$ be a maximal $\epsilon/2$ -packing of Θ , and consider the following brackets:

Claim. The brackets $[m_{\theta_i} \pm \epsilon M/2]$ for $i = 1, \dots, N$ cover \mathcal{F} .

Proof. First, note that a maximal packing set is indeed a covering net with the same ϵ . Therefore, for all $m_\theta \in \mathcal{F}$, there exists i such that $\|\theta - \theta_i\|_2 \leq \epsilon/2$, hence

$$|m_\theta(x) - m_{\theta_i}(x)| \leq M(x)\|\theta - \theta_i\|_2 \leq M(x)\frac{\epsilon}{2},$$

for all $x \in \mathcal{X}$, i.e., $m_\theta \in [m_{\theta_i} \pm \epsilon M/2]$. This means the brackets $[m_{\theta_i} \pm \epsilon M/2]$ cover \mathcal{F} . \otimes

Furthermore, the size of each bracket is $\epsilon \|M\|_{L_2(\mathbb{P})}$, with Proposition 3.3.1, we're done. \blacksquare

Some examples of the parametric Lipschitz function classes are the following.

Example. For $m_\theta(x) = \theta^\top x$, $|\theta_1^\top x - \theta_2^\top x| \leq \|\theta_1 - \theta_2\|_2 \|x\|_2$ from Cauchy-Schwarz. Hence, $M(x) = \|x\|_2$. For Lemma 3.4.2 to apply, consider \mathbb{P} such that $\mathbb{P}\|x\|_2 < \infty$.

Example (Quantile regression). For $m_\theta(x) = |x - \theta|$, $|m_{\theta_1}(x) - m_{\theta_2}(x)| \leq |\theta_1 - \theta_2|$ with $M(x) = 1$. In this case, since any measure \mathbb{P} gives $\mathbb{P}1 < \infty$, so Lemma 3.4.2 applies for all \mathbb{P} .

The above examples extends to essentially all p -norm.

Example. $m_\theta(x) = \|x - \theta\|_p$.

Finally, let's see some standard result on the bracketing numbers.

Example (α -Hölder smooth function class). For \mathcal{S}_α on $[0, 1] \rightarrow [0, 1]$,

$$\log N_{[]}(\mathcal{S}_\alpha, L_2(\mathbb{P}), \epsilon) \leq C \left(\frac{1}{\epsilon}\right)^\alpha.$$

Example. Let \mathcal{M} be the monotone function class on $\mathbb{R} \rightarrow [0, 1]$,

$$\log N_{[]}(\mathcal{M}, L_2(\mathbb{P}), \epsilon) \leq C \left(\frac{1}{\epsilon}\right).$$

Example. Consider \mathcal{C} be the set of convex function class on $[0, 1] \rightarrow [0, 1]$. Let \mathcal{U} be the uniform distribution on $[0, 1]$. Then

$$\log N_{[]}(\mathcal{C}, L_2(\mathcal{U}), \epsilon) \leq C \left(\frac{1}{\sqrt{\epsilon}}\right).$$

More examples are available in [VW96]. To conclude this section, we ask the following:

Problem 3.4.1 (Necessity of VC dimension of boolean function class). Is there a function class for which the uniform entropy bound is infinite, while the bracketing bound is finite?

Answer. Yes! For boolean function class \mathcal{F} , from the [Dudley's theorem](#),

- $\text{VC}(\mathcal{F}) < \infty$: $\sup_{\mu} N(\mathcal{F}, L_2(\mu), \epsilon) < \infty$;
- $\text{VC}(\mathcal{F}) = \infty$: $\sup_{\mu} N(\mathcal{F}, L_2(\mu), \epsilon) = \infty$ for all $\epsilon < 1/2$.

In addition, if $\text{VC}(\mathcal{F}) = \infty$, uniform [Glivenko-Cantelli](#) does not hold, i.e.,

$$\liminf_{n \rightarrow \infty} \sup_{\mathbb{P}} \mathbb{E} \left[\sup_{f \in \mathcal{F}} \mathbb{P}_n f - \mathbb{P} f \right] > 0.$$

However, it's still possible that when $\text{VC}(\mathcal{F}) = \infty$, \mathcal{F} is [Glivenko-Cantelli](#) w.r.t. some \mathbb{P} . ⊗

Example. Let $\mathcal{F} = \{\mathbb{1}_C : C \text{ is compact convex subset of } [0, 1]^d\}$. Then $\text{VC}(\mathcal{F}) = \infty$, and \mathcal{F} is [Glivenko-Cantelli](#) w.r.t. any \mathbb{P} having a density w.r.t. the Lebesgue measure.

Lemma 3.4.3. Let \mathcal{F} be a class of function uniformly bounded by 1. Then for every $\epsilon > 0$,

$$\frac{1}{8} \text{VC}(\mathcal{F}, 4\epsilon) \leq \log \sup_{\mu} N(\mathcal{F}, L_2(\mu), \epsilon)$$

Chapter 4

Applications to M -Estimation

In this chapter, we will focus on M -estimator, and primarily investigate the “rate of convergence” for M -estimators, and look at some examples.

4.1 The M -Estimation Problem

Consider the problem of M -estimation, which formalize subsection 1.2.1:

Problem 4.1.1 (M -estimation). Let Θ be an abstract parameter space,^a and let $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$ be the data. Let $M: \Theta \rightarrow \mathbb{R}$ and $M_n: \Theta \rightarrow \mathbb{R}$ be random functions^b depend on the data. Then, the M -estimation problem tries to estimate the true parameter

$$\theta_0 = \arg \max_{\theta \in \Theta} M(\theta)$$

by minimizing $M_n(\theta)$ instead and find

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} M_n(\theta).$$

^aE.g., \mathbb{R}^d , or some function spaces.

^bOr equivalently, one can view $\{M_n\}_n$ as a stochastic process.

Remark. Typically, for each fixed $\theta \in \Theta$, we have $M_n(\theta) \xrightarrow{P} M(\theta)$.

We have seen some examples in the beginning of the class (subsection 1.2.1).

Example (§1.2.1). Mean, Quantile, and Mode estimation.

Example (Least square estimation). Let $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$, and let $\Theta = \mathcal{F}$. Consider $M(f) = -\mathbb{P}(y - f(x))^2$ and $M_n(f) = -\mathbb{P}_n(y - f(x))^2$.

Example (ERM in classification). Let $M(f) = -\mathbb{P}(y \neq \text{sgn}(f(x)))$ and $M_n(f) = -\mathbb{P}_n(y \neq \text{sgn}(f(x)))$.

Example (MLE). Let $M(\theta) = p \log p_\theta$ and $M_n(\theta) = p_n \log p_\theta$, where $\{p_\theta\}_{\theta \in \Theta}$ is a class of densities w.r.t. some measure.

It only makes sense to look at those M -esimator that are consistent.

Definition 4.1.1 (Consistent). An M -esimator $\hat{\theta}$ is *consistent* if $|M(\hat{\theta}_n) - M(\theta_0)| \xrightarrow{P} 0$ implies $d(\hat{\theta}_n, \theta_0) \xrightarrow{P} 0$ as $n \rightarrow \infty$.

Then, we can ask the following three questions (progressively harder) for an M -estimator.

Problem 4.1.2 (Consistency). Is the M -estimator consistent? I.e., as $|M(\hat{\theta}_n) - M(\theta_0)| \xrightarrow{P} 0$, does $d(\hat{\theta}_n, \theta_0) \xrightarrow{P} 0$ for some metric d ?

Problem 4.1.3 (Rate of convergence). What's the "rate of convergence" for the M -estimator?

Example. The rate for the mean is $O(1/\sqrt{n})$.

We will define the "rate of convergence" precise later (Definition 4.3.1).

Problem 4.1.4 (Limiting distribution). What's the limiting (asymptotic) distribution of $\hat{\theta} - \theta_0$?

4.1.1 Running Example

In the remaining class, we will consider the [sample mode estimation](#) as our running example.

Example (Mode estimation). Let $\chi = \Theta = \mathbb{R}$, and suppose we have $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_{\theta_0}$ supported on χ such that \mathbb{P}_{θ_0} has smooth and bounded density $p_{\theta_0}(x)$ w.r.t. Lebesgue measure, with $p'_{\theta_0}(x) > 0$ for $x < \theta_0$, and $p'_{\theta_0}(x) < 0$ for $x > \theta_0$. The *mode estimation* problem considers

$$M(\theta) = \mathbb{P}_{\theta_0}(\theta - 1 \leq X \leq \theta + 1),$$

so the true parameter $\theta_0 = \arg \max_{\theta \in \Theta} M(\theta)$ is the mode. Let

$$M_n(\theta) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\theta-1 \leq X_i \leq \theta+1},$$

and let $\hat{\theta}_n = \arg \max_{\theta \in \Theta} M_n(\theta)$, i.e., the sample mode.

Remark (Unique optimal). Notice that in the [mode estimation](#) problem, we can check that

- $M'(\theta_0) = 0$;
- $M'(\theta) < 0$ when $\theta < \theta_0$, and $M'(\theta) > 0$ when $\theta > \theta_0$;
- $M''(\theta) > 0$.

These conditions guarantee that θ_0 is the unique maximum.

4.2 Consistency

[Consistency](#) is the easiest task among others. Firstly, we need to show that $|M(\hat{\theta}_n) - M(\theta_0)| \xrightarrow{P} 0$: recall the "basic inequality" step, i.e.,

$$|M(\hat{\theta}_n) - M(\theta_0)| = |M(\hat{\theta}_n) - M_n(\hat{\theta}_n) + \underbrace{M_n(\hat{\theta}_n) - M_n(\theta_0)}_{\leq 0} + M_n(\theta_0) - M(\theta_0)| \leq 2 \sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)|.$$

In the [mode estimation](#) example, denote $m_{\theta}(x) = \mathbb{1}_{\theta-1 \leq x \leq \theta+1}$, and $\mathcal{F} = \{m_{\theta} : \theta \in \mathbb{R}\}$, then consider the following notation.

Notation. $M_n(\theta) = \mathbb{P}_n m_{\theta}$ and $M(\theta) = \mathbb{P} m_{\theta}$.

We then have

$$\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| = \sup_{\theta \in \mathbb{R}} |\mathbb{P}_n m_{\theta} - \mathbb{P} m_{\theta}| \xrightarrow{P} 0$$

since $\text{VC}(\mathcal{F}) = 2$ (see [the previous example](#)). Hence, we conclude that $M(\hat{\theta}) - M(\theta_0) \xrightarrow{P} 0$ as $n \rightarrow \infty$. Now, it remains to answer the following.

Problem. Does $\hat{\theta} \xrightarrow{P} \theta_0$, i.e., does $d(\hat{\theta}, \theta_0) \xrightarrow{P} 0$?

Answer. Need to relate d to M function. *

Lecture 20: A Heuristic Argument for Rate of Convergence

To show $\hat{\theta} \xrightarrow{P} \theta_0$, we need “curvature” or “growth” condition on M at θ_0 .

16 Oct. 9:00



Consider the following.

Definition 4.2.1 (Growth condition). The *growth condition* on $M(\theta)$ is that for all $\epsilon > 0$,

$$\sup_{\theta: d(\theta, \theta_0) > \epsilon} M(\theta) < M(\theta_0).$$

Equivalently, the [growth condition](#) implies that for all $\epsilon > 0$, there exists $\delta > 0$ such that $\delta < \epsilon$ and

$$\inf_{\theta: d(\theta, \theta_0) > \epsilon} M(\theta_0) - M(\theta) > \delta.$$

Hence, as $\epsilon \rightarrow 0$, $\mathbb{P}(d(\hat{\theta}, \theta_0) > \epsilon) \leq \mathbb{P}(M(\theta_0) - M(\hat{\theta}) > \delta) \rightarrow 0$. In our [mode estimation](#) example, since

$$M(\theta) = \mathbb{P}(\theta - 1 \leq X \leq \theta + 1) = \int_{\theta-1}^{\theta+1} p_{\theta_0}(x) dx,$$

if we assume that p_{θ_0} is increasing until θ_0 and decreasing after θ_0 ,¹ one can check that

$$\sup_{\theta: d(\theta, \theta_0) > \epsilon} M(\theta) = \max(M(\theta_0 + \epsilon), M(\theta_0 - \epsilon)) < M(\theta_0)$$

from $p'' < 0$ at θ_0 . More generally, one can refer to the following [[Vaa98](#), Theorem 5.7].

Theorem 4.2.1. Let $\hat{\theta}$ be any minimizer of $M_n(\theta)$, and θ_0 be the unique minimizer of $M(\theta)$. If

- (a) $\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{P} 0$;
- (b) for all $\epsilon > 0$, $\inf_{\theta: d(\theta, \theta_0) > \epsilon} M(\theta) > M(\theta_0)$ for some metric d on Θ ,

then we have $d(\hat{\theta}, \theta_0) \xrightarrow{P} 0$.

Remark. Condition (a) may be too strong. One may need a preliminary “localization”.

4.3 Rate of Convergence

Let’s first see some heuristics before introducing the general theory.

¹Since θ_0 is the unique maximum.

4.3.1 The Heuristic Argument

In this section, we'll first show a heuristic argument for bounding the [rate of convergence](#) in the case of our running example, i.e., [mode estimation](#), where $d(\theta, \theta') = |\theta - \theta'|$.

As previously seen. Recall the previously defined [notation](#) for the [mode estimation](#) example, i.e., $M_n(\theta) = \mathbb{P}_n m_\theta$ and $M(\theta) = \mathbb{P} m_\theta$ where $m_\theta(x) = \mathbb{1}_{\theta-1 \leq x \leq \theta+1}$.

We then see that in the [mode estimation](#) example,

$$M(\theta_0) - M(\hat{\theta}_n) \leq (\mathbb{P}_n - \mathbb{P})(m_{\hat{\theta}_n} - m_{\theta_0}) \leq 2 \sup_{\theta \in \Theta} |\mathbb{P}_n m_\theta - \mathbb{P} m_\theta|,$$

and we can further upper-bound it by c/\sqrt{n} in expectation, i.e., in expectation,

$$M(\theta_0) - M(\hat{\theta}_n) = O_p(1/\sqrt{n}).$$

Now, our goal is to somehow relate $|\hat{\theta}_n - \theta_0|$ with $M(\theta_0) - M(\hat{\theta}_n)$.

Note. We cannot in general, get a rate better than $1/\sqrt{n}$ for this suprema.

Remark. If M is twice differentiable at θ_0 with $M''(\theta_0) < 0$, then there exists a neighborhood I of θ_0 such that for all $\theta \in I$,

$$M(\theta_0) - M(\theta) \geq c \cdot |\theta - \theta_0|^2.$$

Proof. From Taylor expansion, there exists ξ between θ and θ_0 such that

$$M(\theta) - M(\theta_0) = |\theta - \theta_0| M'(\theta_0) + \frac{|\theta - \theta_0|^2}{2} M''(\xi),$$

with $M'(\theta_0) = 0$, $M(\theta) - M(\theta_0) > c \cdot |\theta - \theta_0|^2$ for some c . ⊗

Since $\hat{\theta}_n$ is [consistent](#), we can assume $\hat{\theta}_n \in I$. The upshot is that we want some sorts of “growth condition” (different from [Definition 4.2.1](#))² not for all $\theta \in \mathbb{R}$, but for θ in a fixed neighborhood I of θ_0 .

Intuition. Since $\hat{\theta}_n$ is [consistent](#), the [growth condition](#) on a ball around θ_0 should suffice.

In this case, the [growth condition](#) relating $|\hat{\theta}_n - \theta_0|$ and $M(\theta_0) - M(\hat{\theta}_n)$ is

$$|\hat{\theta}_n - \theta_0|^2 \lesssim M(\theta_0) - M(\hat{\theta}_n).$$

Then by actually showing $M(\theta_0) - M(\hat{\theta}_n) = O_p(1/\sqrt{n})$, we can conclude $|\hat{\theta}_n - \theta_0| = O_p(n^{-1/4})$.

Note. In our example of [mode estimation](#), M'' is constant because p'_{θ_0} is constant, and

$$M''(\theta_0) = p'_{\theta_0}(\theta_0 + 1) - p'_{\theta_0}(\theta_0 - 1) < 0$$

since $p'(\theta_0 + 1) < 0$ and $p'(\theta_0 - 1) > 0$.

Specifically, we now have

$$|\hat{\theta}_n - \theta_0|^2 \lesssim M(\theta_0) - M(\hat{\theta}_n) \leq (\mathbb{P}_n - \mathbb{P})(m_{\hat{\theta}_n} - m_{\theta_0}) \leq 2 \sup_{\theta \in \Theta} |\mathbb{P}_n m_\theta - \mathbb{P} m_\theta| \leq \frac{c}{\sqrt{n}} = O_p\left(\frac{1}{\sqrt{n}}\right)$$

in expectation.

Remark (The right rate). We cannot get better than $O_p(n^{-1/4})$ with this argument. But there exist problems where the rate is better (need to do “localization”!)

- In our example of [mode estimation](#), the right rate is $O_p(n^{-1/3})$.

²This is precisely defined in [Definition 4.3.2](#).

- For “parametric” problems (e.g., mean/quantile estimation) the right rate is $O_p(1/\sqrt{n})$.

Now comes the heuristic part. Let’s first compute the bound for $\mathbb{E}[(\mathbb{P}_n - \mathbb{P})(m_\theta - m_{\theta_0})]$ for a fixed θ close to θ_0 . We have

$$\begin{aligned}\mathbb{E}[(\mathbb{P}_n - \mathbb{P})(m_\theta - m_{\theta_0})] &\leq \sqrt{\text{Var}[\mathbb{P}_n(m_\theta - m_{\theta_0})]} \\ &= \frac{1}{\sqrt{n}} \sqrt{\text{Var}[m_\theta(X_1) - m_{\theta_0}(X_1)]} \leq \frac{1}{\sqrt{n}} \sqrt{\mathbb{E}[(m_\theta(X_1) - m_{\theta_0}(X_1))^2]}.\end{aligned}$$

In our **mode estimation** problem, if θ is close to θ_0 (with $\theta < \theta_0$),

$$m_\theta(x) - m_{\theta_0}(x) = \mathbb{1}_{\theta-1 \leq x \leq \theta_0+1} + \mathbb{1}_{\theta+1 \leq x \leq \theta_0+1},$$

hence

$$\mathbb{E}[(m_\theta(X_1) - m_{\theta_0}(X_1))^2] \leq \mathbb{P}(\theta - 1 \leq X_1 \leq \theta_0 - 1) + \mathbb{P}(\theta + 1 \leq X_1 \leq \theta_0 + 1) \leq 2p_{\theta_0}(\theta_0) \cdot |\theta - \theta_0|,$$

and since p_{θ_0} is bounded, we finally have

$$\mathbb{E}[(m_\theta(X_1) - m_{\theta_0}(X_1))^2] \lesssim |\theta - \theta_0|,$$

which implies

$$(\mathbb{P}_n - \mathbb{P})(m_\theta - m_{\theta_0}) = O_p\left(\sqrt{\frac{|\theta - \theta_0|}{n}}\right).$$

Intuition (Heuristic). Heuristically, we might want to conclude that

$$(\mathbb{P}_n - \mathbb{P})(m_{\hat{\theta}_n} - m_{\theta_0}) = O_p\left(\sqrt{\frac{|\hat{\theta}_n - \theta_0|}{n}}\right),$$

which is a better bound than before.

If this is true, then the overall bound becomes

$$|\hat{\theta}_n - \theta_0|^2 \leq O_p\left(\sqrt{\frac{|\hat{\theta}_n - \theta_0|}{n}}\right).$$

Canceling $\sqrt{|\hat{\theta}_n - \theta_0|}$ from both sides,

$$|\hat{\theta}_n - \theta_0|^{3/2} = O_p\left(\frac{1}{\sqrt{n}}\right) \Rightarrow |\hat{\theta}_n - \theta_0| = O_p(n^{-1/3}),$$

which is the correct rate for $\hat{\theta}_n - \theta_0$.

Remark. In fact, the limiting distribution is also known where we have

$$n^{1/3}(\hat{\theta}_n - \theta_0) \xrightarrow{d} \arg \max_{h \in \mathbb{R}} aB_h - bh^2,$$

where a, b are constants depend on p .

Lecture 21: The General Argument for Rate of Convergence

4.3.2 A General Approach

18 Oct. 9:00

We’re now going to show the general argument for bounding the **rate of convergence**. In this section, we will assume that our **M-estimation problem** is defined for minimum rather than maximum.³

³This can be done without loss of generality since we can simply add a negative sign for M and M_n .

Note. Hence, $M(\theta) \geq M(\theta_0)$ for all $\theta \in \Theta$ now.

Definition 4.3.1 (Rate of convergence). The *rate of convergence* for $\hat{\theta}_n$ is defined as the sequence $\{\delta_n\}$ such that $d(\hat{\theta}_n, \theta_0) = O_p(\delta_n)$.

Recall that instead of using the old [growth condition](#) used when showing [consistency](#), we need an alternative form. Consider the following.

Definition 4.3.2 (Growth condition*). The *growth condition* on $M(\theta)$ is that for all $\theta \in \Theta$,

$$d(\theta, \theta_0)^2 \leq M(\theta) - M(\theta_0).$$

Note. For such [growth condition](#), the canonical choice of d is $d(\theta, \theta_0) = \sqrt{M(\theta_0) - M(\theta)}$.

It suffices to show that given $\epsilon > 0$, there exists M (not depending on n) such that for all n ,

$$\mathbb{P}(d(\hat{\theta}_n, \theta_0) > 2^M \delta_n) \leq \epsilon.$$

Firstly, we apply the *peeling step* to get

$$\mathbb{P}(d(\hat{\theta}_n, \theta_0) > 2^M \delta_n) = \sum_{j>M} \mathbb{P}(2^{j-1} \delta_n < d(\hat{\theta}_n, \theta_0) \leq 2^j \delta_n). \quad (4.1)$$

Note. From the [growth condition](#) and the basic inequality,

$$\begin{aligned} d(\hat{\theta}_n, \theta_0)^2 &\leq M(\hat{\theta}_n) - M(\theta_0) \\ &\leq M(\hat{\theta}_n) - M_n(\hat{\theta}_n) + \underbrace{M_n(\hat{\theta}_n) - M_n(\theta_0)}_{\leq 0} + M_n(\theta_0) - M(\theta_0) \\ &\leq (M_n - M)(\hat{\theta}_n) - (M_n - M)(\theta_0). \end{aligned}$$

With this, we can further upper-bound [Equation 4.1](#) by

$$\begin{aligned} \mathbb{P}(d(\hat{\theta}_n, \theta_0) > 2^M \delta_n) &= \sum_{j>M} \mathbb{P}(2^{j-1} \delta_n < d(\hat{\theta}_n, \theta_0) \leq 2^j \delta_n) && \text{peeling step} \\ &\leq \sum_{j>M} \mathbb{P}((M_n - M)(\theta_0) - (M_n - M)(\hat{\theta}_n) \geq 2^{2j-2} \delta_n^2 \cap d(\hat{\theta}_n, \theta_0) \leq 2^j \delta_n) \\ &\leq \sum_{j>M} \mathbb{P}\left(\sup_{\theta: d(\theta, \theta_0) \leq 2^j \delta_n} (M_n - M)(\theta_0) - (M_n - M)(\theta) \geq 2^{2j-2} \delta_n^2\right) \\ &\leq \sum_{j>M} \mathbb{E} \left[\frac{\sup_{\theta: d(\theta, \theta_0) \leq 2^j \delta_n} (M_n - M)(\theta_0) - (M_n - M)(\theta)}{2^{2j-2} \delta_n^2} \right] \end{aligned}$$

by [Markov's inequality](#), where we need to assume that it's non-negative. Now, we define the following.

Definition 4.3.3 (Localized empirical process). The *localized empirical process* for an [M-estimator problem](#) for $t > 0$ is defined as

$$\mathbb{E} \left[\sup_{\theta: d(\theta, \theta_0) \leq t} (M_n - M)(\theta) - (M_n - M)(\theta_0) \right].$$

Note the following.

Note. For nearly all *M*-estimation problems, the *localized empirical process* can be upper-bounded by a sequence of functions $\phi_n: [0, \infty] \rightarrow [0, \infty]$ such that for all $t > 0$,

$$\mathbb{E} \left[\sup_{\theta: d(\theta, \theta_0) \leq t} (M_n - M)(\theta) - (M_n - M)(\theta_0) \right] \leq \phi_n(t).$$

Assuming ϕ_n 's exist, we then proceed upper-bounding Equation 4.1 as

$$\mathbb{P}(d(\hat{\theta}_n, \theta_0) > 2^M \delta_n) \leq \sum_{j>M} \frac{\mathbb{E} \left[\sup_{\theta: d(\theta, \theta_0) \leq 2^j \delta_n} (M_n - M)(\theta) - (M_n - M)(\theta_0) \right]}{2^{2j-2} \delta_n^2} \leq \sum_{j>M} \frac{\phi_n(2^j \delta_n)}{2^{2j-2} \delta_n^2}.$$

To further bound the right-hand side, consider the following.

Definition 4.3.4 (Sub-quadratic assumption). The *sub-quadratic assumption* assumes that there exists $\alpha < 2$ such that for all n , $c > 1$, and $x > 0$,

$$\phi_n(cx) \leq c^\alpha \cdot \phi_n(x).$$

With *sub-quadratic assumption*, we get

$$\mathbb{P}(d(\hat{\theta}_n, \theta_0) > 2^M \delta_n) \leq \sum_{j>M} \frac{\phi_n(2^j \delta_n)}{2^{2j-2} \delta_n^2} \leq 4 \sum_{j>M} \frac{2^{\alpha j} \phi_n(\delta_n)}{2^{2j} \delta_n^2}.$$

This is basically the final bound we will get by introducing ϕ_n 's. Now, to control them, consider the so-called *rate-determining equation*.

Definition 4.3.5 (Rate-determining equation). Given a sequence $\{\delta_n\}$, the *rate-determining equation* for a *localized empirical process*'s upper-bounds ϕ_n 's is that for some c such that for all n ,

$$\phi_n(\delta_n) \leq c \delta_n^2.$$

Remark. It's important to check that whether such c exists for all n .

Intuition. We want to have $\phi_n(\delta_n) \approx \delta_n^2$.

Finally, assuming the *rate-determining equation* exists for some c , we get

$$\mathbb{P}(d(\hat{\theta}_n, \theta_0) > 2^M \delta_n) \leq 4 \sum_{j>M} \frac{2^{\alpha j} \phi_n(\delta_n)}{2^{2j} \delta_n^2} \leq 4c \cdot \sum_{j>M} \frac{2^{\alpha j}}{2^{2j}},$$

and from the *sub-quadratic assumption*, $\alpha < 2$, so the above sum converges to 0 as $M \rightarrow \infty$.

Remark. We can choose M not depending on n such that the right-hand side is $\leq \epsilon$.

Putting the above together, we have the following.

Theorem 4.3.1 (Non-asymptotic rate of convergence). For an *M*-estimation problem, assume the *growth condition* on M , and the *sub-quadratic assumption* (with parameter $\alpha < 2$) and the *rate-determining equation* are valid for ϕ_n 's arose from bounding the *localized empirical process*,

$$\mathbb{P}(d(\hat{\theta}_n, \theta_0) > 2^M \delta_n) \leq 4c \sum_{j>M} 2^{(\alpha-2)j}.$$

Back to the *mode estimation* example, we now want to formally show that the *rate of convergence* for $|\hat{\theta}_n - \theta_0|$ is $O_p(n^{-1/3})$.

Proposition 4.3.1. For the [mode estimation problem](#), the [rate of convergence](#) for $\hat{\theta}_n$ is $O_p(n^{-1/3})$.

Proof. We check the following.

- The [growth condition](#) is checked in a ball around θ_0 and not for all $\theta \in \Theta$ (since it's [consistent](#)). This is allowed as shown in [Theorem 4.3.2](#).
- Consider the [localized empirical process](#) with notation $m_\theta(x) = \mathbb{1}_{[\theta-1, \theta+1]}$, we have

$$\mathbb{E} \left[\sup_{\theta: |\theta - \theta_0| \leq t} |\mathbb{P}_n(m_\theta - m_{\theta_0}) - \mathbb{P}(m_\theta - m_{\theta_0})| \right].$$

Let $f_\theta(x) = m_\theta(x) - m_{\theta_0}(x)$, and $\mathcal{F} = \{f_\theta: |\theta - \theta_0| \leq t\}$. One can check that

$$f(x) = \mathbb{1}_{\theta_0-1-t \leq x \leq \theta_0-1+t} + \mathbb{1}_{\theta_0+1-t \leq x \leq \theta_0+1+t}$$

is an [envelope](#) for \mathcal{F} . Then, we have

$$\mathbb{P}F^2 \leq \int_{\theta_0-1-t}^{\theta_0-1+t} p_{\theta_0}(x) dx + \int_{\theta_0+1-t}^{\theta_0+1+t} p_{\theta_0}(x) dx \leq p_{\theta_0}(\theta_0) \cdot 4t \leq C_{p_{\theta_0}} t < \infty$$

for some constant $C_{p_{\theta_0}}$ depending on p_{θ_0} . With the [bracketing bound](#), for some constant $C > 0$,

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \sqrt{n} |\mathbb{P}_n f - \mathbb{P} f| \right] \leq C \cdot \|F\|_{L_2(\mathbb{P})} \int_0^1 \sqrt{\log N_{[\cdot]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon \|F\|_{L_2(\mathbb{P})})} d\epsilon.$$

Claim. For all ϵ , there exists some constant $C' > 0$ such that

$$N_{[\cdot]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon) \leq \left(\frac{1}{\epsilon} \right)^{C'} < \infty.$$

With the above claim and $\|F\|_{L_2(\mathbb{P})} \leq \sqrt{C_{p_{\theta_0}} t}$, the integral can be further bounded as

$$\int_0^1 \sqrt{\log N_{[\cdot]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon \|F\|_{L_2(\mathbb{P})})} d\epsilon \leq \int_0^1 \sqrt{C' \log \frac{1}{\epsilon \|F\|_{L_2(\mathbb{P})}}} d\epsilon < \infty,$$

hence, there exists some constant $C > 0$ such that

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{P}_n f - \mathbb{P} f| \right] \leq C \sqrt{\frac{t}{n}}.$$

This motivates us to define $\phi_n(t)$ as $C\sqrt{t/n}$.

- To check the [sub-quadratic assumption](#), for all n and $c > 1$, we have

$$\phi_n(ct) = C\sqrt{\frac{ct}{n}} = \sqrt{c} \cdot \left(C\sqrt{\frac{t}{n}} \right) = \sqrt{c} \cdot \phi_n$$

hence the [sub-quadratic assumption](#) is satisfied with $\alpha = 1/2$.

- Consider the [rate-determining equation](#) $\phi_n(t) \leq t^2$, for $t = \delta_n$,

$$\sqrt{t/n} \leq t^2 \Leftrightarrow \sqrt{\delta_n/n} \leq \delta_n^2 \Rightarrow 1/\sqrt{n} \leq \delta_n^{3/2} \Rightarrow \delta_n \approx 1/n^{1/3}.$$

In all, we have $|\hat{\theta}_n - \theta_0| = O_p(n^{-1/3})$. ■

[Proposition 4.3.1](#) is not fully proven yet, since we only have the [growth condition](#) satisfied in a ball

around θ_0 , not for all $\theta \in \Theta$.

Lecture 22: More Examples on Rate of Convergence

Before we extend [Theorem 4.3.1](#) to handle the local [growth condition](#), we note the following.

20 Oct. 9:00

Remark. The [rate](#) obtained from [Theorem 4.3.1](#) is usually correct; on the other hand, the probability tail bound obtained from [Theorem 4.3.1](#) is

$$\mathbb{P}(d(\hat{\theta}, \theta_0) > t\delta_n) \lesssim \frac{1}{t},$$

with $t = 2^M$ in the argument, which can be weak in the sense that it does not imply $\mathbb{E}[d(\hat{\theta}, \theta_0)] = O(\delta_n)$. Potentially, more sophisticated concentration arguments can be used.

Remark. The main step to apply [Theorem 4.3.1](#) is to bound the expected supremum of [localized empirical process](#), which can be hard.

Finally, as shown in some situations, we cannot expect the [growth condition](#) to hold for all $\theta \in \Theta$; instead, typically we only have $\theta \in B(\theta_0, u^*)$ for some $u^* \in \mathbb{R}$. In this case, a variation of [Theorem 4.3.1](#) still holds [[VW96](#), Theorem 3.2.5].

Theorem 4.3.2. For an [M-estimation problem](#), assume the [growth condition](#) on M holds for $\theta \in B(\theta_0, u^*)$ for some u^* , and the [sub-quadratic assumption](#) (with parameter $\alpha < 2$) and the [rate-determining equation](#) are valid for ϕ_n 's arose from bounding the [localized empirical process](#),

$$\mathbb{P}(d(\hat{\theta}_n, \theta_0) > 2^M \delta_n) \leq 4c \sum_{j>M} 2^{(\alpha-2)j}.$$

Proof. We again start by doing the [peeling step](#), but this time consider

$$\mathbb{P}(d(\hat{\theta}, \theta_0) > 2^M \delta_n) \leq \sum_{\substack{j>M: \\ 2^j \delta_n \leq u^*}} \mathbb{P}(2^{j-1} \delta_n < d(\hat{\theta}, \theta_0) < 2^j \delta_n) + \mathbb{P}\left(d(\hat{\theta}, \theta_0) > \frac{u^*}{2}\right).$$

We then handle the first term as in [Theorem 4.3.1](#), and show the second term goes to 0. ■

4.3.3 More Examples

We see two more examples of using [Theorem 4.3.2](#).

Example (Sample quantile). Let $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$ which has density f w.r.t. Lebesgue measure. Moreover, for $0 < \tau < 1$, let

$$\rho_\tau(x) = \begin{cases} (\tau - 1), & \text{if } x < 0; \\ \tau x, & \text{if } x \geq 0, \end{cases}$$

and $m_\theta(x) = \rho_\tau(x - \theta)$ for all $\theta \in \mathbb{R}$, so

$$M(\theta) = \mathbb{E}[m_\theta(x)], \quad M_n(\theta) = \frac{1}{n} \sum_{i=1}^n \rho_\tau(x_i - \theta)$$

We see that $\theta_0 := \arg \min_\theta M(\theta)$ is the τ^{th} quantile of \mathbb{P} , and let $\hat{\theta} := \arg \min_\theta M_n(\theta)$ be the corresponding [M-estimator](#). The [rate of convergence](#) is $|\hat{\theta}_n - \theta_0| = O_p(1/\sqrt{n})$.

Proof. To show the [rate of convergence](#), consider the following.

- $|m_{\theta_1}(x) - m_{\theta_2}(x)| \leq |\theta_1 - \theta_2|$, i.e., this is a Lipschitz [parametric](#) class.
- To show the [growth condition](#), we need the following.

Lemma 4.3.1. For all $w, v \in \mathbb{R}$,

$$\rho_\tau(w - u) - \rho_\tau(w) = -v(\tau - \mathbb{1}_{w \leq 0}) + \int_0^v [\mathbb{1}_{w \leq z} - \mathbb{1}_{w \leq 0}] dz.$$

Then, we can show the [growth condition](#) satisfy in a neighborhood of θ_0 .

Claim. For $d(\theta, \theta_0) = |\theta - \theta_0|$, for θ in some neighborhood of θ_0 , $|\theta - \theta_0|^2 \lesssim M(\theta) - M(\theta_0)$.

Proof. By denoting F as the CDF of f , we have

$$\begin{aligned} M(\theta_0 + \delta) - M(\theta_0) &= \mathbb{E} [\rho_\tau(x - \theta_0 - \delta) - \rho_\tau(x - \theta_0)] \\ &= \mathbb{E} [-\delta(\tau - \mathbb{1}_{x - \theta_0 \leq 0})] + \mathbb{E} \left[\int_0^\delta (\mathbb{1}_{x - \theta_0 \leq z} - \mathbb{1}_{x - \theta_0 \leq 0}) dz \right] \\ &= \int_0^\delta F(\theta_0 + z) - F(\theta_0) dz \end{aligned}$$

assume there exists a neighborhood of θ_0 such that $f \geq L > 0$, then for some $\xi_z \in (0, \delta)$,

$$\geq \int_0^\delta f(\xi_z) z dz \geq L \cdot \int_0^\delta z dz = \frac{L\delta^2}{2},$$

i.e., in this neighborhood, M grows quadratically in a neighborhood of θ_0 . \otimes

- To bound the [localized empirical process](#), first note that $\hat{\theta}$ is [consistent](#),^a and consider $\mathcal{F} = \{m_\theta - m_{\theta_0} : |\theta - \theta_0| \leq t\}$ where the [localized empirical process](#) is

$$\mathbb{E} \left[\sup_{\theta: |\theta - \theta_0| \leq t} |(\mathbb{P}_n - \mathbb{P})m_\theta - (\mathbb{P}_n - \mathbb{P})m_{\theta_0}| \right].$$

To use the [bracketing bound](#), we first see that $F(x) = t$ is an [envelope](#) with $\|F\|_{L_2(\mathbb{P})} = t$, hence the [localized empirical process](#) can be upper-bounded by

$$\frac{t}{\sqrt{n}} \int_0^1 \sqrt{\log N_{[\cdot]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon t)} d\epsilon \leq \frac{t}{\sqrt{n}} \int_0^1 \sqrt{\log \left(1 + \frac{4}{\epsilon} \right)} d\epsilon,$$

i.e., the [localized empirical process](#) can be upper-bounded by $\phi_n(t) \approx t/\sqrt{n}$.

- It's evident that ϕ_n 's satisfy the [sub-quadratic assumption](#) with $\alpha = 1$.
- By the [rate-determining equation](#),

$$\delta_n/\sqrt{n} \approx \delta_n^2 \Rightarrow \delta_n = 1/\sqrt{n},$$

implying $|\hat{\theta}_n - \theta_0| = O_p(1/\sqrt{n})$.

\otimes

^aThis should be proved beforehand, otherwise things doesn't make sense.

Remark. In this sample quantile example, the classical result for $\tau = 0.5$ shows that

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N} \left(0, \frac{1}{4(f(\theta_0))^2} \right).$$

Another example is the high-dimensional linear regression.

Example (High-dimensional linear regression). Consider $Z = (Y, X_1, \dots, X_p) \in \mathbb{R}^{p+1}$ such that $Z_1, \dots, Z_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$, and we want to predict Y by $\beta^\top X$. Let $M(\beta) = \mathbb{E}[(Y - \beta^\top X)^2]$ with

$$\beta^* = \arg \min_{\beta: \|\beta\|_1 \leq L} M(\beta)$$

for some L , and let $M_n(\beta) = \frac{1}{n} \sum_{i=1}^n (Y^i - \beta^\top X^i)^2$ with

$$\hat{\beta} = \arg \min_{\beta: \|\beta\|_1 \leq L} M_n(\beta).$$

Intuition. We want a *sparse* β^* , which is achieved by controlling L .

Note. We're not assuming the underlying \mathbb{P} to be linear.

The main question of interest for the [high dimensional linear regression](#) is the following.

Problem (Persistency). How large can $L = L(n, p)$ be so $M(\hat{\beta}) - M(\beta^*) \rightarrow 0$ as $n, p \rightarrow \infty$?

Answer. With some assumptions, [Theorem 4.3.3](#) shows that $L \lesssim \sqrt[4]{\log p/n}$. ⊗

Theorem 4.3.3. Under the setup of [high-dimensional linear regression](#), let $Y = X_0$ and define $F(Z) = \max_{0 \leq j, k \leq p} |X_j X_k - \mathbb{E}[X_j X_k]|$. Assume further that $\mathbb{E}[F^2(Z)] < \infty$, then

$$M(\hat{\beta}) - M(\beta^*) = O_p \left(L, \sqrt[4]{\frac{\log p}{n}} \right).$$

Proof. From the basic inequality, $M(\hat{\beta}) - M(\beta^*) \leq 2 \sup_{\beta: \|\beta\|_1 \leq L} |M_n(\beta) - M(\beta)|$. By letting

$$\gamma = (-1, \beta)^\top, \quad \Sigma = (\mathbb{E}[X_j^1 X_k^1])_{j,k=0,\dots,p}, \quad \Sigma_n = \left(\frac{1}{n} \sum_{i=1}^n X_j^i X_k^i \right)_{j,k=0,\dots,p},$$

we may then write $M(\beta) = \gamma^\top \Sigma \gamma$ and $M_n(\beta) = \gamma^\top \Sigma_n \gamma$. Hence,

$$\sup_{\beta: \|\beta\|_1 \leq L} |M_n(\beta) - M(\beta)| = |\gamma^\top \Sigma_n \gamma - \gamma^\top \Sigma \gamma| \leq \|\gamma\|_1^2 \cdot \|\Sigma_n - \Sigma\|_\infty \leq (1+L)^2 \|\Sigma_n - \Sigma\|_\infty,$$

which implies

$$\sup_{\beta: \|\beta\|_1 \leq L} \mathbb{P}(M(\hat{\beta}) - M(\beta^*) > \epsilon) \leq P((1+L)^2 \|\Sigma_n - \Sigma\|_\infty > \epsilon) \leq \frac{(1+L)^2 \mathbb{E}[\|\Sigma_n - \Sigma\|_\infty]}{\epsilon}$$

by [Markov's inequality](#). Finally, we observe that

$$\mathbb{E}[\|\Sigma_n - \Sigma\|_\infty] = \mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{P}_n f - \mathbb{P} f| \right]$$

where $\mathcal{F} = \{f_{jk}: 0 \leq j, k \leq p\}$ with $f_{jk} = X_j X_k - \mathbb{E}[X_j X_k]$. Now, F is clearly an [envelope](#) with $\|F\|_{L_2(\mathbb{P})} < \infty$, and by defining ϵ -brackets to be $[f_{j,k} \pm \epsilon/2]$ for all $j, k = 0, \dots, p$, we have

$$N_{[\cdot]}(\mathcal{F}, L_2(\mathbb{P}), \epsilon) \leq (p+1)^2.$$

By the [bracketing bound](#), $\mathbb{E} [\|\Sigma_n - \Sigma\|_\infty] \leq \frac{1}{\sqrt{n}} \|F\|_{L_2(\mathbb{P})} \sqrt{2 \log(p+1)}$, i.e.,

$$M(\hat{\beta}) - M(\beta^*) \leq \frac{(1+L)^2}{\sqrt{n}} \|F\|_{L_2(\mathbb{P})} \sqrt{2 \log(p+1)}.$$

In order to let this goes to 0, we require $L \lesssim \sqrt[4]{n/\log p}$, which finally implies

$$M(\hat{\beta}) - M(\beta^*) = O_p \left(L, \sqrt[4]{\frac{\log p}{n}} \right).$$

■

Remark. Observe that in the above proof, we do not need to [localize the empirical process](#), i.e., the [bracketing bound](#) can be used for any [empirical process](#) induced by finite class.

Lecture 23: Non-Parametric Regression

4.4 Non-Parametric Regression

23 Oct. 9:00

Now, we can see more advanced applications based on [M-estimations](#).

4.4.1 Fixed Design

First, consider the following problem.

Problem 4.4.1 (Fixed design non-parametric least square). Let $\mathcal{F} = \{f: [0, 1] \rightarrow [-1, 1] \text{ 1-Lipschitz}\}$, and consider $\epsilon_1, \dots, \epsilon_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$ such that for a fixed $f^* \in \mathcal{F}$, for all $i = 1, \dots, n$,

$$y_i = f^*(i/n) + \epsilon_i$$

Then, the problem of *fixed design non-parametric least square* is to find the least square estimate

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \left(y_i - f\left(\frac{i}{n}\right) \right)^2.$$

Notation. [Fixed design non-parametric least square](#) problem is fixed in the sense of the data $(x, y = f(x) + \epsilon)$ is generated at the fixed grid $x \in \{1/n, \dots, 1\}$.

The main idea to solve the [fixed design non-parametric least square](#) is because it appeals to solve an [M-estimation](#) problem.

Remark. \hat{f}_n is an [M-estimator](#).

Proof. We first observe that from the definition of y_i ,^a

$$\hat{f}_n = \arg \max_{f \in \mathcal{F}} \frac{2}{n} \sum_{i=1}^n \epsilon_i \left(f\left(\frac{i}{n}\right) - f^*\left(\frac{i}{n}\right) \right) - \frac{1}{n} \sum_{i=1}^n \left(f\left(\frac{i}{n}\right) - f^*\left(\frac{i}{n}\right) \right)^2,$$

which motivates us to define $M_n(f), M(f): \mathcal{F} \rightarrow \mathbb{R}$ such that

$$M_n(f) = \frac{2}{n} \sum_{i=1}^n \epsilon_i \left(f\left(\frac{i}{n}\right) - f^*\left(\frac{i}{n}\right) \right) - \frac{1}{n} \sum_{i=1}^n \left(f\left(\frac{i}{n}\right) - f^*\left(\frac{i}{n}\right) \right)^2$$

and

$$M(f) = \mathbb{E} [M_n(f)] = -\frac{1}{n} \sum_{i=1}^n \left(f\left(\frac{i}{n}\right) - f^*\left(\frac{i}{n}\right) \right)^2.$$

We see that $f^* = \arg \max_{f \in \mathcal{F}} M(f)$.

⊛

^aNote that the term $-\frac{1}{n} \sum_{i=1}^n \epsilon_i^2$ is omitted since it doesn't depend on $f \in \mathcal{F}$.

Since \hat{f}_n is an *M-estimator*, we want to get the *rate*. We need to verify the following.

Claim. The *growth condition* is satisfied for some d .

Proof. By choosing the canonical d , i.e.,

$$d(f, f^*) := \sqrt{M(f^*) - M(f)} = \sqrt{\frac{1}{n} \sum_{i=1}^n \left(f\left(\frac{i}{n}\right) - f^*\left(\frac{i}{n}\right) \right)^2} = L_2(\mathbb{P}_n)(f, f^*),$$

where \mathbb{P}_n is the empirical measure on the grid $\{1/n, 2/n, \dots, 1\}$, we see that it satisfies the *growth condition* automatically.

⊛

To bound the *localized empirical process* for the *fixed design non-parametric least square*

$$\mathbb{E} \left[\sup_{f \in \mathcal{F} : d(f, f^*) \leq \delta} (M_n - M)(f) - (M_n - M)(f^*) \right].$$

Claim. This *localized empirical process* is bounded by $\phi_n(\delta) = c\sqrt{\delta/n}$ for some $c > 0$.

Proof. We first observe that

$$\mathbb{E} \left[\sup_{\substack{f \in \mathcal{F} : \\ d(f, f^*) \leq \delta}} (M_n - M)(f) - (M_n - M)(f^*) \right] = 2 \cdot \mathbb{E} \left[\sup_{\substack{f \in \mathcal{F} : \\ d(f, f^*) \leq \delta}} \frac{1}{n} \sum_{i=1}^n \epsilon_i \left(f\left(\frac{i}{n}\right) - f^*\left(\frac{i}{n}\right) \right) \right].$$

For $f \in \mathcal{F}$, define $X_f = \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i \left(f\left(\frac{i}{n}\right) - f^*\left(\frac{i}{n}\right) \right)$, so for $f, g \in \mathcal{F}$,

$$X_f - X_g = \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i \left(f\left(\frac{i}{n}\right) - g\left(\frac{i}{n}\right) \right) \sim \mathcal{N}(0, \sigma^2 d^2(f, g)),$$

which implies that X_f is a *sub-Gaussian process* with

$$\mathbb{P}(|X_f - X_g| \geq u) \leq \exp\left(-\frac{u^2}{2d^2(f, g)}\right)$$

when $\sigma^2 = 1$.^a So, by *Dudley's entropy integral bound*, for some $c > 0$,

$$\begin{aligned} \mathbb{E} \left[\sup_{\substack{f \in \mathcal{F} : \\ d(f, f^*) \leq \delta}} \frac{1}{n} \sum_{i=1}^n \epsilon_i \left(f\left(\frac{i}{n}\right) - f^*\left(\frac{i}{n}\right) \right) \right] &\leq \frac{12}{\sqrt{n}} \int_0^\delta \sqrt{\log N(\mathcal{F}, d, \epsilon)} \, d\epsilon \\ &\leq \frac{1}{\sqrt{n}} \int_0^\delta \sqrt{\frac{c}{\epsilon}} \, d\epsilon = c\sqrt{\frac{\delta}{n}} =: \phi_n(\delta) \end{aligned}$$

since $N(\mathcal{F}, d, \epsilon) \leq N(\mathcal{F}, \|\cdot\|_\infty, \epsilon) \leq \exp(c_1/\epsilon)$ from *Theorem 3.3.1*.

⊛

^aWithout this assumption, we just have σ^2 in the final bound, which is still a *sub-Gaussian process*.

Then, we verify the *sub-quadratic assumption* for such ϕ_n 's.

Claim. ϕ_n 's satisfy the [sub-quadratic assumption](#) for $\alpha = 1/2$.

Proof. Since $\phi_n(c\delta) = \sqrt{c} \cdot c\sqrt{\delta/n} = \sqrt{c} \cdot \phi_n(\delta)$. ⊛

With all these, we can finally solve the [rate-determining equation](#), which is

$$\phi_n(\delta) \approx \delta^2 \Rightarrow \sqrt{\frac{\delta}{n}} \approx \delta^2 \Rightarrow \delta \approx n^{-1/3}.$$

In all, we have the [rate of convergence](#)

$$d(\hat{f}_n, f^*) = \sqrt{\frac{1}{n} \sum_{i=1}^n \left(\hat{f}_n\left(\frac{i}{n}\right) - f^*\left(\frac{i}{n}\right) \right)^2} = O_p(n^{-1/3}),$$

or more specifically, from [Theorem 4.3.2](#),

$$\mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n \left(\hat{f}_n\left(\frac{i}{n}\right) - f^*\left(\frac{i}{n}\right) \right)^2 > 2^{2M} n^{-2/3} \right) \leq c \cdot 2^{-M}.$$

Remark. The [rate of convergence](#) in mean-square error (i.e., w.r.t. $d^2 = L_2^2(\mathbb{P})$) for 1-Lipschitz regression is $n^{-2/3}$. Moreover, the “min-max rate” worst case over all $f^* \in \mathcal{F}$ is $n^{-2/3}$, hence the [fixed design non-parametric least square estimator](#) is “min-max rate optimal”. These terms will make sense in the next [remark](#).

More generally, if we assume that $f^* \in \mathcal{S}_\alpha$, we can also solve the [fixed design non-parametric least square](#) over \mathcal{S}_α as follows.

Theorem 4.4.1. Consider the [fixed design non-parametric least square](#) problem over \mathcal{S}_α , the [rate of convergence](#) of $d(\hat{f}_n, f^*)$ for the canonical d is $O_p(n^{-\frac{\alpha}{2\alpha+1}})$ for $\alpha > 1/2$.

Proof. From the same calculation, the corresponding [localized empirical process](#) is bounded by

$$\begin{aligned} \mathbb{E} \left[\sup_{\substack{f \in \mathcal{S}_\alpha: \\ d(f, f^*) \leq \delta}} \frac{1}{n} \sum_{i=1}^n \epsilon_i \left(f\left(\frac{i}{n}\right) - f^*\left(\frac{i}{n}\right) \right) \right] &\leq \frac{c}{\sqrt{n}} \int_0^\delta \sqrt{\log N(\mathcal{S}_\alpha, d, \epsilon)} \, d\epsilon \\ &\leq \frac{c}{\sqrt{n}} \int_0^\delta \sqrt{\log N(\mathcal{S}_\alpha, \|\cdot\|_\infty, \epsilon)} \, d\epsilon \leq \frac{c}{\sqrt{n}} \int_0^\delta \left(\frac{1}{\epsilon} \right)^{\frac{1}{2\alpha}} \, d\epsilon \end{aligned}$$

from [Theorem 3.3.1](#). Since $\alpha > 1/2$, this bound gives $\lesssim \frac{1}{\sqrt{n}} \delta^{1-1/2\alpha}$, so we can take

$$\phi_n(\delta) := \frac{\delta^{1-\frac{1}{2\alpha}}}{\sqrt{n}}.$$

We see that the [sub-quadratic assumption](#) is satisfied since $\phi_n(c\delta) = c^{1-\frac{1}{2\alpha}} \phi_n(\delta)$ with $1 - 1/2\alpha < 2$. Solving the [rate-determining equation](#), we have

$$\phi_n(\delta_n) \approx \delta_n^2 \Rightarrow \frac{\delta_n^{1-\frac{1}{2\alpha}}}{\sqrt{n}} \approx \delta_n^2 \Rightarrow \delta_n \approx n^{-\frac{\alpha}{2\alpha+1}}.$$

It's usually more natural to consider the rate for $d^2(\hat{f}_n, f^*) = L_2^2(\mathbb{P}_n)(\hat{f}_n, f^*) = O_p(n^{-\frac{2\alpha}{2\alpha+1}})$. ■

Note. When $\alpha < 1/2$, we need to use the [modified version of Dudley's entropy integral bound](#). In this case, we will get a slower rate than $n^{-2\alpha/2\alpha+1}$ w.r.t. d^2 .

Remark (Min-max rate optimal). $n^{-2\alpha/2\alpha+1}$ is *min-max rate optimal*, i.e., it's the “right” rate.

Proof. We just showed that $\mathbb{E}[L_2^2(\mathbb{P})(\hat{f}_n, f^*)] \leq cn^{-2\alpha/(2\alpha+1)}$, and it's known that the *min-max rate* is lower-bounded by

$$\inf_{\hat{f}_n} \sup_{f \in \mathcal{F}} \mathbb{E} \left[L_2^2(\mathbb{P})(\hat{f}_n, f^*) \right] \geq cn^{-\frac{2\alpha}{2\alpha+1}}.$$

Hence, the rate for the *fixed design non-parametric least square estimator* is *min-max rate optimal* over \mathcal{S}_α for $\alpha > 1/2$. However, it's not known whether the same conclusion holds for $\alpha < 1/2$. \circledast

We do not know α . A natural question is the following.

Problem. Can we adaptively get $O(n^{-2\alpha/2\alpha+1})$ rate without knowing α ?

Answer. Yes. There are several ways to do this [Ts08]. \circledast

4.4.2 Prediction v.s. Estimation

Now, we want to consider the random design case (i.e., now data does not come from the fixed grid). First, recall the *statistical learning set-up*.

As previously seen. Given $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$, with $f^*(x) = \mathbb{E}[Y \mid X = x]$, i.e.,

$$f^* = \arg \min_{f \text{ measurable}} \mathbb{E}_{(X,Y)} \left[(Y - f(X))^2 \right].$$

Given a class of function \mathcal{F} , we can define

$$f_{\mathcal{F}} = \arg \min_{f \in \mathcal{F}} \mathbb{E} \left[(Y - f(X))^2 \right].$$

For any prediction function $f: \mathcal{X} \rightarrow \mathbb{R}$, its *excess risk* is

$$\mathbb{E} \left[(Y - f(X))^2 \right] - \inf_{h \in \mathcal{F}} \mathbb{E} \left[(Y - h(X))^2 \right].$$

Observe that we can write

$$\begin{aligned} & \mathbb{E} \left[(Y - f(X))^2 \right] - \inf_{h \in \mathcal{F}} \mathbb{E} \left[(Y - h(X))^2 \right] \\ &= \mathbb{E} \left[(Y - f(X) \pm f^*(X))^2 \right] - \inf_{h \in \mathcal{F}} \mathbb{E} \left[(Y - h(X) \pm f^*(X))^2 \right] \\ &= \mathbb{E} \left[(f(X) - f^*(X))^2 \right] - \inf_{h \in \mathcal{F}} \mathbb{E} \left[(h(X) - f^*(X))^2 \right] \quad \text{cross terms are 0} \\ &= \|f - f^*\|_{L_2(\mathbb{P})}^2 - \inf_{h \in \mathcal{F}} \|h - f^*\|_{L_2(\mathbb{P})}^2, \end{aligned}$$

where the last equation is the estimation error in the traditional statistics terminology.

Remark. Prediction and estimation are the same in this content.

Now, if we change f to a predictor \hat{f}_n (depending on the training data), we have

$$\mathbb{E} \left[(Y - \hat{f}_n(X))^2 \right] - \inf_{h \in \mathcal{F}} \mathbb{E} \left[(Y - h(X))^2 \right] = \|\hat{f}_n - f^*\|_{L_2(\mathbb{P})}^2 - \inf_{h \in \mathcal{F}} \|h - f^*\|_{L_2(\mathbb{P})}^2$$

There are two standard scenarios:

- (a) Well specified case: Given \mathcal{F} , assume $f^* \in \mathcal{F}$. In particular, \mathbb{P} is such that $f(x) = \mathbb{E}[Y \mid X = x] \in \mathcal{F}$. In this case, the *excess risk* of \hat{f}_n is just $\|\hat{f}_n - f^*\|_{L_2(\mathbb{P})}^2$ since $\inf_{h \in \mathcal{F}} \|h - f^*\|_{L_2(\mathbb{P})}^2 = 0$, i.e., the estimation error in $L_2(\mathbb{P})$. (this is just the traditional statistics).
- (b) Mis-specified case: do not assume $f^* \in \mathcal{F}$, we need to bound the *oracle inequality* (which bounds both terms) entirely.

Lecture 24

25 Oct. 9:00

Appendix

Bibliography

- [BLM13] S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. OUP Oxford, 2013. ISBN: 978-0-19-953525-5. URL: <https://books.google.com/books?id=5oo4YIz6tR0C>.
- [Han16] Ramon van Handel. “Probability in High Dimensions”. In: (2016). URL: <https://web.math.princeton.edu/~rvan/APC550.pdf>.
- [Nov62] Albert BJ Novikoff. “On convergence proofs on perceptrons”. In: *Proceedings of the Symposium on the Mathematical Theory of Automata*. Vol. 12. 1. New York, NY. 1962, pp. 615–622.
- [Tsy08] A.B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer Series in Statistics. Springer New York, 2008. ISBN: 978-0-387-79052-7. URL: <https://books.google.com/books?id=mwB8rUBsbqoC>.
- [Vaa98] A. W. van der Vaart. *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge: Cambridge University Press, 1998. ISBN: 978-0-521-78450-4. DOI: [10.1017/CB09780511802256](https://doi.org/10.1017/CB09780511802256). URL: <https://www.cambridge.org/core/books/asymptotic-statistics/A3C7DAD3F7E66A1FA60E9C8FE132EE1D> (visited on 10/17/2023).
- [VW96] Aad W. Van Der Vaart and Jon A. Wellner. *Weak Convergence and Empirical Processes*. Springer Series in Statistics. New York, NY: Springer, 1996. ISBN: 978-1-4757-2547-6 978-1-4757-2545-2. DOI: [10.1007/978-1-4757-2545-2](https://doi.org/10.1007/978-1-4757-2545-2). (Visited on 08/21/2023).