

# Uniform Law of Large Numbers

---

**Marcelo Ortiz-Villavicencio**



**February 2, 2024**

Econometrics Reading Group

## 1. Introduction

Uniform convergence of CDF

Uniform law for more general function classes

## 2. A uniform law via Rademacher Complexity

## 3. Upper bounds on the Rademacher complexity

VC dimension

# Introduction

---

- In the previous chapter we pointed out some limitations of asymptotic analysis in high dimensions.
- In this chapter, we will turn our focus to *non-asymptotic analysis*, where we provide convergence guarantees without having the number of observations  $n \rightarrow \infty$ .
- Our focus in this lecture is a set of results called the *uniform law of large numbers*.
- These results represent a strengthening of the usual LLN, which applies to a fixed sequence of RV, to related laws that hold *uniformly* over collections of RV.
- Such uniform laws are of theoretical interest in *Empirical Process Theory*.
- Moreover, they also play an important role in understanding the behavior of different statistical estimators providing guarantees for bounds of the form:

$$\Pr \left[ \sup_{h \in \mathcal{H}} |\hat{\mathcal{L}}(h) - \mathcal{L}(h)| \leq \epsilon \right] \geq 1 - \delta.$$

- A *stochastic process* is a collection of random variables  $\{X(t), t \in T\}$  on the same probability space, indexed by an arbitrary set  $T$ .
- An *empirical process* is a stochastic process based on a random sample.
- Consider a random sample  $X_1, \dots, X_n$  of independent draws from a probability measure  $P$  on an arbitrary sample  $\mathcal{X}$ .
- For a set  $A$ , we define the *empirical measure* (distribution) of  $A$  to be

$$\mathbb{P}_n(A) = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}(A)$$

where  $\delta_x(A)$  is the *Dirac measure* (or point mass) that assigns mass 1 if  $x \in A$  and zero elsewhere

- An *Indicator function* is closely related to a Dirac.

- Given some integrable function  $g$ , we may define the expectation functional  $\gamma_g$  via

$$\gamma_g(P) = \int g(x)dP(x) = \mathbb{E}[g(X)]$$

- We can think about the previous expression as its an *empirical integral*

$$\gamma_g(\mathbb{P}_n) = \int g(x)d\mathbb{P}_n(x) = \mathbb{E}_n[g(X)]$$

- For any class  $\mathcal{F}$  of measurable functions  $f: \mathcal{X} \mapsto \mathbb{R}$ , an *empirical process*  $\{\gamma_f(\mathbb{P}_n), f \in \mathcal{F}\}$  can be defined.
- Our goal is to show "how close"  $\{\gamma_f(\mathbb{P}_n), f \in \mathcal{F}\}$  is to  $\{\gamma_f(P), f \in \mathcal{F}\}$ .

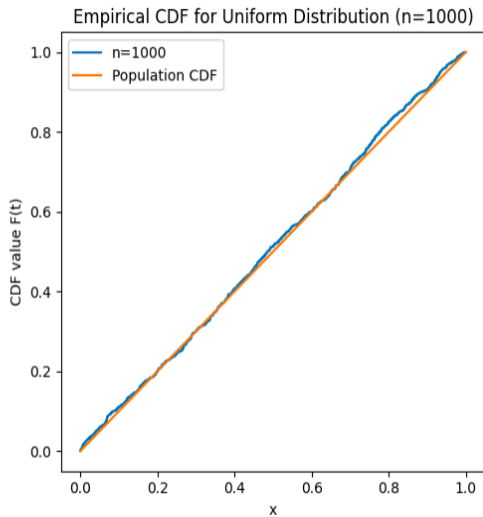
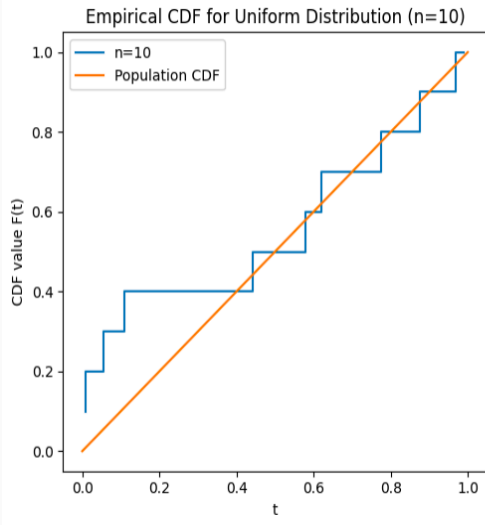
- The law of any scalar random variable  $X$  can be fully specified by its *cumulative distribution function* (CDF), whose value at any point  $t \in \mathbb{R}$  is given by  $F(t) := P(X \leq t)$ .
- Suppose a collection  $\{X_i\}_{i=1}^n$  of  $n$  i.i.d. samples, each drawn according to the law specified by  $F$ .
- A natural estimate of  $F$  is the empirical CDF given by

$$\hat{F}_n(t) := \frac{1}{n} \sum_{i=1}^n 1_{(-\infty, t]} [X_i]$$

where  $1_{(-\infty, t]} [x]$  is a  $\{0, 1\}$ -valued *indicator function* for the event  $\{x \leq t\}$ .

- Since the population CDF can be written as  $F(t) = \mathbb{E} [1_{(-\infty, t]} [X]]$ , the empirical CDF is an *unbiased estimate*.

# Uniform convergence of CDF





- Given a pair of CDFs  $F$  and  $G$ , let us measure the distance between them using the *sup-norm*

$$\|G - F\|_{\infty} := \sup_{t \in \mathbb{R}} |G(t) - F(t)|$$

- We can define then the *continuity* of a functional  $\gamma$  with respect to this norm.
- More precisely, the functional  $\gamma$  is *continuous* at  $F$  in the sup-norm if,  $\forall \epsilon > 0$ , there exists a  $\delta > 0$  such that  $\|G - F\|_{\infty} \leq \delta$  implies that  $|\gamma(G) - \gamma(F)| \leq \epsilon$

## Theorem (Glivenko-Cantelli)

*For any distribution, the empirical CDF  $\hat{F}_n$  is a strongly consistent estimator of the population CDF in the uniform norm, meaning that*

$$\|\hat{F}_n - F\|_{\infty} \xrightarrow{a.s.} 0.$$

- This notion is useful because, for *any* continuous functional, it reduces the consistency question to the issue of whether or not that difference converges to zero.

- Let's focus on a more general consideration of ULLN.
- Let  $\mathcal{F}$  be a class of integrable real-valued functions with domain  $X$ , and let  $\{X_i\}_{i=1}^n$  be a collection of i.i.d. samples from some distribution  $\mathbb{P}$  over  $X$ .
- Consider the random variable

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

which measures the absolute deviation between the sample average  $\frac{1}{n} \sum_{i=1}^n f(X_i)$  and the population average  $\mathbb{E}[f(X)]$ , uniformly over the class  $\mathcal{F}$ .

## Theorem

*We say that  $\mathcal{F}$  is a Glivenko-Cantelli class for  $\mathbb{P}$  if  $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$  converges to zero in probability as  $n \rightarrow \infty$ .*

Note: Not all classes of functions are Glivenko-Cantelli.



- These quantities are one of the main focus of methods based on *empirical risk minimization*.
- Consider an indexed family of probability distributions  $\{\mathbb{P}_\theta \mid \theta \in \Omega\}$ ,
- Suppose we have access to  $n$  samples  $\{X_i\}_{i=1}^n$ , each sample lying in some space  $\mathcal{X}$ .
- Those samples are drawn i.i.d. according to a distribution  $\mathbb{P}_{\theta^*}$ , for some *fixed* but *unknown*  $\theta^* \in \Omega$ . Here the index  $\theta^*$  could lie within a *finite-dimensional space*, such as  $\Omega = \mathbb{R}^d$ , or could lie within some function class  $\Omega = \mathcal{G}$ , in which case the problem is nonparametric (i.e., *infinite-dimensional space*).
- A standard *decision-theoretic* approach to estimating  $\theta^*$  is based on minimizing a **cost function** of the form  $\theta \mapsto \mathcal{L}_\theta(X)$ , which measures the "fit" between a parameter  $\theta \in \Omega$  and the sample  $X \in \mathcal{X}$ .

- Given the collection of  $n$  samples  $\{X_i\}_{i=1}^n$ , the principle of *empirical risk minimization* is based on the objective function (a.k.a. *empirical risk*)

$$\widehat{R}_n(\theta, \theta^*) := \frac{1}{n} \sum_{i=1}^n \mathcal{L}_\theta(X_i)$$

- The empirical risk is contrasted with the *population risk*, where the expectation  $\mathbb{E}_{\theta^*}$  is taken over a sample  $X \sim \mathbb{P}_{\theta^*}$ ,

$$R(\theta, \theta^*) := \mathbb{E}_{\theta^*} [\mathcal{L}_\theta(X)],$$

- In practice, one minimizes the empirical risk over some *subset*  $\Omega_0$  of the full space  $\Omega$ , thereby obtaining some estimate  $\widehat{\theta}$ .
- The question is how to bound the *excess risk*, measured in terms of the population quantities - namely the difference

$$E(\widehat{\theta}, \theta^*) := R(\widehat{\theta}, \theta^*) - \inf_{\theta \in \Omega_0} R(\theta, \theta^*).$$

- Our goal is to develop methods for controlling the *excess risk*.
- Suppose there exists some  $\theta_0 \in \Omega_0$  such that

$$R(\theta_0, \theta^*) = \inf_{\theta \in \Omega_0} R(\theta, \theta^*)$$

- Then, the excess risk can be decomposed as

$$E(\hat{\theta}, \theta^*) = \underbrace{\{R(\hat{\theta}, \theta^*) - \hat{R}_n(\hat{\theta}, \theta^*)\}}_{T_1} + \underbrace{\{\hat{R}_n(\hat{\theta}_0, \theta^*) - \hat{R}_n(\theta_0, \theta^*)\}}_{T_2 \leq 0} + \underbrace{\{\hat{R}_n(\theta_0, \theta^*) - R(\theta_0, \theta^*)\}}_{T_3}$$

- $T_2$  is non-positive, since  $\hat{\theta}$  minimizes the empirical risk over  $\Omega_0$ .
- Because  $\theta_0$  is an unknown but non-random quantity, and recalling the definition of empirical risk,  $T_3$  can be rewritten as

$$T_3 = \left[ \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\theta_0}(X_i) \right] - \mathbb{E}_X[\mathcal{L}_{\theta_0}(X)]$$

corresponding to the deviation of a sample mean from its expectation for the random variable  $\mathcal{L}_{\theta_0}(X)$ .

- $T_1$  can be written in a similar way, namely as the difference

$$T_1 = \mathbb{E}_X [\mathcal{L}_{\hat{\theta}}(X)] - \left[ \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\hat{\theta}}(X_i) \right].$$

- This quantity is more challenging to control, because the parameter  $\hat{\theta}$  (in contrast to  $\theta_0$ ) is now random, and depends on the samples  $\{X_i\}_{i=1}^n$ .
- Hence, controlling  $T_1$  requires a *stronger* result, such as a *uniform law of large numbers* over the cost function class  $\mathcal{L}(\Omega_0) := \{x \mapsto \mathcal{L}_\theta(x), \theta \in \Omega_0\}$ .
- With this notation, we have

$$T_1 \leq \sup_{\theta \in \Omega_0} \left| \frac{1}{n} \sum_{i=1}^n \mathcal{L}_\theta(X_i) - \mathbb{E}_X [\mathcal{L}_\theta(X)] \right| = \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{L}(\Omega_0)}$$

- This same quantity also dominates  $T_3$ , we conclude that the *excess risk* is at most

$$2 \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{L}(\Omega_0)}$$

# **A uniform law via Rademacher Complexity**

---

- Let me now turn to the technical details of deriving such results.
- An important concept that underlies the study of uniform laws is the *Rademacher complexity* of the function class  $\mathcal{F}$ .
- For any fixed collection  $x_1^n := (x_1, x_2, \dots, x_n)$  of points, consider the subset of  $\mathbb{R}^n$  given by

$$\mathcal{F}(x_1^n) := \{(f(x_1), f(x_2), \dots, f(x_n)) \mid f \in \mathcal{F}\}$$

- The set  $\mathcal{F}(x_1^n)$  corresponds to all those vectors in  $\mathbb{R}^n$  that can be realized by applying a function  $f \in \mathcal{F}$  to the collection  $(x_1, x_2, \dots, x_n)$  and the *empirical Rademacher complexity* is given by

$$\mathcal{R}(\mathcal{F}(x_1^n)/n) := \mathbb{E}_\varepsilon \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right].$$

- Given a collection  $X_1^n := \{X_i\}_{i=1}^n$  of random samples, then the *empirical Rademacher complexity*  $\mathcal{R}(\mathcal{F}(X_1^n)/n)$  is a random variable.



- Taking its expectation yields the Rademacher complexity of the function class  $\mathcal{F}$  we get

$$\mathcal{R}_n(\mathcal{F}) := \mathbb{E}_{\mathcal{X}} [\mathcal{R}(\mathcal{F}(X_1^n) / n)] = \mathbb{E}_{\mathcal{X}, \varepsilon} \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right| \right].$$

- Note that the *Rademacher complexity* is the average of the maximum correlation between the vector  $(f(X_1), \dots, f(X_n))$  and the "noise vector"  $(\varepsilon_1, \dots, \varepsilon_n)$ , where the maximum is taken over all functions  $f \in \mathcal{F}$ .
- **Intuition:** a function class is extremely large if we can always find a function that has a high correlation with a randomly drawn noise vector. Conversely, when the Rademacher complexity decays as a function of sample size, then it is impossible to find a function that correlates very highly in expectation with a randomly drawn noise vector.
- **Simple words:** If the  $\mathcal{R}_n(\mathcal{F})$  is small, it suggests that the function class is not very sensitive to random noise in the data. In other words, small  $\mathcal{R}_n(\mathcal{F})$  often implies better generalization performance in a learning algorithm.

- There is a connection between *Rademacher complexity* and the *Glivenko-Cantelli* theorem.
- In particular, any bounded function class  $\mathcal{F}$ , the condition  $\mathcal{R}_n(\mathcal{F}) = o(1)$  implies the *Glivenko-Cantelli* property.

## Theorem

For any  $b$ -uniformly bounded class of functions  $\mathcal{F}$ , any positive integer  $n \geq 1$  and any scalar  $\delta \geq 0$ , we have

$$\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq 2\mathcal{R}_n(\mathcal{F}) + \delta$$

with probability at least  $1 - \exp\left(-\frac{n\delta^2}{2b^2}\right)$ . Consequently, as long as  $\mathcal{R}_n(\mathcal{F}) = o(1)$ , we have

$$\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \xrightarrow{a.s.} 0.$$

- This is nothing more than a tail bound for the probability that the RV  $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$  deviates *substantially above* a multiple of the Rademacher complexity.
- Therefore, we need to obtain upper bounds on the  $\mathcal{R}_n(\mathcal{F})$

# Upper bounds on the Rademacher complexity

---

- To make the previous theorem useful, we require methods for *upper bounding* the Rademacher complexity.
- There are several methods to do so, ranging from simple union bounds (suitable for finite function classes) to more advanced techniques involving *metric entropy* and *chaining* (I will skip this due to time constraints, sorry).
- Instead, we gonna focus on more "elementary" techniques that apply for function classes with *polynomial discrimination* and *Vapnik-Chervonenski* classes.

- It is relatively straightforward to establish uniform laws for function classes with *polynomial discrimination*
- Our interest in function classes for which the *cardinality* grows only as a polynomial function of sample size.

## Definition (Polynomial discrimination)

A class  $\mathcal{F}$  of functions with domain  $\mathcal{X}$  has polynomial discrimination of order  $v \geq 1$  if, for each positive integer  $n$  and collection  $x_1^n = \{x_1, \dots, x_n\}$  of  $n$  points in  $\mathcal{X}$ , the set  $\mathcal{F}(x_1^n)$  has *cardinality upper bounded*

$$\text{card}(\mathcal{F}(x_1^n)) \leq (n + 1)^v$$

- Previous property provides a straightforward approach to controlling the Rademacher complexity.
- For any set  $S \subset \mathbb{R}^n$ , we use  $D(S) := \sup_{x \in S} \|x\|_2$  to denote its *maximal width* in the  $\ell_2$ -norm.

## Lemma (Upper bound on Rademacher Complexity)

Suppose that  $\mathcal{F}$  has polynomial discrimination of order  $v$ . Then for all positive integers  $n$  and any collection of points  $x_1^n = (x_1, \dots, x_n)$ ,

$$\underbrace{\mathbb{E}_\varepsilon \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right]}_{\mathcal{R}(\mathcal{F}(x_1^n)/n)} \leq 4D(x_1^n) \sqrt{\frac{v \log(n+1)}{n}},$$

where  $D(x_1^n) := \sup_{f \in \mathcal{F}} \sqrt{\frac{\sum_{i=1}^n f^2(x_i)}{n}}$  is the  $\ell_2$ -radius of the set  $\mathcal{F}(x_1^n) / \sqrt{n}$ .

- A special simple case is when the function class is  $b$  *uniformly bounded* so that  $D(x_1^n) \leq b$  for all samples.
- Applying the lemma

$$\mathcal{R}_n(\mathcal{F}) \leq 2b\sqrt{\frac{v \log(n+1)}{n}} \text{ for all } n \geq 1$$

- Combined with the Theorem, we conclude that any bounded function class with polynomial discrimination is Glivenko-Cantelli.
- What types of function classes have polynomial discrimination? A good example is based on *indicator functions* of the *left-sided intervals*  $(-\infty, t]$  (e.g., CDFs)

## Corollary (Classical Glivenko-Cantelli)

Let  $F(t) = \mathbb{P}[X \leq t]$  be the CDF of a random variable  $X$ , and let  $\hat{F}_n$  be the empirical CDF based on  $n$  i.i.d. samples  $X_i \sim \mathbb{P}$ . Then

$$\mathbb{P} \left[ \left\| \hat{F}_n - F \right\|_{\infty} \geq 8\sqrt{\frac{\log(n+1)}{n}} + \delta \right] \leq e^{-\frac{n\delta^2}{2}} \quad \text{for all } \delta \geq 0,$$

and hence  $\left\| \hat{F}_n - F \right\|_{\infty} \xrightarrow{a.s.} 0$ .

- In this section, we briefly discuss a classical notion of complexity measure of function class, VC dimension.
- Let us consider a function class  $\mathcal{F}$  in which each function  $f$  is binary-valued, taking the values  $\{0, 1\}$  for concreteness.
- In this case, the set  $\mathcal{F}(x_1^n)$  can have at most  $2^n$  elements.

## Definition (Shattering and VC dimension)

Given a class  $\mathcal{F}$  of binary-valued functions, we say that the set  $x_1^n = (x_1, \dots, x_n)$  is shattered by  $\mathcal{F}$  if  $\text{card}(\mathcal{F}(x_1^n)) = 2^n$ . The VC dimension  $v(\mathcal{F})$  is the largest integer  $n$  for which there is some collection  $x_1^n = (x_1, \dots, x_n)$  of  $n$  points that is shattered by  $\mathcal{F}$ .

- When the quantity  $v(\mathcal{F})$  is finite, then the function class  $\mathcal{F}$  is said to be a VC class.
- Let's finish with an example.



- This example was taken from Ma (2022)
- Will show that VC dimension is an upper bound on the Rademacher complexity.
- The labels belong to the output space  $\mathcal{Y} = \{-1, 1\}$ , each classifier is a function  $h : \mathcal{X} \rightarrow \mathbb{R}$  for all  $h \in \mathcal{H}$ , and the prediction is the sign of the output, i.e.  $\hat{y} = \text{sgn}(h(x))$ .
- We will look at zero-one loss function, i.e.  $\ell_{0-1}((x, y), h) = 1(\text{sgn}(h(x)) \neq y)$ . Note that we can re-express the loss function as

$$\ell_{0-1}((x, y), h) = \frac{1 - \text{sgn}(h(x))y}{2}.$$

- Think about the Rademacher complexity of  $\ell_{0-1}$  loss function, i.e. considering the family of functions  $\mathcal{F} = \{z = (x, y) \mapsto \ell_{0-1}((x, y), h) : h \in \mathcal{H}\}$ .
- Define  $Q$  to be the set of all possible outputs on our dataset:  
 $Q = \{(\text{sgn}(h(x^{(1)})), \dots, \text{sgn}(h(x^{(n)}))) \mid h \in \mathcal{H}\}$ .
- Computing the *Rademacher Complexity* we have

$$\begin{aligned}\mathcal{R}_n(\mathcal{F}) &= \mathbb{E}_{\varepsilon_1, \dots, \varepsilon_n} \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \frac{1 - \text{sgn}(h(x^{(i)}))y_i}{2} \right] \\ &= \frac{1}{2} \mathbb{E}_{\varepsilon_1, \dots, \varepsilon_n} \left[ \sup_{v \in Q} \frac{1}{n} \langle \varepsilon, v \rangle \right]\end{aligned}$$

- For any particular  $v \in Q$ , notice that  $\langle \varepsilon, v \rangle$  is a sum of bounded random variables, so we can use *Hoeffding's inequality* to obtain

$$\Pr \left[ \frac{1}{n} \langle \varepsilon, v \rangle \geq t \right] \leq \exp(-nt^2/2)$$

- Taking the union bound over  $v \in Q$ , we see that

$$\Pr \left[ \exists v \in Q \text{ such that } \frac{1}{n} \langle \varepsilon, v \rangle \geq t \right] \leq |Q| \exp(-nt^2/2).$$

- Thus, with probability at least  $1 - \delta$ , it is true that

$$\sup_{v \in Q} \frac{1}{n} \langle v, \varepsilon \rangle \leq \sqrt{\frac{2(\log |Q| + \log(2/\delta))}{n}}$$

- Similarly, we can show that  $\mathbb{E} \left[ \sup_{v \in Q} \frac{1}{n} \langle v, \varepsilon \rangle \right] \leq O \left( \sqrt{\frac{\log |Q| + \log(2/\delta)}{n}} \right)$  holds.
- *VC dimension* is one way to deal with bounding the size of  $Q$ .
- However, it has some limitations because will always end up with a bound that depends somehow on the dimension.

# Thanks!

✉ [marcelo.ortiz@emory.edu](mailto:marcelo.ortiz@emory.edu)

🔗 [marcelortiz.com](http://marcelortiz.com)

🐦 [@marcelortizv](https://twitter.com/marcelortizv)

## References

- Wainwright, M. J. (2019). *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge University Press.
- Ma, T. (2022). *Lecture Notes for Machine Learning Theory (CS229M/STATS214)*.