

Empirical Processes Reading Group Notes

Manu Navjeevan

March 28, 2021

Contents

1 Math Review	1
1.1 Vector Spaces and Norms	1
1.2 Topology and Continuity	2
1.3 Probability Spaces and Outer Measure	5
2 Weak Convergence	6
2.1 Weak Convergence in Space of Bounded Functions	11
3 Empirical Processes	16
3.1 Maximal Inequality	19
3.2 Chaining and Inequalities for Infinite Classes	22
3.3 Symmetrization	25
3.4 Glivenko-Cantelli	28

1 Math Review

1.1 Vector Spaces and Norms

Definition 1.1 (Vector Space). A vector space X is a set of elements with two operations, addition (+) and scalar multiplication (\cdot), and an additive identity $\mathbf{0} \in X$ satisfying:

1. $x + y = y + x$
2. $(x + y) + z = x + (y + z)$
3. $\mathbf{0} + x = x, \forall x \in X$
4. $\alpha(x + y) = \alpha x + \alpha y$
5. $(\alpha + \beta)x = \alpha x + \beta x$
6. $(\alpha\beta)x = \alpha(\beta x)$
7. $0x = \mathbf{0}$ and $1x = x$

Examples include \mathbb{R}^K and $\mathcal{C}[a, b]$, the set of all continuous functions from $[a, b] \rightarrow \mathbb{R}$.

Definition 1.2 (Norm). Let X be a vector space. A norm is a functional, $\|\cdot\| : X \rightarrow \mathbb{R}$ satisfying

1. $\|x\| \geq 0, \forall x \in X$ and $\|x\| = 0$ if and only if $x = \mathbf{0}$.
2. $\|x + y\| \leq \|x\| + \|y\|$ (Triangle Inequality)
3. $\|\alpha x\| = |\alpha| \|x\|, \forall \alpha \in \mathbb{R}, x \in X$

Examples of norms include the ℓ^p norms on \mathbb{R}^K or the sup-norm on the space of all bounded, real valued, functions. On \mathbb{R}^K all norms are equivalent, which is to say that for any two norms $\|\cdot\|_1, \|\cdot\|_2$ there are constants C_1, C_2 such that $C_1\|\cdot\|_2 \leq \|\cdot\|_1 \leq C_2\|\cdot\|_2$. However, this is not generally the case for functional vector spaces. For example on $\mathcal{C}[a, b]$ there is no constant c such that, for all f :

$$\sup_{x \in [a, b]} f(x) = \|f\|_\infty \leq c\|f\|_2 = \left(\int_a^b f^2(x) dx \right)^{1/2}.$$

Closely related to a norm is the concept of a metric, which is a way of defining a distance on a space.

Definition 1.3 (Metric). Let X be a vector space. A metric (or distance metric) on X is a functional $d(x, y) : X \times X \rightarrow \mathbb{R}$ satisfying:

1. $d(x, y) \geq 0, \forall x, y$ and $d(x, y) = 0 \iff x = y$
2. $d(x, y) = d(y, x)$
3. $d(x, y) \leq d(x, z) + d(z, y), \forall x, y, z$

It is straightforward to verify that, given a norm on a vector space X , we can generate a valid metric:

$$d_{\|\cdot\|}(x, y) := \|x - y\|.$$

We return to these concepts when discussing a topology.

1.2 Topology and Continuity

A topology is a general structure under which we can discuss concepts such as convergence and continuity. We can start with a general structure and then discuss spaces where the topology is generated by a metric (or norm).

Definition 1.4 (Topology). A topology on a set X is a collection of subsets of X , $\tau \subset 2^X$ satisfying:

1. $\emptyset, X \in \tau$.
2. τ is closed under finite intersections, if $\{A_k\}_{k=1}^K \in \tau$ then $\bigcap_{k=1}^K A_k \in \tau$.
3. τ is closed under arbitrary unions. For any index set I , if $\{A_k\}_{k \in I} \in \tau$ then $\bigcup_{k \in I} A_k \in \tau$.

The elements of $A \in \tau$ are called open sets. A set, B , is closed if its complement is in τ , $B^c \in \tau$.

Some simple examples include the trivial topology, $\tau = \{X, \emptyset\}$ and the discrete topology $\tau = 2^X$. Given a topology, we can define some familiar terms:

Definition 1.5 (Interior). For a subset $A \subseteq X$, the interior of A , denoted A° , is the largest open set included in A (where largest is defined under the usual subset ordering). We can also express this as the union of all open sets contained by A .

$$A^\circ = \bigcup \{B : B \in \tau, B \subseteq A\}.$$

Note that a set is open if and only if $A = A^\circ$.

Definition 1.6 (Closure). For a subset $A \subseteq X$, the closure of A , denoted \bar{A} , is the smallest closed set that covers A . We can express this as the intersection of all closed sets containing A :

$$\bar{A} = \bigcap \{B : B^c \in \tau, A \subseteq B\}.$$

By De-Morgan's law and closure of the topology under arbitrary union we can see that this intersection always gives a closed set. A set is closed if and only if $A = \bar{A}$.

Lemma 1.1. Suppose $x \in \bar{A}$, then for every neighborhood of x , V_x , we have that $V_x \cap A \neq \emptyset$.

Proof. Let $x \in \bar{A}$ and suppose for some neighborhood V_x of x we have that $V_x \cap A = \emptyset$. Then we know that $V_x^\circ \cap A = \emptyset$. Take $\tilde{A} = \bar{A} \cap (V_x^\circ)^c$. We can verify that this is a smaller closed set that also contains A . \square

Definition 1.7 (Boundary). The boundary of a set A , denoted δA , is $\bar{A} \setminus A^\circ$.

A useful concept when talking about convergence under a topology is that of a neighborhood of a point $x \in X$.

Definition 1.8 (Neighborhood). For a point $x \in X$ a set V is a neighborhood of x if $x \in V^\circ$.

We can now use the topology to define limit points and convergence.

Definition 1.9 (Limit Point). A point $x \in X$ is a limit point of a set $A \subseteq X$ if, for every neighborhood V of x ,

$$A \cap (V \setminus \{x\}) \neq \emptyset.$$

In other words, every neighborhood of x intersects with A at a point other than x . Let A' be the set of all limit points of $A \subseteq X$.

Lemma 1.2. If S is a subset of X , then $\bar{S} = S \cup S'$.

Proof. First show that $\bar{S} \subseteq S \cup S'$. Let $x \in \bar{S}$. If $x \in S$ then we are done. Otherwise, suppose $x \in \bar{S} \setminus S$. This means that for all V_x we have that $S \cap V_x = S \cap (V_x \setminus \{x\})$. By the result of Lemma 1.1, we have that $V_x \cap S \neq \emptyset$. So, $x \in S'$.

Now suppose that $x \in S \cup S'$. Clearly if $x \in S$ then $x \in \bar{S}$. Suppose then that $x \in S' \setminus S$ but $x \notin \bar{S}$. Let \tilde{S} be any closed set containing S , that is $S \subseteq \tilde{S}$. For sake of contradiction, suppose that $x \notin \tilde{S}$ (x is a limit point of S that is not in \tilde{S}). Because \tilde{S} is closed we know that $\tilde{S}^c \in \tau$. Further, we know that $x \in \tilde{S}^c$ so that \tilde{S}^c is a neighborhood of x . Since x is a limit point of S , we know that $\tilde{S}^c \cap S = \tilde{S}^c \cap S \setminus \{x\} \neq \emptyset$. However, we also know that $S \subseteq \tilde{S}$ so we have a contradiction. Therefore, it must be that $x \in \bar{S}$ which completes the proof. \square

Lemma 1.3 (Characterization of Closed Sets). A set is closed if and only if it contains all of its limit points.

Proof. This is a consequence of Lemma 1.2 and the fact that A is closed if and only if $\bar{A} = A$. \square

Definition 1.10 (Convergence). We say a sequence $\{x_n\}_{n=1}^\infty$ converges to a point $x \in X$ if for every neighborhood V_x of x , there exists a number M such that for all $m \geq M$, $x_m \in V_x$.

Note that under the trivial topology $\tau = \{\emptyset, X\}$ all sequences converge to any point $x \in X$ whereas under the discrete topology on \mathbb{R} , $\tau = 2^\mathbb{R}$, no sequence converges.

Definition 1.11 (Continuity). Let (\mathcal{X}, τ_1) and (\mathcal{Y}, τ_2) be two topological spaces and $f : \mathcal{X} \rightarrow \mathcal{Y}$. We say f is continuous if $f^{-1}(A) \in \tau_1$ for all $A \in \tau_2$. That is, a continuous function maps open sets to open sets.

We can now get ready to combine the notions of continuity and convergence coming from a topology with the notions that we are familiar with from metric spaces. First, we need to define the topology generated by a metric.

Definition 1.12 (Generated Topology). Let \mathcal{A} be a collection of subsets of X . The topology generated by \mathcal{A} , $\langle \mathcal{A} \rangle$ is the smallest topology that contains \mathcal{A} :

$$\langle \mathcal{A} \rangle = \bigcap \{ \tau : \mathcal{A} \subseteq \tau \}.$$

We will then define the topology generated by a metric as the topology generated by the collection of open balls $B(x, \epsilon)$.

Definition 1.13 (Open Ball). Let $d(x, y)$ be a metric on a vector space X . For any point $x \in X$ define the open ball of size ϵ around x as:

$$B(x, \epsilon) = \{ y : d(x, y) \leq \epsilon \}.$$

In a metric space, we consider the topology generated by all the open balls $\tau_d = \langle \{B(x, \epsilon) : x \in X, \epsilon > 0\} \rangle$. In fact, the set of open balls is a basis for this topology, which means that every open set A in τ_d and any point $x \in A$, there is an open ball B such that $x \in B \subseteq A$.¹ Many topological properties such as continuity or convergence can be verified by simply confirming the properties for all members of a basis for the topology. This ties together the “epsilon-delta” notions of continuity and convergence with the more topological versions given above.

For the rest of this subsection we will talk about separability and compactness, but give examples using normed-metric spaces instead of talking in generality about the topology.

Definition 1.14 (Dense Subset). A topological space (X, τ) has a dense subset \mathcal{A} if $\bar{\mathcal{A}} = X$. Equivalent, by Lemma 1.2, every point of X is either in \mathcal{A} or is a limit point of \mathcal{A} .

Informally, all points in X are either in \mathcal{A} or arbitrarily “close” to \mathcal{A} . As an example, in the standard topology on \mathbb{R} generated by the metric $d(x, y) = |x - y|$, the rationals \mathbb{Q} are dense. We also have that, for the set of continuous functions under the sup norm, the set of all polynomials is dense, which means that we can approximate a function arbitrarily well with them.

Definition 1.15 (Seperable Space). We say that a topological space (X, τ) is separable if it has a countable dense subset.

As we went over above, the real line with its standard topology is separable. The $L_p[a, b]$ spaces are also generally separable for $1 \leq p \leq \infty$. However L_∞ is not separable, which will cause issues (this is not the example below).

Example 1.1 (Bounded functions with the sup norm is not seperable). Let $\{f_i\}_{i \in \mathbb{N}}$ be a countable set of functions on $B_\infty[a, b]$. Let $\{q_i\}_{i \in \mathbb{N}}$ be some counting of the rational numbers between a and b . Let \tilde{f} be some function that is equal to 0 except on the rational numbers. For each rational number q_i define

$$\tilde{f}(q_i) = \begin{cases} 1 & \text{if } f_i(q_i) \leq 0 \\ -1 & \text{if } f_i(q_i) > 0 \end{cases}.$$

We can see that \tilde{f} is bounded (and integrates to 0), but it is at least distance one from each function in $\{f_i\}_{i \in \mathbb{N}}$.

Initially I thought this example would work for $L_\infty[a, b]$, but this only forces a difference on a set of measure 0 and I believe L_∞ works with an essential supremum norm.

Another important/useful concept is that of compactness. The general notion is given below:

Definition 1.16 (Compact Set). A set A is compact if for every collection of open sets $\{G_i\}$ such that $A \subset \bigcup G_i$, there is a finite subcollection that also covers A .

Example 1.2. The real-line is not compact. Consider the open cover $\{(n, n+1) | n \in \mathbb{Z}\}$

Example 1.3. The interval $(0, 1]$ is not compact. Consider the open cover $\{(1/n, 1 + 1/n) | n \in \mathbb{N}\}$

Theorem 1.1 (Heine-Borel). For a subset S of the Euclidean Space², \mathbb{R}^n , the following statements are equivalent:

- S is closed and bounded
- S is compact

Compactness is nice because of various extreme value theorems that ensure that a supremum or infimum is attained. Heine-Borel gives a nice way of characterizing compactness for Euclidean Spaces, but in general there is no equivalent result for general metric spaces. We have to strengthen the boundedness assumption.

Definition 1.17 (Totally Bounded). A set \mathcal{A} is totally bounded if for each $\epsilon > 0$ there exists a finite sequence $\{a_1, \dots, a_n\}$ such that for $B_i = \{a \in \mathcal{A} : \|a - a_i\| \leq \epsilon\}$, $\bigcup_{i=1}^N B_i$ covers \mathcal{A} .

¹In fact, the set of all open balls with rational ϵ is a basis for the topology

²That is the space \mathbb{R}^n equipped by the topology generated by the standard distance metric

Intuition: For any precision ϵ , you can find a finite set of points that describe \mathcal{A} arbitrarily well. (much more demanding in infinite dimensions than just bounded).

Theorem 1.2. *In a complete metric space, the following are equivalent:*

- \mathcal{A} is a compact subset
- \mathcal{A} is closed and totally bounded
- Every sequence in \mathcal{A} has a convergent subsequence which converges to a point in \mathcal{A} .

For a compact set T , let $C(T)$ be the set of continuous functions from T to \mathbb{R} equipped with the sup norm. We may want to characterize when a subset K of $C(T)$ is compact.

Definition 1.18 (Equicontinuous). A set of functions $K \subseteq C(T)$ is equicontinuous if for every $t_0 \in T$ and $\epsilon > 0$ there is a $\delta > 0$ such that $|f(t) - f(t_0)| < \epsilon$ whenever $\|t - t_0\| < \delta$ **for all** $f \in K$.

This is a bit like to uniform continuity but adapted a bit to deal with a function space.

Theorem 1.3 (Arzela-Ascoli). *If T is compact, then $K \subseteq C(T)$ is compact (under the sup-norm) if and only if K is bounded and equicontinuous.*

This concludes our discussion of topology and continuity. We now review measurability.

1.3 Probability Spaces and Outer Measure

Definition 1.19 (Sigma Algebra). A collection of subsets \mathcal{F} is a sigma-algebra (or sigma-field) if it contains the whole set and is closed under complement and under countable union.

Definition 1.20 (Borel Sigma Algebra). For any collection of sets \mathcal{A} , we call the smallest sigma algebra containing \mathcal{A} , $\sigma(\mathcal{A})$, the sigma algebra generated by \mathcal{A} . The Borel sigma algebra on a topological space is the sigma algebra generated by all the open sets, $\mathcal{B}(X) = \sigma(\tau)$.

The Borel sigma algebra is useful as it makes all continuous functions measurable (defined below).

Definition 1.21 (Probability Space). A probability space is a triple $(\Omega, \mathcal{F}, \mathbb{P})$ consisting of a set of elements Ω , a sigma algebra on Ω , \mathcal{F} , and a probability measure $\mathbb{P} : \mathcal{F} \rightarrow [0, 1]$ satisfying:

1. $\mathbb{P}(A) \geq \mathbb{P}(\emptyset) = 0$ [Non-negativity]
2. If $A_i \in \mathcal{F}$ is a countable sequence of disjoint sets then $\mathbb{P}(\bigcup_i A_i) = \sum_i \mathbb{P}(A_i)$
3. $\mathbb{P}(\Omega) = 1$.

A measurable function between two spaces equipped with sigma algebra's is simply one that maps measurable sets to measurable sets, similar to the definition of a continuous function.

Definition 1.22 (Measurable Map). A function $f : (\mathcal{X}, \mathcal{A}) \rightarrow (\mathcal{Y}, \mathcal{B})$ is measurable if $f^{-1}(B) \in \mathcal{A}$ for all $B \in \mathcal{B}$

Lemma 1.4 (Lemma 1.3.1 VdV& W). *The Borel σ -field on a metric space \mathbb{D} is the smallest σ -field that makes all elements of $C_b(\mathbb{D})$ measurable (with respect to the Borel sets on \mathbb{R}).¹*

Proof. For any closed set F , F is the null set $\{x : f(x) = 0\}$ of the continuous, bounded function, $x \mapsto d(x, F) \wedge 1$. Since the singleton $\{0\}$ is a closed set in \mathbb{R} (all metric spaces are Hausdorff), F must be in the sigma algebra on \mathbb{D} to make $d(x, F) \wedge 1$ measurable. Since all the closed sets generate the Borel σ -field (because σ -fields are closed under complement), all Borel sets must be included in the sigma-algebra on \mathbb{D} . \square

¹ $C_b(\mathbb{D})$ is the set of all continuous bounded functions from $\mathbb{D} \rightarrow \mathbb{R}$, where \mathbb{R} is endowed with the standard topology on the real line

Given this, we can abstractly think about a random variable as a measurable map from a probability space into another measurable space (typically the real-line). Measurability ensures that things like expectations and probabilities of random variables are well defined.

However, measurability becomes a problem when we are dealing with random functions. For example, if X is a map from a probability space to $L_\infty[a, b]$, the Borel-sigma algebra on $L_\infty[a, b]$ is quite large (its not separable). This means that measurable sets in $L_\infty[a, b]$ may not map back to measurable sets on the probability space $\Omega, \mathcal{F}, \mathbb{P}$.

This is a problem because L_∞ is typically a useful space to work in for empirical process theory. So we have to find a way to relax measurability. This means that we work with outer expectations and probabilities:

Definition 1.23 (Outer Measure and Inner Measure). Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space $T : \Omega \rightarrow \mathbb{R}$. Define the outer expectation:

$$\mathbb{E}^*[T] = \inf \{ \mathbb{E}[U] : T \leq U, U \text{ is measurable} \}.$$

and the inner expectation:

$$\mathbb{E}_*[T] = \sup \{ \mathbb{E}[U] : U \leq T, U \text{ is measurable} \}.$$

We can use this to define inner and outer probability measures by restricting T to be the indicator function for an arbitrary set B . Inner and outer expectations are generally nicely behaved but they require modified versions of dominated and monotone convergence and Fubini's theorem breaks down.

2 Weak Convergence

We can now talk about weak convergence of random variables. Let X_n be a real-valued random variable with cdf $F_n(t)$ and let X be a random variable with cdf $F(t)$. The typical definition of weak convergence is that $X_n \xrightarrow{L} X$ if $F_n(t) \rightarrow F(t)$ pointwise at all continuity points of F . This is not super general for non-real valued random maps.

Theorem 2.1 (Portmanteau). *For real random variables $X_n \xrightarrow{L} X$ is equivalent to:*

- $\mathbb{E}[g(X_n)] \rightarrow \mathbb{E}[g(X)]$ for all bounded continuous functions.
- For all open sets G , $\liminf \mathbb{P}(X_n \in G) \geq \mathbb{P}(X \in G)$.
- For all closed sets K , $\limsup \mathbb{P}(X_n \in K) \leq \mathbb{P}(X \in K)$.

This motivates the theory of weak convergence for general metric spaces. Let \mathbb{D} be a complete metric space with metric d . We can equip \mathbb{D} with its Borel-sigma algebra as defined in Definition 1.20 and a tight probability measure as defined in Definition 2.1. Let $C_b(\mathbb{D})$ be the set of all continuous and bounded real functions on \mathbb{D} . If X is a random variable, $X : (\Omega, \mathcal{F}, \mathbb{P}) \rightarrow \mathbb{D}$ then its law is given $L = \mathbb{P} \circ X^{-1}$.

Definition 2.1 (Tight Probability Measure). A probability measure is tight if for every $\epsilon > 0$ there is a compact set K_ϵ such that $\mathbb{P}(K_\epsilon) \geq 1 - \epsilon$

This is a generalization of bounded in probability I believe.

Definition 2.2 (Borel Law). For a random variable X , we say that X has a Borel Law L if

$$\mathbb{P}(X \in A) = \int_A dL.$$

for all Borel sets A .

Given this setup, we can now define weak convergence:

Definition 2.3 (Weak Convergence). Let $(\Omega_n, \mathcal{F}_n, \mathbb{P}_n)$ be a sequence of probability spaces and $X_n : \Omega_n \rightarrow \mathbb{D}$. Then we say that $X_n \xrightarrow{L} X$ if:

$$\mathbb{E}^* [f(X_n)] \rightarrow \mathbb{E}[f(X)].$$

for every $f \in C_b(\mathbb{D})$

We can characterize this convergence using another Portmanteau theorem.

Theorem 2.2 (Portmanteau). *The following are equivalent:*

1. $X_n \xrightarrow{L} X$
2. $\liminf \mathbb{P}_*(X_n \in G) \geq \mathbb{P}(X \in G)$ for all open sets G .
3. $\limsup \mathbb{P}^*(X_n \in F) \leq \mathbb{P}(X \in F)$ for every closed set F .
4. $\lim P(X_n \in B) = P(X \in B)$ for every Borel set B with $P(X \in \delta B) = 0$.

Question: Is X supposed to have a Borel Law? Otherwise where do open and closed sets get tied into this? Is it from the notion of convergence?

Proof. This proof is in a few steps.

(4) \implies (3): Suppose that $\lim P(X_n \in B) = P(X \in B)$ for every Borel set B with $\mathbb{P}(X \in \delta B) = 0$. Let F be a closed set and let $F^\epsilon = \{x : d(x, F) < \epsilon\}$. The sets δF^ϵ are disjoint for different values of $\epsilon > 0$ (The boundary of this set is $\delta F^\epsilon = \{x : d(x, F) = \epsilon\}$), so at most countably many of them can have nonzero L-measure (otherwise the measure of the entire space would be infinite). Choose a sequence $\epsilon_m \downarrow 0$ with $L(\delta F^{\epsilon_m}) = 0$ for each m (this is possible because only countably many ϵ have $L(F^\epsilon) \neq 0$). For a fixed m , by (4) we have that:

$$\limsup P^*(X_n \in F) \leq \limsup P^*(X_n \in \overline{F^{\epsilon_m}}) = L(\overline{F^{\epsilon_m}}).$$

letting $m \rightarrow \infty$ gives (3).

(3) \iff (2): Take any closed set F . Its complement F^c is open. If

$$\liminf \mathbb{P}_*(X_n \in F^c) \geq \mathbb{P}(X \in F^c).$$

Then

$$\begin{aligned} \limsup \mathbb{P}^*(X_n \in F) &\leq \liminf 1 - \mathbb{P}_*(X_n \in F^c) \\ &\leq 1 - \mathbb{P}(X \in F^c) \\ &= \mathbb{P}(X \in F) \end{aligned}$$

a symmetric argument shows the backwards direction.

(2)+(3) \implies (4): This is straightforward if we recall that, for any set with $L(\delta B) = 0$ we have that $L(B) = L(\bar{B})$. Then we bound the limsup by the liminf:

$$\limsup \mathbb{P}^*(X \in B) \leq \limsup \mathbb{P}(X \in \bar{B}) \leq \mathbb{P}(X \in \bar{B}) = \mathbb{P}(X \in B) \leq \liminf \mathbb{P}_*(X_n \in B).$$

which gives (4).

(1) \implies (2): Take any G open and define the sequence of functions:

$$f_m(x) := \min(1, m \cdot d(x, G^c))$$

Notice that $f_m(x) \in C_b(\mathbb{D})$ and $f_m(x) \leq \mathbb{1}\{x \in G\}$. So, for every m we have that

$$\begin{aligned} \liminf \mathbb{P}_*(X \in G) &= \liminf \mathbb{E}_* [\mathbb{1}\{X \in G\}] \\ &\geq \liminf \mathbb{E}_* [f_m(X)] \\ &\geq \mathbb{E}[f_m(X)] \end{aligned}$$

since $f_m(x) \uparrow \mathbb{1}\{X \in G\}$ by monotone convergence we get the result in (2).

Question: How do we know from weak convergence that this sequence converges in inner expectation?

By VdV and Wellner, weak convergence implies (is equivalent to) $\liminf \mathbb{E}_* [f(X_n)] \geq \mathbb{E} [f(X)]$ for every bounded, Lipschitz continuous, non-negative f . I think the argument for why this is the case goes: Let $f \geq 0$ be bounded and continuous. Then by weak convergence

$$\limsup \mathbb{E}^* [-f(X_n)] = \mathbb{E} [-f(X)].$$

Taking negatives will give:

$$\liminf \mathbb{E}_* [f(X_n)] \geq -\limsup \mathbb{E}^* [-f(X_n)] = \mathbb{E} [f(X)].$$

In any case, $f_m(X)$ is Lipschitz continuous which gives the result.

(2) \implies (1): (SKETCH)

- Suppose $f(x) \geq 0$ is continuous and bounded
- Approximate it from above and below by indicator functions of open sets.

□

Weak convergence is nice because it gives the continuous mapping theorem.

Theorem 2.3 (Continuous Mapping Theorem). *Let $g : \mathbb{D} \rightarrow \mathbb{E}$ be continuous at every point $\mathbb{D}_0 \subseteq \mathbb{D}$. If $X_n \xrightarrow{L} X$ and $\mathbb{P}(X \in \mathbb{D}_0) = 0$ then $g(X_n) \xrightarrow{L} g(X)$.*

Proof. (Without Discontinuity Points): Let $Z_n = g(X_n)$ and $Z = g(X)$. We want to show that $\mathbb{E}^* [f(Z_n)] \rightarrow \mathbb{E} [f(Z)]$ for all $f \in C_b(\mathbb{D}; \mathbb{E})$.

$$\lim_{n \rightarrow \infty} \mathbb{E} [f(Z_n)] = \lim_{n \rightarrow \infty} \mathbb{E} [f(g(X_n))] = \mathbb{E} [f(g(X))] = \mathbb{E} [f(Z)].$$

The main step here is weak convergence of X_n and the stability of $C_b(\mathbb{D}; \mathbb{E})$ under composition.

(With Discontinuity Points, from VdV&W): The set D_g of all points at which g is discontinuous can be written

$$D_g = \bigcup_{m=1}^{\infty} \bigcap_{k=1}^{\infty} \{x : \exists y, z \in B(x, 1/k) \text{ with } d_{\mathbb{E}}(g(y), g(z)) > 1/m\}.$$

Intuition: Recall that g is continuous at x if for every $m \in \mathbb{N}$ there exists a $k \in \mathbb{N}$ such that¹

$$y, z \in B(x, 1/k) \implies d_{\mathbb{E}}(g(y), g(z)) < 1/m$$

If the function is not continuous at x you can find a counterexample for some $k, m \in \mathbb{N}$.

Let $G_k^m = \{x : \exists y, z \in B(x, 1/k) \text{ with } d_{\mathbb{E}}(g(y), g(z)) > 1/m\}$. Every G_k^m is open (if x is in G_k^m the points just around x will be as well so that we can write G_k^m as a union of open balls) so that D_g is a Borel set. For every closed F we then have that:

$$\overline{g^{-1}(F)} \subseteq g^{-1}(F) \cup D_g.$$

By Portmanteau:

$$\begin{aligned} \limsup \mathbb{P}^* (g(X_n) \in F) &\leq \limsup P^* (X_n \in \overline{g^{-1}(F)}) \leq \mathbb{P} (X \in \overline{g^{-1}(F)}) \\ &= \mathbb{P} (X \in g^{-1}(F)) \\ &= \mathbb{P} (g(X) \in F) \end{aligned}$$

Applying Portmanteau again gives weak convergence. □

¹Topologically, this is saying that the inverse map of every open neighborhood of $f(x)$ is an open neighborhood of x

Example 2.1. Take $\mathbb{G}_n \in L^\infty(\mathbb{R})$:

$$\mathbb{G}_n(t) := \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(\mathbb{1}\{X_i \leq t\} - \mathbb{E}[\mathbb{1}\{X \leq t\}] \right)$$

and suppose that $\mathbb{G}_n \xrightarrow{L} \mathbb{G}$ where \mathbb{G} is some other element of $L^\infty(\mathbb{R})$. Let $Z : L^\infty(\mathbb{R}) \rightarrow \mathbb{R}$ be defined as:

$$Z(f) := \sup_t |f(t)|.$$

this function is continuous. Applying the continuous mapping theorem to Z allows us to build uniform confidence intervals.

Let $\gamma_{1-\alpha}$ be the $1 - \alpha$ quantile of $Z := \sup_t |\mathbb{G}(t)|$ and construct a confidence interval (at each t):

$$\left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\} - \gamma_{1-\alpha}/\sqrt{n}, \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\} + \gamma_{1-\alpha}/\sqrt{n} \right].$$

Then:

$$\begin{aligned} & \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\} - \gamma_{1-\alpha}/\sqrt{n} \leq \mathbb{E}[\mathbb{1}\{X \leq t\}] \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\} + \gamma_{1-\alpha}/\sqrt{n} : \text{ for all } t \right) \\ &= \mathbb{P} \left(|\mathbb{G}_n(t)| \leq \gamma_{1-\alpha} \forall t \right) \\ &= \mathbb{P} \left(\sup_t |\mathbb{G}_n(t)| \leq \gamma_{1-\alpha} \right) \end{aligned}$$

But by continuous mapping theorem and Portmanteau, if $\mathbb{P}(\sup_t |\mathbb{G}| = \gamma_{1-\alpha}) = 0$:

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(\sup_t |\mathbb{G}_n(t)| \leq \gamma_{1-\alpha} \right) = \mathbb{P} \left(\sup_t |\mathbb{G}(t)| \leq \gamma_{1-\alpha} \right) = 1 - \alpha.$$

This sort of argument can be applied more generally to functions $\mathbb{G}_n(t) = \hat{m}(t) - m(t)$ to construct uniform confidence intervals.

This shows the usefulness of Portmanteau and Continuous Mapping Theorem. For finite dimension vectors we can use the central limit theorem to establish weak convergence to a normal distribution. However, when X_n is a random element in L^∞ it may be harder to show that $X_n \rightsquigarrow X$ for some other $X \in L^\infty$.

- Don't want to check $\mathbb{E}[f(X_n)] \rightarrow \mathbb{E}[f(X)]$ for all $f \in C_b(L^\infty)$ [There are at least 20 functions in this class]

Instead we will try to use the structure of L^∞ to show the result.

Definition 2.4 (Asymptotic Tightness). A sequence X_n of random maps is asymptotically tight if for every $\epsilon, \delta > 0$ there is a compact K_ϵ such that

$$\liminf P_\star \left(X_n \in K_\epsilon^\delta \right) \geq 0.$$

where $K_\epsilon^\delta = \{y \in \mathbb{D} : d(y, K_\epsilon) < \delta\}$ is the “ δ -enlargement” around K_ϵ .

Definition 2.5 (Asymptotic Measurability). A sequence X_n of random maps is asymptotically measurable if for all $f \in C_b(\mathbb{D})$:

$$\mathbb{E}^\star f(X_n) - \mathbb{E}_\star f(X_n) \rightarrow 0.$$

We would like for a sequence X_n that weakly converges to an element X to inherit some properties from X :

Lemma 2.1 (Lemma 1.3.8 VdV& W). *The following are true:*

1. If $X_n \xrightarrow{L} X$ then X_n is asymptotically measurable
2. If $X_n \xrightarrow{L} X$ then X_n is asymptotically tight if and only if X is tight.

Proof. (1): Take any function $f \in C_b(\mathbb{D})$. By definition of weak convergence we know that

$$\lim \mathbb{E}^* [f(X_n)] = \mathbb{E}[f(X)] \quad \text{and} \quad \lim \mathbb{E}^* [-f(X_n)] = \mathbb{E}[-f(X)].$$

I think we should have that $-\mathbb{E}_* [f(X_n)] \geq \mathbb{E}^* [-f(X_n)]$ for any f which give the result (I think this holds with equality but I leave it as an inequality since this is all we need for the result).

(2): Fix $\epsilon > 0$. If X is tight then there is a compact K with $\mathbb{P}(X \in K) > 1 - \epsilon$. By Portmanteau:

$$\liminf \mathbb{P}_*(X_n \in K^\delta) \geq \mathbb{P}(X \in K^\delta).$$

which is larger than $1 - \epsilon$ for every $\delta > 0$.

Conversely, suppose that X_n is asymptotically tight. Then there exists a compact K with $\liminf \mathbb{P}_*(X_n \in K^\delta) \geq 1 - \epsilon$. By Portmanteau,

$$1 - \epsilon \leq \liminf \mathbb{P}_*(X_n \in K^\delta) \leq \limsup \mathbb{P}^*(X_n \in \overline{K^\delta}) \leq \mathbb{P}(X \in \overline{K^\delta}).$$

Let $\delta \rightarrow 0$ by monotone convergence to complete the result. ² □

The converse is not generally true. Let $X_n = -1$ if n is odd and $X_n = 1$ if n is even. This sequence is asymptotically measurable and asymptotically tight but clearly does not converge. However, it does converge among a subsequence. This is the idea behind the partial converse to this theorem provided by Pohorov's Theorem.

Theorem 2.4 (Pohorev's Theorem, Theorem 1.3.9 VdV& W). *Let X_n be an asymptotically tight and asymptotically measurable sequence. Then there is a subsequence X_{n_j} that converges weakly to a tight Borel law.*

Now a review problem

Example 2.2 (Problem 7; Ch 1.3 VdV& W). Let X_n be a sequence of random elements in \mathbb{D} and $g : \mathbb{D} \rightarrow \mathbb{E}$ a continuous function. Want to show that:

1. If X_n is asymptotically tight then $g(X_n)$ is asymptotically tight.
2. If X_n is asymptotically measurable then $g(X_n)$ is asymptotically measurable.

Proof. 1) Suppose that X_n is asymptotically tight. Fix $\epsilon > 0$. We know that there exists a compact set K such that, $\forall \delta_1 > 0$

$$\liminf \mathbb{P}_*(X_n \in K^{\delta_1}) \geq 1 - \epsilon.$$

The event $\{X_n \in K^{\delta_1}\}$ is a subset of the event that $\{g(X_n) \in g(K^{\delta_1})\}$ so

$$\liminf \mathbb{P}_*(g(X_n) \in g(K^{\delta_1})) \geq \liminf \mathbb{P}_*(X_n \in K^{\delta_1}) \geq 1 - \epsilon.$$

To finish recall that $g(K)$ is a compact set and choose δ_1 such that $g(K^{\delta_1}) \subseteq g(K)^\delta$ (always possible to do so by continuity of g).

²This proof relies on compact sets being closed in metric spaces. The proof of this is as follows: Let A be compact in a metric space. We wish to show that A is closed. Take a point $x \in X \setminus A$. To show that A is closed, we want to show that there is an open neighborhood of x that is not in A (this will show that A contains all of its limit points). For every $a \in A$, let $U_a = B(a, \frac{d(a,x)}{2})$ and $V_a = B(x, \frac{d(a,x)}{2})$. By triangle inequality, U_a and V_a are disjoint. The union of all the sets U_a for all points $a \in A$ is an open cover of A . By compactness of A , we can get a finite subcover U_{a_1}, \dots, U_{a_n} . But then $V_{a_1} \cap \dots \cap V_{a_n}$ is an open neighborhood of x that is disjoint from A . So A is closed. Actually this argument holds in general Hausdorff spaces.

2) Suppose that X_n is asymptotically measurable. This means that, for any $f \in C_b(\mathbb{D})$:

$$\mathbb{E}^* [f(X_n)] - \mathbb{E}_* [f(X_n)] \rightarrow 0.$$

Let $\tilde{f} \in C_B(\mathbb{E})$. For any continuous $g : \mathbb{D} \rightarrow \mathbb{E}$, $f \circ g$ is a continuous and bounded function from $\mathbb{D} \rightarrow \mathbb{R}$. This completes the proof. \square

2.1 Weak Convergence in Space of Bounded Functions

So far, we have defined weak convergence. But, how do we show that $X_n \xrightarrow{L} X$? In \mathbb{R}^K we have the central limit theorem, but no direct analog for random maps into L^∞ .

First, some definitions.

Definition 2.6 (Marginal Random Variable). Let X_n be a random map into $L^\infty(T)$ (the space of all bounded functions from $T \rightarrow \mathbb{R}$). Then, $X_n(t)$ is the marginal distribution of X_n at t . We can view $X_n(t)$ as the composition of X_n with π_t or directly as a real-valued random variable.

A general strategy will be to deal with the marginals directly. By the central limit theorem, we have conditions for the weak convergence of $X_n(t)$. Want to know what these results imply for the random map X_n .

Lemma 2.2 (Lemma 1.5.1, VdV&W). *Let $X_n : \Omega \rightarrow L^\infty(T)$ be asymptotically tight. Then it is asymptotically measurable if and only if $X_n(t)$ is asymptotically measurable for every $t \in T$.*

Lemma 2.3 (Lemma 1.5.3, VdV&W). *Let X and Y be tight Borel measurable maps into $L^\infty(T)$. Then $X \stackrel{L}{=} Y$ if and only if $X(t) \stackrel{L}{=} Y(t)$ for all $t \in T$.*

Theorem 2.5 (Theorem 1.5.4, VdV&W). *Let $X_n : \Omega_n \rightarrow L^\infty(T)$ be arbitrary. Then X_n weakly converges to a tight limit if and only if X_n is asymptotically tight and the marginals $(X_n(t_1), \dots, X_n(t_k))$ converge weakly to a limit for every finite subset t_1, \dots, t_k .*

Proof. Forward direction is simple, backwards direction requires more work:

(\implies) Suppose that $X_n \xrightarrow{L} X$ and X is tight. By Lemma 2.1, this means that X_n is asymptotically tight. Let $T_k : L^\infty(T) \rightarrow \mathbb{R}^k$ be the projection onto the coordinates t_1, \dots, t_k . This is a continuous function so by continuous mapping theorem we have convergence of the marginals for any finite collection.

(\impliedby) Suppose that X_n is asymptotically tight and the marginals converge. Then, by Lemma 2.2, X_n is asymptotically measurable. By Pohorov's theorem, there is a subsequence $X_{n_k} \xrightarrow{L} X$ for some X . Suppose $X_n \not\xrightarrow{L} X$. Then, there is a subsequence $X_{n'_k}$ that stays away from X (in law). However, the marginals converge. This means that the marginals of Y are the same as the marginals of X . By Lemma 2.3, $X \stackrel{L}{=} Y$. \square

Intuition: Why is Tightness + Convergence of Marginals Enough?

- Tightness: $P(X \in K) \geq 1 - \epsilon$ for some compact set K .
 - In a metric space, compact means for any $\epsilon > 0$ there are a finite set of points that approximate the whole set well.
 - * But! For a finite set of points we have convergence of marginals

Showing convergence of marginal distributions is straightforward by CLT. Next, we cover how to show tightness. Then Theorem 2.5 gives convergence of the entire process. To verify tightness we want a better description than the definition of asymptotic tightness. Two approaches

1. Finite Approximation \rightarrow simpler
2. Arzela-Ascoli Theorem \rightarrow larger interest (asymptotic equicontinuity)

2.1.1 Finite Approximation

The general idea here is that, for any $\epsilon > 0$, we can partition the index set T (as in $\ell^\infty(T)$) into a finite number of sets T_i so that the variation in each set is $< \epsilon$. Formally, for any $\eta > 0$,

$$\limsup_{n \rightarrow \infty} \mathbb{P} \left(\max_i \sup_{s, t \in T_i} |X_n(s) - X_n(t)| > \epsilon \right) < \eta.$$

Intuition: Why should we expect this to work?

- Tightness means that you concentrate on a compact set
 - Compact set is well described by a finite # of functions

Theorem 2.6 (Theorem 1.5.6 VdV&W). *A sequence of random maps $X_n \in \ell^\infty(T)$ is asymptotically tight if and only if $X_n(t)$ is asymptotically tight in \mathbb{R} for every t and, for all $\epsilon, \eta > 0$ there is a partition $T = \cup_{i=1}^n T_i$ such that*

$$\limsup_{n \rightarrow \infty} \mathbb{P}^* \left(\max_i \sup_{s, t \in T_i} |X_n(s) - X_n(t)| > \epsilon \right) < \eta \quad (\text{FA-1})$$

Proof. Cover sufficiency. Necessity follows from Theorem 1.5.7 in Van DerVaart and Wellner. Suppose that (FA-1) holds. Fix $\epsilon > 0$ and let the partition $T = \bigcup_{i=1}^k T_i$ satisfy (FA-1) for some $\eta > 0$. We want to show that $\sup_t |X_n(t)|$ is asymptotically tight. Then:

$$\begin{aligned} \limsup \mathbb{P}^* \left(\sup_{t \in T} |X_n(t)| > M \right) &\leq \limsup \mathbb{P}^* \left(\sup_{t \in T} > M, \text{ and (FA-1) holds} \right) \\ &\quad + \limsup \mathbb{P}^* \left(\text{(FA-1) doesn't hold} \right) \\ &\leq \limsup \mathbb{P}^* \left(\max_{1 \leq i \leq k} |x_n(t_i)| + \epsilon > M \right) + \eta \end{aligned}$$

Where in the last line we use the bounded variation within each set T_i and pick some arbitrary elements $t_i \in T_i$. Now note that each $X_n(t_i)$ is asymptotically tight by assumption so that $\max_{1 \leq i \leq k_i} |X_n(t_i)|$ is asymptotically tight.¹. This means that we can pick M so that

$$\limsup \mathbb{P}^* \left(\sup_t |X_n(t)| > M \right) < \eta.$$

or, to put this another way, for every $\eta > 0$ we can show that there is an M such that:

$$\limsup \mathbb{P}^* \left(\sup_t |X_n(t)| > M \right) < \eta.$$

So we have shown that $\sup_t |X_n(t)|$ is bounded in probability. Since $\sup_t |X_n(t)|$ is a map onto the real line, bounded in probability coincides with asymptotic tightness (Heine-Borel).

Now we want to construct a candidate compact set K for the process X_n . Fix $\zeta > 0$ and a sequence $\epsilon_n \downarrow 0$. First, pick an M such that

$$\limsup_{n \rightarrow \infty} \mathbb{P}^* \left(\sup_t |X_n(t)| > M \right) < \zeta.$$

¹Couple of quick arguments to get this one:

1. If each $X_{i,n}$ in $\{X_{i,n}\}_{i=1}^K$ is asymptotically tight then the vector $[X_1 \dots X_K]$ is asymptotically tight. This is because the Cartesian product of a finite number of compact sets is compact (with respect to the product topology).
2. If X_n is asymptotically tight and g is a continuous function then $g(X_n)$ is asymptotically tight. This is shown in Example 2.2 and basically follows from the fact that a continuous function applied to a compact set yields a compact set. The maximum operator is continuous.

we know such an M exists by the above argument. For each ϵ_m partition $T = \bigcup_{i=1}^{K(m)} T_i$ such that

$$\limsup_{n \rightarrow \infty} \mathbb{P}^* \left(\sup_{1 \leq i \leq K(m)} \sup_{s, t \in T_i} |X_n(s) - X_n(t)| > \epsilon_m \right) < \frac{\zeta}{2^m}.$$

For each ϵ_m let $\{z_1, \dots, z_{p(m)}\}$ be the set of functions in $\ell^\infty(T)$ that are constant on T_i and only take values $0, \pm\epsilon_m, \pm 2\epsilon_m, \dots, M$. It is only important for now that, for any m , $p(m)$ is finite (though large). Let

$$K_m = \bigcup_{i=1}^{p(m)} \overline{B}(z_i, \epsilon_m).$$

where $\overline{B}(z_i, \epsilon_m)$ is the closed ball of radius ϵ_m around z_i . Note that if $\sup_t |X_n(t)| \leq M$ and

$$\sup_{1 \leq i \leq p(m)} \sup_{s, t \in T_i} |X_n(s) - X_n(t)| \leq \epsilon_m$$

then $X_n \in K_m$. Let $K = \bigcap_{m=1}^{\infty} K_m$. Then K is closed and totally bounded. Closure follows because each K_m is closed (finite union of closed sets) and an arbitrary intersection of closed sets is closed (because the arbitrary union of open sets is open). To see totally bounded fix $\delta > 0$. Then for each $\epsilon_m < \delta$ we have that $K_m = \bigcup_{i=1}^{p(m)} \overline{B}(z_i, \epsilon_m)$. Since $K_m \supset K$ these balls cover K .

We now have a candidate K . We now want to show that, for every $\delta > 0$, $K^\delta \supset \bigcap_{i=1}^m K_i$ for some m . Suppose not. Then there is a sequence $\{z_m\}$ with $z_m \notin K^\delta$ and $z_m \in \bigcap_{i=1}^m K_i$ for every m .² This sequence has a subsequence contained in one of the balls making up K_1 , this subsequence in one of the balls in K_1 has a further subsequence contained in one of the balls making up K_2 , that subsequence contains a subsequence eventually contained in K_3 , and so on.³ Consider the “diagonal” sequence formed by taking the first element of the first subsequence, the second element of the second sequence, and so on. Eventually, this would be contained in a ball of radius ϵ_m for any m .⁴ Because $\epsilon_m \downarrow 0$ this means the sequence is Cauchy. Since $\ell^\infty(T)$ is a complete (Banach) space this sequence converges and must converge to an element in K . This contradicts the fact that $d(z_m, K) \geq \delta$ for every m .

Finally, combining our previous results, we want to show that $\liminf \mathbb{P}_*(X_n \in K^\delta) \geq 1 - 2\zeta$. for every $\delta > 0$. This is equivalent to saying that $\limsup \mathbb{P}^*(X_n \notin K^\delta) < 2\delta$. Recall that

$$\sup_t |X_n(t)| \leq M \text{ and } \sup_i \sup_{s, t \in T_i} |X_n(s) - X_n(t)| \leq \epsilon_m \implies X_n \in K_m.$$

Then, to show asymptotic tightness:

$$\begin{aligned} \limsup_{n \rightarrow \infty} \mathbb{P}^* \left(X_n \notin \bigcup_{i=1}^n K_i \right) &\leq \limsup \mathbb{P}^* \left(X_n \notin \bigcup_{i=1}^m K_i; \sup_t |X_n(t)| \leq M \right) + \underbrace{\limsup \mathbb{P}^* \left(\sup_t |X_n(t)| > M \right)}_{< \zeta} \\ &\leq \limsup \mathbb{P}^* \left(\sup_i \sup_{s, t \in T_i} |X_n(s) - X_n(t)| > \epsilon_m \text{ for some } m \right) + \zeta \\ &\leq \sum_{j=1}^m \limsup \mathbb{P}^* \left(\sup_i \sup_{s, t \in T_i} |X_n(s) - X_n(t)| > \epsilon_j \right) + \zeta \\ &\leq \sum_{j=1}^m \frac{\zeta}{2^j} + \zeta \\ &< 2\zeta \end{aligned}$$

□

²Pick $z_m \in \bigcap_{i=1}^m K_i \setminus K^\delta$

³Why? Each $\{z_m\}$ is in $\bigcap_{i=1}^m K_i$. Fix some n , then eventually the sequence is contained in $\bigcap_{i=1}^n K_n$ and so is contained in K_n since $K_n \supset \bigcap_{i=1}^n K_i$. This means the sequence $\{z_m\}$ has infinite members in K_n . K_n is the union of a finite number of sets, so one of these sets must contain infinite members

⁴Key here is the boundedness of the functions we are considering.

Proof is involved but useful as it shows the equivalence between asymptotic tightness and a finite approximation notion. The proof also builds some intuition for why tightness is important, at each step we are essentially showing that the whole behavior of the set is well describes (up to a tolerance of size ϵ) by a finite set of marginals. Weak convergence of the marginals is much easier to show.

This being said, the condition in Theorem 2.6 is hard to check. In particular, there is no guidance given on how to select the partition $\{T_i\}_{i=1}^m$. The next way to characterize tightness builds on asymptotic equicontinuity. The idea is the correct way to pick the partition is linked to some form of continuity: pick small T_i so that X_n does not move much on T_i .

Definition 2.7 (Asymptotic ρ -equicontinuity in probability). Suppose ρ is a semimetric on T . Then a sequence of maps $X_n : \Omega_n \rightarrow \ell^\infty(T)$ is asymptotically ρ -equicontinuous if for every $\epsilon, \eta > 0$ there exists a $\delta > 0$ such that

$$\limsup_{n \rightarrow \infty} \mathbb{P}^* \left(\sup_{d(s,t) < \delta} |X_n(s) - X_n(t)| > \epsilon \right) < \eta.$$

Remark. This is basically setting $T_i = \{(s, t) : \rho(s, t) < \delta\}$

Example. Let $X_n(t) = \frac{1}{\sqrt{n}} \sum_{i=1}^n [\mathbb{1}\{X_i \leq t\} - \mathbb{P}(X \leq t)]$. Then $|X_n(t) - X_n(t')| \approx 0$ for all $|t - t'| < \delta$. Note that here, for every n , $X_n(t)$ is still a discontinuous function of t , it's just that the jumps get closer together or smaller.

Example. Suppose that $\gamma = g(X, \beta_0) + \epsilon$ with $\mathbb{E}[\epsilon|X] = 0$. By the vector LLN, we can say that $\hat{\beta} - \beta_0 \rightarrow_p 0$.

In contrast, *asymptotic equicontinuity* will allow to say that:

$$\hat{\beta} \rightarrow_p \beta_0 \implies \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \left\{ \left(g(x_i, \hat{\beta}) - \mathbb{E}[g(x, \hat{\beta})] \right) - \left(g(x_i, \beta_0) - \mathbb{E}[g(x, \beta_0)] \right) \right\} \right| = o_p(1).$$

which is a more powerful result.

Theorem 2.7 (Theorem 1.5.7 VdV&W). *A sequence of random maps, $X_n : \Omega_n \rightarrow \ell^\infty(T)$ is asymptotically tight if and only if $X_n(t)$ is asymptotically tight in \mathbb{R} for each t and there exists a semimetric ρ on T such that (T, ρ) is totally bounded and X_n is asymptotically uniformly ρ -equicontinuous.*

Proof. First prove sufficiency then necessity:

(\Leftarrow) Fix $\epsilon, \eta > 0$. Then, there is a $\delta > 0$ such that

$$\limsup \mathbb{P}^* \left(\sup_{\rho(s,t) < \delta} |X_n(s) - X_n(t)| > \epsilon \right) < \eta.$$

Since T is totally bounded, then there are finitely many balls of radius δ that cover T , $B_1, \dots, B_{K(\delta)}$. Make these balls disjoint by taking successive “set-minuses” and then we have a partition of T . Then

$$\limsup \mathbb{P}^* \left(\max_i \sup_{s,t \in T_i} |X_n(s) - X_n(t)| > \epsilon \right) \leq \limsup \mathbb{P}^* \left(\sup_{\rho(s,t) < \delta} |X_n(s) - X_n(t)| > \epsilon \right) < \eta$$

and we can apply the results of Theorem 2.6.

(\Rightarrow) If X_n is asymptotically tight, then $g(X_n)$ is asymptotically tight for each continuous function g . Let $K_1 \subset K_2 \subset \dots$ be compact sets with:

$$\liminf \mathbb{P}_* (X_n \in K_m^\epsilon) \geq 1 - 1/m.^5$$

⁵We can choose nested compact sets with this property because the union of a finite number of compact sets is compact and the probability functional is increasing with respect to the subset ordering.

For each m define a semimetric ρ_m on T by:

$$\rho_m(s, t) = \sup_{z \in K_m} |z(s) - z(t)|.$$

Then (T, ρ_m) is totally bounded. How? Cover K_m by finitely many balls of arbitrarily small radius η centered at z_1, \dots, z_k .⁶ Partition \mathbb{R}^k into cubes of edge η and for every cube pick at most one $t \in T$ such that $(z_1(1), \dots, z_k(t))$ is in the cube. Since z_1, \dots, z_k are uniformly bounded,⁷ this gives finitely many points t_1, \dots, t_p . Now, the balls $\{t : \rho_m(t, t_i) < 3\eta\}$ cover T : t is in the ball around t_i for which $(z_1(t), \dots, z_k(t))$ and $(z_1(t_i), \dots, z_k(t_i))$ fall in the same cube. This in turn follows from the fact that $\rho_m(t, t_i)$ can be bounded by $2 \sup_{z \in K_m} \inf_i \|z - z_i\|_T + \sup_j |z_j(t_i) - z_j(t)|$.⁸

Now set

$$\rho(s, t) = \sum_{m=1}^{\infty} 2^{-m} (\rho_m(s, t) \wedge 1).$$

Fix some $\eta > 0$. Take a natural number m with $2^{-m} < \eta$. Cover T with finitely many ρ_m -balls of radius m .⁹ Let t_1, \dots, t_p be their centers. Since $\rho_1 \leq \rho_2 \leq \dots$,¹⁰ there is for every t a t_i with

$$\rho(t, t_i) \leq \sum_{k=1}^m 2^{-k} \rho_k(t, t_i) + 2^{-m} < 2\eta.$$
¹¹

So (T, ρ) is totally bounded as well. It is clear from definitions that $|z(s) - z(t)| \leq \rho_m(s, t)$ for every $z \in K_m$ and that $(\rho_m(s, t) \wedge 1) \leq 2^m \rho(s, t)$.¹² Further, if $\|z_0 - z\|_T < \epsilon$ for $z \in K_m$, then $|z_0(s) - z_0(t)| < 2\epsilon + |z(s) - z(t)|$ for any pair s, t .¹³ This gives us that

$$K_m^\epsilon \subset \left\{ z : \sup_{\rho(s, t) < 2^{-m}\epsilon} |z(s) - z(t)| \leq 3\epsilon \right\}.$$

⁶This is possible by compactness. Cover K_m by balls of radius η and then take a finite subcover.

⁷Recall that each z_i is in $\ell^\infty(T)$ which is the space of all bounded functions from $T \rightarrow \mathbb{R}$. A finite collection of bounded functions is uniformly bounded

⁸Recall that $\|f\|_T = \sup_{t \in T} |f(t)|$, $\rho_m(t, t_i) = \sup_{z \in K_m} |z(t) - z(t_i)|$, z_1, \dots, z_K are the points (bounded functions of T) around which balls of radius η cover K_m , and t_1, \dots, t_p are points of T such that the vector valued function $(z_1(\cdot), \dots, z_k(\cdot))$ takes values only in cubes of edge length η of which one of t_1, \dots, t_p is an element. Then, applying the triangle inequality and the above statements:

$$\begin{aligned} \rho_m(t, t_i) &= \sup_{z \in K_m} |z(t) - z(t_i)| \\ &\leq \sup_{z \in K_m} |z(t) - z_j(t_i)| + |z_j(t_i) - z(t)| \\ &\leq \sup_{z \in K_m} |z(t) - z_j(t_i)| + |z_j(t_i) - z_j(t)| + |z_j(t) - z(t)| \\ &\leq 2 \sup_{z \in K_m} \|z - z_j\|_T + |z_j(t_i) - z_j(t)| \end{aligned}$$

Since this holds for all j , we obtain

$$\rho_m(t, t_i) \leq 2 \sup_{z \in K_m} \inf_j \|z - z_j\|_T + \sup_j |z_j(t) - z_j(t_i)|.$$

For any t such that $(z_1(t), \dots, z_k(t))$ falls in the same cube as $(z_1(t_i), \dots, z_k(t_i))$, the first quantity is (strictly) bounded by 2η by the definition of z_1, \dots, z_k whereas the second quantity is bounded by η because t falls in the same cube as t_i . Now, since, for each $t \in T$, $(z_1(t), \dots, z_k(t)) \in T$ must fall in the same cube as $(z_1(t_i), \dots, z_k(t_i))$ for some $i \in \{1, \dots, p\}$ we have that $t \in \{\tilde{t} : \rho_m(t_i, \tilde{t}) < 3\eta\}$ for some $i \in \{1, \dots, p\}$. Since η is arbitrary, this shows that (T, ρ_m) is totally bounded.

⁹This is possible because (T, ρ_m) is totally bounded by the above argument

¹⁰Because $K_1 \subseteq K_2 \subseteq K_3 \dots$

¹¹If t is distance at most η from t_i under ρ_m , it is also distance at most η from t_i under ρ_k for $k \leq m$

¹²In the definition of ρ multiply left side and right side by 2^m . A semimetric is always (weakly) positive.

¹³Same triangle inequality decomposition as above:

$$\begin{aligned} |z_0(s) - z_0(t)| &\leq |z_0(s) - z(s)| + |z(s) - z_0(t)| \\ &\leq |z_0(s) - z(s)| + |z(s) - z(t)| + |z(t) - z_0(t)| \\ &\leq 2 \|z_0 - z\|_T + |z(s) - z(t)| \end{aligned}$$

The system of implications to get this is: if $z \in K_m$ and $\epsilon < 1$ then $\rho(s, t) < 2^{-m}\epsilon \implies \rho_m(s, t) \leq \epsilon \implies |z(s) - z(t)| \leq \epsilon$. That this holds for all $z \in K_m$ gives that for $z \in K_m^\epsilon$, $\rho(s, t) < 2^{-m}\epsilon \implies |z(s) - z(t)| \leq 3\epsilon$. Taking $\epsilon \leq 1$ is without loss of generality. To finish not that this gives us that, for given ϵ and m and for $\delta < 2^{-m}\epsilon$

$$\liminf \mathbb{P}_\star \left(\sup_{\rho(s, t) < \delta} |X_n(s) - X_n(t)| < 3\epsilon \right) \geq 1 - \frac{1}{m}.$$

This shows the backwards direction of Theorem 2.6 as well. As a note, this whole argument can be used with nets instead of sequences. \square

Remark. Important not to forget the totally bounded part of the theorem. For example, in the example of the empirical CDF case, we need to show that \mathbb{R} is totally bounded. The good news is we have choice of semi-metric.

Remark (Connection to Arzela-Ascoli). Arzela-Ascoli: Let T be a set with metric ρ that is compact. Let $C(T)$ be the set of all real valued continuous functions on T . Then $A \subset C(T)$ is compact under $|\cdot|_\infty$ if and only if it is equicontinuous and bounded.

We can think of Theorem 2.7 as a stochastic version of this. That is for

$$\liminf \mathbb{P}_\star \left(\sup_{\rho(s, t) < \delta} |X_n(s) - X_n(t)| \leq \epsilon \right) \geq 1 - \eta.$$

The set of functions satisfying this condition is equicontinuous. So then, if X_n falls here it is in a compact set by Arzela-Ascoli (Theorem 1.3). Showing this is a focus later.

3 Empirical Processes

These notes follow Section 2 in VdV&W. So far, we have discussed theory for $X_n \xrightarrow{L} X$ where both X_n and X are random elements in $\ell^\infty(T)$. The classic example that we have kept in mind is convergence of the empirical CDF process, $X_n(t) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (\mathbb{1}\{X_i \leq t\} - \mathbb{P}(X \leq t))$. In this next section we will build on the theory developed to show the convergence of some empirical processes on ℓ^∞ .

Definition 3.1 (Empirical Measure). For a random sample $\{X_i\}_{i=1}^n$, the empirical measure \mathbb{P}_n is the measure constructed from the sample (putting mass $1/n$ at each X_i). That is, for any set C :

$$\mathbb{P}_n(C) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \in C\}.$$

We can also write this in terms of the degenerate measures on each X_i :

$$\mathbb{P}_n := \frac{1}{n} \sum_{i=1}^n \delta_{X_i}.$$

Definition 3.2 (Empirical Process). For a random sample $\{X_i\}_{i=1}^n$ drawn from common distribution P , the empirical process \mathbb{G}_n is the scaled and demeaned measure on X given by:

$$\mathbb{G}_n(C) := \frac{1}{\sqrt{n}} \sum_{i=1}^n (\mathbb{1}\{X_i \in C\} - P(X_i \in C)).$$

This is often related to the empirical measure in Definition 3.1 by

$$\mathbb{G}_n = \sqrt{n} (\mathbb{P}_n - P).$$

Or written in terms of the degenerate measures on each X_i :

$$\mathbb{G}_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n (\delta_{X_i} - P).$$

Remark (Notation). We will make the following notations to save space later on. For a measure \mathbb{Q} on a space let $\mathbb{Q}f = \mathbb{E}_{\mathbb{Q}}[f(X)]$. E.j: $\mathbb{P}_n f = \mathbb{E}_n[f(X)] = \frac{1}{n} \sum_{i=1}^n f(X_i)$ and $\mathbb{G}_n f = \frac{1}{\sqrt{n}} \sum_{i=1}^n (f(X_i) - \mathbb{P}_n f)$.

With this notation:

$$\begin{aligned} \mathbb{P}_n f \xrightarrow{\text{a.s.}} P f \text{ is just saying } \frac{1}{n} \sum_{i=1}^n f(X_i) &\xrightarrow{\text{a.s.}} \mathbb{E}[f(X)] \\ \mathbb{G}_n f \xrightarrow{L} N(0, \sigma^2) \text{ is just saying } \frac{1}{\sqrt{n}} \sum_{i=1}^n (f(X_i) - \mathbb{E}[f(X)]) &\xrightarrow{L} N(0, \sigma^2) \end{aligned}$$

By LLN and CLT we have that for any function f , $\mathbb{P}_n f \rightarrow_{a.s.} P f$ and $\mathbb{G}_n f \xrightarrow{L} N(0, P(f - P f)^2)$

Example 3.1 (Classes of Functions). LLN and CLT establish the behavior of the empirical measure $\mathbb{P}_n f$ and the empirical process $\mathbb{G}_n f$ for a fixed function f (which could even be vector valued). However, we often want to study the behavior of the empirical measure of empirical process over a class of functions \mathcal{F} . In this case we can think of $\mathbb{G}_n(\mathcal{F})$ or $\mathbb{P}_n(\mathcal{F})$ as random maps onto $\ell^\infty(\mathcal{F})$. The marginal, $\mathbb{G}_n f$ or $\mathbb{P}_n f$, is then the behavior of the empirical measure/process for a single function $f \in \mathcal{F}$.

Mapping this back to the empirical CDF example of before let $\mathcal{F} = \{f_t : \mathbb{R} \rightarrow \mathbb{R} \mid f_t(x) = \mathbb{1}\{x \leq t\}, t \in T\}$. Before, we considered convergence of the whole CDF through the map $X_n : \Omega_n \rightarrow \ell^\infty(T)$ with the marginals $X_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\}$. With these new definitions/notations, we equivalently consider convergence of the entire CDF through the map $\mathbb{P}_n(\mathcal{F}) : \Omega_n \rightarrow \ell^\infty(\mathcal{F})$ with marginals $\mathbb{P}_n f_t = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\}$.

This sort of notation/generalizability is useful as we can consider the behavior of the empirical measure or empirical process over a larger class of functions. For example, if we wanted to study an entire semiparametric model we may consider the behavior of $\mathbb{G}_n(\mathcal{F})$ where

$$\mathcal{F} = \{f(x; \theta) \text{ for some } \theta \in \Theta\}.$$

Or, if we wanted to consider convergence after imposing some shape restriction, we may take

$$\mathcal{F} = \{f : X \rightarrow \mathbb{R} \mid f \text{ is monotonic}\}.$$

Remark (Notation). Sometimes we use \rightsquigarrow to denote weak convergence/convergence in law instead of \xrightarrow{L} .

Remark (Definition of ℓ^∞ Space). It is useful to review the $\ell^\infty(T)$ space for an arbitrary index space T . Define:

$$\ell^\infty(T) = \left\{ f : T \rightarrow \mathbb{R} : \sup_{t \in T} |f(t)| < \infty \right\} \quad (3.1)$$

and equip this space with the sup-norm, $\|f\|_T = \sup_{t \in T} |f(t)|$. Note that, for any \mathcal{F} , $\mathbb{G}_n(\mathcal{F})$ can be viewed as a random map into $\ell^\infty(\mathcal{F})$ for each n . Boundedness comes from the finiteness of the sample. We will sometimes make the notation $\|\mathbb{Q}\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} |\mathbb{Q}f|$ for a given measure \mathbb{Q} .

Now make some important definitions and then talk about how they relate to what we want to show.

Definition 3.3 (Glivenko-Cantelli Class). A class of functions, \mathcal{F} , for which

$$\|\mathbb{P}_n - P\|_{\mathcal{F}} \rightarrow_p 0 \quad (3.2)$$

is called a Glivenko-Cantelli class, or a P -Glivenko-Cantelli class to emphasize the dependence on the underlying measure P from which the sample is drawn.

Definition 3.4 (Donsker Class). A class of functions, \mathcal{F} , for which

$$\mathbb{G}_n(\mathcal{F}) \xrightarrow{L} \mathbb{G}(\mathcal{F}) \quad (3.3)$$

where \mathbb{G} is a tight, Borel measurable element in $\ell^\infty(\mathcal{F})$, is called a Donsker class, or P -Donsker class to emphasize the dependence on the underlying measure P from which the sample is drawn.

A Donsker class is trivially Glivenko-Cantelli.

Example 3.2 (Some Donsker Classes). Some examples of function classes:

1. If \mathcal{F} consists of a single function with finite variance then \mathcal{F} is Donsker by the Central Limit Theorem. That is $\mathbb{G}_n \xrightarrow{L} \mathbb{G}$ where \mathbb{G} is a tight element on $\ell^\infty(\mathcal{F}) = \ell^\infty(\{f\})$
2. The class of functions $\mathcal{F} = \{f(x) = x'\beta : \beta \in \mathcal{B}\}$ is Donsker if \mathcal{B} is bounded.
3. The class of monotonic densities on $[0, 1]$ is Donsker.
4. The class of square integrable functions is not Donsker (too large).

How do we know if $\mathbb{G}_N \rightsquigarrow \mathbb{G}$ where \mathbb{G} is a tight, Borel measurable element on $\ell^\infty(\mathcal{F})$? By Theorem 2.5 we know that X_n weakly converges if and only if X_n is asymptotically tight and the marginals $(X_n(t_1), \dots, X_n(t_k))$ converge weakly to a limit for every finite subset. Moreover, by Lemma 2.2 asymptotic measurability of the process is equivalent to asymptotic measurability of the marginals. By the Central Limit Theorem, we typically have weak convergence and asymptotic measurability of the marginals, what remains is to show asymptotic tightness.

Theorem 2.7 characterizes asymptotic tightness in terms of ρ -equicontinuity. Much of the work in showing tightness will be to find some semimetric ρ on \mathcal{F} such that for any $\epsilon, \eta > 0$ there is a $\delta > 0$ such that

$$\limsup_{n \rightarrow \infty} \mathbb{P}^* \left(\sup_{\rho(f,g) < \delta} |\mathbb{G}_n(f) - \mathbb{G}_n(g)| > \epsilon \right) < \eta. \quad (3.4)$$

A typical approach will be to let $\mathcal{F}_\delta = \{f, g \in \mathcal{F}, \rho(f, g) < \delta\}$. If we can show that, for some $M(\delta)$ that goes to 0 as $\delta \downarrow 0$

$$\begin{aligned} \mathbb{E} \left[\|\mathbb{G}_n\|_{\mathcal{F}_\delta} \right] &= \mathbb{E} \left[\sup_{\rho(f,g) < \delta} |\mathbb{G}_n(f) - \mathbb{G}_n(g)| \right] \\ &= \mathbb{E} \left[\sup_{\rho(f,g) < \delta} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \{f(X_i) - \mathbb{E}[f(X_i)] - g(X_i) + \mathbb{E}[g(X_i)]\} \right| \right] \\ &\leq M(\delta) \end{aligned}$$

Then, we would get the result in (3.4) by Markov's inequality. This type of result, that $\mathbb{E} \left[\|\mathbb{G}_n\|_{\mathcal{F}_\delta} \right] \leq M(\delta)$ is called a maximal inequality and is immensely useful.

Obtaining such a maximal inequality/establishing asymptotic tightness is dependent on the space not being “too large” (loosely speaking). In the example above, the class $\mathcal{F} = \{f(x) = x'\beta \mid \beta \in \mathcal{B}\}$ is Donsker so long as \mathcal{B} is bounded. To illustrate, see in the single dimensional case that

$$\sup_{b \in \mathcal{B}} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i b - \mathbb{E}[xb] \right| = \sup_{b \in \mathcal{B}} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i - \mathbb{E}[x] \right| |b|.$$

If we don't impose $|b| \leq M$ then this will blow up to $+\infty$ with probability 1, whereas if we do we have that this is $O_p(1)$. For more involved function classes, we want a way of measuring whether \mathcal{F} is large or not. This motivates the definitions of bracketing and covering numbers below.

Definition 3.5 (Covering Number). The covering number, $\mathcal{N}(\epsilon, \mathcal{F}, \|\cdot\|)$ of a class of functions \mathcal{F} is the smallest number of balls of radius ϵ under $\|\cdot\|$ needed to cover the set \mathcal{F} .

Definition 3.6 (Bracketing Number). Given two functions, ℓ and u , the bracket $[\ell, u]$ is the set of all functions f with $\ell(x) \leq f \leq u(x)$ for all x . An ϵ -bracket is a bracket $[\ell, u]$ with $\|u - \ell\| < \epsilon$. The bracketing number $\mathcal{N}_{[]}(\epsilon, \mathcal{F}, \|\cdot\|)$ is the minimum number of ϵ -brackets needed to cover \mathcal{F} .

Example 3.3 (Covering Number). Let $A = [0, 1]$ and $\|\cdot\|$ be the standard Euclidean norm¹.

1. If $\epsilon \geq 1/2$, then a ball centered at $1/2$ covers the entire interval so $\mathcal{N}(\epsilon, A, \|\cdot\|) = 1$.
2. If $\epsilon < 1/2$, then we need $\lceil \frac{1}{2\epsilon} \rceil$ balls to cover A .

Note that (i) in this example the covering number coincides with the bracketing number (ii) in general the balls needed to cover \mathcal{F} need not be centered at points in \mathcal{F} (iii) (in general) as $\epsilon \downarrow 0$ we have that $\mathcal{N}(\epsilon, \mathcal{F}, \|\cdot\|) \uparrow \infty$.

Example 3.4 (Bracketing Number). Suppose x takes values in $[0, 1]$ and let $\mathcal{F} = \{f(x) = x\beta, \text{ for } \beta \in [0, 1]\}$. Then, if $\beta_i < \beta_{i+1}$, $[x\beta_i, x\beta_{i+1}]$ forms a bracket containing all functions $f(x) = x\beta$ with $\beta_i \leq \beta \leq \beta_{i+1}$. Further note that

$$\|x\beta_i - x\beta_{i+1}\| = \sup_{x \in [0,1]} |x||\beta_i - \beta_{i+1}| = |\beta_i - \beta_{i+1}|.$$

For any $\epsilon > 0$ break up $[0, 1]$ into $[0, \epsilon, 2\epsilon, \dots]$ and take $\beta_i = (i-1)\epsilon$ to get brackets $[x\beta_i, x\beta_{i+1}]$ of size ϵ . We need $\lceil 1/\epsilon \rceil$ of these brackets to cover \mathcal{F} so that $\mathcal{N}_{[]}(\epsilon, \mathcal{F}, \|\cdot\|_\infty) \leq \lceil 1/\epsilon \rceil < 2/\epsilon$.

Remark (Bracketing vs. Covering Numbers). In general we have that $\mathcal{N}(\epsilon, \mathcal{F}, \|\cdot\|) \leq \mathcal{N}_{[]} (2\epsilon, \mathcal{F}, \|\cdot\|)$, but no opposite relationship. This shows that bracketing numbers are in general stronger than covering numbers and give you better control over the class of functions.

We will see conditions for Glivenko-Cantelli and Donsker properties under both, but in general proving Glivenko-Cantelli involves using bracketing numbers whereas proving Donsker involves using covering numbers.

In general, finding the covering/bracketing number will be difficult but we will learn some tips. Verifying that a set is Donsker will often come down to showing that the covering/bracketing number does not go to infinity “too fast.”

3.1 Maximal Inequality

For an arbitrary set of functions, \mathcal{F} , want to develop an inequality that looks something like:

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

Or, rewriting in the notation of above:

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{G}_n f| \right] \leq \text{size}(\mathcal{F}).$$

This sort of inequality is useful as it can be used to show the uniform law of large numbers:

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} |(\mathbb{P}_n - P) f| \right] \leq \frac{1}{\sqrt{n}} \text{size}(\mathcal{F}) + \text{Markov's Inequality}.$$

Or show asymptotic tightness through stochastic equicontinuity:

$$\mathbb{E} \left[\sup_{\rho(f,g) < \delta} |\mathbb{G}_n(f - g)| \right] \leq \text{size}(\mathcal{F}_\delta) + \text{Theorem 2.7}.$$

However, often we may need to change the exact application of these maximal inequalities. We will work out where these come from as we go along. The inequality will be presented for general stochastic processes (for our purposes, a stochastic process is a random map into $\ell^\infty(T)$). To build the maximal inequality, we will need to define a new norm which generalizes the L_p norms. We do so quickly below.

¹If we want to view this as a function class we can equivalently say A is the set of constant functions taking values in the interval $[0, 1]$ and consider any L_p norm on this class

3.1.1 Orlicz Norm

Definition 3.7 (Orlicz Norm). Let ψ be a non-decreasing, convex function with $\psi(0) = 0$ and X a random variable. Then, the Orlicz norm $\|X\|_\psi$ is defined as

$$\|X\|_\psi = \inf \left\{ C > 0 : \mathbb{E} \psi \left(\frac{|X|}{C} \right) \leq 1 \right\} \quad (3.5)$$

Where here the infimum over the empty set is taken to be $+\infty$.

Remark (Orlicz norms generalize L_p). Note that for any $p \geq 1$ the function $f(x) = x^p$ is convex and non-decreasing. With this in mind we can view the Orlicz norms as a generalization of the L_p norms to general convex and non-decreasing functions.

Remark (Orlicz p-norms). Of particular interest will be the Orlicz norms generated by the functions

$$\psi_p = e^{x^p} - 1.$$

for $p \geq 1$. The Orlicz norm in this case is often denoted $\|\cdot\|_{\psi_p}$. These norms give more weight to the tails of X than the standard L_p norms. It is not the case that these norms are uniformly larger than all L_p norms, however, we do have the inequalities

$$\begin{aligned} \|X\|_{\psi_p} &\leq \|X\|_{\psi_q} (\log 2)^{p/q} \\ \|X\|_p &\leq p! \|X\|_{\psi_1} \end{aligned}$$

Remark (Orlicz Norms and Markov's Inequality). Any Orlicz norm can be used to bound tail probabilities. Using Markov's inequality:

$$\mathbb{P}(|X| > x) \leq \mathbb{P} \left(\psi \left(|X| / \|X\|_\psi \right) \geq \psi \left(x / \|X\|_\psi \right) \right) \leq \frac{1}{\psi \left(x / \|X\|_\psi \right)}.$$

For $\psi_p(x) = e^{x^p} - 1$ this leads to tail estimates like $\exp(-Cx^p)$ for any random variable with a finite ψ_p -norm. Conversely, an exponential tail bound of this type shows that $\|X\|_{\psi_p}$ is finite.

Lemma 3.1 (Lemma 2.2.1 VdV&W). *Let X be a random variable with $\mathbb{P}(|X| > x) \leq Ke^{-Dx^p}$ for every x and some (fixed) constants K and D and for some $p \geq 1$. Then, the Orlicz norm of X satisfies*

$$\|X\|_{\psi_p} \leq ((1 + K)/D)^{1/p}$$

In particular, this will mean that for $C = ((1 + K)/D)^{1/p}$

$$\mathbb{E} \left[\psi \left(\frac{|X|}{C} \right) \right] \leq 1.$$

Proof. By Fundamental Theorem of Calculus and Tonelli's Theorem, for any constant B :

$$\mathbb{E} \left[e^{B|X|^p} - 1 \right] = \mathbb{E} \int_0^{|X|^p} B e^{Bs} ds = \int_0^\infty \mathbb{P}(|X| > s^{1/p}) B e^{Bs} ds$$

Now use the inequality on the tails of $|X|$, plug in $B = C^{-p} = D/(1 + K)$, and see that the final equality is bounded by 1. \square

Using the fact that $\max |X_i|^p \leq \sum |X_i|^p$ we obtain for the L_p norms, the result that

$$\left\| \max_{1 \leq i \leq m} X_j \right\|_p = \left(\mathbb{E} \max_{1 \leq i \leq m} |X_i|^p \right)^{1/p} \leq m^{1/p} \max_{1 \leq i \leq m} \|X_i\|_p.$$

We can generalize this for the Orlicz norm.

Lemma 3.2 (Lemma 2.2.2 VdV&W). *Let ψ be a convex, non-decreasing, nonzero function with $\psi(0) = 0$ and $\limsup_{x,y \rightarrow \infty} \psi(x)\psi(y)/\psi(cxy) < \infty$ for some constant c . Then, for any random variables X_1, \dots, X_m ,*

$$\left\| \max_{1 \leq i \leq m} X_i \right\|_{\psi} \leq K\psi^{-1}(m) \max_{1 \leq i \leq m} \|X_i\|_{\psi} \quad (3.6)$$

For a constant K depending only on ψ .

Proof. Without loss of generality, assume that $\psi(x)\psi(y) \leq \psi(cxy)$ for all $x, y \geq 1$ and that $\psi(1) \leq 1/2$.¹ In this case, $\psi(x/y) \leq \psi(cx)/\psi(y)$ for all $x \geq y \geq 1$.² Thus, for $y \geq 1$ and any D ;

$$\begin{aligned} \max_{1 \leq i \leq m} \psi \left(\frac{|X_i|}{Dy} \right) &\leq \max_{1 \leq i \leq m} \left[\frac{\psi(c|X_i|/D)}{\psi(y)} + \psi \left(\frac{|X_i|}{Dy} \right) \mathbb{1} \left\{ \frac{|X_i|}{Dy} < 1 \right\} \right] \\ &\leq \sum_{i=1}^m \frac{\psi(c|X_i|D)}{\psi(y)} + \psi(1) \end{aligned}$$

Let $D = c \max_{1 \leq i \leq m} \|X_i\|_{\psi}$, and take expectations to get:

$$\mathbb{E} \psi \left(\frac{\max |X_i|}{Dy} \right) \leq \frac{m}{\psi(y)} + \psi(1).$$

When $\psi(1) \leq 1/2$ take $y = \psi^{-1}(2m)$. Then:

$$\left\| \max_{1 \leq i \leq m} |X_i| \right\|_{\psi} \leq \psi^{-1}(2m) c \max_{1 \leq i \leq m} \|X_i\|_{\psi}.$$

By the convexity of ψ and the fact that $\psi(0) = 0$, it follows that $\psi^{-1}(2m) \leq 2\psi^{-1}(m)$. This gives the result. \square

To review, we have established the following inequalities above:

1. For maximums of a finite number of random variables

$$\mathbb{E} \left[\max_{1 \leq i \leq m} |X_i| \right] \leq m \max_{1 \leq i \leq m} \mathbb{E} [|X_i|].$$

2. Then, generalized this to the L_p norms

$$\left\| \max_{1 \leq i \leq m} |X_i| \right\|_{L_p} \leq m^{1/p} \max_{1 \leq i \leq m} \|X_i\|_{L_p}.$$

3. Then, generalized this using the Orlicz norm (Definition 3.7)

$$\left\| \max_{1 \leq i \leq m} |X_i| \right\|_{\psi} \leq K\psi^{-1}(m) \max_{1 \leq i \leq m} \|X_i\|_{\psi}.$$

In particular, taking $\psi(a) = e^{a^2} - 1$, we have that $\mathbb{E} [\max_{1 \leq i \leq m} |X_i|] \leq C\sqrt{\log(m+1)}$ for any C such that $\max_{1 \leq i \leq m} \mathbb{E} \left[\psi \left(\frac{|X_i|}{C} \right) \right] \leq 1$. Lemma 3.1 gives a condition for the existence of such a C .

¹If this is not the case there are constants $\sigma \leq 1$ and $\tau > 0$ such that $\phi(x) = \sigma\psi(\tau x)$ satisfies these conditions. Apply the inequality to ϕ and note that

$$\|X\|_{\psi} \leq \|X\|_{\phi}/(\sigma\tau) \leq \|X\|_{\psi}/\sigma.$$

² $x/y \geq 1$ so $\psi(x/y)\psi(y) \leq \psi(c(x/y)y)$

3.2 Chaining and Inequalities for Infinite Classes

So far, we have developed inequalities that deal with finite number of random variables. These inequalities are useful for showing Donsker/Glivenko-Cantelli property for finite classes of functions, $|\mathcal{F}| < \infty$, just set $X_i = \mathbb{G}_n f_i$. However, we often want to show uniform convergence for (uncountably) infinite classes of sets, $|\mathcal{F}| = |\mathbb{Q}|$ or $|\mathcal{F}| = |\mathbb{R}|$. To do this, we will use a technique called *chaining*.

Roughly speaking, this will work whenever our class of functions \mathcal{F} is “separable”, with respect to the empirical process \mathbb{G}_n (or empirical measure \mathbb{P}_n). This means there is a countable subset $\tilde{\mathcal{F}}$ of \mathcal{F} such that $\sup_{\mathcal{F}} |\mathbb{G}_n(f)| = \sup_{\tilde{\mathcal{F}}} |\mathbb{G}_n(f)|$. What does this buy us? If $\tilde{\mathcal{F}}_0 \subset \tilde{\mathcal{F}}_1 \subset \tilde{\mathcal{F}}_2 \cdots \subset \tilde{\mathcal{F}}$ is an infinite sequence of sets whose union is $\tilde{\mathcal{F}}$ and where each $\tilde{\mathcal{F}}_i$ is finite, then:

$$\lim_{k \rightarrow \infty} \sup_{\tilde{\mathcal{F}}_k} |\mathbb{G}_n(f)| \stackrel{a.s.}{=} \sup_{\tilde{\mathcal{F}}} |\mathbb{G}_n(f)| \xrightarrow{\text{monotone convergence}} \lim_{k \rightarrow \infty} \mathbb{E} \left[\sup_{\tilde{\mathcal{F}}_k} |\mathbb{G}_n(f)| \right] = \mathbb{E} \left[\sup_{\tilde{\mathcal{F}}} |\mathbb{G}_n(f)| \right].$$

and by separability, the last expectation is equal to the expectation of the supremum over the whole class \mathcal{F} . To make this work, we want to make sure that we can apply the inequalities that we developed in the past section. Specifically, we want to make sure that the conditions of Lemma 3.1 hold. To do so, make a definition.

Definition 3.8 (Subgaussian Process). Let \mathbb{G} be a stochastic process on a space \mathcal{F} equipped with a metric $d(\cdot, \cdot)$. Then \mathbb{G} is subgaussian if

$$\mathbb{P} \left(|\mathbb{G}(f) - \mathbb{G}(g)| > x \right) \leq 2e^{-1/2x^2/d^2(f,g)} \quad (3.7)$$

for all $f, g \in \mathcal{F}$ and any $x \geq 0$.

Also define a separable function as an analytic concept and then extend this to the case of stochastic processes.

Definition 3.9 (Separable Function). A function $f : A \rightarrow B$ from a topological space A into a topological space B is separable if there is a countable, dense, subset $S \subset A$ such that for any closed $F \subset B$ and any open $I \subset A$, if $f(t) \in F$ for all $t \in I \cap S$ then $f(t) \in F$ for all $t \in I$. This is often denoted as an S -separable function to emphasize the dependence on the countable, dense subset S .

Lemma 3.3 (Continuity and Separability). A continuous function $f : \mathcal{X} \rightarrow \mathcal{Y}$ from a separable space \mathcal{X} onto \mathcal{Y} is separable.

Definition 3.10 (Separable Process; Shalizi 2007). A stochastic process on a topological space \mathcal{F} , $\mathbb{G}(\cdot, \omega) : \Omega \rightarrow \ell^\infty(\mathcal{F})$, is separable if there is a countable, dense, subset of \mathcal{F} , $\tilde{\mathcal{F}}$, and a measure zero set N such that for all $\omega \notin N$, $\mathbb{G}(\cdot, \omega)$ is $\tilde{\mathcal{F}}$ -separable.¹

Separability can be roughly interpreted as ensuring that the behavior of the function (and therefore the stochastic process) can be well described by its behavior on countable subset. This ensures some of the properties that we’ve seen above, namely that $\sup_{f \in \tilde{\mathcal{F}}} |\mathbb{G}(f)| \stackrel{a.s.}{=} \sup_{f \in \mathcal{F}} |\mathbb{G}(f)|$. We are now ready for the main theorem of this subsection, the proof of which will rely on the chaining argument roughly discussed above.

Theorem 3.1 (Theorem 2.2.4 VdV&W). Let \mathbb{G} be a separable subgaussian process on a space \mathcal{F} equipped with a metric $d(\cdot, \cdot)$ and let $\text{diam}(\mathcal{F}) = \sup_{f,g \in \mathcal{F}} d(f, g)$. Then

$$\mathbb{E} \sup_{f,g \in \mathcal{F}} |\mathbb{G}(f) - \mathbb{G}(g)| \leq K \int_0^{\text{diam}(\mathcal{F})} \sqrt{\log \mathcal{N}(\epsilon, \mathcal{F}, d)} d\epsilon \quad (3.8)$$

$$\mathbb{E} \sup_{f \in \mathcal{F}} |\mathbb{G}(f)| \leq \mathbb{E} |\mathbb{G}(f_0)| + K \int_0^{\text{diam}(\mathcal{F})} \sqrt{\log \mathcal{N}(\epsilon, \mathcal{F}, d)} d\epsilon, \quad \forall f_0 \in \mathcal{F} \quad (3.9)$$

¹Note that this requires a topology on \mathcal{F} . In the applications we will be talking about \mathcal{F} will be equipped with a metric d . This will generate a topology.

Proof. Proof proceeds in steps. Let $M = \text{diam}(\mathcal{F}) = \sup_{f,g \in \mathcal{F}} d(f,g)$. For any $f_0, f \in \mathcal{F}$ we have that $d(f_0, g) \leq M$. First step will be to build a “chain” to almost any point in \mathcal{F} . Further, let $\tilde{\mathcal{F}}$ be the dense subset as described in Definitions 3.9 and 3.10.

Step 1: Building a Chain. Pick any $f_0 \in \tilde{\mathcal{F}}$ and let $\tilde{\mathcal{F}}_0 = \{f_0\}$. Build nesting sets, $\mathcal{F}_0 \subset \tilde{\mathcal{F}}_1 \subset \tilde{\mathcal{F}}_2 \subset \dots \subset \tilde{\mathcal{F}}$, such that for each $k \in \mathbb{N}$ $\tilde{\mathcal{F}}_k = \{f_1, \dots, f_{m(k)}\}$ is a maximal collection of points such that $d(f_k, g_k) > \frac{M}{2^k}$ for any $f_k, g_k \in \mathcal{F}_k$. By definition of the packing numbers we know that $\mathcal{N}\left(\frac{M}{2^{k+1}}, \mathcal{F}, d\right)$ balls cover \mathcal{F} . Putting a point at the center of each of these balls creates points that are at least distance $\frac{M}{2^k}$ from each other. Similarly, if we could fit more points at least distance $\frac{M}{2^k}$ distance away from each other than we could pack more balls of radius $\frac{M}{2^{k+1}}$ into \mathcal{F} by centering a ball at each point. So, $|\tilde{\mathcal{F}}_k| \leq \mathcal{N}\left(\frac{M}{2^{k+1}}, \mathcal{F}, d\right)$ (Inequality comes because each $\tilde{\mathcal{F}}_k$ has to contain all previous sets).

Finally, link each point $f_k \in \tilde{\mathcal{F}}_k$ to a unique point $f_{k-1} \in \tilde{\mathcal{F}}_{k-1}$ such that $d(f_k, f_{k-1}) \leq \frac{M}{2^{k-1}}$.²

Step 2: Use the chain to build a bound. Using these links, for any $f_k, g_k \in \tilde{\mathcal{F}}_k$ we can build a chain back to f_0 :

$$\begin{aligned} |\mathbb{G}(f_k) - \mathbb{G}(g_k)| &= |(\mathbb{G}(f_k) - \mathbb{G}(f_0)) - (\mathbb{G}(g_k) - \mathbb{G}(f_0))| \\ &= \left| \sum_{j=0}^k (\mathbb{G}(f_j) - \mathbb{G}(f_{j-1})) - \sum_{j=0}^k (\mathbb{G}(g_j) - \mathbb{G}(g_{j-1})) \right| \end{aligned}$$

By the triangle inequality:

$$\mathbb{E} \left[\max_{g_k, f_k \in \tilde{\mathcal{F}}_k} |\mathbb{G}(f_k) - \mathbb{G}(g_k)| \right] \leq 2 \sum_{j=0}^K \mathbb{E} \left[\max_{s_i \in \tilde{\mathcal{F}}_i} |\mathbb{G}(s_i) - \mathbb{G}(s_{i-1})| \right] \quad (\text{P-1})$$

With this setup, we can use the maximal inequalities developed above, applying them to the finite sets $\tilde{\mathcal{F}}_k$.

Step 3: Try to control the jumps. Recall that there are at most $\mathcal{N}\left(\frac{M}{2^{k+1}}, \mathcal{F}, d\right)$ points in \mathcal{F}_k and that $d(s_k, s_{k-1}) \leq \frac{M}{2^{k-1}}$. By our maximal inequality in Lemma 3.2, taking $\psi(a) = e^{a^2} - 1$ we have that

$$\mathbb{E} \left[\max_{s_j \in \tilde{\mathcal{F}}_j} |\mathbb{G}(s_j) - \mathbb{G}(s_{j-1})| \right] \leq C_j \sqrt{\log \left(\mathcal{N} \left(\frac{M}{2^{j+1}}, \mathcal{F}, d \right) + 1 \right)}.$$

For any constant C_j such that

$$\mathbb{E} \left[\exp \left(\frac{(\mathbb{G}(s_j) - \mathbb{G}(s_{j-1}))^2}{c_j^2} \right) - 1 \right] \leq 1, \quad \forall s_j \in \tilde{\mathcal{F}}_j.$$

Since \mathbb{G} is subgaussian we know that $\mathbb{P}(|\mathbb{G}(f) - \mathbb{G}(g)| > x) \leq 2e^{-\frac{1}{2}x^2/d^2(f,g)}$. By construction, we know that $d(s_j, s_{j-1}) \leq \frac{M}{2^{j-1}}, \forall s_j \in \tilde{\mathcal{F}}_j$. So

$$\mathbb{P}(|\mathbb{G}(s_j) - \mathbb{G}(s_{j-1})| > x) \leq 2e^{-\frac{1}{2} \frac{x^2}{\lceil M/2^{j-1} \rceil^2}}.$$

By Lemma 3.1 we can take $C_j = \frac{\sqrt{3}M}{2^{j-1}}$ and combine with the other results in this section to get

$$\mathbb{E} \left[\max_{s_j \in \tilde{\mathcal{F}}_j} |\mathbb{G}(s_j) - \mathbb{G}(s_{j-1})| \right] \leq \frac{\sqrt{3}M}{2^{j-1}} \sqrt{\log \left(\mathcal{N} \left(\frac{M}{2^{j+1}}, \mathcal{F}, d \right) + 1 \right)} \quad (\text{P-2})$$

²I found it helpful to remember here that $\tilde{\mathcal{F}}_{k-1} \subset \tilde{\mathcal{F}}_k$. If no such f_{k-1} exists we could add f_k to $\tilde{\mathcal{F}}_{k-1}$, a contradiction. If $f_k \in \tilde{\mathcal{F}}_{k-1}$ we can link it to itself.

Step 4: Combine Results of Previous Steps. Combine the inequalities from (P-1) and (P-2) to get

$$\mathbb{E} \left[\max_{g_k, f_k \in \tilde{\mathcal{F}}_k} |\mathbb{G}(f_g) - \mathbb{G}(g_k)| \right] \leq \sqrt{12}M \sum_{j=0}^k \frac{1}{2^{j-1}} \sqrt{\log \left(\mathcal{N} \left(\frac{M}{2^{j+1}}, \mathcal{F}, d \right) + 1 \right)}$$

With some complex rearranging of squares, we can bound the sum in the display by it's integral up to a constant scale, dropping the added 1 in the log in the process.³ That is, we ultimately obtain for some constant K :

$$\mathbb{E} \left[\max_{g_k, f_k \in \tilde{\mathcal{F}}_k} |\mathbb{G}(f_g) - \mathbb{G}(g_k)| \right] \leq K \int_0^M \sqrt{\log \mathcal{N}(\epsilon, \mathcal{F}, d)} d\epsilon \quad (\text{P-3})$$

Step 5: Conclude by Separability. $\{\tilde{\mathcal{F}}_k\}_{k=1}^\infty$ is an increasing sequence of sets that approaches $\tilde{\mathcal{F}}$ and note that the bound in (P-3) does not depend on (little) k . So, invoking monotone convergence and separability of \mathbb{G} :

$$\begin{aligned} \mathbb{E} \left[\sup_{f, g \in \mathcal{F}} |\mathbb{G}(f) - \mathbb{G}(g)| \right] &= \mathbb{E} \left[\sup_{f, g \in \tilde{\mathcal{F}}} |\mathbb{G}(f) - \mathbb{G}(g)| \right] \\ &= \lim_{k \rightarrow \infty} \mathbb{E} \left[\max_{f, g \in \tilde{\mathcal{F}}_k} |\mathbb{G}(f) - \mathbb{G}(g)| \right] \\ &\leq K \int_0^M \sqrt{\log \mathcal{N}(\epsilon, \mathcal{F}, d)} d\epsilon \end{aligned}$$

This is the inequality in equation (3.8). To get equation (3.9) fix any f_0 and apply triangle inequality. \square

Remark (Comments on Theorem 3.1). Theorem 3.1 is an involved result. Some remarks below.

1. We have shown that if \mathbb{G}_n is a separable subgaussian process then

$$\mathbb{E} \left[\sup_{f, g \in \mathcal{F}} |\mathbb{G}_n(f) - \mathbb{G}_n(g)| \right] \leq K \int_0^{\text{diam}(\mathcal{F})} \sqrt{\log \mathcal{N}(\epsilon, \mathcal{F}, d)} d\epsilon.$$

Note that the right hand side does not depend on \mathbb{G}_n at all! Only on the “size” of \mathcal{F} .

2. So, suppose we want to show that \mathbb{G}_n is an asymptotically tight process on \mathcal{F} . By Theorem 2.7 it is sufficient (and necessary) to show that for every $\epsilon, \eta > 0$ there is a $\delta > 0$ such that:

$$\limsup_{n \rightarrow \infty} \mathbb{P} \left(\sup_{\rho(f, g) \leq \delta} |\mathbb{G}_n(f) - \mathbb{G}_n(g)| > \epsilon \right) < \eta.$$

Let $\mathcal{F}_\delta = \{s = f - g : f, g \in \mathcal{F} \text{ and } \rho(f, g) \leq \delta\}$. The above can be restated as showing that $\exists \delta > 0$ such that:

$$\limsup_{n \rightarrow \infty} \mathbb{P} \left(\sup_{s \in \mathcal{F}_\delta} |\mathbb{G}_n(s)| > \epsilon \right) < \eta.$$

By Markov's inequality we can bound the probability in the above display by

$$\frac{1}{\epsilon} \mathbb{E} \left[\sup_{s \in \mathcal{F}_\delta} |\mathbb{G}_n(s)| \right] \leq \frac{K}{\epsilon} \int_0^\delta \sqrt{\log \mathcal{N}(\epsilon, \mathcal{F}, d)} d\epsilon$$

And then we can sent the RHS to 0 by sending $\delta \downarrow 0$ as long as the integral on the RHS is finite. Asymptotic tightness plus convergence of marginals will give convergence to a tight element in $\ell^\infty(\mathcal{F})$ by Theorem 2.5. What remains is to show the conditions on \mathbb{G} , separability and subgaussian.

³Here we use the fact that $\log(1+m) \leq 2\log(m)$ for $m \geq 2$

3.3 Symmetrization

Symmetrization is a technique that will allow us to get/show(?) a subgaussian process. This follows the discussion in Chapter 2.2.1 and 2.3 in VanDerVaart and Wellner.

What sort of variables are subgaussian? A classic example below.

Definition 3.11 (Rademacher Random Variable). Random variable $\epsilon_i : \Omega_i \rightarrow \mathbb{R}$ is a Rademacher random variable if $\mathbb{P}(\epsilon_i = 1) = \mathbb{P}(\epsilon_i = -1) = 1/2$.

The following lemma shows that a particular process consisting of Rademacher random variables is subgaussian.

Lemma 3.4 (Hoeffding's Inequality). *Let a_1, \dots, a_n be constants and $\epsilon_1, \dots, \epsilon_n$ be independent Rademacher random variables. Then*

$$\mathbb{P}\left(\left|\sum_{i=1}^n a_i \epsilon_i\right| > x\right) \leq 2e^{\frac{1}{2} \frac{x^2}{\|a\|^2}}.$$

where $\|a\|$ denotes the Euclidean norm of a .

Proof. (From VdV&W, Lemma 2.2.7) For any λ and any Rademacher random variable ϵ one has that $\mathbb{E}e^{\lambda\epsilon} = (e^\lambda + e^{-\lambda})/2$. By power series expansion:

$$\begin{aligned} e^\lambda &= 1 + \lambda + \frac{\lambda^2}{2!} + \frac{\lambda^3}{3!} + \dots \\ e^{-\lambda} &= 1 - \lambda + \frac{\lambda^2}{2!} - \frac{\lambda^3}{3!} + \dots \\ \implies (e^\lambda + e^{-\lambda})/2 &= 1 + \frac{\lambda^2}{2!} + \frac{\lambda^4}{4!} + \frac{\lambda^6}{6!} \dots \\ &\leq 1 + \frac{\lambda^2}{2} + \frac{\lambda^4}{2^2 \cdot 2!} + \frac{\lambda^6}{2^3 \cdot 3!} \\ &= e^{\lambda^2/2} \end{aligned}$$

where in the last inequality we use that $2^k \cdot k! \leq (2k)!$ so that in total we have that $\mathbb{E}e^{\lambda\epsilon} = (e^\lambda + e^{-\lambda})/2 \leq e^{\lambda^2/2}$. Take $\lambda = x/\|a\|^2$ and apply Markov's inequality to get the result. \square

Example. For any functions f, g we have that

$$\mathbb{P}\left(\left|\sum_{i=1}^n \frac{\epsilon_i}{\sqrt{n}} (f(x_i) - g(x_i))\right| > x \mid \{X_i\}\right) \leq 2e^{-\frac{1}{2} \frac{x^2}{d_n^2(f,g)}}.$$

where $d_n^2(f, g) := \frac{1}{n} \sum_{i=1}^n (f(x_i) - g(x_i))^2$ is the square of the prediction norm.

We would like to use the maximal inequality in Theorem 3.1 to control $\mathbb{E}[\sup_{f,g} |\mathbb{G}_n(f) - \mathbb{G}_n(g)|]$, but the problem is that \mathbb{G}_n is not (in general), subgaussian. However, from Lemma 3.4 we know that, at least conditional on our data, $\mathbb{G}_n^\circ := \frac{1}{\sqrt{n}} \sum \epsilon_i (f(x_i) - g(x_i))$ is. Strategy will be to relate the two processes, \mathbb{G}_n and \mathbb{G}_n° .

Before starting, it is useful to formally define the probability space that we are working with. Let $\epsilon_1, \dots, \epsilon_n$ be i.i.d Rademacher random variables that are generated independent of (X_1, \dots, X_n) , our observed data. Define the symmetrized process:

$$\mathbb{P}_n^\circ f = \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i).$$

Because \mathbb{P}_n° is subgaussian, conditional on X_1, \dots, X_n , it can be easier to study. We want to bound supremum of the process $\mathbb{P}_n - P$ by that of the symmetrized process. To formalize these bounds, we have to be careful about the non-measurability of supremum like $\|\mathbb{P}_n - P\|_{\mathcal{F}}$.¹

¹Even if \mathcal{F} is a class of measurable functions, the supremum may not be measurable.

In the following discussion, outer expectations of functions of X_1, \dots, X_n are assumed to be taken with respect to the coordinate projection of the infinite product space $(\mathcal{X}^{\mathbb{N}}, \mathcal{A}^{\mathbb{N}}, P^{\mathbb{N}})$ onto its first n coordinates, $(\mathcal{X}^n, \mathcal{A}^n, P^n)$.² When auxiliary variables, independent of the X 's are involved, as in the next lemma, we can use a similar convention. The underlying probability space is assumed to be of the form $(\mathcal{X}^n, \mathcal{A}^n, P^n) \times (\mathcal{Z}, \mathcal{C}, Q)$. Independence is understood in terms of a product probability space.³ To manage all this, we take advantage of a modified Fubini's theorem for outer expectations, stated here without proof.

Lemma 3.5 (Fubini's Theorem, Lemma 1.2.6 VdV&W). *Let T be defined on a product probability space. Then*

$$\mathbb{E}_{\star} T \leq \mathbb{E}_{1\star} \mathbb{E}_{2\star} T \leq \mathbb{E}_1^{\star} \mathbb{E}_2^{\star} T \leq \mathbb{E}^{\star} T.$$

Proof. For the last inequality, we can assume that $\mathbb{E}^{\star} T < \infty$ so that $\mathbb{E}^{\star} T = \mathbb{E} T^{\star}$. Since T^{\star} is jointly measurable with respect to the product σ -field, the map $\omega_2 \mapsto T^{\star}(\omega_1, \omega_2)$ is a measurable majorant of $\omega_2 \mapsto T(\omega_1, \omega_2)$ for P_1 almost all ω_1 . Hence $\int T^{\star}(\omega_1, \omega_2) dP_2(\omega_2) \geq (\mathbb{E}_2^{\star} T)(\omega_1)$ for P_1 almost all ω_1 . Further, by Fubini's theorem for standard integrals, this is a measurable function of ω_1 . Thus the integral of this with respect to P_1 is an upper bound for $\mathbb{E}_1^{\star} \mathbb{E}_2^{\star} T$. Since T^{\star} is jointly measurable, by another application Fubini's theorem for standard integrals:

$$\mathbb{E}^{\star} T = \mathbb{E} T^{\star} = \int \left(\int T^{\star}(\omega_1, \omega_2) dP_2(\omega_2) \right) dP_1(\omega_1) \geq \mathbb{E}_1^{\star} \mathbb{E}_2^{\star} T.$$

The inequalities for inner expectations hold by considering $-T$. □

Lemma 3.6 (Symmetrization). *For every non-decreasing, convex, $\Phi : \mathbb{R} \rightarrow \mathbb{R}$ and class of measurable functions \mathcal{F} :*

$$\mathbb{E}^{\star} \Phi (\|\mathbb{P}_n - P\|_{\mathcal{F}}) \leq \mathbb{E}^{\star} \Phi (2 \|\mathbb{P}_n^{\circ}\|_{\mathcal{F}}).$$

Where outer expectations are calculated as described above.

Proof. Let Y_1, \dots, Y_n be independent copies of X_1, \dots, X_n (independently drawn from the same joint distribution as X_1, \dots, X_n , defined formally as the coordinate projections on the last n coordinates in the product space $(\mathcal{X}^n, \mathcal{A}^n, P^n) \times (\mathcal{Z}, \mathcal{C}, Q) \times (\mathcal{X}^n, \mathcal{A}^n, P^n)$).

For fixed values X_1, \dots, X_n applying Jensen's inequality to the absolute value gives:

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

where \mathbb{E}_Y^{\star} is the outer expectation with respect to Y_1, \dots, Y_n computed for P^n . Again applying Jensen's inequality gives:

$$\Phi (\|\mathbb{P}_n - P\|_{\mathcal{F}}) \leq \mathbb{E}_Y^{\star} \Phi \left(\left\| \frac{1}{n} \sum_{i=1}^n [f(X_i) - f(Y_i)] \right\|_{\mathcal{F}}^{\star Y} \right).$$

where $f^{\star Y}$ is the minimal measurable majorant of f with respect to the distribution of Y . Because Φ is non-decreasing and continuous, the $\star Y$ inside Φ can be moved to \mathbb{E}_Y^{\star} . In total then:

$$\Phi (\|\mathbb{P}_n - P\|_{\mathcal{F}}) \leq \mathbb{E}_Y^{\star} \Phi \left(\left\| \frac{1}{n} \sum_{i=1}^n [f(X_i) - f(Y_i)] \right\|_{\mathcal{F}} \right).$$

²That is the outer expectation is taken relative to P^n where P^n is defined from the projection of the infinite product space onto its first n coordinates

³Two sub-sigma algebras, $\mathcal{A}_1, \mathcal{A}_2 \subset \mathcal{A}$ are considered independent if $\mathbb{P}(A_1 A_2) = \mathbb{P}(A_1) \mathbb{P}(A_2)$ for any $A_1 \in \mathcal{A}_1, A_2 \in \mathcal{A}_2$. The sigma algebra generated by a random map $X : (\Omega, \mathcal{A}, \mathbb{P}) \rightarrow (\mathcal{X}, \mathcal{B})$ is the smallest sigma algebra on Ω that makes X measurable,

$$\sigma(X) := \{X^{-1}(B) : B \in \mathcal{B}\}.$$

Two random variables, X, Y , defined on the same probability space are independent if their generated sigma algebras, $\sigma(X), \sigma(Y)$, are independent. In the context of having independent draws X_1, \dots, X_n we can think of this as the projection mappings $\pi_i(\mathcal{X}^n)$ being independent.

Next, take the expectation with respect to X_1, \dots, X_n of the above quantity to get:

$$\mathbb{E}^* \Phi (\|\mathbb{P}_n - P\|_{\mathcal{F}}) \leq \mathbb{E}_X^* \mathbb{E}_Y^* \Phi \left(\frac{1}{n} \left\| \sum_{i=1}^n [f(X_i) - f(Y_i)] \right\| \right).$$

Adding a minus sign in front of the term $[f(X_i) - f(Y_i)]$ has the effect of exchanging X_i and Y_i . By construction of the underlying probability space this does not change the expectation. Hence, the expression

$$\mathbb{E}^* \Phi \left(\frac{1}{n} \left\| \sum_{i=1}^n e_i [f(X_i) - f(Y_i)] \right\| \right).$$

is the same for any n -tuple $(e_1, \dots, e_n) \in \{-1, 1\}^n$. So:

$$\mathbb{E}^* \Phi (\|\mathbb{P}_n - P\|_{\mathcal{F}}) \leq \mathbb{E}_{\epsilon} \mathbb{E}_{X,Y}^* \Phi \left(\left\| \sum_{i=1}^n \epsilon_i [f(X_i) - f(Y_i)] \right\|_{\mathcal{F}} \right).$$

where each ϵ_i is an independent Rademachar random variable and $\epsilon = (\epsilon_1, \dots, \epsilon_n)$. By triangle inequality and convexity of the Φ :

$$\begin{aligned} & \mathbb{E}_{\epsilon} \mathbb{E}_{X,Y}^* \Phi \left(\left\| \sum_{i=1}^n \epsilon_i [f(X_i) - f(Y_i)] \right\|_{\mathcal{F}} \right) \\ & \leq \mathbb{E}_{\epsilon} \mathbb{E}_{X,Y}^* \Phi \left(\left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right\|_{\mathcal{F}} + \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(Y_i) \right\|_{\mathcal{F}} \right) \\ & \leq \frac{1}{2} \mathbb{E}_{\epsilon} \mathbb{E}_{X,Y}^* \Phi \left(2 \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right\|_{\mathcal{F}} \right) + \frac{1}{2} \mathbb{E}_{\epsilon} \mathbb{E}_{X,Y}^* \Phi \left(2 \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(Y_i) \right\|_{\mathcal{F}} \right) \\ & \leq \mathbb{E}^* \Phi (2 \|\mathbb{P}_n^{\circ}\|_{\mathcal{F}}) \end{aligned}$$

where we use the fact that a repeated outer expectation can be bounded above by a joint outer expectation, $\mathbb{E}_{\epsilon} \mathbb{E}_{X,Y}^* \leq \mathbb{E}_{\epsilon,X,Y}^* (= \mathbb{E}^*)$ using Lemma 3.5. \square

Corollary 3.1 (Symmetrization of Empirical Process, Andres' Notes). *For real valued processes as described above:*

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

Proof. Take $\Phi(x) = |x|$. All norms on \mathbb{R} are equivalent to $|\cdot|$. Lemma 3.6 then gives us that $\mathbb{E}^* \|\mathbb{P}_n - P\|_{\mathcal{F}} \leq 2 \mathbb{E}^* \|\mathbb{P}_n^{\circ}\|_{\mathcal{F}}$. Expand this out and scale by \sqrt{n} to get the result. For non real-valued processes (vector valued, function valued, etc.) we can replace $|\cdot|$ with $\|\cdot\|$ above. \square

Remark. The proof of Corollary 3.1 uses the fact that Lemma 3.6 is not an asymptotic bound, it holds in every finite sample.

We now have the pieces to show a class of functions \mathcal{F} is either

- Glivenko-Cantelli, i.e that $\|\mathbb{P}_n - P\|_{\mathcal{F}} = o_p(1)$. We will do this by placing conditions on the bracketing/covering numbers.
- Donsker, i.e that $\mathbb{G}_n(\mathcal{F}) \xrightarrow{L} \mathbb{G}(\mathcal{F})$ for some tight \mathbb{G} . To do so, we will use covering numbers. The system of arguments needed to show this is usually as follows:

- By Theorem 2.5 weak convergence to a tight limit is equivalent to asymptotic tightness and weak convergence of the marginals.
- Weak convergence of the marginals is generally provided by CLT. Theorem 2.7 shows that asymptotic tightness is equivalent to uniform ρ -equicontinuity (Definition 1.18)
- Asymptotic equicontinuity holds if $\mathbb{E} \left[\sup_{f \in \mathcal{F}_\delta} |\mathbb{G}_n(f)| \right]$ goes to 0 as $\delta \downarrow 0$. Theorem 3.1 gives conditions where this is possible for separable, subgaussian processes.
- Lemma 3.3 suggests separability if \mathcal{F} is separable. Lemma 3.4 gives us that the Rademachar process is subgaussian conditional on X_1, \dots, X_n . Combining with Theorem 3.1 gives

$$\mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}_\delta} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right] \leq \int_0^{\text{diam}(\mathcal{F}_\delta)} \sqrt{\log \mathcal{N}(s, \mathcal{F}_\delta, L_2(\mathbb{P}_n))} ds$$

where $L_2(\mathbb{P}_n) = \frac{1}{n} \sum_{i=1}^n f(X_i)^2$ is the L_2 norm with respect to the empirical measure. Since X is random this norm will also end up random. This seems like it will make dealing with

$$\mathbb{E} \left[\int_0^{\text{diam}(\mathcal{F}_\delta)} \sqrt{\log \mathcal{N}(s, \mathcal{F}_\delta, L_2(\mathbb{P}_n))} ds \right]$$

painful, but we end up having good bounds for this.

- Lemma 3.6, and in particular Corollary 3.1, relates the empirical process to the Rademachar process. Take expectations with respect to X in the above bound to bound the and apply the symmetrization lemma to get bounds on the empirical process of interest.

We next move to verifying the various conditions and applying them to show that some specific processes are Glivenko-Cantelli or Donsker.

3.4 Glivenko-Cantelli

This subsection follows Section 2.4 in Van DerVaart and Wellner. Goal is to establish conditions for a uniform law of large numbers using bracketing and covering numbers.

Theorem 3.2 (Theorem 2.4.1 VdV&W). *Let \mathcal{F} be a class of measurable functions such that*

$$\mathcal{N}_{[]}(\epsilon, \mathcal{F}, L_1(P)) < \infty$$

for every $\epsilon > 0$. Then \mathcal{F} is Glivenko-Cantelli.

Proof. Fix $\epsilon > 0$. Choose finitely many ϵ -brackets $[l_i, u_i]$ whose union contains \mathcal{F} and such that $P(u_i - l_i) < \epsilon$ for every i . Then, for every $f \in \mathcal{F}$ there is a bracket, $l_i \leq f \leq u_i$, such that:

$$\begin{aligned} (\mathbb{P}_n - P)f &\leq \mathbb{P}_n u_i - Pf \leq (\mathbb{P}_n - P)u_i + P(u_i - f) \leq (\mathbb{P}_n - P)u_i + \epsilon \\ (\mathbb{P}_n - P)f &\geq \mathbb{P}_n l_i - Pf \geq (\mathbb{P}_n - P)l_i + P(l_i - f) \geq (\mathbb{P}_n - P)l_i - \epsilon \end{aligned}$$

Consequently,

$$\begin{aligned} \sup_{f \in \mathcal{F}} (\mathbb{P}_n - P)f &\leq \max_i (\mathbb{P}_n - P)u_i + \epsilon \\ \inf_{f \in \mathcal{F}} (\mathbb{P}_n - P)f &\geq \min_i (\mathbb{P}_n - P)l_i - \epsilon \end{aligned}$$

By the strong law of large numbers, both the maximums and the minimums on the right hand side of the inequalities above converge almost surely to 0. Combination these yields that $\limsup \|\mathbb{P}_n - P\|_{\mathcal{F}}^* \leq \epsilon$ almost surely for every $\epsilon > 0$. Take $\epsilon \downarrow 0$ to see that the \limsup must be 0 almost surely. \square

Remark. Some comments on Theorem 3.2:

1. Proof is really quite straightforward. Bracketing gives pointwise control so just use the upper and lower bounds.
2. No measurability condition is needed and no requirements on the rate of growth of $\mathcal{N}_{[]}(\epsilon, \cdot, \cdot)$ as $\epsilon \downarrow 0$.

Example (Empirical CDF is Glivenko-Cantelli). Let X be a scalar random variable.¹ We want to show that

$$\sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\} - P(X_i \leq t) \right| = o_p(1).$$

Let $\mathcal{F} = \{f(x) = \mathbb{1}\{X_i \leq t\} : t \in \mathbb{R}\}$. Partition \mathbb{R} into grids $-\infty = t_0 < t_1 < \dots < t_m = \infty$ such that $\mathbb{P}(t_i \leq X \leq t_{i+1}) < \epsilon$ for each i . Then the finitely many brackets $[\mathbb{1}\{X_i \leq t_i\}, \mathbb{1}\{X_i \leq t_{i+1}\}]$ cover \mathcal{F} and are “size” ϵ under P . So, $\mathcal{N}_{[]}(\epsilon, \mathcal{F}, L_1(P)) < \infty$ for every $\epsilon > 0$. So \mathcal{F} is Glivenko-Cantelli (i.e, we have a uniform law of large numbers).

The requirement on the bracketing numbers can in general be hard. Would like a result for the covering numbers as well. This will make showing that some classes are Glivenko-Cantelli easier later on. Before doing so, we need to make a couple definitions:

Definition 3.12 (Envelope). A class \mathcal{F} has envelope F if $|f(x)| \leq F(x)$ for all x and all $f \in \mathcal{F}$.

Definition 3.13 (Truncated Class). Let \mathcal{F} be a class of functions. Then the truncated class \mathcal{F}_M is given

$$\mathcal{F}_M = \{f(x)\mathbb{1}\{f \leq M\} : f \in \mathcal{F}\}.$$

Definition 3.14 (P-Measurable Class). A class \mathcal{F} is P -measurable if $\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i)\epsilon_i \right|$ is measurable with respect to the product measure on $(\mathcal{X}^m, \mathcal{A}^m, P^n) \times (\mathcal{Z}, \mathcal{C}, Q)$, where $(\mathcal{Z}, \mathcal{C}, Q)$ denotes the probability space that the Rademachar random variables are defined on.

Definition 3.15 ($L_P(\mathbb{P}_n)$ -norm). We have that $\|f - g\|_{L_1(P)} = \mathbb{E}_P [|f(x) - g(x)|]$, similarly we can define

$$\|f - g\|_{L_1(\mathbb{P}_n)} = \mathbb{E}_{\mathbb{P}_n} [|f(x) - g(x)|].$$

and through this define $\mathcal{N}_{[]}(\epsilon, \mathcal{F}, L_1(\mathbb{P}_n))$.

Theorem 3.3 (Theorem 2.4.3, VdV&W). Let \mathcal{F} be a P -measurable class of measurable functions with envelope F such that $\mathbb{P}^* F < \infty$. If $\log \mathcal{N}_{[]}(\epsilon, \mathcal{F}_M, L_1(\mathbb{P}_n)) = o_{P^*}(n)$ for every ϵ and $M > 0$, then $\|\mathbb{P}_n - P\|_{\mathcal{F}}^* \rightarrow 0$ almost surely and in mean.

Proof. Idea will be to apply the maximal inequality in Theorem 3.1.

Step 1: Symmetrization. First, we will apply symmetrization (Corollary 3.1)

$$\mathbb{E}^* \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - Pf(X_i) \right| \right] \leq 2 \cdot \mathbb{E}^* \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right]$$

And then truncate the functions, $f = f\mathbb{1}\{f \leq M\} + f\mathbb{1}\{f > M\}$, apply triangle inequality, and bound the functions not in \mathcal{F}_M with the envelope F .

$$\begin{aligned} &\leq 2\mathbb{E}_X \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}_M} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right] + 2\mathbb{E}^* [\epsilon_i F(X_i)\mathbb{1}\{F \geq M\}] \\ &= 2\mathbb{E}_X \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}_M} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right] + 2P^* F(X_i)\mathbb{1}\{F \geq M\} \end{aligned}$$

¹This generalizes easily for a vector valued random variable

Note the argument that allows us to replace the first outer expectation with iterated expectations over X and ϵ : each of the functions in \mathcal{F} are measurable and \mathcal{F}_M is bounded, which means that the supremum will be measurable and bounded with probability 1 in any finite sample (with respect to the empirical measure/conditional on the X data).

Since $P^*F < \infty$ we can choose M so that the term on the right is arbitrarily small.² That is, for any $\delta > 0$ we can pick M such that

$$\mathbb{E}^* \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - Pf(X_i) \right| \right] \leq \mathbb{E}_X \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}_M} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right] + \delta.$$

Step 2: Deal with the term that is conditional on $\{X_i\}$. Let $\mathcal{G}_\delta = \{g_1, \dots, g_{K(\delta)}\}$ be such that for every $f \in \mathcal{F}_M$ there is a $g \in \mathcal{G}_\delta$ such that $\|f - g\|_{L_1(\mathbb{P}_n)} < \delta$. Since $\log \mathcal{N}(\delta, \mathcal{F}_M, L_1(\mathbb{P}_n)) = o_p(n)$, we know that it is possible to pick a \mathcal{G}_δ in this fashion with probability approaching 1. Note that:

- Cardinality of \mathcal{G}_δ : $|\mathcal{G}_\delta| = \mathcal{N}(\delta, \mathcal{F}_M, \|\cdot\|_{L_1(\mathbb{P}_n)})$.
- Envelope of \mathcal{G}_δ : by construction $\mathcal{F}_M \leq M$ so we can assume that $\mathcal{G}_\delta \leq M$.

Then, for all $f \in \mathcal{F}_M$ we have that, for some $g \in \mathcal{G}_\delta$:

$$\begin{aligned} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| &\leq \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i g(X_i) \right| + \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(X_i) - g(X_i)) \right| \\ &\leq \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i g(X_i) \right| + \delta \end{aligned}$$

This gives us that

$$\mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}_M} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right] \leq \mathbb{E}_\epsilon \left[\sup_{g \in \mathcal{G}_\delta} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i g(X_i) \right| \right] + \delta.$$

Step 3: Apply the Maximal Inequality. We bound the first term in the last display using the maximal inequality in Lemma 3.2, for the particular case of the ψ_2 Orlicz norm³: if $D = \{f_1, \dots, f_M\}$ then

$$\mathbb{E} \left[\sup_{f \in D} |f(X_i)| \right] \leq C \sqrt{1 + \log m}$$

for any C with $\mathbb{E} \left[\exp \left(\frac{f(X_i)}{C^2} \right) - 1 \right] \leq 1$. In our setting we will apply this to the functions $\frac{1}{n} \sum_{i=1}^n \epsilon_i g(X_i)$ for $g \in \mathcal{G}_\delta$, with the $g(X_i)$ treated as fixed so that these are considered random variables in ϵ_i . In our setting we can bound:

$$\mathbb{E}_\epsilon \left[\sup_{g \in \mathcal{G}_\delta} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i g(X_i) \right| \right] \leq C \sqrt{1 + \log \mathcal{N}(\delta, \mathcal{F}_M, \|\cdot\|_{L_1(\mathbb{P}_n)})}.$$

for such a C such that

$$\mathbb{E}_\epsilon \left[\exp \left(\left(\frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right)^2 / C^2 \right) - 1 \right] \leq 1.$$

□

²I am sort of using P^* and \mathbb{E}_X^* interchangeably here, which I apologize for

³We know that, in any finite sample, this Orlicz norm exists because our empirical expectation is bounded.