

STA 711: Probability & Measure Theory

Robert L. Wolpert

7 Convergence in \mathbb{R}^d and in Metric Spaces

A sequence of elements a_n of \mathbb{R}^d converges to a limit a if and only if, for each $\epsilon > 0$, the sequence a_n eventually lies within a ball of radius ϵ centered at a . It's okay if the first few (or few million) terms lie outside that ball—and the number of terms that do lie outside the ball may depend on how big ϵ is (if ϵ is small enough it typically *will* take millions of terms before the remaining sequence lies inside the ball). This can be made mathematically precise by introducing a letter (say, N_ϵ) for how many initial terms we have to throw away, so that $a_n \rightarrow a$ if and only if there is an $N_\epsilon < \infty$ so that, for each $n \geq N_\epsilon$, $|a_n - a| < \epsilon$: only finitely many a_n can be farther than ϵ from a .

The same notion of convergence really works in any *metric space*, where we require that some measure of the *distance* $d(a_n, a)$ from a_n to a tend to zero in the sense that it exceeds each number $\epsilon > 0$ for at most some finite number N_ϵ of terms.

Points a_n in d -dimensional Euclidean space will converge to a limit $a \in \mathbb{R}^d$ if and only if each of their coordinates converges in \mathbb{R} ; and, since there are only finitely many coordinates, if they all converge then they do so *uniformly* (i.e., for each ϵ we can take the same N_ϵ for all d of the coordinate sequences), so all notions of convergence in \mathbb{R}^d are equivalent. For example,

$$\max_{1 \leq i \leq d} |x_i - y_i| \leq \left[\sum (x_i - y_i)^2 \right]^{\frac{1}{2}} \leq \sum_{1 \leq i \leq d} |x_i - y_i| \leq d \max_{1 \leq i \leq d} |x_i - y_i|$$

so convergence is identical for all three of these metrics. Convergence is much more complex and interesting for random variables.

7.1 Convergence of Random Variables

For *random variables* X_n the idea of convergence to a limiting random variable X is more delicate, since each X_n is a *function* of $\omega \in \Omega$ and usually there are infinitely many points $\omega \in \Omega$. What should we mean in saying that a sequence X_n converges to a limit X ? That $X_n(\omega)$ converges to $X(\omega)$ for each fixed ω ? Or that $X_n(\omega)$ converges *uniformly* in $\omega \in \Omega$? Or that some notion of the distance $d(X_n, X)$ between X_n and the limit X decreases to zero? Should the probability measure \mathbf{P} be involved in some way?

Here are a few different choices of what we *might* mean by the statement that “ X_n converges to X ,” for a sequence of random variables X_n and a random variable X , all defined on the same probability space $(\Omega, \mathcal{F}, \mathbf{P})$:

pw: The real numbers $X_n(\omega) \rightarrow X(\omega)$ for every $\omega \in \Omega$ (pointwise cgce):

$$(\forall \epsilon > 0) (\forall \omega \in \Omega) (\exists N_{\epsilon, \omega} < \infty) (\forall n \geq N_{\epsilon, \omega}) |X_n(\omega) - X(\omega)| < \epsilon.$$

uni: The sequences of real numbers $X_n(\omega) \rightarrow X(\omega)$ *uniformly* for $\omega \in \Omega$:

$$(\forall \epsilon > 0) (\exists N_\epsilon < \infty) (\forall \omega \in \Omega) (\forall n \geq N_\epsilon) |X_n(\omega) - X(\omega)| < \epsilon.$$

a.s.: Outside some null event $\mathcal{N} \in \mathcal{F}$, each sequence of real numbers $X_n(\omega) \rightarrow X(\omega)$ (Almost-Sure convergence, or convergence “almost everywhere” (*a.e.*)): for some $\mathcal{N} \in \mathcal{F}$ with $\mathbf{P}[\mathcal{N}] = 0$,

$$(\forall \epsilon > 0) (\forall \omega \notin \mathcal{N}) (\exists N_{\epsilon, \omega} < \infty) (\forall n \geq N_{\epsilon, \omega}) |X_n(\omega) - X(\omega)| < \epsilon,$$

$$i.e., \quad \mathbf{P}\left\{ \bigcup_{\epsilon > 0} \bigcap_{N < \infty} \bigcup_{n \geq N} |X_n(\omega) - X(\omega)| \geq \epsilon \right\} = 0.$$

L_∞ : Outside some null event $\mathcal{N} \in \mathcal{F}$, the sequences of real numbers $X_n(\omega) \rightarrow X(\omega)$ converge *uniformly* (“almost-uniform” or “ L_∞ ” convergence): for some $\mathcal{N} \in \mathcal{F}$ with $\mathbf{P}[\mathcal{N}] = 0$,

$$(\forall \epsilon > 0) (\exists N_\epsilon < \infty) (\forall \omega \notin \mathcal{N}) (\forall n \geq N_\epsilon) |X_n(\omega) - X(\omega)| < \epsilon.$$

pr.: For each $\epsilon > 0$, the probabilities $\mathbf{P}[|X_n - X| > \epsilon] \rightarrow 0$ (convergence “in probability”, or “in measure”):

$$(\forall \epsilon > 0) (\forall \eta > 0) (\exists N_{\epsilon, \eta} < \infty) (\forall n \geq N_{\epsilon, \eta}) \mathbf{P}[|X_n - X| > \epsilon] < \eta.$$

L_1 : The expectation $\mathbf{E}[|X_n - X|]$ converges to zero (convergence “in L_1 ”):

$$(\forall \epsilon > 0) (\exists N_\epsilon < \infty) (\forall n \geq N_\epsilon) \mathbf{E}[|X_n - X|] < \epsilon.$$

L_p : For some fixed number $p > 0$, the expectation of the p^{th} absolute power $\mathbf{E}[|X_n - X|^p]$ converges to zero (convergence “in L_p ,” sometimes called “in the p^{th} mean”):

$$(\forall \epsilon > 0) (\exists N_\epsilon < \infty) (\forall n \geq N_\epsilon) \mathbf{E}[|X_n - X|^p] < \epsilon.$$

dist.: The *distributions* of X_n converge to the distribution of X , *i.e.*, the measures $\mathbf{P} \circ X_n^{-1}$ converge in some way to $\mathbf{P} \circ X^{-1}$ (“vague” or “weak” convergence, or “convergence in distribution”, sometimes written $X_n \Rightarrow X$):

$$(\forall \epsilon > 0) (\forall \phi \in \mathcal{C}_b(\mathbb{R})) (\exists N_{\epsilon, \phi} < \infty) (\forall n \geq N_{\epsilon, \phi}) \mathbf{E}[|\phi(X_n) - \phi(X)|] < \epsilon.$$

Which of these eight notions of convergence is right for random variables? The answer is that *all* of them are useful in probability theory for one purpose or another. You will want to know which ones imply which other ones, under what conditions. All but the first two (pointwise, uniform) notions depend upon the measure \mathbf{P} ; it is possible for a sequence X_n to converge to X in any of these senses for one probability measure \mathbf{P} , but to fail to converge for another \mathbf{P}' . Most of them can be phrased as metric convergence for some notion of distance between random variables:

pr.: $X_n \rightarrow X$ in probability if and only if $d_0(X, X_n) \rightarrow 0$ as real numbers, where:

$$d_0(X, Y) := \mathbb{E} \left(\frac{|X - Y|}{1 + |X - Y|} \right)$$

L_1 : $X_n \rightarrow X$ in L_1 if and only if $d_1(X, X_n) := \|X - X_n\|_1 \rightarrow 0$ as real numbers, where:

$$\|Z\|_1 := \mathbb{E}|Z|$$

L_p : $X_n \rightarrow X$ in L_p if and only if $d_p(X, X_n) \rightarrow 0$ as real numbers, where:

$$d_p(X, Y) := \begin{cases} (\mathbb{E}|X - Y|^p)^{1/p} & p \geq 1 \\ \mathbb{E}|X - Y|^p & 0 < p < 1. \end{cases}$$

L_∞ : $X_n \rightarrow X$ almost uniformly if and only if $d_\infty(X, X_n) := \|X - X_n\|_\infty \rightarrow 0$ as real numbers, where:

$$\|Z\|_\infty := \sup\{r \geq 0 : \mathbb{P}[|Z| > r] > 0\}$$

As the notation suggests, convergence in probability and in L_∞ are in some sense limits of convergence in L_p as $p \rightarrow 0$ and $p \rightarrow \infty$, respectively. Almost-sure convergence is an exception: there is no metric notion of distance $d(X, Y)$ for which $X_n \rightarrow X$ almost surely if and only if $d(X, X_n) \rightarrow 0$ (unless Ω is countable or \mathbb{P} is atomic).

7.1.1 Almost-Sure Convergence

Let $\{X_n\}$ and X be a collection of RVs on some $(\Omega, \mathcal{F}, \mathbb{P})$. The set of points ω for which $X_n(\omega)$ does converge to $X(\omega)$ is just

$$\bigcap_{\epsilon > 0} \bigcup_{m=1}^{\infty} \bigcap_{n=m}^{\infty} [\omega : |X_n(\omega) - X(\omega)| \leq \epsilon],$$

the points which, for all $\epsilon > 0$, have $|X_n(\omega) - X(\omega)|$ less than ϵ for all but finitely-many n . The sequence X_n is said to converge “almost everywhere” (*a.e.*) to X , or to converge to X “almost surely” (*a.s.*), if this set of ω has probability one, or (conversely) if its complement is a null set:

$$\mathbb{P} \left[\bigcup_{\epsilon > 0} \bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} [\omega : |X_n(\omega) - X(\omega)| > \epsilon] \right] = 0.$$

Despite their appearance these intersections and unions over $\epsilon > 0$ are only countable, since we need include only rational ϵ (or, for that matter, any sequence ϵ_k tending to zero, such as $\epsilon_k = 1/k$). Thus $X_n \rightarrow X$ *a.e.* if and only if, for each $\epsilon > 0$,

$$\mathbf{P}\left[\bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} [\omega : |X_n(\omega) - X(\omega)| > \epsilon]\right] = 0. \quad (a.e.)$$

This combination of intersection and union occurs frequently in probability, and has a name; for any sequence E_n of events, $[\bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} E_n]$ is called the *lim sup* of the $\{E_n\}$, and is sometimes described more colorfully as $[E_n \text{ i.o.}]$, the set of points in E_n “infinitely often.” Its complement is the *lim inf* of the sets $F_n = E_n^c$, $[\bigcup_{m=1}^{\infty} \bigcap_{n=m}^{\infty} F_n]$: the set of points in all but finitely many of the F_n . Since \mathbf{P} is countably additive, and since the intersection in the definition of *lim sup* is *decreasing* and the union in the definition of *lim inf* is *increasing*, always we have

$\mathbf{P}[\bigcup_{n=m}^{\infty} E_n] \searrow \mathbf{P}[\bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} E_n]$ and $\mathbf{P}[\bigcap_{n=m}^{\infty} F_n] \nearrow \mathbf{P}[\bigcup_{m=1}^{\infty} \bigcap_{n=m}^{\infty} F_n]$ as $m \rightarrow \infty$. Thus,

Theorem 1 $X_n \rightarrow X$ *P*-a.s. if and only if for every $\epsilon > 0$,

$$\lim_{m \rightarrow \infty} \mathbf{P}[|X_n - X| > \epsilon \text{ for some } n \geq m] = 0.$$

In particular, $X_n \rightarrow X$ *P*-a.s. if $\sum_{n < \infty} \mathbf{P}[|X_n - X| > \epsilon] < \infty$ for each $\epsilon > 0$ (why?).

7.1.2 Convergence In Probability

The sequence X_n is said to converge to X “in probability” (*pr.*) if, for each $\epsilon > 0$,

$$\mathbf{P}[\omega : |X_n(\omega) - X(\omega)| > \epsilon] \rightarrow 0. \quad (pr.)$$

If we denote by E_n the event $[\omega : |X_n(\omega) - X(\omega)| > \epsilon]$ we see that convergence *almost surely* requires that $\mathbf{P}[\bigcup_{m \geq n} E_m] \rightarrow 0$ as $n \rightarrow \infty$, while convergence *in probability* requires only that $\mathbf{P}[E_n] \rightarrow 0$. For another proof, if $X_n \rightarrow X$ (*a.s.*) then the DCT implies that $\mathbf{P}[|X_n - X| > \epsilon] = \mathbf{E}\mathbf{1}_{\{|X_n - X| > \epsilon\}} \rightarrow 0$ for each $\epsilon > 0$. Thus:

Theorem 2 If $X_n \rightarrow X$ a.s. then $X_n \rightarrow X$ *pr.*

Here is a partial converse:

Theorem 3 If $X_n \rightarrow X$ *pr.*, then there is a subsequence n_k such that $X_{n_k} \rightarrow X$ a.s.

Proof. Set $n_0 := 0$ and, for each integer $k \geq 1$, set

$$n_k := \inf \left\{ n > n_{k-1} : \mathbf{P} \left[\omega : |X_n(\omega) - X(\omega)| > \frac{1}{k} \right] \leq 2^{-k} \right\}.$$

For any $\epsilon > 0$ we have $\frac{1}{k} \leq \epsilon$ eventually (namely, for $k \geq k_0 := \lceil \frac{1}{\epsilon} \rceil$) and for each $m \geq k_0$,

$$\begin{aligned} \mathbf{P}\left[\bigcup_{k=m}^{\infty} \left\{\omega : |X_{n_k}(\omega) - X(\omega)| > \epsilon\right\}\right] &\leq \mathbf{P}\left[\bigcup_{k=m}^{\infty} \left\{\omega : |X_{n_k}(\omega) - X(\omega)| > \frac{1}{k}\right\}\right] \\ &\leq \sum_{k=m}^{\infty} \mathbf{P}\left\{\omega : |X_{n_k}(\omega) - X(\omega)| > \frac{1}{k}\right\} \\ &\leq \sum_{k=m}^{\infty} 2^{-k} = 2^{1-m} \rightarrow 0 \text{ as } m \rightarrow \infty. \end{aligned}$$

□

In fact, this characterizes convergence in probability:

Theorem 4 *Let $\{X_n\}$, X be random variables on $(\Omega, \mathcal{F}, \mathbf{P})$. Then $X_n \rightarrow X$ pr. if and only if every sequence $\mathbb{N} \ni n_k \nearrow \infty$ has a subsequence n_{k_i} such that $X_{n_{k_i}} \rightarrow X$ a.s. as $i \rightarrow \infty$.*

Proof. The “only if” (\Rightarrow) direction is just Theorem 3. Suppose (for contradiction) that $X_n \not\rightarrow X$ pr.; then for some $\epsilon > 0$ and $\delta > 0$ there are infinitely-many n for which $\mathbf{P}[|X_n - X| > \epsilon] > \delta$. Let n_k be an increasing sequence satisfying this bound. By hypothesis, there is a subsequence along which $X_{n_{k_i}} \rightarrow X$ a.s.; but by Theorem 2, also $X_{n_{k_i}} \rightarrow X$ pr, so $\mathbf{P}[|X_n - X| > \epsilon] \rightarrow 0$, a contradiction. □

Theorem 5 *Let $\{X_n\}$, X be real-valued random variables on $(\Omega, \mathcal{F}, \mathbf{P})$ with $X_n \rightarrow X$ pr. and let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ be continuous. Then $Y_n := \phi(X_n) \rightarrow Y := \phi(X)$ pr.*

Proof. For an easy but indirect proof, simply apply Theorem 4. For a more direct approach, begin by selecting any $\epsilon > 0$ and $\delta > 0$. Find a compact set $K_\epsilon \subset \mathbb{R}$ with $\mathbf{P}[X \in K_\epsilon] \geq 1 - \epsilon/2$; since ϕ is uniformly continuous on K_ϵ , find $\eta > 0$ such that

$$(\forall x \in K_\epsilon)(\forall y \in \mathbb{R}) \quad |x - y| \leq \eta \quad \Rightarrow \quad |\phi(x) - \phi(y)| < \delta.$$

Now, since $X_n \rightarrow X$ pr., find $N \in \mathbb{N}$ such that

$$(\forall n \geq N) \quad \mathbf{P}[|X_n - X| > \eta] \leq \epsilon/2.$$

Then for $n \geq N$,

$$\begin{aligned} \mathbf{P}[|Y_n - Y| > \delta] &\leq \mathbf{P}[X \notin K_\epsilon] + \mathbf{P}[X \in K_\epsilon, |\phi(X_n) - \phi(X)| > \delta] \\ &\leq \mathbf{P}[X \notin K_\epsilon] + \mathbf{P}[X \in K_\epsilon, |X_n - X| > \eta] \\ &\leq \epsilon/2 + \epsilon/2 = \epsilon. \end{aligned}$$

□

The same result (with the same proof) is true for random variables X_n, X taking values in any σ -compact complete separable metric space \mathcal{X} and, in particular, for $\mathcal{X} = \mathbb{R}^d$.

7.1.3 A Counter-Example

If $X_n \rightarrow X$ *a.s.* implies $X_n \rightarrow X$ *pr.*, and if the converse holds at least along subsequences, are the two notions really identical? Or is it possible for RVs X_n to converge to X *pr.*, but not *a.s.*? The answer is that the two notions *are* different, and that *a.s.* convergence is strictly stronger than convergence *pr.* Here's an example:

First notice that every integer $n \in \mathbb{N}$ can be written uniquely in the form $n = i + 2^j$ for integers $j \geq 0$ and $0 \leq i < 2^j$ (set $j := \lfloor \log_2 n \rfloor$ and $i := n - 2^j$). Let $(\Omega, \mathcal{F}, \mathbf{P})$ be the unit interval with Borel sets and Lebesgue measure. Define a sequence of random variables $X_n : \Omega \rightarrow \mathbb{R}$ by

$$X_n(\omega) = \begin{cases} 1 & \text{if } \frac{i}{2^j} < \omega \leq \frac{i+1}{2^j} \\ 0 & \text{otherwise} \end{cases} \quad \text{where } n = i + 2^j, 0 \leq i < 2^j.$$

SO,

X_1 is one on 1 interval of length 1,
 X_2, X_3 are one on 2 intervals each of length 1/2,
 X_4, \dots, X_7 are one on 4 intervals each of length 1/4,
 X_8, \dots, X_{15} are one on 8 intervals each of length 1/8,

and, in general, each X_n is one on an interval of length 2^{-j} . Since $\frac{1}{n} \leq \frac{1}{2^j} < \frac{2}{n}$,

$$\mathbf{P}[|X_n| > \epsilon] = 2^{-j} < \frac{2}{n} \rightarrow 0$$

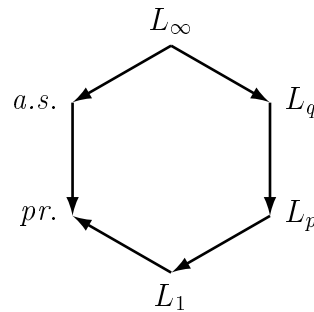
for each $0 < \epsilon < 1$ and $X_n \rightarrow 0$ *pr.* On the other hand, for every $j > 0$ we have each $\omega \in \Omega$ in one of the 2^j intervals of length 2^{-j} where some X_n is one,

$$(\forall j \in \mathbb{N}) \quad \Omega = \bigcup_{i=0}^{2^j-1} \left(\frac{i}{2^j}, \frac{i+1}{2^j} \right] = \bigcup_{n=2^j}^{2^{j+1}-1} [\omega : X_n(\omega) = 1]$$

For every $\omega \in \Omega$ and $j \in \mathbb{N}$ there is some $n \geq 2^j$ with $X_n(\omega) = 1$ (and so there are infinitely-many such X_n); thus $[\omega : X_n(\omega) \rightarrow 0]$ is *empty*, not a set of probability one! Obviously X_n does not converge *a.s.*, although it does converge *pr.*

This example is a building-block for several examples to come, so getting to know it well is worthwhile. Try to verify that $X_n \rightarrow 0$ in probability (how large must n be to ensure $\mathbf{P}[|X_n| > \delta] < \epsilon$?) and in L_p (how large must n be to ensure $\|X_n\|_p < \epsilon$?) but *not* almost surely (what is $\{\omega : X_n \rightarrow 0\}$? Why?). Find $\|X_n\|_p$ explicitly. Why *doesn't* $X_n \rightarrow 0$ *a.s.*? What would happen if we multiplied X_n by n ? By n^2 ? By $j := \lfloor \log_2 n \rfloor$? What about the subsequence $Y_n := X_{2^n}$? Does X_n converge in L_∞ ?

In summary, for $1 < p < q < \infty$ the convergence implications are:



with *partial* converses (*pr.* \rightarrow *a.s.* along subsequences, *pr.* $\rightarrow L_1, L_p, L_q$ under UI— see the next section). All of these imply convergence in distribution, which we'll consider later.

7.2 Uniform Integrability

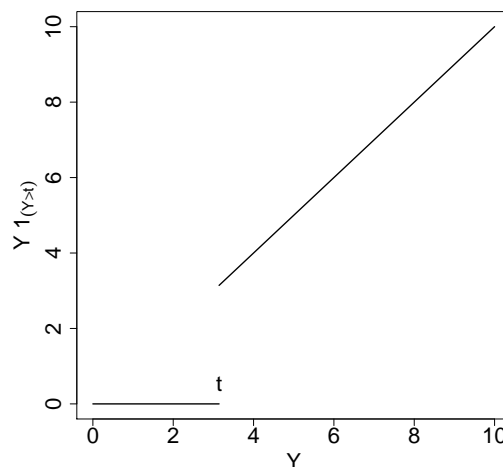
Let $Y \geq 0$ be integrable on some probability space $(\Omega, \mathcal{F}, \mathbf{P})$,

$$\mathbf{E}[Y] = \int_{\Omega} Y \, d\mathbf{P} < \infty.$$

By Lebesgue's DCT it follows that

$$\lim_{t \rightarrow \infty} \mathbf{E}[Y \mathbf{1}_{\{Y > t\}}] = \lim_{t \rightarrow \infty} \int_{[\omega: Y(\omega) > t]} Y \, d\mathbf{P} = 0$$

since $\{Y \mathbf{1}_{\{Y > t\}}\}$ is dominated by $Y \in L_1$ and converges to zero *a.s.* as $t \rightarrow \infty$.



Consequently, for any sequence of random variables X_n that are dominated by Y in the sense that $|X_n| \leq Y$ *a.s.*,

$$\mathbf{E}[|X_n| \mathbf{1}_{\{|X_n| > t\}}] \leq \mathbf{E}[Y \mathbf{1}_{\{Y > t\}}] \rightarrow 0, \text{ uniformly in } n.$$

Call the sequence X_n *uniformly integrable* (or simply UI) if $\mathbf{E}[|X_n| \mathbf{1}_{\{|X_n|>t\}}] \rightarrow 0$ uniformly in n as $t \rightarrow \infty$, even if it isn't dominated by some integrable $Y \in L_1$. The big result is:

Theorem 6 *If $X_n \rightarrow X$ pr. and if X_n is UI then $X_n \rightarrow X$ in L_1 .*

Proof. Without loss of generality take $X \equiv 0$. Fix any $\epsilon > 0$; find (by UI) $t_\epsilon > 0$ such that $\mathbf{E}[|X_n| \mathbf{1}_{\{|X_n|>t_\epsilon\}}] \leq \epsilon$ for all n . Now (by $X_n \rightarrow X$ pr.) find $N_\epsilon \in \mathbb{N}$ such that, for $n \geq N_\epsilon$, $\mathbf{P}[|X_n| > \epsilon] < \epsilon/t_\epsilon$. Then:

$$\begin{aligned} \mathbf{E}[|X_n|] &= \int_{\{|X_n| \leq \epsilon\}} |X_n| d\mathbf{P} + \int_{\{\epsilon < |X_n| \leq t_\epsilon\}} |X_n| d\mathbf{P} + \int_{\{t_\epsilon < |X_n|\}} |X_n| d\mathbf{P} \\ &\leq \int_{\{|X_n| \leq \epsilon\}} \epsilon d\mathbf{P} + \int_{\{\epsilon < |X_n| \leq t_\epsilon\}} t_\epsilon d\mathbf{P} + \int_{\{t_\epsilon < |X_n|\}} |X_n| d\mathbf{P} \\ &\leq \epsilon + t_\epsilon \times \mathbf{P}[|X_n| > \epsilon] + \epsilon \\ &\leq 3\epsilon. \end{aligned}$$

□

Theorem 6 has a partial converse— if $\{X_n\} \subset L_1$ and $X_n \rightarrow X$ pr. and if $\mathbf{E}|X_n| \rightarrow \mathbf{E}|X| < \infty$, then $\{X_n\}$ is UI (see Theorem 9). For another proof, first note $\mathbf{E}[|X_n| \mathbf{1}_{\{|X_n|>t\}}] = \mathbf{E}[|X_n|] - \mathbf{E}[|X_n| \mathbf{1}_{\{|X_n| \leq t\}}]$. For $\epsilon > 0$ first pick $t < \infty$ s.t. $\mathbf{E}[|X| \mathbf{1}_{\{|X|>t\}}] < \epsilon$ (possible, since $X \in L_1$); then (by DCT) find N s.t. for $n \geq N$, $\|X_n \mathbf{1}_{\{|X_n| \leq t\}} - X \mathbf{1}_{\{|X| \leq t\}}\|_1 < \epsilon$ and (by assumption) $|\mathbf{E}|X_n| - \mathbf{E}|X|| < \epsilon$. Then $\mathbf{E}[|X_n| \mathbf{1}_{\{|X_n|>t\}}] \leq 3\epsilon$ for $n \geq N$; since $X_j \in L_1$ for $j < N$, the result follows (with $t^* = \max(t, t_1, \dots, t_{N-1})$).

Similarly, for any $p > 0$, $X_n \rightarrow X$ (pr.) and $|X_n|^p$ UI (for example, $|X_n| \leq Y \in L_p$, or $\|X_n\|_q \leq B < \infty$ for some $q > p$) gives $X_n \rightarrow X$ (L_p). In the special case of $|X_n| \leq Y \in L_1$ this is just Lebesgue's **Dominated Convergence Theorem** (DCT), in a little stronger version than we proved in Week 4 because here we require only $X_n \rightarrow X$ (pr.), while then we required $X_n \rightarrow X$ (a.s.).

We have seen that $\{X_n\}$ is UI whenever $|X_n| \leq Y \in L_1$, but UI is more general than that. Here are two more criteria. The first is called the “crystal ball condition”:

Theorem 7 *If $\{X_n\}$ is uniformly bounded in L_p for some $p > 1$ then $\{X_n\}$ is UI.*

Proof. Let $B \in \mathbb{R}_+$ be an upper bound for $\|X_n\|_p$ and set $q := \frac{p}{p-1}$. By Hölder's inequality

$$\begin{aligned} \mathbf{E}[|X_n| \mathbf{1}_{\{|X_n|>t\}}] &\leq \|X_n\|_p \|\mathbf{1}_{\{|X_n|>t\}}\|_q \\ &\leq B \{\mathbf{P}[|X_n| > t]\}^{1/q} \\ &\leq B \{\mathbf{E}[|X_n|^p/t^p]\}^{1/q} \quad \text{by Markov's inequality} \\ &= B \{\|X_n\|_p^p/t^p\}^{1/q} \\ &= B \|X_n\|_p^{p-1}/t^{p-1} \quad \text{since } p/q = p-1 \\ &\leq B^p t^{1-p} \rightarrow 0 \end{aligned}$$

uniformly in $\{n\}$.

□

Theorem 8 *The random variables $\{X_\alpha\}$ are UI if and only if they are uniformly bounded in L_1 and*

$$(\forall \epsilon > 0)(\exists \delta > 0)(\forall A \in \mathcal{F}) \quad \mathbf{P}(A) < \delta \Rightarrow \mathbf{E}[|X_\alpha| \mathbf{1}_A] < \epsilon.$$

If $(\Omega, \mathcal{F}, \mathbf{P})$ is non-atomic, the condition “ $\{X_\alpha\}$ is uniformly bounded in L_1 ” is unnecessary.

Proof. Straightforward. If $(\Omega, \mathcal{F}, \mathbf{P})$ is non-atomic, for each $\epsilon > 0$ find $\delta > 0$ by hypothesis and then cover Ω with $N := \lfloor 1 + 1/\delta \rfloor$ sets $A_i \in \mathcal{F}$ with $\mathbf{P}(A_i) < \delta$ to see $\mathbf{E}|X_\alpha| < N\epsilon$ uniformly. Try to offer an example to illustrate what can go wrong if \mathbf{P} has an atom $\omega^* \in \Omega$ with $c = \mathbf{P}\{\omega^*\} > 0$.

In fact, $X_n \rightarrow X$ in L_1 **if and only if** $\{X_i\}$ are UI:

Theorem 9 *Let $\{X_n\}, X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ with $X_n \rightarrow X$ pr. Then the following are equivalent:*

1. $\{X_n\}$ is UI
2. $X_n \rightarrow X$ in L_1
3. $\mathbf{E}|X_n| \rightarrow \mathbf{E}|X|$

Proof. Since $1 \Rightarrow 2$ by Theorem 6 and $2 \Rightarrow 3$ by the triangle inequality, it remains only to show that $3 \Rightarrow 1$. Let $X_n \rightarrow X$ pr. and $\mathbf{E}|X_n| \rightarrow \mathbf{E}|X|$. Fix $\epsilon > 0$ and, for each $t > 0$, set

$$f_t(x) := \begin{cases} |x| & 0 \leq |x| \leq t \\ t(t+1-|x|) & t < |x| \leq t+1 \\ 0 & t+1 < |x| < \infty \end{cases}$$

and note that $0 \leq |x| \mathbf{1}_{\{|x| \leq t\}} \leq f_t(x) \leq |x| \mathbf{1}_{\{|x| \leq t+1\}}$ and $0 \leq f_t(x) \leq t$.

Since $X \in L_1$, $\mathbf{E}|X| \mathbf{1}_{\{|X| > t\}} \rightarrow 0$ as $t \rightarrow \infty$ by DCT so we can find t sufficiently large that

$$\mathbf{E}|X| \mathbf{1}_{\{|X| > t\}} < \epsilon. \tag{1}$$

Since f_t is continuous and bounded, $f_t(X_n) \rightarrow f_t(X)$ as $n \rightarrow \infty$ in probability by Theorem 5 and also in L_1 by DCT. For large enough n (say, $n \geq N'$), $\|f_t(X_n) - f_t(X)\|_1 \leq \epsilon$ and so

$$\mathbf{E}|X_n| \mathbf{1}_{\{|X_n| \leq t+1\}} \geq \mathbf{E}f_t(X_n) \geq \mathbf{E}f_t(X) - \epsilon \geq \mathbf{E}|X| \mathbf{1}_{\{|X| \leq t\}} - \epsilon. \tag{2}$$

Since $\mathbf{E}|X_n| \rightarrow \mathbf{E}|X|$, for sufficiently large n (say, $n \geq N \geq N'$)

$$\mathbf{E}|X_n| \leq \mathbf{E}|X| + \epsilon. \tag{3}$$

Upon subtracting (2) from (3) we have

$$\mathbf{E}|X_n| \mathbf{1}_{\{|X_n| > t+1\}} \leq \mathbf{E}|X| \mathbf{1}_{\{|X| > t\}} + 2\epsilon.$$

Applying (1) we get

$$\mathbf{E}|X_n| \mathbf{1}_{\{|X_n| > t+1\}} < 3\epsilon$$

completing the proof that $\{X_n\}$ is UI.

□

7.3 Cauchy Convergence

Every form of convergence we have considered for random variables *except* a.s. can be represented as convergence in a metric space. A sequence X_n converges to a limit X in a metric space if every ball centered at X contains all but finitely-many of $\{X_n\}$. Sometimes we wish to consider a sequence X_n that converges to *some* limit, without knowing the limit in advance. The concept of *Cauchy Convergence* is ideal for this— we insist that for each $\epsilon > 0$, all but finitely-many of the ϵ -balls centered *at points X_m of the sequence* contain all but finitely-many of the points. For any of distance measures d_p of Section (7.1), with $0 \leq p \leq \infty$, say “ X_n is a Cauchy sequence in L_p ” if

$$(\forall \epsilon > 0)(\exists N_\epsilon < \infty)(\forall m, n \geq N_\epsilon) \quad d_p(X_m, X_n) < \epsilon.$$

The spaces L_p for $0 \leq p \leq \infty$ are all *complete* in the sense that *if X_n is Cauchy for d_p then there exists $X \in L_p$ for which $d_p(X_n, X) \rightarrow 0$* . To see this, take an increasing subsequence N_k along which $d_p(X_m, X_n) < 2^{-k}$ for $n \geq m \geq N_k$, and set $X_0 := 0$ and $N_0 = 0$; for $k \in \mathbb{N}$ set $Y_k := X_{N_k} - X_{N_{k-1}}$. Check to confirm that $X := \sum_{k=1}^{\infty} Y_k$ is well-defined and $X \in L_p$ a.s., and that $d_p(X_n, X) \rightarrow 0$.

7.4 Convergence in Distribution

Some of the most famous results in probability theory concern the limiting *distribution* of sequences of random variables. For example,

- **Central Limit Theorem:** For iid $\{X_n\} \subset L_2$ and large n , $(S_n - n\mu)/\sqrt{n\sigma^2}$ has approximately the Normal $\text{No}(0, 1)$ distribution;
- **Law of Small Numbers:** For independent \mathbb{N}_0 -valued random variables X_i with small means μ_i , S_n has approximately the Poisson $\text{Po}(\lambda_n)$ distribution with mean $\lambda_n = \sum \mu_i$;
- **Extreme Value or Three Types Theorem:** For iid $\{X_n\}$ and suitable sequences $\{a_n\}$, $\{b_n\}$ of constants, $\max\{(X_j - a_n)/b_n : j \leq n\}$ has approximately the Gumbel, Fréchet, or reversed Weibull distribution.

But what does it mean to have a distribution “approximately”???

If X_n has the discrete uniform distribution on the points $\{i/n\}$ for $1 \leq i \leq n$, and if n is huge (say, $2^{51} \approx 10^{17}$) then most of us would regard X_n as “approximately” a standard uniform $\text{Un}(0, 1)$ random variable— in fact, this is the best one can hope for in a 64-bit double-precision floating-point representation. Yet each X_n has a *discrete* distribution while the limit is a *continuous* distribution, so whatever we mean by “converging in distribution” it can’t involve either density functions or probability mass functions, and it is too much to ask for the distributions $\mu_n(B) = \mathbf{P}[X_n \in B]$ to converge for every Borel set $B \subset \mathbb{R}$.

In this example the distribution functions $F_n(x) = \mathbf{P}[X_n \leq x]$ do in fact converge at every $x \in \mathbb{R}$, but that too is a bit too much to ask— consider constant random variables $Y_n :=$

$-1/n$, for example, which converge almost-surely to $Y := 0$, but the distribution functions do not converge at the point $y = 0$. This is closer to the right idea, however:

Definition 1 Let $\{\mu_n(dx)\}$ and $\mu(dx)$ be distributions on a measurable space $(\mathcal{X}, \mathcal{E})$. Say “ μ_n converges in distribution to μ ” or write “ $\mu_n \Rightarrow \mu$ ” if

$$\int_{\mathcal{X}} \phi(x) \mu_n(dx) \rightarrow \int_{\mathcal{X}} \phi(x) \mu(dx)$$

as $n \rightarrow \infty$ for every continuous bounded function $\phi : \mathcal{X} \rightarrow \mathbb{R}$.

If $\{X_n\}$ are random variables with distributions $\{\mu_n\}$, we also say that “ X_n converges in distribution to μ ” and write “ $X_n \Rightarrow \mu$.” This definition works not only for real-valued random variables where $\mathcal{X} = \mathbb{R}$ and $\mathcal{E} = \mathcal{B}$ is the Borel sets, but also in \mathbb{R}^d or any other complete separable metric space \mathcal{X} with Borel sets \mathcal{E} . The connection with distribution functions is given by:

Proposition 1 Let $\{\mu_n(dx)\}$ and $\mu(dx)$ be distributions on the real line $(\mathbb{R}, \mathcal{B})$, and let $F_n(x) := \mu_n((-\infty, x])$ and $F(x) := \mu((-\infty, x])$ be the associated DFs. Then the following are equivalent:

- $\mu_n \Rightarrow \mu$;
- $F_n(x) \rightarrow F(x)$ on any dense set of $\{x\}$ in \mathbb{R} ;
- $F_n(x) \rightarrow F(x)$ at every $x \in \mathbb{R}$ where $F(x-) = F(x)$ is continuous.

A similar but more general result for complete separable metric spaces \mathcal{X} (like \mathbb{R}^n) with their Borel sets is that $\mu_n \Rightarrow \mu$ if and only if $\mu_n(B) \rightarrow \mu(B)$ for every Borel set B whose boundary $\partial B := \overline{B} \cap \overline{B^c}$ has zero measure for the limiting distribution, $\mu(\partial B) = 0$. Convergence in distribution is weaker than any of the other forms of convergence we have considered:

Proposition 2 Let $\{X_n\}$, X be random variables on some probability space $(\Omega, \mathcal{F}, \mathbf{P})$ with distributions $\{\mu_n\}$, μ , and suppose $X_n \rightarrow X$ pr. Then $X_n \Rightarrow \mu$.

The proof is a simple application of the UI Convergence Theorem to the (uniformly bounded and hence UI) random variables $\phi(X_n)$.

Convergence in distribution is *strictly* weaker than any other notion, because it even applies if the $\{X_n\}$ are defined on different probability spaces $(\Omega_n, \mathcal{F}_n, \mathbf{P}_n)$, where none of the other forms of convergence even makes sense. As usual there is a *partial* converse, however:

Proposition 3 Let $\mu_n \Rightarrow \mu$ on the real line $(\mathbb{R}, \mathcal{B})$. Then there exists a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and random variables $\{X_n\}$, X with distributions $\{\mu_n\}$, μ for which $X_n \rightarrow X$ a.s.

It is enough to take $\Omega = (0, 1]$, $\mathcal{F} = \mathcal{B}$, $\mathbf{P} = \lambda$, and for $\omega \in \Omega$ set

$$X_n(\omega) := \inf \{x \in \bar{\mathbb{R}}_+ : F_n(x) \geq \omega\} \quad X(\omega) := \inf \{x \in \bar{\mathbb{R}}_+ : F(x) \geq \omega\}$$

and verify that $X_n(\omega) \rightarrow X(\omega)$ for all but countably-many $\omega \in \Omega$.

7.5 Metrics on Distributions [Optional]

Convergence in distribution can be metrized by the *Lévy-Prokhorov* metric. For Borel probability measures μ, ν on a complete separable metric space (\mathcal{X}, d) , set

$$d_{LP}(\mu, \nu) := \inf \left\{ \epsilon > 0 : \begin{aligned} &\mu(A) \leq \nu(A^\epsilon) + \epsilon \\ &\text{and } \nu(A) \leq \mu(A^\epsilon) + \epsilon \text{ for all } A \in \mathcal{B}(\mathcal{X}) \end{aligned} \right\} \quad (4a)$$

where $A^\epsilon := \{y \in \mathcal{X} : (\exists x \in A) d(x, y) < \epsilon\} = \cup_{x \in A} B_\epsilon(x)$, the union of open ϵ -balls around all points in A . In $d = 1$ dimension this reduces to the *Lévy metric* for the DFs F and G of μ and ν ,

$$d_{LP}(F, G) := \inf \{ \epsilon > 0 \mid F(x - \epsilon) \leq G(x) \leq F(x + \epsilon) \text{ for all } x \in \mathbb{R} \}$$

While convergence in probability (and metric (4a)) is the most important notion of closeness and convergence of probability distributions, several others arise. Many of these are described and compared in the Central Limit Theorem notes for Week 9 at

<http://www.stat.duke.edu/courses/Fall116/sta711/lec/topics/dstn.pdf>.

Some of the most important are:

Total Variation:

$$\begin{aligned} d_{TV}(\mu, \nu) &:= \sup \{ |\mu(A) - \nu(A)| : A \in \mathcal{B}(\mathcal{X}) \} \\ &= \inf \{ \mathbf{P}[X \neq Y] : X \sim \mu, Y \sim \nu \}. \end{aligned} \quad (4b)$$

If μ and ν have density functions with respect to some σ -finite reference measure λ (like Lebesgue measure, for continuous distributions, or counting measure, for discrete ones), this can be evaluated as half the L_1 distance between their densities,

$$= \frac{1}{2} \int_{\mathcal{X}} |f(x) - g(x)| \lambda(dx),$$

Hellinger:

Again if $\mu(dx) = f(x)\lambda(dx)$ and $\nu(dx) = g(x)\lambda(dx)$ have densities, the *Hellinger Distance* is

$$\begin{aligned} d_H(\mu, \nu) &:= \left\{ \frac{1}{2} \int_{\mathcal{X}} \left(\sqrt{f(x)} - \sqrt{g(x)} \right)^2 \lambda(dx) \right\}^{1/2} \\ &= \left\{ 1 - \int_{\mathcal{X}} \sqrt{f(x)g(x)} \lambda(dx) \right\}^{1/2} \end{aligned}$$

Hellinger and Total Variation determine the same topology, *i.e.*, so a sequence converges in one if and only if it does in the other.

Kolmogorov-Smirnov:

The *Kolmogorov-Smirnov* distance between two probability measures on the real line $\mathcal{X} = \mathbb{R}$ is the L_∞ distance between their DFs:

$$d_{KS}(\mu, \nu) := \sup_{x \in \mathbb{R}} \{ |\mu((-\infty, x]) - \nu((-\infty, x])| \} \quad (4c)$$

Wasserstein W_1 :

The *Wasserstein* W_1 distance between two distributions on a complete separable metric space (\mathcal{X}, d) is

$$\begin{aligned} d_W(\mu, \nu) &:= \inf_{\gamma \in \Gamma(\mu, \nu)} \iint_{\mathcal{X} \times \mathcal{X}} d(x, y) \gamma(dx dy) \\ &= \inf \{ \|X - Y\|_1 : X \sim \mu, Y \sim \nu \} \end{aligned} \quad (4d)$$

where $\Gamma(\mu, \nu)$ is the space of probability measures on \mathcal{X}^2 with marginals μ and ν . It can also be evaluated in terms of unit Lipschitz functions $f(x)$ as:

$$= \sup_{f: \mathcal{X} \rightarrow \mathbb{R}} \left\{ \left| \int_{\mathcal{X}} f(x) \mu(dx) - \int_{\mathcal{X}} f(x) \nu(dx) \right| : \text{Lip}(f) \leq 1 \right\}.$$

Sometimes called the “transportation metric,” d_W can be interpreted as the minimum cost of moving mass distributed on \mathcal{X} according to μ to mass distributed according to ν if moving cost is proportional to the product of mass times distance. There are also L_p versions of this metric, the Wasserstein W_p distance, but they are encountered less frequently.

Kullback-Leibler Divergence:

The “Kullback-Leibler divergence” (Kullback and Leibler, 1951), also called Relative Entropy, from distribution μ (with pdf $f(x)$) to ν (with pdf $g(x)$) on a Polish space \mathcal{X} is:

$$\text{KL}(\mu \| \nu) := \int_{\mathcal{X}} -\log \left[\frac{g(x)}{f(x)} \right] f(x) \lambda(dx). \quad (4e)$$

It is nonnegative, because $\log y \leq y - 1$ for all $y > 0$ (or by Jensen’s inequality), but it is not actually a distance metric because it’s not symmetric in μ and ν and doesn’t satisfy the triangle inequality. It does determine a topology, though, and hence a notion of convergence.

7.6 Summary: UI and Convergence Concepts**I. Uniform Integrability (UI)**

- A. $|X_n| \leq Y \in L_r$, $r > 0$, implies $\int_{\{|X_n| > t\}} |X_n|^r d\mathbf{P} \leq \int_{\{Y > t\}} Y^r d\mathbf{P} \rightarrow 0$ as $t \rightarrow \infty$, uniformly in n , so $\{|X_n|^r\}$ is UI. This is the *definition* of UI for $r = 1$.
- B. If $(\Omega, \mathcal{F}, \mathbf{P})$ nonatomic, X_n UI iff $(\forall \epsilon)(\exists \delta) \ni \mathbf{P}[\Lambda] \leq \delta \Rightarrow \int_{\Lambda} |X_n| d\mathbf{P} < \epsilon$ (take $\delta = \frac{\epsilon}{2t}$). If $(\Omega, \mathcal{F}, \mathbf{P})$ has atoms, also need uniform bound $\mathbf{E}|X_n| \leq B$.

C. $\mathbf{E}|X_n|^q \leq c^q < \infty$ implies $|X_n|^p$ UI for each $p < q$.

1. Remark: *not* for $p = q$ (counter-eg: $X_n := n\mathbf{1}_{(0,1/n]}$)

D. Main result: If $X_n \rightarrow X$ *pr.*, then

$$|X_n|^p \text{ is UI} \iff X_n \rightarrow X \text{ in } L_p \iff \mathbf{E}|X_n|^p \rightarrow \mathbf{E}|X|^p.$$

II. Convergence in Distribution (aka *Vague Convergence*)

A. $X_n \rightarrow X$ *pr.* iff $(\forall n_k \uparrow)(\exists n_{k_i} \uparrow) \ni X_{n_{k_i}} \rightarrow X$ *a.s.* (by contradiction)

B. $X_n \rightarrow X$ *a.s.* and $g(x)$ continuous implies $g(X_n) \rightarrow g(X)$ *a.s.*

C. $X_n \rightarrow X$ *pr.* and $g(x)$ continuous $\Rightarrow g(X_n) \rightarrow g(X)$ *pr.* (use A, B)

D. Definition: $X_n \Rightarrow X$ if $(\forall \phi \in C_b(\mathbb{R})) \mathbf{E}\phi(X_n) \rightarrow \mathbf{E}\phi(X)$

1. Prop: $X_n \rightarrow X$ *pr.* implies $X_n \Rightarrow X$ (use II.C)

2. Prop: $X_n \Rightarrow X$ implies $F_n(r) \rightarrow F(r)$ for all r s.t. $F(r) = F(r-)$.

a. Remark: Even if $X_n \Rightarrow X$, $F_n(r)$ may not converge where $F(r)$ jumps;

b. Remark: Even if $X_n \Rightarrow X$, $f_n(r) := F'_n(r)$ may not converge to $f(r) := F'(r)$; in fact, either may fail to exist.

III. Implications among cgce notions: *a.s.*, *pr.*, L_p , L_q , L_∞ , *dist.* ($0 < p < q < \infty$):

A. *a.s.* \implies *pr.* (by Easy Borel-Cantelli)

1. *pr.* \implies *a.s.* along subsequences

2. *pr.* $\not\Rightarrow$ *a.s.* (counter-eg: $X_n(\omega) = \mathbf{1}_{(i/2^j, (i+1)/2^j]}(\omega)$, $n = i + 2^j$)

B. $L_p \implies$ *pr.* (by Markov's inequality)

1. *pr.* $\implies L_p$ under Uniform Integrability

2. *pr.* $\not\Rightarrow L_p$ (counter-eg: $X_n = n^{1/p}\mathbf{1}_{(0,1/n]}$)

C. $L_q \implies L_p$ for $p < q$ (by Jensen's inequality)

1. $L_p \not\Rightarrow L_q$ (counter-eg: $X_n = n^{1/q}\mathbf{1}_{(0,1/n]}$)

D. $L_\infty \implies L_p$ (simple estimate, or $\|X\|_p \nearrow$ as $p \nearrow$ by Jensen)

1. $L_p \not\Rightarrow L_\infty$ (counter-eg: $X_n = \mathbf{1}_{(0,1/n]}$)

E. $L_\infty \implies$ *a.s.* (almost-uniform cgce implies almost-pointwise cgce)

F. *pr.* \implies *dist.* (II.D.1 above)

1. *dist.* $\not\Rightarrow$ *pr.* (counter-eg: X_n, X on different spaces or iid $\{X_n\}$)

2. *dist.* \implies *a.s.* ($(\exists(\Omega, \mathcal{F}, \mathbf{P}), X_n, X) \ni X_n \rightarrow X$ *a.s.*)

References

Kullback, S. and Leibler, R. A. (1951), "On Information and sufficiency," *Annals of Mathematical Statistics*, 22, 79–86.