

STA 711: Probability & Measure Theory

Robert L. Wolpert

6 Independence (cont)

6.5 B/C + Independence Illustration

Let $X \in L_1(\Omega, \mathcal{F}, \mathbb{P})$ and let $\{A_n\} \subset \mathcal{F}$ be events that satisfy $\mathbb{P}[A_n] \rightarrow 0$. Does it follow that $X_n := X\mathbf{1}_{A_n}$ converges almost-surely to 0?

If $\sum_n \mathbb{P}[A_n] < \infty$, then *yes*— by the Borel-Cantelli lemma,

$$\mathbb{P}[X_n \not\rightarrow 0] \leq \mathbb{P}[\limsup A_n] = 0,$$

so $X_n \rightarrow 0$ *a.s.*

BUT, if $\{\mathbb{P}[A_n]\}$ is not summable, then *a.s.* convergence can fail. For example, if $\{A_n\}$ are independent and $X \equiv 1$, then

$$\mathbb{P}[X_n \not\rightarrow 0] \geq \mathbb{P}[\limsup A_n] = 1,$$

so $\mathbb{P}[X_n \rightarrow 0] = 0$. In Week 7 we will find a new sense of convergence called “convergence in probability” that is weaker than almost-sure convergence, and we’ll show that $X_n \rightarrow 0$ *pr.*

6.6 Another Zero-One Law: Kolmogorov’s

For any collection $\{X_n\}$ of random variables on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, define two sequences of σ -algebras (“past” and “future”) by:

$$\mathcal{F}_n := \sigma\{X_i : i \leq n\} \quad \mathcal{T}_n := \sigma\{X_i : i \geq n+1\}$$

and, from them, construct the π -system \mathcal{P} and “tail” σ -algebra \mathcal{T} by

$$\mathcal{P} := \bigcup_{n=1}^{\infty} \mathcal{F}_n \quad \mathcal{T} := \bigcap_{n=1}^{\infty} \mathcal{T}_n.$$

In general \mathcal{P} will not be a σ -algebra, because it will not be closed under countable unions or intersections, but it is a field and hence a π -system, and generates the σ -algebra $\vee \mathcal{F}_n := \sigma(\mathcal{P}) \subseteq \mathcal{F}$.

The class \mathcal{T} , called the *tail* σ -field, includes those events that depend only on what happens *eventually*, regardless of what happens for the first few (or few million) $\{X_n\}$. These include such events as “ $\{X_n \text{ converges}\}$ ” or “ $\{\limsup X_n \leq 1\}$ ” or, with $S_n := \sum_1^n X_j$, “ $\{\frac{1}{n}S_n \text{ converges}\}$ ” or “ $\{\frac{1}{n}S_n \rightarrow 0\}$,” but not events like “ $\{\min X_n \leq c\}$.”

Theorem 1 (Kolmogorov's Zero-One Law) *For independent random variables X_n , the tail σ -field \mathcal{T} is “almost trivial” in the sense that every event $\Lambda \in \mathcal{T}$ has probability $P[\Lambda] = 0$ or $P[\Lambda] = 1$.*

Proof. Let $A \in \mathcal{P} = \cup \mathcal{F}_n$, and $\Lambda \in \mathcal{T}$. Then for some $n \in \mathbb{N}$, $A \in \mathcal{F}_n$ and $\Lambda \in \mathcal{T}_n$, so $A \perp\!\!\!\perp \Lambda$; thus \mathcal{P} and \mathcal{T} are independent. Since \mathcal{P} is a π -system, it follows from the Basic Criterion that $\sigma(\mathcal{P})$ and \mathcal{T} are also independent. But $\mathcal{F}_n \subset \mathcal{P}$ so each X_n is $\sigma(\mathcal{P})$ -measurable, hence $\mathcal{T} \subset \sigma(\mathcal{P})$ and each $\Lambda \in \mathcal{T}$ must also be in $\sigma(\mathcal{P}) \perp\!\!\!\perp \mathcal{T}$. It follows that:

$$P[\Lambda] = P[\Lambda \cap \Lambda] = P[\Lambda]P[\Lambda] = P[\Lambda]^2,$$

so $0 = P[\Lambda](1 - P[\Lambda])$ proving that $P[\Lambda]$ must be zero or one. \square

6.7 Product Spaces

Do independent random variables exist, with arbitrary (marginal) specified distributions? How can they be constructed? One way is to build *product probability spaces*; let's see how to do that.

Let $(\Omega_j, \mathcal{F}_j, P_j)$ be a probability space for $j = 1, 2$ and set:

$$\Omega = \Omega_1 \times \Omega_2 := \{(\omega_1, \omega_2) : \omega_j \in \Omega_j\}$$

$$\mathcal{F} = \mathcal{F}_1 \times \mathcal{F}_2 := \sigma\{A_1 \times A_2 : A_j \in \mathcal{F}_j\}$$

$$P := P_1 \otimes P_2, \text{ the unique extension to } \mathcal{F} \text{ satisfying:}$$

$$P(A_1 \times A_2) = P_1(A_1) \cdot P_2(A_2) \text{ for } A_1 \in \mathcal{F}_1, A_2 \in \mathcal{F}_2.$$

Any random variables X_1 on $(\Omega_1, \mathcal{F}_1, P_1)$ and X_2 on $(\Omega_2, \mathcal{F}_2, P_2)$ can be extended to the common space (Ω, \mathcal{F}, P) by defining $X_1^*(\omega_1, \omega_2) := X_1(\omega_1)$ and $X_2^*(\omega_1, \omega_2) := X_2(\omega_2)$; it's easy to show that $\{X_j^*\}$ are independent and have the same marginal distributions as $\{X_j\}$. Thus, independent random variables do exist with arbitrary distributions. The same construction extends to countable families.

6.8 Fubini's Theorem

We now consider how to evaluate probabilities and integrals on product spaces.

For any \mathcal{F} -measurable random variable $X : \Omega_1 \times \Omega_2 \rightarrow \mathcal{S}$ (\mathcal{S} would be \mathbb{R} , for real-valued RVs, but could also be \mathbb{R}^n or any complete separable metric space), and for any $\omega_2 \in \Omega_2$, the (second) *section* of X is the $(\Omega_1, \mathcal{F}_1, P_1)$ random variables $X_{\omega_2} : \Omega_1 \rightarrow \mathcal{S}$ defined by

$$X_{\omega_2}(\omega_1) := X(\omega_1, \omega_2).$$

It is not *quite* obvious, but true, that X_{ω_2} is \mathcal{F}_1 -measurable— show this first for indicator random variables $X = \mathbf{1}_{A_1 \times A_2}$ of product sets, then extend by a π - λ argument to indicators

$X = \mathbf{1}_A$ of sets $A \in \mathcal{F}$, then to simple RVs by linearity, then to the nonnegative RVs X_+ and X_- for an arbitrary \mathcal{F} -measurable X by monotone limits. Similarly, the first section $X_{\omega_1}(\cdot) := X(\omega_1, \cdot)$ is \mathcal{F}_2 -measurable for each $\omega_1 \in \Omega_1$.

Finally: Fubini's theorem gives conditions (namely, that either $X \geq 0$ or $E|X| < \infty$) to guarantee that these three integrals are meaningful and equal:

$$\int_{\Omega_2} \left\{ \int_{\Omega_1} X_{\omega_2} dP_1 \right\} dP_2 \stackrel{?}{=} \iint_{\Omega} X dP \stackrel{?}{=} \int_{\Omega_1} \left\{ \int_{\Omega_2} X_{\omega_1} dP_2 \right\} dP_1 \quad (1)$$

To prove this, first note that it's true for indicators $X = \mathbf{1}_{A_1 \times A_2}$ of the π -system of measurable rectangles ($A_1 \times A_2$) with each $A_j \in \mathcal{F}_j$; then verify that the class \mathcal{C} of events $A \in \mathcal{F}$ for which it holds for $X = \mathbf{1}_A$ is a λ -system. By Dynkin's π - λ theorem it follows that $\mathcal{F} \subset \mathcal{C}$ so (1) holds for all indicators $X = \mathbf{1}_A$ of events $A \in \mathcal{F}$, hence for all nonnegative simple functions in \mathcal{E}_+ , and finally for all \mathcal{F} -measurable $X \geq 0$ by the MCT. For $X \in L_1$, apply this result separately to X_+ and X_- .

Fubini's theorem applies more generally. Each probability measure P_j may be replaced by an arbitrary σ -finite¹ measure:

Theorem 2 (Fubini) *Let $(\mathcal{X}, \mathcal{F}, \mu)$ and $(\mathcal{Y}, \mathcal{G}, \nu)$ be two σ -finite measure spaces, and let $f(x, y)$ be a real-valued measurable function on the product space $(\mathcal{X} \times \mathcal{Y}, \mathcal{F} \times \mathcal{G}, \mu \otimes \nu)$. Then*

$$\int_{\mathcal{Y}} \left\{ \int_{\mathcal{X}} f(x, y) \mu(dx) \right\} \nu(dy) = \iint_{\mathcal{X} \times \mathcal{Y}} f(x, y) (\mu \otimes \nu)(dx dy) = \int_{\mathcal{X}} \left\{ \int_{\mathcal{Y}} f(x, y) \nu(dy) \right\} \mu(dx)$$

if either

- $f(x, y) \geq 0$ for $(x, y) \in \mathcal{N}^c$ for some $\mathcal{N} \in \mathcal{F} \times \mathcal{G}$ with $(\mu \otimes \nu)(\mathcal{N}) = 0$, or
- $f \in L_1(\mathcal{X} \times \mathcal{Y}, \mathcal{F} \times \mathcal{G}, \mu \otimes \nu)$.

Also, one of the measures (say, P_2) may be replaced by a measurable kernel² $K(\omega_1, d\omega_2)$ that is a σ -finite measure $K(\omega_1, \cdot)$ on \mathcal{F}_2 in its second variable for each fixed ω_1 , and an \mathcal{F}_1 -measurable function $K(\cdot, B_2)$ in its first variable for each fixed $B_2 \in \mathcal{F}_2$. Now Fubini's Theorem asserts (under positivity or L_1 conditions) the equality of integrals of X wrt the measure $P(d\omega_1 d\omega_2) = P_1(d\omega_1)K(\omega_1, d\omega_2)$ to the iterated integrals

$$\int_{\Omega_2} \nu_X(d\omega_2) = \iint_{\Omega} X dP = \int_{\Omega_1} \left\{ \int_{\Omega_2} X_{\omega_1}(\omega_2) K(\omega_1, d\omega_2) \right\} P_1(d\omega_1)$$

for the measure on \mathcal{F}_2 given by $\nu_X(d\omega_2) := \int_{\Omega_1} X_{\omega_2}(\omega_1) K(\omega_1, d\omega_2) P_1(d\omega_1)$.

¹Recall that a measure μ on a measurable space $(\mathcal{X}, \mathcal{F})$ is " σ -finite" if there are countably-many sets $\{\Lambda_j\} \subset \mathcal{F}$ with $\mu(\Lambda_j) < \infty$ for each j , and $\mathcal{X} = \cup_j \Lambda_j$. Evidently any finite measure (including probability measures) is also σ -finite, but the converse is false. Lebesgue measure is σ -finite on \mathbb{R}^n , for example.

²Measurable kernels come up all the time when studying conditional distributions (as you'll see in week 9 of this course) and, in particular, Markov chains and processes.

As an easy consequence of Theorem 2 (take $(\mathcal{X}, \mathcal{F}, \mu) = (\Omega, \mathcal{F}, \mathbb{P})$, let $(\mathcal{Y}, \mathcal{G}) = (\mathbb{N}, 2^{\mathbb{N}})$, and let $\nu(A) := \#\{A\}$ be counting measure for $A \subset \mathbb{N}$), for any sequence of random variables we may exchange summation and expectation and conclude that equality holds in

$$\mathbb{E} \left\{ \sum_{n=1}^{\infty} X_n \right\} \stackrel{?}{=} \sum_{n=1}^{\infty} \{\mathbb{E} X_n\}$$

whenever each $X_n \geq 0$ or when $\sum_{n=1}^{\infty} \mathbb{E}|X_n| < \infty$, but otherwise equality may fail.

6.8.1 A Counter-example

For an example where interchanging integration order fails, integrate by parts to verify:

$$\begin{aligned} \int_0^1 \left\{ \int_0^1 \frac{y^2 - x^2}{(x^2 + y^2)^2} dx \right\} dy &= \int_0^1 \left\{ \frac{-1}{1 + y^2} \right\} dy = \frac{-\pi}{4} \\ \int_0^1 \left\{ \int_0^1 \frac{y^2 - x^2}{(x^2 + y^2)^2} dy \right\} dx &= \int_0^1 \left\{ \frac{+1}{1 + x^2} \right\} dx = \frac{+\pi}{4}. \end{aligned}$$

As expected in light of Fubini's Theorem, the integrand isn't nonnegative nor is it in L_1 :

$$\begin{aligned} \iint_{[0,1]^2} \left| \frac{y^2 - x^2}{(x^2 + y^2)^2} \right| dx dy &\geq \int_0^{\pi/2} \int_0^1 \frac{r^2 |\sin^2 \theta - \cos^2 \theta|}{r^4} r dr d\theta \\ &= \left(\int_0^{\pi/2} |\cos(2\theta)| d\theta \right) \left(\int_0^1 r^{-1} dr \right) \\ &= (1)(\infty). \end{aligned}$$

6.8.2 A Simple but Useful Consequence of Fubini

For any $p > 0$ and any random variable X ,

$$\begin{aligned} \mathbb{E}|X|^p &= \mathbb{E} \left[\int_0^{|X|} p x^{p-1} dx \right] = \mathbb{E} \left[\int_0^{\infty} \mathbf{1}_{\{|X| > x\}} p x^{p-1} dx \right] \\ &= \int_0^{\infty} [\mathbb{E} \mathbf{1}_{\{|X| > x\}}] p x^{p-1} dx = \int_0^{\infty} p x^{p-1} \mathbb{P}[|X| > x] dx. \end{aligned}$$

It follows that

$$X \in L_p \Leftrightarrow \mathbb{E}|X|^p < \infty \Leftrightarrow \sum_{n=1}^{\infty} n^{p-1} \mathbb{P}[|X| > n] < \infty.$$

To see this, set $Y := |X|$ and fix $p \geq 1$. Then, since $\lceil \frac{n+1}{n} \rceil \leq 2$ for $n \in \mathbb{N}$, for $p \geq 1$:

$$\begin{aligned}
 \mathbf{E}Y^p &\leq \mathbf{E}(\lceil Y \rceil^p) = \int_0^\infty p y^{p-1} \mathbf{P}[\lceil Y \rceil > y] dy \\
 &= \sum_{n=0}^\infty \int_n^{n+1} p y^{p-1} \mathbf{P}[\lceil Y \rceil > y] dy \\
 &= \sum_{n=0}^\infty \int_n^{n+1} p y^{p-1} \mathbf{P}[Y > n] dy \\
 &\leq \sum_{n=0}^\infty p(n+1)^{p-1} \mathbf{P}[Y > n] \\
 &\leq p + \sum_{n=1}^\infty p 2^{p-1} n^{p-1} \mathbf{P}[Y > n]
 \end{aligned}$$

and hence $Y \in L_p$ if $\sum n^{p-1} \mathbf{P}[Y > n]$ converges. Conversely, if $0 \leq Y \in L_p$ for $p \geq 1$, then:

$$\begin{aligned}
 \mathbf{E}Y^p &\geq \mathbf{E}(\lfloor Y \rfloor^p) = \int_0^\infty p y^{p-1} \mathbf{P}[\lfloor Y \rfloor > y] dy \\
 &= \sum_{n=0}^\infty \int_n^{n+1} p y^{p-1} \mathbf{P}[Y > n] dy \\
 &\geq \sum_{n=0}^\infty \int_n^{n+1} p n^{p-1} \mathbf{P}[Y > n] dy \\
 &= \sum_{n=1}^\infty p n^{p-1} \mathbf{P}[Y > n],
 \end{aligned}$$

so the sum converges if $Y \in L_p$. The argument for $0 < p < 1$ is very similar, but differs slightly because now y^{p-1} is decreasing on each interval $(n, n+1]$ instead of increasing.

The case $p = 1$ is easiest and most important: if $S := \sum_{n=0}^\infty \mathbf{P}[|X| > n] < \infty$, then $X \in L_1$ with $\mathbf{E}|X| \leq S \leq \mathbf{E}|X| + 1$. If X takes on only nonnegative integer values then $\mathbf{E}X = S$. For any $\epsilon > 0$, apply this to $Y := X/\epsilon$ to see

$$\epsilon \sum_{n=1}^\infty \mathbf{P}[|X| > n\epsilon] < \mathbf{E}|X| \leq \epsilon \sum_{n=0}^\infty \mathbf{P}[|X| > n\epsilon]$$

6.9 Hoeffding's Inequality

If $\{X_j\}$ are independent and (individually) bounded, so $(\forall j \in \mathbb{N}) (\exists \{a_j, b_j\})$ for which $P[a_j \leq X_j \leq b_j] = 1$, then $(\forall c > 0)$, $S_n := \sum_{j=1}^n X_j$ satisfies

$$P[(S_n - ES_n) \geq c] \leq \exp\left(-2c^2 / \sum_{j=1}^n |b_j - a_j|^2\right).$$

If X_j are iid and bounded by $\|X_j\|_\infty \leq 1$, e.g., then take $a_j = -1$, $b_j = 1$, and $c = n\epsilon$ to see

$$P[(\bar{X}_n - \mu) \geq \epsilon] \leq e^{-n\epsilon^2/2}. \quad (2)$$

Wassily Hoeffding proved this improvement on Chebychev's inequality for L_∞ random variables in 1963 at UNC. It follows from Hoeffding's Lemma:

$$E[e^{\lambda(X_j - \mu_j)}] \leq \exp(\lambda^2(b_j - a_j)^2/8),$$

proved in turn from Jensen's ineq and Taylor's theorem (with remainder). The importance of (2) is that the bound decreases *exponentially* in n as $n \rightarrow \infty$, while the Chebychev bound $P[|\bar{X}_n - \mu| \geq \epsilon] \leq \sigma^2/n\epsilon^2$ decreases only like $1/n$. The price for this better bound is that the $\{X_j\}$ must be bounded in L_∞ , not merely in L_2 . See also related and earlier **Bernstein's** inequality (1937), **Chernoff** bounds (1952), and **Azuma's** inequality (1967).

Here's a proof for the important special case of $X_j = \pm 1$ with probability $1/2$ each (and hence $\mu = 0$):

$$\begin{aligned} P[\bar{X}_n \geq \epsilon] &= P[S_n \geq n\epsilon] \\ &= P[e^{\lambda S_n} \geq e^{n\lambda\epsilon}] && \text{for any } \lambda > 0 \\ &\leq E[e^{\lambda S_n}] e^{-n\lambda\epsilon} && \text{by Markov's inequality} \\ &= \left\{\frac{1}{2}e^\lambda + \frac{1}{2}e^{-\lambda}\right\}^n e^{-n\lambda\epsilon} && \text{by independence} \\ &\leq \left\{e^{\lambda^2/2}\right\}^n e^{-n\lambda\epsilon} && \text{see footnote}^3 \\ &= \exp(n\lambda^2/2 - n\lambda\epsilon) \end{aligned}$$

The exponent is minimized at $\lambda = \epsilon$, so the tightest bound is:

$$P[\bar{X}_n \geq \epsilon] \leq \exp(n\epsilon^2/2 - n\epsilon\epsilon) = e^{-n\epsilon^2/2}.$$

The general case isn't much harder, but proving $Ee^{\lambda(X-\mu)} \leq e^{\lambda^2/2}$ is a bit more delicate.

By Borel/Cantelli it follows from Hoeffding's inequality that $(\bar{X}_n - \mu) > \epsilon$ only finitely-many times for each $\epsilon > 0$, if $\{X_n\} \subset L_\infty$ are iid, leading to our first **Strong Law of Large Numbers**: $P[\bar{X}_n \rightarrow \mu] = 1$ (why does this follow?).

³ $\cosh(\lambda) = \{\frac{1}{2}e^\lambda + \frac{1}{2}e^{-\lambda}\} = \sum \frac{\lambda^{2k}}{(2k)!} \leq \sum \frac{(\lambda^2)^k}{2^k(k!)} = e^{\lambda^2/2}$

Note that Chebychev's inequality only guarantees the algebraic bound $P[\bar{X}_n \geq \epsilon] \leq 1/n\epsilon^2$, instead of Hoeffding's exponential bound. Since $1/n\epsilon^2$ isn't summable in n , Chebychev's bound isn't strong enough to prove a strong LLN, but Hoeffding's is.

Hoeffding's inequality is now used commonly in computer science and machine learning, applied to indicators of Bernoulli events (or, equivalently, to binomial random variables). It gives the bound

$$P[|\bar{X}_n - p| \leq \epsilon] = P[(p - \epsilon)n \leq S_n \leq (p + \epsilon)n] \geq 1 - 2e^{-2\epsilon^2 n}$$

for $S_n := \sum_{j \leq n} X_j \sim \text{Bi}(n, p)$ for iid Bernoulli $X_j \stackrel{\text{iid}}{\sim} \text{Bi}(1, p)$ variables, showing exponential concentration of probability around the mean. This is far stronger than Chebychev's bound of $P[|\bar{X}_n - p| \leq \epsilon] \geq 1 - p(1 - p)/\epsilon^2 n$, since Hoeffding's bound for $P[|\bar{X}_n - p| > \epsilon]$ is *exponentially* small as $n \rightarrow \infty$ and hence is summable.