

# Conditional Expectation

Robert L. Wolpert  
Department of Statistical Science  
Duke University, Durham, NC, USA

## 10 Conditioning

Frequently in probability and (especially Bayesian) statistics we wish to find the probability of some event  $A$  or the expectation of some random variable  $X$ , *conditionally* on some body of information— such as the occurrence of another event  $B$  or the value of another random variable  $Z$  (or collection of them  $\{Z_\alpha\}$ ). In elementary probability we encounter the usual formulas for conditional probabilities and expectations

$$\mathbb{P}[A | B] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B]} \quad \mathbb{E}[X | Z] = \begin{cases} \frac{\int x f(x,Z) dx}{\int f(x,Z) dx} & X, Z \text{ jointly continuous} \\ \frac{\sum x f(x,Z)}{\sum f(x,Z)} & X, Z \text{ discrete} \end{cases}$$

but this notion breaks down either for distributions which are *not* jointly absolutely continuous or discrete, and also when we wish to condition on the value of infinitely-many (even uncountably-many) random variables  $\{Z_\alpha\}$ , as we will when we consider stochastic processes. There simply is no such thing as a joint density function for an infinite collection  $\{Z_\alpha\}$ , even if each finite set has an absolutely continuous joint distribution.

Since “information” in probability theory is represented by  $\sigma$ -algebras (here  $\sigma\{B\}$  or  $\sigma\{Z_\alpha\}$ ), what we need are ways to express, interpret, and compute *conditional* probabilities of events and *conditional* expectations of random variables, given  $\sigma$ -algebras. As a bonus, this will unify the notions of conditional probability and conditional expectation, for distributions that are discrete or continuous or neither. First, a tool to help us.

### 10.1 Lebesgue’s Decomposition

Let  $\mu$  and  $\lambda$  be two positive  $\sigma$ -finite measures on the same measurable space  $(\Omega, \mathcal{F})$ . Call  $\mu$  and  $\lambda$  *equivalent*, and write  $\mu \equiv \lambda$ , if they have the same null sets— so the notion of “a.e.” is the same for both. More generally, we call  $\lambda$  *absolutely continuous* (AC) w.r.t.  $\mu$ , and write  $\lambda \ll \mu$ , if  $\mu(A) = 0$  implies  $\lambda(A) = 0$ , *i.e.*, if every  $\mu$ -null set is also  $\lambda$ -null (so  $\lambda \equiv \mu$  if and only if  $\lambda \ll \mu$  and  $\mu \ll \lambda$ ). We call  $\mu$  and  $\lambda$  *mutually singular*, and write  $\mu \perp \lambda$ , if for some

disjoint sets  $A, B \in \mathcal{F}$  we have  $\mu(A^c) = 0$  and  $\lambda(B^c) = 0$ , so  $\mu$  and  $\lambda$  are “concentrated” on disjoint sets.

For example— if  $\lambda(A) = \int_A f(x)\mu(dx)$  for some  $\mathcal{F}$ -measurable function  $f : \Omega \rightarrow \mathbb{R}_+$ , then  $\lambda \ll \mu$ ; if  $f > 0$   $\mu$ -a.s, then also  $\mu(A) = \int_A f(x)^{-1}\lambda(dx)$  and  $\mu \equiv \lambda$ . If for some other  $\sigma$ -finite measure  $\nu$  and some  $\mathcal{F}$ -measurable  $f, g : \Omega \rightarrow \mathbb{R}_+$  we set

$$\mu(A) := \int_A f(x)\nu(dx) \quad \lambda(A) := \int_A g(x)\nu(dx)$$

then  $\mu \perp \lambda$  if  $f(x)g(x) = 0$  for  $\nu$ -a.e.  $x \in \Omega$ . The functions  $f$  and  $g$  are called the densities of  $\mu$  and  $\lambda$  with respect to  $\nu$ , generalizing the familiar idea of density functions w.r.t. Lebesgue measure.

**Theorem 1 (Lebesgue Decomposition)** *Let  $\mu, \lambda$  be two  $\sigma$ -finite measures on a measurable space  $(\Omega, \mathcal{G})$ . Then there exist a unique pair  $\lambda_a, \lambda_s$  of  $\sigma$ -finite measures on  $(\Omega, \mathcal{G})$  and a unique  $\mathcal{G}$ -measurable function  $Y$  such that:*

$$\begin{aligned} \lambda &= \lambda_a + \lambda_s \\ \lambda_a &\ll \mu, \quad \lambda_s \perp \mu \\ \lambda_a(G) &= \int_G Y(\omega)\mu(d\omega), \quad G \in \mathcal{G}. \end{aligned}$$

**Proof Sketch.** Set

$$\mathcal{H} := \{h \in L_1(\Omega, \mathcal{G}, \mu) : h \geq 0, (\forall G \in \mathcal{G}) \int_G h d\mu \leq \lambda(G)\}$$

Show that  $\mathcal{H}$  is closed under maxima, then find  $\{h_n\}$  such that

$$\sup \left\{ \int h_n d\mu : n \in \mathbb{N} \right\} = \sup \left\{ \int h d\mu : h \in \mathcal{H} \right\}$$

and set  $h := \sup h_n$  and  $Y := h\mathbf{1}_{\{h < \infty\}}$ . Now verify the statement of the Theorem.  $\square$

If  $\mu(dx) = dx$  is Lebesgue measure on  $\mathbb{R}^d$ , for example, then this decomposes any probability distribution  $\lambda$  into an absolutely continuous part  $\lambda_a(dx) = Y(x) dx$  with pdf  $Y$  and a singular part  $\lambda_s(dx)$  (the sum of the singular-continuous and discrete components). When  $\lambda \ll \mu$  (so  $\lambda_a = \lambda$  and  $\lambda_s = 0$ ) the Radon-Nikodym derivative is often denoted  $Y = \frac{d\lambda}{d\mu}$  or  $\frac{\lambda(d\omega)}{\mu(d\omega)}$ , and extends the idea of “density” from densities with respect to Lebesgue measure to those with respect to an arbitrary “reference” (or “base” or “dominating”) measure  $\mu$ . For example, the pmf  $f(x) = \mathbf{P}[X = x]$  of an integer-valued random variable  $X$  may now be viewed as its pdf with respect to counting measure on  $\mathbb{Z}$ , so families of discrete distributions now have pdf’s (if they take values in a common countable set), and random variables with mixed distributions (truncated normals, for example) have density functions with respect to

a dominating measure that includes point masses where the distributions have atoms, and Lebesgue measure where they are absolutely continuous. With respect to the finite base measure  $\lambda(A) := \sum\{1/k! : k \in A\}$  on the nonnegative integers  $\mathbb{N}_0$ , for example, the  $\text{Po}(\lambda)$  distribution has pdf  $f(k) = \lambda^k e^{-\lambda}$ .

To explore further conditioning we apply Lebesgue's decomposition in a quite different way, with  $\mu = \mathbb{P}$  a probability measure on  $(\Omega, \mathcal{F}, \mathbb{P})$  and  $\lambda(d\omega) = X d\mathbb{P}$  for some  $X \in L_1$  a  $\sigma$ -finite measure to prove the important:

## 10.2 The Radon-Nikodym Theorem

**Theorem 2 (Radon-Nikodym)** *Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space,  $X \in L_1(\Omega, \mathcal{F}, \mathbb{P})$ , and  $\mathcal{G} \subset \mathcal{F}$  a sub- $\sigma$ -algebra. Then there exists a unique  $Y \in L_1(\Omega, \mathcal{G}, \mathbb{P})$ , which we will denote  $Y = \mathbb{E}[X | \mathcal{G}]$  and call a “conditional expectation of  $X$ , given  $\mathcal{G}$ ,” that satisfies for every  $G \in \mathcal{G}$ :*

$$(\forall G \in \mathcal{G}) \quad \mathbb{E} (X - Y)\mathbf{1}_G = 0$$

The important feature to notice is that  $Y$  must be  $\mathcal{G}$ -measurable, which may be hard to achieve if  $\mathcal{G}$  is much smaller than  $\mathcal{F}$ . In some sense  $Y$  is the best possible  $\mathcal{G}$ -measurable approximation to  $X$ .

**Proof.** First take  $X$  to be non-negative,  $X \geq 0$ . The measure  $\mathbb{P}$ , initially defined on all of  $\mathcal{F}$ , can also be viewed as a probability measure on the smaller  $\sigma$ -algebra  $\mathcal{G} \subset \mathcal{F}$ . Define another measure  $\lambda$  on  $\mathcal{G}$  (not on all of  $\mathcal{F}$ ) by

$$\lambda(G) := \mathbb{E} X \mathbf{1}_G = \int_G X(\omega) \mathbb{P}(d\omega), \quad G \in \mathcal{G}.$$

This is bounded (since  $X \in L_1(\Omega, \mathcal{F}, \mathbb{P})$ ) and positive (since  $X \geq 0$ ), so by Theorem 1 (applied on  $(\Omega, \mathcal{G}, \mathbb{P})$ , not  $(\Omega, \mathcal{F}, \mathbb{P})$ ) we can write  $\lambda = \lambda_a + \lambda_s$  with  $\lambda_a \ll \mathbb{P}$ ,  $\lambda_s \perp \mathbb{P}$ , and  $\lambda_a(G) = \int_G Y d\mathbb{P}$  for some  $Y \in L_1(\Omega, \mathcal{G}, \mathbb{P})$ . But  $\lambda \ll \mathbb{P}$  by construction, so (by uniqueness)  $\lambda_s = 0$ ,  $\lambda_a = \lambda$ , and the Theorem follows.

For general  $X$ , consider separately the positive and negative parts  $X_+ := \max(X, 0)$  and  $X_- := \max(-X, 0)$  and set  $Y := Y_+ - Y_-$ .  $\square$

For events  $A \in \mathcal{F}$  and sub- $\sigma$ -algebras  $\mathcal{G} \subseteq \mathcal{F}$  we denote the *conditional probability of  $A$ , given  $\mathcal{G}$*  by

$$\mathbb{P}[A | \mathcal{G}] = \mathbb{E} [\mathbf{1}_A | \mathcal{G}],$$

a  $\mathcal{G}$ -measurable random variable (not a numerical constant) taking values in  $[0, 1]$ .

Of course  $X$  itself has the property that its integrals over events  $G \in \mathcal{G}$  coincide with those of  $X$ —the point is that  $Y = \mathbb{E}[X | \mathcal{G}]$  is a  $\mathcal{G}$ -measurable approximation to  $X$  (i.e., one that depends only on the “information” encoded in  $\mathcal{G}$ ) with this property. As we'll see below, if  $\mathcal{F} \subseteq \mathcal{G}$  (or, more generally, if  $X$  is  $\mathcal{G}$ -measurable, so  $\sigma(X) \subseteq \mathcal{G}$ ) then the best  $\mathcal{G}$ -measurable approximation is  $\mathbb{E}[X | \mathcal{G}] = X$  itself. At the other extreme, if  $X$  is independent of  $\mathcal{G}$ , then one can do no better than the constant random variable  $\mathbb{E}[X | \mathcal{G}] \equiv \mathbb{E}X$ .

### 10.2.1 Key Example: Countable Partitions

If  $\mathcal{G} = \sigma\{\Lambda_n\}$  for a finite or countable partition  $\{\Lambda_n\} \subset \mathcal{F}$  (so  $\Lambda_m \cap \Lambda_n = \emptyset$  for  $m \neq n$  and  $\Omega = \cup \Lambda_n$ ), then for any  $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ ,

$$\mathbf{E}[X | \mathcal{G}] = \sum \mathbf{1}_{\Lambda_n} \mathbf{E}_{\Lambda_n}[X] = \sum \mathbf{1}_{\{\Lambda_n\}}(\omega) \frac{1}{\mathbf{P}[\Lambda_n]} \mathbf{E}[X \mathbf{1}_{\{\Lambda_n\}}]$$

is constant on partition elements and equal there to the  $\mathbf{P}$ -weighted average value of  $X$  (omit from the sum any term with  $\mathbf{P}[\Lambda_n] = 0$ ).

In particular— let  $(\Omega, \mathcal{F}, \mathbf{P})$  be the unit interval with Lebesgue measure, and let  $\mathcal{G}_n = \sigma\{(i/2^n, j/2^n]\}$ ,  $0 \leq i < j \leq 2^n$ . Note that  $\mathcal{G}_n \subset \mathcal{G}_m$  for  $n \leq m$  and that  $\mathcal{F} = \bigvee \mathcal{G}_n$ . Then for any  $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ ,

$$X_n := \mathbf{E}[X | \mathcal{G}_n] = 2^n \int_{i/2^n}^{(i+1)/2^n} X(v) dv, \quad i/2^n < \omega \leq (i+1)/2^n, \quad 0 \leq i < 2^n.$$

This is our first example of a *martingale*, a sequence of random variables  $X_n \in L_1(\Omega, \mathcal{F}, \mathbf{P})$  with the property that  $X_n = \mathbf{E}[X_m | \mathcal{G}_n]$  for  $n \leq m$ ; we'll see more soon. What happens as  $n \rightarrow \infty$ ?

### 10.2.2 Properties:

- The conditional expectation is *almost* unique: if  $Y_1$  and  $Y_2$  are each  $\mathcal{G}$ -measurable and for some  $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$  and all  $G \in \mathcal{G}$  satisfy

$$\mathbf{E}(X - Y_1)\mathbf{1}_G = 0 = \mathbf{E}(X - Y_2)\mathbf{1}_G,$$

then each may be called “ $\mathbf{E}[X | \mathcal{G}]$ ” but they may not be equal for all  $\omega \in \Omega$ . The difference  $(Y_1 - Y_2)$  is  $\mathcal{G}$  measurable and is zero almost-surely (can you prove that?), but still may not vanish for all  $\omega \in \Omega$ . Thus one speaks of “a” conditional expectation rather than “the” conditional expectation.

- If  $X = \mathbf{1}_A$  and if  $\mathcal{G} = \sigma\{B\}$  for some  $A, B \in \mathcal{F}$  with  $0 < \mathbf{P}[B] < 1$ ,

$$\mathbf{P}[A | \mathcal{G}](\omega) = \mathbf{E}[\mathbf{1}_A | \sigma(B)](\omega) = \begin{cases} \mathbf{P}[A \cap B] / \mathbf{P}[B] & \omega \in B \\ \mathbf{P}[A \cap B^c] / \mathbf{P}[B^c] & \omega \notin B \end{cases}$$

Thus, conditional expectation (given a  $\sigma$ -algebra  $\mathcal{G}$ ) generalizes the notion of the conditional probability of one event  $A$  given another  $B$  (or its complement  $B^c$ ).

- More generally, If  $X \in L_1$  and if  $\mathcal{G} = \sigma\{G_i\}$  for some (finite or countable) measurable partition  $\{G_i\} \subset \mathcal{F}$ , then

$$\mathbf{E}[X | \mathcal{G}](\omega) = \sum \mathbf{1}_{G_i}(\omega) \frac{1}{\mathbf{P}(G_i)} \int_{G_i} X(\omega') P(d\omega')$$

is the weighted average of  $X$  over the partition element that contains  $\omega$ .

- If  $X, Z \sim f(x, z)$  are jointly absolutely-continuous and if  $\mathcal{G} = \sigma(Z)$ ,

$$\mathbb{E}[X | Z] := \mathbb{E}[X | \sigma(Z)] = \frac{\int x f(x, Z) dx}{\int f(x, Z) dx}.$$

Thus, conditional expectation (given a  $\sigma$ -algebra  $\mathcal{G}$ ) generalizes the elementary notion of conventional expectation (given an RV  $Z$ ). What if  $X$  and  $Z$  are both discrete? What if just one is discrete? What if  $Z$  is a vector?

To prove this property, first show that a random variable is  $\mathcal{G}$  measurable if and only if it is a Borel function of  $Z$  (obvious for simple RVs, then take monotone limits). Apply this to write  $\mathbb{E}[X | \mathcal{G}] = \phi(Z)$ ; then for  $G \in \mathcal{G}$ , solve the equation  $0 = \mathbb{E}\mathbf{1}_G[\phi(Z) - X]$  for  $\phi(Z)$ .

- If  $X \in L_1(\Omega, \mathcal{F}, \mathbb{P})$  and if  $X \perp\!\!\!\perp \mathcal{G}$  then

$$\mathbb{E}[X | \mathcal{G}] \equiv \mathbb{E}X.$$

In particular,  $\mathbb{E}[X | \{\Omega, \emptyset\}] = \mathbb{E}X$ . Thus, conditional expectation (given a  $\sigma$ -algebra  $\mathcal{G}$ ) generalizes the elementary notion of expectation.

- If  $X \in L_1(\Omega, \mathcal{F}, \mathbb{P})$  and if  $\mathcal{H} \subset \mathcal{G} \subset \mathcal{F}$ , then

$$\mathbb{E}[X | \mathcal{H}] = \mathbb{E}[\mathbb{E}[X | \mathcal{G}] | \mathcal{H}]$$

This is called the “tower” (or sometimes “smoothing” or “telescoping”) property of conditional expectation. It’s especially useful when we have entire nested families (called *filtrations*) of  $\sigma$ -algebras  $\{\mathcal{F}_n\}$  with  $n \leq m \Rightarrow \mathcal{F}_n \subseteq \mathcal{F}_m$ ; for example,  $\mathcal{F}_n := \sigma\{X_j : j \leq n\}$  for a family  $\{X_n\}$  of (non-necessarily-independent) random variables.

- A common use of the tower property is the calculation for  $\mathcal{G}$ -measurable  $X$  with  $Y$ ,  $XY \in L_1(\Omega, \mathcal{F}, \mathbb{P})$ ,

$$\mathbb{E}[XY] = \mathbb{E}[\mathbb{E}[XY | \mathcal{G}]] = \mathbb{E}[X \mathbb{E}[Y | \mathcal{G}]]$$

*i.e.*,  $\mathcal{G}$ -measurable random variables can be pulled out of conditional expectations just like constants. The case  $\mathcal{G} = \sigma(X)$  is most common:  $\mathbb{E}[XY] = \mathbb{E}[X \mathbb{E}[Y | X]]$ .

- If  $X$  and  $\{Y_n\}$  are jointly Gaussian, then  $\mathbb{E}[X | \sigma\{Y_n\}]$  is the orthogonal projection of  $X$  onto the linear span of  $\{Y_n\}$  in the Hilbert space  $L_2(\Omega, \mathcal{F}, \mathbb{P})$ . This is usually the easiest way to compute conditional expectations in multivariate normal examples. Thus, conditional expectation (given a  $\sigma$ -algebra  $\mathcal{G}$ ) generalizes the notion of orthogonal projection. Similarly,
- $L_2$  prediction: For  $X \in L_2$ ,  $\mathbb{E}[X | \mathcal{G}]$  minimizes  $\|X - Y\|_2^2$  among all  $\mathcal{G}$ -measurable  $Y$ .

- Let  $\{X_n\} \subset L_1(\Omega, \mathcal{F}, \mathbb{P})$  be iid with means  $\mu = \mathbb{E}[X_n]$  and set  $S_n := \sum_{j \leq n} X_j$  and  $\mathcal{G}_n := \sigma\{X_1, \dots, X_n\}$ . Then for  $n \leq m$ ,

$$\mathbb{E}[S_m | \mathcal{G}_n] = S_n + (m - n)\mu;$$

in particular,  $(S_n - n\mu)$  is a *martingale*. If  $\{X_n\} \subset L_2(\Omega, \mathcal{F}, \mathbb{P})$ , set  $\sigma^2 := \mathbb{V}X_n$  and check that  $(S_n - n\mu)^2 - n\sigma^2$  is also a martingale.

- **Monotonicity:** If  $X \geq Z$  *a.s.*, then  $\mathbb{E}[X | \mathcal{G}] \geq \mathbb{E}[Z | \mathcal{G}]$  *a.s.* for any  $\mathcal{G} \subset \mathcal{F}$ . To see this, note first it suffices to take  $Z \equiv 0$ ; then set  $Y := \mathbb{E}[X | \mathcal{G}]$  and  $G = \{\omega : Y < 0\}$ . Since  $G \in \mathcal{G}$ , note  $\mathbb{E}[Y \mathbf{1}_{\{Y < 0\}}] = \mathbb{E}[X \mathbf{1}_{\{Y < 0\}}] \geq 0$  since  $X \geq 0$  *a.s.*, so  $\mathbb{P}[Y < 0] = 0$ .
- **Conditional Mean/Variance Formula:** If  $X \in L_2(\Omega, \mathcal{F}, \mathbb{P})$  and  $Y := \mathbb{E}[X | \mathcal{G}]$ ,

$$\mathbb{V}[X] = \mathbb{E}\left\{\mathbb{E}[(X - Y)^2 | \mathcal{G}]\right\} + \mathbb{V}(Y).$$

Thus the variance of  $X$  is the mean of the conditional variance plus the variance of the conditional mean. This elegant formula is worth remembering.

- All the usual integration tools and inequalities— DCT, MCT, Fatou, Jensen, Hölder, Minkowski, Markov, Chebychev, *etc.*— have *conditional* versions as well. For example, for  $X \in L_1$  and convex  $\phi(\cdot)$  with  $\phi(X) \in L_1$ ,

$$\phi(\mathbb{E}[X | \mathcal{G}]) \leq \mathbb{E}[\phi(X) | \mathcal{G}] \text{ a.s.}$$

(note both sides are  $\mathcal{G}$ -measurable *random variables* now, not constants as in the familiar Jensen inequality, so the “almost surely” qualification is needed).

If  $0 \leq X_n \uparrow X$  in probability, for another example, then

$$\mathbb{E}[X_n | \mathcal{G}] \rightarrow \mathbb{E}[X | \mathcal{G}] \text{ a.s.},$$

and also  $\mathbb{E}[|X_n - X| | \mathcal{G}] \rightarrow 0$  *a.s.*, a conditional generalization of Lebesgue’s MCT.

This MCT can be used to prove a conditional Fatou’s Lemma is available for  $X_n \geq 0$ :

$$\mathbb{E}\left[\liminf_{n \rightarrow \infty} X_n | \mathcal{G}\right] \leq \liminf_{n \rightarrow \infty} \mathbb{E}[X_n | \mathcal{G}] \text{ a.s.}$$

If  $X_n \rightarrow X$  *pr.* and  $|X_n| \leq Y \in L_1$  *a.s.*, then again

$$\mathbb{E}[X_n | \mathcal{G}] \rightarrow \mathbb{E}[X | \mathcal{G}] \text{ a.s.},$$

a conditional version of Lebesgue’s DCT. To prove this, just apply the conditional Fatou’s lemma to the nonnegative RVs  $Y + X_n$  and  $Y - X_n$  to see

$$\begin{aligned} \mathbb{E}[Y + X | \mathcal{G}] &= \mathbb{E}\left[\liminf_{n \rightarrow \infty} (Y + X_n) | \mathcal{G}\right] \leq \liminf_{n \rightarrow \infty} \mathbb{E}[Y + X_n | \mathcal{G}] \text{ and} \\ \mathbb{E}[Y - X | \mathcal{G}] &= \mathbb{E}\left[\liminf_{n \rightarrow \infty} (Y - X_n) | \mathcal{G}\right] \leq \liminf_{n \rightarrow \infty} \mathbb{E}[Y - X_n | \mathcal{G}]. \end{aligned}$$

Upon subtracting  $\mathbb{E}[Y | \mathcal{G}]$  and changing signs for the second equation, we conclude

$$\mathbb{E}[X | \mathcal{G}] \leq \liminf_{n \rightarrow \infty} \mathbb{E}[X_n | \mathcal{G}] \text{ and } \mathbb{E}[X | \mathcal{G}] \geq \limsup_{n \rightarrow \infty} \mathbb{E}[X_n | \mathcal{G}].$$

### 10.3 Borel's Paradox

Let  $(X, Y)$  be the longitude,  $0 \leq X < 2\pi$ , and latitude,  $-\pi/2 \leq Y \leq \pi/2$ , of a point drawn uniformly from a sphere  $\mathcal{S}$  (perhaps the globe). What is its *conditional* distribution of  $(X, Y)$ , given that it lies on a great circle  $\mathcal{C}$ ? This famously ill-posed question helps motivate a careful consideration of conditioning. If the “great circle” is the equator  $Y = 0$ , the answer is the (perhaps expected) uniform distribution, with longitude  $X \sim \text{Un}([0, 2\pi))$ . But if the great circle is, say, the prime meridian  $X = 0$ , then the point is much more likely to be near the equator (where an interval of  $Y = 0 \pm 1$  degree latitude has a large area) than near either pole (where it doesn't); in that case the conditional distribution of  $Y$  has density  $f(y | x) = \frac{1}{2} \cos(y) \mathbf{1}_{[-\pi/2, \pi/2]}(y)$  for any  $0 \leq x < 2\pi$ .

We simply cannot condition meaningfully on the null event that  $(X, Y)$  lies on a set of zero probability, such as a great circle. We *can* condition on events of positive probability, or on the  $\sigma$ -algebra generated by a random variable.

In *Radon spaces* (which include  $\mathbb{R}^d$  and all complete separable metric spaces) these notions are closely related: in particular, we can always compute a version of the conditional expectation of one random-variable  $X$  given another  $Z$  as  $\mathbf{E}[X | Z] = \phi_X(z)$  for the limit

$$\phi_X(z) = \limsup_{\epsilon \rightarrow 0} \mathbf{E}[X | \{|Z - z| < \epsilon\}].$$

Let's use this to try to answer the question: *What is the conditional distribution of the horizontal component  $X$  of a point drawn from the unit square, given that the point lies on the bottom edge?* Let  $(X, Y)$  be the coordinates of a point drawn uniformly from the unit square and  $0 < \epsilon < 1$ , and let  $\Delta$  denote the bottom edge of the square. For  $0 < x < 1$  we can compute

$$\mathbf{P}[X \leq x | 0 \leq Y \leq \epsilon] = \frac{\epsilon x}{\epsilon} = x$$

and conclude (taking  $\epsilon \rightarrow 0$ ) that the conditional *distribution* of  $X$ , given  $Y = 0$ , is the standard uniform, and hence the conditional expectation  $\mathbf{E}[X | Y = 0] = 1/2$ . Similarly if we let  $R = Y/X$  be the ratio of  $Y$  to  $X$ , we can also compute

$$\mathbf{P}[X \leq x | 0 \leq R \leq \epsilon] = \frac{\epsilon x^2/2}{\epsilon/2} = x^2,$$

so the conditional distribution of  $X$ , given  $R = 0$ , is  $\text{Be}(2, 1)$ , with conditional density  $f(x | R=0) = 2x$  on  $[0, 1]$  and conditional mean  $\mathbf{E}[X | R=0] = 2/3$ . Note that both of these “events” on which we condition are identical—the null event that  $(X, Y)$  lies on the bottom edge  $\Delta$  of the square, another example of Borel's paradox. Really these two different results were answers to different questions: one found the values of  $\mathbf{P}[X \leq x | \sigma\{Y\}]$  and  $\mathbf{E}[X | \sigma\{Y\}]$ , the other found  $\mathbf{P}[X \leq x | \sigma\{R\}]$  and  $\mathbf{E}[X | \sigma\{R\}]$ . Geometrically, what do events in  $\sigma\{Y\}$  and those in  $\sigma\{R\}$  look like in the square? For an arbitrary density  $f(x)$  on the unit interval, can you find a random variable  $Z$  (a function of  $X$  and  $Y$ ) such that  $\{Z = 0\}$  is the bottom edge of the square and the conditional distribution of  $X$  given  $Z = 0$  is  $f(x) dx$ ? Are any conditions on  $f(x)$  needed?

**A little more generally...**

Let  $f(x)$  be any strictly-positive bounded pdf on the unit interval, and set  $Z := Y/f(X)$ . Then  $(X, Y) \in \Delta = [0, 1] \times \{0\}$  if and only if  $Z = 0$  and, for  $\epsilon > 0$ ,

$$\begin{aligned} \mathbb{P}[X \leq x \mid Z \leq \epsilon] &= \frac{\mathbb{P}[X \leq x \cap Y \leq \epsilon f(X)]}{\mathbb{P}[Y \leq \epsilon f(X)]} \\ &= \frac{\int_0^x \min[1, \epsilon f(\xi)] d\xi}{\int_0^1 \min[1, \epsilon f(\xi)] d\xi} = \frac{\int_0^x \min[1/\epsilon, f(\xi)] d\xi}{\int_0^1 \min[1/\epsilon, f(\xi)] d\xi} \\ &\rightarrow \int_0^x f(\xi) d\xi \end{aligned}$$

as  $\epsilon \rightarrow 0$  by LMCT, so the limiting distribution of  $X$  conditional on  $Z \leq \epsilon$  is the completely arbitrary distribution with density  $f(x)$ . Thus, in a very strong way, the “conditional distribution of  $X$  given that  $(X, Y) \in \Delta$ ” is not determined. We can find conditional probabilities and distributions given *random variables* or non-null *events* or (more generally than either) *sigma algebras*, but not given events of probability zero.

**Be careful out there...**

Borel’s paradox isn’t just an academic puzzle. Naïve attempts to “condition” on null events (for example, by trying to impose Bayesian prior distributions on both the inputs and outputs of deterministic models, as in *Inference from a Deterministic Population Dynamics Model for Bowhead Whales* by Raftery, Givens & Zeh, JASA 1995) pop up every year or two in the literature, and sometimes aren’t caught in the review process. That one (I kid you not) led to discussions about Borel’s Paradox at meetings of the International Whaling Commission, and in the 1995 IWC Annual Report (try googling “bowhead whale borel paradox”).

Be careful!