

Contents

7	Conditioning	61
7a	What is the problem	61
7b	Discrete case	63
7c	Conditional expectation	64
7d	A convergence theorem	66
7e	Conditional measures	66
7f	Markov property as conditional independence	74

7 Conditioning

7a What is the problem

First, three results announced. They are rather trivial in such special cases as discrete distributions or absolutely continuous distributions, but nontrivial for singular (and mixed) distributions. We'll see that the general theory holds for *all* distributions, and generalizes to high and even infinite dimension.

7a1 Proposition. Every 2-dimensional random variable (X, Y) is distributed like $(f(U), g(U, V))$ for some Lebesgue measurable functions $f : (0, 1) \rightarrow \mathbb{R}$ and $g : (0, 1) \times (0, 1) \rightarrow \mathbb{R}$ such that f is increasing on $(0, 1)$, and $g(u, \cdot)$ is increasing on $(0, 1)$ for each $u \in (0, 1)$; here U, V are independent random variables distributed uniformly on $(0, 1)$.

Such f is called the quantile function of X ; and $g(u, \cdot)$ is the conditional quantile function of Y given $X = f(u)$.

7a2 Proposition. Let (X_1, Y_1) be a 2-dimensional random variable (on some probability space), and (Y_2, Z_2) another 2-dimensional random variable (on another probability space) such that Y_1 and Y_2 are identically distributed. Then there exists (on some probability space) a 3-dimensional random variable (X, Y, Z) such that (X, Y) is distributed like (X_1, Y_1) , and (Y, Z) is distributed like (Y_2, Z_2) .

The idea is simple: X and Z are conditionally independent given Y . But what exactly does it mean?

7a3 Theorem (disintegration of measure). For every probability measure μ on \mathbb{R}^2 there exist a probability measure ν on \mathbb{R} and a family $(\mu_x)_{x \in \mathbb{R}}$ of probability measures μ_x on \mathbb{R} such that

$$\mu(A) = \int \mu_x(A_x) \nu(dx)$$

for all Borel sets $A \subset \mathbb{R}^2$; here $A_x = \{y : (x, y) \in A\}$. (The function $x \mapsto \mu_x(A_x)$ is thus claimed to be ν -measurable.)

Now, the problem in general.

The elementary definition

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$$

works if and only if $\mathbb{P}(B) \neq 0$. However, in many cases a limiting procedure gives us a useful result when $\mathbb{P}(B) = 0$.

7a4 Example. Let $(S_n)_n$ be the simple one-dimensional random walk and $B = \{\forall n, S_n > -10\}$ (a zero-probability event). We introduce $B_n = \{S_1 > -10, \dots, S_n > -10\}$, observe that $B_n \downarrow B$ and let $\mathbb{P}(A|B) = \lim_n \mathbb{P}(A|B_n)$ for “simple” A . In fact, we get a Markov chain with the transition probability $p_{k,k-1} = \frac{k+9}{2k+20}$, $p_{k,k+1} = \frac{k+11}{2k+20}$. However, the formula $\mathbb{P}(A|B) = \lim_n \mathbb{P}(A|B_n)$ should not be used for all A ; otherwise, trying $A = B$, we get a paradox: $\mathbb{P}(B|B) = 0$.

Similarly we may define the self-avoiding random walk on \mathbb{Z}^2 (assuming convergence); in fact, no one knows the joint distribution of the first two moves (even roughly)!

Sometimes different “reasonable” sequences $B_n \downarrow B$ lead to different results, which is known as Borel’s paradox or Borel-Kolmogorov paradox. For example,

$$\lim_{\varepsilon \rightarrow 0^+} \mathbb{P}(X \leq 1 \mid -\varepsilon < Y < \varepsilon) \neq \lim_{\varepsilon \rightarrow 0^+} \mathbb{P}(X \leq 1 \mid -\varepsilon|X| < Y < \varepsilon|X|).$$

Also, meridians (lines of longitude) and parallels (circles of latitude) on a sphere.

Sometimes conditioning is *really* impossible.

7a5 Example. Let $(S_n)_n$ be the simple one-dimensional random walk and $B = \{S_n \rightarrow +\infty\}$ (a zero-probability event). We note that $B = \{S_n - S_{10} \rightarrow +\infty\}$, “therefore” B is independent of S_1, \dots, S_{10} ; conditionally, given B , the walk behaves as usual, and we get the paradox, $\mathbb{P}(B|B) = 0$, once again.

7a6 Example. Let $X \sim U(0, 1)$ and $B = \{X \in \mathbb{Q}\}$. We consider $e^{2\pi i X}$; by symmetry, all rational points “must” get equal probabilities, which is impossible.

The conditional density formula

$$f_{Y|X=x}(y) = \frac{f_{X,Y}(x, y)}{f_X(x)}$$

works well whenever the joint distribution of X, Y is absolutely continuous. Neither this formula, nor the similar discrete formula

$$\mathbb{P}(Y = y | X = x) = \frac{\mathbb{P}(X = x, Y = y)}{\mathbb{P}(X = x)}$$

covers the (practically important) case of discrete Y but absolutely continuous X .¹ A still more complicated case is, conditioning of Y on $X = g(Y)$ when Y is absolutely continuous and g is not one-to-one (especially when g behaves like the Weierstrass function).

Bad news: we have no conditioning theory that covers under a single umbrella all “good” cases listed above. Good news (for those who do not fear of measure theory): we have a conditioning theory that covers conditioning of Y on X for an *arbitrary* joint distribution of random variables (or random vectors) X, Y , and this theory includes both the discrete case and the absolutely continuous case.

7b Discrete case

Let Ω be (at most) countable, $\mathcal{F} = 2^\Omega$, and $\mathcal{F}_1 \subset \mathcal{F}$ a sub- σ -algebra. Clearly, $\mathcal{F}_1 = \sigma(X)$ for some $X : \Omega \rightarrow \mathbb{R}$ (X just indexes the equivalence classes with some real numbers). Here is the elementary conditioning:

$$\begin{aligned} \mathbb{P}(A | X = x) &= \frac{\mathbb{P}(A \cap X^{-1}(x))}{\mathbb{P}(X^{-1}(x))} = \frac{\sum_{\omega \in A \cap X^{-1}(x)} p(\omega)}{\sum_{\omega \in X^{-1}(x)} p(\omega)} = P_x(A) = f(x), \\ \mathbb{P}(A | \mathcal{F}_1) &= \mathbb{P}(A | X) = f(X) : \Omega \rightarrow \mathbb{R}; \\ \mathbb{E}(Y | X = x) &= \frac{\sum_{\omega \in A \cap X^{-1}(x)} Y(\omega) p(\omega)}{\sum_{\omega \in X^{-1}(x)} p(\omega)} = \int Y \, dP_x = g(x), \\ \mathbb{E}(Y | \mathcal{F}_1) &= \mathbb{E}(Y | X) = g(X) : \Omega \rightarrow \mathbb{R}; \end{aligned}$$

each *conditional measure* P_x is a probability measure on Ω , concentrated on $X^{-1}(x)$; the map $x \rightarrow P_x$ depends on the choice of X , but the map $P_X : \omega \mapsto P_{X(\omega)}$ does not. Similarly, *regression functions* f and g depend on the choice of X , but the random variables $f(X)$ and $g(X)$ do not. The conditional probability is a special case of the conditional expectation: $Y = \mathbb{1}_A$.

¹See 7e18.

Convergence of the series is ensured if Y is integrable (except for negligible x , if any).

Note that $\mathbb{E}(Y|\mathcal{F}) = Y$ and $\mathbb{P}(A|\mathcal{F}) = \mathbb{1}_A$. On the other extreme, $\mathbb{E}(Y|\{\emptyset, \Omega\}) = \mathbb{E}Y$ (a constant function), and $\mathbb{P}(A|\{\emptyset, \Omega\}) = \mathbb{P}(A)$.

The *total probability formula* and *total expectation formula*

$$\begin{aligned}\mathbb{E}(\mathbb{P}(A|\mathcal{F}_1)) &= \mathbb{P}(A), \\ \mathbb{E}(\mathbb{E}(Y|\mathcal{F}_1)) &= \mathbb{E}Y\end{aligned}$$

boil down to the *decomposition of measure*,

$$P = \sum_x \mathbb{P}(X = x) \cdot P_x = \sum_{\omega} p(\omega) P_{X(\omega)} = \mathbb{E} P_X,$$

the latter expectation being taken in the linear space of signed measures...¹

The function $a \mapsto \mathbb{E}((Y-a)^2) = a^2 - 2a\mathbb{E}Y + \mathbb{E}(Y^2)$ reaches its minimum at $a = \mathbb{E}Y$ (assuming $Y \in L_2$, which evidently holds for $Y = \mathbb{1}_A$). That is, $\mathbb{E}Y$ is the orthogonal projection of Y to the one-dimensional space of constants, in $L_2(\Omega, \mathcal{F}, P)$. The same holds in each $L_2(P_x)$, thus, the regression function g minimizes $\mathbb{E}(Y - g(X))^2 = \mathbb{E}(\mathbb{E}((Y - g(X))^2|X))$. That is, $\mathbb{E}(Y|\mathcal{F}_1)$ is the orthogonal projection of Y to $L_2(\Omega, \mathcal{F}_1, P) \subset L_2(\Omega, \mathcal{F}, P)$.

7c Conditional expectation

We turn to the general case: (Ω, \mathcal{F}, P) is an arbitrary probability space, and $\mathcal{F}_1 \subset \mathcal{F}$ a sub- σ -algebra. We assume that all null sets belong to \mathcal{F} and also to \mathcal{F}_1 .

The Hilbert space $L_2(\mathcal{F}_1) = L_2(\Omega, \mathcal{F}_1, P)$ is a subspace of $L_2(\mathcal{F}) = L_2(\Omega, \mathcal{F}, P)$. We consider the orthogonal projection $L_2(\mathcal{F}) \rightarrow L_2(\mathcal{F}_1)$ and denote it $Y \mapsto \mathbb{E}(Y|\mathcal{F}_1)$. Note that $\mathbb{E}(Y|\mathcal{F}_1)$ is an equivalence class. Orthogonality means that $\langle Y - \mathbb{E}(Y|\mathcal{F}_1), X \rangle = 0$, that is,

$$\mathbb{E}(X \cdot \mathbb{E}(Y|\mathcal{F}_1)) = \mathbb{E}(XY) \quad \text{for all } X \in L_2(\mathcal{F}_1);$$

this property characterizes $\mathbb{E}(Y|\mathcal{F}_1)$ among $L_2(\mathcal{F}_1)$. In particular,

$$\mathbb{E}(\mathbb{E}(Y|\mathcal{F}_1)) = \mathbb{E}(Y),$$

the total expectation formula (not a characterization, of course). Moreover,

$$(7c1) \quad \mathbb{E}(\mathbb{E}(Y|\mathcal{F}_1); B) = \mathbb{E}(Y; B) \quad \text{for all } B \in \mathcal{F}_1;$$

¹Recall Sect. 5c.

also a characterization, since indicators and their linear combinations are dense in L_2 . A projection operator is always linear:

$$\begin{aligned}\mathbb{E}(aX | \mathcal{F}_1) &= a\mathbb{E}(X | \mathcal{F}_1), \\ \mathbb{E}(X + Y | \mathcal{F}_1) &= \mathbb{E}(X | \mathcal{F}_1) + \mathbb{E}(Y | \mathcal{F}_1), \\ \text{if } X_n \rightarrow X \text{ in } L_2 \text{ then } \mathbb{E}(X_n | \mathcal{F}_1) &\rightarrow \mathbb{E}(X | \mathcal{F}_1) \text{ in } L_2.\end{aligned}$$

But the subspace $L_2(\mathcal{F}_1)$ is special:

$$\text{if } X \in L_2(\mathcal{F}_1) \text{ then } X^+ \in L_2(\mathcal{F}_1).$$

It follows easily that the projection is positive:

$$\text{if } X \geq 0 \text{ a.s. then } \mathbb{E}(X | \mathcal{F}_1) \geq 0 \text{ a.s.}$$

Thus, the projection operator is continuous (of norm 1) also in the L_1 norm (apply the total expectation formula to X^+ and X^-), and therefore extends to $L_1(\Omega, \mathcal{F}, P)$ by continuity. It is still positive.

The “tower property”

$$\mathbb{E}(X | \mathcal{F}_1) = \mathbb{E}(\mathbb{E}(X | \mathcal{F}_2) | \mathcal{F}_1) \text{ a.s. if } \mathcal{F}_1 \subset \mathcal{F}_2 \subset \mathcal{F}$$

holds in L_2 for a simple geometric reason, and extends to L_1 by continuity.

7c2 Example. Let $Y \sim U(0, 1)$ and $X = f(Y)$,

$$f(y) = \begin{cases} 3y & \text{for } 0 \leq y \leq 1/3, \\ 1.5(1 - y) & \text{for } 1/3 \leq y \leq 2/3, \\ 0.5 & \text{for } 2/3 \leq y \leq 1. \end{cases}$$

Then $\mathbb{E}(Y | X) = g(X)$,

$$g(x) = \begin{cases} x/3 & \text{for } 0 < x < 0.5, \\ 5/6 & \text{for } x = 0.5, \\ (2 - x)/3 & \text{for } 0.5 < x < 1. \end{cases}$$

7c3 Exercise. Do it twice. Namely, (a) check it via (7c1); (b) derive it by minimization.

7c4 Exercise. Two σ -algebras $\mathcal{F}_1, \mathcal{F}_2 \subset \mathcal{F}$ are independent if and only if $\mathbb{E}(Y | \mathcal{F}_1) = \mathbb{E}Y$ a.s. for all $Y \in L_2(\mathcal{F}_2)$.

Formulate it accurately, and prove.

7c5 Exercise. Let $\mathbb{E}(Y|\mathcal{F}_1) = \mathbb{E}Y$ a.s. for a given $Y \in L_2(\mathcal{F})$.

- (a) Does it follow that Y is independent of \mathcal{F}_1 ?
- (b) Reconsider (a) assuming in addition that Y takes on only 3 values.
- (c) The same for only two values.

MORE ON MARTINGALES

Now we may use Definition 6a3 (rather than 6a4) in full generality.

A martingale bounded in L_2 converges in L_2 (recall Lemma 6c5 and the paragraph after it), $M_n \rightarrow M_\infty$ in L_2 . We have $M_\infty = M_\infty^+ - M_\infty^-$ and $M_n = \mathbb{E}(M_\infty|\mathcal{F}_n) = \mathbb{E}(M_\infty^+|\mathcal{F}_n) - \mathbb{E}(M_\infty^-|\mathcal{F}_n)$, the difference between two L_2 -bounded positive martingales. Proposition 6c6 is thus generalized.

7c6 Proposition. A martingale bounded in L_2 converges both in L_2 and almost surely.

7d A convergence theorem

In the end of Sect. 6c, when proving $\mathbb{P}(1 \leq Z_n = o((2p)^n)) = 0$, we used the relation $\mathbb{P}(A|\mathcal{F}_n) \rightarrow \mathbb{1}_A$ a.s. (for $A \in \mathcal{F}_\infty$). This relation is proved below.

Recall 3b13: if $\mathcal{F}_n \uparrow \mathcal{F}_\infty$ then $\cup_n L_2(\mathcal{F}_n)$ is dense in $L_2(\mathcal{F}_\infty)$. It follows that $\mathbb{E}(X|\mathcal{F}_n) \rightarrow X$ in L_2 for all $X \in L_2(\mathcal{F}_\infty)$, and $\mathbb{E}(X|\mathcal{F}_n) \rightarrow \mathbb{E}(X|\mathcal{F}_\infty)$ in L_2 for all $X \in L_2 = L_2(\mathcal{F})$. But this is a martingale, bounded in L_2 ; 7c6 ensures a.s. convergence, and we get the following.

7d1 Theorem. Let (Ω, \mathcal{F}, P) be a probability space and $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots \subset \mathcal{F}_\infty \subset \mathcal{F}$ sub- σ -algebras such that $\mathcal{F}_n \uparrow \mathcal{F}_\infty$. Then

$$\begin{aligned} \mathbb{E}(X|\mathcal{F}_n) &\rightarrow X \text{ a.s. and in } L_2 \text{ for all } X \in L_2(\mathcal{F}_\infty); \\ \mathbb{E}(X|\mathcal{F}_n) &\rightarrow \mathbb{E}(X|\mathcal{F}_\infty) \text{ a.s. and in } L_2 \text{ for all } X \in L_2(\mathcal{F}); \\ \mathbb{P}(A|\mathcal{F}_n) &\rightarrow \mathbb{1}_A \text{ a.s. and in } L_2 \text{ for all } A \in \mathcal{F}_\infty; \\ \mathbb{P}(A|\mathcal{F}_n) &\rightarrow \mathbb{P}(A|\mathcal{F}_\infty) \text{ a.s. and in } L_2 \text{ for all } A \in \mathcal{F}. \end{aligned}$$

7e Conditional measures

Let (Ω, \mathcal{F}, P) be a probability space, and $\mathcal{F}_1 \subset \mathcal{F}$ a sub- σ -algebra. The *conditional probability*

$$\mathbb{P}(A|\mathcal{F}_1) = \mathbb{E}(\mathbb{1}_A|\mathcal{F}_1)$$

satisfies

$$\begin{aligned} 0 \leq \mathbb{P}(A|\mathcal{F}_1) \leq 1 & \quad \text{a.s.}, \\ \mathbb{P}(A_1 \uplus A_2 \uplus \dots | \mathcal{F}_1) &= \mathbb{P}(A_1 | \mathcal{F}_1) + \mathbb{P}(A_2 | \mathcal{F}_1) + \dots \quad \text{a.s.}, \\ \forall B \in \mathcal{F}_1 \quad \mathbb{P}(B | \mathcal{F}_1) &= \mathbb{1}_B; \\ \mathbb{E}(\mathbb{P}(A | \mathcal{F}_1)) &= \mathbb{P}(A), \end{aligned}$$

which, however, does not mean that we can define conditional measures just by $P_\omega(A) = \mathbb{P}(A | \mathcal{F}_1)(\omega)$.

In the discrete case (at most countable Ω) we may get rid of all negligible points (if any) and define $P_\omega(A) = \mathbb{P}(A | \mathcal{F}_1)(\omega)$, getting

$$\begin{aligned} 0 \leq P_\omega(A) \leq 1; \\ P_\omega(A_1 \uplus A_2 \uplus \dots) &= P_\omega(A_1) + P_\omega(A_2) + \dots; \\ \forall B \in \mathcal{F}_1 \quad P_\omega(B) &= \mathbb{1}_B(\omega); \quad \text{especially, } P_\omega(\Omega) = 1; \\ \forall A \in \mathcal{F} \quad \int P_\omega(A) P(d\omega) &= P(A) \quad \left(\text{in this sense, } \int P_\omega P(d\omega) = P \right). \end{aligned}$$

It is a *disintegration* of P into probability measures P_ω localized on corresponding parts of the partition.

In general it does not go...

7e1 Definition. Let (Ω, \mathcal{F}, P) be a probability space and $\mathcal{F}_1 \subset \mathcal{F}$ a sub- σ -algebra. A *regular conditional probability* (given \mathcal{F}_1) is a family $(P_\omega)_{\omega \in \Omega}$ of probability measures P_ω on (Ω, \mathcal{F}) such that for every $A \in \mathcal{F}$ the function $\omega \mapsto P_\omega(A)$ belongs to the equivalence class $\mathbb{P}(A | \mathcal{F}_1)$.

The function $\omega \mapsto P_\omega(A)$ is not required to be \mathcal{F}_1 -measurable; rather, it must be equal a.s. to some \mathcal{F}_1 -measurable function (of the given equivalence class).

It is usual to write " $P_\omega(A) = \mathbb{P}(A | \mathcal{F}_1)(\omega)$ a.s." treating $\mathbb{P}(A | \mathcal{F}_1)$ as an (arbitrary) element of the equivalence class.

Only the equivalence class of the map $\omega \mapsto P_\omega$ matters; but the exceptional set should not depend on A .

Generally, a regular conditional probability need not exist (see 7e3).

7e2 Theorem. For $(\Omega, \mathcal{F}) = (\mathbb{R}, \mathcal{B})$ a regular conditional probability exists and is unique (up to equivalence).

Note that (a) P is an *arbitrary* Borel probability measure on \mathbb{R} ; (b) \mathcal{F} is the Borel σ -algebra; P -null sets are *not* added; (c) $\mathcal{F}_1 \subset \mathcal{F}$ is an arbitrary sub- σ -algebra.

7e3 Example. There exists¹ $Z \subset [0, 1]$ of interior Lebesgue measure 0 and exterior Lebesgue measure 1. We take $\Omega = [0, 1]$, $\mathcal{F} = \{(A \cap Z) \uplus (B \setminus Z) : A, B \in \mathcal{B}\}$, $\mathcal{F}_1 = \mathcal{B}$, and define a probability measure P on (Ω, \mathcal{F}) by $P((A \cap Z) \uplus (B \setminus Z)) = 0.5 \text{ mes } A + 0.5 \text{ mes } B$.

If $(P_\omega)_\omega$ is a regular conditional probability, then for almost every² $x \in [0, 1]$ the measure P_x must be equal to δ_x (the atom at x), since P_x is concentrated on every rational interval containing x . Thus $P_x(Z) = \mathbb{1}_Z(x)$ a.s., which contradicts to its measurability w.r.t. \mathcal{F}_1 .

7e4 Exercise. Show that $\mathbb{P}(Z | \mathcal{F}_1) = 0.5$ a.s.

PROOF OF THEOREM 7E2

The uniqueness part of Theorem 7e2 is easy: almost every x satisfies $P_x(I) = P'_x(I)$ for all rational intervals I , which implies $P_x = P'_x$.

The existence part needs more effort.

7e5 Definition. A measurable space is a pair (Ω, \mathcal{F}) of a set Ω and a σ -algebra \mathcal{F} on it.

Do not confuse “measure space” and “measurable space”! A probability measure on a measurable space turns it into a probability space.

Elements of \mathcal{F} are called measurable sets. A map between two measurable spaces is called measurable, if the inverse image of every measurable set is measurable. Two measurable spaces are called isomorphic, if there exists an isomorphism between them, that is, a measurable bijection with measurable inverse.

The disjoint union of two measurable spaces is a measurable space, $(\Omega, \mathcal{F}) = (\Omega', \mathcal{F}') \uplus (\Omega'', \mathcal{F}'')$.

A measurable part of a measurable space is itself a measurable space (and the original measurable space becomes the disjoint union).

By an embedding of one measurable space into another we mean an isomorphism between the former and a measurable part of the latter.

The following technical definition is introduced temporarily, for this proof only.

7e6 Definition. A measurable space (Ω, \mathcal{F}) is *good*, if a regular conditional probability exists for every probability measure on (Ω, \mathcal{F}) and every sub- σ -algebra of \mathcal{F} .

¹Using the choice axiom, of course.

²W.r.t. Lebesgue measure.

Existence of regular conditional probability on $(\mathbb{R}, \mathcal{B})$ becomes the claim that $(\mathbb{R}, \mathcal{B})$ is good. It follows from the three lemmas below (since a measurable space isomorphic to a good one is good).

7e7 Lemma. The Cantor set (with its Borel σ -algebra) is a good measurable space.

7e8 Lemma. The real line (with its Borel σ -algebra) is embeddable (as a measurable space) into the Cantor set (with its Borel σ -algebra).

7e9 Lemma. A measurable part of a good measurable space is a good measurable space. In other words, if $(\Omega, \mathcal{F}) = (\Omega', \mathcal{F}') \uplus (\Omega'', \mathcal{F}'')$ is good then (Ω', \mathcal{F}') is good.

Given $(\Omega, \mathcal{F}) = (\Omega', \mathcal{F}') \uplus (\Omega'', \mathcal{F}'')$ and two sub- σ -algebras, $\mathcal{F}'_1 \subset \mathcal{F}'$ on Ω' and $\mathcal{F}''_1 \subset \mathcal{F}''$ on Ω'' , we get the corresponding sub- σ -algebra $\mathcal{F}_1 \subset \mathcal{F}$ on Ω (namely, $\mathcal{F}_1 = \{A \uplus B : A \in \mathcal{F}'_1, B \in \mathcal{F}''_1\}$).

Proof of 7e9. Given a probability measure P on (Ω', \mathcal{F}') , we extend it to a probability measure on (Ω, \mathcal{F}) ($P(\Omega'') = 0$, necessarily). Further, given a sub- σ -algebra $\mathcal{F}'_1 \subset \mathcal{F}'$, we choose some sub- σ -algebra $\mathcal{F}''_1 \subset \mathcal{F}''$ (no matter which one) and get the corresponding sub- σ -algebra $\mathcal{F}_1 \subset \mathcal{F}$ on Ω such that $L_2(\Omega, \mathcal{F}_1) = L_2(\Omega', \mathcal{F}'_1)$ and therefore $\mathbb{E}(\cdot | \mathcal{F}_1) = \mathbb{E}(\cdot | \mathcal{F}'_1)$.

We take a regular conditional probability $(P_\omega)_{\omega \in \Omega}$ for \mathcal{F}_1 and restrict it to $(P_\omega)_{\omega \in \Omega'}$. We have $P_\omega(\Omega') = 1$ for almost all $\omega \in \Omega'$, since $P_\omega(\Omega') = \mathbb{P}(\Omega' | \mathcal{F}_1)(\omega) = \mathbb{1}_{\Omega'}(\omega) = 1$. Thus we may treat P_ω as a probability measure on (Ω', \mathcal{F}') , getting a regular conditional probability for \mathcal{F}'_1 . \square

Proof of 7e8. First, \mathbb{R} is isomorphic to $(0, 1)$ (as a topological space, the more so, as a measurable space).

Second, we embed $(0, 1)$ into the Cantor set via binary digits and observe that the image is a Borel set. \square

7e10 Lemma. If the function $\omega \mapsto P_\omega(A)$ belongs to the equivalence class $\mathbb{P}(A | \mathcal{F}_1)$ for all A of an algebra \mathcal{E} that generates \mathcal{F} then it holds for all $A \in \mathcal{F}$ (and therefore $(P_\omega)_{\omega \in \Omega}$ is a regular conditional probability).

Proof. Recall Sect. 1b: $P^*(A) = P(A)$ for all $A \in \mathcal{F}$, where P^* is the outer measure (defined via \mathcal{E}).

Let $B \in \mathcal{F}$; we have to prove that $\mathbb{P}(B | \mathcal{F}_1) = P_\bullet(B)$ a.s. (denoting the function $\omega \mapsto P_\omega(\dots)$ by $P_\bullet(\dots)$). Given $\varepsilon > 0$, we take $A_1, A_2, \dots \in \mathcal{E}$ such that $\cup_n A_n \supset B$ and $\sum_n P(A_n) \leq P(B) + \varepsilon$. Introducing a measurable function (and equivalence class) $h_\varepsilon = \sum_n P_\bullet(A_n) = \sum_n \mathbb{P}(A_n | \mathcal{F}_1)$ we get, on one hand, $P_\bullet(B) \leq h_\varepsilon$ (since $P_\omega(\cup_n A_n) \leq \sum_n P_\omega(A_n)$) and

$\mathbb{E} h_\varepsilon \leq P(B) + \varepsilon$ (since $\mathbb{E} P_\bullet(A_n) = \mathbb{E} \mathbb{P}(A_n | \mathcal{F}_1) = P(A_n)$), and on the other hand, $\mathbb{P}(B | \mathcal{F}_1) \leq h_\varepsilon$ (since $\mathbb{1}_B \leq \sum_n \mathbb{1}_{A_n}$).

Applying the same argument to $\Omega \setminus B$ we get a measurable function g_ε such that $\mathbb{E} g_\varepsilon \geq P(B) - \varepsilon$, $P_\bullet(B) \geq g_\varepsilon$ and $\mathbb{P}(B | \mathcal{F}_1) \geq g_\varepsilon$.

We have $|\mathbb{P}(B | \mathcal{F}_1) - P_\bullet(B)| \leq h_\varepsilon - g_\varepsilon$ and $\mathbb{E}(h_\varepsilon - g_\varepsilon) \leq 2\varepsilon$ for arbitrary ε ; therefore $\mathbb{P}(B | \mathcal{F}_1) = P_\bullet(B)$ a.s. \square

Proof of 7e7. The Borel σ -algebra of the Cantor set is generated by the countable algebra \mathcal{E} of clopen (that is, both closed and open) sets. Every finitely additive set function on this algebra is automatically σ -additive, due to compactness: if $A = A_1 \uplus A_2 \uplus \dots$ then $A_k = \emptyset$ for all k large enough. By Theorem 1b3, every finitely additive set function on this algebra extends to a measure.

We define P_ω by

$$P_\omega(A) = \mathbb{P}(A | \mathcal{F}_1)(\omega) \quad \text{for all } A \in \mathcal{A},$$

choosing a function in each equivalence class (countably many choices) and extend P_ω to a probability measure. Additivity holds a.s. (countably many equalities!) and is easily ensured everywhere. It remains to use Lemma 7e10. \square

Theorem 7e2 is thus proved.

7e11 Exercise. Generalize Theorem 7e2 to:

- (a) \mathbb{R}^n (with the Borel σ -algebra);
- (b) \mathbb{R}^∞ (all infinite sequences of reals with the σ -algebra generated by the coordinates).

Hint: only Lemma 7e8 needs to be generalized; embedding into the Cantor set ensures both existence and uniqueness.

Measurable spaces embeddable into the Cantor set are called *standard*. Observe that (a) Theorem 7e2 holds for all standard measurable spaces; (b) all \mathbb{R}^n and \mathbb{R}^∞ (and all their Borel subsets) are standard.¹

The requirement that the function $\omega \mapsto P_\omega(A)$ belongs to the equivalence class $\mathbb{P}(A | \mathcal{F}_1)$ can be reformulated using (7c1):

$$\int_B P_\omega(A) P(d\omega) = P(A \cap B) \quad \text{for all } B \in \mathcal{F}_1, A \in \mathcal{F}.$$

¹In fact, all Polish spaces (with Borel σ -algebras), and all their Borel subsets, are standard.

Or, equivalently,

$$(7e12) \quad P_\omega(B) = \mathbb{1}_B(\omega) \text{ a.s. for each } B \in \mathcal{F}_1 \text{ (separately),}$$

$$(7e13) \quad \int_{\Omega} P_\omega(A) P(d\omega) = P(A) \text{ for each } A \in \mathcal{F} \text{ (separately).}$$

Indeed, (7e12) implies $\int_B P_\omega(A) P(d\omega) = \int_{\Omega} P_\omega(A \cap B) P(d\omega)$.

7e14 Exercise. A given $A \in \mathcal{F}$ is independent of \mathcal{F}_1 if and only if $P_\omega(A) = P(A)$ for almost all ω .

Prove it.

Proof of Theorem 7a3. Theorem 7e2 (the existence part), applied to $(\Omega, \mathcal{F}, P) = (\mathbb{R}^2, \mathcal{B}_2, \mu)$ and \mathcal{F}_1 generated by the first coordinate, gives

$$\mu(A) = \int P_{x,y}(A) \mu(dxdy).$$

Measurability of $P_{x,y}$ w.r.t. \mathcal{F}_1 means that $P_{x,y} = P_x$. Property (7e12) means that P_x is concentrated on $\{x\} \times \mathbb{R}$. Thus, $P_{x,y}(A) = P_x(A) = \mu_x(A_x)$. The first projection of μ is ν . \square

7e15 Exercise. Prove uniqueness of $(\mu_x)_{x \in \mathbb{R}}$ up to a ν -negligible set.

The same holds for many other spaces.

7e16 Exercise. Disintegrate the joint distribution μ of the random variables X, Y of Example 7c2. Namely, guess the measures μ_x and check the equality $\mu(A) = \int \mu_x(A_x) \nu(dx)$.

7e17 Exercise. Disintegrate a measure that has a density (w.r.t. the 2-dimensional Lebesgue measure). Namely, guess the measures μ_x and check the equality $\mu(A) = \int \mu_x(A_x) \nu(dx)$.

7e18 Exercise. Derive formulas for conditioning of a discrete random variable $Y : \Omega \rightarrow \mathbb{Z}$ on a continuous random variable $X : \Omega \rightarrow \mathbb{R}$ that has a density f_X .

7e19 Exercise. Every random variable X is distributed like $f(U)$ for some increasing function $f : (0, 1) \rightarrow \mathbb{R}$; here U is distributed uniformly on $(0, 1)$.

Prove it.

Hint: first, prove it for a discrete X ; second, take discrete X_n such that $X_n \uparrow X$.

The function f is unique except for its values at discontinuity points (at most countable set).

Proof of Proposition 7a1. By 7e19, X is distributed like some $f(U)$. Theorem 7a3 gives us the conditional distributions μ_x of Y given $X = x$. By 7e19 (again), $\mu_{f(u)}$ is the distribution of some $g(u, V)$. For every y ,

$$m(\{v : g(u, v) \leq y\}) = \mu_{f(u)}((-\infty, y])$$

is measurable in u , and

$$\{(u, v) : v < \mu_{f(u)}((-\infty, y])\} \subset \{(u, v) : g(u, v) \leq y\} \subset \{(u, v) : v \leq \mu_{f(u)}((-\infty, y])\}$$

by monotonicity of $g(u, \cdot)$; it follows that g is measurable. Finally,

$$\begin{aligned} \mathbb{P}((f(U), g(U, V)) \in A) &= \int_0^1 du \int_0^1 dv \mathbb{1}_A(f(u), g(u, v)) = \\ &= \int_0^1 du \int_0^1 dv \mathbb{1}_{A_{f(u)}}(g(u, v)) = \int_0^1 du \mu_{f(u)}(A_{f(u)}) = \\ &= \int \mu_x(A_x) \nu(dx) = \mu(A) = \mathbb{P}((X, Y) \in A). \end{aligned}$$

□

Proof of Proposition 7a2. We have the one-dimensional distribution ν of Y_1 (and Y_2 as well), the two-dimensional distribution μ of (Y_1, X_1) and the two-dimensional distribution ξ of (Y_2, Z_2) . Theorem 7a3 gives us μ_y and ξ_y such that¹ $\mu(A) = \int \mu_y(A_y) \nu(dy)$ and $\xi(B) = \int \xi_y(B_y) \nu(dy)$ for all Borel sets $A, B \subset \mathbb{R}^2$; here $A_y = \{x : (y, x) \in A\}$ and $B_y = \{z : (y, z) \in B\}$. We introduce a measure η on \mathbb{R}^3 by²

$$\eta(C) = \int (\mu_y \times \xi_y)(C_y) \nu(dy)$$

for all Borel sets $C \subset \mathbb{R}^3$; here $C_y = \{(x, z) : (x, y, z) \in C\}$. For a three-dimensional random variable (X, Y, Z) distributed η we have

$$\begin{aligned} \mathbb{P}((Y, X) \in A) &= \eta(\{(x, y, z) : (y, x) \in A\}) = \int (\mu_y \times \xi_y)(A_y \times \mathbb{R}) \nu(dy) = \\ &= \int \mu_y(A_y) \nu(dy) = \mu(A) = \mathbb{P}((Y_1, X_1) \in A); \end{aligned}$$

and the same for (Y, Z) . □

¹ μ_y is the conditional distribution of X_1 given $Y_1 = y$; and ξ_y is the conditional distribution of Z_2 given $Y_2 = y$.

²Think, is the integrand measurable?

CONDITIONAL MEASURES AND CONDITIONAL EXPECTATIONS

7e20 Proposition. Let (Ω, \mathcal{F}, P) be a probability space, $\mathcal{F}_1 \subset \mathcal{F}$ a sub- σ -algebra, $(P_\omega)_{\omega \in \Omega}$ a regular conditional probability (given \mathcal{F}_1), and $Y \in L_1(\Omega, \mathcal{F}, P)$. Then the function $\omega \mapsto \int Y dP_\omega$ belongs to the equivalence class $\mathbb{E}(Y | \mathcal{F}_1)$.

Proof. The relation holds for indicators, and (by linearity) for their linear combinations. Let it hold for each Y_n , and $0 \leq Y_n \uparrow Y$ pointwise (not just a.s.), and $Y \in L_1(\Omega, \mathcal{F}, P)$; it is sufficient to prove that the relation holds for Y . We have $Y_n \rightarrow Y$ in L_1 , therefore $\mathbb{E}(Y_n | \mathcal{F}_1) \rightarrow \mathbb{E}(Y | \mathcal{F}_1)$ in L_1 . On the other hand, $\int Y_n dP_\omega \uparrow \int Y dP_\omega \in [0, \infty]$ for each ω by the monotone convergence theorem; the rest is easy. \square

It is tempting to extend $\mathbb{E}(\cdot | \mathcal{F}_1)$ to all Y such that $\int |Y| dP_\omega < \infty$ for almost all ω (even if $\int |Y| dP = \infty$). Then, however, strange things happen. For example, it may be that $\mathbb{E}(Y | \mathcal{F}_1) > 0$ a.s., but $\mathbb{E}(Y | \mathcal{F}_2) < 0$ a.s.¹

If X is \mathcal{F}_1 -measurable then P_ω is concentrated on $X^{-1}(X(\omega))$ for almost every ω . That is, $X(\omega') = X(\omega)$ for P_ω -almost all ω' .

“Taking out what is known”:

$$\mathbb{E}(XY | \mathcal{F}_1) = X\mathbb{E}(Y | \mathcal{F}_1) \quad \text{for } X \in L_\infty(\mathcal{F}_1), Y \in L_1(\mathcal{F}).$$

Conditional versions of many inequalities follow immediately from existence of regular conditional probability. Conditional Markov inequality:

$$\mathbb{P}(Y > X | \mathcal{F}_1) \leq \frac{\mathbb{E}(Y | \mathcal{F}_1)}{X} \quad \text{a.s. for } X \in L_0^+(\mathcal{F}_1), Y \in L_0^+(\mathcal{F}).$$

Conditional Jensen’s inequality:

$$\mathbb{E}(h(Y) | \mathcal{F}_1) \geq h(\mathbb{E}(Y | \mathcal{F}_1))$$

for convex $h(\cdot)$, provided that $\mathbb{E}(|X| | \mathcal{F}_1) < \infty$ a.s. Choosing $h(x) = |x|^p$ with $p \in (1, \infty)$ and taking the (unconditional) expectation we get

$$\|\mathbb{E}(Y | \mathcal{F}_1)\|_p \leq \|Y\|_p.$$

Conditional Cauchy-Schwartz inequality:

$$|\mathbb{E}(YZ | \mathcal{F}_1)| \leq \sqrt{\mathbb{E}(Y^2 | \mathcal{F}_1)} \sqrt{\mathbb{E}(Z^2 | \mathcal{F}_1)}.$$

And so on.

¹A counterexample (sketch): $\mathbb{P}(X = n, Y = n + 1) = \mathbb{P}(X = n + 1, Y = n) = 0.5p^n(1 - p)$ for $n = 0, 1, 2, \dots$; then $\mathbb{E}(a^Y | X = x) = \frac{pa + a^{-1}}{1 + p} a^x$ for $x = 1, 2, \dots$; we take $ap > 1$ and get $\mathbb{E}(a^Y | X) > a^X$ a.s., but also $\mathbb{E}(a^X | Y) > a^Y$ a.s.

7f Markov property as conditional independence

The “shortest” Markov process consists of 3 random variables X, Y, Z such that for every Borel set $B \subset \mathbb{R}$,

$$(7f1) \quad \mathbb{P}(Z \in B | X, Y) = \mathbb{P}(Z \in B | Y) \quad \text{a.s.}$$

In other words: in order to predict the future (Z) knowing the past (X) and the present (Y), use only the present; the past is irrelevant.

7f2 Exercise. A finite Markov chain (S_0, S_1, S_2, \dots) (recall Sect. 5b) satisfies

$$\mathbb{P}(S_{n+1} \in B | S_0, \dots, S_n) = \mathbb{P}(S_{n+1} \in B | S_n) \quad \text{a.s.}$$

for every n .

Prove it.

Note that 7f2 may be treated as (7f1) generalized to multi-dimensional $X = (S_0, \dots, S_{n-1})$.

It may seem that (7f1) is not time-symmetric; that is, this condition on (X, Y, Z) does not imply the same condition on (Z, Y, X) . However, let us consider it in the light of conditional distributions. Denoting $\mathbb{P}_y(\dots) = \mathbb{P}(\dots | Y = y)$ we rewrite (7f1) as $\mathbb{P}_y(Z \in B | X) = \mathbb{P}_y(Z \in B)$ a.s., which is just independence of Z and X (recall 7e14)... but conditionally, given $Y = y$.

In order to make this argument rigorous we need the following.

7f3 Proposition. Let (Ω, \mathcal{F}) be a standard measurable space, (Ω, \mathcal{F}, P) a probability space, $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \mathcal{F}$ sub- σ -algebras; $(P_\omega^{(1)})_{\omega \in \Omega}$ a regular conditional probability on (Ω, \mathcal{F}, P) given \mathcal{F}_1 ; and $(P_\omega^{(2)})_{\omega \in \Omega}$ a regular conditional probability on (Ω, \mathcal{F}, P) given \mathcal{F}_2 . Then for almost every ω_1 , $(P_\omega^{(2)})_{\omega \in \Omega}$ is also a regular conditional probability on $(\Omega, \mathcal{F}, P_{\omega_1}^{(1)})$ given \mathcal{F}_2 .

Proof. By 7e10 (and standardness) it is sufficient to check (7e12) and (7e13) on $(\Omega, \mathcal{F}, P_{\omega_1}^{(1)})$ for almost all ω_1 ; the exceptional set of ω_1 may depend on A, B . That is, we need

$$P_\omega^{(2)}(B) = \mathbb{1}_B(\omega) \quad P_{\omega_1}^{(1)}\text{-a.s.} \quad \text{for each } B \in \mathcal{F}_2 \text{ (separately),}$$

$$\int_{\Omega} P_\omega^{(2)}(A) P_{\omega_1}^{(1)}(d\omega) = P_{\omega_1}^{(1)}(A) \quad P\text{-a.s.} \quad \text{for each } A \in \mathcal{F} \text{ (separately).}$$

The former: the equality holds P -a.s. therefore it holds $P_{\omega_1}^{(1)}$ -a.s. for P -almost every ω_1 .

The latter: the function $f : \omega \rightarrow P_\omega^{(2)}(A)$ belongs to the equivalence class $\mathbb{P}(A | \mathcal{F}_2)$; thus, by 7e20, the function $\omega_1 \mapsto \int_\Omega f dP_{\omega_1}^{(1)} = \int_\Omega P_\omega^{(2)}(A) P_{\omega_1}^{(1)}(d\omega)$ belongs to the equivalence class $\mathbb{E}(f | \mathcal{F}_1) = \mathbb{E}(\mathbb{P}(A | \mathcal{F}_2) | \mathcal{F}_1) = \mathbb{P}(A | \mathcal{F}_1)$ (the tower property); and the function $\omega_1 \mapsto P_{\omega_1}^{(1)}(A)$ belongs to the same equivalence class. \square

We may apply it to a probability measure μ on \mathbb{R}^3 and its disintegrations:

$$\begin{aligned} \mu(A) &= \int \mu_y(A_y) \nu_1(dy), \quad A_y = \{(x, z) : (x, y, z) \in A\}; \\ \mu(A) &= \int \mu_{x,y}(A_{x,y}) \nu_2(dx dy), \quad A_{x,y} = \{z : (x, y, z) \in A\} \end{aligned}$$

for Borel $A \subset \mathbb{R}^3$. By 7f3, for ν_1 -almost every y we have the disintegration

$$\mu_y(A) = \int \mu_{x,y}(A_x) \mu_y(dx dz) = \int \mu_{x,y}(A_x) \nu_y(dx), \quad A_x = \{z : (x, z) \in A\}$$

for Borel $A \subset \mathbb{R}^2$; here ν_y is the marginal of μ_y .

In other words, for arbitrary random variables X, Y, Z we have

$$\begin{aligned} \mathbb{P}((X, Z) \in A | Y = y) &= \int \mathbb{P}((X, Z) \in A | X = x, Y = y) \nu_y(dx) = \\ &= \mathbb{E}(\mathbb{P}((X, Z) \in A | X, Y = y) | Y = y) \end{aligned}$$

for Borel $A \subset \mathbb{R}^2$ (since ν_y is the conditional distribution of X given $Y = y$). We get the conditioned (on $Y = y$) version of the formula $\mathbb{P}(\dots) = \mathbb{E}\mathbb{P}(\dots | X)$ mentioned in Sect. 7e.

7f4 Definition. Random variables X, Z are *conditionally independent* given Y , if for P_Y -almost every y the conditional distribution of (X, Z) given $Y = y$ is a product,

$$P_{X,Z|Y=y} = P_{X|Y=y} \times P_{Z|Y=y}.$$

An equivalent definition without conditional distributions:

$$\mathbb{P}(X \in A, Z \in B | Y) = \mathbb{P}(X \in A | Y) \mathbb{P}(Z \in B | Y) \quad \text{a.s.}$$

for all Borel sets $A, B \subset \mathbb{R}$.

7f5 Proposition. The Markov property (7f1) is equivalent to the conditional independence of X, Z given Y .

Now we can define easily the Markov property for random functions on a graph, on the lattice \mathbb{Z}^n , etc.