

CONDITIONAL EXPECTATION

STEVEN P. LALLEY

1. CONDITIONAL EXPECTATION: L^2 -THEORY

Definition 1. Let (Ω, \mathcal{F}, P) be a probability space and let \mathcal{G} be a σ -algebra contained in \mathcal{F} . For any real random variable $X \in L^2(\Omega, \mathcal{F}, P)$, define $E(X|\mathcal{G})$ to be the *orthogonal projection* of X onto the closed subspace $L^2(\Omega, \mathcal{G}, P)$.

This definition may seem a bit strange at first, as it seems not to have any connection with the naive definition of conditional probability that you may have learned in elementary probability. However, there is a compelling rationale for Definition 1: the orthogonal projection $E(X|\mathcal{G})$ minimizes the expected squared difference $E(X - Y)^2$ among all random variables $Y \in L^2(\Omega, \mathcal{G}, P)$, so in a sense it is the *best predictor* of X based on the information in \mathcal{G} . It may be helpful to consider the special case where the σ -algebra \mathcal{G} is generated by a single random variable Y , i.e., $\mathcal{G} = \sigma(Y)$. In this case, every \mathcal{G} -measurable random variable is a Borel function of Y (exercise!), so $E(X|\mathcal{G})$ is the unique Borel function $h(Y)$ (up to sets of probability zero) that minimizes $E(X - h(Y))^2$. The following exercise indicates that the special case where $\mathcal{G} = \sigma(Y)$ for some real-valued random variable Y is in fact very general.

Exercise 1. Show that if \mathcal{G} is *countably generated* (that is, there is some countable collection of set $B_j \in \mathcal{G}$ such that \mathcal{G} is the *smallest* σ -algebra containing all of the sets B_j) then there is a \mathcal{G} -measurable real random variable Y such that $\mathcal{G} = \sigma(Y)$.

The following exercise shows that in the special case where the σ -algebra \mathcal{G} is finite, Definition 1 is equivalent to the naive definition of conditional expectation.

Exercise 2. Suppose that the σ -algebra \mathcal{G} is finite, that is, suppose that there is a finite measurable partition B_1, B_2, \dots, B_n of Ω such that \mathcal{G} is the σ -algebra generated by the sets B_i . Show that for any $X \in L^2(\Omega, \mathcal{F}, P)$,

$$E(X|\mathcal{G}) = \sum_{i=1}^n \frac{E(X\mathbf{1}_{B_i})}{P(B_i)} \mathbf{1}_{B_i} \quad a.s.$$

Because conditional expectation is defined by orthogonal projection, all of the elementary properties of orthogonal projection operators translate to corresponding properties of conditional expectation. Following is a list of some of the more important of these.

Properties of Conditional Expectation: Let $X \in L^2(\Omega, \mathcal{F}, P)$ and let \mathcal{G} be a σ -algebra contained in \mathcal{F} . Then

(0) **Linearity:** $E(aX_1 + bX_2|\mathcal{G}) = aE(X_1|\mathcal{G}) + bE(X_2|\mathcal{G})$.

(1) **Orthogonality:** $X - E(X|\mathcal{G}) \perp L^2(\Omega, \mathcal{G}, P)$.

- (2) **Best Prediction:** $E(X|\mathcal{G})$ minimizes $E(X - Y)^2$ among all $Y \in L^2(\Omega, \mathcal{G}, P)$.
(3) **Tower Property:** If \mathcal{H} is a σ -algebra contained in \mathcal{G} , so that $\mathcal{H} \subseteq \mathcal{G} \subseteq \mathcal{F}$, then

$$E(X|\mathcal{H}) = E(E(X|\mathcal{G})|\mathcal{H}).$$

- (4) **Covariance Matching:** $E(X|\mathcal{G})$ is the unique random variable $Z \in L^2(\Omega, \mathcal{G}, P)$ such that for every $Y \in L^2(\Omega, \mathcal{G}, P)$,

$$E(XY) = E(ZY).$$

Property (4) is just a re-statement of the orthogonality law. It is usually taken to be the *definition* of conditional expectation (as in Billingsley). Observe that for the equation in (4) to hold for *all* $Y \in L^2(\Omega, \mathcal{G}, P)$ it is necessary that Z be square-integrable (because the equation must hold for $Y = Z$). That there is only one such random variable Z (up to change on events of probability 0) follows by an easy argument using indicators of events $B \in \mathcal{G}$: if there were two \mathcal{G} -measurable random variables Z_1, Z_2 such that the equation in (4) were valid for both $Z = Z_1$ and $Z = Z_2$, and all Y , then for every $B \in \mathcal{G}$,

$$E(Z_1 - Z_2)\mathbb{1}_B = 0.$$

But any \mathcal{G} -measurable random variable that integrates to 0 on *every* event $B \in \mathcal{G}$ must equal 0 a.s. (why?).

2. CONDITIONAL EXPECTATION: L^1 -THEORY

The major drawback of Definition 1 is that it applies only to square-integrable random variables. Thus, our next objective will be to extend the definition and basic properties of conditional expectation to all *integrable* random variables. Since an integrable random variable X need not be square-integrable, its conditional expectation $E(X|\mathcal{G})$ on a σ -algebra \mathcal{G} cannot be defined by orthogonal projection. Instead, we will use the covariance property (4) as a basis for a general definition.

Definition 2. Let (Ω, \mathcal{F}, P) be a probability space and let \mathcal{G} be a σ -algebra contained in \mathcal{F} . For any real random variable $X \in L^1(\Omega, \mathcal{F}, P)$, define $E(X|\mathcal{G})$ to be the unique random variable $Z \in L^1(\Omega, \mathcal{G}, P)$ such that for every *bounded*, \mathcal{G} -measurable random variable Y ,

$$(5) \quad E(XY) = E(ZY).$$

In particular, equation (5) must hold for every indicator $Y = \mathbf{1}_G$, where $G \in \mathcal{G}$. We have already seen that there can be at most one such random variable $Z \in L^1(\Omega, \mathcal{G}, P)$; thus, to verify that Definition 2 is a valid definition, we must prove that there is *at least one* random variable $Z \in L^1(\Omega, \mathcal{G}, P)$ satisfying equation (5).

Proposition 1. For every $X \in L^1(\Omega, \mathcal{F}, P)$ there exists $Z \in L^1(\Omega, \mathcal{G}, P)$ such that equation (5) holds for all bounded, \mathcal{G} -measurable random variables Y .

First Proof. There is a short, easy proof using the Radon-Nikodym theorem. It is enough to consider the case where X is nonnegative, because in general an integrable random variable can be split into positive and negative parts. Assume, then, that $X \geq 0$. Define a finite, positive measure ν on (Ω, \mathcal{G}) by

$$\nu(G) = EX\mathbf{1}_G.$$

It is easily checked that $\nu \ll P$ on \mathcal{G} . Therefore, the Radon-Nikodym theorem implies that there exists a nonnegative, \mathcal{G} -measurable random variable Z such that for every bounded, \mathcal{G} -measurable Y ,

$$\int Y d\nu = E(ZY).$$

□

Although it is short and elegant, the preceding proof relies on a deep theorem, the Radon-Nikodym theorem. In fact, the use of the Radon-Nikodym theorem is superfluous; the fact that every L^1 random variable can be arbitrarily approximated by L^2 random variables makes it possible to construct a solution to (5) by approximation. For this, we need several more properties of the conditional expectation operator on L^2 .

(6) **Normalization:** $E(1|\mathcal{G}) = 1$ almost surely.

(7) **Positivity:** For any nonnegative, bounded random variable X ,

$$E(X|\mathcal{G}) \geq 0 \quad \text{almost surely.}$$

(8) **Monotonicity:** If X, Y are bounded random variables such that $X \leq Y$ a.s., then

$$E(X|\mathcal{G}) \leq E(Y|\mathcal{G}) \quad \text{almost surely.}$$

The normalization property (6) is almost trivial: it holds because any constant random variable c is measurable with respect to *any* σ -algebra, and in particular \mathcal{G} ; and any random variable $Y \in L^2(\Omega, \mathcal{G}, P)$ is its own projection. The positivity property (7) can be easily deduced from the covariance matching property (4): since $X \geq 0$ a.s., for any event $G \in \mathcal{G}$,

$$E(X\mathbf{1}_G) = E(E(X|\mathcal{G})\mathbf{1}_G) \geq 0;$$

consequently, $E(X|\mathcal{G})$ must be nonnegative almost surely, because its integral on any event is nonnegative. The monotonicity property (8) follows directly from *linearity* and *positivity*.

Second Proof of Proposition 1. First, observe again that it suffices to consider the case where X is nonnegative. Next, recall that $X = \lim \uparrow X \wedge n$. Each of the random variables $X \wedge n$ is bounded, hence in L^2 , and so its conditional expectation is well-defined by orthogonal projection. Moreover, by the positivity and monotonicity laws (7) – (8), the conditional expectations $E(X \wedge n|\mathcal{G})$ are nonnegative and non-decreasing with n . Consequently,

$$Z := \lim \uparrow E(X \wedge n|\mathcal{G})$$

exists and is \mathcal{G} -measurable. Now by the Monotone Convergence Theorem, for any bounded, nonnegative, \mathcal{G} -measurable random variable Y ,

$$\begin{aligned} E(XY) &= \lim \uparrow E((X \wedge n)Y) \\ &= \lim \uparrow E(E(X \wedge n|\mathcal{G})Y) \\ &= E(ZY). \end{aligned}$$

This proves equation (5) for *nonnegative* Y ; it then follows for arbitrary bounded Y by linearity, using $Y = Y_+ - Y_-$. □

2.1. Properties of Conditional Expectation. Henceforth we shall take Definition 2 to be the definition of conditional expectation. By the covariance matching property (4), this definition agrees with Definition 1 for $X \in L^2$. Given Definition 2, the following properties are all easily established. (Exercise: Check any that you do not find immediately obvious.)

- (1) **Definition:** $EXY = E(E(X|\mathcal{G})Y)$ for all bounded, \mathcal{G} -measurable random variables Y .
- (2) **Linearity:** $E(aU + bV|\mathcal{G}) = aE(U|\mathcal{G}) + bE(V|\mathcal{G})$ for all scalars $a, b \in \mathbb{R}$.
- (3) **Positivity:** If $X \geq 0$ then $E(X|\mathcal{G}) \geq 0$.
- (4) **Stability:** If X is \mathcal{G} -measurable, then $E(XZ|Y) = XE(Z|Y)$.
- (5) **Independence Law:** If X is independent of \mathcal{G} then $E(X|\mathcal{G}) = EX$ is constant a.s.
- (6) **Tower Property:** If $\mathcal{H} \subseteq \mathcal{G}$ then $E(E(X|\mathcal{G})|\mathcal{H}) = E(X|\mathcal{H})$.
- (7) **Expectation Law:** $E(E(X|\mathcal{G})) = EX$.
- (8) **Constants:** For any scalar a , $E(a|\mathcal{G}) = a$.
- (9) **Jensen Inequalities:** If $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is convex and $E|X| < \infty$ then

$$E(\varphi(X)) \geq \varphi(EX) \text{ and}$$

$$E(\varphi(X)|Y) \geq \varphi(E(X|Y)).$$

In all of these statements, the relations $=$ and \leq are meant to hold *almost surely*. Properties (3)–(7) extend easily to nonnegative random variables X with infinite expectation. Observe that, with the exceptions of the Stability, Tower, and Independence Properties, all of these correspond to basic properties of *ordinary* expectation. Later, we will see a deeper reason for this. Following is another property of conditional expectation that generalizes a corresponding property of ordinary expectation.

Proposition 2. (*Jensen Inequality*) *If $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is convex and $E|X| < \infty$ then*

$$(9) \quad E(\varphi(X)|\mathcal{G}) \geq \varphi(E(X|\mathcal{G})).$$

REMARK. Since (9) holds for the trivial σ -algebra $\{\emptyset, \Omega\}$, the usual Jensen inequality $E\varphi(X) \geq \varphi(EX)$ is a consequence.

Proof of the Jensen Inequalities. One of the basic properties of convex functions is that every point on the graph of a convex function φ has a *support line*: that is, for every argument $x_* \in \mathbb{R}$ there is a linear function $y_{x_*}(x) = ax + b$ such that

$$\varphi(x_*) = y_{x_*}(x_*) \quad \text{and}$$

$$\varphi(x) \geq y_{x_*}(x) \quad \text{for all } x \in \mathbb{R}.$$

Let X be a random variable such that $E|X| < \infty$, so that the expectation EX is well-defined and finite. Let $y_{EX}(x) = ax + b$ be a support line to the convex function at the point $(EX, \varphi(EX))$. Then by definition of a support line, $y_{EX}(EX) = \varphi(EX)$; also, $y_{EX}(X) \leq \varphi(X)$, and so

$$Ey_{EX}(X) \leq E\varphi(X).$$

But because $y_{EX}(x) = ax + b$ is a linear function of x ,

$$Ey_{EX}(X) = y_{EX}(EX) = \varphi(EX).$$

This proves the Jensen inequality for ordinary expectation. The proof for conditional expectation is similar. Let $y_{E(X|\mathcal{G})}(x)$ be the support line at the point $(E(X|\mathcal{G}), \varphi(E(X|\mathcal{G})))$. Then

$y_{E(X|\mathcal{G})}(E(X|\mathcal{G})) = \varphi(E(X|\mathcal{G}))$, and for every value of X , $y_{E(X|\mathcal{G})}(X) \leq \varphi(X)$. Consequently, by the linearity and positivity properties of conditional expectation,

$$\begin{aligned}\varphi(E(X|\mathcal{G})) &= y_{E(X|\mathcal{G})}(E(X|\mathcal{G})) \\ &= E(y_{E(X|\mathcal{G})}(X)|\mathcal{G}) \\ &\leq E(\varphi(X)|\mathcal{G}).\end{aligned}$$

□

Proposition 3. *Let $\{X_\lambda\}_{\lambda \in \Lambda}$ be a (not necessarily countable) uniformly integrable family of real random variables on (Ω, \mathcal{F}, P) . Then the set of all conditional expectations $\{E(X_\lambda|\mathcal{G})\}_{\lambda \in \Lambda, \mathcal{G} \subset \mathcal{F}}$, where \mathcal{G} ranges over all σ -algebras contained in \mathcal{F} , is uniformly integrable.*

Proof. We will use the equivalence of the following two characterizations of uniform integrability: a family $\{Y_\lambda\}_{\lambda \in \Lambda}$ of nonnegative random variables is uniformly integrable if either of the following holds.

- (a) For each $\varepsilon > 0$ there exists $\alpha < \infty$ such that $EY_\lambda \mathbf{1}_{\{Y_\lambda \geq \alpha\}} < \varepsilon$.
- (b) The expectations EY_λ are bounded (that is, $\sup_{\lambda \in \Lambda} EY_\lambda < \infty$), and for each $\varepsilon > 0$ there exists $\delta > 0$ such that for every event A of probability less than δ ,

$$\sup_{\lambda \in \Lambda} EY_\lambda \mathbf{1}_A < \varepsilon.$$

We may assume without loss of generality in proving the proposition that all of the random variables X_λ are nonnegative. Since the family $\{X_\lambda\}_{\lambda \in \Lambda}$ is uniformly integrable, for each $\varepsilon > 0$ there exists $\delta > 0$ such that for every event A satisfying $P(A) < \delta$ and every $\lambda \in \Lambda$,

$$(10) \quad EX_\lambda \mathbf{1}_A < \varepsilon.$$

Moreover, since the set $\{X_\lambda\}_{\lambda \in \Lambda}$ is uniformly integrable, the expectations are uniformly bounded, and so the set $\{E(X_\lambda|\mathcal{G})\}_{\lambda \in \Lambda, \mathcal{G} \subset \mathcal{F}}$ of all conditional expectations is bounded in L^1 . Hence, by the Markov inequality,

$$\lim_{\alpha \rightarrow \infty} \sup_{\lambda \in \Lambda} \sup_{\mathcal{G} \subset \mathcal{F}} P\{E(X_\lambda|\mathcal{G}) \geq \alpha\} = 0.$$

To show that the set of all conditional expectations $\{E(X_\lambda|\mathcal{G})\}_{\lambda \in \Lambda, \mathcal{G} \subset \mathcal{F}}$ is uniformly integrable it suffices, by criterion (a), to prove that for every $\varepsilon > 0$ there exists a constant $0 < \alpha < \infty$ such that for each $\lambda \in \Lambda$ and each σ -algebra $\mathcal{G} \subset \mathcal{F}$,

$$(11) \quad EE(X_\lambda|\mathcal{G}) \mathbf{1}_{\{E(X_\lambda|\mathcal{G}) > \alpha\}} < \varepsilon.$$

But by the definition of conditional expectation,

$$EE(X_\lambda|\mathcal{G}) \mathbf{1}_{\{E(X_\lambda|\mathcal{G}) > \alpha\}} = EX_\lambda \mathbf{1}_{\{E(X_\lambda|\mathcal{G}) > \alpha\}}.$$

If $\alpha < \infty$ is sufficiently large then by the preceding paragraph $P\{E(X_\lambda|\mathcal{G}) > \alpha\} < \delta$, where $\delta > 0$ is so small that the inequality (10) holds for all events A with probability less than δ . The desired inequality (11) now follows. □

3. CONVERGENCE THEOREMS FOR CONDITIONAL EXPECTATION

Just as for ordinary expectations, there are versions of Fatou's lemma and the monotone and dominated convergence theorems.

Monotone Convergence Theorem . *Let X_n be a nondecreasing sequence of nonnegative random variables on a probability space (Ω, \mathcal{F}, P) , and let $X = \lim_{n \rightarrow \infty} X_n$. Then for any σ -algebra $\mathcal{G} \subset \mathcal{F}$,*

$$(12) \quad E(X_n | \mathcal{G}) \uparrow E(X | \mathcal{G}).$$

Fatou's Lemma . *Let X_n be a sequence of nonnegative random variables on a probability space (Ω, \mathcal{F}, P) , and let $X = \liminf_{n \rightarrow \infty} X_n$. Then for any σ -algebra $\mathcal{G} \subset \mathcal{F}$,*

$$(13) \quad E(X | \mathcal{G}) \leq \liminf E(X_n | \mathcal{G}).$$

Dominated Convergence Theorem . *Let X_n be a sequence of real-valued random variables on a probability space (Ω, \mathcal{F}, P) such that for some integrable random variable Y and all $n \geq 1$,*

$$(14) \quad |X_n| \leq Y.$$

Then for any σ -algebra $\mathcal{G} \subset \mathcal{F}$,

$$(15) \quad \lim_{n \rightarrow \infty} E(X_n | \mathcal{G}) = E(X | \mathcal{G}) \quad \text{and} \quad \lim_{n \rightarrow \infty} E(|X_n - X| | \mathcal{G}) = 0$$

As in Properties 1–9 above, the limiting equalities and inequalities in these statements hold *almost surely*. The proofs are easy, given the corresponding theorems for *ordinary* expectations; I'll give the proof for the Monotone Convergence Theorem and leave the other two, which are easier, as exercises.

Proof of the Monotone Convergence Theorem. This is essentially the same argument as used in the second proof of Proposition 1 above. By the Positivity and Linearity properties of conditional expectation,

$$E(X_n | \mathcal{G}) \leq E(X_{n+1} | \mathcal{G}) \leq E(X | \mathcal{G})$$

for every n . Consequently, the limit $V := \lim_{n \rightarrow \infty} \uparrow E(X_n | \mathcal{G})$ exists with probability one, and $V \leq E(X | \mathcal{G})$. Moreover, since each conditional expectation is \mathcal{G} -measurable, so is V . Set $B = \{V < E(X | \mathcal{G})\}$; we must show that $P(B) = 0$. Now $B \in \mathcal{G}$, so by definition of conditional expectation,

$$E(X \mathbf{1}_B) = E(E(X | \mathcal{G}) \mathbf{1}_B) \quad \text{and} \quad E(X_n \mathbf{1}_B) = E(E(X_n | \mathcal{G}) \mathbf{1}_B).$$

But the Monotone Convergence Theorem for ordinary expectation implies that

$$\begin{aligned} E X \mathbf{1}_B &= \lim_{n \rightarrow \infty} E X_n \mathbf{1}_B \quad \text{and} \\ E V \mathbf{1}_B &= \lim_{n \rightarrow \infty} E E(X_n | \mathcal{G}) \mathbf{1}_B, \end{aligned}$$

so $E X \mathbf{1}_B = E V \mathbf{1}_B$. Since $V < X$ on B , this implies that $P(B) = 0$. □

4. REGULAR CONDITIONAL DISTRIBUTIONS

Let X be a real random variable defined on a probability space (Ω, \mathcal{F}, P) and let \mathcal{G} be a σ -algebra contained in \mathcal{F} . For each Borel set $B \subset \mathbb{R}$ define

$$(16) \quad P(X \in B | \mathcal{G}) = E(\mathbf{1}_{\{X \in B\}} | \mathcal{G})$$

to be the *conditional probability* of the event $\{X \in B\}$ given \mathcal{G} . Observe that the conditional probability $P(X \in B | \mathcal{G})$ is a *random variable*, not a scalar. Nevertheless, conditional probability, like unconditional probability, is countably additive, in the sense that for any sequence B_n of pairwise disjoint Borel sets,

$$(17) \quad P(\cup_{n \geq 1} B_n | \mathcal{G}) = \sum_{n \geq 1} P(B_n | \mathcal{G}) \quad \text{a.s.}$$

(This follows by the linearity of conditional expectation and the monotone convergence theorem, as you should check.) Unfortunately, there are *uncountably* many sequences of Borel sets, so we cannot be sure *a priori* that the null sets on which (17) fails do not add up to an event of positive probability. The purpose of this section is to show that in fact they do not.

Definition 3. Let $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$ and $(\mathcal{Y}, \mathcal{F}_{\mathcal{Y}})$ be measurable spaces. A *Markov kernel* (or *Markov transition kernel*) from $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$ to $(\mathcal{Y}, \mathcal{F}_{\mathcal{Y}})$ is a family $\{\mu_x(dy)\}_{x \in \mathcal{X}}$ of probability measures on $(\mathcal{Y}, \mathcal{F}_{\mathcal{Y}})$ such that for each event $F \in \mathcal{F}_{\mathcal{Y}}$ the random variable $x \mapsto \mu_x(F)$ is $\mathcal{F}_{\mathcal{X}}$ -measurable.

Definition 4. Let X be a random variable defined on a probability space (Ω, \mathcal{F}, P) that takes values in a measurable space $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$, and let \mathcal{G} be a σ -algebra contained in \mathcal{F} . A *regular conditional distribution* for X given \mathcal{G} is a Markov kernel μ_ω from (Ω, \mathcal{G}) to $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$ such that for every set $F \in \mathcal{F}_{\mathcal{X}}$,

$$(18) \quad \mu_\omega(F) = P(X \in F | \mathcal{G})(\omega) \quad \text{a.s.}$$

Theorem 1. *If X is a real random variable defined on a probability space (Ω, \mathcal{F}, P) then for every σ -algebra $\mathcal{G} \subset \mathcal{F}$ there is a regular conditional distribution for X given \mathcal{G} .*

The proof will use the *quantile transform* method of constructing a random variable with a specified cumulative distribution function F from a uniform-[0,1] random variable. Given a c.d.f. F , define

$$(19) \quad F^-(t) = \inf\{x : F(x) \geq t\} \quad \text{for } t \in (0, 1).$$

Lemma 1. *If U has the uniform distribution on $[0, 1]$ then $F^-(U)$ has cumulative distribution function F , that is, for every $x \in \mathbb{R}$,*

$$P(F^-(U) \leq x) = F(x).$$

Proof. Exercise. □

Corollary 1. *For any cumulative distribution function F on \mathbb{R} (that is, any right-continuous, nondecreasing function satisfying $F \rightarrow 0$ at $-\infty$ and $F \rightarrow 1$ at $+\infty$) there is a Borel probability measure μ_F on \mathbb{R} such that for every $x \in \mathbb{R}$,*

$$(20) \quad \mu_F(-\infty, x] = F(x).$$

Proof. This follows directly from the following more general fact: if (Ω, \mathcal{F}, P) is any probability space and if $X : \Omega \rightarrow \mathcal{X}$ is a measurable transformation taking values in $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$ then $\mu = P \circ X^{-1}$ is a probability measure on $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$. \square

Thus, the quantile transform together with the existence of Lebesgue measure on $[0, 1]$ implies the existence of a measure satisfying (20). In proving Theorem 1 we shall use this to reduce the problem of constructing probability measures $\mu_{\omega}(dx)$ to the simpler problem of constructing c.d.f.s $F_{\omega}(x)$ such that

$$(21) \quad F_{\omega}(x) = P(X \leq x | \mathcal{G}) \quad \text{a.s.}$$

Proof of Theorem 1. For each rational x define $F_{\omega}(x) = P(X \leq x | \mathcal{G})(\omega)$ to be some version of the conditional probability. Consider the following events:

$$\begin{aligned} B_1 &= \{F_{\omega} \text{ not monotone on } \mathbb{Q}\}; \\ B_2 &= \{\exists x \in \mathbb{Q} : F_{\omega}(x) \neq \lim_{y \rightarrow x^+} F_{\omega}(y)\}; \\ B_3 &= \{\lim_{n \rightarrow \infty} F_{\omega}(n) \neq 1\}; \\ B_4 &= \{\lim_{n \rightarrow \infty} F_{\omega}(-n) \neq 0\}; \\ B &= B_1 \cup B_2 \cup B_3 \cup B_4. \end{aligned}$$

(In B_2 , the limit is through the set of rationals, and in B_3 and B_4 the limit is through the positive integers n .) These are all events of probability 0. (Exercise: why?)

On the event B redefine $F_{\omega}(x) = \Phi(x)$, where Φ is the standard normal c.d.f. On B^c , and on B^c extend F_{ω} to all real arguments by setting

$$F_{\omega}(x) = \inf_{y \in \mathbb{Q}, y > x} F_{\omega}(y).$$

Note that for rational x this definition coincides with the original choice of $F_{\omega}(x)$ (because the event $B^c \subset B_2^c$). Then $F_{\omega}(x)$ is, for every $\omega \in \Omega$, a cumulative distribution function, and so there exists a Borel probability measure $\mu_{\omega}(dx)$ on \mathbb{R} with c.d.f. F_{ω} . The proof will be completed by establishing the following two claims.

Claim 1: $\{\mu_{\omega}(dx)\}_{\omega \in \Omega}$ is a Markov kernel.

Claim 2: $\{\mu_{\omega}(dx)\}_{\omega \in \Omega}$ is a regular conditional distribution for X given \mathcal{G} .

Proof of Claim 1. Each μ_{ω} is by construction a Borel probability measure, so what must be proved is that for each Borel set B the mapping $\omega \mapsto \mu_{\omega}(B)$ is \mathcal{G} -measurable. This is certainly true for each interval $(-\infty, x]$ with $x \in \mathbb{Q}$, since $F_{\omega}(x)$ is a version of $P(X \leq x | \mathcal{G})$, which by definition is \mathcal{G} -measurable. Moreover, by countable additivity, the set of Borel sets B such that $\omega \mapsto \mu_{\omega}(B)$ is \mathcal{G} -measurable is a λ -system (why?). Since the set of intervals $(-\infty, x]$ with $x \in \mathbb{Q}$ is a π -system, the claim follows from the π - λ theorem. \square

Proof of Claim 2. To show that $\{\mu_{\omega}(dx)\}_{\omega \in \Omega}$ is a regular conditional distribution for X , we must show that equation (18) holds for every Borel set F . By construction, it holds for every $F = (-\infty, x]$ with x rational, and it is easily checked that the collection of all F for which (18) holds is a λ -system (why?). Therefore, the claim follows from the π - λ theorem. \square

\square

Regular conditional distributions are useful in part because they allow one to reduce many problems concerning conditional expectations to problems concerning only *ordinary* expectations. For such applications the following *disintegration formula* for conditional expectations is essential.

Theorem 2. Let $\mu_\omega(dx)$ be a regular conditional distribution for X given \mathcal{G} , let Y be \mathcal{G} -measurable, and let $f(x, y)$ be a jointly measurable real-valued function such that $E|f(X, Y)| < \infty$. Then

$$(22) \quad E(f(X, Y) | \mathcal{G}) = \int f(x, Y(\omega)) \mu_\omega(dx) \quad \text{a.s.}$$

Proof. By the usual arguments (linearity plus approximation by simple functions plus monotone convergence) it suffices to prove this in the special case where $f = \mathbf{1}_F$ is an indicator function. By the usual $\pi - \lambda$ argument, it suffices to consider sets F of the form $F = A \times B$, so $f(X, Y) = \mathbf{1}_A(X)\mathbf{1}_B(Y)$. Thus, our task is to show that

$$E(\mathbf{1}_A(X)\mathbf{1}_B(Y) | \mathcal{G}) = \int \mathbf{1}_A(x)\mathbf{1}_B(Y(\omega)) \mu_\omega(dx).$$

Since the indicator $\mathbf{1}_B(Y)$ factors out on both sides (by the stability property of conditional expectation – see section 2.1 above) it is enough to show that for any measurable set A ,

$$E(\mathbf{1}_A(X) | \mathcal{G})(\omega) = \int \mathbf{1}_A(x) \mu_\omega(dx).$$

But this follows directly from the definition of regular conditional distribution. □

Exercise 3. Let X and Y be real random variables defined on a probability space (Ω, \mathcal{F}, P) , and let $\mathcal{G} = \sigma(Y)$ be the σ -algebra generated by Y . Show that there is a Markov kernel $\mu_Y(dx)$ from $(\mathbb{R}, \mathcal{B})$ to $(\mathbb{R}, \mathcal{B})$ such that a regular conditional distribution for X given \mathcal{G} is given by

$$\mu_Y(B) = P(X \in B | \mathcal{G}) \quad \text{a.s.}$$

□

As an example of the use of regular conditional distributions, consider the Jensen inequality (9)

$$E(\varphi(X) | \mathcal{G}) \geq \varphi(E(X | \mathcal{G})).$$

The Jensen inequality for ordinary expectations is easy: see the proof of Proposition 2 above. Now let $\mu_\omega(dx)$ be a regular conditional distribution for X given \mathcal{G} . Then

$$\begin{aligned} E(\varphi(X) | \mathcal{G})(\omega) &= \int \varphi(x) \mu_\omega(dx) \quad \text{a.s.} \\ &\geq \varphi\left(\int x \mu_\omega(dx)\right) \\ &= \varphi(E(X | \mathcal{G})) \quad \text{a.s.} \end{aligned}$$

The inequality in the middle line follows from the Jensen inequality for ordinary expectations, since each μ_ω is a Borel probability.

5. BOREL SPACES

It is not unusual – especially in the theory of Markov processes, and also in decision theory – that a random variable X of interest will take its values not in \mathbb{R} but in some other space of objects. These might be *paths*, or *configurations*, or (in decision problems) *actions*. The argument given in section 4 to establish the existence of regular conditional distributions used in an essential way the hypothesis that the random variable X was real-valued. Are there regular conditional distributions for random variables that do not take real values? In this section we will show that if X takes values in a *Borel space* then rclds always exist.

Definition 5. A *Borel space*¹ is a measurable space $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$ such that $\mathcal{F}_{\mathcal{X}}$ is *countably generated*, that is, such that for some countable collection $\{F_n\}_{n \geq 1}$,

$$(23) \quad \mathcal{F}_{\mathcal{X}} = \sigma(\{F_n\}_{n \geq 1})$$

is the minimal σ -algebra containing each set F_n .

Example 1. The space $(\mathbb{R}^d, \mathcal{B}_d)$, where \mathcal{B}_d is the σ -algebra of d -dimensional Borel sets, is a Borel space, because it is generated by the rectangles $\times_{i \leq d} (a_i, b_i]$.

Example 2. If $(\mathcal{X}_n, \mathcal{F}_n)$ is a sequence of Borel spaces then the Cartesian product $(\times_{n=1}^{\infty} \mathcal{X}_n, \times_{n=1}^{\infty} \mathcal{F}_n)$ is a Borel space. Here $\times_{n=1}^{\infty} \mathcal{F}_n$ is the product σ -algebra, that is, the minimal σ -algebra containing all *finite-dimensional* cylinder sets. (A generating set is the collection of all cylinders $F_1 \times F_2 \times \cdots \times F_n$, where F_i is taken from the countable set of generators of \mathcal{F}_n .)

Thus, in particular, the sequence space \mathbb{R}^{∞} with the usual Borel field is a Borel space. This is perhaps the most important example in probability theory.

Example 3. A *Polish space* is a complete, separable metric space, together with its Borel field. (The Borel field of a metric space is the minimal σ -algebra containing all open balls of rational radii.) Every Polish space is a Borel space. (Exercise!)

A particular case of interest is the space $C[0, 1]$ of continuous paths parametrized by time $t \in [0, 1]$. The metric on $C[0, 1]$ is induced by the supremum norm. That $C[0, 1]$ is separable in the sup norm metric follows because the set of *polygonal paths* with rational endpoints is dense.

Proposition 4. *If $(\mathcal{X}, \mathcal{F})$ is a Borel space then there is a real-valued random variable $Y : \mathcal{X} \rightarrow \mathbb{R}$ such that $\mathcal{F} = \sigma(Y)$.*

Proof. In brief, the idea is *coding by indicators*. If \mathcal{F} is countably generated then it is generated by a countable collection $\{F_n\}_{n \geq 1}$. Use the indicator functions of these sets F_n to construct a mapping to the Cantor set:

$$Y = \sum_{n=1}^{\infty} \frac{2\mathbf{1}_{F_n}}{3^n}.$$

This mapping is clearly measurable, and the value of Y uniquely determines the value of each indicator $\mathbf{1}_{F_n}$, because the ternary expansion provides a bijective mapping between the Cantor

¹The terminology varies. In some places what I have defined to be a Borel space is called a *standard* Borel space, and in some places (see, especially, the fundamental paper of Mackey on the subject) it is also required that the σ -algebra *separates points*.

set and the set of all infinite sequences of 0s and 1s. Consider the σ -algebra consisting of all events of the form $\{Y \in B\}$, where B is a one-dimensional Borel set. This collection contains the events F_n , so it must also contain every member of \mathcal{F} , since \mathcal{F} is generated by $\{F_n\}_{n \geq 1}$. Since Y is measurable, it follows that

$$\mathcal{F} = Y^{-1}(\mathcal{B}_1).$$

□

Theorem 3. *Let (Ω, \mathcal{F}, P) be a probability space and X a random variable on Ω taking values in a Borel space $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$. Then for any σ -algebra \mathcal{G} contained in \mathcal{F} there exists a regular conditional distribution for X given \mathcal{G} .*

Proof. By Proposition 4 there is a measurable mapping $Y : \mathcal{X} \rightarrow \mathbb{R}$ such that $\mathcal{F}_{\mathcal{X}} = Y^{-1}(\mathcal{B}_1)$. By Theorem 1, there is a regular conditional distribution $\nu_{\omega}(dy)$ for the real random variable $Y \circ X$ given \mathcal{G} . Define a Markov kernel $\mu_{\omega}(dx)$ from (Ω, \mathcal{F}) to $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$ as follows: for every Borel set $B \in \mathcal{B}_1$, set

$$\mu_{\omega}(Y^{-1}(B)) = \nu_{\omega}(B).$$

Since $\mathcal{F}_{\mathcal{X}} = Y^{-1}(\mathcal{B}_1)$, this defines $\mu_{\omega}(F)$ for every $F \in \mathcal{F}_{\mathcal{X}}$. That $\mu_{\omega}(dx)$ is a Markov kernel follows directly from the fact that $\nu_{\omega}(dy)$ is a Markov kernel. That $\mu_{\omega}(dx)$ is a regular conditional distribution for X given \mathcal{G} can be seen as follows. For each $F \in \mathcal{F}_{\mathcal{X}}$ there is a Borel set B such that $F = Y^{-1}(B)$, and so

$$\begin{aligned} P(X \in F | \mathcal{G})(\omega) &= P(Y \circ X \in B | \mathcal{G})(\omega) \\ &= \nu_{\omega}(B) \\ &= \mu_{\omega}(F). \end{aligned}$$

□

Exercise 4. Use the existence of regular conditional distributions for random variables valued in $(\mathbb{R}^{\infty}, \mathcal{B}_{\infty})$ to deduce the dominated convergence theorem for conditional expectations from the dominated convergence theorem for ordinary expectations. NOTE: You may also need the disintegration formula of Theorem 2.