

# Conditional expectation

Jason Swanson

April 17, 2009

## 1 Conditioning on $\sigma$ -algebras

Let  $(\Omega, \mathcal{F}, P)$  be a probability space and let  $A \in \mathcal{F}$  with  $P(A) > 0$ . Define

$$Q(B) = P(B | A) = \frac{P(B \cap A)}{P(A)}, \quad \text{for all } B \in \mathcal{F}.$$

It is easy to check that  $Q$  is a probability measure on  $(\Omega, \mathcal{F})$ . If  $X$  is a random variable, we define the **conditional expectation of  $X$  given  $A$**  as

$$E[X | A] = \int X dQ, \tag{1.1}$$

whenever this integral is well-defined. Note that  $E[1_B | A] = P(B | A)$ .

**Theorem 1.1.** *If  $E[|X|1_A] < \infty$ , then  $X$  is  $Q$ -integrable. If  $X \geq 0$  or  $E[|X|1_A] < \infty$ , then*

$$E[X | A] = \frac{E[X1_A]}{P(A)}. \tag{1.2}$$

**Remark 1.2.** Note that (1.2) may be written as

$$E[X | A] = \frac{\alpha(A)}{P(A)}, \tag{1.3}$$

where  $d\alpha = X dP$ . Also note that (1.2) gives us the formula  $E[X1_A] = P(A)E[X | A]$ . If  $X = 1_B$ , then this reduces to the familiar multiplication rule,  $P(A \cap B) = P(A)P(B | A)$ .

**Proof of Theorem 1.1.** Note that if  $P(B) = 0$ , then  $Q(B) = 0$ . Hence  $Q \ll P$ . Also note that

$$Q(B) = \int_B \frac{1_A}{P(A)} dP, \quad \text{for all } B \in \mathcal{F}.$$

Thus,  $dQ/dP = 1_A/P(A)$ . It follows that if  $X \geq 0$ , then

$$E[X | A] = \int X dQ = \int X \frac{dQ}{dP} dP = E \left[ X \frac{1_A}{P(A)} \right] = \frac{E[X1_A]}{P(A)}.$$

Therefore, if  $E[|X|1_A] < \infty$ , then  $X$  is  $Q$ -integrable, and the same formula holds.  $\square$

**Lemma 1.3.** Any finite  $\sigma$ -algebra  $\mathcal{F}$  on  $\Omega$  can be written as  $\mathcal{F} = \sigma(\{A_j\}_{j=1}^n)$ , where  $\{A_j\}_{j=1}^n$  is a partition of  $\Omega$ , that is,  $\Omega = \dot{\bigcup}_{j=1}^n A_j$ . Moreover, the partition  $\{A_j\}_{j=1}^n$  is unique.

**Proof.** For each  $\omega \in \Omega$ , let  $A_\omega$  be the smallest measurable set containing  $\omega$ . That is,  $A_\omega = \bigcap \mathcal{F}_\omega$ , where  $\mathcal{F}_\omega = \{A \in \mathcal{F} : \omega \in A\}$ . Since this is a finite intersection,  $A_\omega \in \mathcal{F}$ . In particular,  $\mathcal{E} = \{A_\omega : \omega \in \Omega\}$  is a finite set. We claim that  $\mathcal{E}$  is a partition of  $\Omega$  and that  $\mathcal{F} = \sigma(\mathcal{E})$ .

Clearly,  $\Omega = \bigcup_{\omega \in \Omega} A_\omega$ , so to show that  $\mathcal{E}$  is a partition, it suffices to show that this is a disjoint union. More specifically, we wish to show that if  $\omega, \omega' \in \Omega$ , then either  $A_\omega = A_{\omega'}$  or  $A_\omega \cap A_{\omega'} = \emptyset$ . Let  $\omega, \omega' \in \Omega$ . Note that for any  $A \in \mathcal{F}$ , if  $\omega \in A$ , then  $A \in \mathcal{F}_\omega$ , which implies  $A_\omega \subset A$ . Hence, if  $\omega \in A_{\omega'}$ , then  $A_\omega \subset A_{\omega'}$ ; and if  $\omega \in A_{\omega'}^c$ , then  $A_\omega \subset A_{\omega'}^c$ . That is, either  $A_\omega \subset A_{\omega'}$  or  $A_\omega \subset A_{\omega'}^c$ . By symmetry, either  $A_{\omega'} \subset A_\omega$  or  $A_{\omega'} \subset A_\omega^c$ . Taken together, this shows that either  $A_\omega = A_{\omega'}$  or  $A_\omega \cap A_{\omega'} = \emptyset$ .

To see that  $\mathcal{F} = \sigma(\mathcal{E})$ , simply note that any  $A \in \mathcal{F}$  can be written as  $A = \bigcup_{\omega \in A} A_\omega$ , and that this is a finite union.

For uniqueness, suppose that  $\mathcal{F} = \sigma(\{B_j\}_{j=1}^n)$ , where  $\Omega = \dot{\bigcup}_{j=1}^n B_j$ . If  $\omega \in B_j$ , then  $A_\omega = B_j$ . Therefore,  $\mathcal{E} = \{B_j\}_{j=1}^n$ .  $\square$

**Exercise 1.4.** Show that every infinite  $\sigma$ -algebra is uncountable.

Let  $(\Omega, \mathcal{F}, P)$  be a probability space and  $X$  an integrable random variable. Let  $\mathcal{G} \subset \mathcal{F}$  be a finite  $\sigma$ -algebra. Write  $\mathcal{G} = \sigma(\{A_j\}_{j=1}^n)$ , where  $\{A_j\}_{j=1}^n$  is a partition of  $\Omega$ . The **conditional expectation of  $X$  given  $\mathcal{G}$** , written  $E[X | \mathcal{G}]$ , is a random variable defined by

$$E[X | \mathcal{G}](\omega) = \begin{cases} E[X | A_1] & \text{if } \omega \in A_1, \\ E[X | A_2] & \text{if } \omega \in A_2, \\ \vdots & \\ E[X | A_n] & \text{if } \omega \in A_n. \end{cases}$$

Note that we may write

$$E[X | \mathcal{G}] = \sum_{j=1}^n E[X | A_j] 1_{A_j}. \quad (1.4)$$

We also define the **conditional probability of  $A$  given  $\mathcal{G}$**  as  $P(A | \mathcal{G}) = E[1_A | \mathcal{G}]$ . Note that  $P(A | \mathcal{G})(\omega) = P(A | A_j)$  if  $\omega \in A_j$ .

To illustrate this concept, consider the following heuristic example. Imagine my friend is at the local bar, and is about to throw a dart at a dartboard. If I model the dart board by a unit circle, which I call  $\Omega$ , then his dart will land at some point  $\omega \in \Omega$ .

Unfortunately, I am not there with him and will not be able to observe the exact location of  $\omega$ . But after he throws the dart, he is going to call me on the phone and tell me what his score for that shot was. This information will not be enough for me to determine  $\omega$ . It will, however, narrow it down. Before I receive his call, I can partition the dartboard  $\Omega$  into several pieces,  $A_1, \dots, A_n$ , with each piece corresponding to a unique score. Once he calls me, I will know which piece contains his dart.

Let  $X$  be the distance from his dart to the bullseye. Suppose he calls me and I determine that his dart is somewhere inside  $A_j$ . I can then compute  $E[X | A_j]$ . However, before

he calls, I can get prepared by computing  $E[X | A_j]$  for all  $j$ , and then encoding all this information into the single random variable  $E[X | \mathcal{G}]$ .

In probability theory, we model information by  $\sigma$ -algebras. In this example, the  $\sigma$ -algebra  $\mathcal{G}$  generated by the partition  $\{A_j\}$  models the information I will receive from my friend's phone call. Imagine that while I am waiting for my friend's phone call, an interviewer starts asking me questions. For various events  $A$ , the interviewer asks me, "After your friend calls, will you know with certainty whether or not  $A$  has occurred?" Depending on the event  $A$ , I will have to answer either yes, no, or maybe. The events  $A \in \mathcal{G}$  are precisely those events for which I can answer yes.

**Theorem 1.5.** *Let  $(\Omega, \mathcal{F}, P)$  be a probability space and  $X$  an integrable random variable. Let  $\mathcal{G} \subset \mathcal{F}$  be a finite  $\sigma$ -algebra, and let  $d\alpha = X dP$ . Then*

$$E[X | \mathcal{G}] = \frac{d(\alpha|_{\mathcal{G}})}{d(P|_{\mathcal{G}})}. \quad (1.5)$$

Or equivalently,  $Z = E[X | \mathcal{G}]$  is the unique random variable such that

- (i)  $Z$  is  $\mathcal{G}$ -measurable, and
- (ii)  $E[X1_A] = E[Z1_A]$ , for all  $A \in \mathcal{G}$ .

**Remark 1.6.** First, note the similarity between (1.5) and (1.3). Second, it should be mentioned that  $Z = E[X | \mathcal{G}]$  is unique in the following sense. If  $\tilde{Z}$  is another  $\mathcal{G}$ -measurable random variable such that  $E[X1_A] = E[\tilde{Z}1_A]$  for all  $A \in \mathcal{G}$ , then  $Z = \tilde{Z}$  a.s.

**Proof of Theorem 1.5.** Let  $Z = E[X | \mathcal{G}]$ . To prove (1.5), we must show that  $Z$  is  $\mathcal{G}$ -measurable and that

$$(\alpha|_{\mathcal{G}})(A) = \int_A Z d(P|_{\mathcal{G}}), \quad \text{for all } A \in \mathcal{G}.$$

If  $A \in \mathcal{G}$ , then  $(\alpha|_{\mathcal{G}})(A) = \alpha(A) = E[X1_A]$ ; and if  $Z$  is  $\mathcal{G}$ -measurable, then  $\int_A Z d(P|_{\mathcal{G}}) = \int_A Z dP = E[Z1_A]$ . It therefore follows that (1.5) is equivalent to (i) and (ii). Uniqueness follows from the uniqueness for the Radon-Nikodym derivative.

Write  $\mathcal{G} = \sigma(\{A_j\}_{j=1}^n)$ , where  $\{A_j\}_{j=1}^n$  is a partition of  $\Omega$ . By (1.4), it is clear that  $Z$  is  $\mathcal{G}$ -measurable, so we need only verify (ii). By the linearity of the expected value, it will suffice to show that  $E[X1_{A_j}] = E[Z1_{A_j}]$  for all  $j$ . Using (1.4), we have

$$\begin{aligned} E[Z1_{A_j}] &= E\left[\left(\sum_{i=1}^n E[X | A_i]1_{A_i}\right)1_{A_j}\right] = E\left[\sum_{i=1}^n E[X | A_i]1_{A_i \cap A_j}\right] \\ &= E[E[X | A_j]1_{A_j}] = E[X | A_j]P(A_j) = E[X1_{A_j}], \end{aligned}$$

where the last equality follows from (1.2). □

We now extend the definition of  $E[X | \mathcal{G}]$  to infinite  $\sigma$ -algebras. Let  $(\Omega, \mathcal{F}, P)$  be a probability space and  $X$  an integrable random variable. Let  $\mathcal{G} \subset \mathcal{F}$  be a  $\sigma$ -algebra. The **conditional expectation of  $X$  given  $\mathcal{G}$**  is the unique random variable  $Z = E[X | \mathcal{G}]$  satisfying (i) and (ii) of Theorem 1.5. Equivalently,  $E[X | \mathcal{G}]$  can be defined by (1.5), where  $d\alpha = X dP$ . We also define the **conditional probability of  $A$  given  $\mathcal{G}$**  as  $P(A | \mathcal{G}) = E[1_A | \mathcal{G}]$ .

## 2 Conditioning on random variables

Let  $(\Omega, \mathcal{F}, P)$  be a probability space and  $X$  an integrable random variable. Let  $Y$  be a discrete random variable taking values in a finite set  $S = \{k_1, \dots, k_n\}$ , and assume that  $P(Y = k_j) > 0$  for all  $j$ . Define  $f(k) = E[X | Y = k]$  for all  $k \in S$ . Since  $\{Y = k\}$  is an event with positive probability, this is the simplest kind of conditional expectation we have used so far, namely, the one defined in (1.1).

We now define the **conditional expectation of  $X$  given  $Y$**  as  $E[X | Y] = f(Y)$ . In other words,  $E[X | Y]$  is a random variable whose outcome depends on  $Y$ . Specifically, for all  $\omega \in \{Y = k\}$ , we have  $E[X | Y](\omega) = f(k)$ .

**Theorem 2.1.** *Under the above definitions,  $E[X | Y] = E[X | \sigma(Y)]$ .*

**Proof.** Under the above definitions,  $E[X | Y] = f(Y)$ . In order to show that  $f(Y) = E[X | \sigma(Y)]$ , we must show that (i)  $f(Y)$  is  $\sigma(Y)$ -measurable, and (ii)  $E[X1_A] = E[f(Y)1_A]$  for all  $A \in \sigma(Y)$ .

Recall that a random variable  $Z$  is  $\sigma(Y)$ -measurable if and only if  $Z = g(Y)$  for some measurable function  $g$ . Hence, (i) clearly holds. For (ii), let  $A \in \sigma(Y)$  be arbitrary. Then  $A = \{Y \in B\}$  for some  $B \in \mathcal{R}$ . This implies there is some subset  $C \subset S$  such that  $A = \bigcup_{k \in C} \{Y = k\}$ . By the linearity of the expected value, we may assume without loss of generality that  $A = \{Y = k\}$ . In this case,

$$\begin{aligned} E[f(Y)1_A] &= E[f(Y)1_{\{Y=k\}}] = E[f(k)1_{\{Y=k\}}] = f(k)P(Y = k) \\ &= E[X|Y = k]P(Y = k) = E[X1_{\{Y=k\}}] = E[X1_A]. \end{aligned}$$

Note that we have again used (1.2). □

We now extend the definition of  $E[X | Y]$  to arbitrary random variables. Let  $(\Omega, \mathcal{F}, P)$  be a probability space and  $X$  an integrable random variable. Let  $Y$  be an arbitrary random variable. The **conditional expectation of  $X$  given  $Y$**  is  $E[X | Y] = E[X | \sigma(Y)]$ . We also define the **conditional probability of  $A$  given  $Y$**  as  $P(A | Y) = E[1_A | Y]$ . Note that since  $E[X | Y]$  is  $\sigma(Y)$ -measurable, there exists a measurable function  $g$  (which depends on  $X$ ) such that  $E[X | Y] = g(Y)$ .

**Example 2.2.** Suppose  $X$  and  $Y$  are random variables with a joint density function  $f(x, y)$ . Note that  $P(X \in A | Y = k)$  is undefined, since  $P(Y = k) = 0$ . Nonetheless, it should be intuitively clear that the following equation ought to hold:

$$P(X \in A | Y = k) = \frac{\int_A f(x, k) dx}{\int_{\mathbb{R}} f(x, k) dx}. \quad (2.1)$$

The integral in the denominator is necessary in order to make the function  $x \mapsto f(x, k)$  a probability density function. In this example, we will explore the sense in which this formula is rigorously valid. (In an undergraduate class, it may be rigorously valid *by definition*. But for us, as usual, it is a special case of something more general.)

**Theorem 2.3.** Let  $X$  and  $Y$  have joint density  $f(x, y)$ . Let  $g$  be a measurable function such that  $E|g(X)| < \infty$ . Define

$$h(k) = \frac{\int_{\mathbb{R}} g(x)f(x, k) dx}{\int_{\mathbb{R}} f(x, k) dx},$$

whenever  $\int_{\mathbb{R}} f(x, k) dx > 0$ , and  $h(k) = 0$  otherwise. Then  $E[g(X) | Y] = h(Y)$ .

**Remark 2.4.** If  $g(x) = 1_A(x)$ , then  $h(k)$  agrees with the right-hand side of (2.1). Also, as can be seen from the proof below, we could have defined  $h(k)$  arbitrarily when  $\int_{\mathbb{R}} f(x, k) dx = 0$ .

**Proof of Theorem 2.3.** Since  $h(Y)$  is  $\sigma(Y)$ -measurable, we only need to show that  $E[h(Y)1_A] = E[g(X)1_A]$  for all  $A \in \sigma(Y)$ . Let  $A \in \sigma(Y)$ , so that  $A = \{Y \in B\}$  for some  $B \in \mathcal{R}$ . We now have

$$\begin{aligned} E[h(Y)1_A] &= E[h(Y)1_B(Y)] = \int_B \int_{\mathbb{R}} h(y)f(x, y) dx dy \\ &= \int_B \left( h(y) \int_{\mathbb{R}} f(x, y) dx \right) dy = \int_{B \cap C} \left( h(y) \int_{\mathbb{R}} f(x, y) dx \right) dy, \end{aligned}$$

where  $C = \{y : \int_{\mathbb{R}} f(x, y) dx > 0\}$ . Note that for all  $y \in C$ , we have  $h(y) \int_{\mathbb{R}} f(x, y) dx = \int_{\mathbb{R}} g(x)f(x, y) dx$ . Also, for all  $y \in C^c$ , we have  $f(x, y) = 0$  for Lebesgue almost every  $x$ . Thus,  $y \in C^c$  implies  $\int_{\mathbb{R}} g(x)f(x, y) dx = 0$ . It therefore follows that

$$\begin{aligned} E[h(Y)1_A] &= \int_{B \cap C} \int_{\mathbb{R}} g(x)f(x, y) dx dy = \int_B \int_{\mathbb{R}} g(x)f(x, y) dx dy \\ &= E[g(X)1_B(Y)] = E[g(X)1_A], \end{aligned}$$

which was what we needed to prove.  $\square$

In general, we interpret  $E[X | Y = y]$  to mean  $g(y)$ , where  $g$  is a measurable function such that  $E[X | Y] = g(Y)$ . Some caution is needed in these cases, though, since such a function  $g$  is only defined  $\mu_Y$ -a.s., where  $\mu_Y$  is the distribution of  $Y$ .

## 3 Properties of conditional expectation

### 3.1 Basic properties

Let  $(\Omega, \mathcal{F}, P)$  be a probability space,  $X$  an integrable random variable, and  $\mathcal{G} \subset \mathcal{F}$ .

**Proposition 3.1.**  $E[E[X | \mathcal{G}]] = EX$ .

**Proof.** By the definition of conditional expectation, we have  $E[E[X | \mathcal{G}]1_A] = E[X1_A]$  for all  $A \in \mathcal{G}$ . Take  $A = \Omega$ .  $\square$

**Proposition 3.2.** If  $X$  is  $\mathcal{G}$ -measurable, then  $E[X | \mathcal{G}] = X$ .

**Proof.** It follows trivially, since  $X$  is  $\mathcal{G}$ -measurable and  $E[X1_A] = E[X1_A]$  for all  $A \in \mathcal{G}$ .  $\square$

Recall that  $X$  and  $\mathcal{G}$  are independent if  $\sigma(X)$  and  $\mathcal{G}$  are independent, which in turn means that  $P(A \cap B) = P(A)P(B)$  whenever  $A \in \sigma(X)$  and  $B \in \mathcal{G}$ . Hence,  $X$  and  $\mathcal{G}$  are independent if and only if  $P(\{X \in C\} \cap B) = P(X \in C)P(B)$  for all  $C \in \mathcal{R}$  and  $B \in \mathcal{G}$ .

**Proposition 3.3.** *If  $X$  and  $\mathcal{G}$  are independent, then  $E[X | \mathcal{G}] = E[X]$ . In particular,  $E[X | \{\emptyset, \Omega\}] = E[X]$ .*

**Proof.** A constant random variable is measurable with respect to every  $\sigma$ -algebra, so  $E[X]$  is trivially  $\mathcal{G}$ -measurable. Also, for all  $A \in \mathcal{G}$ , we have  $E[X1_A] = E[X]E[1_A] = E[E[X]1_A]$ . The final claim holds since every random variable is independent of the trivial  $\sigma$ -algebra.  $\square$

**Theorem 3.4.** *If  $\mathcal{G}_1 \subset \mathcal{G}_2 \subset \mathcal{F}$ , then  $E[E[X | \mathcal{G}_1] | \mathcal{G}_2] = E[E[X | \mathcal{G}_2] | \mathcal{G}_1] = E[X | \mathcal{G}_1]$ .*

**Remark 3.5.** In words, this says that in a battle between nested  $\sigma$ -algebras, the smallest  $\sigma$ -algebra always wins.

**Proof of Theorem 3.4.** Since  $\mathcal{G}_1 \subset \mathcal{G}_2$  and  $E[X | \mathcal{G}_1]$  is  $\mathcal{G}_1$ -measurable, it is also  $\mathcal{G}_2$ -measurable. Hence, by Proposition 3.2,  $E[E[X | \mathcal{G}_1] | \mathcal{G}_2] = E[X | \mathcal{G}_1]$ . The other equality holds since  $E[X | \mathcal{G}_1]$  is  $\mathcal{G}_1$ -measurable and, for all  $A \in \mathcal{G}_1 \subset \mathcal{G}_2$ , we have  $E[E[X | \mathcal{G}_1]1_A] = E[X1_A] = E[E[X | \mathcal{G}_2]1_A]$ .  $\square$

**Theorem 3.6.** *If  $Y$  and  $XY$  are integrable, and  $X$  is  $\mathcal{G}$ -measurable, then*

$$E[XY | \mathcal{G}] = XE[Y | \mathcal{G}] \text{ a.s.}$$

**Proof.** Let  $d\alpha = Y dP$ , so that  $E[Y | \mathcal{G}] = d(\alpha|_{\mathcal{G}})/d(P|_{\mathcal{G}})$ . Let  $d\beta_1 = X d\alpha_+$ ,  $d\beta_2 = X d\alpha_-$ , and  $\beta = \beta_1 - \beta_2$ , so that  $d\beta = XY dP$  and  $E[XY | \mathcal{G}] = d(\beta|_{\mathcal{G}})/d(P|_{\mathcal{G}})$ . Since  $X$  is  $\mathcal{G}$ -measurable, we have  $d(\beta_1|_{\mathcal{G}})/d(\alpha_+|_{\mathcal{G}}) = d(\beta_2|_{\mathcal{G}})/d(\alpha_-|_{\mathcal{G}}) = X$ . Hence,

$$\begin{aligned} E[XY | \mathcal{G}] &= \frac{d(\beta|_{\mathcal{G}})}{d(P|_{\mathcal{G}})} = \frac{d(\beta_1|_{\mathcal{G}})}{d(P|_{\mathcal{G}})} - \frac{d(\beta_2|_{\mathcal{G}})}{d(P|_{\mathcal{G}})} \\ &= \frac{d(\beta_1|_{\mathcal{G}})}{d(\alpha_+|_{\mathcal{G}})} \cdot \frac{d(\alpha_+|_{\mathcal{G}})}{d(P|_{\mathcal{G}})} - \frac{d(\beta_2|_{\mathcal{G}})}{d(\alpha_-|_{\mathcal{G}})} \cdot \frac{d(\alpha_-|_{\mathcal{G}})}{d(P|_{\mathcal{G}})} \\ &= X \left( \frac{d(\alpha_+|_{\mathcal{G}})}{d(P|_{\mathcal{G}})} - \frac{d(\alpha_-|_{\mathcal{G}})}{d(P|_{\mathcal{G}})} \right) = X \frac{d(\alpha|_{\mathcal{G}})}{d(P|_{\mathcal{G}})} = XE[Y | \mathcal{G}], \text{ } P\text{-a.s.}, \end{aligned}$$

and we are done.  $\square$

**Theorem 3.7. (linearity)**  $E[aX + Y | \mathcal{G}] = aE[X | \mathcal{G}] + E[Y | \mathcal{G}]$ .

**Proof.** The right-hand side is clearly  $\mathcal{G}$ -measurable. Let  $A \in \mathcal{G}$ . Then

$$\begin{aligned} E[(aE[X | \mathcal{G}] + E[Y | \mathcal{G}])1_A] &= aE[E[X | \mathcal{G}]1_A] + E[E[Y | \mathcal{G}]1_A] \\ &= aE[X1_A] + E[Y1_A] = E[(aX + Y)1_A], \end{aligned}$$

and we are done.  $\square$

**Lemma 3.8.** *Suppose  $U$  and  $V$  are  $\mathcal{H}$ -measurable random variables. If  $E[U1_A] \leq E[V1_A]$  for all  $A \in \mathcal{H}$ , then  $U \leq V$  a.s. If  $E[U1_A] = E[V1_A]$  for all  $A \in \mathcal{H}$ , then  $U = V$  a.s.*

**Proof.** By reversing the roles of  $U$  and  $V$ , the second claim follows from the first. To prove the first, suppose  $E[U1_A] \leq E[V1_A]$  for all  $A \in \mathcal{H}$ . Let  $A = \{U > V\} \in \mathcal{H}$  and define  $Z = (U - V)1_A$ , so that  $Z \geq 0$ . Note that  $EZ = E[U1_A] - E[V1_A] \leq 0$ . Hence,  $EZ = 0$ , so  $Z = 0$  a.s., which implies  $P(A) = 0$ .  $\square$

**Theorem 3.9. (monotonicity)** *If  $X \leq Y$  a.s., then  $E[X | \mathcal{G}] \leq E[Y | \mathcal{G}]$  a.s.*

**Proof.** For all  $A \in \mathcal{G}$ , we have  $E[E[X | \mathcal{G}]1_A] = E[X1_A] \leq E[Y1_A] = E[E[Y | \mathcal{G}]1_A]$ . Hence, by Lemma 3.8,  $E[X | \mathcal{G}] \leq E[Y | \mathcal{G}]$  a.s.  $\square$

**Theorem 3.10.** *Suppose  $X$  and  $Y$  are independent and  $\varphi$  is a measurable function such that  $E|\varphi(X, Y)| < \infty$ , then  $E[\varphi(X, Y) | X] = g(X)$ , where  $g(x) = E[\varphi(x, Y)]$ .*

**Remark 3.11.** It is important here that  $X$  and  $Y$  are independent. This result is not true when  $X$  and  $Y$  are dependent.

**Proof of Theorem 3.10.** Clearly,  $g(X)$  is  $\sigma(X)$ -measurable. Let  $A \in \mathcal{R}$ . Then

$$\begin{aligned} E[\varphi(X, Y)1_{\{X \in A\}}] &= \int \int \varphi(x, y)1_A(x)\mu_Y(dy)\mu_X(dx) \\ &= \int 1_A(x) \left( \int \varphi(x, y)\mu_Y(dy) \right) \mu_X(dx) = \int 1_A(x)g(x)\mu_X(dx) = E[g(X)1_{\{X \in A\}}], \end{aligned}$$

and we are done.  $\square$

**Example 3.12.** Let  $X, Y, Z$  be iid, uniformly distributed on  $(0, 1)$ . We shall compute the distribution of  $(XY)^Z$ . We begin by computing the distribution of  $W = XY$ . Let  $w \in (0, 1)$ . Then

$$P(W \leq w) = P(XY \leq w) = E[1_{\{XY \leq w\}}] = E[E[1_{\{XY \leq w\}} | X]].$$

By Theorem 3.10,  $E[1_{\{XY \leq w\}} | X] = f(X)$ , where

$$f(x) = E[1_{\{xY \leq w\}}] = P(xY \leq w) = P\left(Y \leq \frac{w}{x}\right) = 1_{\{x < w\}} + \frac{w}{x}1_{\{x \geq w\}}.$$

Thus,

$$P(W \leq w) = E[f(X)] = E\left[1_{\{X < w\}} + \frac{w}{X}1_{\{X \geq w\}}\right] = w + \int_w^1 \frac{w}{x} dx = w - w \log w.$$

Differentiating, we find that  $W$  has density  $f_W(w) = (-\log w)1_{(0,1)}(w)$ .

Similarly, for  $x \in (0, 1)$ , we now compute

$$P((XY)^Z \leq x) = E[P(W^Z \leq x | W)] = E[g(W)],$$

where

$$g(w) = P(w^Z \leq x) = P\left(Z \geq \frac{\log x}{\log w}\right) = \left(1 - \frac{\log x}{\log w}\right)1_{\{w \leq x\}}.$$

Thus,

$$\begin{aligned} P((XY)^Z \leq x) &= E \left[ \left( 1 - \frac{\log x}{\log W} \right) 1_{\{W \leq x\}} \right] = \int_0^x \left( 1 - \frac{\log x}{\log w} \right) (-\log w) dw \\ &= - \int_0^x \log w dw + x \log x = x. \end{aligned}$$

In other words,  $(XY)^Z$  is uniformly distributed on  $(0, 1)$ .

## 3.2 Limit theorems and inequalities

**Theorem 3.13. (monotone convergence)** *If  $0 \leq X_n \uparrow X$  a.s. and  $X$  is integrable, then  $E[X_n | \mathcal{G}] \uparrow E[X | \mathcal{G}]$  a.s.*

**Proof.** By monotonicity, there exists a  $\mathcal{G}$ -measurable random variable  $Z$  such that  $E[X_n | \mathcal{G}] \uparrow Z$  a.s. Let  $A \in \mathcal{G}$ . Using monotone convergence,

$$E[Z 1_A] = \lim_{n \rightarrow \infty} E[E[X_n | \mathcal{G}] 1_A] = \lim_{n \rightarrow \infty} E[X_n 1_A] = E[X 1_A],$$

which shows  $Z = E[X | \mathcal{G}]$ . □

**Theorem 3.14. (Fatou's lemma)** *If  $X_n \geq 0$  a.s., each  $X_n$  is integrable, and  $\liminf_{n \rightarrow \infty} X_n$  is integrable, then*

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

**Proof.** Let  $\bar{X}_n = \inf_{j \geq n} X_j$  and  $X = \liminf_{n \rightarrow \infty} X_n$ . Note that  $0 \leq \bar{X}_n \uparrow X$ . In particular,  $\bar{X}_n$  is integrable. For each  $j \geq n$ , we have  $\bar{X}_n \leq X_j$  a.s. Hence, by monotonicity,  $E[\bar{X}_n | \mathcal{G}] \leq E[X_j | \mathcal{G}]$  a.s. It follows that

$$E[\bar{X}_n | \mathcal{G}] \leq \inf_{j \geq n} E[X_j | \mathcal{G}] \text{ a.s.}$$

Monotone convergence implies

$$E[X | \mathcal{G}] = \lim_{n \rightarrow \infty} E[\bar{X}_n | \mathcal{G}] \leq \liminf_{n \rightarrow \infty} E[X_n | \mathcal{G}] \text{ a.s.,}$$

and we are done. □

**Theorem 3.15. (dominated convergence)** *Let  $X_n$  be random variables with  $X_n \rightarrow X$  a.s. Suppose there exists an integrable random variable  $Y$  such that  $|X_n| \leq Y$  a.s. for all  $n$ . Then*

$$\lim_{n \rightarrow \infty} E[X_n | \mathcal{G}] = E[X | \mathcal{G}] \text{ a.s.}$$

**Exercise 3.16.** Prove Theorem 3.15 by mimicking the proof of the ordinary dominated convergence theorem.

**Exercise 3.17.** Show that if  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  is convex, then the left-hand derivative,

$$\varphi'_-(c) = \lim_{h \downarrow 0} \frac{\varphi(c) - \varphi(c-h)}{h}$$

exists for all  $c$ . Moreover,

$$\varphi(x) - \varphi(c) - (x-c)\varphi'_-(c) \geq 0, \quad (3.1)$$

for all  $x$  and  $c$ .

**Theorem 3.18. (Jensen's inequality)** *If  $\varphi$  is convex and  $X$  and  $\varphi(X)$  are integrable, then  $\varphi(E[X | \mathcal{G}]) \leq E[\varphi(X) | \mathcal{G}]$ .*

**Proof.** Let  $Z = (X - E[X | \mathcal{G}])\varphi'_-(E[X | \mathcal{G}])$ , so that by (3.1),  $\varphi(X) - \varphi(E[X | \mathcal{G}]) - Z \geq 0$ , which implies

$$0 \leq E[\varphi(X) - \varphi(E[X | \mathcal{G}]) - Z | \mathcal{G}] = E[\varphi(X) | \mathcal{G}] - \varphi(E[X | \mathcal{G}]) - E[Z | \mathcal{G}].$$

It therefore suffices to show that  $E[Z | \mathcal{G}] = 0$ . To see this, we calculate

$$\begin{aligned} E[Z | \mathcal{G}] &= E[(X - E[X | \mathcal{G}])\varphi'_-(E[X | \mathcal{G}]) | \mathcal{G}] \\ &= \varphi'_-(E[X | \mathcal{G}])E[X - E[X | \mathcal{G}] | \mathcal{G}] \\ &= \varphi'_-(E[X | \mathcal{G}])E[X | \mathcal{G}] - E[E[X | \mathcal{G}] | \mathcal{G}] \\ &= \varphi'_-(E[X | \mathcal{G}])E[X | \mathcal{G}] - E[X | \mathcal{G}] = 0, \end{aligned}$$

and we are done. □

**Theorem 3.19. (Hölder's inequality)** *Let  $p, q \in (1, \infty)$  be conjugate exponents, so that  $1/p + 1/q = 1$ . Suppose that  $|X|^p$  and  $|Y|^q$  are integrable. Then*

$$E[|XY| | \mathcal{G}] \leq (E[|X|^p | \mathcal{G}])^{1/p} (E[|Y|^q | \mathcal{G}])^{1/q} \text{ a.s.}$$

**Proof.** Note that by the ordinary Hölder's inequality,  $XY$  is integrable, so that  $E[|XY| | \mathcal{G}]$  is well-defined. Let  $U = (E[|X|^p | \mathcal{G}])^{1/p}$  and  $V = (E[|Y|^q | \mathcal{G}])^{1/q}$ . Note that both  $U$  and  $V$  are  $\mathcal{G}$ -measurable. Observe that

$$E[|X|^p 1_{\{U=0\}}] = E[E[|X|^p 1_{\{U=0\}} | \mathcal{G}]] = E[1_{\{U=0\}} E[|X|^p | \mathcal{G}]] = E[1_{\{U=0\}} U^p] = 0.$$

Hence,  $|X| 1_{\{U=0\}} = 0$  a.s., which implies

$$E[|XY| | \mathcal{G}] 1_{\{U=0\}} = E[|XY| 1_{\{U=0\}} | \mathcal{G}] = 0.$$

Similarly,  $E[|XY| | \mathcal{G}] 1_{\{V=0\}} = 0$ . It therefore suffices to show that  $E[|XY| | \mathcal{G}] 1_H \leq UV$ , where  $H = \{U > 0, V > 0\}$ . For this, we will use Lemma 3.8 to prove that

$$\frac{E[|XY| | \mathcal{G}]}{UV} 1_H \leq 1 \text{ a.s.}$$

Note that the left-hand side is defined to be zero on  $H^c$ .

Let  $A \in \mathcal{G}$  be arbitrary and define  $G = H \cap A$ . Then

$$\begin{aligned}
E \left[ \frac{E[|XY| \mid \mathcal{G}]}{UV} 1_H 1_A \right] &= E \left[ E \left[ \frac{|XY|}{UV} 1_G \mid \mathcal{G} \right] \right] \\
&= E \left[ \frac{|X|}{U} 1_G \cdot \frac{|Y|}{V} 1_G \right] \\
&\leq \left( E \left[ \frac{|X|^p}{U^p} 1_G \right] \right)^{1/p} \left( E \left[ \frac{|Y|^q}{V^q} 1_G \right] \right)^{1/q} \\
&= \left( E \left[ \frac{E[|X|^p \mid \mathcal{G}]}{U^p} 1_G \right] \right)^{1/p} \left( E \left[ \frac{E[|Y|^q \mid \mathcal{G}]}{V^q} 1_G \right] \right)^{1/q} \\
&= (E[1_G])^{1/p} (E[1_G])^{1/q} = E[1_G] \leq E[1_A].
\end{aligned}$$

Applying Lemma 3.8 finishes the proof.  $\square$

### 3.3 Minimizing the mean square error

We say that a random variable  $X$  is **square integrable** if  $E|X|^2 < \infty$ . Let  $X$  be square integrable and consider the function  $f(a) = E|X - a|^2 = a^2 - 2(EX)a + E|X|^2$ . This function has a minimum at  $a = EX$ . In other words, if we wish to approximate  $X$  by a constant, then the constant  $EX$  is the one which minimizes our mean square error.

The conditional expectation has a similar property. If we wish to approximate  $X$  by a square integrable,  $\mathcal{G}$ -measurable random variable, then  $E[X \mid \mathcal{G}]$  is the random variable which minimizes our mean square error. This is made precise in the following theorem.

**Theorem 3.20.** *Let  $(\Omega, \mathcal{F}, P)$  be a probability space and let  $X$  be square integrable. Let  $\mathcal{G} \subset \mathcal{F}$  and define  $Z = E[X \mid \mathcal{G}]$ . If  $Y$  is any square integrable,  $\mathcal{G}$ -measurable random variable, then  $E|X - Z|^2 \leq E|X - Y|^2$ .*

**Proof.** First note that by Jensen's inequality,

$$|Z|^2 = |E[X \mid \mathcal{G}]|^2 \leq E[|X|^2 \mid \mathcal{G}] \text{ a.s.}$$

Hence,  $E|Z|^2 \leq E[E[|X|^2 \mid \mathcal{G}]] = E|X|^2 < \infty$  and  $Z$  is square integrable. Let  $W = Z - Y$ . Since  $W$  is  $\mathcal{G}$ -measurable,

$$E[WZ] = E[WE[X \mid \mathcal{G}]] = E[E[WX \mid \mathcal{G}]] = E[WX].$$

Hence,  $E[W(X - Z)] = 0$ , which implies

$$\begin{aligned}
E|X - Y|^2 &= E|X - Z + W|^2 = E|X - Z|^2 + 2E[W(X - Z)] + E|W|^2 \\
&= E|X - Z|^2 + E|W|^2 \geq E|X - Z|^2,
\end{aligned}$$

and we are done.  $\square$

**Remark 3.21.** In the language of Hilbert spaces and  $L^p$  spaces, this theorem says the following:  $X$  is an element of the Hilbert space  $L^2(\Omega, \mathcal{F}, P)$ , and  $E[X \mid \mathcal{G}]$  is the orthogonal projection of  $X$  onto the subspace  $L^2(\Omega, \mathcal{G}, P)$ .

## 4 A preview of stochastic processes

A **stochastic process** is a collection of random variable  $\{X(t) : t \in T\}$  indexed by some set  $T$ . We usually think of  $T$  as time. A **discrete time stochastic process** is where  $T = \mathbb{N}$ , in which case the process is just a sequence of random variables.

Let  $\{X_n : n \in \mathbb{N}\}$  be a discrete time stochastic process. Define  $\mathcal{F}_n = \sigma(X_1, \dots, X_n)$ . The  $\sigma$ -algebra  $\mathcal{F}_n$  represents all the information at time  $n$  that we would have from observing the values  $X_1, \dots, X_n$ . Note that  $\mathcal{F}_n \subset \mathcal{F}_{n+1}$ .

More generally, a **filtration** is a sequence of  $\sigma$ -algebras  $\{\mathcal{F}_n\}_{n=1}^{\infty}$  such that  $\mathcal{F}_n \subset \mathcal{F}_{n+1}$ . A stochastic process  $\{X_n : n \in \mathbb{N}\}$  is said to be **adapted** to the filtration  $\{\mathcal{F}_n\}_{n=1}^{\infty}$  if  $X_n$  is  $\mathcal{F}_n$ -measurable for all  $n$ . The special case  $\mathcal{F}_n = \sigma(X_1, \dots, X_n)$  is called the **filtration generated by  $X$** , and is denoted by  $\{\mathcal{F}_n^X\}_{n=1}^{\infty}$ .

An important class of discrete time stochastic processes is the martingales. A stochastic process  $\{X_n : n \in \mathbb{N}\}$  is a **martingale** with respect to the filtration  $\{\mathcal{F}_n\}_{n=1}^{\infty}$  if

- (i)  $X_n$  is integrable for all  $n$ ,
- (ii)  $\{X_n : n \in \mathbb{N}\}$  is adapted to  $\{\mathcal{F}_n\}_{n=1}^{\infty}$ , and
- (iii)  $E[X_{n+1} | \mathcal{F}_n] = X_n$  for all  $n$ .

The critical item is (iii). Imagine that  $X_n$  models our cumulative wealth as we play a sequence of gambling games. Condition (iii) says that, given all the information up to time  $n$ , our expected wealth at time  $n + 1$  is the same as our wealth at time  $n$ . In other words, a martingale models a “fair” game.

Another important class of discrete time stochastic processes is the Markov chains. A stochastic process  $\{X_n : n \in \mathbb{N}\}$  is a **Markov chain** with respect to the filtration  $\{\mathcal{F}_n\}_{n=1}^{\infty}$  if

- (i)  $\{X_n : n \in \mathbb{N}\}$  is adapted to  $\{\mathcal{F}_n\}_{n=1}^{\infty}$ , and
- (ii)  $P(X_{n+1} \in B | \mathcal{F}_n) = P(X_{n+1} \in B | X_n)$  for all  $B \in \mathcal{R}$ .

Here, the critical item is (ii). It is called the **Markov property**. In words, it says that the conditional distribution of  $X_{n+1}$  given all the information up to time  $n$  is the same as if we were only given  $X_n$ . In other words, the future behavior of a Markov chain depends only on the present location of the chain, and not on how it got there.

The canonical example of a Markov chain is a random walk. If  $\{X_j\}_{j=1}^{\infty}$  are iid and  $S_n = X_1 + \dots + X_n$ , then  $\{S_n : n \in \mathbb{N}\}$  is a **random walk**. The random walk is a Markov chain with respect to the filtration generated by  $S$ . Moreover, if each  $X_j$  is integrable with mean zero, then the random walk is also a martingale.

A continuous time stochastic process has the form  $\{X(t) : t \in [0, \infty)\}$ . Examples include the Poisson process and Brownian motion. Concepts such as filtrations, adaptedness, martingales, and the Markov property can all be extended to continuous time. Care is needed however, because (for one thing) the time domain is uncountable. Brownian motion is the continuous time analog of a random walk. It is the canonical example in continuous time of both a martingale and a Markov process. It can be realized as the limit of a sequence of

random walks, where the step sizes are becoming smaller and the steps are occurring more frequently.

More specifically, let  $\{S_n : n \in \mathbb{N}\}$  be a mean zero random walk. Let  $S(t) = S_{\lfloor t \rfloor}$ , where  $\lfloor \cdot \rfloor$  denotes the greatest integer function. Then the sequence of processes

$$\left\{ \frac{S(nt)}{\sqrt{n}} : t \in [0, \infty) \right\}$$

converges (in a certain sense) as  $n \rightarrow \infty$  to a continuous time stochastic process called Brownian motion. This is the conclusion of Donsker's theorem, which is a kind of central limit theorem for stochastic processes.

Differential equations that involve Brownian motion are referred to as stochastic differential equations (SDEs). SDEs are used to model dynamical systems that involve randomness, and are very common in scientific applications. In order to understand SDEs, one must first understand the stochastic integral (with respect to Brownian motion), which behaves quite differently from the ordinary Lebesgue-Stieltjes integral. In particular, the classical fundamental theorem of calculus no longer applies when one is working with stochastic integrals. It must be replaced by a new rule called Itô's rule. Itô's rule gives rise to a whole new calculus called stochastic calculus.