

## 1 Conditional expectation given an event

Let  $Y$  be a random variable, and  $A$  an event. We assume we simulate  $n$  values  $\{y_1, y_2, \dots, y_n\}$  independently from the distribution of  $Y$ . The expectation  $E(Y)$  can be numerically estimated as the average of the simulated values (this is called a Monte Carlo estimate):

$$E(Y) = \frac{y_1 + y_2 + \dots + y_n}{n}.$$

The conditional expectation  $E(Y | A)$  is numerically approximated by considering only those values among  $\{y_1, y_2, \dots, y_n\}$  where  $A$  occurred, and taking their average:

$$E(Y | A) = \frac{\sum_{\{j: A \text{ occurred for } y_j\}} y_j}{\#\{j : A \text{ occurred for } y_j\}}.$$

**Definition** (Conditional expectation given an event). Let  $A$  be an event with positive probability,  $P(A) > 0$ . If  $Y$  is a discrete random variable, then the *conditional expectation of  $Y$  given  $A$*  is

$$E(Y | A) = \sum_y yP(Y = y | A),$$

where the sum is over the support of  $Y$ . If  $Y$  is a continuous random variable with PDF  $f$ , then

$$E(Y | A) = \int_{-\infty}^{\infty} yf(y | A) dy,$$

where the conditional PDF  $f(y | A)$  is defined as the derivative of the conditional CDF

$$F(y | A) = P(Y \leq y | A),$$

and can also be computed by a hybrid version of Bayes' rule:

$$f(y | A) = \frac{P(A | Y = y)f(y)}{P(A)}.$$

**Theorem** (Law of total expectation). Let  $A_1, A_2, \dots, A_n$  be a partition of a sample space, with  $P(A_i) > 0$  for all  $i = 1, 2, \dots, n$ . Let  $Y$  be a random variable on this sample space. Then

$$E(Y) = \sum_{i=1}^n E(Y | A_i)P(A_i).$$

The law of total probability (LOTP) is a particular case of the law of total expectation. Let  $B$  be an event, and let  $Y = I_B$  be its indicator. The law of total expectation says:

$$E(I_B) = \sum_{i=1}^n E(I_B | A_i)P(A_i).$$

But, by the fundamental bridge, we have

$$E(I_B) = P(B), \quad E(I_B | A_i) = P(B | A_i).$$

Thus we obtain LOTP:

$$P(B) = \sum_{i=1}^n P(B | A_i)P(A_i).$$

**Example:** Geometric expectation

Let  $X \sim \text{Geom}(p)$ . Then  $X$  represents the number of failures before the first successful trial in a sequence of independent Bernoulli trials, each with the same probability of success  $p \in (0, 1)$ . We denote  $q = 1 - p$ . We calculate  $E(X)$  by conditioning on the outcome of the first trial: if this outcome is a success,  $X = 0$ . Otherwise, if this outcome is a failure, we are back where we started by memorylessness (this property comes from the independence of the Bernoulli trials). Thus

$$E(X) = E(X \mid \text{outcome of first trial is success}) \cdot p + E(X \mid \text{outcome of first trial is failure}) \cdot q = 0 \cdot p + (1 + E(X))q.$$

We solve the equation  $E(X) = (1 + E(X))q$ , and obtain  $E(X) = q/p$ .

## 2 Conditional expectation given a random variable

The key to understanding the conditional expectation of a random variable  $Y$  given another random variable  $X$ , denoted by  $E(Y \mid X)$ , is first to understand the conditional expectation  $E(Y \mid X = x)$  of  $Y$  given the event  $X = x$ . If  $Y$  is discrete, we have:

$$E(Y \mid X = x) = \sum_y y P(Y = y \mid X = x).$$

If  $Y$  is continuous, we have:

$$E(Y \mid X = x) = \int_{-\infty}^{\infty} y f_{Y|X}(y \mid x) dy.$$

**Definition** (Conditional expectation given a random variable). Let  $g(x) = E(Y \mid X = x)$ . The *conditional expectation of  $Y$  given  $X$* , denoted by  $E(Y \mid X)$ , is defined to be the random variable  $g(X)$ . This random variable can have an expectation  $E(E(Y \mid X))$ , a variance  $\text{Var}(E(Y \mid X))$ , and also higher order moments.

**Example:**

Suppose we have a stick of length 1, and break the stick at a random point  $X$  chosen uniformly at random. Given that  $X = x$ , we then choose another breakpoint  $Y$  uniformly on the interval  $[0, x]$ . Find the conditional expectation  $E(Y \mid X)$ , its mean  $E(E(Y \mid X))$ , and its variance  $\text{Var}(E(Y \mid X))$ .

*Solution:* We have  $X \sim \text{Unif}(0, 1)$  and  $Y \mid X = x \sim \text{Unif}(0, x)$ . Then

$$g(x) = E(Y \mid X = x) = \frac{x}{2}.$$

Thus the conditional mean of  $Y$  given  $X$  is

$$E(Y \mid X) = g(X) = \frac{X}{2}.$$

Since  $X \sim \text{Unif}(0, 1)$ , we have  $\frac{X}{2} \sim \text{Unif}\left(0, \frac{1}{2}\right)$ . Thus the mean and the variance of the conditional expectation  $E(Y \mid X)$  are:

$$E(E(Y \mid X)) = E\left(\frac{X}{2}\right) = \frac{1}{4}, \quad \text{Var}(E(Y \mid X)) = \text{Var}\left(\frac{X}{2}\right) = \frac{1}{48}.$$

### 3 Properties of conditional expectation

**Theorem** (Dropping what is independent). *If  $X$  and  $Y$  are independent, then  $E(Y | X) = E(Y)$ . That is, the random variable  $E(Y | X)$  is the constant  $E(Y)$ .*

*Proof.* Independence implies  $E(Y | X = x) = E(Y)$  for all  $x$ . □

**Theorem** (Taking out what is known). *For any function  $h(\cdot)$ , we have*

$$E(h(X)Y | X) = h(X)E(Y | X).$$

*Note that the above equality means that the random variable  $g_1(X) = E(h(X)Y | X)$  is equal with the random variable  $g_2(X) = h(X)E(Y | X)$ .*

*Proof.* We know that, for any constant  $c \in \mathbb{R}$ , we have

$$E(cX) = cE(X).$$

Once we know the value  $x$  of  $X$ , the function  $h(X)$  becomes a constant  $h(x)$ . □

**Theorem** (Linearity). *For any random variables  $Y_1, Y_2, \dots, Y_n$  and  $X$ , we have*

$$E(Y_1 + Y_2 + \dots + Y_n | X) = E(Y_1 | X) + E(Y_2 | X) + \dots + E(Y_n | X).$$

**Example:**

Let  $Y_1, \dots, Y_n$  be i.i.d., and  $S_n = Y_1 + \dots + Y_n$ . Find  $E(X_1 | S_n)$ .

*Solution:* By symmetry,

$$E(Y_1 | S_n) = E(Y_2 | S_n) = \dots = E(Y_n | S_n).$$

By linearity and by taking out what is known,

$$E(Y_1 | S_n) + E(Y_2 | S_n) + \dots + E(Y_n | S_n) = E(S_n | S_n) = S_n.$$

Thus

$$E(Y_1 | S_n) = \frac{S_n}{n}.$$

**Theorem** (Adam's law: connecting conditional expectation with unconditional expectation). *For any random variables  $X$  and  $Y$ , we have*

$$E(E(Y | X)) = E(Y).$$

**Theorem** (Adam's law with extra conditioning). *For any random variables  $X, Y$  and  $Z$ , we have*

$$E(E(Y | X, Z) | Z) = E(Y | Z).$$

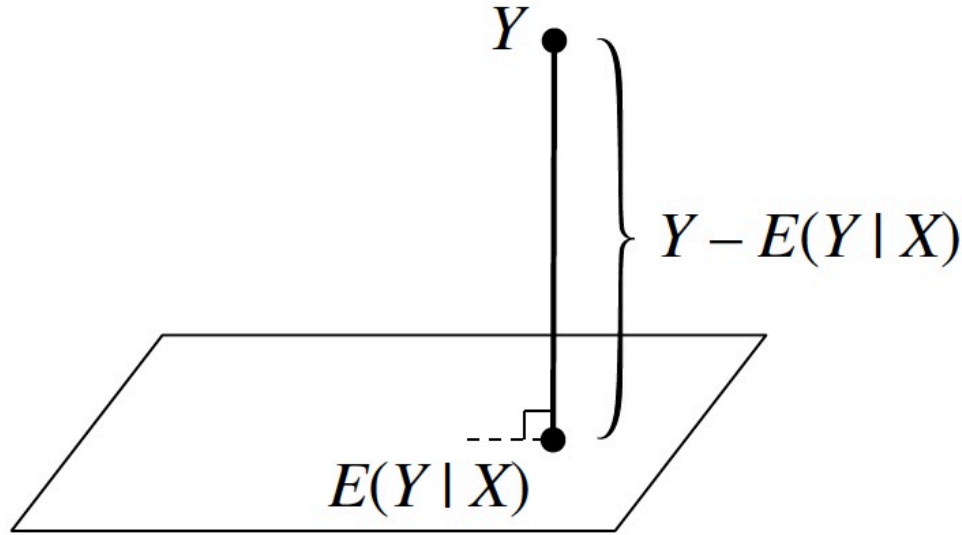


Figure 1: In this figure, we consider the vector space of all the random variables with zero mean and finite variance on a certain probability space. Each random variable corresponds with a point (vector) in this space. The subspace of random variables of the form  $h(X)$  is represented as a plane. The conditional expectation  $E(Y | X)$  belongs to this plane. To obtain  $E(Y | X)$ , we project  $Y$  onto this plane and obtain  $E(Y | X)$  as the function of  $X$  that is closest to  $Y$ . The line from  $Y$  to  $E(Y | X)$  is orthogonal (perpendicular) to the plane since any other route from  $Y$  to  $E(Y | X)$  would be longer. Then the residual  $Y - E(Y | X)$  is orthogonal to  $h(X)$  for all functions  $h(\cdot)$ , and  $E(Y | X)$  is the function of  $X$  and best predicts  $Y$ , where “best” means that the mean squared error  $E((Y - h(X))^2)$  is minimized by choosing  $h(X) = E(Y | X)$ .

**Theorem** (Projection interpretation). *For any function  $h(\cdot)$ , the random variable  $Y - E(Y | X)$  is uncorrelated with  $h(X)$ . Since*

$$E(Y - E(Y | X)) = E(Y) - E(E(Y | X)) = E(Y) - E(Y) = 0,$$

*this is equivalent with*

$$E((Y - E(Y | X))h(X)) = 0.$$

*The geometric interpretation of this result is given in Figure 1.*

*Proof.*

$$E((Y - E(Y | X))h(X)) = E(Yh(X)) - E(E(Y | X)h(X)) = h(X)E(Y) - h(X)E(E(Y | X)) = h(X)E(Y) - h(X)E(Y) = 0.$$

□

**Example:** Linear regression

In its simplest form, the linear regression model uses a single explanatory variable  $X$  to predict a response variable  $Y$ . It assumes that the conditional expectation of  $Y$  given  $X$  is linear in  $X$ :

$$E(Y | X) = a + bX.$$

(a) Show that an equivalent way to express this is to write:

$$Y = a + bX + \epsilon,$$

where  $\epsilon$  is a random variable (called the error) with  $E(\epsilon | X) = 0$ .

(b) Express the constants  $a$  and  $b$  in terms of  $E(X)$ ,  $E(Y)$ ,  $\text{Cov}(X, Y)$ , and  $\text{Var}(X)$ .

*Solution:* (a) Let  $Y = a + bX + \epsilon$  with  $E(\epsilon | X) = 0$ . Then

$$E(Y | X) = E(a | X) + E(bX | X) + E(\epsilon | X) = a + bE(X | X) = a + bX.$$

Conversely, assume that  $E(Y | X) = a + bX$ , and define  $\epsilon = Y - (a + bX)$ . Then

$$E(\epsilon | X) = E(Y | X) - E(a + bX | X) = E(Y | X) - (a + bX) = 0.$$

(b) By Adam's law, the unconditional mean of  $Y$  is

$$E(Y) = E(E(Y | X)) = E(a + bX) = a + bE(X).$$

The unconditional mean of the error  $\epsilon$  is

$$E(\epsilon) = E(E(\epsilon | X)) = E(0) = 0.$$

$X$  and  $\epsilon$  are uncorrelated because:

$$E(\epsilon X) = E(E(\epsilon X | X)) = E(XE(\epsilon | X)) = E(X \cdot 0) = 0.$$

Thus

$$\text{Cov}(X, Y) = \text{Cov}(X, a + bX + \epsilon) = \text{Cov}(X, a) + b\text{Cov}(X, X) + \text{Cov}(X, \epsilon) = b\text{Var}(X).$$

Thus  $b = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$ , and

$$a = E(Y) - bE(X) = E(Y) - \frac{\text{Cov}(X, Y)}{\text{Var}(X)}E(X)$$

## 4 Conditional variance

**Definition** (Conditional variance). The *conditional variance of  $Y$  given  $X$*  is

$$\text{Var}(Y | X) = E([Y - E(Y | X)]^2 | X).$$

This is equivalent to

$$\text{Var}(Y | X) = E(Y^2 | X) - (E(Y | X))^2.$$

The conditional variance  $\text{Var}(Y | X)$  is a random variable, and it is a function of  $X$ .

**Example:**

Let  $Z \sim N(0, 1)$  and  $Y = Z^2$ . Find  $E(Y | Z)$ ,  $\text{Var}(Y | Z)$ ,  $E(Z | Y)$ , and  $\text{Var}(Z | Y)$ .

*Solution:* We apply the taking out what is known theorem with  $h(z) = z^2$ , and obtain:

$$E(Y | Z) = E(Z^2 | Z) = E(h(Z) \cdot 1 | Z) = h(Z) \cdot E(1 | Z) = h(Z) = Z^2.$$

We can see this directly: conditional on  $Z = z$ ,  $Y = Z^2 = z^2$  a constant. Since the expectation of a constant is the constant itself, and the variance of a constant is 0, we have

$$E(Y | Z) = Z^2, \quad \text{Var}(Y | Z) = 0.$$

On the other hand, if  $Y = Z^2 = y$ , then  $Z \in \{-\sqrt{t}, \sqrt{t}\}$ . From the Bayes' rule we have:

$$P(Z = -\sqrt{t} | Y = t) \propto \phi(-\sqrt{t}), \quad P(Z = \sqrt{t} | Y = t) \propto \phi(\sqrt{t}),$$

where  $\phi(\cdot)$  is the standard Normal PDF. By the symmetry of the standard Normal, we have  $\phi(-\sqrt{t}) = \phi(\sqrt{t})$ .

Hence

$$P(Z = -\sqrt{t} | Y = t) = P(Z = \sqrt{t} | Y = t) = \frac{1}{2}.$$

Thus

$$E(Z | Y = t) = -\sqrt{t}P(Z = -\sqrt{t} | Y = t) + \sqrt{t}P(Z = \sqrt{t} | Y = t) = 0,$$

which implies  $E(Z | Y) = 0$ . Moreover

$$E(Z^2 | Y = t) = tP(Z = -\sqrt{t} | Y = t) + tP(Z = \sqrt{t} | Y = t) = t,$$

which implies  $E(Z^2 | Y) = Y$ . By definition:

$$\text{Var}(Z | Y) = E(Z^2 | Y) - (E(Z | Y))^2 = Y.$$

**Theorem** (Eve's law: connecting conditional variance to unconditional variance). *For any random variables  $X$  and  $Y$ , we have*

$$\text{Var}(Y) = E(\text{Var}(Y | X)) + \text{Var}(E(Y | X)).$$

*This relation is known as the law of total variance, or as the variance decomposition formula.*

*Proof.* Let  $g(X) = E(Y | X)$ . By Adam's law,  $E(g(X)) = E(Y)$ . Then

$$E(\text{Var}(Y | X)) = E(E(Y^2 | X) - g(X)^2) = E(Y^2) - E(g(X)^2),$$

and

$$\text{Var}(E(Y | X)) = E(g(X)^2) - (Eg(X))^2 = E(g(X)^2) - (EY)^2.$$

Thus

$$E(\text{Var}(Y | X)) + \text{Var}(E(Y | X)) = E(Y^2) - (EY)^2 = \text{Var}(Y).$$

□

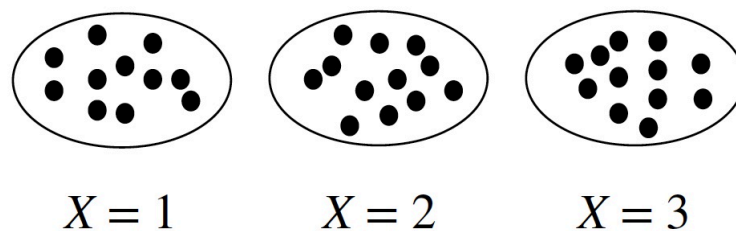


Figure 2: Illustration of Eve's law. Consider a population where each person has a value of  $X$  (e.g., age) and a value of  $Y$  (height). Assume that we divide this population into subpopulations or groups: here, we have three age groups associated with  $X = 1$ ,  $X = 2$  and  $X = 3$ . There are two sources contributing to the variation in people's heights in the overall population: (1) *within-group variation*  $E(\text{Var}(Y | X))$  represents the average amount of variation in  $Y$  (height) within (conditional on) each age group; and (2) *between-group variation*  $\text{Var}(E(Y | X))$  (here, the variance of the group means  $E(Y | X = 1)$ ,  $E(Y | X = 2)$  and  $E(Y | X = 3)$ ) represents the variance of average heights across age groups. Eve's law says that total variance is the sum of within-group and between-group variation.