

1 Conditional Expectation

1.1 Conditional distribution

Consider two random variables X and Y with joint mass function or joint density function denoted by $f_{X,Y}$, i.e.,

$$f_{X,Y}(x,y) = \begin{cases} \Pr(X=x, Y=y), & X \text{ and } Y \text{ are discrete at points } x \text{ and } y \text{ respectively} \\ \frac{\partial^2}{\partial x \partial y} \Pr(X \leq x, Y \leq y), & X \text{ and } Y \text{ are continuous at points } x \text{ and } y \text{ respectively} \end{cases}$$

We define the following concepts.

- the **marginal mass or density function of X**

$$f_X(x) = \sum_y f_{X,Y}(x,y) \quad \text{or} \quad f_X(x) = \int_{\mathbb{R}} f_{X,Y}(x,y) \, dy.$$

- the **marginal mass or density function of Y**

$$f_Y(y) = \sum_x f_{X,Y}(x,y) \quad \text{or} \quad f_Y(y) = \int_{\mathbb{R}} f_{X,Y}(x,y) \, dx.$$

- the **conditional mass or density function of X given $Y = y$**

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}, \quad \text{provided } f_Y(y) > 0. \quad (1)$$

Using the conditional distribution of X given Y , the marginal mass or density function of X can be expressed as

$$f_X(x) = \int_{\mathbb{R}} f_{X|Y}(x|y) f_Y(y) \, dy \quad \text{or} \quad f_X(x) = \sum_{y \in \mathbb{R}} f_{X|Y}(x|y) f_Y(y) \quad (2)$$

Proposition 1. *If the random variables X and Y are independent, we have*

$$f_{X,Y}(x,y) = f_X(x) f_Y(y).$$

As an immediate consequence, we have

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

1.2 Conditional expectation w.r.t. random variables

Throughout this section, we assume that X given $Y = y$ is a continuous random variable with density function $f_{X|Y}(\cdot|y)$ (if $X|Y$ is discrete, replace all the integral signs by summation signs). The conditional expectation of X given $Y = y$ is given by the expected value with respect to the conditional density function

$$\mathbb{E}[X|Y=y] = \int_{\mathbb{R}} x f_{X|Y}(x|y) \, dx.$$

This motivates the following definition:

Definition 1. The **conditional expectation of X given Y** is the random variable

$$\mathbb{E}[X|Y] = \int_{\mathbb{R}} x f_{X|Y}(x|Y) \, dx.$$

Remark 1. The conditional expectation is a random variable since it takes elements in the range of Y and assigns it to a number. In other words, if we define the function g through

$$g(y) = \mathbb{E}[X|Y = y] = \int_{\mathbb{R}} x f_{X|Y}(x|y) dx,$$

then

$$\mathbb{E}[X|Y] = g(Y).$$

We can interpret the conditional expected value as the “best” estimate for the value of X given a realization of Y (see Problem 1.6).

The conditional expectation obeys the following useful properties.

Proposition 2. *The conditional expectation has the following properties:*

1. *Law of total expectation:* $\mathbb{E}[\mathbb{E}[X|Y]] = \mathbb{E}[X]$
2. *Pulling out known factors:* If h is a function, then

$$\mathbb{E}[h(Y)X|Y] = h(Y)\mathbb{E}[X|Y]$$

Proof. The properties follow directly from the definition

(a) We define $g(y) = \mathbb{E}[X|Y = y] = \int_{\mathbb{R}} x f_{X|Y}(x|y) dx$. By the definition of the expected value,

$$\begin{aligned} \mathbb{E}[\mathbb{E}[X|Y]] &= \mathbb{E}[g(Y)] = \int_{\mathbb{R}} g(y) f_Y(y) dy = \int_{\mathbb{R}} \left(\int_{\mathbb{R}} x f_{X|Y}(x|y) dx \right) f_Y(y) dy \\ &= \int_{\mathbb{R}} \int_{\mathbb{R}} x f_{X|Y}(x|y) f_Y(y) dx dy \\ &= \int_{\mathbb{R}} \int_{\mathbb{R}} x f_{X,Y}(x, y) dx dy \\ &= \int_{\mathbb{R}} x \left(\int_{\mathbb{R}} f_{X,Y}(x, y) dy \right) dx \\ &= \int_{\mathbb{R}} x f_X(x) dx = \mathbb{E}[X]. \end{aligned}$$

(b) For any y in the support of Y ,

$$g(y) = \mathbb{E}[h(Y)X|Y = y] = \int_{\mathbb{R}} h(y) x f_{X|Y}(x|y) dx = h(y) \int_{\mathbb{R}} x f_{X|Y}(x|y) dx = h(y) \mathbb{E}[X|Y = y].$$

Therefore,

$$\mathbb{E}[h(Y)X|Y] = g(Y) = h(Y)\mathbb{E}[X|Y].$$

□

Likewise, one can define the conditional variance in the obvious way.

Definition 2. The **conditional variance of X given Y** is defined as

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

The conditional variance satisfies the following useful properties.

Proposition 3. *We have*

1. $\text{Var}(X|Y) = \mathbb{E}[X^2|Y] - (\mathbb{E}[X|Y])^2$
2. *Law of total variance:* $\text{Var}(X) = \mathbb{E}[\text{Var}(X|Y)] + \text{Var}(\mathbb{E}[X|Y])$

Proof. (a) With $g(Y) = \mathbb{E}[X|Y]$ we have from Proposition 2 (b) that

$$\begin{aligned}
 \text{Var}(X|Y) &= \mathbb{E}[X^2 - 2X\mathbb{E}[X|Y] + (\mathbb{E}[X|Y])^2 | Y] \\
 &= \mathbb{E}[X^2 | Y] - 2\mathbb{E}[X\mathbb{E}[X|Y] | Y] + \mathbb{E}[(\mathbb{E}[X|Y])^2 | Y] \\
 &= \mathbb{E}[X^2 | Y] - 2\mathbb{E}[Xg(Y) | Y] + \mathbb{E}[(g(Y))^2 | Y] \\
 &= \mathbb{E}[X^2 | Y] - 2g(Y) \cdot \mathbb{E}[X | Y] + (g(Y))^2 \mathbb{E}[1 | Y] \quad (\text{by Proposition 2 (b)}) \\
 &= \mathbb{E}[X^2 | Y] - 2\mathbb{E}[X | Y] \cdot \mathbb{E}[X | Y] + (\mathbb{E}[X|Y])^2 \\
 &= \mathbb{E}[X^2 | Y] - (\mathbb{E}[X | Y])^2
 \end{aligned}$$

(b) It follows from (a) and Proposition 2 (a) that

$$\begin{aligned}
 \mathbb{E}[\text{Var}(X|Y)] &= \mathbb{E}[\mathbb{E}[X^2|Y]] - \mathbb{E}[(\mathbb{E}[X|Y])^2] \\
 &= \mathbb{E}[X^2] - \mathbb{E}[(\mathbb{E}[X|Y])^2].
 \end{aligned}$$

On the other hand,

$$\begin{aligned}
 \text{Var}(\mathbb{E}[X|Y]) &= \mathbb{E}[(\mathbb{E}[X|Y])^2] - (\mathbb{E}[\mathbb{E}[X|Y]])^2 \\
 &= \mathbb{E}[(\mathbb{E}[X|Y])^2] - (\mathbb{E}[X])^2.
 \end{aligned}$$

Combining the preceding two relations implies

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

□

1.3 Example Problems

Problem 1.1. Suppose a fair coin is tossed 3 times. Define the random variables X = “number of Heads”, and

$$Y = \begin{cases} 1 & \text{Head occurs on the first toss,} \\ 0 & \text{Tail occurs on the first toss.} \end{cases}$$

1. Find the joint PMF for (X, Y) .
2. Are X and Y independent?
3. What is the conditional distribution of X given Y ?
4. What is the probability that $X + Y = 2$?

Solution 1.1.

Part 1: We can compute all the probabilities one by one and encode the joint PMF of X and Y in the table

$f_{X,Y}(x,y)$		x				$f_Y(y)$
		0	1	2	3	
y	0	1/8	2/8	1/8	0	1/2
	1	0	1/8	2/8	1/8	1/2
$f_X(x)$		1/8	3/8	3/8	1/8	1

Part 2: We can see

$$f_{X,Y}(0,1) = 0 \neq \frac{1}{8} \cdot \frac{1}{2} = f_X(0)f_Y(1)$$

which implies that X and Y are not independent (which makes perfect sense, as the number of heads we have should depend on whether we had heads in the first toss).

Part 3: Using the formula $f_{X|Y}(x|y) = f_{X,Y}(x,y)/f_Y(y)$ we find

	x			
	0	1	2	3
$f_{X Y}(x y=0)$	2/8	4/8	2/8	0
$f_{X Y}(x y=1)$	0	2/8	4/8	2/8

Part 4: We have $X + Y = 2$ if and only if $X = 2, Y = 0$ or $X = 1, Y = 1$. We can sum these terms up in the joint PMF

$$\mathbb{P}(X + Y = 2) = f(2,0) + f(1,1) + f(0,2) = \frac{1}{8} + \frac{1}{8} = \frac{1}{4}.$$

Problem 1.2. Suppose that X and Θ are two random variables such that X given $\Theta = \theta$ is Poisson distributed with mean θ , i.e.,

$$f_{X|\Theta}(k|\theta) = e^{-\theta} \frac{\theta^k}{k!}, \quad k = 0, 1, 2, \dots$$

and Θ is Gamma distributed with parameters $\alpha, \beta > 0$. That is, Θ has the density function

$$f_{\Theta}(\theta) = \frac{\beta^{\alpha} \theta^{\alpha-1} e^{-\beta\theta}}{\Gamma(\alpha)}, \quad \theta > 0,$$

where Γ denotes the Gamma function,

$$\Gamma(\alpha) = \int_0^{\infty} \theta^{\alpha-1} e^{-\theta} d\theta.$$

Compute the marginal mass function of X .

Solution 1.2. The marginal mass function of X is given by

$$\begin{aligned}
 \mathbb{P}(X = k) &= \int_0^\infty f_{X|\Theta}(k|\theta) f_\Theta(\theta) d\theta \\
 &= \int_0^\infty \frac{\theta^k e^{-\theta}}{k!} \cdot \frac{\beta^\alpha \theta^{\alpha-1} e^{-\beta\theta}}{\Gamma(\alpha)} d\theta \\
 &= \frac{\beta^\alpha}{k! \Gamma(\alpha)} \int_0^\infty \theta^{k+\alpha-1} e^{-(\beta+1)\theta} d\theta \\
 &= \frac{\beta^\alpha}{k! \Gamma(\alpha)} \cdot \frac{1}{(\beta+1)^{k+\alpha}} \int_0^\infty x^{k+\alpha-1} e^{-x} dx \\
 &= \frac{1}{k! \Gamma(\alpha)} \left(\frac{\beta}{\beta+1} \right)^\alpha \left(\frac{1}{\beta+1} \right)^k \Gamma(k+\alpha) \\
 &= \frac{(k+\alpha-1)(k+\alpha-2) \cdots (\alpha+1)\alpha}{k!} \left(1 - \frac{1}{\beta+1} \right)^\alpha \left(\frac{1}{\beta+1} \right)^k \\
 &= \binom{k+\alpha-1}{k} \left(1 - \frac{1}{\beta+1} \right)^\alpha \left(\frac{1}{\beta+1} \right)^k.
 \end{aligned}$$

Therefore, X follows a negative binomial distribution with parameters α and $\frac{1}{\beta+1}$.

Problem 1.3. Suppose that X given $\Theta = \theta$ is Poisson distributed with mean θ and Θ is Gamma distributed with density function

$$f_\Theta(\theta) = \frac{\beta^\alpha \theta^{\alpha-1} e^{-\beta\theta}}{\Gamma(\alpha)}, \quad \theta > 0.$$

1. Compute $\mathbb{E}[X]$.
2. Compute $\text{Var}[X]$.

Recall that if $X \sim \text{Poi}(\lambda)$ then $\mathbb{E}[X] = \lambda$ and $\text{Var}(X) = \lambda$ and if $X \sim \Gamma(\alpha, \beta)$ then $\mathbb{E}[X] = \frac{\alpha}{\beta}$ and $\text{Var}(X) = \frac{\alpha}{\beta^2}$.

Solution 1.3.

(a) Using the law of total expectation,

$$\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X|\Theta]] = \mathbb{E}[\Theta] = \frac{\alpha}{\beta}.$$

We used the fact that the expected value of a Poisson distributed random variable is equal to its mean parameter.

(b) By the law of total variance

$$\begin{aligned}
 \text{Var}(X) &= \mathbb{E}[\text{Var}(X|\Theta)] + \text{Var}(\mathbb{E}[X|\Theta]) \\
 &= \mathbb{E}[\Theta] + \text{Var}(\Theta) \\
 &= \frac{\alpha}{\beta} + \frac{\alpha}{\beta^2} = \frac{\alpha(\beta+1)}{\beta^2}.
 \end{aligned}$$

We used the fact that the expected value of a Poisson distributed random variable is equal to its mean parameter and its variance is equal to its mean parameter.

Problem 1.4. Suppose that

$$X = \begin{cases} \sum_{i=1}^N Y_i, & \text{if } N > 0, \\ 0, & \text{if } N = 0, \end{cases}$$

where N is Poisson distributed with mean λ and Y_1, Y_2, \dots is a sequence of iid random variables with mean μ and variance σ^2 that is independent of N . We say that X is a **compound Poisson random variable**.

1. Compute $\mathbb{E}[X]$.
2. Compute $\text{Var}[X]$.

Recall that if $X \sim \text{Poi}(\lambda)$ then $\mathbb{E}[X] = \lambda$ and $\text{Var}(X) = \lambda$.

Solution 1.4.

(a) By the law of total expectation

$$\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X|N]] = \mathbb{E}[N\mu] = \lambda\mu,$$

(b) By the law of total variance

$$\begin{aligned} \text{Var}(X) &= \mathbb{E}[\text{Var}(X|N)] + \text{Var}(\mathbb{E}[X|N]) \\ &= \mathbb{E}[N\sigma^2] + \text{Var}(N\mu) \\ &= \sigma^2\mathbb{E}[N] + \mu^2\text{Var}(N) \\ &= \lambda(\sigma^2 + \mu^2). \end{aligned}$$

We used the fact that the expected value of a Poisson distributed random variable is equal to its mean parameter and its variance is equal to its mean parameter.

Problem 1.5. For any constant c , show that

$$\mathbb{E}[(X - c)^2] \geq \mathbb{E}[(X - \mathbb{E}[X])^2].$$

In particular, the expected value is the constant that minimizes the mean squared error.

Solution 1.5. This proof follows directly from the properties of the expected value. By adding and subtracting $\mathbb{E}[X]$, we see that

$$\begin{aligned} \mathbb{E}[(X - c)^2] &= \mathbb{E}[(X - \mathbb{E}[X] + \mathbb{E}[X] - c)^2] \\ &= \mathbb{E}[(X - \mathbb{E}[X])^2] + \mathbb{E}[(\mathbb{E}[X] - c)^2] + 2\mathbb{E}[(X - \mathbb{E}[X])(\mathbb{E}[X] - c)] \end{aligned}$$

Since $\mathbb{E}[X] - c$ is not random, we see that the cross terms vanish

$$\mathbb{E}[(X - \mathbb{E}[X])(\mathbb{E}[X] - c)] = (\mathbb{E}[X] - c)\mathbb{E}[(X - \mathbb{E}[X])] = (\mathbb{E}[X] - c)(\mathbb{E}[X] - \mathbb{E}[X]) = 0.$$

Since $\mathbb{E}[(\mathbb{E}[X] - c)^2] \geq 0$, we conclude that

$$\mathbb{E}[(X - c)^2] = \mathbb{E}[(X - \mathbb{E}[X])^2] + \mathbb{E}[(\mathbb{E}[X] - c)^2] \geq \mathbb{E}[(X - \mathbb{E}[X])^2]$$

as required.

Problem 1.6. For any measurable function f , show that

$$\mathbb{E}[(X - f(Y))^2] \geq \mathbb{E}[(X - \mathbb{E}[X | Y])^2].$$

In particular, the conditional expectation minimizes the mean squared error.

Solution 1.6. This proof follows directly from the properties of the conditional expected value. By adding and subtracting $\mathbb{E}[X | Y]$, we see that

$$\begin{aligned} \mathbb{E}[(X - f(Y))^2] &= \mathbb{E}[(X - \mathbb{E}[X | Y] + \mathbb{E}[X | Y] - f(Y))^2] \\ &= \mathbb{E}[(X - \mathbb{E}[X | Y])^2] + \mathbb{E}[(\mathbb{E}[X | Y] - f(Y))^2] + 2\mathbb{E}[(X - \mathbb{E}[X | Y])(\mathbb{E}[X | Y] - f(Y))] \end{aligned}$$

Applying the law of total expectation and using the fact that $\mathbb{E}[X | Y]$ and $f(Y)$ are measurable functions of Y , we see that the cross terms vanish

$$\begin{aligned} \mathbb{E}[(X - \mathbb{E}[X | Y])(\mathbb{E}[X | Y] - f(Y))] &= \mathbb{E}[\mathbb{E}[(X - \mathbb{E}[X | Y])(\mathbb{E}[X | Y] - f(Y)) | Y]] \\ &= \mathbb{E}[(\mathbb{E}[X | Y] - f(Y)) \mathbb{E}(X - \mathbb{E}[X | Y]) | Y]] \\ &= \mathbb{E}[(\mathbb{E}[X | Y] - f(Y))(\mathbb{E}[X | Y] - \mathbb{E}[X | Y])] \\ &= 0. \end{aligned}$$

Since $\mathbb{E}[(\mathbb{E}[X | Y] - f(Y))^2] \geq 0$, we conclude that

$$\mathbb{E}[(X - f(Y))^2] = \mathbb{E}[(X - \mathbb{E}[X | Y])^2] + \mathbb{E}[(\mathbb{E}[X | Y] - f(Y))^2] \geq \mathbb{E}[(X - \mathbb{E}[X | Y])^2]$$

as required.

Remark 2. Notice that this proof only uses the law of total expectation and the trick that allows us to pull known factors. This will imply that this proof carries over to the general conditional expectation with respect to σ -algebras.

2 Conditional expectations w.r.t. σ -algebras

We now introduce general definition of conditional expectation that will allow us to condition on more general forms of (random) information. We will use σ -algebra $\mathcal{F}_0 \subset \mathcal{F}$ as a **model of information** and define the general notation of the conditional expectation of X given information \mathcal{F}_0

$$\mathbb{E}[X|\mathcal{F}_0].$$

A σ -algebra is a natural model for the information because it contains both the negation and union of outcomes, which can easily deduced from existing information.

2.1 Constructing σ -algebras

We first take a closer look at possible constructions of σ -algebras.

Definition 3. Given a collection of sets \mathcal{A} of Ω , the σ -**algebra** generated by the collection of sets \mathcal{A} is the smallest σ -algebra containing \mathcal{A} and is often denoted by $\sigma(\mathcal{A})$.

Example 1. On $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4\}$ consider the following two partitions:

$$\begin{aligned}\mathcal{P}_1 &= \{\{\omega_1\}, \{\omega_2\}, \{\omega_3\}, \{\omega_4\}\} \\ \mathcal{P}_2 &= \{\{\omega_1, \omega_2\}, \{\omega_3, \omega_4\}\}.\end{aligned}$$

In the first, we are able to distinguish between all elements of Ω . In the second, we cannot distinguish between ω_1 and ω_2 and between ω_3 and ω_4 . Thus, \mathcal{P}_1 is *finer* than \mathcal{P}_2 . The σ -algebra $\sigma(\mathcal{P}_1)$ is equal to the power set of Ω , i.e., it contains all subsets of Ω . On the other hand,

$$\sigma(\mathcal{P}_2) = \{\emptyset, \{\omega_1, \omega_2\}, \{\omega_3, \omega_4\}, \Omega\}.$$

2.1.1 σ -algebras generated by random variables

Suppose the σ -algebra \mathcal{F}_0 corresponds to the information from observing the values of a collection Y_1, \dots, Y_n of \mathcal{F} -measurable random variables. Informally, \mathcal{F}_0 then consists of all events that can be described through the random variables Y_1, \dots, Y_n .

Definition 4. The σ -algebra \mathcal{F}_0 generated by Y_1, \dots, Y_n is the σ -algebra generated by events of the form $\{Y_i \leq x\}$ for all $x \in \mathbb{R}$ and $i = 1, \dots, n$. We write

$$\mathcal{F}_0 := \sigma(Y_1, \dots, Y_n).$$

Remark 3. Let X be a random variable on $(\Omega, \mathcal{F}, \mathbb{P})$. One can prove that the σ -algebra $\sigma(X)$ generated by X is equivalent to

$$\sigma(X) = \{X^{-1}(B) \mid B \in \mathcal{B}(\mathbb{R})\},$$

where we recall that $\mathcal{B}(\mathbb{R})$ is the Borel σ -algebra on \mathbb{R} and

$$X^{-1}(B) = \{\omega \in \Omega \mid X(\omega) \in B\} = \{X \in B\}$$

is the pre-image of B .

Example 2. Let X be the number of heads obtained for a coin tossed twice. In this case, $\Omega = \{HH, HT, TH, TT\}$. Clearly, $X(HH) = 2$, $X(HT) = X(TH) = 1$ and $X(TT) = 0$. We have

$$\sigma(X) = \{\emptyset, \{HH\}, \{TT\}, \{TT, HH\}, \{HT, TH\}, \{HT, TH, HH\}, \{HT, TH, TT\}, \Omega\}.$$

Notice that this set is not equal to the power set of Ω . In particular, the set $\{HT\}$ is not in $\sigma(X)$ since knowing the number of heads does not allow you to determine that $\{HT\}$ happened since it is indistinguishable from the event $\{TH\}$, while $\{HT, TH\}$ is in the set, since the events you flipped HT or TH corresponds to the event of flipping exactly 1 heads.

2.2 Independent σ -algebras

Consider a probability space $(\Omega, \mathcal{F}, \mathbb{P})$. Recall that two **events** $A, B \in \mathcal{F}$ are called **independent** under \mathbb{P} if

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

The notion of independence can be extended to σ -algebras in the obvious way.

Definition 5. Two σ -algebras $\mathcal{F}_1, \mathcal{F}_2 \subset \mathcal{F}$ are **independent** if

$$\mathbb{P}(A_1 \cap A_2) = \mathbb{P}(A_1)\mathbb{P}(A_2), \quad \text{for any } A_1 \in \mathcal{F}_1 \text{ and } A_2 \in \mathcal{F}_2.$$

The notation of independence of random variables can also be stated with respect to σ -algebras.

Definition 6. Two **random variables** X_1 and X_2 on $(\Omega, \mathcal{F}, \mathbb{P})$ are **independent** if $\sigma(X_1)$ and $\sigma(X_2)$ are independent.

Remark 4. This notion of independence is equivalent to the earlier notation defined in Week 1. That is the following statements are equivalent

1. X_1 and X_2 are independent,
2. The probabilities satisfy

$$\mathbb{P}(X_1 \in B_1, X_2 \in B_2) = \mathbb{P}(X_1 \in B_1)\mathbb{P}(X_2 \in B_2),$$

for any $B_1, B_2 \in \mathcal{B}(\mathbb{R})$.

3. The CDFs satisfy

$$\mathbb{P}(X_1 \leq x_1, X_2 \leq x_2) = \mathbb{P}(X_1 \leq x_1)\mathbb{P}(X_2 \leq x_2) = F_{X_1}(x_1)F_{X_2}(x_2) \quad \forall x_1, x_2$$

The independence between a random variable and σ -algebra is also defined in the natural way.

Definition 7. A random variable X is independent of a σ -algebra $\mathcal{F}_1 \subset \mathcal{F}$ if $\sigma(X)$ and \mathcal{F}_1 are independent.

2.3 Conditional expectations with respect to general σ -algebras

Definition 8. Consider a random variable X on $(\Omega, \mathcal{F}, \mathbb{P})$ and a σ -algebra $\mathcal{F}_0 \subset \mathcal{F}$. We define the **conditional expectation** of X given \mathcal{F}_0 as a random variable $\mathbb{E}[X|\mathcal{F}_0]$ satisfying the following two conditions:

1. $\mathbb{E}[X|\mathcal{F}_0]$ is a \mathcal{F}_0 -measurable random variable.
2. $\mathbb{E}[\mathbb{E}[X|\mathcal{F}_0]\mathbb{1}_A] = \mathbb{E}[X\mathbb{1}_A]$ for any $A \in \mathcal{F}_0$.

Remark 5. The conditional expectation is unique, so there is only one random variable that can satisfy the second condition (see Problem 2.3).

The first condition is natural because we want to be able to define the conditional expectation with respect to the outcome of a random events: your best guess for a random variable should be able to adapt to a random event in \mathcal{F}_0 . Since the expected value is the “best” (see Problem 1.5) guess for X , the second condition means that the conditional expected value is always the best guess for X given the outcome of any $A \in \mathcal{F}_0$.

Example 3. One can show that the preceding definition gives the following special cases:

- Consider the case $\mathcal{F}_0 = \sigma(Y)$. In general, a random variable Z is \mathcal{F}_0 -measurable if and only if there is a function h such that

$$Z = h(Y_1, \dots, Y_n).$$

The conditional expectation is the function given by

$$\mathbb{E}[X|\mathcal{F}_0] = \mathbb{E}[X|Y]$$

where the right-hand side is the function of Y defined in the same way as in Section 1.2.

- Consider the case $\mathcal{F}_0 = \sigma(Y_1, \dots, Y_n)$. In general, a random variable Z is \mathcal{F}_0 measurable if and only if there is a function h such that

$$Z = h(Y_1, \dots, Y_n).$$

The conditional expectation is the function given by

$$\mathbb{E}[X|\mathcal{F}_0] = \mathbb{E}[X|Y_1, \dots, Y_n] = g(Y_1, \dots, Y_n).$$

where the function g can be defined in the same way as in Section 1.2. We denote by f_{Y_1, \dots, Y_n} the joint probability density (or probability mass function) of Y_1, \dots, Y_n and define

$$f_{X|Y_1, \dots, Y_n}(x|y_1, \dots, y_n) := \frac{f_{X, Y_1, \dots, Y_n}(x, y_1, \dots, y_n)}{f_{Y_1, \dots, Y_n}(y_1, \dots, y_n)},$$

where f_{X, Y_1, \dots, Y_n} is the joint density of X, Y_1, \dots, Y_n . Then we let

$$g(y_1, \dots, y_n) = \int_{\mathbb{R}} x f_{X|Y_1, \dots, Y_n}(x|y_1, \dots, y_n) dx.$$

- Let $\mathcal{P} = \{A_1, A_2, \dots\}$ be a partition of Ω and let $\mathcal{F}_0 = \sigma(\mathcal{P})$. In general, a random variable Z is \mathcal{F}_0 -measurable if and only if Z is of the form

$$Z = \sum_{i=1}^{\infty} z_i \mathbb{1}_{A_i}$$

for some real numbers z_1, z_2, \dots . The conditional expectation is the function given by

$$\mathbb{E}[X|\mathcal{F}_0] = \sum_{i=1}^{\infty} \mathbb{E}[X|A_i] \mathbb{1}_{A_i}$$

where the coefficients are given by the **(elementary) conditional expectation**

$$\mathbb{E}[X|A_i] = \frac{\mathbb{E}[X \mathbb{1}_{A_i}]}{\mathbb{P}(A_i)}$$

whenever $\mathbb{P}(A_i) > 0$ and 0 if $\mathbb{P}(A_i) = 0$

The following proposition lists many useful propositions of the conditional expectation that will be used to do computations with conditional expectations.

Proposition 4. For a random variable X on $(\Omega, \mathcal{F}, \mathbb{P})$ and a σ -algebra $\mathcal{F}_0 \subset \mathcal{F}$:

1. **Stability:** If X is \mathcal{F}_0 -measurable, then $\mathbb{E}[X|\mathcal{F}_0] = X$
2. **Trivial σ -algebra:** If \mathcal{G} is the trivial σ -algebra, i.e., $\mathcal{G} = \{\emptyset, \Omega\}$, then

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

3. **Law of total expectation:** $\mathbb{E}[\mathbb{E}[X|\mathcal{F}_0]] = \mathbb{E}[X]$

4. **Linearity:** $\mathbb{E}[aX + bY|\mathcal{F}_0] = a\mathbb{E}[X|\mathcal{F}_0] + b\mathbb{E}[Y|\mathcal{F}_0]$

5. **Pulling out known factors:** If Y is \mathcal{F}_0 -measurable, then

$$\mathbb{E}[XY|\mathcal{F}_0] = Y\mathbb{E}[X|\mathcal{F}_0]$$

6. **Tower property:** If $\mathcal{F}_0 \subset \mathcal{F}_1 \subset \mathcal{F}$ are σ -algebras, then

$$\mathbb{E}[X|\mathcal{F}_0] = \mathbb{E}[\mathbb{E}[X|\mathcal{F}_1]|\mathcal{F}_0]$$

7. **(Interchanged) Tower Property:** If $\mathcal{F}_0 \subset \mathcal{F}_1 \subset \mathcal{F}$ are σ -algebras, then

$$\mathbb{E}[X|\mathcal{F}_0] = \mathbb{E}[\mathbb{E}[X|\mathcal{F}_0]|\mathcal{F}_1]$$

8. **Independence:** If X is independent of \mathcal{F}_0 , then

$$\mathbb{E}[X|\mathcal{F}_0] = \mathbb{E}[X].$$

9. **Monotonicity:** If $X \leq Y$ then

$$\mathbb{E}[X|\mathcal{F}_0] \leq \mathbb{E}[Y|\mathcal{F}_0].$$

10. **Jensen's inequality:** If $\phi: \mathbb{R} \rightarrow \mathbb{R}$ is convex, then

$$\phi(\mathbb{E}[X|\mathcal{F}_0]) \leq \mathbb{E}[\phi(X)|\mathcal{F}_0]$$

2.4 Example Problems

Problem 2.1. Show that $\mathbb{E}[1|\mathcal{F}_0] = 1$ for any $\mathcal{F}_0 \subseteq \mathcal{F}$.

Solution 2.1. We require that for every $A \in \mathcal{F}_0$,

$$\mathbb{E}[\mathbb{1}_A] = \mathbb{E}[\mathbb{1}_A \mathbb{E}[1|\mathcal{F}_0]].$$

Clearly this equality holds when $\mathbb{E}[1|\mathcal{F}_0] = 1$, so we conclude by uniqueness. An alternative proof is to realize that 1 is independent of \mathcal{F}_0 since $\sigma(1) = \{\emptyset, \mathcal{F}\}$ which is trivially independent with \mathcal{F} .

Problem 2.2. We roll two fair dice and record their values D_1 and D_2 . Let $X = D_1 + D_2$ denote the sum of both dice, and let $Y = D_1$ denote the value of the first dice. Compute $\mathbb{E}[X]$ and $\mathbb{E}[X|Y]$.

Solution 2.2. The expected value of X by linearity is

$$\mathbb{E}[X] = \mathbb{E}[D_1] + \mathbb{E}[D_2] = \sum_{x=1}^6 \frac{x}{6} + \sum_{y=1}^6 \frac{y}{6} = \frac{7}{2} + \frac{7}{2} = 7.$$

Our expected value of X will adapt if we know the value of D_2 . In particular, we have

$$\mathbb{E}[X|Y] = \mathbb{E}[D_1 + Y|Y] = \mathbb{E}[D_1|Y] + \mathbb{E}[Y|Y] = \mathbb{E}[D_1|Y] + Y\mathbb{E}[1|Y] = \mathbb{E}[D_1] + Y = \frac{7}{2} + Y.$$

where we used linearity in the second equality, pulling out known factors in the third and independence in the fourth equality.

Problem 2.3. Show that the conditional expectation is unique, i.e. there is only one \mathcal{F}_0 measurable random variable denoted by $\mathbb{E}[X|\mathcal{F}_0]$ that satisfies

$$\mathbb{E}[\mathbb{E}[X|\mathcal{F}_0]\mathbb{1}_A] = \mathbb{E}[X\mathbb{1}_A]$$

for any $A \in \mathcal{F}_0$.

Solution 2.3. Suppose there exists random variables $Y = \mathbb{E}[X | \mathcal{F}_0]$ and $Y' = \mathbb{E}[X | \mathcal{F}_0]$ such that $\mathbb{P}(Y = Y') \neq 1$. We must have that either $\mathbb{P}(Y > Y') > 0$ or $\mathbb{P}(Y < Y') > 0$. Suppose without loss of generality that the first holds. Since both Y and Y' are measurable, the event $A = \{Y > Y'\} \in \mathcal{F}$ and $\mathbb{P}(Y - Y' > 0) = \mathbb{P}(Y > Y') > 0$, so

$$\mathbb{E}[(Y - Y')\mathbb{1}_A] > 0$$

since $Y > Y'$ on A and $\mathbb{P}(A) > 0$. However, by the second property of conditional expectations,

$$\mathbb{E}[(Y - Y')\mathbb{1}_A] = \mathbb{E}[(\mathbb{E}[X | \mathcal{F}_0] - \mathbb{E}[X | \mathcal{F}_0])\mathbb{1}_A] = \mathbb{E}[X\mathbb{1}_A] - \mathbb{E}[X\mathbb{1}_A] = 0$$

which is a contradiction.

Problem 2.4. Prove all the statements in Proposition 4

Solution 2.4. These properties will allow from manipulations of the definition of conditional expectation.

1. **Stability:** Clearly X is \mathcal{F}_0 measurable. Furthermore, the random variable X satisfies the second condition for the conditional expectation $\mathbb{E}[X | \mathcal{F}_0]$ trivially, i.e. for any $A \in \mathcal{F}_0$,

$$\mathbb{E}[X\mathbb{1}_A] = \mathbb{E}[\mathbb{E}[X | \mathcal{F}_0]\mathbb{1}_A]$$

so $\mathbb{E}[X | \mathcal{F}_0] = X$ by uniqueness.

Alternative Proof: This is also a special case of the the pulling out known factors property. Suppose that X is \mathcal{F}_0 measurable. By pulling known factors (proved later) we have that

$$\mathbb{E}[X | \mathcal{F}_0] = X \mathbb{E}[1 | \mathcal{F}_0] = X.$$

2. **Trivial σ -algebra:** Since $\mathbb{E}[X | \mathcal{G}]$ is \mathcal{G} measurable it must be a constant function. To find what this constant $\mathbb{E}[X | \mathcal{G}] = c$ must be, we can take $A = \Omega$ see that

$$c = \mathbb{E}[c] = \mathbb{E}[\mathbb{E}[X | \mathcal{G}]] = \mathbb{E}[\mathbb{E}[X | \mathcal{G}]\mathbb{1}_\Omega] = \mathbb{E}[X\mathbb{1}_\Omega] = \mathbb{E}[X].$$

3. **Law of Total Expectation:** Taking $A = \Omega$, we have that

$$\mathbb{E}[\mathbb{E}[X | \mathcal{F}_0]] = \mathbb{E}[\mathbb{E}[X | \mathcal{F}_0]\mathbb{1}_\Omega] = \mathbb{E}[X\mathbb{1}_\Omega] = \mathbb{E}[X].$$

4. **Linearity:** This implies that the conditional expectation behaves like the usual expectation. We check that the right hand side satisfies the definition of conditional expectation for the random variable $\mathbb{E}[aX + bY | \mathcal{F}_0]$, i.e. for any $A \in \mathcal{F}_0$

$$\begin{aligned} \mathbb{E}[(a\mathbb{E}[X | \mathcal{F}_0] + b\mathbb{E}[Y | \mathcal{F}_0])\mathbb{1}_A] &= a\mathbb{E}[\mathbb{E}[X | \mathcal{F}_0]\mathbb{1}_A] + b\mathbb{E}[\mathbb{E}[Y | \mathcal{F}_0]\mathbb{1}_A] \\ &= a\mathbb{E}[X\mathbb{1}_A] + b\mathbb{E}[Y\mathbb{1}_A] \\ &= \mathbb{E}[(aX + bY)\mathbb{1}_A] \\ &= \mathbb{E}[\mathbb{E}[(aX + bY) | \mathcal{F}_0]\mathbb{1}_A]. \end{aligned}$$

Therefore, by uniqueness $\mathbb{E}[aX + bY | \mathcal{F}_0] = a\mathbb{E}[X | \mathcal{F}_0] + b\mathbb{E}[Y | \mathcal{F}_0]$.

5. **Pulling out known factors:** The statement itself is very natural, since Y being \mathcal{F}_0 measurable means it is in some sense not random once given information \mathcal{F}_0 , so we can factor out.

The proof is a bit technical for this course, but we will include it here for completeness. It suffices to show this result for constant random variables $Y = \mathbb{1}_B \in \mathcal{F}_0$, because more general random variables can be approximated by linear combinations of such random variables, so the result will follow from linearity. Therefore, it remains to show that

$$\mathbb{E}[\mathbb{1}_B X \mid \mathcal{F}_0] = \mathbb{1}_B \mathbb{E}[X \mid \mathcal{F}_0].$$

Notice that for any $A \in \mathcal{F}_0$, $A \cap B \in \mathcal{F}_0$ by the definition of a σ -algebra. We check that the random variable $\mathbb{1}_B \mathbb{E}[X \mid \mathcal{F}_0]$ (which is clearly \mathcal{F}_0 measurable since it is the product of measurable functions) satisfies the definition of the conditional expectation $\mathbb{E}[\mathbb{1}_B X \mid \mathcal{F}_0]$, i.e. for any $A \in \mathcal{F}_0$

$$\mathbb{E}[(\mathbb{1}_B \mathbb{E}[X \mid \mathcal{F}_0]) \mathbb{1}_A] = \mathbb{E}[\mathbb{E}[X \mid \mathcal{F}_0] \mathbb{1}_{A \cap B}] = \mathbb{E}[X \mathbb{1}_{A \cap B}] = \mathbb{E}[(X \mathbb{1}_B) \mathbb{1}_A] = \mathbb{E}[\mathbb{E}[\mathbb{1}_B X \mid \mathcal{F}_0] \mathbb{1}_A].$$

By uniqueness, $\mathbb{E}[\mathbb{1}_B X \mid \mathcal{F}_0] = \mathbb{1}_B \mathbb{E}[X \mid \mathcal{F}_0]$.

6. **Tower Property:** This proof is part of Assignment #2.
7. **(Interchanged) Tower Property:** This proof is part of Assignment #2.
8. **Independence:** Clearly $\mathbb{E}[X]$ is \mathcal{F}_0 measurable since it is a constant function. We show that $\mathbb{E}[X]$ satisfies the definition of conditional expectation of $\mathbb{E}[X \mid \mathcal{F}_0]$, i.e. for any $A \in \mathcal{F}_0 \subset \mathcal{F}_1$ we have

$$\mathbb{E}[\mathbb{E}[X] \mathbb{1}_A] = \mathbb{E}[X] \mathbb{E}[\mathbb{1}_A] \underbrace{=}_{\star} \mathbb{E}[X \mathbb{1}_A] = \mathbb{E}[\mathbb{E}[X \mid \mathcal{F}_0] \mathbb{1}_A]$$

where in \star we used the fact that if X and Y are independent then $\mathbb{E}[X] \mathbb{E}[Y] = \mathbb{E}[XY]$. We conclude that $\mathbb{E}[X] = \mathbb{E}[X \mid \mathcal{F}_0]$ by uniqueness.

9. **Monotonicity:** We proceed by contradiction. Suppose that $X \leq Y$, but assume for the sake of contradiction that the set $A = \{\mathbb{E}[X \mid \mathcal{F}_0] > \mathbb{E}[Y \mid \mathcal{F}_0]\} = \{\mathbb{E}[X \mid \mathcal{F}_0] - \mathbb{E}[Y \mid \mathcal{F}_0] > 0\}$ is satisfies $\mathbb{P}(A) > 0$. Notice that $A \in \mathcal{F}_0$ since $\mathbb{E}[X \mid \mathcal{F}_0]$ and $\mathbb{E}[Y \mid \mathcal{F}_0]$ are \mathcal{F}_0 measurable.

Since $X \leq Y$ implies $X - Y \leq 0$, the properties of conditional probability implies that

$$0 \geq \mathbb{E}[(X - Y) \mathbb{1}_A] = \mathbb{E}[\mathbb{E}[(X - Y) \mid \mathcal{F}_0] \mathbb{1}_A] = \mathbb{E}[(\mathbb{E}[X \mid \mathcal{F}_0] - \mathbb{E}[Y \mid \mathcal{F}_0]) \mathbb{1}_A] > 0,$$

since we have that $\mathbb{E}[X \mid \mathcal{F}_0] - \mathbb{E}[Y \mid \mathcal{F}_0] > 0$ on A . However, we arrived at a contradiction $0 > 0$, so we must have $\mathbb{P}(A) = 0$. Therefore, we must have that

$$\mathbb{E}[X \mid \mathcal{F}_0] \leq \mathbb{E}[Y \mid \mathcal{F}_0]$$

almost surely.

10. **Jensen's Inequality:** This proof is a bit too complicated for this course. Instead, we know that Jensen's inequality holds for expectation $\mathbb{E}[\phi(X)] \geq \phi(\mathbb{E}[X])$, and it is reasonable that the same result holds for conditional expectations. The proof classical proof of Jensen's uses linearity and monotonicity, which also holds for conditional expectation.