

**32.11.** 32.2↑ Absolute continuity of a set function  $\varphi$  with respect to a measure  $\mu$  is defined just as if  $\varphi$  were itself a measure:  $\mu(A) = 0$  must imply that  $\varphi(A) = 0$ . Show that, if this holds and  $\mu$  is  $\sigma$ -finite, then  $\varphi(A) = \int_A f d\mu$  for some integrable  $f$ . Show that  $A^+ = [\omega : f(\omega) \geq 0]$  and  $A^- = [\omega : f(\omega) < 0]$  give a Hahn decomposition for  $\varphi$ . Show that the three variations satisfy  $\varphi^+(A) = \int_A f^+ d\mu$ ,  $\varphi^-(A) = \int_A f^- d\mu$ , and  $|\varphi|(A) = \int_A |f| d\mu$ . Hint: To construct  $f$ , start with (32.2).

**32.12.** ↑ A *signed measure*  $\varphi$  is a set function that satisfies (32.1) if  $A_1, A_2, \dots$  are disjoint and may assume one of the values  $+\infty$  and  $-\infty$  but not both. Extend the Hahn and Jordan decompositions to signed measures

**32.13.** 31.22↑ Suppose that  $\mu$  and  $\nu$  are a probability measure and a  $\sigma$ -finite measure on the line and that  $\nu \ll \mu$ . Show that the Radon–Nikodym derivative  $f$  satisfies

$$\lim_{h \rightarrow 0} \frac{\nu(x-h, x+h)}{\mu(x-h, x+h)} = f(x)$$

on a set of  $\mu$ -measure 1.

**32.14.** Find on the unit interval uncountably many probability measures  $\mu_p$ ,  $0 < p < 1$ , with supports  $S_p$  such that  $\mu_p\{x\} = 0$  for each  $x$  and  $p$  and the  $S_p$  are disjoint in pairs.

**32.15.** Let  $\mathcal{F}_0$  be the field consisting of the finite and the cofinite sets in an uncountable  $\Omega$ . Define  $\varphi$  on  $\mathcal{F}_0$  by taking  $\varphi(A)$  to be the number of points in  $A$  if  $A$  is finite, and the negative of the number of points in  $A^c$  if  $A$  is cofinite. Show that (32.1) holds (this is not true if  $\Omega$  is countable). Show that there are no negative sets for  $\varphi$  (except the empty set), that there is no Hahn decomposition, and that  $\varphi$  does not have bounded range.

## SECTION 33. CONDITIONAL PROBABILITY

The concepts of conditional probability and expected value with respect to a  $\sigma$ -field underlie much of modern probability theory. The difficulty in understanding these ideas has to do not with mathematical detail so much as with probabilistic meaning, and the way to get at this meaning is through calculations and examples, of which there are many in this section and the next.

### The Discrete Case

Consider first the conditional probability of a set  $A$  with respect to another set  $B$ . It is defined of course by  $P(A|B) = P(A \cap B)/P(B)$ , unless  $P(B)$  vanishes, in which case it is not defined at all.

It is helpful to consider conditional probability in terms of an observer in possession of partial information.<sup>†</sup> A probability space  $(\Omega, \mathcal{F}, P)$  describes

<sup>†</sup>As always, *observer*, *information*, *know*, and so on are informal, nonmathematical terms; see the related discussion in Section 4 (p. 57).

the working of a mechanism, governed by chance, which produces a result  $\omega$  distributed according to  $P$ ;  $P(A)$  is for the observer the probability that the point  $\omega$  produced lies in  $A$ . Suppose now that  $\omega$  lies in  $B$  and that the observer learns this fact and no more. From the point of view of the observer, now in possession of this partial information about  $\omega$ , the probability that  $\omega$  also lies in  $A$  is  $P(A|B)$  rather than  $P(A)$ . This is the idea lying back of the definition.

If, on the other hand,  $\omega$  happens to lie in  $B^c$  and the observer learns of this, his probability instead becomes  $P(A|B^c)$ . These two conditional probabilities can be linked together by the simple function

$$(33.1) \quad f(\omega) = \begin{cases} P(A|B) & \text{if } \omega \in B, \\ P(A|B^c) & \text{if } \omega \in B^c. \end{cases}$$

The observer learns whether  $\omega$  lies in  $B$  or in  $B^c$ ; his new probability for the event  $\omega \in A$  is then just  $f(\omega)$ . Although the observer does not in general know the argument  $\omega$  of  $f$ , he can calculate the value  $f(\omega)$  because he knows which of  $B$  and  $B^c$  contains  $\omega$ . (Note conversely that from the value  $f(\omega)$  it is possible to determine whether  $\omega$  lies in  $B$  or in  $B^c$ , unless  $P(A|B) = P(A|B^c)$ —that is, unless  $A$  and  $B$  are independent, in which case the conditional probability coincides with the unconditional one anyway.)

The sets  $B$  and  $B^c$  partition  $\Omega$ , and these ideas carry over to the general partition. Let  $B_1, B_2, \dots$  be a finite or countable partition of  $\Omega$  into  $\mathcal{F}$ -sets, and let  $\mathcal{G}$  consist of all the unions of the  $B_i$ . Then  $\mathcal{G}$  is the  $\sigma$ -field generated by the  $B_i$ . For  $A$  in  $\mathcal{F}$ , consider the function with values

$$(33.2) \quad f(\omega) = P(A|B_i) = \frac{P(A \cap B_i)}{P(B_i)} \quad \text{if } \omega \in B_i, \quad i = 1, 2, \dots$$

If the observer learns which element  $B_i$  of the partition it is that contains  $\omega$ , then his new probability for the event  $\omega \in A$  is  $f(\omega)$ . The partition  $\{B_i\}$ , or equivalently the  $\sigma$ -field  $\mathcal{G}$ , can be regarded as an experiment, and to learn which  $B_i$  it is that contains  $\omega$  is to learn the outcome of the experiment. For this reason the function or random variable  $f$  defined by (33.2) is called the *conditional probability of  $A$  given  $\mathcal{G}$*  and is denoted  $P[A|\mathcal{G}]$ . This is written  $P[A|\mathcal{G}]_\omega$  whenever the argument  $\omega$  needs to be explicitly shown.

Thus  $P[A|\mathcal{G}]$  is the function whose value on  $B_i$  is the ordinary conditional probability  $P(A|B_i)$ . This definition needs to be completed, because  $P(A|B_i)$  is not defined if  $P(B_i) = 0$ . In this case  $P[A|\mathcal{G}]$  will be taken to have any constant value on  $B_i$ ; the value is arbitrary but must be the same over all of the set  $B_i$ . If there are nonempty sets  $B_i$  for which  $P(B_i) = 0$ ,  $P[A|\mathcal{G}]$  therefore stands for any one of a family of functions on  $\Omega$ . A specific such function is for emphasis often called a *version* of the conditional

probability. Note that any two versions are equal except on a set of probability 0.

**Example 33.1.** Consider the Poisson process. Suppose that  $0 \leq s \leq t$ , and let  $A = [N_s = 0]$  and  $B_i = [N_t = i]$ ,  $i = 0, 1, \dots$ . Since the increments are independent (Section 23),  $P(A|B_i) = P[N_s = 0]P[N_t - N_s = i]/P[N_t = i]$ , and since they have Poisson distributions (see (23.9)), a simple calculation reduces this to

$$(33.3) \quad P[N_s = 0|\mathcal{G}]_\omega = \left(1 - \frac{s}{t}\right)^i \quad \text{if } \omega \in B_i, \quad i = 0, 1, 2, \dots$$

Since  $i = N_t(\omega)$  on  $B_i$ , this can be written

$$(33.4) \quad P[N_s = 0|\mathcal{G}]_\omega = \left(1 - \frac{s}{t}\right)^{N_t(\omega)}.$$

Here the experiment or observation corresponding to  $\{B_i\}$  or  $\mathcal{G}$  determines the number of events—telephone calls, say—occurring in the time interval  $[0, t]$ . For an observer who knows this number but not the locations of the calls within  $[0, t]$ , (33.4) gives his probability for the event that none of them occurred before time  $s$ . Although this observer does not know  $\omega$ , he knows  $N_t(\omega)$ , which is all he needs to calculate the right side of (33.4). ■

**Example 33.2.** Suppose that  $X_0, X_1, \dots$  is a Markov chain with state space  $S$  as in Section 8. The events

$$(33.5) \quad [X_0 = i_0, \dots, X_n = i_n]$$

form a finite or countable partition of  $\Omega$  as  $i_0, \dots, i_n$  range over  $S$ . If  $\mathcal{G}_n$  is the  $\sigma$ -field generated by this partition, then by the defining condition (8.2) for Markov chains,  $P[X_{n+1} = j|\mathcal{G}_n]_\omega = p_{i_n j}$  holds for  $\omega$  in (33.5). The sets

$$(33.6) \quad [X_n = i]$$

for  $i \in S$  also partition  $\Omega$ , and they generate a  $\sigma$ -field  $\mathcal{G}_n^0$  smaller than  $\mathcal{G}_n$ . Now (8.2) also stipulates  $P[X_{n+1} = j|\mathcal{G}_n^0]_\omega = p_{ij}$  for  $\omega$  in (33.6), and the essence of the Markov property is that

$$(33.7) \quad P[X_{n+1} = j|\mathcal{G}_n] = P[X_{n+1} = j|\mathcal{G}_n^0]. \quad \blacksquare$$

### The General Case

If  $\mathcal{G}$  is the  $\sigma$ -field generated by a partition  $B_1, B_2, \dots$ , then the general element of  $\mathcal{G}$  is a disjoint union  $B_{i_1} \cup B_{i_2} \cup \dots$ , finite or countable, of certain of the  $B_i$ . To know which set  $B_i$  it is that contains  $\omega$  is the same thing

as to know which sets in  $\mathcal{G}$  contain  $\omega$  and which do not. This second way of looking at the matter carries over to the general  $\sigma$ -field  $\mathcal{G}$  contained in  $\mathcal{F}$ . (As always, the probability space is  $(\Omega, \mathcal{F}, P)$ .) The  $\sigma$ -field  $\mathcal{G}$  will not in general come from a partition as above.

One can imagine an observer who knows for each  $G$  in  $\mathcal{G}$  whether  $\omega \in G$  or  $\omega \in G^c$ . Thus the  $\sigma$ -field  $\mathcal{G}$  can in principle be identified with an experiment or observation. This is the point of view adopted in Section 4; see p. 57. It is natural to try and define conditional probabilities  $P[A|\mathcal{G}]$  with respect to the experiment  $\mathcal{G}$ . To do this, fix an  $A$  in  $\mathcal{F}$  and define a finite measure  $\nu$  on  $\mathcal{G}$  by

$$\nu(G) = P(A \cap G), \quad G \in \mathcal{G}.$$

Then  $P(G) = 0$  implies that  $\nu(G) = 0$ . The Radon–Nikodym theorem can be applied to the measures  $\nu$  and  $P$  on the measurable space  $(\Omega, \mathcal{G})$  because the first one is absolutely continuous with respect to the second.<sup>†</sup> It follows that there exists a function or random variable  $f$ , measurable  $\mathcal{G}$  and integrable with respect to  $P$ , such that<sup>†</sup>  $P(A \cap G) = \nu(G) = \int_G f dP$  for all  $G$  in  $\mathcal{G}$ .

Denote this function  $f$  by  $P[A|\mathcal{G}]$ . It is a random variable with two properties:

- (i)  $P[A|\mathcal{G}]$  is measurable  $\mathcal{G}$  and integrable.
- (ii)  $P[A|\mathcal{G}]$  satisfies the functional equation

$$(33.8) \quad \int_G P[A|\mathcal{G}] dP = P(A \cap G), \quad G \in \mathcal{G}.$$

There will in general be many such random variables  $P[A|\mathcal{G}]$ , but any two of them are equal with probability 1. A specific such random variable is called a *version* of the conditional probability.

If  $\mathcal{G}$  is generated by a partition  $B_1, B_2, \dots$  the function  $f$  defined by (33.2) is measurable  $\mathcal{G}$  because  $[\omega: f(\omega) \in H]$  is the union of those  $B_i$  over which the constant value of  $f$  lies in  $H$ . Any  $G$  in  $\mathcal{G}$  is a disjoint union  $G = \bigcup_k B_{i_k}$ , and  $P(A \cap G) = \sum_k P(A|B_{i_k})P(B_{i_k})$ , so that (33.2) satisfies (33.8) as well. Thus the general definition is an extension of the one for the discrete case.

Condition (i) in the definition above in effect requires that the values of  $P[A|\mathcal{G}]$  depend only on the sets in  $\mathcal{G}$ . An observer who knows the outcome of  $\mathcal{G}$  viewed as an experiment knows for each  $G$  in  $\mathcal{G}$  whether it contains  $\omega$  or not; for each  $x$  he knows this in particular for the set  $[\omega': P[A|\mathcal{G}]_{\omega'} = x]$ ,

<sup>†</sup>Let  $P_0$  be the restriction of  $P$  to  $\mathcal{G}$  (Example 10.4), and find on  $(\Omega, \mathcal{G})$  a density  $f$  for  $\nu$  with respect to  $P_0$ . Then, for  $G \in \mathcal{G}$ ,  $\nu(G) = \int_G f dP_0 = \int_G f dP$  (Example 16.4). If  $g$  is another such density, then  $P[f \neq g] = P_0[f \neq g] = 0$ .

and hence he knows in principle the functional value  $P[A|\mathcal{G}]_\omega$  even if he does not know  $\omega$  itself. In Example 33.1 a knowledge of  $N_t(\omega)$  suffices to determine the value of (33.4)— $\omega$  itself is not needed.

Condition (ii) in the definition has a gambling interpretation. Suppose that the observer, after he has learned the outcome of  $\mathcal{G}$ , is offered the opportunity to bet on the event  $A$  (unless  $A$  lies in  $\mathcal{G}$ , he does not yet know whether or not it occurred). He is required to pay an entry fee of  $P[A|\mathcal{G}]$  units and will win 1 unit if  $A$  occurs and nothing otherwise. If the observer decides to bet and pays his fee, he gains  $1 - P[A|\mathcal{G}]$  if  $A$  occurs and  $-P[A|\mathcal{G}]$  otherwise, so that his gain is

$$(1 - P[A|\mathcal{G}])I_A + (-P[A|\mathcal{G}])I_{A^c} = I_A - P[A|\mathcal{G}].$$

If he declines to bet, his gain is of course 0. Suppose that he adopts the strategy of betting if  $G$  occurs but not otherwise, where  $G$  is some set in  $\mathcal{G}$ . He can actually carry out this strategy, since after learning the outcome of the experiment  $\mathcal{G}$  he knows whether or not  $G$  occurred. His expected gain with this strategy is his gain integrated over  $G$ :

$$\int_G (I_A - P[A|\mathcal{G}]) dP.$$

But (33.8) is exactly the requirement that this vanish for each  $G$  in  $\mathcal{G}$ . Condition (ii) requires then that each strategy be fair in the sense that the observer stands neither to win nor to lose on the average. Thus  $P[A|\mathcal{G}]$  is the just entry fee, as intuition requires.

**Example 33.3.** Suppose that  $A \in \mathcal{G}$ , which will always hold if  $\mathcal{G}$  coincides with the whole  $\sigma$ -field  $\mathcal{F}$ . Then  $I_A$  satisfies conditions (i) and (ii), so that  $P[A|\mathcal{G}] = I_A$  with probability 1. If  $A \notin \mathcal{G}$ , then to know the outcome of  $\mathcal{G}$  viewed as an experiment is in particular to know whether or not  $A$  has occurred. ■

**Example 33.4.** If  $\mathcal{G}$  is  $\{\emptyset, \Omega\}$ , the smallest possible  $\sigma$ -field, every function measurable  $\mathcal{G}$  must be constant. Therefore,  $P[A|\mathcal{G}]_\omega = P(A)$  for all  $\omega$  in this case. The observer learns nothing from the experiment  $\mathcal{G}$ . ■

According to these two examples,  $P[A|\{\emptyset, \Omega\}]$  is identically  $P(A)$ , whereas  $I_A$  is a version of  $P[A|\mathcal{F}]$ . For any  $\mathcal{G}$ , the function identically equal to  $P(A)$  satisfies condition (i) in the definition of conditional probability, whereas  $I_A$  satisfies condition (ii). Condition (i) becomes more stringent as  $\mathcal{G}$  decreases, and condition (ii) becomes more stringent as  $\mathcal{G}$  increases. The two conditions work in opposite directions and between them delimit the class of versions of  $P[A|\mathcal{G}]$ .

**Example 33.5.** Let  $\Omega$  be the plane  $R^2$  and let  $\mathcal{F}$  be the class  $\mathcal{R}^2$  of planar Borel sets. A point of  $\Omega$  is a pair  $(x, y)$  of reals. Let  $\mathcal{G}$  be the  $\sigma$ -field consisting of the vertical strips, the product sets  $E \times R^1 = [(x, y) : x \in E]$ , where  $E$  is a linear Borel set. If the observer knows for each strip  $E \times R^1$  whether or not it contains  $(x, y)$ , then, as he knows this for each one-point set  $E$ , he knows the value of  $x$ . Thus the experiment  $\mathcal{G}$  consists in the determination of the first coordinate of the sample point. Suppose now that  $P$  is a probability measure on  $\mathcal{R}^2$  having a density  $f(x, y)$  with respect to planar Lebesgue measure:  $P(A) = \iint_A f(x, y) dx dy$ . Let  $A$  be a horizontal strip  $R^1 \times F = [(x, y) : y \in F]$ ,  $F$  being a linear Borel set. The conditional probability  $P[A|\mathcal{G}]$  can be calculated explicitly.

Put

$$(33.9) \quad \varphi(x, y) = \frac{\int_F f(x, t) dt}{\int_{R^1} f(x, t) dt}.$$

Set  $\varphi(x, y) = 0$ , say, at points where the denominator here vanishes; these points form a set of  $P$ -measure 0. Since  $\varphi(x, y)$  is a function of  $x$  alone, it is measurable  $\mathcal{G}$ . The general element of  $\mathcal{G}$  being  $E \times R^1$ , it will follow that  $\varphi$  is a version of  $P[A|\mathcal{G}]$  if it is shown that

$$(33.10) \quad \int_{E \times R^1} \varphi(x, y) dP(x, y) = P(A \cap (E \times R^1)).$$

Since  $A = R^1 \times F$ , the right side here is  $P(E \times F)$ . Since  $P$  has density  $f$ , Theorem 16.11 and Fubini's theorem reduce the left side to

$$\begin{aligned} \int_E \left\{ \int_{R^1} \varphi(x, y) f(x, y) dy \right\} dx &= \int_E \left\{ \int_F f(x, t) dt \right\} dx \\ &= \iint_{E \times F} f(x, y) dx dy = P(E \times F). \end{aligned}$$

Thus (33.9) does give a version of  $P[R^1 \times F|\mathcal{G}]$ . ■

The right side of (33.9) is the classical formula for the conditional probability of the event  $R^1 \times F$  (the event that  $y \in F$ ) given the event  $\{x\} \times R^1$  (given the value of  $x$ ). Since the event  $\{x\} \times R^1$  has probability 0, the formula  $P(A|B) = P(A \cap B)/P(B)$  does not work here. The whole point of this section is the systematic development of a notion of conditional probability that covers conditioning with respect to events of probability 0. This is accomplished by conditioning with respect to *collections* of events—that is, with respect to  $\sigma$ -fields  $\mathcal{G}$ .

**Example 33.6.** The set  $A$  is by definition independent of the  $\sigma$ -field  $\mathcal{G}$  if it is independent of each  $G$  in  $\mathcal{G}$ :  $P(A \cap G) = P(A)P(G)$ . This being the same thing as  $P(A \cap G) = \int_G P(A) dP$ ,  $A$  is independent of  $\mathcal{G}$  if and only if  $P[A|\mathcal{G}] = P(A)$  with probability 1. ■

The  $\sigma$ -field  $\sigma(X)$  generated by a random variable  $X$  consists of the sets  $[\omega: X(\omega) \in H]$  for  $H \in \mathcal{R}^1$ ; see Theorem 20.1. The conditional probability of  $A$  given  $X$  is defined as  $P[A|\sigma(X)]$  and is denoted  $P[A|X]$ . Thus  $P[A|X] = P[A|\sigma(X)]$  by definition. From the experiment corresponding to the  $\sigma$ -field  $\sigma(X)$ , one learns which of the sets  $[\omega': X(\omega') = x]$  contains  $\omega$  and hence learns the value  $X(\omega)$ . Example 33.5 is a case of this: take  $X(x, y) = x$  for  $(x, y)$  in the sample space  $\Omega = \mathbb{R}^2$  there.

This definition applies without change to random vector, or, equivalently, to a finite set of random variables. It can be adapted to arbitrary sets of random variables as well. For any such set  $[X_t, t \in T]$ , the  $\sigma$ -field  $\sigma[X_t, t \in T]$  it generates is the smallest  $\sigma$ -field with respect to which each  $X_t$  is measurable. It is generated by the collection of sets of the form  $[\omega: X_t(\omega) \in H]$  for  $t$  in  $T$  and  $H$  in  $\mathcal{R}^1$ . The *conditional probability*  $P[A|X_t, t \in T]$  of  $A$  with respect to this set of random variables is by definition the conditional probability  $P[A|\sigma[X_t, t \in T]]$  of  $A$  with respect to the  $\sigma$ -field  $\sigma[X_t, t \in T]$ .

In this notation the property (33.7) of Markov chains becomes

$$(33.11) \quad P[X_{n+1} = j | X_0, \dots, X_n] = P[X_{n+1} = j | X_n].$$

The conditional probability of  $[X_{n+1} = j]$  is the same for someone who knows the present state  $X_n$  as for someone who knows the present state  $X_n$  and the past states  $X_0, \dots, X_{n-1}$  as well.

**Example 33.7.** Let  $X$  and  $Y$  be random vectors of dimensions  $j$  and  $k$ , let  $\mu$  be the distribution of  $X$  over  $\mathbb{R}^j$ , and suppose that  $X$  and  $Y$  are independent. According to (20.30),

$$P[X \in H, (X, Y) \in J] = \int_H P[(x, Y) \in J] \mu(dx)$$

for  $H \in \mathcal{R}^j$  and  $J \in \mathcal{R}^{j+k}$ . This is a consequence of Fubini's theorem; it has a conditional-probability interpretation. For each  $x$  in  $\mathbb{R}^j$  put

$$(33.12) \quad f(x) = P[(x, Y) \in J] = P[\omega': (x, Y(\omega')) \in J].$$

By Theorem 20.1(ii),  $f(X(\omega))$  is measurable  $\sigma(X)$ , and since  $\mu$  is the distribution of  $X$ , a change of variable gives

$$\int_{[X \in H]} f(X(\omega)) P(d\omega) = \int_H f(x) \mu(dx) = P([(X, Y) \in J] \cap [X \in H]).$$

Since  $[X \in H]$  is the general element of  $\sigma(X)$ , this proves that

$$(33.13) \quad f(X(\omega)) = P[(X, Y) \in J | X]_{\omega}$$

with probability 1. ■

The fact just proved can be written

$$\begin{aligned} P[(X, Y) \in J | X]_{\omega} &= P[(X(\omega), Y) \in J] \\ &= P[\omega' : (X(\omega), Y(\omega')) \in J]. \end{aligned}$$

Replacing  $\omega'$  by  $\omega$  on the right here causes a notational collision like the one replacing  $y$  by  $x$  causes in  $\int_a^b f(x, y) dy$ .

Suppose that  $X$  and  $Y$  are independent random variables and that  $Y$  has distribution function  $F$ . For  $J = [(u, v) : \max\{u, v\} \leq m]$ , (33.12) is 0 for  $m < x$  and  $F(m)$  for  $m \geq x$ ; if  $M = \max\{X, Y\}$ , then (33.13) gives

$$(33.14) \quad P[M \leq m | X]_{\omega} = I_{[X \leq m]}(\omega) F(m)$$

with probability 1. All equations involving conditional probabilities must be qualified in this way by the phrase *with probability 1*, because the conditional probability is unique only to within a set of probability 0.

The following theorem is useful for checking conditional probabilities.

**Theorem 33.1.** *Let  $\mathcal{P}$  be a  $\pi$ -system generating the  $\sigma$ -field  $\mathcal{G}$ , and suppose that  $\Omega$  is a finite or countable union of sets in  $\mathcal{P}$ . An integrable function  $f$  is a version of  $P[A | \mathcal{G}]$  if it is measurable  $\mathcal{G}$  and if*

$$(33.15) \quad \int_G f dP = P(A \cap G)$$

*holds for all  $G$  in  $\mathcal{P}$ .*

**PROOF.** Apply Theorem 10.4. ■

The condition that  $\Omega$  is a finite or countable union of  $\mathcal{P}$ -sets cannot be suppressed; see Example 10.5.

**Example 33.8.** Suppose that  $X$  and  $Y$  are independent random variables with a common distribution function  $F$  that is positive and continuous. What is the conditional probability of  $[X \leq x]$  given the random variable  $M = \max\{X, Y\}$ ? As it should clearly be 1 if  $M \leq x$ , suppose that  $M > x$ . Since  $X \leq x$  requires  $M = Y$ , the chance of which is  $\frac{1}{2}$  by symmetry, the conditional probability of  $[X \leq x]$  should by independence be  $\frac{1}{2}F(x)/F(m) = \frac{1}{2}P[X \leq x | X \leq m]$  with the random variable  $M$  substituted

for  $m$ . Intuition thus gives

$$(33.16) \quad P[X \leq x | M]_{\omega} = I_{[M \leq x]}(\omega) + \frac{1}{2} I_{[M > x]}(\omega) \frac{F(x)}{F(M(\omega))}.$$

It suffices to check (33.15) for sets  $G = [M \leq m]$ , because these form a  $\pi$ -system generating  $\sigma(M)$ . The functional equation reduces to

$$(33.17) \quad P[M \leq \min\{x, m\}] + \frac{1}{2} \int_{x < M \leq m} \frac{F(x)}{F(M)} dP = P[M \leq m, X \leq x].$$

Since the other case is easy, suppose that  $x < m$ . Since the distribution of  $(X, Y)$  is product measure, it follows by Fubini's theorem and the assumed continuity of  $F$  that

$$\begin{aligned} \int_{x < M \leq m} \frac{1}{F(M)} dP &= \iint_{\substack{u \leq t \\ x < u \leq m}} \frac{1}{F(v)} dF(u) dF(v) \\ &\quad + \iint_{\substack{t < u \\ x < u \leq m}} \frac{1}{F(u)} dF(u) dF(v) = 2(F(m) - F(x)), \end{aligned}$$

which gives (33.17). ■

**Example 33.9.** A collection  $[X_t : t \geq 0]$  of random variables is a *Markov process in continuous time* if for  $k \geq 1$ ,  $0 \leq t_1 \leq \dots \leq t_k \leq u$ , and  $H \in \mathcal{R}^1$ ,

$$(33.18) \quad P[X_u \in H | X_{t_1}, \dots, X_{t_k}] = P[X_u \in H | X_{t_k}]$$

holds with probability 1. The analogue for discrete time is (33.11). (The  $X_n$  there have countable range as well, and the transition probabilities are constant in time, conditions that are not imposed here.)

Suppose that  $t \leq u$ . Looking on the right side of (33.18) as a version of the conditional probability on the left shows that

$$(33.19) \quad \int_G P[X_u \in H | X_t] dP = P([X_u \in H] \cap G)$$

if  $0 \leq t_1 \leq \dots \leq t_k = t \leq u$  and  $G \in \sigma(X_{t_1}, \dots, X_{t_k})$ . Fix  $t$ ,  $u$ , and  $H$ , and let  $k$  and  $t_1, \dots, t_k$  vary. Consider the class  $\mathcal{F} = \bigcup \sigma(X_{t_1}, \dots, X_{t_k})$ , the union extending over all  $k \geq 1$  and all  $k$ -tuples satisfying  $0 \leq t_1 \leq \dots \leq t_k = t$ . If  $A \in \sigma(X_{t_1}, \dots, X_{t_k})$  and  $B \in \sigma(X_{s_1}, \dots, X_{s_j})$ , then  $A \cap B \in \sigma(X_{r_1}, \dots, X_{r_j})$ , where the  $r_\alpha$  are the  $s_\beta$  and the  $t_\gamma$  merged together. Thus  $\mathcal{F}$  is a  $\pi$ -system. Since  $\mathcal{F}$  generates  $\sigma[X_s : s \leq t]$  and  $P[X_u \in H | X_t]$  is measurable with respect to this  $\sigma$ -field, it follows by (33.19) and Theorem 33.1 that  $P[X_u \in H | X_t]$  is a version of  $P[X_u \in H | X_s, s \leq t]$ :

$$(33.20) \quad P[X_u \in H | X_s, s \leq t] = P[X_u \in H | X_t], \quad t \leq u,$$

with probability 1.

This says that for calculating conditional probabilities about the future, the present  $\sigma(X_t)$  is equivalent to the present and the *entire* past  $\sigma[X_s: s \leq t]$ . This follows from the apparently weaker condition (33.18). ■

**Example 33.10.** The Poisson process  $[N_t: t \geq 0]$  has independent increments (Section 23). Suppose that  $0 \leq t_1 \leq \dots \leq t_k \leq u$ . The random vector  $(N_{t_1}, N_{t_2} - N_{t_1}, \dots, N_{t_k} - N_{t_{k-1}})$  is independent of  $N_u - N_{t_k}$ , and so (Theorem 20.2)  $(N_{t_1}, N_{t_2}, \dots, N_{t_k})$  is independent of  $N_u - N_{t_k}$ . If  $J$  is the set of points  $(x_1, \dots, x_k, y)$  in  $R^{k+1}$  such that  $x_k + y \in H$ , where  $H \in \mathcal{R}^1$ , and if  $\nu$  is the distribution of  $N_u - N_{t_k}$ , then (33.12) is  $P[(x_1, \dots, x_k, N_u - N_{t_k}) \in J] = P[x_k + N_u - N_{t_k} \in H] = \nu(H - x_k)$ . Therefore, (33.13) gives  $P[N_u \in H | N_{t_1}, \dots, N_{t_k}] = \nu(H - N_{t_k})$ . This holds also if  $k = 1$ , and hence  $P[N_u \in H | N_{t_1}, \dots, N_{t_k}] = P[N_u \in H | N_{t_k}]$ . The Poisson process thus has the Markov property (33.18); this is a consequence solely of the independence of the increments. The extended Markov property (33.20) follows. ■

### Properties of Conditional Probability

**Theorem 33.2.** *With probability 1,  $P[\emptyset | \mathcal{G}] = 0$ ,  $P[\Omega | \mathcal{G}] = 1$ ; and*

$$(33.21) \quad 0 \leq P[A | \mathcal{G}] \leq 1$$

for each  $A$ . If  $A_1, A_2, \dots$  is a finite or countable sequence of disjoint sets, then

$$(33.22) \quad P\left[\bigcup_n A_n | \mathcal{G}\right] = \sum_n P[A_n | \mathcal{G}]$$

with probability 1.

**PROOF.** For each version of the conditional probability,  $\int_G P[A | \mathcal{G}] dP = P(A \cap G) \geq 0$  for each  $G$  in  $\mathcal{G}$ ; since  $P[A | \mathcal{G}]$  is measurable  $\mathcal{G}$ , it must be nonnegative except on a set of  $P$ -measure 0. The other inequality in (33.21) is proved the same way.

If the  $A_n$  are disjoint and if  $G$  lies in  $\mathcal{G}$ , it follows (Theorem 16.6) that

$$\begin{aligned} \int_G \left( \sum_n P[A_n | \mathcal{G}] \right) dP &= \sum_n \int_G P[A_n | \mathcal{G}] dP = \sum_n P(A_n \cap G) \\ &= P\left(\left(\bigcup_n A_n\right) \cap G\right). \end{aligned}$$

Thus  $\sum_n P[A_n | \mathcal{G}]$ , which is certainly measurable  $\mathcal{G}$ , satisfies the functional equation for  $P[\bigcup_n A_n | \mathcal{G}]$ , and so must coincide with it except perhaps on a set of  $P$ -measure 0. Hence (33.22). ■

Additional useful facts can be established by similar arguments. If  $A \subset B$ , then

$$(33.23) \quad P[B - A \parallel \mathcal{G}] = P[B \parallel \mathcal{G}] - P[A \parallel \mathcal{G}], \quad P[A \parallel \mathcal{G}] \leq P[B \parallel \mathcal{G}].$$

The inclusion-exclusion formula

$$(33.24) \quad P\left[\bigcup_{i=1}^n A_i \parallel \mathcal{G}\right] = \sum_i P[A_i \parallel \mathcal{G}] - \sum_{i < j} P[A_i \cap A_j \parallel \mathcal{G}] + \dots$$

holds. If  $A_n \uparrow A$ , then

$$(33.25) \quad P[A_n \parallel \mathcal{G}] \uparrow P[A \parallel \mathcal{G}],$$

and if  $A_n \downarrow A$ , then

$$(33.26) \quad P[A_n \parallel \mathcal{G}] \downarrow P[A \parallel \mathcal{G}].$$

Further,  $P(A) = 1$  implies that

$$(33.27) \quad P[A \parallel \mathcal{G}] = 1,$$

and  $P(A) = 0$  implies that

$$(33.28) \quad P[A \parallel \mathcal{G}] = 0.$$

Of course (33.23) through (33.28) hold with probability 1 only.

### Difficulties and Curiosities

This section has been devoted almost entirely to examples connecting the abstract definition (33.8) with the probabilistic idea lying back of it. There are pathological examples showing that the interpretation of conditional probability in terms of an observer with partial information breaks down in certain cases.

**Example 33.11.** Let  $(\Omega, \mathcal{F}, P)$  be the unit interval  $\Omega$  with Lebesgue measure  $P$  on the  $\sigma$ -field  $\mathcal{F}$  of Borel subsets of  $\Omega$ . Take  $\mathcal{G}$  to be the  $\sigma$ -field of sets that are either countable or cocountable. Then the function identically equal to  $P(A)$  is a version of  $P[A \parallel \mathcal{G}]$ : (33.8) holds because  $P(G)$  is either 0 or 1 for every  $G$  in  $\mathcal{G}$ . Therefore,

$$(33.29) \quad P[A \parallel \mathcal{G}]_\omega = P(A)$$

with probability 1. But since  $\mathcal{G}$  contains all one-point sets, to know which

elements of  $\mathcal{G}$  contain  $\omega$  is to know  $\omega$  itself. Thus  $\mathcal{G}$  viewed as an experiment should be completely informative—the observer given the information in  $\mathcal{G}$  should know  $\omega$  exactly—and so one might expect that

$$(33.30) \quad P[A|\mathcal{G}]_\omega = \begin{cases} 1 & \text{if } \omega \in A, \\ 0 & \text{if } \omega \notin A. \end{cases}$$

This is Example 4.10 in a new form. ■

The mathematical definition gives (33.29); the heuristic considerations lead to (33.30). Of course, (33.29) is right and (33.30) is wrong. The heuristic view breaks down in certain cases but is nonetheless illuminating and cannot, since it does not intervene in proofs, lead to any difficulties.

The point of view in this section has been “global.” To each fixed  $A$  in  $\mathcal{F}$  has been attached a function (usually a family of functions)  $P[A|\mathcal{G}]_\omega$  defined over all of  $\Omega$ . What happens if the point of view is reversed—if  $\omega$  is fixed and  $A$  varies over  $\mathcal{F}$ ? Will this result in a probability measure on  $\mathcal{F}$ ? Intuition says it should, and if it does, then (33.21) through (33.28) all reduce to standard facts about measures.

Suppose that  $B_1, \dots, B_r$  is a partition of  $\Omega$  into  $\mathcal{F}$ -sets, and let  $\mathcal{G} = \sigma(B_1, \dots, B_r)$ . If  $P(B_1) = 0$  and  $P(B_i) > 0$  for the other  $i$ , then one version of  $P[A|\mathcal{G}]$  is

$$P[A|\mathcal{G}]_\omega = \begin{cases} 1 & \text{if } \omega \in B_1, \\ \frac{P(A \cap B_i)}{P(B_i)} & \text{if } \omega \in B_i, i = 2, \dots, r. \end{cases}$$

With this choice of version for each  $A$ ,  $P[A|\mathcal{G}]_\omega$  is, as a function of  $A$ , a probability measure on  $\mathcal{F}$  if  $\omega \in B_2 \cup \dots \cup B_r$ , but not if  $\omega \in B_1$ . The “wrong” versions have been chosen. If, for example,

$$P[A|\mathcal{G}]_\omega = \begin{cases} P(A) & \text{if } \omega \in B_1, \\ \frac{P(A \cap B_i)}{P(B_i)} & \text{if } \omega \in B_i, i = 2, \dots, r, \end{cases}$$

then  $P[A|\mathcal{G}]_\omega$  is a probability measure in  $A$  for each  $\omega$ . Clearly, versions such as this one exist if  $\mathcal{G}$  is finite.

It might be thought that for an arbitrary  $\sigma$ -field  $\mathcal{G}$  in  $\mathcal{F}$  versions of the various  $P[A|\mathcal{G}]$  can be so chosen that  $P[A|\mathcal{G}]_\omega$  is for each fixed  $\omega$  a probability measure as  $A$  varies over  $\mathcal{F}$ . It is possible to construct a

counterexample showing that this is not so.<sup>†</sup> The example is possible because the exceptional  $\omega$ -set of probability 0 where (33.22) fails depends on the sequence  $A_1, A_2, \dots$ ; if there are uncountably many such sequences, it can happen that the union of these exceptional sets has positive probability whatever versions  $P[A \parallel \mathcal{G}]$  are chosen.

The existence of such pathological examples turns out not to matter. Example 33.9 illustrates the reason why. From the assumption (33.18) the notably stronger conclusion (33.20) was reached. Since the set  $[X_u \in H]$  is fixed throughout the argument, it does not matter that conditional probabilities may not, in fact, be measures. What does matter for the theory is Theorem 33.2 and its extensions.

Consider a point  $\omega_0$  with the property that  $P(G) > 0$  for every  $G$  in  $\mathcal{G}$  that contains  $\omega_0$ . This will be true if the one-point set  $\{\omega_0\}$  lies in  $\mathcal{F}$  and has positive probability. Fix any versions of the  $P[A \parallel \mathcal{G}]$ . For each  $A$  the set  $\{\omega : P[A \parallel \mathcal{G}]_\omega < 0\}$  lies in  $\mathcal{G}$  and has probability 0; it therefore cannot contain  $\omega_0$ . Thus  $P[A \parallel \mathcal{G}]_{\omega_0} \geq 0$ . Similarly,  $P[\Omega \parallel \mathcal{G}]_{\omega_0} = 1$ , and, if the  $A_n$  are disjoint,  $P[\bigcup_n A_n \parallel \mathcal{G}]_{\omega_0} = \sum_n P[A \parallel \mathcal{G}]_{\omega_0}$ . Therefore,  $P[A \parallel \mathcal{G}]_{\omega_0}$  is a probability measure as  $A$  ranges over  $\mathcal{F}$ .

Thus conditional probabilities behave like probabilities at points of positive probability. That they may not do so at points of probability 0 causes no problem because individual such points have no effect on the probabilities of sets. Of course, sets of points individually having probability 0 do have an effect, but here the global point of view reenters.

### Conditional Probability Distributions

Let  $X$  be a random variable on  $(\Omega, \mathcal{F}, P)$ , and let  $\mathcal{G}$  be a  $\sigma$ -field in  $\mathcal{F}$ .

**Theorem 33.3.** *There exists a function  $\mu(H, \omega)$ , defined for  $H$  in  $\mathcal{R}^1$  and  $\omega$  in  $\Omega$ , with these two properties:*

- (i) *For each  $\omega$  in  $\Omega$ ,  $\mu(\cdot, \omega)$  is a probability measure on  $\mathcal{R}^1$ .*
- (ii) *For each  $H$  in  $\mathcal{R}^1$ ,  $\mu(H, \cdot)$  is a version of  $P[X \in H \parallel \mathcal{G}]$ .*

The probability measure  $\mu(\cdot, \omega)$  is a *conditional distribution* of  $X$  given  $\mathcal{G}$ . If  $\mathcal{G} = \sigma(Z)$ , it is a conditional distribution of  $X$  given  $Z$ .

**PROOF.** For each rational  $r$ , let  $F(r, \omega)$  be a version of  $P[X \leq r \parallel \mathcal{G}]_\omega$ . If  $r \leq s$ , then by (33.23),

$$(33.31) \quad F(r, \omega) \leq F(s, \omega)$$

<sup>†</sup>The argument is outlined in Problem 33.11. It depends on the construction of certain nonmeasurable sets.

for  $\omega$  outside a  $\mathcal{G}$ -set  $A_{rs}$  of probability 0. By (33.26),

$$(33.32) \quad F(r, \omega) = \lim_n F(r + n^{-1}, \omega)$$

for  $\omega$  outside a  $\mathcal{G}$ -set  $B_r$  of probability 0. Finally, by (33.25) and (33.26),

$$(33.33) \quad \lim_{r \rightarrow -\infty} F(r, \omega) = 0, \quad \lim_{r \rightarrow \infty} F(r, \omega) = 1$$

outside a  $\mathcal{G}$ -set  $C$  of probability 0. As there are only countably many of these exceptional sets, their union  $E$  lies in  $\mathcal{G}$  and has probability 0.

For  $\omega \notin E$  extend  $F(\cdot, \omega)$  to all of  $R^1$  by setting  $F(x, \omega) = \inf\{F(r, \omega) : x < r\}$ . For  $\omega \in E$  take  $F(x, \omega) = F(x)$ , where  $F$  is some arbitrary but fixed distribution function. Suppose that  $\omega \notin E$ . By (33.31) and (33.32),  $F(x, \omega)$  agrees with the first definition on the rationals and is nondecreasing; it is right-continuous; and by (33.33) it is a probability distribution function. Therefore, there exists a probability measure  $\mu(\cdot, \omega)$  on  $(R^1, \mathcal{R}^1)$  with distribution function  $F(\cdot, \omega)$ . For  $\omega \in E$ , let  $\mu(\cdot, \omega)$  be the probability measure corresponding to  $F(x)$ . Then condition (i) is satisfied.

The class of  $H$  for which  $\mu(H, \cdot)$  is measurable  $\mathcal{G}$  is a  $\lambda$ -system containing the sets  $H = (-\infty, r]$  for rational  $r$ ; therefore  $\mu(H, \cdot)$  is measurable  $\mathcal{G}$  for  $H$  in  $\mathcal{R}^1$ .

By construction,  $\mu((-\infty, r], \omega) = P[X \leq r \mid \mathcal{G}]_\omega$  with probability 1 for rational  $r$ ; that is, for  $H = (-\infty, r]$  as well as for  $H = R^1$ ,

$$\int_G \mu(H, \omega) P(d\omega) = P([X \in H] \cap G)$$

for all  $G$  in  $\mathcal{G}$ . Fix  $G$ . Each side of this equation is a measure as a function of  $H$ , and so the equation must hold for all  $H$  in  $\mathcal{R}^1$ . ■

**Example 33.12.** Let  $X$  and  $Y$  be random variables whose joint distribution  $\nu$  in  $R^2$  has density  $f(x, y)$  with respect to Lebesgue measure:  $P[(X, Y) \in A] = \nu(A) = \iint_A f(x, y) dx dy$ . Let  $g(x, y) = f(x, y) / \int_{R^1} f(x, t) dt$ , and let  $\mu(H, x) = \int_H g(x, y) dy$  have probability density  $g(x, \cdot)$ ; if  $\int_{R^1} f(x, t) dt = 0$ , let  $\mu(\cdot, x)$  be an arbitrary probability measure on the line. Then  $\mu(H, X(\omega))$  will serve as the conditional distribution of  $Y$  given  $X$ . Indeed, (33.10) is the same thing as  $\int_{E \times R^1} \mu(F, x) d\nu(x, y) = \nu(E \times F)$ , and a change of variable gives  $\int_{[X \in E]} \mu(F, X(\omega)) P(d\omega) = P[X \in E, Y \in F]$ . Thus  $\mu(F, X(\omega))$  is a version of  $P[Y \in F \mid X]_\omega$ . This is a new version of Example 33.5. ■