

# You Can't Always Get What You Want\*

## Some considerations regarding conditional probabilities

Wayne C. Myrvold  
Department of Philosophy  
The University of Western Ontario  
wmyrvold@uwo.ca

June 13, 2014

### **Abstract**

The standard treatment of conditional probability leaves conditional probability undefined when the conditioning proposition has zero probability. Nonetheless, some find the option of extending the scope of conditional probability to include zero-probability conditions attractive or even compelling. This article reviews some of the pitfalls associated with this move, and concludes that, for the most part, probabilities conditional on zero-probability propositions are more trouble than they are worth.

---

\*But if you try, sometimes, you might find you get what you need.

# 1 Introduction

Let  $\mathcal{A}$  be a set of propositions, closed under Boolean operations, let  $P$  be a probability function on  $\mathcal{A}$ , and, for some proposition  $C$ , let  $P_C$  be another probability function on  $\mathcal{A}$ , to be thought of as yielding probabilities conditional on  $C$ .<sup>1</sup> It is uncontroversial that, if  $C$  is in  $\mathcal{A}$ , these should satisfy

$$P(AC) = P_C(A) P(C). \tag{1}$$

If  $P(C) > 0$ , then the unconditional probability function  $P$ , together with the requirement that (1) hold, uniquely determines  $P_C(A)$ , for any  $A \in \mathcal{A}$ :

$$P_C(A) = \frac{P(AC)}{P(C)}. \tag{2}$$

If, however,  $P(C) = 0$ , then  $P(AC)$  is also equal to zero, and (1) is satisfied for any value whatsoever of  $P_C(A)$ , and so (1) leaves  $P_C(A)$  completely undetermined.

One reaction, the standard one, is to leave  $P_C(A)$  undefined except for  $C \in \mathcal{A}$  with  $P(C) > 0$ . But, since (1) places no constraints whatsoever on the function  $P_C$  when  $P(C) = 0$ , for such propositions we are free, without fear of violating this condition, to define  $P_C$  to be any probability function whatsoever on  $\mathcal{A}$ . Instead of relying on (2) to define conditional probability functions in terms of the unconditional probability function  $P$ , we can take conditional probability as primitive. This is a route that has been recommended by a number of authors over the years (see, *e.g.*, Rényi 1955, Popper 1938, 1955, 1959, Carnap 1950, Harper 1975, Harper and Hájek 1997, Hájek 2003, van Fraassen 1976, and Dorr 2010). In support of this, cases are sometimes adduced that suggest that there are probabilities conditional on zero-probability propositions that have clearly defined values (see §2, below). Moreover, it might seem that we *have* to regard some probabilities conditional on zero-probability propositions as well-defined, in order to do justice to statistical practice, since statistical practice invokes likelihood functions, which ascribe probabilities to data as a function of some continuously varying parameter, and these are well-defined for all parameter values even if every point value of the parameter is ascribed zero probability. We do not want to eschew the use of such functions; does this not commit us to probabilities conditional on zero-probability propositions?

In this essay, I hope to convince the reader that things are not so straightforward. The examples that purport to show that there are clear-cut answers to requests for probabilities conditional on propositions of probability zero are misleading. We can *give* such questions answers by requiring that the conditional probability functions possess certain symmetry properties, but this is our choice, not dictated by the nature of the problem, and we should not let the intuitive appeal of such symmetry properties blind us to the fact that we must stipulate that the conditional probabilities have

---

<sup>1</sup>We will also use the notation  $P(A|C)$ , when convenient.

them, in order for the questions to acquire determinate answers. Moreover, there will be cases in which symmetry conditions that we may wish to impose will clash with each other, or may clash with the *desideratum* of countable additivity, illustrating Jagger’s Theorem: You Can’t Always Get What You Want.<sup>2</sup>

Furthermore, the consequences of taking the standard route, and leaving undefined probabilities conditional on null propositions (that is, propositions of unconditional probability zero), are not as dire as some would make them out to be. Though there are cases (such as the likelihood functions already mentioned) in which quantities appear that can unproblematically be taken to be probabilities conditional on null propositions, they need not be. We can take them to be no more than auxiliary functions, useful for calculating probabilities. The theory goes along straightforwardly if we take all conditional probabilities to have conditions with positive probability.

Some will be undaunted, and will insist on introducing a host of null-condition conditional probabilities. This can be done, but, if it is done, it should be *done* and not merely gestured at: those who invoke probabilities conditional on null propositions should specify which pairs of propositions  $A$ ,  $B$  they take the conditional probability  $P_B(A)$  to be defined for, and specify the values of these conditional probabilities.

A comment, before we begin, on the word “probability.” It has become common, in the philosophical literature on probability, to acknowledge that the word is used in (at least) two senses. There is an epistemic sense, having to do with degrees of belief, and a physical sense, having to do with characteristics of chance set-ups. For the most part, our considerations will bear equally on probability in either sense. When it matters, we will use the “credence” for the former, epistemic sense, “chance” for the physical sense, and “probability” when we want to be neutral between the two senses.

## 2 Examples

Consider the following examples.

---

<sup>2</sup> The fact that, in probability theory, we can’t always get what we want, is a familiar fact. We might want our probability function to be defined on arbitrary subsets of our probability space, but, as is well-known, we can’t always do so while satisfying *desiderata* such as symmetry conditions and countable additivity. Consider, for example, the task of defining a uniform distribution—that is, a distribution invariant under all rotations—on the unit circle. There can be no distribution that is invariant under rotations, is countably additive, and is defined on all subsets of the unit circle. The proof is found in many probability texts, *e.g.* Billingsley (2012, p. 47). The standard response is to preserve countable additivity and to restrict the domain of definition of the probability function to certain subsets of the probability space, the measurable sets, leaving the probability of other sets undefined. In one and two dimensions, as Banach (1923) showed, one can extend the probability function to one defined on arbitrary subsets, if one is willing to give up countable additivity. The well-known Banach-Tarski paradox shows that we can’t do so in three-dimensional space; there is no finitely additive set function that is defined on all subsets and invariant under translations and rotations.

**Example 1.**<sup>3</sup> A number is chosen, with uniform probability, from the interval  $[0, 1]$ . Conditional on the supposition that the chosen number is either  $1/4$  or  $3/4$ , what is the probability that it is  $1/4$ ?

**Example 2.** (Borel-Kolmogorov).<sup>4</sup> A point is chosen, with uniform probability, on the surface of the earth, which we treat as a perfect sphere.

- a). What is the probability that the chosen point is in the Western Hemisphere, given that that it lies on the equator?
- b). Conditional on the chosen point lying on the great circle containing the Greenwich meridian, what is the probability that it lies closer to the equator than to a pole?

**Example 3.** A number is chosen, with uniform probability, from the interval  $[0, 1]$ . Conditional on the supposition that the chosen number is rational, what is the probability that it is greater than  $1/2$ ?

For many, perhaps most, readers, each of the above questions will have an obvious answer. This should give us pause.

In each case, the set-up of the problem mentions a uniform probability distribution. There is a standard way of making this mathematically precise, insofar as unconditional probabilities are concerned. These unconditional probabilities determine probabilities conditional on propositions of nonzero probability. They do not determine conditional probabilities when the conditioning propositions have probability zero. We may extend our domain of conditioning to include some of these, but there is no canonical way of doing so. Moreover, when we venture into to the project of extending the domain of conditioning, if we are not careful, we run the risk of imposing conditions that seem to be intuitively compelling, but which cannot be jointly satisfied. We must proceed with caution, if we proceed at all.

## 3 Symmetry Conditions

### 3.1 Example 1.

Example 1 seems beguilingly simple. It may seem that the symmetry of the problem dictates the answer  $1/2$ , on pain of irrationality. Nothing at all in the set-up of the problem favours either  $1/4$  or  $3/4$ .

But consider this variant on the question. Suppose that the number is chosen from the unit interval, not with uniform distribution, but according to a distribution given by the density function

$$f(x) = 2x. \tag{3}$$

---

<sup>3</sup>Adapted from Hájek (2003).

<sup>4</sup>Based on Kolmogorov (1950, §V.2), which in turn is based on Borel (1909, §45) (§8.6 of Borel 1965). See also Jaynes (2003, §15.7), Hájek (2003, §4.4).

Now ask the question: conditional on the number chosen being either  $1/4$  or  $3/4$ , what is the probability that it is  $1/4$ ?

Here, I suspect, intuitions will vary. To some, the answer might still be, obviously,  $1/2$ . Others, reflecting on the fact that the number chosen is more likely to be greater than  $1/2$  than less than  $1/2$ , might regard  $3/4$  as the more probable value. This intuition can be given a numerical value by considering, that, for any sufficiently small positive  $\epsilon$ ,

$$\frac{Pr(X \in [\frac{3}{4} - \epsilon, \frac{3}{4} + \epsilon])}{Pr(X \in [\frac{1}{4} - \epsilon, \frac{1}{4} + \epsilon])} = 3, \quad (4)$$

which suggests that the number  $3/4$  is 3 times as probable as  $1/4$ .

Suppose, now, that we change the question only slightly, and ask: if the number is chosen from the unit interval according to a distribution with density (3), what is the probability, conditional on the number chosen being either  $1/2$  or  $\sqrt{3}/2$ , that it is  $1/2$ ? Similar considerations suggest that  $\sqrt{3}/2$  is more probable than  $1/2$ .

If we give this answer, we have thereby achieved incoherence, because this last question is just our first question rephrased.  $X$  being uniformly distributed on the unit interval is the same as  $\sqrt{X}$  being distributed with density (3), and so we have two ways of asking one and the same question. Taking  $Y = \sqrt{X}$ , we can ask the same question via either of:

- i). If  $X$  is chosen from the unit interval, with uniform distribution, then, conditional on the chosen number being either  $1/4$  or  $3/4$ , what is the probability that it is  $1/4$ ?
- ii). If  $Y$  is chosen from the unit interval, with a probability distribution given by density function (3), then, conditional on the chosen number being either  $1/2$  or  $\sqrt{3}/2$ , what is the probability that it is  $1/2$ ?

We can escape incoherence by requiring that, if a number is chosen according to *any* probability distribution on the unit interval that assigns probability zero to all singleton sets, then, for any finite subset of the unit interval, the probability conditional on the number being in that subset is the same for every member of the set. And, if we are to have equiprobability when the distribution is uniform, this is the *only* way to escape incoherence, since, for any random variable  $X$ , there will always be *some* function of  $X$  that is uniformly distributed.<sup>5</sup> But this convention may seem odd to some: consider

---

<sup>5</sup>To see this: let  $X$  be any random variable, with distribution  $\mu_X$ , and take  $g(X)$  to be the function of  $X$  given by

$$g(X) = \int_{-\infty}^X d\mu_X(x). \quad (5)$$

Then  $g$  has range in  $[0, 1]$ , and, for any  $a \in [0, 1]$ ,

$$P(g(X) \leq a) = a. \quad (6)$$

That is,  $g(X)$  is uniformly distributed on  $[0, 1]$ .

a density function that is very sharply peaked around  $1/2$ . On the convention under consideration, conditional on the supposition that the chosen number is either  $1/2$  or  $9/10$  (which could be as many standard deviations away from the peak as we like),  $1/2$  and  $9/10$  are equally probable.

These considerations will, I hope, lead some readers who initially regarded question 1 as having an obvious answer to conclude: things aren't as straightforward as they seemed.

### 3.2 The Sphere

Consider, again, Examples 2. A point is chosen, with uniform probability, on the surface of a sphere, and we are asked to reflect on the questions: a) What is the probability that the chosen point is in the Western Hemisphere, given that that it lies on the equator? b) Conditional on the chosen point lying on the great circle containing the Greenwich meridian, what is the probability that it lies closer to the equator than to a pole?

For question 2(a), the seemingly obvious answer is  $1/2$ . For 2(b), the obvious answer might seem to be  $1/2$ , again, as half of the length of any meridian consists of points that are closer to the equator than to a pole.

But consider this: it is *not* true that  $1/2$  of the earth's surface is closer to the equator than it is to a pole; more of it is closer to the equator. The probability that a point chosen with uniform probability is closer to the equator than to a pole is  $1/\sqrt{2} \approx 0.707$ . Since every point lies on some meridian, we might want to say that the probability, conditional on our point lying on the Greenwich (or any other) meridian, of being closer to the equator than to a pole, is  $1/\sqrt{2}$ .

Any reader who is wondering whether the *correct* answer to 2(b) is  $1/2$  or  $1/\sqrt{2}$  or some other number is reminded: the setup of the problem does not determine *any* answer. The answer of  $1/2$  seems to rely on some intuition that the conditional probabilities should share relevant symmetries with the unconditional distribution. An intuition is a dangerous thing; we would do well to replace the intuition with an explicit requirement regarding symmetries.

In Appendix 1 we define probability spaces and conditional probability spaces. Given a probability space, we define the associated *standard conditional probability space* as the conditional probability space that takes, as its domain of conditioning, all and only the propositions with nonzero probability. If a probability space is invariant under a transformation  $T$ , then *ipso facto* so is the standard conditional probability space. We may want to use symmetry considerations to extend the standard conditional probability space to one that includes conditionalization on null propositions. As a first pass, we might be tempted to require that our conditional probability space be invariant under all transformations—that is, one-one mappings that take measurable sets to measurable sets—that leave the unconditional probability space that we started with invariant. This, we might speculate, is the requirement needed to underwrite the

“obvious” answer to question 2(a). A moment’s reflection, however, reveals that this is unreasonably strong. Let  $C$  be any probability-zero subset of the sphere, and, for *any* one-one mapping  $T_C$  of  $C$  onto itself, consider a transformation of the sphere that consists of performing  $T_C$  on  $C$  and doing nothing elsewhere. Since  $P(C) = 0$ , this transformation does not change the unconditional probability of any set, and so our unconditional probability space is invariant under this transformation. Thus, to require invariance under arbitrary transformations that leave the unconditional probability space invariant entails that probabilities conditional on a null set  $C$  be invariant under arbitrary permutations of  $C$ , a requirement that is satisfiable when  $C$  is a finite set but not otherwise.<sup>6</sup>

The set  $S$  of events might have additional structure that we can require our transformations to preserve. In the sphere case, the elementary events are choices of points on a sphere, and these points have distances between them. We can restrict our attention to transformations of our probability space that preserve these distances. These are just the rigid rotations of the sphere. Requiring invariance under all rigid rotations entails that the conditional probability function, conditional on the chosen point lying on a circle, be invariant under the subgroup of rotations that leave the circle invariant. This is uniquely satisfied by a uniform distribution on the circle.

If the intuition that the obvious answer to the sphere questions 2(a) and 2(b) is  $1/2$  rests on an implicit assumption that probabilities, conditional on the point lying on a circle, should be invariant under rotations that leave the circle invariant, then, rather than leave this implicit, we should place it as an explicit condition on our conditional probability space. Can we do this? If we’re not too demanding about the extent of the set  $\mathcal{B}$  on which we conditionalize, then it is easy to show that we can. This is done in Appendix 2, where we construct a conditional probability space that includes conditionalization on all circles and subsets of circles of nonzero length, and is invariant under rigid rotations of the sphere.

We might want more than this in our domain of conditionalization. Can our conditional probability space be extended in such a way that it includes conditionalization on *all* measurable subsets of the sphere, and preserves symmetry under rotations?

If we demand countable additivity, then the answer is easy: no, we can’t. Given a coordinatization of the sphere by latitude and longitude, consider  $E_Q$ , the set of points on the equator whose longitudes are rational numbers. This set is invariant under rational rotations of the sphere about its axis. Invariance under such rotations requires that the probability, conditional on  $E_Q$ , ascribed to any interval of the equator be proportional to the length of the interval, and this in turn requires the probability assigned to single points on the equator be zero. But  $P_{E_Q}(E_Q)$  must be equal to one, and so the conditional probability function  $P_{E_Q}$  cannot be countably additive.

Similar considerations apply, of course, to Example 3. Our unconditional probabil-

---

<sup>6</sup>If  $C$  is a finite set, we can have a probability function that always assigns equal probabilities to sets of equal cardinality. This is not possible if  $C$  is infinite. In the infinite case, there must be measurable sets  $A$ ,  $B$ , of equal cardinality, with  $P(A) \neq P(B)$ . We can then choose some mapping that takes  $A$  to  $B$ .

ity function is invariant under translations of the unit interval (modulo 1). The set of rationals in the unit interval is invariant under the subgroup consisting of translations through a rational distance. Imposing translation symmetry on probabilities conditional on the number chosen being rational gives the expected answer: conditional on the number being rational, the probability that it lies in any interval is equal to the length of that interval. But this comes at the cost of violating countable additivity. If we conditionalize on the rationals we are faced with a choice between a symmetry condition that may be desired, and preserving countable additivity. This is something that we do not have to face when conditioning on sets of nonzero probability; if  $P$  is countably additive, and  $P(C) > 0$ , then  $P_C$  is also countably additive.

Suppose we're willing to give up countable additivity. Is there a conditional probability space that permits conditionalization on arbitrary measurable subsets of the sphere, and is invariant under rotations? Since this will include conditionalization on measure-zero subsets of  $S$  that are neither invariant under rotations nor contained in nontrivial subsets that are invariant under rotations, it is likely that, if such conditional probability spaces do exist, rotational symmetry will not suffice for uniqueness. We should expect that, if there are any, there are many such spaces, and that it would not be a trivial task to specify one. It is, as far as I know, an open question whether such conditional probability spaces exist. Philosophers who write as if one can blithely assume that such conditional probability spaces exist are kindly requested to show that they do, and, if there is more than one, to specify which one they have in mind.

### 3.3 The Eternal Coin

In the case of the sphere, things worked out (reasonably) well. We were able to identify a natural group of symmetries, and imposition of these symmetries entailed one of the 'obvious' answers to our questions. In other cases, we will not be so lucky. Symmetries that we may wish to impose can come into conflict.

An interesting example of this is provided by Cian Dorr (2010), in the set-up that he calls "The Eternal Coin." The Eternal Coin is a fair coin that is flipped every day, throughout an infinite past, and will continue to be flipped every day into an infinite future. In the absence of any other information about the coin, we are invited to consider credences in propositions such as

*H*: The Coin lands Heads today.

*P*: The Coin landed Heads on every day in the past.

*F*: The Coin will land Heads on every day in the future.

All credences—including those conditional on propositions with probability zero—will be taken to be predicated on the setup being as we have described it.<sup>7</sup>

---

<sup>7</sup>This is necessary because, if one has nonzero credence that the coin is not fair, or that the tosses are not independent, then conditionalization on either *F* or *P* will send credence that the setup is as described



We construct a probability space as follows. Our set  $\Upsilon$  of elementary events is the set of bi-infinite sequences of Heads and Tails. To form a  $\sigma$ -algebra  $\mathcal{F}$  of measurable sets, we proceed as follows. For any finite set of integers  $K$ , and any  $u \in \Upsilon$ , we form a *cylinder set*  $C_K(u)$  consisting of all elements of  $\Upsilon$  that agree with  $u$  on the set  $K$ . That is, a cylinder set is the set of all events that agree on some finite subset of integers. We take  $\mathcal{C}$  to be the smallest  $\sigma$ -algebra containing all cylinder sets.<sup>8</sup>

To define a probability measure  $Pr$  on  $\langle \Upsilon, \mathcal{C} \rangle$ , it suffices to specify the probabilities of cylinder sets.<sup>9</sup> To do this, we assign, for any  $k$ -element set  $K$ , the probability  $2^{-k}$  to each cylinder set  $C_K(u)$ . This function has a unique countably additive extension to  $\mathcal{C}$ , which we will take to be our probability measure  $Pr$ . This gives us a probability space  $\langle \Upsilon, \mathcal{C}, Pr \rangle$ .

This probability measure has, as expected, the following features:

- i). Each individual flip has equal probability  $1/2$  for  $H$  and  $T$ .
- ii). Outcomes of distinct flips are independent: if  $K, L$  are disjoint sets, then, for all  $u, v \in \Upsilon$ ,

$$Pr(C_K(u) \cap C_L(v)) = Pr(C_K(u)) \cdot Pr(C_L(v)).$$

For any set of integers  $L$ , let  $F_L : \Upsilon \rightarrow \Upsilon$  be the ‘bit flip’ transformation on  $L$ , that is, the transformation that consists of exchanging  $H$  and  $T$  at each place in  $L$ . Our probability space is invariant under all such transformations.

Our probability space is also invariant under permutations of the integers. For any bijection  $\pi : \mathbb{Z} \rightarrow \mathbb{Z}$ , let  $T_\pi : \Upsilon \rightarrow \Upsilon$  be the operation whose action on a bi-sequence  $u$  permutes the values of  $u$ ,

$$(T_\pi u)_k = u_{\pi(k)}. \tag{7}$$

Permutations that will be of particular interest are the shift operations. For any integer  $n$ , let  $S_n : \Upsilon \rightarrow \Upsilon$  be the operation of shifting everything  $n$  places:

$$(S_n u)_k = u_{k-n}. \tag{8}$$

Invariance under shift operations means that, although our coordinatization has a distinguished origin (the day 0, which we are calling “today”), our probability space is invariant under shift of this origin.

If  $P_n$  is the proposition that the coin landed Heads on the past  $n$  days, then  $Pr(P_n) = 2^{-n}$ . Since  $P$  entails  $P_n$  for each  $n$ , it follows that  $Pr(P) = 0$ . Similarly,  $Pr(F) = 0$ .

The function  $Pr$ , of course, uniquely determines probabilities conditional on propositions with non-zero probability. Dorr invites us to consider probabilities conditional on some zero-probability propositions, such as  $P$ ,  $F$ , and  $P \vee F$ . It is, of course,

---

to zero.

<sup>8</sup>See Appendix 1 for definitions of any terms that might be unfamiliar.

<sup>9</sup>In this section, we use  $Pr$  for our probability function to avoid confusion with the proposition  $P$ .

possible to extend our probability assignments to include probabilities conditional on propositions such as these, and this can be done in a variety of ways.<sup>10</sup>

Here's one way to do it. For any  $n$ , let  $K_n = [-n, n]$ , and, for any  $A \in \mathcal{C}$ , let  $A_n$  be the proposition that commits only to what  $A$  says about coin flips in  $K_n$ , and says nothing about what happens outside this interval.<sup>11</sup>

Our set  $\mathcal{B}$  of conditions will consist of all nonempty  $B \in \mathcal{C}$ . For  $B \in \mathcal{B}$ , let  $\mathcal{A}_B$  be the set of  $A \in \mathcal{C}$  such that the sequence  $P(A_n|B_n)$  converges to a limit as  $n \rightarrow \infty$ , and, for  $A \in \mathcal{A}_B$ , take

$$Pr(A|B) = \lim_{n \rightarrow \infty} Pr(A_n|B_n). \quad (10)$$

A few of the conditional probabilities that we thereby obtain are,

$$\begin{aligned} Pr(H|F) &= Pr(T|F) = Pr(H|P) = Pr(T|P) = 1/2; \\ Pr(P|P \vee F) &= Pr(F|P \vee F) = 1/2; \\ Pr(P|P \vee HF) &= Pr(F|HP \vee F) = 2/3; \\ Pr(HF|P \vee HF) &= Pr(HP|HP \vee F) = 1/3; \\ Pr(P \vee HF|P \vee F) &= Pr(HP \vee F|P \vee F) = 3/4. \end{aligned} \quad (11)$$

The limiting procedure we have sketched is, of course, only one possible limiting procedure, and no claim is made for priority of this over other procedures. We have made a frankly arbitrary choice, and have obtained the above conditional probabilities; other choices will yield other values.

The conditional probabilities we have obtained preserve independence and bit-flip symmetry. The limiting procedure we have chosen manifestly breaks shift symmetry. Unsurprisingly, the conditional probabilities we obtain from it also violate shift symmetry. To see this, consider the one-day shift  $S_1$ . We have,

$$S_1(P) = HP \quad S_1(HF) = F \quad (12)$$

However,

$$Pr(P|P \vee HF) \neq Pr(HP|HP \vee F). \quad (13)$$

We therefore have extended our probability function in a way that respects independence of distinct flips, and also bit-flip symmetry, but violates shift symmetry. We

---

<sup>10</sup>In this section, we will find ourselves conditionalizing on some fairly complex propositions, and so it will be convenient to switch from the subscript notation for conditional probabilities used in the rest of the paper to the slash notation.

<sup>11</sup>That is, take

$$A_n = \bigcup_{u \in A} C_{K_n}(u). \quad (9)$$

should ask whether we can do better, and extend our probability function in such a way that all of the above conditional probabilities are defined so as to respect all of these symmetries.

Dorr shows that, counterintuitively,<sup>12</sup> the answer is no. Provided that  $P(P|P \vee F)$  and  $P(F|P \vee F)$  are defined and are both positive, shift invariance entails that

$$Pr(H|F) = Pr(H|P) = 1. \tag{14}$$

Proof is given in Appendix 3.

A similar argument yields a violation of countable additivity. Let  $P^+$  be the proposition that the coin has landed Heads every day in the past but will land Tails sometime, either today or in the future, and let  $F^+$  be the proposition that the coin will land Heads every day in the future, but landed Tails today or sometime in the past. Shift invariance, together with the conditions that  $Pr(P^+|P^+ \vee F^+)$  and  $Pr(F^+|P^+ \vee F^+)$  are defined and are both nonzero, entails that, for every  $n$ , the probability conditional on  $P^+$  that the coin lands Heads today and every day for  $n$  days into the future is one. This in turn entails (letting  $H_n$  be the proposition that the coin will land Heads  $n$  days from now and  $T_n$ , the proposition that it will land Tails), that, for each  $n$ .

$$\begin{aligned} Pr(H_n|P^+) &= 1; \\ Pr(T_n|P^+) &= 0; \end{aligned} \tag{15}$$

even though the probability, conditional on  $P^+$ , that, for some  $n$ ,  $T_n$  is true, is unity.

Including the propositions  $P$ ,  $F$ , and  $P \vee F$  in the set of propositions on which we can conditionalize, and imposing shift symmetry, is possible, but it comes at a high cost: we lose independence; it is no longer true that conditionalization on a proposition that specifies outcomes on a set of days not including today leaves the probability of the coin landing Heads today unchanged. Symmetry conditions that we would like our conditional probability space to respect clash; we can't get all that we want.

Depending on our purpose, we might prefer to preserve one or the other of the symmetries. If the Eternal Coin is being considered as an idealization of a situation in which a coin is tossed a large but finite number of times, then shifts will not be symmetries of the finite system, which is our real object of interest, and so it will not be important for our purposes to demand shift invariance of the conditional probability space. There might be other purposes for which shift invariance is of such paramount importance that it would be worth abandoning independence (though it is hard to see why it would not be preferable to simply leave those conditional probabilities undefined).

If we think of the setup as involving an actual bi-infinite sequence of coin tosses, not an idealization of a finite set-up, then, as Dorr convincingly argues, violation of shift invariance is bizarre. Dorr invites us to imagine ourselves causally isolated from

---

<sup>12</sup>Perhaps. The more one thinks about what is required to give values to these conditional probabilities, the less clear it becomes that we have intuitions about them at all.

the Eternal Coin. I learn nothing about the outcomes of its flips as the days pass. Now, consider the following:  $HP$ , the proposition that the coin lands Heads today and landed Heads every day in the past, is the proposition that, tomorrow, I will express by the words, “The coin landed Heads every day in the past,” the same sentence that I use today to express the proposition  $P$ . Similarly,  $HP \vee F$  is the proposition that I will express tomorrow using the same words I use today to express  $P \vee HF$ . Today, when I say “My credence that the coin landed every day in the past, conditional on the supposition that it either landed Heads every day in the past or will land Heads today and every day in the future,” I denote  $Pr(P|P \vee HF)$ ; tomorrow, the same phrase denotes  $Pr(HP|HP \vee F)$ . Does it make sense for these to have different values? To do so involves distinguishing between today and tomorrow in a way that seems unwarranted by the setup of the problem. Shift invariance, it seems, is a requirement of rationality.

On the other hand, it is stipulated in the setup that coin tosses on distinct days are independent of each other.  $Pr(H|P_n)$  is equal to  $1/2$ , for every  $n$ , no matter how large. The toss today is independent of every past toss; should it not also be independent of *all* the past tosses? Recall that all of these probabilities are meant to be predicated on the supposition that the setup is as described, which includes stipulation of independent tosses. For our credences, conditional on this setup, to violate independence, setting  $Pr(H|P)$  equal to 1, seems no less irrational than violation of shift invariance.

Violation of either symmetry, shift invariance or independence, is a high price to pay for probabilities conditional on null events. Dorr bites the bullet and preserves shift invariance at the price of independence, but it is not clear that this is preferred over the alternative. Better still, it would seem, would to be preserve both symmetries, which we can do, of course, by restricting the domain of conditionalization to propositions with positive probability.

## 4 Probabilities conditional on a $\sigma$ -algebra

Consider, once again, Example 2. As noted, an “obvious” answer to the question 2(a) of the probability that a point chosen with uniform probability on the sphere lies in the Western hemisphere, conditional on the supposition that it lies on the equator, is  $1/2$ . For the question 2(b) of the probability that the point lies closer to an equator than a pole, conditional on the supposition that it lies on the Greenwich meridian, both  $1/2$  and  $1/\sqrt{2}$  seem to have merit.

One way to think about question 2(a) is to imagine that, first, a circle of latitude is chosen, and then a point is chosen on that circle according a probability distribution conditional on the point lying on the circle. Taking the total area of the surface of the sphere to be 1, the area between two circles of latitude, at angles  $a$ ,  $b$ , measured from the equator, is equal to

$$\frac{1}{2} \int_a^b \cos \phi \, d\phi.$$

This means that the latitude  $\Phi$  must be distributed according to

$$P(\Phi \in A) = \frac{1}{2} \int_A \cos \phi \, d\phi. \quad (16)$$

That is,  $\Phi$  has density function

$$f_\Phi(\phi) = \frac{1}{2} \cos \phi. \quad (17)$$

The longitude  $\Theta$  is distributed with uniform probability on  $[-\pi, \pi]$ , and so has density function

$$f_\Theta(\theta) = \frac{1}{2\pi}. \quad (18)$$

Latitude and longitude are independent random variables. That is,

$$P(\Phi \in A \ \& \ \Theta \in B) = P(\Phi \in A) P(\Theta \in B) = \int_A f_\Phi(\phi) \, d\phi \int_B f_\Theta(\theta) \, d\theta. \quad (19)$$

for all measurable  $A \subseteq [-\pi/2, \pi/2]$  and  $B \subseteq [-\pi, \pi]$ .

What should the conditional distribution of the longitude  $\Theta$  be taken to be, conditional on a given circle of latitude? We may want the conditional probabilities to mesh with the unconditional probabilities in a nice way, and demand

$$P(\Phi \in A \ \& \ \Theta \in B) = \int_A P(\Theta \in B | \Phi = \phi) f_\Phi(\phi) \, d\phi \quad (20)$$

for all measurable  $A, B$ . The simplest way to do this, which is also the way that is naturally suggested by the independence of  $\Theta$  and  $\Phi$ , is to take  $P(\Theta \in B | \Phi = \phi)$ , for each  $B$ , to have the constant value  $P(\Theta \in B)$ , independent of  $\phi$ . But it's not the only way. We can take any set of latitudes of measure zero, and choose distributions for  $\Theta$ , conditional on  $\Phi = \phi$  in that set, any way we want, and still satisfy the meshing condition (20). That means that (20) is compatible with *any* answer to question 2(a).

It is natural, however, to take  $P(\Theta \in B | \Phi = \phi)$  to be, for each  $B$ , a continuous function of  $\phi$ . This condition, together with the meshing condition (20), uniquely fixes

$$P(\Theta \in B | \Phi = \phi) = P(\Theta \in B). \quad (21)$$

Similarly, we can define conditional distributions of latitude, conditional on meridian lines (lines of constant longitude), and demand that these also mesh with the unconditional probabilities:

$$P(\Phi \in A \ \& \ \Theta \in B) = \int_B P(\Phi \in A | \Theta = \theta) f_\Theta(\theta) \, d\theta \quad (22)$$

This, together with the requirement that for each  $A$ ,  $P(\Phi \in A | \Theta = \theta)$  be a continuous function of  $\theta$ , uniquely fixes

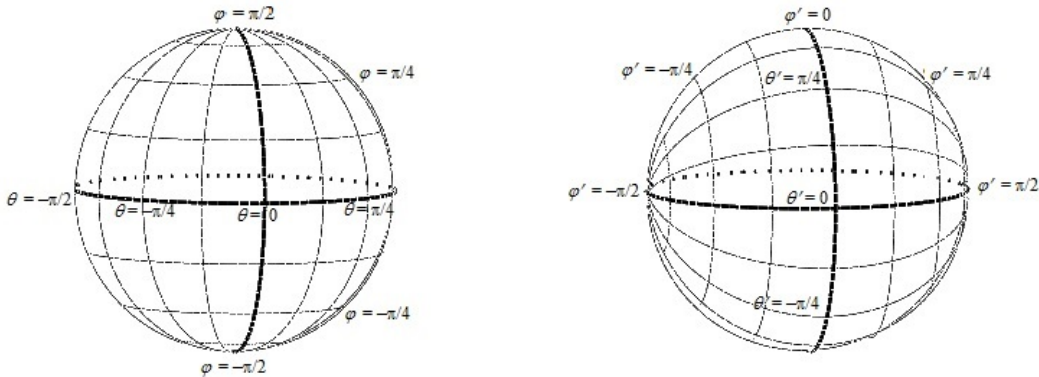
$$P(\Phi \in A | \Theta = \theta) = P(\Phi \in A), \quad (23)$$

corresponding to conditional density functions

$$f_{\Phi}(\phi | \Theta = \theta) = \frac{1}{2} \cos \phi. \quad (24)$$

Consider, now, question 2(b). What is the probability that the point lies closer to an equator than a pole, conditional on the supposition that it lies on the Greenwich meridian? We can imagine that a meridian is first chosen, and then a point chosen on that meridian. Using (23) yields the result that, conditional on any meridian, the probability is  $1/\sqrt{2}$  that the chosen point is closer to the equator than to a pole.

On the other hand, since we are only imagining these things, we can also imagine the sphere partitioned by circles parallel to the circle containing the Greenwich meridian (see Figure 1) and imagine that first one of these circles is chosen, and then a point chosen on that circle according to a probability distribution conditional on the circle. If this procedure is to yield uniform probabilities on the sphere, we must have the distributions on almost all of these circles be uniform, and this plus continuity militates a uniform distribution on all of them. This yields the answer  $1/2$  to question 2(b).



**Figure 1: Two coordinatizations of the sphere.**

Which answer is correct? If the point on the sphere is, in fact, chosen according to one of the two-step procedures we have imagined, then such a set-up privileges one of the answers. But if the point is simply chosen, with uniform probability, on the sphere, then the set-up privileges neither answer, and, if one or the other has greater intuitive appeal, this may be because one is implicitly assuming one or the other scenario.

A circle is just a circle,<sup>13</sup> and the great circle containing the Greenwich meridian, *qua* circle on the sphere, is an element of many different partitions of the sphere. If we really think that  $f_{\Phi}(\phi | \Theta = 0)$ , as given by (24), is a conditional density function yielding the distribution of the random variable  $\Phi$  conditional on the supposition that  $\Theta = 0$ , then it shouldn't matter how this supposition is described. The supposition can equally well be described using coordinates that take circles parallel to the great circle containing the Greenwich meridian as lines of latitude  $\phi'$ . Then the great circle containing our original Greenwich meridian is the set of points for which  $\phi' = 0$ . On this circle the new longitude  $\theta'$  differs from the old latitude  $\phi$  by a constant, and, if we choose the zero-point of our new longitude as our old equator, we will have  $\theta'$  equal to  $\phi$  on the circle. But a uniform distribution of the new longitude on circles of constant  $\phi'$  requires a conditional density function

$$f_{\Theta'}(\theta' | \Phi' = \phi') = \frac{1}{2\pi} \quad (25)$$

It can't be the case that, conditional on the chosen point lying on the circle that is the great circle containing the Greenwich meridian of our first coordinatization and is the equator of our second, we have different conditional distributions depending on how we describe the circle. Taking (24) to yield the conditional distribution of  $\Phi$  on this circle is incompatible with a uniform distribution of  $\Theta'$ , as given by (25).

All of this suggests that perhaps we are better off leaving probabilities conditional on null propositions undefined. Nonetheless, we can still write the probability function on the sphere in the form (20) or (22), and this can be a useful thing to do, whether or not we regard the quantities  $P(\Phi \in A | \Theta = \theta)$  and  $P(\Theta' \in B | \Phi' = \phi')$  as conditional probabilities.

Whether or not we think of them as genuine conditional probabilities, quantities such as  $P(\Phi \in A | \Theta = \theta)$  are always there if we want them—that is, they can be shown to always exist. Given a probability space  $\langle \Omega, \mathcal{A}, P \rangle$ , and a random variable  $X$ , with distribution  $\mu$ , we can always find a function  $f_A : \mathbb{R} \rightarrow \mathbb{R}$  such that, for all Borel sets  $\Delta$ ,

$$P(A \ \& \ X \in \Delta) = \int_{\Delta} f_A(x) \, d\mu(x). \quad (26)$$

The existence of such functions is guaranteed by the Radon-Nikodym theorem (see, *e.g.*, Billingsley (2012, §32–33).) The condition (26) defines the function  $f_A$  only up to

---

<sup>13</sup>Oddly enough, this has been disputed. In connection with this example, E.T. Jaynes (2003, p. 470) writes,

Nearly everybody feels that he knows perfectly well what a great circle is; so it is difficult to get people to see that the term 'great circle' is ambiguous until we specify what limiting operation is to produce it.

This strikes me as confused. One and the same great circle can be the limit of many different decreasing sequences of subsets of the sphere, but the circle is not itself *produced* by the limiting operation. Not so with probabilities conditional on a great circle, which, unless stipulated as primitive, are obtained via some limiting operation.

sets of probability zero. If  $f_A$  is any function satisfying (26), then any function that differs from  $f_A$  only on a set of probability zero will also satisfy it, and any functions that satisfy the condition will differ at most on a set of probability zero.

Such functions always exist, and can be useful calculational tools. But, as we have seen, it may be problematic to regard  $f_A(x)$  as yielding a conditional probability, namely, the probability of  $A$ , conditional on the supposition that  $X = x$ . The problem, illustrated by the example of the sphere, is the existence of some other random variable,  $X'$ , such that one and the same set of events can be equivalently picked out by two conditions  $X = x$  and  $X' = x'$ . Continuity or other considerations might lead to functions  $f_A$  and  $f'_A$  that differ on the set  $C$  picked out by either the condition  $X = x$  or  $X' = x'$ . Yet the probability of  $A$ , conditional on an event  $C$ , should depend only on  $A$  and  $C$ , and not how we happen to describe the event  $C$ . The standard view, which goes back to Kolmogorov (1950, p. 51), is that, useful as such functions are, we ought not to regard them as yielding probabilities conditional on an event of probability zero.<sup>14</sup>

Generalizing: for any probability space  $\langle \Omega, \mathcal{A}, P \rangle$ , and any  $\sigma$ -algebra  $\mathcal{G} \subseteq \mathcal{A}$ , and any  $A \in \mathcal{A}$ , a function  $g_A : \Omega \rightarrow \mathbb{R}$  is said to be a *conditional probability of  $A$  with respect to  $\mathcal{G}$*  iff it is a  $\mathcal{G}$ -measurable function such that

$$P(AG) = \int_G g_A dP \quad (27)$$

for all  $G \in \mathcal{G}$ . We will write  $g_A = P(A|\mathcal{G})$ .<sup>15</sup> Condition (27) then becomes

$$P(AG) = \int_G P(A|\mathcal{G}) dP. \quad (28)$$

Hájek (2003, 291) calls this “Kolmogorov’s elaboration of the ratio formula.”

Functions of the sort appearing in (26) are yielded as special cases. For any random variable  $X$ , let  $\sigma(X)$  be the  $\sigma$ -algebra consisting of the sets  $X^{-1}(B)$ , where  $B$  ranges over Borel subsets of the real line. Let  $P(A|\sigma(X))$  be a conditional probability of  $A$  with respect to  $\sigma(X)$ . Recall, this is a function from our sample space  $\Omega$  to the reals. The condition that it be a  $\sigma(X)$ -measurable function entails that it is constant on sets of constant  $X$ . That is, if, for two points  $\omega, \omega' \in \Omega$ , we have  $X(\omega) = X(\omega')$ , then the function  $P(A|\sigma(X))$  must take on the same value at these points:  $P(A|\sigma(X))(\omega) = P(A|\sigma(X))(\omega')$ . We can define a function  $f_A : \mathbb{R} \rightarrow \mathbb{R}$  via

$$f_A(X(\omega)) = P(A|\sigma(X))(\omega). \quad (29)$$

If  $P(A|\sigma(X))$  satisfies (28), then  $f_A$  will satisfy (26).

---

<sup>14</sup>In his discussion of the Borel paradox, Kolmogorov writes, “This shows that the concept of a probability conditional on an isolated given hypothesis whose probability equals 0 is inadmissible” (Kolmogorov, 1950, p. 51).

<sup>15</sup>The notation is intended to be both reminiscent of, and distinct from, the notation used for conditional probabilities.



Conditional probabilities, as usually conceived, that is, as defined by ratios of unconditional probabilities, are a special case of these conditional probabilities with respect to a  $\sigma$ -algebra. Let  $\{G_i\}$  be a countable partition, and let  $\mathcal{G}$  be the  $\sigma$ -algebra generated by this partition. Since the elements of the partition  $\{G_i\}$  are atoms of this  $\sigma$ -algebra, and  $P(A|\mathcal{G})$  is required to be a  $\mathcal{G}$ -measurable function, it must be a constant function on each  $G_i$ . Let  $P(A|G_i)$  be the value that  $P(A|\mathcal{G})(\omega)$  takes on for  $\omega \in G_i$ . Then the condition that (28) hold for all  $G \in \mathcal{G}$  is equivalent to the condition that

$$P(AG_i) = P(A|G_i)P(G_i), \tag{30}$$

which, of course, yields the familiar ratio formula for  $P(A|G_i)$  whenever  $P(G_i) > 0$ . In this sense, we have a generalization of conditional probabilities.

When an agent learns which element of  $\{G_i\}$  is true, she at the same time learns the truth value of each proposition in the  $\sigma$ -algebra  $\mathcal{G}$ . The heuristic idea behind the introduction of probabilities conditional on more general  $\sigma$ -algebras is to mimic this. A random variable  $X$  partitions the space  $\Omega$  of events into sets of constant  $X$ . If a point  $\omega$  is chosen from  $\Omega$ , learning the valuing of  $X(\omega)$  provides partial information about  $\omega$ ; it tells us for every set  $\Delta \in \sigma(X)$ , whether or not  $\omega \in \Delta$ . Similarly for other  $\sigma$ -algebras, whether or not generated by a random variable; a specification, for each set in a  $\sigma$ -algebra  $\mathcal{G}$ , whether or not  $\omega$  is in that set, provides some information about  $\omega$ , information that is partial unless singleton sets are among the members of  $\mathcal{G}$ .

Let  $\mathcal{G} \subseteq \mathcal{A}$  be a  $\sigma$ -algebra that contains atoms—that is, elements of  $\mathcal{G}$  with no non-empty proper subsets in  $\mathcal{G}$ —that cover  $\Omega$ . For  $A \in \mathcal{A}$ , let  $P(A|\mathcal{G})$  be a conditional probability of  $A$  with respect to  $\mathcal{G}$ . If  $G$  is an atom of  $\mathcal{G}$ , then  $P(A|\mathcal{G})$  must take on a constant value on  $G$ . Should we regard this value, the value of  $P(A|\mathcal{G})$  for  $\omega \in G$ , as the probability of  $A$  conditional on the proposition  $G$ ?

There are two sorts of problems with this. The first is technical and local, in that it applies only to certain  $\sigma$ -algebras that we might dismiss as pathological. Nonetheless, it should give us pause, as it shows that the heuristic motivation of the characterization of  $P(A|\mathcal{G})$ , namely, as conditional probabilities resulting from information specifying, for each element  $G$  of a  $\sigma$ -algebra  $\mathcal{G}$ , whether or not  $\omega \in G$ , can break down. The second sort of problem is conceptual and global, and poses a serious objection to taking the value of  $P(A|\mathcal{G})$  for  $\omega \in G$ , as the probability of  $A$  conditional on the proposition  $G$  (except in special circumstances, to be discussed in the next section).

The first problem is this. On the heuristic view that  $P(A|\mathcal{G})$ , evaluated on some atom  $G$  of  $\mathcal{G}$ , yields the probability of  $A$  appropriate to learning that  $\omega \in G$ , we would expect that, if  $\mathcal{G}$  is a  $\sigma$ -algebra whose atoms are all the singleton sets, then  $P(A|\mathcal{G})$  would be equal to 1 if  $\omega$  is in  $A$  and 0 if not, since learning which atom of  $\mathcal{G}$  obtains is complete information about  $\omega$ . But this won't always be the case. Let our probability space be the unit interval with Lebesgue measure. Let  $\mathcal{G}$  consist of the smallest  $\sigma$ -algebra containing all of the singleton sets; this consists of the countable sets and their complements. Now let  $A$  be any set with  $P(A) \in (0, 1)$ . It is easy to see that  $P(A|\mathcal{G})(\omega)$  must be equal to  $P(A)$  for almost all  $\omega$ , violating our expectation

that it will everywhere be equal to 0 or 1.<sup>16</sup>

This problem can be thought of one of being excessively permissive about the sub- $\sigma$ -algebras on which we may conditionalize. Easwaran (2008, §8.1) makes a well-motivated proposal on which this problem does not arise. Instead of conditionalizing on arbitrary sub- $\sigma$ -algebras, we consider only those that consist of all the measurable sets that are unions of elements of some partition  $\mathcal{E}$ . Call the  $\sigma$ -algebra consisting of the elements of  $\mathcal{A}$  that are unions of elements of a partition  $\mathcal{E}$ ,  $\mathcal{A}_{\mathcal{E}}$ . If the information on which we are to update consists of a specification of which element of  $\mathcal{E}$  obtains, then  $\mathcal{A}_{\mathcal{E}}$  is the relevant sub- $\sigma$ -algebra, since specifying which element of the partition  $\mathcal{E}$  contains  $\omega$  is equivalent to specifying, for every  $F \in \mathcal{A}_{\mathcal{E}}$ , whether or not  $\omega \in F$ . If the partition  $\mathcal{E}$  contains all singleton sets, then  $\mathcal{A}_{\mathcal{E}}$  is just  $\mathcal{A}$ , and, for any  $A \in \mathcal{A}$ ,  $P(A|\mathcal{A}_{\mathcal{E}})$  must be equal, on all but a set of measure zero, to the characteristic function of  $A$ .

The second problem is the one we have already been discussing, and it is more serious. Let  $G$  be an atom of a  $\sigma$ -algebra  $\mathcal{G}$ . Though, for any  $G$  with  $P(G) = 0$ , the condition (28) leaves the value of  $P(A|\mathcal{G})$  on  $G$  undetermined, the condition together with other natural constraints, such as requiring  $P(A|\mathcal{G})$  to be a continuous function, can, as we have seen, determine the value of  $P(A|\mathcal{G})$  on  $G$ . But this is not enough to warrant taking this value as the probability of  $A$ , conditional on  $G$ , since the same set  $G$  will be an atom of other  $\sigma$ -algebras, and the same considerations might dictate that, for some other  $\sigma$ -algebra  $\mathcal{G}'$  containing  $G$ , the value that  $P(A|\mathcal{G}')$  has on  $G$  be different from the value that  $P(A|\mathcal{G})$  on  $G$ . In cases, such as the sphere example, in which the set-up privileges neither  $\sigma$ -algebra, it would be a mistake to take either of these values (or any other) as *the* probability of  $A$  conditional on  $G$ .

This is, as mentioned, the standard view. Taking up this suggestion, Easwaran concludes,

this means we must view conditional probability as (in general) a three-place function, depending not only on  $A$  and  $G$ , but also the partition  $\mathbf{G}$  defining the set of “relevant alternatives” to  $G$ . In particular cases, this partition will be specified by the experiment an agent is considering  $G$  as an outcome to, or the set of alternative hypotheses under consideration, or some other contextual factor. Thus, we must think of conditional degree of belief as a function  $P(A|G, \mathbf{G})$  rather than just  $P(A|G)$  (Easwaran, 2011, pp. 143–44).

We should ask: under what conditions will there be a set of relevant alternatives that is uniquely picked out by the set-up?

It is frequently suggested, as in the quotation from Easwaran, that it is the experiment that yields the data that determines a relevant partition (see also the discussion in Rényi 2007a, §2.1). On this rationale, though, it is hard to see that we would ever need to go beyond a finite partition. Unless we are entertaining the fiction of agents with infinite powers of discrimination, there are only finitely many distinguishable al-

---

<sup>16</sup>This is example 33.11 of Billingsley (2012).

ternatives as to the outcome of any experiment.<sup>17</sup> Even if we do imagine agents with infinite powers of discrimination, the set of alternatives they could record, using a finite alphabet, in a lab notebook of finite capacity, is a finite set.

Unproblematic null-condition conditional probabilities are not as commonplace as some of the literature might suggest. However, there are cases in which the set-up of a problem *does* permit one to speak unambiguously of the the probability of an event conditional on a null proposition. In those cases, null-condition conditional probabilities are unobjectionable, and they can be useful, though they are not indispensable. These are the subject of the next section.

## 5 Unproblematic Null-Condition Probabilities

### 5.1 Likelihood functions

It is common, in statistical practice, to regard outcomes of some experiment as being generated by an incompletely known probability distribution characteristic of the experimental set-up. Data gathered is used to gain information about that distribution. We commonly consider a family of candidate distributions; typically this family is indexed by some set of parameters. For instance, we might regard an experimentally measurable variable as being normally distributed with unknown mean  $\mu$  and unknown variance  $\sigma^2$ . A data-set is generated, and is used to gain information about the values of the parameters.

Let  $\Omega$  be the set of possible outcomes of an experiment, and let  $\mathcal{F}$  be the set of measurable subsets of  $\Omega$ . Suppose that the candidate probability distributions are characterized by specifying the values of  $n$  parameters  $(\theta_1, \theta_2, \dots, \theta_n)$ . We take our parameter space  $\Gamma$  to be the set of all such ordered  $n$ -tuples of parameters. For every  $n$ -tuple  $\theta = (\theta_1, \theta_2, \dots, \theta_n)$ , let  $P_\theta$  be the corresponding probability distribution on  $\langle \Omega, \mathcal{F} \rangle$ . For each  $E \in \mathcal{F}$ , let  $\mathcal{L}_E$  be the function from  $\Gamma$  to the real numbers defined by

$$\mathcal{L}_E(\theta) = P_\theta(E). \tag{31}$$

These functions are called *likelihood functions*. Given a proposition  $E$  about the experimental outcome (which might, for example, be a specification of the data gathered), we can, for example, ask which  $n$ -tuple  $\theta$  of parameters yields the highest value of  $\mathcal{L}_E(\theta)$ ; this will be the *maximum likelihood* estimate of the parameters.

In standard, frequentist statistics, the parameter space is not itself subject to probabilistic considerations; it is regarded as nonsensical to ascribe probabilities, prior or posterior, to propositions regarding values of the parameters. Hence,  $P_\theta$  is not regarded as a conditional probability distribution, conditional on a proposition of probability 0.

---

<sup>17</sup>This is even easier to see in these days in which laboratory equipment has digital readout than it was in the old days of pointers and dials!

On a Bayesian approach, on the other hand, one also ascribes probabilities to propositions regarding the values of the parameters, and the process of gaining information about the parameter values is modelled by conditionalization on the experimental result. Let  $\mathcal{G}$  be a  $\sigma$ -algebra of subsets of the parameter space  $\Gamma$ . Let  $\mathcal{H}$  be the smallest  $\sigma$ -algebra containing all sets of the form  $F \times G$ , with  $F \in \mathcal{F}$  and  $G \in \mathcal{G}$ . Suppose that, for each  $E \in \mathcal{F}$ , the likelihood function  $\mathcal{L}_E(\theta)$  is a  $\mathcal{G}$ -measurable function. Then, given a probability measure  $Q$  on  $\langle \Gamma, \mathcal{G} \rangle$ , we can form a new probability space whose event space is the Cartesian product  $\Omega \times \Gamma$  of the experimental outcome space and the parameter space, and whose measurable sets are the sets in  $\mathcal{H}$ : we define a probability measure  $P$  as the unique countably additive extension to  $\mathcal{H}$  of the function that, for  $F \in \mathcal{F}$  and  $G \in \mathcal{G}$ , is given by

$$P(F \times G) = \int_G \mathcal{L}_F(\theta) dQ. \quad (32)$$

We now have a probability space  $\langle \Omega \times \Gamma, \mathcal{H}, P \rangle$ . The experimental outcome  $X$ , and parameter value  $\Theta$ , are random variables on this probability space. The  $\sigma$ -algebra  $\sigma(X)$  that consists of propositions about the experimental outcome is  $\mathcal{F} \times \Gamma$ —that is, the set whose elements are sets of the form  $F$ , for all  $F \in \mathcal{F}$ —, and  $\sigma(\Theta)$ , the  $\sigma$ -algebra that consists of propositions about parameter values, is  $\Omega \times \mathcal{G}$ . One can readily verify that a version of conditional probability with respect to  $\sigma(\Theta)$  is obtained by setting

$$P(E||\Theta)(\omega) = \mathcal{L}_E(\Theta(\omega)). \quad (33)$$

Any version of conditional probability with respect to  $\sigma(\Theta)$  will have to agree with (33) with probability 1.

Is it permissible to regard the values of  $P(E||\Theta)$  on the atoms of  $\sigma(\Theta)$  as probabilities conditional on null propositions? There is a natural one-one correspondence between the atoms of  $\sigma(\Theta)$  and the points in the parameter space  $\Gamma$ . In this case, we have a  $\sigma$ -algebra that is picked out as special by the set-up of the problem; the random variable  $\Theta$  represents the parameters of the system about which we are trying to gain information, and the atoms of the  $\sigma$ -algebra  $\sigma(\Theta)$  correspond to maximal specification of these parameters. In this case, there seems no threat of ambiguity due to variant choices of  $\sigma$ -algebra to conditionalize on, and we can, relatively unproblematically, regard these values as null-condition probabilities.

We can form a conditional probability space by taking the set  $\mathcal{B}$  of permissible conditions to include, in addition to all propositions with positive probability, also the atoms of  $\sigma(\Theta)$ , corresponding to point values of our parameters. This still leaves us with a set of conditions that, though it goes beyond the standard set, is still fairly sparse compared to the full set  $\mathcal{H}$  of measurable sets.

It is not uncommon to deal with nested families of models, in which the parameter space of one model is a lower-dimensional subspace of the parameter space of another. This might come about, for example, by considering a model in which the value of some parameter is fixed, or two parameters are constrained to be equal. We will want to

retain the same likelihoods in the reduced model. We will also want a probability distribution over the reduced space. Here the issue illustrated by the Borel-Kolmogorov paradox resurfaces; the probability distribution on the higher-dimensional space does not determine a distribution on the lower-dimensional space, and defining one via a limiting procedure will lead to differing results, depending on the procedure chosen. It is a mistake to regard the lower-dimensional model as being obtained, in a straightforward way, from the higher-dimensional model via conditionalization.

Though, in this case, it is permissible to treat the likelihoods as probabilities conditional on point-values of the parameters, it is by no means necessary to do so. All of our standard statistical reasoning goes through if we restrict our domain of conditionalization to the traditional choice of sets with positive probability.

If we obtain evidence  $E$  about the experimental outcome, then we can update our credences about parameter values via conditionalization. For any measurable subset  $\Delta$  of parameter-space,

$$P(\Theta \in \Delta) \rightarrow P_E(\Theta \in \Delta) = \frac{P(E \& \Theta \in \Delta)}{P(E)}. \quad (34)$$

If the prior distribution of  $\Theta$  is given by a density function  $\mu$ , then the process of conditionalization yields a new density function  $\mu_E$ . In order for (34) to be satisfied,

we must have, for almost all  $\theta$ ,<sup>18</sup>

$$\mu_E(\theta) = \frac{\mathcal{L}_E(\theta) \mu(\theta)}{P(E)}. \quad (40)$$

This is often called the *continuous form of Bayes' theorem*. Thinking of (40) as a form of Bayes' theorem invites to think of  $\mathcal{L}_E(\theta)$  as the probability of  $E$  conditional on a point value of the parameter  $\theta$ . But the use of the new density  $\mu_E$  is to generate the new probability distribution  $P_E$ , and this can be done directly via (34), and there is no need to invoke probabilities conditional on null subsets of the parameter space.

All that we need for Bayesian statistical inference is the probability space  $\langle \Omega \times \Gamma, \mathcal{H}, P \rangle$ , and operations on this, including conditionalization on new evidence, can go through in the standard way, without invoking any conditional probabilities conditional on null subsets of the parameter space. We can, if convenient, work with the likelihood functions  $\mathcal{L}_E(\theta)$ , whose existence is guaranteed by the Radon-Nikodym theorem. But there is no *need* to regard these as *bona fide* conditional probabilities, and their usefulness as calculational tools does not depend on any such interpretation.

## 5.2 Stochastic processes

In the theory of stochastic processes, we deal with a set  $\{X_t \mid t \in T\}$  of random variables, where the index  $t$  is to be thought of as a time index (which may be continuous or

---

<sup>18</sup>To see this: suppose that the probability distribution  $Q$  over the parameters is given by a density function  $\mu$ . That is, for all Borel subsets  $\Delta$  of the parameter space,

$$Q(\Theta \in \Delta) = \int_{\Delta} \mu(\theta) d\theta. \quad (35)$$

Then, by (32),

$$P(E \& \Theta \in \Delta) = \int_{\Delta} \mathcal{L}_E(\theta) \mu(\theta) d\theta. \quad (36)$$

Combing this with (34), we get

$$P_E(\Theta \in \Delta) = \frac{P(E \& \Theta \in \Delta)}{P(E)} = \int_{\Delta} \frac{\mathcal{L}_E(\theta) \mu(\theta)}{P(E)} d\theta. \quad (37)$$

Suppose, now, that  $P_E$  is given by a density function  $\mu_E$ .

$$P_E(\Theta \in \Delta) = \int_{\Delta} \mu_E(\theta) d\theta. \quad (38)$$

Then, setting (74) and (75) equal to each other, we get

$$\int_{\Delta} \mu_E(\theta) d\theta = \int_{\Delta} \frac{\mathcal{L}_E(\theta) \mu(\theta)}{P(E)} d\theta. \quad (39)$$

Since this must be true for every Borel set  $\Delta$ , the integrands must be equal almost everywhere.

discrete). As an example, consider the following simple two-step process, adapted from Bayes (1763). A ball is thrown onto a square table  $ABCD$ , with unit sides, with uniform probability on the square for its landing place. A line is drawn through its landing point, parallel to  $AD$ . We then throw a second ball, again with uniform probability, and are provided with a report of whether the second ball landed to the left or the right of the line we drew. In this case, it is unproblematic to say that, conditional on the first ball's landing at a distance  $x$  from the left side of the table, the chance of the second ball landing to the left of the line is  $x$ .

But we don't have to; everything we need to say about the process can be said without invocation of null-condition conditional probabilities. Let  $X_1$  be the random variable that represents the distance of the landing place of the first ball from the left side of the table, and let  $X_2$  be the random variable that takes on the value  $L$  or  $R$  depending on whether the second ball lands to the left or right of the line through the landing-place of the first ball. We can specify the two-step process by saying that  $X_1$  is uniformly distributed on  $(0, 1)$ , and that conditional probabilities for  $X_2$  are given by

$$\begin{aligned} P(X_2 = L | X_1 = x) &= x \\ P(X_2 = R | X_1 = x) &= 1 - x. \end{aligned} \tag{41}$$

But we can also specify the same probability distribution over possible outcomes of the two-step process by saying that joint probabilities regarding  $X_1$  and  $X_2$  satisfy

$$\begin{aligned} P(X_2 = L \ \&\ X_1 \in \Delta) &= \int_{\Delta} x \, dx \\ P(X_2 = R \ \&\ X_1 \in \Delta) &= \int_{\Delta} (1 - x) \, dx \end{aligned} \tag{42}$$

for every Borel set  $\Delta \subseteq (0, 1)$ . Null-condition conditional probabilities, though they may provide a useful way of talking, are not needed to specify the stochastic process.

More generally, given a stochastic process involving random variables  $\{X_t \mid t \in T\}$ , for any time  $t_0$  we can consider the set of random variables with  $t \leq t_0$ , and form a  $\sigma$ -algebra  $\mathcal{T}_0$  generated by this set of random variables. For any proposition of the form  $X_r \in \Delta$ , we will have conditional probabilities with respect to  $\mathcal{T}_0$ ,  $P(X_r \in \Delta \mid \mathcal{T}_0)$ . The values these take on the atoms of  $\mathcal{T}_0$  may be regarded as probabilities conditional on a full specification of events up to  $t_0$ , even if these atoms have zero probability.

Cautions that by now are familiar are in place: though the set-up gives us a privileged  $\sigma$ -algebra, namely, the  $\sigma$ -algebra corresponding to a full specification of events up to  $t_0$ , including these events in our set of admissible conditions for conditional probability still leaves us with a rather sparse set of conditions, and problems and ambiguities may arise if we seek to include in our set of conditions null propositions with less than complete information about the past. Secondly, the stochastic process only specifies these conditional probabilities for almost all histories; different versions of the conditional probabilities may differ on probabilities conditional on past histories

in some set of measure zero. These are not taken as corresponding distinct stochastic processes, as they yield the same probability for any set of events.

The Eternal Coin example of §3.3 illustrates this latter point. Let  $\mathcal{P}$  be the  $\sigma$ -algebra consisting of propositions about results of coin tosses to the past of today. The atoms of this  $\sigma$ -algebra comprise all possible complete specifications of the past; the proposition  $P$ , that the coin landed heads every day in the past, is one such. The proposition  $H$ , that the coin lands heads today, is independent of the  $\sigma$ -algebra  $\mathcal{P}$ . That is,

$$P(AH) = P(A)P(H) \tag{43}$$

for all  $A \in \mathcal{P}$ . This entails that we must have

$$P(H|\mathcal{P})(u) = P(H) = 1/2 \tag{44}$$

for almost all  $u$  in our event space. But this doesn't preclude Dorr, or anyone else so inclined, from assigning the value 1, or any other value, to the probability of  $H$  conditional on the proposition  $P$ , or on any set of propositions comprising a set of measure zero. Distinct choices of this sort yield the same probabilities for all propositions about histories.

This agreement of conditional probabilities up to a set of conditions of total probability zero is, arguably, all that matters when it comes to formulation of a stochastic process. Suppose that we are formulating a physical theory with stochastic dynamics, and formulate the theory in terms of transition probabilities, that is, probabilities about future events conditional on past events. Suppose that we have two specifications of such transition probabilities, that agree for almost all conditions, but assign different conditional probabilities to some events, on a set of conditions of total probability zero. Do we have here two distinct physical theories, or two different formulations of the same theory?

Both specifications share the same possibility space, that is, the same set of possible histories. They agree on which propositions about histories are to be assigned probabilities, and they agree on what those probabilities are. These things are all that matters, in formulating a stochastic theory; the means we use to specify them is inessential. They should be counted as alternate formulations of essentially the same theory.

### 5.3 The Principal Principle

The Principal Principle, so named by Lewis (1980), is the prescription that your credence at time  $t$  in a proposition  $A$ , conditional on the supposition that the chance at  $t$  of  $A$  is  $x$  and any admissible proposition, be  $x$ . That is,

$$Cr_t(A | ch_t(A) = x \ \& \ E) = x. \tag{45}$$

for any admissible  $E$ , where “[a]dmissible propositions are the sort of information whose impact on credence about outcomes comes entirely by way of credence about



the chances of those outcomes” (Lewis, 1980, p. 272). This is to be true for every value of  $x$  in  $[0, 1]$ . Any credences about the chance of  $A$  will assign zero credence to uncountably many singleton sets. Thus, it looks as if the Principal Principle *commits* us to conditionalizing on null propositions.

This, again, is unobjectionable, as we have a distinguished  $\sigma$ -algebra, consisting of propositions of the form  $ch_t(A) \in \Delta$ , where  $\Delta$  ranges over Borel subsets of  $[0, 1]$ . But use of the Principal Principle itself does not by itself commit us to null-condition probabilities.

The work that the Principal Principle does is to ensure that our credences about  $A$  mesh with our credences about the chance of  $A$ . Suppose that  $\mu_E$  is the probability function that represents an agent’s credences about the chance of  $A$ , conditional on some admissible evidence  $E$ . That is,

$$Cr_t(ch_t(A) \in \Delta | E) = \int_{\Delta} d\mu_E(x). \quad (46)$$

Then we will have (assuming that the conditional probabilities are defined),

$$Cr_t(A | E) = \int Cr_t(A | ch_t(A) = x \ \& \ E) d\mu_E(x). \quad (47)$$

Imposing the Principal Principle entails that

$$Cr_t(A | E) = \int x d\mu_E(x). \quad (48)$$

This condition encapsulates the effect of Principal Principle on the agent’s credence: it ensures that the agent’s credences about  $A$  mesh properly with credence about the chance of  $A$ .

We can achieve the same effect—that is, ensure satisfaction of (48)—without reference to probabilities conditional on null propositions. All need do is prescribe that, for every interval  $\Delta \subseteq [0, 1]$  with  $Cr_t(ch_t(A) \in \Delta) > 0$ , and any admissible  $E$ ,

$$Cr_t(A | E \ \& \ ch_t(A) \in \Delta) \in \Delta. \quad (49)$$

Satisfaction of this condition entails that joint credences about a proposition and the chance of that proposition satisfy (48). Thus, even without primitive probabilities conditional on null propositions, we get what we need.

We can readily extend (49) to credences about multiple propositions and their chances. For any finite set  $\mathbf{A} = \{A_1, \dots, A_n\}$  of propositions, we require that, for all measurable  $\Delta \subseteq [0, 1]^n$ , with  $Cr_t(ch_t(\mathbf{A} \in \Delta)) > 0$ ,

$$Cr_t(A_1 \ \& \ A_2 \ \& \dots \ \& \ A_n | E \ \& \ ch_t(\mathbf{A}) \in \Delta) \in \text{Conv}(\Delta), \quad (50)$$

where  $\text{Conv}(\Delta)$  is the convex hull of  $\Delta$ .

## 6 Conclusion

Talk of probabilities conditional on zero-probability propositions is common in the philosophical literature. There is nothing *necessarily* incoherent in such talk, and we may, for certain purposes, find it convenient to include such propositions in the stock of proposition on which we conditionalize. But the motivations for doing so have been exaggerated.

Moreover, though symmetry considerations may guide us in choice of probability distribution conditional on null propositions, such considerations can be less than reliable guides. Imposing the requirement that the conditional probability space be invariant under all symmetries of the unconditional probability space is excessively restrictive. If we want to extend our conditional probability space to include conditionalization on null propositions, we will have to be selective about which symmetries of the unconditional probability space we impose on the conditional probability space. In some cases—such as the sphere—there may be a natural choice of which symmetries to impose. In other cases, of which Dorr’s Eternal Coin is a striking example, symmetry considerations will lead us in opposing directions, without a clear choice to be made.

If, nonetheless, you want to include null propositions in your set of conditions: proceed with caution, and with care to state explicitly how your conditional probability space is to be constructed.

## 7 Acknowledgments

I thank Alan Hájek, Bill Harper and Joshua Luczak for helpful discussions on these matters. I also thank two anonymous referees for *Erkenntnis*, who read the manuscript with extraordinary care and saved me from numerous errors (some minor, some less so). This work was sponsored by a grant from the Social Sciences and Humanities Research Council of Canada (SSHRC).

# Appendix 1 Terminology

## 1.1 Probability Spaces

For any set  $S$ , an *algebra* of subsets of  $S$  is a set of subsets of  $S$  that contains  $S$  and is closed under complementation and unions. A  $\sigma$ -*algebra* of subsets of  $S$  is an algebra that is closed under countable unions. For the real line  $\mathbb{R}$ , we define the Borel sets as the smallest  $\sigma$ -algebra containing all open intervals.

If  $\mathcal{A}$  is an algebra of subsets of  $S$ , a function  $P : \mathcal{A} \rightarrow \mathbb{R}$  is *additive* iff, for any disjoint  $A, B \in \mathcal{A}$ ,

$$P(A \cup B) = P(A) + P(B).$$

If  $\mathcal{A}$  is a  $\sigma$ -algebra of subsets of  $S$ , a function  $P : \mathcal{A} \rightarrow \mathbb{R}$  is *countably additive* iff, for any sequence  $\{A_i\}$  of disjoint sets in  $\mathcal{A}$ ,

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i).$$

A *probability space* is a triple  $\langle S, \mathcal{A}, P \rangle$ , where  $S$  is a set, to be thought of as the set of elementary events,  $\mathcal{A}$  is an algebra of subsets of  $S$ , which are the sets of events (propositions) to which probabilities will be ascribed, and  $P : \mathcal{A} \rightarrow \mathbb{R}$  is a probability function, that is, a positive, additive set function with  $P(S) = 1$ . Since we will have reasons to consider probability functions that are not countably additive, we depart from tradition in not assuming countable additivity unless explicitly stated. If we require countable additivity, then  $\mathcal{A}$  is required to be a  $\sigma$ -algebra, and we will refer to  $P$  as a *probability measure*.

If  $\langle S, \mathcal{A}, P \rangle$  is a probability space, a *random variable* is a measurable function  $X : S \rightarrow \mathbb{R}$ , that is, a function such that, for any Borel set  $B$ , the set

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

is in  $\mathcal{A}$ . A random variable  $X$  generates a subalgebra of  $\mathcal{A}$ , called  $\sigma(X)$ , which is the set of all  $X^{-1}(B)$ , as  $B$  ranges over Borel subsets of the real line.

## 1.2 Conditional Probability Spaces

Following Rényi (1955, 2007b),<sup>19</sup> we define a *conditional probability space* as a quadruplet  $\langle S, \mathcal{A}, \mathcal{B}, P \rangle$ , where  $S$  is a set of events,  $\mathcal{A}$  an algebra of subsets of  $S$ ,  $\mathcal{B}$  a subset

---

<sup>19</sup>Though inspired and instructed by Rényi's treatment, this definition departs from Rényi in two ways. First, Rényi requires  $P_B(A)$  to be defined for every  $A \in \mathcal{A}$ . This may be undesirable; see Appendix 2. Second, Rényi requires countable additivity, and we leave open the possibility of conditional probability functions that are merely finitely additive.

In his later book, (2007a), Rényi revises the definition of a conditional probability space to exclude zero-probability conditions, and further requires that the set  $\mathcal{B}$  of conditions be closed under finite disjunctions, and that it contain a sequence  $\{B_n\}$  that covers  $\Omega$  (see §2.2). A subset of a  $\sigma$ -algebra  $\mathcal{A}$  satisfying these two

of  $\mathcal{A}$ , to be thought of as the set of events on which we may conditionalize, and  $P$  is a function that takes  $B \in \mathcal{B}$  to a function  $P_B : \mathcal{A}_B \rightarrow \mathbb{R}$ , where, for each  $B$ ,  $\mathcal{A}_B$  is a subalgebra of  $\mathcal{A}$ , and

- i). For each  $B \in \mathcal{B}$ 
  - (a)  $P_B(A) \geq 0$  for all  $A \in \mathcal{A}_B$ .
  - (b) For all  $A \in \mathcal{A}$ , if  $B \subseteq A$ , then  $A \in \mathcal{A}_B$  and  $P_B(A) = 1$ .
  - (c) For disjoint  $A, A' \in \mathcal{A}_B$ ,  $P_B(A \cup A') = P_B(A) + P_B(A')$ .
- ii). For all  $B, C \in \mathcal{B}$  and  $A, B \in \mathcal{A}_C$ , if  $BC \in \mathcal{B}$  then

$$P_C(AB) = P_{BC}(A) P_C(B),$$

provided that  $B \in \mathcal{A}_C$  and  $A \in \mathcal{A}_{BC}$ .

A conditional probability space can be thought of as a family of probability spaces  $\{\langle S, \mathcal{A}_B, P_B \rangle \mid B \in \mathcal{B}\}$ , required to mesh with each other via (ii).

It is an immediate consequence of (ii) that, for any  $C \in \mathcal{B}$  and  $B \subseteq C$ , if  $B \in \mathcal{A}_C$  and  $P_C(B) > 0$ , then, for all  $A \in \mathcal{A}_C$ ,

$$P_B(A) = \frac{P_C(AB)}{P_C(B)} \quad (51)$$

provided that  $B \in \mathcal{B}$  and  $A \in \mathcal{A}_B$ . This allows us to define probabilities conditional on  $B$ , provided they don't clash with those yielded by some other  $D \in \mathcal{B}$  such that  $B \subseteq D$ ,  $B \in \mathcal{A}_D$ , and  $P_D(B) > 0$ . For this reason, we will usually assume the further condition,

- iii). For all  $C \in \mathcal{B}$  and  $B \subseteq C$ , if  $B \in \mathcal{A}_C$  and  $P_C(B) > 0$ , then  $B \in \mathcal{B}$  and  $\mathcal{A}_C \subseteq \mathcal{A}_B$ .

Given a probability space  $\langle S, \mathcal{A}, P \rangle$ , let  $\mathcal{A}^*$  be the subset of  $\mathcal{A}$  consisting of sets  $B$  with  $P(B) > 0$ . Let  $P^*$  be the function that maps  $B \in \mathcal{A}^*$  to the probability function  $P_B : \mathcal{A} \rightarrow [0, 1]$ , given by

$$P_B(A) = \frac{P(AB)}{P(B)}. \quad (52)$$

Then  $\langle S, \mathcal{A}, \mathcal{A}^*, P^* \rangle$  is a conditional probability space, corresponding to the standard choice of having conditional probability defined only when the condition has nonzero probability.

We will say that a probability space  $\langle S, \mathcal{A}, P \rangle$  is invariant under a bijection  $\mathbb{T} : S \rightarrow S$  if and only if  $\mathbb{T}$  leaves the set  $\mathcal{A}$  of measurable sets invariant, and, for all  $A \in \mathcal{A}$ ,  $P(\mathbb{T}(A)) = P(A)$ , where  $\mathbb{T}(A)$  is  $\{\mathbb{T}(x) \mid x \in A\}$ . Similarly, a conditional probability space  $\langle S, \mathcal{A}, \mathcal{B}, P \rangle$  is invariant under  $\mathbb{T}$  if and only if  $\mathbb{T}(\mathcal{A}) = \mathcal{A}$ ,  $\mathbb{T}(\mathcal{B}) = \mathcal{B}$ , and, for all  $B \in \mathcal{B}$ ,  $\mathcal{A}_{\mathbb{T}(B)} = \mathbb{T}(\mathcal{A}_B)$  and  $P_{\mathbb{T}(B)}(\mathbb{T}(A)) = P_B(A)$  for all  $A \in \mathcal{A}_B$ .

---

conditions, and not containing the null set, Rényi calls a *bunch* of sets. In this work, Rényi calls conditional probability spaces in which conditionalization on null propositions is permitted *generalized conditional probability spaces* (see Problem 2.8).

### 1.3 Lebesgue measure

Consider the unit interval  $I = (0, 1]$ . Let  $\mathcal{B}$  be the the smallest  $\sigma$ -algebra that contains all intervals  $(a, b]$ . These are the *Borel sets*. Extension to the entire real line, or to  $\mathbb{R}^n$ , is straightforward; the Borel subsets of  $\mathbb{R}^n$  are the elements of the smallest  $\sigma$ -algebra containing all rectangles of the form

$$\{(x_1, \dots, x_n) \mid x_i \in (a_i, b_i], i = 1, \dots, n\}.$$

The *uniform measure*, or *Lebesgue measure*, on  $\langle I, \mathcal{B} \rangle$  is the unique countably additive measure that assigns measure  $b - a$  to each interval  $(a, b]$ . Call this measure  $\lambda$ . Though the Borel sets include all sets of interest for most purposes, we can extend our measure to a wider  $\sigma$ -algebra, called the Lebesgue measurable sets. We define an outer measure  $P^*$ , defined for any subset  $A$  of the unit interval, by

$$P^*(A) = \inf \sum_i \lambda(A_i), \tag{53}$$

where the infimum is taken over all countable collections  $\{A_i\}$  of intervals such that  $A \subseteq \bigcup_i A_i$ . The set  $\mathcal{L}$  of Lebesgue measurable sets consists of all sets  $A \subseteq I$  such that

$$P^*(AE) + P^*(A^cE) = P^*(E) \tag{54}$$

for all  $E \subseteq I$ , where  $A^c$  is the complement of  $A$  in  $I$ . It can be shown that this is a  $\sigma$ -algebra, and that  $\lambda$  has a unique countably additive extension to  $\mathcal{L}$ ; this extension is also called Lebesgue measure. Again, extension of this concept to  $\mathbb{R}^n$  is straightforward.

It can be shown that, if  $A$  is a Lebesgue-measurable subset of  $\mathbb{R}^n$  with Lebesgue measure zero, then every subset of  $A$  is Lebesgue measurable (and, of course, also has Lebesgue measure zero).

## Appendix 2 A rotationally invariant conditional probability space

On the  $n$ -sphere  $S^n$  (that is, the  $n$ -dimensional space of points at unit distance from a fixed point in  $n + 1$ -dimensional Euclidean space, which is a circle for  $n = 1$ , and the surface of a sphere for  $n = 2$ ), we construct a uniform spherical measure  $\sigma_n$ . One way to characterize these measures is in terms of Lebesgue measure on the ambient  $(n + 1)$ -dimensional space. A subset of  $A$  of  $S_n$  is  $\sigma_n$ -measurable if and only the wedge subtended by  $A$ —that is, the set of points on straight lines between points of  $A$  and the origin—is Lebesgue measurable, and we take the measure of  $A$  to be proportional to the Lebesgue measure of the wedge it subtends. Let  $\mathcal{L}_n$  be the set of all  $\sigma_n$ -measurable subsets of  $S_n$ .

Let  $S$  be a 2-sphere—that is, the surface of a sphere in 3D space—and let  $\mathcal{L}_S$  be the set of all  $\sigma_2$ -measurable subsets of  $S$ , and let  $\sigma_S$  be  $\sigma_2$  measure on  $\langle S, \mathcal{L}_S \rangle$ . Let  $\mathcal{C}$

be the set of all circles on the sphere  $S$ . For each circle  $C \in \mathcal{C}$ , let  $\mathcal{L}_C$  be the set of  $\sigma_1$ -measurable subsets of  $C$ , and let  $\sigma_C$  be  $\sigma_1$  measure on  $\langle C, \mathcal{L}_C \rangle$ .

We can construct standard conditional probability spaces  $\langle S, \mathcal{L}_S, \mathcal{L}_S^*, P_S^* \rangle$  and  $\langle C, \mathcal{L}_C, \mathcal{L}_C^*, P_C^* \rangle$ , where, as in 1.2,  $\mathcal{L}_S^*$  is the set of sets in  $\mathcal{L}_S$  with positive probability. We want to extend  $\langle S, \mathcal{L}_S, \mathcal{L}_S^*, P_S^* \rangle$  to encompass, in our set of conditions, other sets, including, at minimum, all circles, in such a way that conditionalizing on any circle yields uniform probabilities on that circle

Let us take  $\mathcal{B}$  to be

$$\mathcal{B} = \mathcal{L}_S^* \cup \bigcup_{C \in \mathcal{C}} \mathcal{L}_C^*. \quad (55)$$

Note that each element of  $\mathcal{B}$  is either in  $\mathcal{L}_S^*$  or is a subset of a *unique* circle  $C$  (this is because the members of  $\mathcal{L}_C^*$  have positive measure, and hence contain infinitely many points, and no two circles share more than two points). Take  $\mathcal{A}_B$  to be  $\mathcal{L}_S$  for  $B \in \mathcal{L}_S^*$ . For  $B \in \mathcal{L}_C^*$ , take  $A \in \mathcal{A}_B$  iff  $AB \in \mathcal{L}_C$ . Define

$$P_B(A) = \begin{cases} \frac{\lambda_S(AB)}{\lambda_S(B)}, & B \in \mathcal{L}_S^*; \\ \frac{\lambda_C(AB)}{\lambda_C(B)}, & B \in \mathcal{L}_C^*. \end{cases} \quad (56)$$

We have constructed a conditional probability space that is invariant under all rigid rotations, and includes conditionalization on circles and some subsets of circles.

For any  $A \in \mathcal{L}_S$  with  $\lambda_S(A) = 0$ , every subset of  $A$  is a measurable set, and is assigned measure 0. Since each circle  $C$  has  $\lambda_S(C) = 0$ , this means that every subset of  $C$  is in  $\mathcal{L}_S$ . Since *not* every subset of  $C$  is in  $\mathcal{L}_C$ ,  $P_C(A)$  is not defined for arbitrary  $A \in \mathcal{L}_S$ . We might want to extend  $P_C$  so that it is defined on all  $A \in \mathcal{L}_S$ . But, as already mentioned (see footnote 2), we can do so, and preserve rotational invariance, only at the price of sacrificing countable additivity. We can't get all that we want.

## Appendix 3 The Eternal Coin: Proof of Dorr's theorem

We will speak in general terms, but readers should think of the example at hand, that of the Eternal Coin. We assume Axiom (iii) of Appendix 1.2. In this appendix, we will find ourselves conditionalizing on complex propositions, and it will be convenient to shift from the subscript notation  $Pr_C$  for probabilities conditional on conditional probabilities to the notation  $Pr(\cdot | C)$ .

Suppose there is a proposition  $P$ , and a transformation  $\mathbb{T}$ , such that  $\mathbb{T}(P) \models P$ . If  $Pr(P) > 0$ , then

$$Pr(\mathbb{T}(P) | P) = \frac{Pr(\mathbb{T}(P))}{Pr(P)}, \quad (57)$$

and so  $\mathbb{T}$ -invariance would entail that  $Pr(\mathbb{T}(P) | P) = 1$ . Furthermore, if there exists a proposition  $Z$  such that  $\mathbb{T}(Z) = Z$ ,  $P \models Z$ , and  $Pr(P|Z) > 0$ , then

$$Pr(\mathbb{T}(P) | P) = \frac{Pr(\mathbb{T}(P) | Z)}{Pr(P | Z)}, \quad (58)$$

and so, once again,  $\mathbb{T}$ -invariance would entail that  $P(\mathbb{T}(P) | P) = 1$ .

But  $\mathbb{T}$ -invariant propositions of the right sort may be hard to come by, and there may be no such  $Z$ . Suppose, however, that there exist propositions  $X$ ,  $Z$ , such that  $P \models X \models Z$ , and  $\mathbb{T}(P) \models \mathbb{T}(X) \models Z$ . Then, if  $P$ ,  $\mathbb{T}(P)$ ,  $X$ , and  $\mathbb{T}(X)$  are all in  $\mathcal{A}_Z$ , and  $Pr(P|Z) > 0$ , and if  $Pr(\mathbb{T}(P) | \mathbb{T}(X))$  is defined, we have

$$Pr(P | Z) = Pr(P | X) Pr(X | Z); \quad (59)$$

$$Pr(\mathbb{T}(P) | Z) = Pr(\mathbb{T}(P) | \mathbb{T}(X)) Pr(\mathbb{T}(X) | Z),$$

and so, still assuming that  $\mathbb{T}(P) \models P$ , we have

$$Pr(\mathbb{T}(P) | P) = \frac{Pr(\mathbb{T}(P) | Z)}{Pr(P | Z)} = \frac{Pr(\mathbb{T}(P) | \mathbb{T}(X))}{Pr(P | X)} \frac{Pr(\mathbb{T}(X) | Z)}{Pr(X | Z)}. \quad (60)$$

Now suppose that there is also a proposition  $F$  such that  $\mathbb{T}^{-1}(F) \models F$ , with  $Pr(F|Z) > 0$ . Suppose, also, that  $\mathbb{T}^{-1}(F) \models X$ . From this it follows that  $F \models \mathbb{T}(X)$ , and we have

$$Pr(\mathbb{T}^{-1}(F) | F) = \frac{Pr(\mathbb{T}^{-1}(F) | X)}{Pr(F | \mathbb{T}(X))} \frac{Pr(X | Z)}{Pr(\mathbb{T}(X) | Z)}. \quad (61)$$

Multiplying (60) and (61) gives us,

$$Pr(\mathbb{T}(P) | P) Pr(\mathbb{T}^{-1}(F) | F) = \frac{Pr(\mathbb{T}(P) | \mathbb{T}(X))}{Pr(P | X)} \frac{Pr(\mathbb{T}^{-1}(F) | X)}{Pr(F | \mathbb{T}(X))}. \quad (62)$$

So far, we haven't invoked any symmetry assumptions. If we impose  $\mathbb{T}$ -invariance, we have

$$\begin{aligned} Pr(\mathbb{T}(P) | \mathbb{T}(X)) &= Pr(P | X); \\ Pr(\mathbb{T}^{-1}(F) | X) &= Pr(F | \mathbb{T}(X)), \end{aligned} \quad (63)$$

and (62) becomes

$$Pr(\mathbb{T}(P) | P) Pr(\mathbb{T}^{-1}(F) | F) = 1, \quad (64)$$

from which it follows that

$$Pr(\mathbb{T}(P) | P) = Pr(\mathbb{T}^{-1}(F) | F) = 1. \quad (65)$$

Now, since we have assumed that  $\mathbb{T}(P) \models P$  and  $\mathbb{T}^{-1}(F) \models F$ , there always do exist  $Z, X$  satisfying the conditions stipulated. Take  $Z$  to be  $P \vee F$ , and take  $X$  to be  $P \vee \mathbb{T}^{-1}(F)$ . Then  $\mathbb{T}(X)$  is  $\mathbb{T}(P) \vee F$ .

To sum up: we have established

**Proposition 1** *Let  $\langle \Omega, \mathcal{A}, \mathcal{B}, Pr \rangle$  be a conditional probability space that satisfies condition (iii) and is invariant under a transformation  $\mathbb{T}$ . Suppose there are propositions  $P, F$ , such that  $Z = P \vee F \in \mathcal{B}$  and  $P, F \in \mathcal{A}_Z$ , such that*

- i). (a)  $\mathbb{T}(P) \models P$ ;*  
*(b)  $\mathbb{T}^{-1}(F) \models F$ ;*
- ii). (a)  $Pr(P|Z) > 0$ ;*  
*(b)  $Pr(F|Z) > 0$ .*

*Then*

$$Pr(\mathbb{T}(P) | P) = Pr(\mathbb{T}^{-1}(F) | F) = 1.$$

Applied to the Eternal Coin, let  $\mathbb{T}$  be  $S_1$ , which shifts everything forward one day.  $P$ , as before, is the proposition that the coin landed Heads every day in the past, and  $F$ , the proposition that the coin will land Heads every day in the future. Let  $H$  be the proposition that the coin lands Heads today. Then  $S_1(P)$  is  $HP$ , and  $S_1^{-1}(F)$  is  $HF$ . Clearly,  $HP \models P$  and  $HF \models F$ . If

$$\begin{aligned} Pr(P|P \vee F) &> 0; \\ Pr(F|P \vee F) &> 0; \end{aligned} \tag{66}$$

and if

$$\begin{aligned} Pr(HP|HP \vee F) &= Pr(P|P \vee HF); \\ Pr(F|HP \vee F) &= Pr(HF|P \vee HF), \end{aligned} \tag{67}$$

then

$$Pr(H|P) = Pr(H|F) = 1. \tag{68}$$

We can run the same argument with  $S_k$ , for any positive  $k$ , yielding the conclusion that, for every  $n \geq 0$ , the probability conditional on  $P$  that the coin lands Heads today and  $n$  days into the future is 1, as is the probability, conditional on  $F$ , that the coin lands Heads today and  $n$  days into the past.



## Notes

<sup>1</sup>We will also use the notation  $P(A|C)$ , when convenient.

<sup>2</sup> The fact that, in probability theory, we can't always get what we want, is a familiar fact. We might want our probability function to be defined on arbitrary subsets of our probability space, but, as is well-known, we can't always do so while satisfying *desiderata* such as symmetry conditions and countable additivity. Consider, for example, the task of defining a uniform distribution—that is, a distribution invariant under all rotations—on the unit circle. There can be no distribution that is invariant under rotations, is countably additive, and is defined on all subsets of the unit circle. The proof is found in many probability texts, *e.g.* Billingsley (2012, p. 47). The standard response is to preserve countable additivity and to restrict the domain of definition of the probability function to certain subsets of the probability space, the measurable sets, leaving the probability of other sets undefined. In one and two dimensions, as Banach (1923) showed, one can extend the probability function to one defined on arbitrary subsets, if one is willing to give up countable additivity. The well-known Banach-Tarski paradox shows that we can't do so in three-dimensional space; there is no finitely additive set function that is defined on all subsets and invariant under translations and rotations.

<sup>3</sup>Adapted from Hájek (2003).

<sup>4</sup>Based on Kolmogorov (1950, §V.2), which in turn is based on Borel (1909, §45) (§8.6 of Borel 1965). See also Jaynes (2003, §15.7), Hájek (2003, §4.4).

<sup>5</sup>To see this: let  $X$  be any random variable, with distribution  $\mu_X$ , and take  $g(X)$  to be the function of  $X$  given by

$$g(X) = \int_{-\infty}^X d\mu_X(x). \quad (69)$$

Then  $g$  has range in  $[0, 1]$ , and, for any  $a \in [0, 1]$ ,

$$P(g(X) \leq a) = a. \quad (70)$$

That is,  $g(X)$  is uniformly distributed on  $[0, 1]$ .

<sup>6</sup>If  $C$  is a finite set, we can have a probability function that always assigns equal probabilities to sets of equal cardinality. This is not possible if  $C$  is infinite.

In the infinite case, there must be measurable sets  $A, B$ , of equal cardinality, with  $P(A) \neq P(B)$ . We can then choose some mapping that takes  $A$  to  $B$ .

<sup>7</sup>This is necessary because, if one has nonzero credence that the coin is not fair, or that the tosses are not independent, then conditionalization on either  $F$  or  $P$  will send credence that the setup is as described to zero.

<sup>8</sup>See Appendix 1 for definitions of any terms that might be unfamiliar.

<sup>9</sup>In this section, we use  $Pr$  for our probability function to avoid confusion with the proposition  $P$ .

<sup>10</sup>In this section, we will find ourselves conditionalizing on some fairly complex propositions, and so it will be convenient to switch from the subscript notation for conditional probabilities used in the rest of the paper to the slash notation.

<sup>11</sup>That is, take

$$A_n = \bigcup_{u \in A} C_{K_n}(u). \quad (71)$$

<sup>12</sup>Perhaps. The more one thinks about what is required to give values to these conditional probabilities, the less clear it becomes that we have intuitions about them at all.

<sup>13</sup>Oddly enough, this has been disputed. In connection with this example, E.T. Jaynes (2003, p. 470) writes,

Nearly everybody feels that he knows perfectly well what a great circle is; so it is difficult to get people to see that the term ‘great circle’ is ambiguous until we specify what limiting operation is to produce it.

This strikes me as confused. One and the same great circle can be the limit of many different decreasing sequences of subsets of the sphere, but the circle is not itself *produced* by the limiting operation. Not so with probabilities conditional on a great circle, which, unless stipulated as primitive, are obtained via some limiting operation.

<sup>14</sup>In his discussion of the Borel paradox, Kolmogorov writes, “This shows that the concept of a probability conditional on an isolated given hypothesis whose probability equals 0 is inadmissible” (Kolmogorov, 1950, p. 51).

<sup>15</sup>The notation is intended to be both reminiscent of, and distinct from, the notation used for conditional probabilities.

<sup>16</sup>This is example 33.11 of Billingsley (2012).

<sup>17</sup>This is even easier to see in these days in which laboratory equipment has digital readout than it was in the old days of pointers and dials!

<sup>18</sup>To see this: suppose that the probability distribution  $Q$  over the parameters is given by a density function  $\mu$ . That is, for all Borel subsets  $\Delta$  of the parameter space,

$$Q(\Theta \in \Delta) = \int_{\Delta} \mu(\theta) d\theta. \quad (72)$$

Then, by (32),

$$P(E \& \Theta \in \Delta) = \int_{\Delta} \mathcal{L}_E(\theta) \mu(\theta) d\theta. \quad (73)$$

Combing this with (34), we get

$$P_E(\Theta \in \Delta) = \frac{P(E \& \Theta \in \Delta)}{P(E)} = \int_{\Delta} \frac{\mathcal{L}_E(\theta) \mu(\theta)}{P(E)} d\theta. \quad (74)$$

Suppose, now, that  $P_E$  is given by a density function  $\mu_E$ .

$$P_E(\Theta \in \Delta) = \int_{\Delta} \mu_E(\theta) d\theta. \quad (75)$$

Then, setting (74) and (75) equal to each other, we get

$$\int_{\Delta} \mu_E(\theta) d\theta = \int_{\Delta} \frac{\mathcal{L}_E(\theta) \mu(\theta)}{P(E)} d\theta. \quad (76)$$

Since this must be true for every Borel set  $\Delta$ , the integrands must be equal almost everywhere.

<sup>19</sup>Though inspired and instructed by Rényi's treatment, this definition departs from Rényi in two ways. First, Rényi requires  $P_B(A)$  to be defined for every  $A \in \mathcal{A}$ . This may be undesirable; see Appendix 2. Second, Rényi requires countable additivity, and we leave open the possibility of conditional probability functions that are merely finitely additive.

In his later book, (2007a), Rényi revises the definition of a conditional probability space to exclude zero-probability conditions, and further requires that the set  $\mathcal{B}$  of conditions be closed under finite disjunctions, and that it contain a sequence  $\{B_n\}$  that covers  $\Omega$  (see §2.2). A subset of a  $\sigma$ -algebra  $\mathcal{A}$  satisfying these two conditions, and not containing the null set, Rényi calls a *bunch* of sets. In this work, Rényi calls conditional probabilities spaces in which conditionalization on null propositions is permitted *generalized conditional probability spaces* (see Problem 2.8).

## References

- Banach, S. (1923). Sur le problème de la mesure. *Fundamenta Mathematicae* 4, 7–33.
- Bayes, T. (1763). An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions* 53, 370–418.
- Billingsley, P. (2012). *Probability and Measure, Anniversary Edition*. Hoboken, NJ: Wiley.
- Borel, E. (1909). *Éléments de la Théorie des Probabilités*. Paris: Librairie Scientifique A. Hermann & Fils. English translation in Borel (1965).
- Borel, E. (1965). *Elements of the Theory of Probability*. Englewood Cliffs, NJ: Prentice-Hall. English translation of Borel (1909).
- Carnap, R. (1950). *The Logical Foundations of Probability*. The University of Chicago Press.
- Dorr, C. (2010). The Eternal Coin: A puzzle about self-locating conditional credence. *Philosophical Perspectives* 24, 189–205.
- Easwaran, K. (2011). The varieties of conditional probability. In P. Bandyopadhyay and M. Forster (Eds.), *Handbook of the Philosophy of Science. Philosophy of Statistics*, pp. 137–148. Amsterdam: North-Holland.
- Easwaran, K. K. (2008). The foundations of conditional probability. Doctoral Dissertation, Department of Philosophy, University of California, Berkeley. Available at <http://www.kennyeaswaran.org/research>.
- Hájek, A. (2003). What conditional probability could not be. *Synthese* 137, 273–323.
- Harper, W. and A. Hájek (1997). Full belief and probability: Comments on van Fraassen. *Dialogue* 36, 91–100.
- Harper, W. L. (1975). Rational belief change, Popper functions, and counterfactuals. *Synthese* 30, 221–262.
- Jaynes, E. T. (2003). *Probability Theory: The Logic of Science*. Cambridge: Cambridge University Press.
- Kolmogorov, A. (1950). *Foundations of the Theory of Probability*. New York: Chelsea Publishing Company. Tr. Nathan Morrison.
- Lewis, D. (1980). A subjectivist’s guide to objective chance. In R. C. Jeffrey (Ed.), *Studies in Inductive Logic and Probability*, Volume II, pp. 263–93. University of California Press.

- Popper, K. R. (1938). A set of independent axioms for probability. *Mind* 47, 275–277.
- Popper, K. R. (1955). Two autonomous axiom systems for the calculus of probabilities. *The British Journal for the Philosophy of Science* 6, 51–57.
- Popper, K. R. (1959). *The Logic of Scientific Discovery*. London: Hutchinson.
- Rényi, A. (1955). On a new axiomatic theory of probability. *Acta Mathematica Hungarica* 6, 265–333.
- Rényi, A. ([1970] 2007a). *Foundations of Probability*. Mineola, NY: Dover Publications, Inc. Reprint of edition published by Holden-Day, 1970.
- Rényi, A. ([1970] 2007b). *Probability Theory*. Mineola, NY: Dover publications, Inc. Reprint of edition published by North-Holland, 1970.
- van Frassen, B. C. (1976). Representation of conditional probabilities. *Journal of Philosophical Logic* 5, 417–430.