

Paradoxes of probability theory

I protest against the use of infinite magnitude as something accomplished, which is never permissible in mathematics. Infinity is merely a figure of speech, the true meaning being a limit.

C. F. Gauss

The term ‘paradox’ appears to have several different common meanings. Székely (1986) defines a paradox as anything which is true but surprising. By that definition, every scientific fact and every mathematical theorem qualifies as a paradox for someone. We use the term in almost the opposite sense; something which is absurd or logically contradictory, but which appears at first glance to be the result of sound reasoning. Not only in probability theory, but in all mathematics, it is the careless use of infinite sets, and of infinite and infinitesimal quantities, that generates most paradoxes.

In our usage, there is no sharp distinction between a paradox and an error. A paradox is simply an error out of control; i.e. one that has trapped so many unwary minds that it has gone public, become institutionalized in our literature, and taught as truth. It might seem incredible that such a thing could happen in an ostensibly mathematical field; yet we can understand the psychological mechanism behind it.

15.1 How do paradoxes survive and grow?

As we stress repeatedly, from a false proposition – or from a fallacious argument that leads to a false proposition – all propositions, true and false, may be deduced. But this is just the danger; if fallacious reasoning always led to absurd conclusions, it would be found out at once and corrected. But once an easy, shortcut mode of reasoning has led to a few correct results, almost everybody accepts it; those who try to warn against it are not listened to.

When a fallacy reaches this stage, it takes on a life of its own, and develops very effective defenses for self-preservation in the face of all criticisms. Mathematicians of the stature of Henri Poincaré and Hermann Weyl tried repeatedly to warn against the kind of reasoning used in infinite-set theory, with zero success. For details, see Appendix B and Kline (1980). The writer was also guilty of this failure to heed warnings for many years, until

absurd results that could no longer be ignored finally forced him to see the error in an easy mode of reasoning.

To remove a paradox from probability theory will require, at the very least, detailed analysis of the result and the reasoning that leads to it, showing that:

- (1) the result is indeed absurd;
- (2) the reasoning leading to it violates the rules of inference developed in Chapter 2;
- (3) when one obeys those rules, the paradox disappears and we have a reasonable result.

There are too many paradoxes contaminating the current literature for us to analyze separately. Therefore we seek here to study a few representative examples in some depth, in the hope that the reader will then be on the alert for the kind of reasoning which leads to them.

15.2 Summing a series the easy way

As a kind of introduction to fallacious reasoning with infinite sets, we recall an old parlor game by which you can prove that any given infinite series $S = \sum_i a_i$ converges to any number x that your victim chooses. The sum of the first n terms is $s_n = a_1 + a_2 + \cdots + a_n$. Then, defining $s_0 \equiv 0$, we have

$$a_n = (s_n - x) - (s_{n-1} - x), \quad 1 \leq n < \infty, \quad (15.1)$$

so that the series becomes

$$\begin{aligned} S &= (s_1 - x) + (s_2 - x) + (s_3 - x) + \cdots \\ &\quad - (s_0 - x) - (s_1 - x) - (s_2 - x) - \cdots. \end{aligned} \quad (15.2)$$

The terms $(s_1 - x)$, $(s_2 - x)$, \dots all cancel out, so the sum of the series is

$$S = -(s_0 - x) = x \quad QED. \quad (15.3)$$

The reader for whom this reasoning appears at first glance to be valid has a great deal of company, and is urged to study this example carefully. Such fallacious arguments are avoided if we follow this advice, repeated from Chapter 2:

Apply the ordinary processes of arithmetic and analysis only to expressions with a finite number n of terms. Then after the calculation is done, observe how the resulting finite expressions behave as the parameter n increases indefinitely.

Put more succinctly, passage to a limit should always be the last operation, not the first. In case of doubt, this is the only safe way to proceed. Our present theory of convergence of infinite series could never have been achieved if its founders – Abel, Cauchy, d'Alembert, Dirichlet, Gauss, Weierstrasz, and others – had not followed this advice meticulously. In pre-Bourbakist mathematics (such as Whittaker and Watson, 1927) this policy was considered so obvious that there was no need to stress it. The results thus obtained have never been found defective.

Had we followed this advice above, we would not have tried to cancel out an infinite number of terms in a single stroke; we would have found that at any finite n th stage, instead

of the s_i cancelling out and one x remaining, the x values would have cancelled out and the last s remains, leading to the correct summation of the series.

Yet today, reasoning essentially equivalent to what we did in (15.2) is found repeatedly where infinite sets are used in probability theory. As an example, we examine another of the consequences of ignoring this advice, which has grown into far more than a parlor game.

15.3 Nonconglomerability

If (C_1, \dots, C_n) denote a finite set of mutually exclusive, exhaustive propositions on prior information I , then for any proposition A the sum and product rules of probability theory give

$$P(A|I) = \sum_{i=1}^n P(AC_i|I) = \sum_{i=1}^n P(A|C_i I)P(C_i|I) \quad (15.4)$$

in which the prior probability $P(A|I)$ is written as a weighted average of the conditional probabilities $P(A|C_i I)$. Now, it is a very elementary theorem that a weighted average of a set of real numbers cannot lie outside the range spanned by those numbers; if

$$L \leq P(A|C_i I) \leq U, \quad (1 \leq i \leq n) \quad (15.5)$$

then necessarily

$$L \leq P(A|I) \leq U, \quad (15.6)$$

a property which de Finetti (1972) called ‘conglomerability’ or, more precisely, ‘conglomerability in the partition $\{C_i\}$ ’, although it may seem too trivial to deserve a name. Obviously, nonconglomerability cannot arise from a correct application of the rules of probability theory on finite sets. It cannot, therefore, occur in an infinite set which is approached as a well-defined limit of a sequence of finite sets.

Yet nonconglomerability has become a minor industry, with a large and growing literature. There are writers who believe that it is a real phenomenon, and that they are proving theorems about the circumstances in which it occurs, which are important for the foundations of probability theory. Nonconglomerability has become, quite literally, institutionalized in our literature and taught as truth.

In spite of its mathematical triviality, then, we need to examine some cases where nonconglomerability has been claimed. Rather than trying to cite all of this vast literature, we draw upon a single reference (Kadane, Schervish and Seidenfeld, 1986), hereafter denoted by KSS, where several examples of nonconglomerability and some references to other work may be found.

Example 1: Rectangular array. Firstly, we note the typical way in which nonconglomerability is manufactured, and the illustrative example most often cited. We start from a two-dimensional $(M \times N)$ set of probabilities:

$$p(i, j), \quad 1 \leq i \leq M, \quad 1 \leq j \leq N, \quad (15.7)$$

and think of i plotted horizontally, j vertically, so that the sample space is a rectangular array of MN points in the first quadrant. It will suffice to take some prior information I for which these probabilities are uniform: $p(i, j) = (1/MN)$. Then the probability of the event $(A : i < j)$ is found by direct counting to be

$$P(A|I) = \begin{cases} (2N - M - 1)/2N & M \leq N \\ (N - 1)/2M & N \leq M. \end{cases} \quad (15.8)$$

Let us resolve this in the manner of (15.4), into probabilities conditional on the set of propositions (C_1, \dots, C_M) , where C_i is the statement that we are on the i th column of the array: then $P(C_i|I) = (1/M)$, and

$$P(A|C_i I) = \begin{cases} (N - i)/N & 1 \leq i \leq M \leq N \\ (N - i)/N & 1 \leq i \leq N \leq M \\ 0 & N \leq i \leq M. \end{cases} \quad (15.9)$$

These conditional probabilities reach the upper and lower bounds

$$U = (N - 1)/N \quad \text{all } M, N, \\ L = \begin{cases} 1 - R & M \leq N \\ 0 & N \leq M. \end{cases} \quad (15.10)$$

where R denotes the ratio $R = M/N$. Substituting (15.8) and (15.10) into (15.6), it is evident that the condition for conglomerability is always satisfied, as it must be, whatever the values of (M, N) . How, then, can one possibly create a nonconglomerability out of this?

Just pass to the limit $M \rightarrow \infty$, $N \rightarrow \infty$, and ask for the probabilities $P(A|C_i I)$ for $i = 1, 2, \dots$. But instead of examining the limiting form of (15.9), which gives the exact values for all (M, N) , we try to evaluate these probabilities directly on the infinite set.

Then, it is argued that, for any given i , there are an infinite number of points where A is true and only a finite number where it is false. *Ergo*, the conditional probability $P(A|C_i I) = 1$ for all i ; yet $P(A|I) < 1$. We see here the same kind of reasoning that we used in (15.2); we are trying to carry out very simple arithmetic operations (counting), but directly on an infinite set.

Now consider the set of propositions (D_1, \dots, D_N) , where D_j is the statement that we are on the j th row of the array, counting from the bottom. Now, by the same argument, for any given j , there are an infinite number of points where A is false, and only a finite number where A is true. *Ergo*, the conditional probability $P(A|D_j I) = 0$ for all j ; yet $P(A|I) > 0$. By this reasoning, we have produced two nonconglomerabilities, in opposite directions, from the same model (i.e. the same infinite set).

It is even more marvellous than that. In (15.8), it is true that if we pass to the limit holding i fixed, the conditional probability $P(A|C_i B)$ tends to one for all i ; but if instead we hold $(N - i)$ fixed, it tends to zero for all i . Therefore, if we consider the cases $(i = 1, i = 2, \dots)$ in increasing order, the probabilities $P(A|C_i B)$ appear to be one for all i . But it is equally

valid to consider them in decreasing order ($i = N, i = N - 1, \dots$); then, by the same reasoning, they would appear to be zero for all i . (Note that we could redefine the labels by subtracting $N + 1$ from each one, thus numbering them ($i = -N, \dots, i = -1$) so that as $N \rightarrow \infty$ the upper indices stay fixed; this would have no effect on the validity of the reasoning.)

Thus, to produce two opposite nonconglomerabilities we need not introduce two different partitions $\{C_i\}$, $\{D_j\}$; they can be produced by two equally valid arguments from a single partition. What produces them is that one supposes the infinite limit already accomplished *before* doing the arithmetic, reversing the policy of Gauss which we recommended above. But if we follow that policy and do the arithmetic first, then an arbitrary redefinition of the labels $\{i\}$ has no effect; the counting for any N is the same.

Once one has understood the fallacy in (15.2), then whenever someone claims to have proved some result by carrying out arithmetic or analytical operations directly on an infinite set, it is hard to shake off a feeling that he could have proved the opposite just as easily and by an equally sound argument, had he wished to. Thus there is no reason to be surprised by what we have just found.

Suppose that instead we had done the calculation by obeying our rules strictly, doing first the arithmetic operations on finite sets to obtain the exact solution (15.8); then passing to the limit. However the infinite limit is approached, the conditional probabilities take on values in a wide interval whose lower bound is zero or $1 - R$, and whose upper bound tends to one. The condition (15.5) is always satisfied, and a nonconglomerability could never have been found.

The reasoning leading to this nonconglomerability contains another fallacy. Clearly, one cannot claim to have produced a nonconglomerability on the infinite set until the 'unconditional' probability $P(A|I)$ has also been calculated on that set, not merely bounded by a verbal argument. But as M and N increase, from (15.8) the limiting $P(A|I)$ depends only on the ratio $R = M/N$:

$$P(A|I) \rightarrow \begin{cases} 1 - R/2 & R \leq 1 \\ 1/(2R) & R \geq 1. \end{cases} \quad (15.11)$$

If we pass to the infinite limit without specifying the limiting ratio, the unconditional probability $P(A|I)$ becomes indeterminate; we can get any value in $[0, 1]$ depending on how the limit is approached. Put differently, the ratio R contains all the information relevant to the probability of A ; yet it was thrown away in passing to the limit too soon. The unconditional probability $P(A|I)$ could not have been evaluated directly on the infinite set, any more than could the conditional probabilities.

Thus, nonconglomerability on a rectangular array, far from being a phenomenon of probability theory, is only an artifact of failure to obey the rules of probability theory as developed in Chapter 2. But from studying a single example we cannot see the common feature underlying all claims of nonconglomerability.

15.4 The tumbling tetrahedra

We now examine a claim that nonconglomerability can occur even in a one-dimensional infinite set $n \rightarrow \infty$ where there does not appear to be any limiting ratio like the above M/N to be ignored. Also we now consider a problem of inference, instead of the above sampling distribution example. The scenario (Stone, 1979) appears to be equivalent to the ‘strong inconsistency’ problem (Stone, 1976). We follow the KSS notation for the time being – until we see why we must not.

A regular tetrahedron with faces labeled e^+ (positron), e^- (electron), μ^+ (muon), μ^- (antimuon), is tossed repeatedly. A record is kept of the result of each toss, except that, whenever a record contains e^+ followed immediately by e^- (or e^- by e^+ , or μ^+ by μ^- , or μ^- by μ^+), the particles annihilate each other, erasing that pair from the record. At some arbitrary point in the sequence, the player (who is ignorant of what has happened to date) calls for one more toss, and then is shown the final record $x \in X$, after which he must place bets on the truth of the proposition $A \equiv$ ‘annihilation occurred at the final toss’. What probability $P(A|x)$ should he assign?

When we try to answer this by application of probability theory, we come up immediately against the difficulty that, in the problem as stated, the solution depends on a nuisance parameter, the unspecified length n of the original sequence of tosses. This was pointed out by Hill (1980), but KSS take no note of it. In fact, they do not mention n at all except by implication, in a passing remark that the die is ‘rolled a very large number of times’. We infer that they meant the limit $n \rightarrow \infty$, from later phrases such as ‘the countable set S ’ and ‘every finite subset of S ’.

In other words, once again an infinite set is supposed to be something already accomplished, and one is trying to find relations between probabilities by reasoning directly on the infinite set. Nonconglomerability enters through asking whether the prior probability $P(A)$ is conglomerable in the partition x , corresponding to the equation

$$P(A) = \sum_{x \in X} P(A|x)P(x). \quad (15.12)$$

KSS denote by $\theta \in S$ the record just before the final toss (thought of as a ‘parameter’ not known by the player), where S is the set of all possible such records, and conclude by verbal arguments that:

- (a) $0 \leq p(A|\theta) \leq 1/4$, all $\theta \in S$;
- (b) $3/4 \leq p(A|x) \leq 1$, all $x \in X$.

It appears that another violent nonconglomerability has been produced; for if $P(A)$ is conglomerable in the partition $\{x\}$ of final records, it must be true that $3/4 \leq P(A) \leq 1$, while if it is conglomerable in the partition $\{\theta\}$ of previous records, we require $0 \leq P(A) \leq 1/4$; it cannot be conglomerable in both. So where is the error this time?

We accept statement (a); indeed, given the independence of different tosses, knowing anything whatsoever about the earlier tosses gives us no information about the final one, so

the uniform prior assignment $1/4$ for the four possible results of the final toss still holds. Therefore, $p(A|\theta) = 1/4$, except when the record θ is blank, in which case there is nothing to annihilate, and so $p(A|\theta) = 0$. But this argument does not hold for statement (b); since the result of the final toss affects the final record x , it follows that knowing x must give some information about the final toss, invalidating the uniform $1/4$ assignment.

Also, the argument that KSS gave for statement (b) supposed prior information different from that used for statement (a). This was concealed from view by the notation $p(A|\theta)$, $p(A|x)$ which fails to indicate prior information I . Let us repeat (15.12) with adequate notation:

$$P(A|I) = \sum_{x \in X} P(A|xI)P(x|I). \quad (15.13)$$

Now as I varies, all these quantities will in general vary. By ‘conglomerability’ we mean, of course, ‘conglomerability with some particular fixed prior information I .’ Recognizing this, we repeat statements (a) and (b) in a notation adequate to show this difference:

$$\begin{aligned} \text{(a)} \quad & 0 \leq p(A|\theta I_a) \leq 1/4, & \theta \in S; \\ \text{(b)} \quad & 3/4 \leq p(A|x I_b) \leq 1, & x \in X. \end{aligned}$$

From reading KSS we find that prior information I_a , in effect, assigned uniform probabilities on the set T of 4^n possible outcomes of n tosses, as is appropriate for the case of ‘independent repetitions of a random experiment’ assumed in the statement of the problem. But I_b assigned uniform probabilities on the set S of different previous records θ . This is very different; an element of S (or X) may correspond to one element of T , or to many millions of elements of T , so a probability assignment uniform on the set of tosses is very nonuniform on the set of records. Therefore it is not evident whether there is any contradiction here; they are statements about two quite different problems.

Exercise 15.1. In $n = 40$ tosses there are $4^n = 1.21 \times 10^{24}$ possible sequences of results in the set T . Show that, if those tosses give the expected number $m = 10$ of annihilations leading to a record $x \in X$ of length 20, the specific record x corresponds to about 10^{14} elements of T . On the other hand, if there are no annihilations, the resulting record x of length 40 corresponds to only one element of T .

Perhaps this makes clearer the reason for our seemingly fanatical insistence on indicating the prior information I explicitly in every formal probability symbol $P(A|BI)$. Those who fail to do this may be able to get along without disaster for a while, judging the meaning of an equation from the surrounding context rather than from the equation as written. But eventually they are sure to find themselves writing nonsense, when they start inadvertently using probabilities conditional on different prior information in the same equation, or the same argument, and their notation conceals that fact. We shall see presently a more famous

and more serious error (the marginalization paradox) caused by failure to indicate the fact that two probabilities are conditional on different prior information.

To show the crucial role that n plays in the problem, let I agree with I_a in assigning equal prior probabilities to each of the 4^n outcomes of n tosses. Then, if n is known, calculations of $p(A|nI)$, $p(x|nI)$, $p(A|nxI)$ are determinate combinatorial problems on finite sets (i.e. in each case there is one and only one correct answer), and the solutions obviously depend on n . So let us try to calculate $P(A|xI)$; denoting summation over all n in $(0 \leq n < \infty)$ by \sum , we have for the prior probabilities

$$\begin{aligned} p(A|I) &= \sum p(A|nI) = \sum p(A|nI)p(n|I) \\ p(x|I) &= \sum p(x|nI) = \sum p(x|nI)p(n|I) \end{aligned} \quad (15.14)$$

and for the conditional one

$$p(A|xI) = \sum p(A|nxI)p(n|xI) = \frac{\sum p(A|nxI)p(x|nI)p(n|I)}{\sum p(x|nI)p(n|I)}, \quad (15.15)$$

where we expanded $p(n|xI)$ by Bayes' theorem. It is evident that the problem is indeterminate until the prior probabilities $p(n|I)$ are assigned. Quite generally, failure to specify the prior information makes a problem of inference just as ill-posed as does failure to specify the data.

Passage to infinite n then corresponds to taking the limit of prior probabilities $p(n|I)$ that are nonzero only for larger and larger n . Evidently, this can be done in many different ways, and the final results will depend on which limiting process we use unless $p(A|nI)$, $p(x|nI)$, $p(A|nxI)$ all approach limits independent of n .

The number of different possible records x is less than 4^n (asymptotically, about 3^n) because many different outcomes with annihilation may produce the same final record, as the above exercise shows. Therefore, for any $n < \infty$, there is a finite set X of different possible final records x , and *a fortiori* a finite set S of previous records θ , so the prior probability of final annihilation can be written in either of the forms:

$$p(A|nI) = \sum_{x \in X} p(A|xnI)p(x|nI) = \sum_{\theta \in S} p(A|\theta nI)p(\theta|nI), \quad (15.16)$$

and the general theorem on weighted averages guarantees that nonconglomerability cannot occur in either partition for any finite n , or for an infinite set generated as a well-behaved limit of a sequence of these finite sets.

A few things about the actual range of variability of the conditional probabilities $p(A|nxI)$ can be seen at once without any calculation. For any n , there are possible records of length n for which we know that no annihilation occurred; the lower bound is always reached for some x , and it is $p(A|nxI) = 0$, not $3/4$. The lower bound in statement (b) could never have been found for any prior information, had the infinite set been approached as a limit of a sequence of finite sets. Furthermore, for any even n there are possible records of length zero for which we know that the final toss was annihilated; the upper bound is always reached for some x , and it is $p(A|nxI) = 1$.

Likewise, for even n it is not possible for θ to be blank, so from (15.16) we have $p(A|nI) = p(A|\theta nI) = 1/4$ for all $\theta \in S$. Therefore, if n is even, there is no need to invoke even the weighted average theorem; there is no possibility for nonconglomerability in either the partition $\{x\}$ or $\{\theta\}$.

At this point it is clear that the issue of nonconglomerability is disposed of in the same way as in our first example; it is an artifact of trying to calculate probabilities directly on an infinite set without considering any limit from a finite set. Then it is not surprising that KSS never found any specific answer to their original question: ‘What can we infer about final annihilation from the final record x ?’ But we would still like to see the answer (particularly since it reveals an even more startling feature of the problem).

15.5 Solution for a finite number of tosses

If n is known, we can get the exact analytical solution easily from valid application of our rules. It is a straightforward Bayesian inference in which we are asking only for the posterior probability for final annihilation A . But this enables us to simplify the problem; there is no need to draw inferences about every detail of the previous record θ .

If there is annihilation at the n th toss, then the length of the record decreases by one: $y(n) = y(n-1) - 1$. If there is no annihilation at the n th toss, the length increases by one: $y(n) = y(n-1) + 1$. The only exception is that $y(n)$ is not permitted to become negative; if $y(n-1) = 0$, then the n th toss cannot give annihilation. Therefore, since the available record x tells us the length $y(n)$ but not $y(n-1)$, any reasoning about final annihilation may be replaced immediately by reasoning about $\alpha \equiv y(n-1)$, which is the sole parameter needed in the problem.

Likewise, any permutations of the symbols $\{e^\pm, \mu^\pm\}$ in $x(n)$ which keep the same $y(n)$ will lead to just the same inferences about A . But then n and $y \equiv y(n)$ are sufficient statistics; all other details of the record x are irrelevant to the question being asked. Thus the scenario of the tetrahedrons is more complicated than it needs to be in order to define the mathematical problem (in fact, so complicated that it seems to have prevented recognition that it is a standard textbook random walk problem).

At each n th toss we have the sampling probability $1/4$ of annihilating, independently of what happened earlier (with a trivial exception if $y(n-1) = 0$). Therefore if we plot n horizontally, $y(n)$ vertically, we have the simplest random walk problem in one dimension, with a perfectly reflecting boundary on the horizontal axis $y = 0$. At each horizontal step, if $y > 0$ there is probability $3/4$ of moving up one unit, $1/4$ of moving down one unit; if $y = 0$, we can move only up. Starting with $y(0) = 0$, annihilation cannot occur on step 1, and immediately after the n th step, if there have been m annihilations, the length of the record is $y(n) = n - 2m$.

After the n th step we have a prior probability distribution for $y(n)$ to have the value i :

$$p_i^{(n)} \equiv p(i|nI), \quad 0 \leq i \leq n, \quad (15.17)$$

with the initial vector

$$p_i^{(0)} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ \vdots \end{pmatrix}, \quad (15.18)$$

and successive distributions are connected by the Markov chain relation

$$p_i^{(n)} = \sum_{j=0}^{n-1} M_{ij} p_j^{(n-1)} \quad \begin{matrix} 0 \leq i \leq n \\ 1 \leq n < \infty, \end{matrix} \quad (15.19)$$

with the transition matrix (number the rows and columns starting with zero)

$$M \equiv \begin{pmatrix} 0 & 1/4 & 0 & 0 & \dots \\ 1 & 0 & 1/4 & 0 & \dots \\ 0 & 3/4 & 0 & 1/4 & \dots \\ 0 & 0 & 3/4 & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix}. \quad (15.20)$$

The reflecting boundary at $y = 0$ is indicated by the element $M_{10} = 1$, which would be $3/4$ without the reflection.

The matrix M is, in principle, infinite-dimensional, but for the n th step only the first $n + 1$ rows and columns are needed. The vector $p^{(n)}$ is also, in principle, infinite-dimensional, but $p_i^{(n)} = 0$ when $i > n$. Then the exact solution for the prior probabilities $p_i^{(n)}$ is the first column of M^n :

$$p_i^{(n)} = M_{i,0}^n \quad (15.21)$$

(note that this is intended to represent $(M^n)_{i,0}$, not $(M_{i,0})^n$).

Now let us see how this prior is to be used in our Bayesian inference problem. Denote the data and the hypothesis being tested by

$$D \equiv y(n) = i, \quad H \equiv y(n-1) = \alpha, \quad (15.22)$$

which are the only parts of the data x and the parameter θ that are relevant to our problem. From the above their prior probabilities are

$$p(D|I) = M_{i,0}^n, \quad p(H|I) = M_{\alpha,0}^{n-1}. \quad (15.23)$$

The sampling distribution is

$$p(D|HI) = \begin{cases} 3/4 \delta(i, \alpha + 1) + 1/4 \delta(i, \alpha - 1) & \alpha > 0 \\ \delta(i, 1) & \alpha = 0. \end{cases} \quad (15.24)$$

So, Bayes' theorem gives the posterior probability for α as

$$p(H|DI) = p(H|I) \frac{p(D|HI)}{p(D|I)} = \frac{M_{\alpha,0}^{n-1}}{M_{i,0}^n} \begin{cases} 3/4 \delta(i, \alpha + 1) + 1/4 \delta(i, \alpha - 1) & \alpha > 0 \\ \delta(i, 1) & \alpha = 0. \end{cases} \quad (15.25)$$

Now, final annihilation A occurs if and only if $\alpha = i + 1$, so the exact solution for finite n is

$$p(A|DnI) = \frac{M_{i+1,0}^{n-1}}{4 M_{i,0}^n}, \quad (15.26)$$

in which $i = y(n)$ is a sufficient statistic. Another way of writing this is to note that the denominator of (15.26) is

$$4M_{i,0}^n = 4 \sum_j M_{i,j} M_{j,0}^{n-1} = 3M_{i-1,0}^{n-1} + M_{i+1,0}^{n-1}, \quad (15.27)$$

and so the posterior odds on A are

$$o(A|DnI) \equiv \frac{p(A|xnI)}{p(\bar{A}|xnI)} = \frac{1}{3} \frac{M_{i+1,0}^{n-1}}{M_{i-1,0}^{n-1}}, \quad (15.28)$$

and it would appear, from their remarks, that the exact solution to the problem that KSS had in mind is the limit of (15.26) or (15.28) as $n \rightarrow \infty$.

This solution for finite n is complicated because of the reflecting boundary. Without it, the aforementioned matrix element $M_{1,0}$ would be $3/4$ and the problem would reduce to the simplest of all random walk problems. That solution gives us a very good approximation to (15.26), which actually yields the exact solution to our problem in the limit. Let us examine this alternative formulation because its final result is very simple and the derivation is instructive about a point that is not evident from the above exact solution.

The problem where at each step there is probability p to move up one unit, $q = 1 - p$ to move down one unit, is defined by the recursion relation in which $f(i|n)$ is the probability to move a total distance i in n steps:

$$f(i|n+1) = pf(i-1|n) + qf(i+1|n). \quad (15.29)$$

With initial conditions $f(i|n=0) = \delta(i, 0)$, the standard textbook solution is the binomial for r successes in n trials; $f_0(i|n) = b(r|np)$, with $r = (n+i)/2$. In our problem we know that on the first step we necessarily move up, $y(1) = 1$, so our initial conditions are $f(i|n=1) = \delta(i, 1)$, and using the binomial recursion (15.29) after that the solution would be $f(i|n) = f_0(i-1|n-1) = b(r|n-1, p)$, with again $r = (n+i)/2$.

But with $p = 3/4$, this is not exactly the same as (15.19) because it neglects the reflecting boundary. If too many 'failures' (i.e. annihilations) occur early in the sequence, this could reduce the length of the record to zero, forcing the upward probability for the next step to be one rather than $3/4$; and (15.19) is taking all that into account. Put differently, in the solution to (15.29), when n is small, some probability drifts into the region $y < 0$; but if $p = 3/4$ the amount is almost negligibly small, and it all returns eventually to $y > 0$.

When n is very large, the solution drifts arbitrarily far away from the reflecting boundary, putting practically all the probability into the region $(\hat{y} - \sqrt{n} < y < \hat{y} + \sqrt{n})$, where $\hat{y} \equiv (p - q)n = n/2$. So conclusions drawn from (15.29) become highly accurate (in the limit, exact).

The sampling distribution (15.24) is unchanged, but we need binomial approximations to the priors for i and α . The latter is the length of the record after $n - 1$ steps, or tosses. No annihilation is possible at the first toss, so after $n - 1$ tosses we know that there were $n - 2$ tosses at which annihilation could have occurred, with probability $1/4$ at each, so the prior probability for m annihilations in the first $n - 1$ tosses is the binomial $b(m|n - 2, 1/4)$:

$$f(m) \equiv p(m|n) = \binom{n-2}{m} \left(\frac{1}{4}\right)^m \left(\frac{3}{4}\right)^{n-2-m}, \quad 0 \leq m \leq n-2. \quad (15.30)$$

Then the prior probability for α , replacing the numerator in (15.28), is

$$p(\alpha|n) = f\left(\frac{n-1-\alpha}{2}\right), \quad (15.31)$$

from which we find the prior expectation $E(\alpha|I) = n/2$. Likewise in the denominator we want the prior for $y(n) = i$. This is just (15.31) with the replacements $n - 1 \rightarrow n, \alpha \rightarrow i$.

Given y , the possible values of α are $\alpha = y \pm 1$, so the posterior odds on final annihilation are, writing $m \equiv (n - y)/2$,

$$o = \frac{p(A|yn)}{p(\bar{A}|yn)} = \frac{p(\alpha = y + 1|yn)}{p(\alpha = y - 1|yn)} = \frac{(1/4)^{\binom{n-2}{m-1}} (1/4)^{m-1} (3/4)^{n-1-m}}{(3/4)^{\binom{n-2}{m}} (1/4)^m (3/4)^{n-2-m}}. \quad (15.32)$$

But, at first sight astonishing, the factors $(1/4)$, $(3/4)$ cancel out, so the result depends only on the factorials:

$$o = \frac{m! (n - 2 - m)!}{(m - 1)! (n - 1 - m)!} = \frac{n - y}{n - 2 + y}, \quad (15.33)$$

and the posterior probability of final annihilation reduces simply to

$$p(A|yn) = \frac{o}{1 + o} = \frac{n - y}{2(n - 1)}, \quad (15.34)$$

which does not bear any resemblance to any of the solutions proposed by those who tried to solve the problem by reasoning directly on infinite sets. The sampling probabilities $p = 3/4$, $q = 1/4$, which figured so prominently in previous discussions, do not appear at all in the solution.

But now *think* about it. Given n and $y(n)$, we know that annihilation might have occurred in any of $n - 1$ tosses, but that in fact it did occur in exactly $(n - y)/2$ tosses. But we have no information about which tosses, so the posterior probability for annihilation at the final toss (or at any toss after the first) is, of course,

$$\frac{n - y}{2(n - 1)}. \quad (15.35)$$

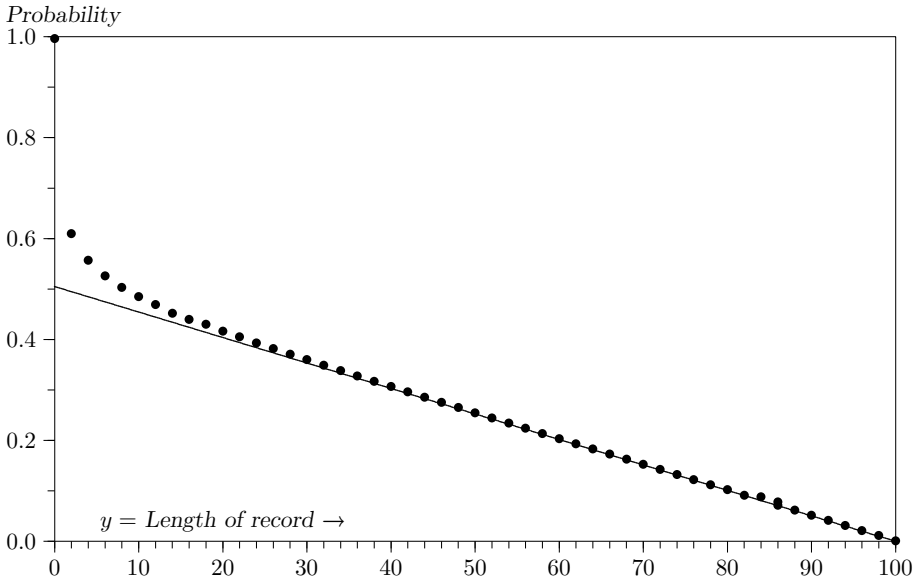


Fig. 15.1. Solution to the ‘strong inconsistency’ problem for $n = 100$ tosses. Solid line = approximation, Eq. (15.34); dots = exact solution, Eq. (15.26).

We derived (15.34) directly from the principles of probability theory by a rather long calculation; but with a modicum of intuitive understanding of the problem, we could have reasoned it out in our heads without any calculation at all!

In Figure 15.1 we compare the exact solution (15.26) with the asymptotic solution (15.34). The difference is negligible numerically when $n > 20$. So then, why did so many people think the answer should be $1/4$? Perhaps it helps to note that the prior expectation for y is $E(y|I) = (n + 1)/2$, so the *predictive probability* of final annihilation is

$$p(A|nI) = \frac{n - E(y|I)}{2(n - 1)} = \frac{1}{4}. \quad (15.36)$$

The posterior probability of final annihilation is indeed $1/4$, *if the observed record length y is the expected value*. If new information is only what we already expected, it does not change our estimates; it only makes us more confident of them. But if y is observed to be different from its prior expectation, this tells us the *actual number* of annihilations, and of course this information takes precedence over whatever initial probability assignments ($1/4$, $3/4$) we might have made. That is why they cancelled out in the posterior odds.¹ In spite of our initial surprise, then, Bayes’ theorem is doing exactly the right thing here; and the exact solution of the problem originally posed is given also by

¹ This cancellation is the thing that is not evident at all in the exact solution (15.26), although it is still taking place out of sight.

the limit of (15.35) as $n \rightarrow \infty$:

$$p(A|xI) = \frac{1}{2}(1 - z) \quad (15.37)$$

where $z \equiv \lim y(n)/n$.

In summary, the common feature of these two claims of nonconglomerability is now apparent. In the first scenario, there was no mention of the existence of the finite numbers M , N whose ratio M/N is the crucial quantity on which the solution depends. In the second scenario, essentially the same thing was done; failure to introduce the length n of the sequence and, incredibly, even the length $y(n)$ of the observed record, likewise causes one to lose the crucial thing – in this case, the sufficient statistic y/n – on which the solution depends. In both cases, by supposing the infinite limit as something already accomplished at the start, *one is throwing away the very information required to find the solution*.

This has been a very long discussion, but it is hard to imagine a more instructive lesson in how and why one must carry out probability calculations where infinite sets are involved, or a more horrible example of what can happen if we fail to heed the advice of Gauss.

15.6 Finite vs. countable additivity

At this point, the reader will be puzzled and asking, ‘Why should anybody care about nonconglomerability? What difference does it make?’ Nonconglomerability is, indeed, of little interest in itself; it is only a kind of red herring that conceals the real issue. A follower of de Finetti would say that the underlying issue is the technical one of ‘finite additivity’. To which we would reply that ‘finite additivity’ is also a red herring, because it is used for a purpose almost the opposite of what it sounds like.

In Chapter 2 we derived the sum rule (2.85) for mutually exclusive propositions: if as a statement of Boolean algebra, $A \equiv A_1 + A_2 + \cdots + A_n$ is a disjunction of a finite number of mutually exclusive propositions, then

$$p(A|C) = \sum_{i=1}^n p(A_i|C). \quad (15.38)$$

Then it is a trivial remark that our probabilities have ‘finite additivity’. As $n \rightarrow \infty$ it seems rather innocuous to suppose that the sum rule goes in the limit into a sum over a countable number of terms, forming a convergent series; whereupon our probabilities would be called countably additive. Indeed (although we do not see how it could happen in a real problem), if this should ever fail to yield a convergent series we would conclude that the infinite limit does not make sense, and we would refuse to pass to the limit at all. In our formulation of probability theory, it is difficult to see how one could make any substantive issue out of this perfectly straightforward situation.

The conventional formulations, reversing our policy, suppose the infinite limit already accomplished at the beginning, before such questions as additivity are raised; and then are

concerned with additivity over propositions about intervals on infinite sets. To quote Feller (1966, 1971 edn, p. 107):

Let F be a function assigning to each interval I a finite value $F\{I\}$. Such a function is called (finitely) *additive* if for every partition of an interval I into finitely many non-overlapping intervals $I_1 \cdots I_n$, $F\{I\} = F\{I_1\} + \cdots + F\{I_n\}$.

Then (p. 108) Feller gives an example showing why he wishes to replace finite additivity by countable additivity:

In R^1 put $F\{I\} = 0$ for any interval $I = (a, b)$ with $b < \infty$ and $F\{I\} = 1$ when $I = (a, \infty)$. This interval function is additive but weird because it violates the natural continuity requirement that $F\{(a, b)\}$ should tend to $F\{(a, \infty)\}$ as $b \rightarrow \infty$.

This last example shows the desirability of strengthening the requirement of finite additivity. We shall say that an interval function F is countably additive, or σ -additive, if for every partitioning of an interval I into countably many intervals I_1, I_n, \dots , $F\{I\} = \sum F\{I_k\}$.

He then adds that the condition of countable additivity is ‘manifestly violated’ in the above weird example (let it be an exercise for the reader to explain clearly *why* this is manifest).

What is happening in that weird example? Surely, the weirdness does not lie in lack of continuity (since continuity is quite unnecessary in any event), but in something far worse. Supposing those intervals occupied by some variable x and the interval function $F\{I\}$ to be the probability $p(x \in I)$, one is assigning zero probability to any finite range of x , but unit probability to the infinite range. This is almost impossible to comprehend when we suppose the infinite interval already accomplished, but we can understand what is happening if we heed the advice of Gauss and think in terms of passage to a limit. Suppose we have a properly normalized pdf:

$$p(x|r) = \begin{cases} 1/r & 0 \leq x < r \\ 0 & r \leq x < \infty. \end{cases} \quad (15.39)$$

As long as $0 < r < \infty$, there is nothing strange, and we could describe this by an interval function

$$F(a, b) \equiv \int_a^b dx \, p(x|r) = \begin{cases} (b-a)/r & 0 \leq a \leq b \leq r < \infty \\ (r-a)/r & 0 \leq a \leq r \leq b < \infty \\ 0 & 0 \leq r \leq a \leq b < \infty, \end{cases} \quad (15.40)$$

which is, rather trivially, countably additive and *a fortiori* finitely additive. As r increases, the density function becomes smaller and spread over a wider interval; but as long as $r < \infty$ we have a well-defined and nonparadoxical mathematical situation.

If we try to describe the limit of $p(x|r)$ as something already accomplished *before* discussing additivity, then we have created Feller’s weird example. We are trying to make a probability density that is everywhere zero, but which integrates to unity. But *there is no such thing*, according not only to all the warnings of classical mathematicians from Gauss on, but according to our own elementary common sense.

Invoking finite additivity is a sneaky way of approaching the real issue. To see why the kind of additivity matters in the conventional formulation, let us note what happens when one carries out the order of operations corresponding to our advice above. We assign a continuous monotonic increasing cumulative probability function $G(x)$ on the real line, with the natural continuity property that

$$G(x) \rightarrow \begin{cases} 1 & x \rightarrow +\infty \\ 0 & x \rightarrow -\infty; \end{cases} \quad (15.41)$$

then, the interval function F for the interval $I = (a, b)$ may be taken as $F\{I\} = G(b) - G(a)$, and it is ‘manifest’ that this interval function is countably additive in the sense defined. That is, we can choose x_k satisfying $a < x_1 < x_2 < \dots < b$ so as to break the interval (a, b) into as many nonoverlapping subintervals $\{I_0, I_1, \dots, I_n\} = \{(a, x_1), (x_1, x_2), \dots, (x_n, b)\}$ as we please, and it will be true that $F\{I\} = \sum F\{I_k\}$. If $G(x)$ is differentiable, then its derivative $f(x) \equiv G'(x)$ may be interpreted as a normalized probability density: $\int dx f(x) = 1$.

We see, finally, what the point of all this is: ‘finite additivity’ is a euphemism for ‘reversing the proper order of approaching limits, and thereby getting into trouble with non-normalizable probability distributions’. Feller saw this instantly, warned the reader against it, and proceeded to develop his own theory in a way that avoids the many useless and unnecessary paradoxes that arise from it.²

As we saw in Chapter 6, passage to the limit $r \rightarrow \infty$ at the end of a calculation can yield useful results; some other probability derived from $p(x|r)$ might approach a definite, finite, and simple limiting value. We have now seen that trying to pass to the limit at the beginning of a calculation can generate nonsense because crucial information is lost before we have a chance to use it.

The real issue here is: do we admit such things as uniform probability distributions on infinite sets into probability theory as legitimate mathematical objects? Do we believe that an infinite number of zeroes can add up to one? In the strange language in which these things are discussed, to advocate ‘finite additivity’, as de Finetti and his followers do, is a devious way of answering ‘yes’ without seeming to do so. To advocate ‘countable additivity’, as Kolmogorov and Feller did, is an equally devious way to answer ‘no’ in the spirit of Gauss.

The terms are red herrings because ‘finite additivity’ sounds colloquially as if it were a cautious assumption, ‘countable additivity’ a bit more adventurous. de Finetti does indeed seem to think that finite additivity is the weaker assumption; and he rails against those who, as he sees it, are intellectually dishonest when they invoke countable additivity only for ‘mathematical convenience’, instead of for a compelling reason. As we see it, jumping directly into an infinite set at the very beginning of a problem is a vastly greater error of judgment, which has far worse consequences for probability theory; there is a little more than just ‘mathematical convenience’ at stake here.

² Since we disagree with Feller so often on conceptual issues, we are glad to be able to agree with him on nearly all technical ones. He was, after all, a very great contributor to the technical means for solving sampling theory problems, and practically everything he did is useful to us in our wider endeavors.

We noted the same psychological phenomenon in Chapter 3, when we introduced the binomial distribution for sampling with replacement; those who committed the sin of throwing away relevant information invented the term ‘randomization’ to conceal that fact and make it sound like they were doing something respectable. Those who commit the sin of doing reckless, irresponsible things with infinity often invoke the term ‘finite additivity’ to make it sound as if they are being *more* careful than others with their mathematics.

15.7 The Borel–Kolmogorov paradox

For the most part, the transition from discrete to continuous probabilities is uneventful, proceeding in the obvious way with no surprises. However, there is one tricky point concerning continuous densities that is not at all obvious, but can lead to erroneous calculations unless we understand it. The following example continues to trap many unwary minds.

Suppose I is prior information according to which (x, y) are assigned a bivariate normal pdf with variance unity and correlation coefficient ρ :

$$p(dx dy|I) = \frac{\sqrt{1-\rho^2}}{2\pi} \exp \left\{ \frac{1}{2}(x^2 + y^2 - 2\rho xy) \right\} dx dy. \quad (15.42)$$

We can integrate out either x or y to obtain the marginal pdfs (to prepare for integrating out x , write $x^2 + y^2 - 2\rho xy = (x - \rho y)^2 + (1 - \rho^2)y^2$, etc.):

$$p(dx|I) = \sqrt{\left(\frac{1-\rho^2}{2\pi}\right)} \exp \left\{ -\frac{1}{2}(1-\rho^2)x^2 \right\} dx \quad (15.43)$$

$$p(dy|I) = \sqrt{\left(\frac{1-\rho^2}{2\pi}\right)} \exp \left\{ -\frac{1}{2}(1-\rho^2)y^2 \right\} dy. \quad (15.44)$$

Thus far, all is routine. But now, what is the conditional pdf for x , given that $y = y_0$? We might think that we need only set $y = y_0$ in (15.42) and renormalize:

$$p(dx|y = y_0 I) = A \exp \left\{ -\frac{1}{2}(x^2 + y_0^2 - 2\rho xy_0) \right\} dx, \quad (15.45)$$

where A is a normalizing constant. But there is no guarantee that this is valid, because we have obtained (15.45) by an intuitive *ad hoc* device; we did not derive it from (15.42) by applying the basic rules of probability theory, which we derived in Chapter 2 for the discrete case:

$$p(AB|X) = p(A|BX)p(B|X), \quad (15.46)$$

from which a discrete conditional probability is given by the usual rule

$$p(A|BX) = \frac{p(AB|X)}{p(B|X)} \quad (15.47)$$

often taken as the definition of a conditional probability. But we can do the calculation by strict application of our rules if we define the discrete propositions

$$\begin{aligned} A &\equiv x \text{ in } dx \\ B &\equiv y \text{ in } (y_0 < y < y_0 + dy). \end{aligned} \quad (15.48)$$

Then we should write instead of (15.45), using (15.42) and (15.44),

$$p(A|BI) = p(dx|dy I) = \frac{p(dx dy|I)}{p(dy|I)} = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2}(x - \rho y_0)^2 \right\} dx. \quad (15.49)$$

Since dy cancels out, taking the limit $dy \rightarrow 0$ does nothing.

On working out the normalizing constant in (15.45), we find that (15.45) and (15.49) are in fact identical. So, why all this agony? Didn't the quick argument leading to (15.45) give us the right answer?

This is a good example of our opening remarks that a fallacious argument may lead to correct or incorrect results. The reasoning that led us to (15.45) happened to give a correct result here; but it can equally well yield any result we please instead of (15.45). It depends on the particular form in which you or I choose to write our equations. To show this, and therefore generate a paradox, suppose that we had used instead of (x, y) the variables (x, u) , where

$$u \equiv \frac{y}{f(x)} \quad (15.50)$$

with $0 < f(x) < \infty$; for example, $f(x) = 1 + x^2$ or $f(x) = \cosh(x)$, etc. The Jacobian is

$$\frac{\partial(x, u)}{\partial(x, y)} = \left(\frac{\partial u}{\partial y} \right)_x = \frac{1}{f(x)} \quad (15.51)$$

so the pdf (15.42), expressed in the new variables, is

$$p(dx du|I) = \frac{\sqrt{1 - \rho^2}}{2\pi} \exp \left\{ -\frac{1}{2}(x^2 + u^2 f^2(x) - 2\rho u f(x)) \right\} f(x) dx du. \quad (15.52)$$

Again, we can integrate out u or x , leading to a marginal distribution $p(dx|I)$, which is easily seen to be identical with (15.43), and $p(du|I)$, which is found to be identical with (15.44) transformed to the variable u , as it should be; so far, so good.

But now, what is the conditional pdf for x , given that $u = 0$? If we follow the reasoning that led us to (15.45); i.e. simply set $u = 0$ in (15.52) and renormalize, we find

$$p(dx|u = 0 I) = A \exp \left\{ -\frac{1}{2}x^2 \right\} f(x) dx. \quad (15.53)$$

Now from (15.50) the condition $u = 0$ is the same as $y = 0$; so it appears that this should be the same as (15.45) with $y_0 = 0$. But (15.53) differs from that by an extra factor $f(x)$, which could be arbitrary!

Many find this astonishing and unbelievable; they repeat over and over: 'But the condition $u = 0$ is *exactly the same condition* as $y = 0$; how can there be a different result?' We warned against this phenomenon briefly, and perhaps too cryptically, in Chapter 4; but there it did

not actually cause an error because we had only one parameter in the problem. Now we need to examine it carefully to see the error and the solution.

We noted in Chapter 1 that we shall make no attempt to define any probability conditional on contradictory premises; there could be no unique solution to such a problem. We start each problem by defining a ‘sample space’ or ‘hypothesis space’ which sets forth the range of conditions we shall consider *in that problem*. In the present problem, our discrete hypotheses were of the form ‘ $a \leq y \leq b$ ’, placing y in an interval of positive measure $b - a$. Then what could we mean by the proposition ‘ $y = 0$ ’, which has measure zero? We could mean only the limit of some sequence of propositions referring to positive measure, such as

$$A_\epsilon \equiv |y| < \epsilon \quad (15.54)$$

as $\epsilon \rightarrow 0$. The propositions A_ϵ confine the point (x, y) to successively narrower horizontal strips, but for any $\epsilon > 0$, A_ϵ is a discrete proposition with a definite positive probability, so by the product rule the conditional probability of any hypothesis $H \equiv 'x \text{ in } dx'$,

$$p(H|A_\epsilon I) = \frac{p(HA_\epsilon|I)}{p(A_\epsilon|I)} \quad (15.55)$$

is well-defined, and the limit of this as $\epsilon \rightarrow 0$ is also a well-defined quantity. Perhaps that limit is what one meant by $p(H|y = 0 I)$.³

But the proposition ‘ $y = 0$ ’ may be defined equally well as the limit of the sequence

$$B_\epsilon \equiv |y| < \epsilon|x| \quad (15.56)$$

of successively thinner wedges, and $p(H|B_\epsilon I)$ is also unambiguously defined as in (15.55) for all $\epsilon > 0$. Although the sequences $\{A_\epsilon\}$, $\{B_\epsilon\}$ tend to the same limit $y = 0$, the conditional densities tend to different limits:

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} p(H|A_\epsilon) &\propto g(x), \\ \lim_{\epsilon \rightarrow 0} p(H|B_\epsilon) &\propto |x|g(x), \end{aligned} \quad (15.57)$$

and in place of $|x|$ we could put an arbitrary non-negative function $f(x)$. As we see from this, merely to specify ‘ $y = 0$ ’ without any qualifications is ambiguous; it tells us to pass to a measure-zero limit, but does not tell us which of any number of limits is intended.

We have here one more example showing why the rules of inference derived in Chapter 2 must be obeyed *strictly, in every detail*. Intuitive shortcuts have a potential for disaster, which is particularly dangerous just because of the fact that it strikes only intermittently. An intuitive *ad hockery* that violates those rules will probably lead to a correct result in some cases; but it will surely lead to disaster in others. Whenever we have a probability density on one space, and we wish to generate from it one on a subspace of measure zero, the only safe procedure is to pass to an explicitly defined limit by a process like (15.55). In general, the final result will and must depend on which limiting operation was specified.

³ Note again what we belabor constantly: the rules of probability theory tell us unambiguously that it is the limit of the ratio, not the ratio of the limits, that is to be taken in (15.55). The former quantity remains finite and well-behaved in conditions where the latter does not exist.

This is extremely counter-intuitive at first hearing; yet it becomes obvious when the reason for it is understood.

A famous puzzle based on this paradox concerns passing from the surface of a sphere to a great circle on it. Given a uniform probability density over the surface area, what is the corresponding conditional density on any great circle? Intuitively, everyone says immediately that, from geometrical symmetry, it must be uniform also. But if we specify points by latitude ($-\pi/2 \leq \theta \leq \pi/2$) and longitude ($-\pi < \phi \leq \pi$), we do not seem to get this result. If that great circle is the equator, defined by $|\theta| < \epsilon$ as $\epsilon \rightarrow 0$, we have the expected uniform distribution $p(\phi) = (2\pi)^{-1}$ ($-\pi < \phi \leq \pi$). But if it is the meridian of Greenwich defined by $|\phi| < \epsilon$ as $\epsilon \rightarrow 0$, we have $p(\theta) = (1/2)\cos(\theta)$ ($-\pi/2 \leq \theta \leq \pi/2$), with the density reaching a maximum on the equator and zero at the poles.

Many quite futile arguments have raged – between otherwise competent probabilists – over which of these results is ‘correct’. The writer has witnessed this more than once at professional meetings of scientists and statisticians. 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. The intuitive symmetry argument presupposes unconsciously the equatorial limit; yet one eating slices of an orange might presuppose the other.

15.8 The marginalization paradox

The tumbling tetrahedrons problem flared up into an even more spectacular case of probability theory gone crazy, with the work of Dawid, Stone, and Zidek (1973), hereafter denoted by DSZ, which for a time seemed to threaten the consistency of all probability theory. The marginalization paradox is more complicated than the ones discussed above, because it arises not from a single error, but from a combination of errors of logic and intuition, insidious because they happened to support each other. When first propounded it seems to have fooled every expert in the field, with the single exception of D. A. S. Fraser, who, as discussant of the DSZ paper, saw that the conclusions were erroneous and put his finger correctly on the cause of this; but he was not listened to.

The marginalization paradox also differs from the others in that it received the immediate, enthusiastic endorsement of the establishment, and therefore it has been able to do far more damage to the cause of scientific inference than any, other; yet, when properly understood, the phenomenon has useful applications in scientific inference. Marginalization as a potentially useful means of constructing uninformative priors is discussed incompletely in Jaynes (1980); this rather deep subject still has the status of ongoing research, in which the main theorems are probably not yet known.

In the present chapter we are concerned with the marginalization story only as a weird episode of history which accomplished one good thing by forcing Bayesians to revise some easy, shortcut inference procedures. We illustrate the original paradox by the scenario of DSZ, again following their notation until we see why we must not. It starts as a conventional, and seemingly harmless, nuisance parameter problem.

A conscientious Bayesian B_1 studies a problem with data $x \equiv (x_1, \dots, x_n)$ and a multidimensional parameter θ which he partitions into two components, $\theta = (\eta, \zeta)$, being interested only in inferences about ζ . Thus his model is defined by some specified sampling distribution $p(x|\eta\zeta)$ supposed given in the statement of the problem, and η is a nuisance parameter to be integrated out. With a prior $\pi(\eta, \zeta)$, B_1 thus obtains the marginal posterior pdf for ζ :

$$p(\zeta|x) = \int d\eta p(\eta\zeta|x) = \frac{\int d\eta p(x|\eta\zeta)\pi(\eta, \zeta)}{\int d\zeta \int d\eta p(x|\eta\zeta)\pi(\eta, \zeta)}, \quad (15.58)$$

the standard result, which summarizes everything B_1 knows about ζ . The issue now turns on what class of priors $\pi(\eta, \zeta)$ we may assign for this purpose. Our answer is, of course:

Any proper prior, or any limit of a sequence of such priors, such that the ratio of integrals in (15.58) converges to yield a proper posterior pdf for ζ , may be admitted into our theory as representing a conceivable state of prior knowledge about the parameters. Eq. (15.58) will then yield the correct conclusions that follow from that state of knowledge.

This need not be qualified by any special circumstances of the particular problem; we believe that this policy, followed strictly, cannot generate ambiguities or contradictions. But *failure* to follow it can lead to almost anything.

However, DSZ did not see it that way at all. They concentrate on a special circumstance, noting that in many cases the data x may be partitioned into two components: $x = (y, z)$ in such a way that ‘the sampling distribution for z is independent of the nuisance parameter η ’, which property they write (DSZ, Eq. (1.2)) as

$$p(z|\eta\zeta) = \int dy p(yz|\eta\zeta) = p(z|\zeta), \quad (15.59)$$

which, by itself, would appear rather generally possible, but without any very deep significance. For example, if η is a location parameter, then any function $z(x)$ of the data that is invariant under rigid translations will have a sampling distribution independent of η . If η is a scale parameter, then any function $z(x)$ invariant under scale changes will have this property. If η is a rotation angle, then any component of the data that is invariant under those rotations will qualify. DSZ proceed to discover cases in which, when (15.59) holds and B_1 assigns an improper prior to η , he finds that his marginal posterior pdf for ζ ‘is a function of z only’, which, in view of (15.59), DSZ would presumably write as

$$p(\zeta|yz) = p(\zeta|z). \quad (15.60)$$

At this point there enters a lazy Bayesian B_2 , who ‘always arrives late on the scene of inference’, and the combination of (15.59) and (15.60) sets off for him a curious train of thought. From (15.60) as written it appears that the component y of the data can be discarded as irrelevant to inferences about ζ . The appearance of (15.59) then suggests that η might also be removed from the model as irrelevant. So he proposes to simplify the calculation; his intuitive judgment is that, given (15.59) and (15.60), we should be able to derive the marginal pdf for ζ more easily by direct application of Bayes’ theorem in a reduced model

$p(z|\zeta)$ in which (y, η) do not appear at all. Thus if B_2 assigns the prior $\pi(\zeta)$, he obtains the posterior distribution

$$p(\zeta|z) = \frac{p(z|\zeta)\pi(\zeta)}{\int d\zeta p(z|\zeta)\pi(\zeta)}. \quad (15.61)$$

But he finds to his dismay that he cannot reproduce B_1 's result (15.58) whatever prior he assigns to ζ . What conclusions should we draw from this?

For DSZ, the reasoning of B_2 seemed compelling; on grounds of this intuitive 'reduction principle' they considered it obvious that B_1 and B_2 ought to get the same results, and therefore that one of them must be guilty of some transgression. They point the accusing finger at B_1 thus: ' B_2 's intervention has revealed the paradoxical unBayesianity of B_1 's posterior distribution for ζ .' They place the blame on his use of an improper prior for η .

For us, the situation appears very different; B_2 's result was not derived by application of our rules. Eq. (15.61) was only an intuitive guess; as the reader may verify, it does not follow mathematically from (15.58), (15.59) and (15.60). Therefore, (15.61) is *not a valid application of probability theory to B_1 's problem*. If intuition suggests otherwise, then that intuition needs educating – just as it did in the other paradoxes.

At this stage we are faced not just with one confusion, but with three. The notation used above conceals from view some crucial points:

- (1) While the result (15.60) is 'a function of z only' in the sense that y does not appear explicitly in (15.60), it is a *different* function of z for different η -priors. That is, it is still a functional of the η -prior, as is clear from a glance at (15.58); through this dependence, probability theory is telling us that prior information about η still matters. As soon as we realize this, we see that B_2 comes to a different conclusion than B_1 not because B_1 is committing a transgression, but for just the opposite reason: B_1 is taking into account relevant prior information that B_2 is ignoring.
- (2) But the real trouble starts farther back than that. We need to be aware that current orthodox notation has a more basic ambiguity that makes the meaning of (15.59) and (15.60) undefined, and this is corrected only by the notation introduced by Harold Jeffreys (1939) and expounded in our Chapter 2 and Appendix B. Thus, we understand that the symbol $p(yz|\eta\zeta)$ stands for the joint probability (density) for y, z conditional on *specific numerical values* for the two parameters η, ζ that are present in our model. But then what does $p(z|\zeta)$ stand for? Presumably this is not intended to say that η has no numerical value at all!

Indeed, if he wished to refer to a *different model* in which η is not present at all, the orthodoxian would use the same notation $p(z|\zeta)$. So it seems that, strictly speaking, we should always interpret the symbol $p(z|\zeta)$ as referring to that different model. But that is not the intention in (15.59); reference is being made to a model in which η is still present, but the probability for z is independent of its numerical value. It seems that the only way this could be expressed in orthodox notation is to rewrite (15.59) as

$$\frac{\partial}{\partial \eta} p(z|\eta\zeta) = 0. \quad (15.62)$$

- (3) This ambiguity, and still another one, is present in (15.60); here the intention is only to indicate that $p(\zeta|yz)$ is independent of the numerical value of y ; but the symbol $p(\zeta|z)$, strictly speaking, must be held to refer to a different model in which the datum y was not given at all. Now we

have the additional ambiguity that any posterior probability depends necessarily on the prior

information; yet the notation in (15.60) makes no reference to any prior information.⁴ We begin to see why the marginalization paradox was so confusing!

There is a better way of looking at this, which avoids all the above confusions while using the mathematics that was intended by DSZ; we may take a more charitable view of B_2 if we put these equations in a different scenario. The lazy Bayesian, B_2 , was introduced as a fellow who invents a shortcut method that violates the rules of probability theory. But we may suppose equally well that, through no fault of his own, he is only an uninformed fellow who was given only the reduced model $p(z|\zeta)$ in which η is not present; and he is unaware of the existence of (η, y) . Then (15.61) is a valid inference for the *different state of knowledge* that B_2 has; and it is valid whether or not the separation property (15.60) holds.⁵ Although the equations are the same because we defined B_2 's model by B_1 's marginal sampling distribution $p(z|\zeta)$, this avoids much confusion; viewed in this way, B_1 and B_2 are both making valid inferences, but about two different problems.

Both of these new ambiguities arise from the fact that orthodox notation fails to indicate which model is being considered. But both are corrected by including the prior information symbol I , understood to be a proposition defined somewhere in the surrounding context, that includes full specification of the model. If we follow the example of Jeffreys and write the right-hand sides of (15.58) and (15.61) correctly as $p(\zeta|yzI_1)$ and $p(\zeta|zI_2)$, thereby making this difference in the problems clear, there can be no appearance of paradox. The prior information I_1 specifies the full sampling distribution $p(yz|\eta\zeta)$, while I_2 specifies a model only by $p(z|\zeta)$, which makes no reference to (η, y) . That B_1 and B_2 came to different conclusions from different prior information is no more strange than if they had come to different conclusions from different data.

Exercise 15.2. Consider the intermediate case of a third Bayesian, B_3 , who has the same prior information as B_1 about η, ζ but is not given the data component y . Then y never appears in B_3 's equations at all; his model is the marginal sampling distribution $p(z|\eta\zeta I_3)$. Show that, nevertheless, if (15.59) still holds (in the interpretation intended, as indicated by (15.62)), then B_2 and B_3 are always in agreement, $p(\zeta|zI_3) = p(\zeta|zI_2)$, and that to prove this it is not necessary to appeal to (15.60). Merely withholding the datum y automatically makes any prior knowledge about η irrelevant to inference about ζ . Ponder this until you can explain in words why it is, after all, intuitively obvious.

⁴ Yet, as we stress again, if you fail to specify the prior information, a problem of inference is just as ill-posed as if you had failed to specify the data. In practice, orthodoxy is able to function in spite of this in some problems, by the tacit assumption that an uninformative prior is to be used. Of course, the dedicated orthodoxian will deny vehemently that he is making any such assumption; nevertheless, it is a mathematical fact that his conclusions are what a Bayesian would obtain *from an uninformative prior*. This was demonstrated already by Jeffreys (1939).

⁵ The fact that (15.60) is not essential to the problem was not yet clearly seen in Jaynes (1980); the marginalization problem was more subtle than that any Bayesians had faced up to that time. Because DSZ laid so much stress on (15.60), we followed them in concentrating on finding conditions for its validity. Today, with the benefit of hindsight, it is clear that there is in general no reason to expect (15.60) to hold, so it loses its supposed importance. This deeper understanding enables us to find useful solutions to current problems of inference far more subtle than marginalization, as demonstrated by Bretthorst (1988). But the secret of success here is, as always, simply: *absolutely strict adherence* to the rules of conduct derived in Chapter 2. As these paradoxes show, the slightest departure from them can generate gross absurdities.

15.8.1 On to greater disasters

Up to this point, we had only a misreading of equations through inadequate notation; but now a comedy of mutually reinforcing errors commenced. In support of their contention that B_1 is the guilty party, DSZ offered a proof that this paradox (i.e. the discrepancy in the results of B_1 and B_2) ‘could not have arisen if B_1 had employed proper prior distributions’. Let us examine their proof of this, still using their notation. With a general joint proper prior $\pi(\eta, \zeta)$ the integrals in (15.58) are separately convergent and positive, so if we multiply through by the denominator, we are neither multiplying nor dividing by zero. Then

$$p(x|\eta\zeta) = p(yz|\eta\zeta) = p(y|z\eta\zeta)p(z|\eta\zeta) = p(y|z\eta\zeta)p(z|\zeta), \quad (15.63)$$

where we used the product rule and (15.59). Then (15.58) becomes

$$p(\zeta|yz) \int d\zeta \int d\eta p(y|z\eta\zeta)p(z|\zeta)\pi(\eta, \zeta) = \int d\eta p(y|z\eta\zeta)p(z|\zeta)\pi(\eta, \zeta). \quad (15.64)$$

But now we assume that (15.60) still holds; because the integrals are absolutely convergent, we may integrate out y from both sides of (15.64), whereupon $\int d\eta \pi(\eta, \zeta) = \pi(\zeta)$ and (15.64) reduces to

$$p(\zeta|z) \int d\zeta p(z|\zeta)\pi(\zeta) = p(z|\zeta)\pi(\zeta), \quad (15.65)$$

which is identical with (15.61). DSZ concluded that, if B_1 uses a proper prior, then B_1 and B_2 are necessarily in agreement – from which it would follow again, in agreement with their intuition, that the paradox must be caused by B_1 ’s use of improper priors.

But this proof of (15.65) has used mutually contradictory assumptions. As Fraser recognized, if B_1 uses a proper prior, then (15.60) *cannot* be true and (15.65) does not follow; it is no accident that DSZ had found (15.60) only with improper priors. This is easiest to see in terms of a specific example, after which it will become obvious why it is true in general. In the following we use the full notation of Jeffreys so that we always distinguish between the two problems.

The change-point problem

Observations have been made of n successive, independent, positive real, ‘exponentially distributed’ quantities $\{x_1, \dots, x_n\}$. It is known (definition of the model) that the first ζ of these have expectations $1/\eta$ and the remaining $(n - \zeta)$ have expectations $1/(c\eta)$, where c is known and $c \neq 1$, while η and ζ are unknown. From the data, we want to estimate at what point in the sequence the change occurred. The sampling density for $x \equiv (x_1, \dots, x_n)$ is

$$p(x|\eta\zeta I_1) = c^{n-\zeta} \eta^n \exp \left\{ -\eta \left(\sum_{i=1}^{\zeta} x_i + c \sum_{i=\zeta+1}^n x_i \right) \right\}, \quad 1 \leq \zeta \leq n. \quad (15.66)$$

If $\zeta = n$, then there is no change, the last sum in (15.66) is absent, and c disappears from the model. Since η is a scale parameter, the sampling distribution for ratios of observations

$z_i \equiv x_i/x_1$ should be independent of η . Indeed, separating the data $x = (y, z)$ into $y \equiv x_1$, which sets the scale and the ratios (z_2, \dots, z_n) , and noting that the volume element transforms as $dx_1 \cdots dx_n = y^{n-1} dy dz_2 \cdots dz_n$, we find that the joint sampling distribution for $z \equiv (z_2, \dots, z_n)$ depends only on ζ :

$$p(z_2 \cdots z_n | \eta \zeta I_1) = \int_0^\infty dy c^{n-\zeta} \eta^n y^{n-1} \exp\{\eta y Q(\zeta, z)\} = \frac{c^{n-\zeta} (n-1)!}{Q(\zeta, z)^n} = p(z | \zeta I_1), \quad (15.67)$$

where $z_1 \equiv 1$ and

$$Q(\zeta, z) \equiv \sum_1^\zeta z_i + c \sum_{\zeta+1}^n z_i \quad (15.68)$$

is a function that is known from the data. Let B_1 choose a properly normalized discrete prior $\pi(\zeta)$ in $(1 \leq \zeta \leq n)$, and independently a prior $\pi(\eta) d\eta$ in $(0 < \eta < \infty)$. Then B_1 's marginal posterior distribution for ζ is, from (15.66),

$$p(\zeta | y z I_1) \propto \pi(\zeta) c^{n-\zeta} \int_0^\infty d\eta \exp\{-\eta y Q\} \pi(\eta) \eta^n, \quad (15.69)$$

and, from (15.67), B_2 's posterior distribution (15.61) for ζ is now

$$p(\zeta | z I_2) \propto \pi(\zeta) p(z | \zeta) = \frac{\pi(\zeta) c^{-\zeta}}{[Q(\zeta, z)]^n}, \quad (15.70)$$

which takes no note of $\pi(\eta)$. But, as expected from the above discussion, not only does B_1 's knowledge about ζ depend on both y and z , it depends just as strongly on what prior $\pi(\eta)$ he assigned to the nuisance parameter.

On meditation, we see that a little common sense would have anticipated this result at once. If we know absolutely nothing about η except that it is positive, then the only evidence we can have about the change point ζ must come from noting the relative values of the x_i ; for example, at which i does the ratio x_i/x_1 appear to change? On the other hand, suppose that we knew η exactly; then clearly not only the ratios x_i/x_1 , but also the absolute values of the x_i , would be relevant to inference about ζ . Then, whether x_i is closer to $1/\eta$ or to $1/(c\eta)$ tells us something about whether $(i < \zeta)$ or $(i > \zeta)$ that the ratio x_i/x_1 does not tell us, this extra information would enable us to make better estimates of ζ . If we had only partial prior knowledge of η , then knowledge of the absolute values of the x_i would be less helpful, but still relevant, so, as Fraser noted, (15.60) could not be valid.

But now B_1 discovers that use of the improper prior

$$\pi(\eta) = \eta^{-k}, \quad 0 < \eta < \infty, \quad (15.71)$$

where k is any real number for which the integral (15.69) converges, leads to the separation property (15.60), and to the posterior pdf

$$p(\zeta | z I_1) \propto \frac{\pi(\zeta) c^{-\zeta}}{[Q(\zeta, z)]^{n-k+1}}, \quad (15.72)$$

which still depends, through k , on the prior assigned to η . We see that for no prior $\pi(\zeta)$ can B_2 agree with B_1 , except when $k = 1$, in which case B_2 and B_1 find themselves in agreement after all, and with the same prior $\pi(\zeta)$. But this result is not peculiar to the change-point model; it holds quite generally, as the following Exercise shows.

Exercise 15.3. Prove that the $k = 1$ prior is always uninformative in this sense whenever η is a scale parameter for y . That is, if the sampling distribution has the functional form

$$p(yz|\eta\zeta) = \eta^{-1} h(z, \zeta; y/\eta), \quad (15.73)$$

then (15.59) follows at once, and B_1 and B_2 agree if and only if we use a prior $\pi(\eta) \propto \eta^{-1}$.

It seems to us that this is an eminently satisfactory result without any trace of paradox. For the case $k = 1$ is just the Jeffreys prior, which we have already seen to be ‘completely uninformative’ about any scale parameter η , by several different criteria. Then, of course, with this prior B_1 has no extra information after all, and should, indeed, find himself in agreement with B_2 .

DSZ did not see it that way at all, and persisted in their intuitive judgment that there is a serious paradox and that B_1 was at fault for using an improper prior; so the story continues. DSZ proceed to exhibit many more examples in which this ‘paradox’ appears – invariably when an improper prior was used. The totality of all these demonstrations appeared to mount up into overwhelming evidence that to use any improper prior is to generate inconsistencies. But, in the belief that their proof of (15.65) had already dealt with it, they failed to examine what happens in those examples in the case of proper priors, and so they managed to get through a long string of examples without discovering the error in that proof.⁶

To correct this omission, and reveal the error in (15.65) clearly, we need only to examine any of the DSZ examples, to see what happens in the case of proper priors $\pi(\eta)$. In the change-point problem, whatever this prior, B_1 ’s result (15.69) depends on y and z through a function of the product $yQ(\zeta, z)$. Then for what functions $f(yQ)$ will the separation property (15.60) hold? Evidently, the necessary and sufficient condition for this is that y and ζ appear in separate factors: in the case where the integrals in (15.58) converge, we

⁶ Another reason for this was their tendency to write the priors in terms of the ‘wrong’ parameters. Usually, a model was defined initially with certain parameters α, β . The parameters η, ζ for which the relations (15.59), (15.60) held were certain functions of them: $\eta = \eta(\alpha, \beta)$, etc. But DSZ continued to write the priors in terms of α, β , which made it seem that the Jeffreys prior has no particular significance; a wide variety of different priors appeared to ‘avoid the paradox’ in various different problems. In Jaynes (1980) we showed that, had they transformed their parameters to the relevant ones η, ζ , they would have found in every such case except one that η was a scale parameter for y and the ‘paradox’ disappeared for and only for the Jeffreys prior $\pi(\eta)$. Thus Exercise 15.3 includes, in effect, all their examples except the infamous Example #5, which requires a separate treatment given below.

require the integral to have the functional form

$$\int_0^\infty d\eta \exp\{-\eta y Q\} \pi(\eta) \eta^n = f(yQ) = g(y, z)h(\zeta, z), \quad (15.74)$$

for then and only then will y cancel out upon normalization of $p(\zeta|yz)$. The answer is obvious: if a function of $[\log(y) + \log Q(\zeta)]$ has the form $[\log g(y) + \log h(\zeta)]$, the only possibility is a linear function: $\log f(yQ) = a[\log(y) + \log(Q)]$ or $f(yQ) = (yQ)^a$, where $a(z, n)$ may depend on z and n . But then, noting that the Laplace transform is uniquely invertible, and that

$$\int_0^\infty d\eta \exp\{-\eta y Q\} \eta^{a-1} = \frac{(a-1)!}{(yQ)^a}, \quad (15.75)$$

we see that, contrary to the assumption of DSZ, (15.60) *cannot hold unless the prior is of the improper form* $\pi(\eta) = \eta^{-k}$, $0 < \eta < \infty$.

Exercise 15.4. Show that this result is also general; that is, not only in the change-point problem, but in any problem like that of Exercise 15.3 where η is a scale parameter for y , a prior of the form $\pi(\eta) = \eta^{-k}$ will lead to a factorization of the form $\int d\eta p(yz|\eta\zeta)\pi(\eta) = g(y, z)h(\zeta, z)$ for some functions g, h , whereupon (15.60) will hold. For this reason, the many later examples of DSZ are essentially repetitious; they are only making the same point over and over again.

Evidently, any value of k which makes the integral (15.74) converge will lead to a well-behaved posterior distribution for ζ ; but a still wider class of values of k may do so if the improper prior is approached, as it should be, as the limit of a sequence of proper priors, as explained previously.

But use of a proper prior $\pi(\eta)$ necessarily means that the separation property (15.60) cannot hold. For example, choose the prior $\pi(\eta) \propto \eta^a \exp\{-b\eta\}$. Then (15.69) becomes

$$p(\zeta|yzI_1) \propto \frac{\pi(\zeta)c^{-\zeta}}{(b+yQ)^{n+a+1}}, \quad (15.76)$$

and as long as the prior is proper (that is, $b > 0$), the datum y cannot be disentangled, but remains relevant; and so (15.60) does not hold, as we expected from (15.75). The ‘paradox’ disappears, not because B_1 and B_2 agree, but because B_2 cannot invoke his ‘reduction principle’ at all. Indeed, in any of the DSZ examples, inserting any proper prior $\pi(\eta)$ for which we can do the integrals will yield an equally good counter-example to (15.65); how could this have gone undetected for years? We note some of the circumstances that led to this.

15.9 Discussion

Some have denied that there is any such thing as ‘complete ignorance’, much less any ‘completely uninformative’ prior. From their introductory remarks, it appears that to demonstrate this was the original goal of DSZ, and several discussants continued to emphasize the point in agreement with them. But their arguments were verbal, expressing only intuitive feelings; the mathematical facts confirm the sense of the idea of ‘complete ignorance’ after all. The Jeffreys prior is doing here what we should naturally suppose an uninformative prior ought to do, and it does this quite generally (whenever η is a scale parameter).

Technically, the concurrence of many different results like that of Exercise 15.3 shows us that the notion of complete ignorance is consistent and useful; the fact that the same Jeffreys prior emerges uniquely from many different and independent lines of reasoning shows how impossible it would be to modify it or abandon it. As is invariably the case in this field, past difficulties with the ideas of Jeffreys signified not any defects in his ideas, but only misapplications of probability theory by his critics.

Exercise 15.3 shows another sense in which our previous conclusion (that the prior $d\eta/\eta$ is uninformative about a scale parameter η) is quite literally true; not as an intuitive judgment, but now as a definite theorem that follows from the rules of probability theory. Of course, our ultimate goal is always to represent honestly the prior information that we actually have. But, both conceptually and mathematically, the notion of ‘complete ignorance’ is a valid and necessary part of this program, as the starting point from which all inference proceeds; just as the notion of zero is a necessary part of arithmetic.

In the discussion following the DSZ paper, nobody noticed that there was a counter-example to their proof of (15.65) already in plain sight in the DSZ article (their Example #5, where it is evident by inspection that B_1 and B_2 remain in disagreement for all priors, proper or improper), and only Fraser expressed any doubts about the DSZ conclusions. He noted that DSZ

... propose that the confusion can be avoided by a restriction to *proper* priors. This is a strange proposal as a resolution of the difficulties – for it means in the interesting cases that one cannot eliminate a variable, and hence cannot go to the marginal likelihood.

But it seems that these words were, like the prophecies of Nostradamus, too cryptic for anyone to understand until he had first located the error for himself. Fraser’s point – and ours above – is that when B_1 uses a proper prior, then in general B_2 ’s ‘reduction principle’ cannot be applied because (15.60) ceases to be true. In other words, when B_1 uses proper priors, this does not bring B_1 and B_2 into agreement. In (15.74) and (15.75) we have demonstrated that in the change-point problem, agreement of B_1 and B_2 *requires* that B_1 uses an improper prior; just the opposite of the DSZ conclusion.

It is evident, to one who has understood the above analysis, that the situation found in the change-point problem is actually quite general. For, if one knew both y and η , that information must be relevant to the inference about ζ unless the sampling distributions are completely independent; that is, unless $p(yz|\eta\zeta) = p(y|\eta)p(z|\zeta)$. Except in this trivial

case, if one knows y , any partial information about η must still be relevant for inference about ζ or, similarly, if one knew η , any partial information about y would be relevant.

But common sense should have told us that any proper prior $\pi(\eta)$ on an infinite domain is necessarily informative about η , for it determines finite upper and lower bounds within which η is almost certain to lie. Seen in this way, Fraser's cryptic remark becomes obvious – and in full generality.

In any event, what happened was that nearly everybody accepted the DSZ conclusions uncritically, without careful examination of their argument. Anti-Bayesians, who very much wanted the DSZ conclusion to be true, seized upon it eagerly as sounding the death-knell of all Bayesianity. Under this pressure the prominent Bayesian D. V. Lindley broke down and confessed to sins of which he was not guilty, and the Royal Statistical Society bestowed a warm vote of thanks upon DSZ for this major contribution to our understanding of inference.

As a result, since 1973 a flood of articles has appeared, rejecting the use of improper priors under any and all circumstances, on the grounds that they have been proved by DSZ to generate inconsistencies. Incredibly, the fact that proper priors never 'correct' the supposed inconsistencies never came out in all this discussion. Thus the marginalization paradox became, like nonconglomerability, quite literally institutionalized in the literature of this field, and taught as truth. Scientific inference thus suffered a setback from which it will require decades to recover.

Nobody noted that this same 'paradox' had been found and interpreted correctly long before by Harold Jeffreys (1939, Sect. 3.8) in connection with estimating the correlation coefficient ρ in a bivariate normal distribution, in which the location parameters are the uninteresting nuisance parameters. He gives two examples of B_1 's result, corresponding to different prior information about the nuisance parameters, in his equations (10) and (24), their difference indicating the effect of that prior information. Then he gives B_2 's result in (28), the agreement with (24) indicating that a uniform prior for the location parameters is uninformative about ρ .

This was seen again independently by Geisser and Cornfield (1963) in connection with priors for multivariate normal distributions. They perceived that the difference between the results of B_1 and B_2 , their equations (3.10) and (3.26), was not a paradox, because B_2 's result was not a valid solution to the problem; they termed it, very properly, a 'pseudoposterior distribution'. DSZ refer to this work, but when faced with this discrepancy they still place more confidence in the 'reduction principle' than in the rules of probability theory.

In all these examples except one – that Example #5 again – an interesting phenomenon occurred. While the paradox was present for general improper priors in some infinite class C , there was always one particular improper prior in that class for which the paradox disappeared; B_1 and B_2 found themselves in agreement after all. DSZ noted this curious fact, but do not appear to have noticed its significance. We suggest that this was by far the most important fact uncovered in all the marginalization work.

Any prior $\pi(\eta)$ which leaves B_1 and B_2 in agreement must be *completely uninformative* about η (and, *a fortiori*, about ζ). This means that, far from casting doubt on the notion of complete ignorance, in the marginalization phenomena we have for the first time a purely

objective definition of complete ignorance that springs directly out of the product and sum rules of probability theory without appeal to any other notions like entropy or group invariance.

This is, again, an eminently satisfactory result; but why does it seem not to be true in DSZ's Example #5? There is still something new and important to be learned here.

15.9.1 The DSZ Example #5

We have data $D = \{x_1, \dots, x_n\}$ consisting of n observations from the standard normal sampling distribution $N(\mu, \sigma)$. With prior information I described by the proper prior pdf

$$p(d\mu d\sigma | I) = f(\mu, \sigma) d\mu d\sigma, \quad (15.77)$$

we have the usual joint posterior pdf for the parameters:

$$p(d\mu d\sigma | DI) = g(\mu, \sigma) d\mu d\sigma \quad (15.78)$$

with

$$g(\mu, \sigma) = \frac{f(\mu, \sigma) L(\mu, \sigma)}{\int d\mu \int d\sigma f(\mu, \sigma) L(\mu, \sigma)} \quad (15.79)$$

and the likelihood function

$$L(\mu, \sigma) = \sigma^{-n} \exp \left\{ -\frac{n}{2\sigma^2} [s^2 + (\mu - \bar{x})^2] \right\}, \quad (15.80)$$

in which, as usual, $\bar{x} \equiv n^{-1} \sum x_i$ and $s^2 \equiv n^{-1} \sum (x_i - \bar{x})^2$ are the sufficient statistics. Although we suppose the prior $f(\mu, \sigma)$ normalizable, it need not be actually normalized in (15.79) because any normalization constant appears in both numerator and denominator, and cancels out.

As long as $s^2 > 0$, the likelihood is bounded throughout the region of integration $-\infty < \mu < \infty, 0 \leq \sigma < \infty$, and therefore with a proper prior the integral in (15.79) is guaranteed to converge, leading to a proper posterior pdf. Furthermore, if the prior has moments of order m, k ,

$$\int_{-\infty}^{\infty} d\mu \int_0^{\infty} d\sigma \mu^m \sigma^k f(\mu, \sigma) < \infty, \quad (15.81)$$

the posterior distribution is guaranteed to have moments of higher order (in fact, all orders for μ and at least as high as order $k + n$ for σ). The solution is therefore very well-behaved mathematically.

But now we throw the proverbial monkey-wrench into this by declaring that we are interested only in the quantity

$$\zeta \equiv \frac{\mu}{\sigma}. \quad (15.82)$$

Making the change of variables $(\mu, \sigma) \rightarrow (\zeta, \sigma)$, the volume element transforms as $d\mu d\sigma = \sigma d\zeta d\sigma$, so writing $p(d\zeta|DI_1) = h_1(\zeta)d\zeta$, B_1 's marginal posterior pdf is

$$h_1(\zeta) = \int_0^\infty \sigma d\sigma g(\sigma\zeta, \sigma), \quad (15.83)$$

and in view of the high moments of g there are no convergence problems here, as long as $n > 1$. Thus far, there is no hint of trouble.

Now we examine the solution for a specific proper prior that can approach an improper prior. Consider the conjugate prior probability element

$$f(\mu, \sigma) d\mu d\sigma \propto \sigma^{-\gamma-1} \exp\{-\beta/\sigma - \alpha\mu^2\} d\mu d\sigma, \quad (15.84)$$

which is proper when $(\alpha, \beta, \gamma) > 0$, and tends to the Jeffreys uninformative prior $d\mu d\sigma/\sigma$ as $(\alpha, \beta, \gamma) \rightarrow 0$. This leads to the joint posterior pdf, $p(d\mu d\sigma|DI) = g(\mu, \sigma) d\mu d\sigma$ with density function

$$g(\mu, \sigma) \propto \sigma^{-n-\gamma-1} \exp\left\{-\frac{\beta}{\sigma} - \alpha\mu^2 - \frac{n}{2\sigma^2}[s^2 + (\mu - \bar{x})^2]\right\}, \quad (15.85)$$

from which we are to calculate the marginal posterior pdf for ζ alone by the integration (15.83). The result depends on both sufficient statistics (\bar{x}, s) , but is most easily written in terms of a different set. The quantities R, r , where

$$R^2 \equiv n(\bar{x}^2 + s^2) = \sum x_i^2, \quad r \equiv \frac{n\bar{x}}{R} = \frac{\sum x_i}{\sqrt{\sum x_i^2}}, \quad (15.86)$$

also form a set of jointly sufficient statistics, and from (15.85) and (15.83) we find the functional form $p(d\zeta|DI_1) = h_1(\zeta|r, R)d\zeta$, where

$$h_1(\zeta|r, R) \propto \exp\left\{-\frac{n\zeta^2}{2}\right\} \int_0^\infty d\omega \omega^{n+\gamma-1} \exp\left\{-\frac{1}{2}\omega^2 + r\zeta\omega - \beta R^{-1}\omega - \alpha\zeta^2 R^2\omega^{-2}\right\}. \quad (15.87)$$

As long as α or β is positive, the result depends on both sufficient statistics, as Fraser predicted; but, as α, β tend to zero and we approach an improper prior, the statistic R becomes less and less informative about ζ , and when α, β both vanish the dependence on R drops out altogether:

$$h_1(\zeta|r, R) \rightarrow h_1(\zeta|r) \propto \exp\left\{-\frac{n\zeta^2}{2}\right\} \int_0^\infty d\omega \omega^{n+\gamma-1} \exp\left\{-\frac{1}{2}\omega^2 + r\zeta\omega\right\}. \quad (15.88)$$

If then one were to look only at the limiting case $\alpha = \beta = 0$ and not at the limiting process, it might appear that just r alone is a sufficient statistic for ζ , as it did in (15.60). This supposition is encouraged by noting that the sampling distribution for r in turn depends only on ζ , not on μ and σ separately:

$$p(r|\mu\sigma) \propto (n-r^2)^{(n-3)/2} \int_0^\infty d\omega \omega^{n-1} \exp\left\{\frac{1}{2}\omega^2 + r\zeta\omega\right\}. \quad (15.89)$$

It might then seem that, in view of (15.88) and (15.89), we should be able to derive the same result by applying Bayes' theorem to the reduced sampling distribution (15.89). But one who supposes this finds, to his dismay, that (15.89) is not a factor of (15.88); that is, the ratio $h_1(\zeta|r)/p(r|\zeta)$ depends on r as well as ζ . The Jeffreys uninformative prior $\gamma = 0$ does indeed make the two integrals equal, but there remains an uncompensated factor with $(n - r^2)$, and so even the uninformative Jeffreys prior for (μ, σ) cannot bring about agreement of B_1 and B_2 . There is no prior $p(\zeta|I_2)$ that can yield B_1 's posterior distribution (15.88) from B_2 's sampling distribution (15.89).

Since the paradox is still present for a proper prior, this is another counter-example to (15.65); but it has a deeper meaning for us. What is now the information being used by B_1 but ignored by B_2 ? It is not the prior probability for the nuisance parameter; the new feature is that in this model the mere qualitative fact of the existence of the nuisance parameter in the model *already constitutes prior information relevant to B_1 's inference*, which B_2 is ignoring.

Recognizing this, we suddenly see the whole subject in a much broader light. We found above that (15.60) is not essential to the marginalization phenomenon; now we see that concentration on the nuisance parameter η is not an essential feature either! If there is any prior information whatsoever that is relevant to ζ , *whether or not it refers to η* , that B_1 is taking into account but B_2 is not, then we are in the same situation, and our two Bayesians come, necessarily, to different conclusions. In other words, DSZ considered only a very special case of the real phenomenon.

This situation is discussed in Jaynes (1980, following Eq. (79)), where the phenomenon is called 'ζ-overdetermination'. Reverting to our original notation in (15.58) and denoting B_1 's prior information by I_1 , it is shown that the general necessary and sufficient condition for agreement of B_1 and B_2 is that

$$\int d\eta p(y|z\eta\zeta I_1)\pi(\eta) = p(y|z\zeta I_1) \quad (15.90)$$

shall be independent of ζ for all possible samples y, z . Denoting the parameter space and our partitioning into subspaces by $S_\theta = S_\zeta \otimes S_\eta$, we may write this as

$$\int_{S_\eta} d\eta p(yz|\eta\zeta)\pi(\eta) = p(y|zI_1)p(z|\zeta) \quad \left\{ \begin{array}{l} \zeta \in S_\zeta \\ \text{all } y, z \end{array} \right. \quad (15.91)$$

or, more suggestively,

$$\int_{S_\eta} d\eta K(\zeta, \eta)\pi(\eta) = \lambda f(\zeta). \quad (15.92)$$

This is a Fredholm integral equation in which the kernel is B_1 's likelihood, $K(\zeta, \eta) = p(yz|\zeta\eta)$, the 'driving force' is B_2 's likelihood $f(\zeta) = p(z|\zeta)$, and $\lambda(y, z) \equiv p(y|zI_1)$ is an unknown function to be determined from (15.92). But now we see the meaning of 'uninformative' much more deeply; for every different data set (y, z) there is a different integral equation. Therefore, for a single prior $\pi(\eta)$ to qualify as 'uninformative', it must satisfy many different (in general, an uncountable number) of these integral equations simultaneously.

At first glance, it seems almost beyond belief that any prior could do this; from a mathematical standpoint the condition seems hopelessly overdetermined, casting doubt on the notion of an uninformative prior. Yet we have many examples where such a prior does exist. In Jaynes (1980) we analyzed the structure of these integral equations in some detail, showing that the different status of Example #5 is due to the ‘incompleteness’ of the kernel.

More specifically, the set of all L^2 functions on S_ζ forms a Hilbert space H_ζ . For any specified data set $x = (y, z)$, as η ranges over S_η , the functions $K(\zeta, \eta)$, in their dependence on ζ , span a certain subspace $H'_\zeta(y, z) \in H_\zeta$. The kernel is said to be *complete* if $H'_\zeta = H_\zeta$. If it is incomplete, then if there is any data set (y, z) for which $f(\zeta)$ does not lie in H'_ζ , there can be no solution of (15.92). In such cases, the mere qualitative fact of the *existence* of the components (y, η) – irrespective of their numerical values – already constitutes prior information relevant to B_1 ’s inference, because introducing them into the model restricts the space of B_1 ’s possible likelihood functions (from different data sets y, z) from H_ζ to H'_ζ . In this case the shrinkage of H_ζ cannot be restored by any prior on S_η , and there is no possibility for agreement of B_1 and B_2 .

In general, the point is that the integral equation for any one data set x imposes only very weak conditions on $\pi(\eta)$, determining its projection on only a tiny subspace $H(x) \in H_\zeta$. As we consider different data sets, the $H(x)$ are scattered about, like stars in the sky, within the full Hilbert space H_ζ . There is room for all of them, so the system of integral equations has nontrivial solutions after all.

15.9.2 Summary

Looking at the above equations with all this in mind, we now see that there was never any paradox or inconsistency after all; one should not have expected (15.88) to be derivable from (15.89) by Bayes’ theorem because they are the posterior distribution and sampling distribution for two different problems, in which the model has different parameters. Eq. (15.88) is the correct marginal posterior pdf for ζ in a problem P_1 with two parameters (ζ, σ) ; but, although σ is integrated out to form the marginal pdf, the result still depends on what prior we have assigned to σ – as it should, since, if σ is known, it is highly relevant to the inference; if it is unknown, any partial prior information we have about it must still be relevant.

In contrast, (15.89) can be interpreted as a valid sampling distribution for a problem P_2 in which ζ is the only parameter present; the prior information does not even include the existence of the parameter σ which was integrated out in P_1 . With a prior density $f_2(\zeta)$ it would yield a posterior pdf

$$h_2(\zeta) \propto f_2(\zeta) \int d\omega \omega^{n-1} \exp \left\{ -\frac{1}{2} \omega^2 + r\zeta\omega \right\} \quad (15.93)$$

of a different functional form than (15.88). In view of the earlier work of Jeffreys and of Geisser and Cornfield, one could hardly claim that the situation was new and startling, much less paradoxical.

Forty years earlier, Harold Jeffreys was immune from such errors because (1) he perceived that the product and sum rules of probability theory are adequate to conduct inference and they take precedence over intuitive *ad hoc* devices like the reduction principle; (2) he had recognized from the start that all inferences are necessarily conditional not only on the data, but also on the prior information – therefore his formal probability symbols $P(A|BI)$ always indicated the prior information I , which included specification of the model.

Today, it seems to us incredible that anyone could have examined even one problem of inference without perceiving this necessary role of prior information; what kind of logic could they have been using? Nevertheless, those trained in the ‘orthodox’ tradition of probability theory did not recognize it. They did not have a term for prior information in their vocabulary, much less a symbol for it in their equations; and *a fortiori* no way of indicating when two probabilities are conditional on different prior information.⁷ So they were helpless when prior information mattered.

15.10 A useful result after all?

In most paradoxes there is something of value to be salvaged from the debris, and we think (Jaynes, 1980) that the marginalization paradox may have made an important and useful contribution to the old problem of ‘complete ignorance’. How is the notion to be defined, and how is one to construct priors expressing complete ignorance? We have discussed this from the standpoint of entropy and symmetry (transformation groups) in previous chapters; now marginalization suggests still another principle for constructing uninformative priors.

Many cases are known, of which we have seen examples in DSZ, where a problem has a parameter of interest ζ and an uninteresting nuisance parameter η . Then the marginal posterior pdf for ζ will depend on the prior assigned to η as well as on the sufficient statistics. Now for certain particular priors $p(\eta|I)$ one of the sufficient statistics may drop out of the marginal distribution $p(\zeta|DI)$, as R did in (15.88). It is at first glance surprising that the sampling distribution for the remaining sufficient statistics may in turn depend only on ζ as in (15.89).

Put differently, suppose a problem has a set of sufficient statistics (t_1, t_2) for the parameters (ζ, η) . Now, if there is some function $r(t_1, t_2)$ whose sampling distribution depends only on ζ , so that $p(r|\zeta\eta I) = p(r|\zeta I)$, this defines a pseudoproblem with different prior information I_2 , in which η is never present at all. Then there may be a prior $p(\eta|I)$ for which the posterior marginal distribution $p(\zeta|DI) = p(\zeta|rI)$ depends only on the component r of the sufficient

⁷ Indeed, in the period 1930–1960 nearly all orthodoxians, under the influence of R. A. Fisher, scorned Jeffreys’ work, and some took a militant stand against prior information, teaching their students that it is not only intellectually foolish, but also morally reprehensible – a deliberate breach of ‘scientific objectivity’ – to allow one’s self to be influenced by prior information at all! This did little damage in the very simple problems considered in the orthodox literature, where there was no significant prior information anyway. And it did relatively little damage in physical science where prior information is important, because scientists ignored orthodox teaching and persisted in doing, qualitatively, the Bayesian reasoning using prior information that their own common sense told them was the right thing to do. But we think it was a disaster for fields such as econometrics and artificial intelligence, where adoption of the orthodox view of probability had the automatic consequence that the significant problems could not even be formulated, much less solved, because the orthodox view of probability theory does not recognize probability as expressing information at all.

statistic. This happened in the example studied above; but now, more may be true. It may be that for that prior on η the pseudoposterior pdf for ζ is identical with the marginal pdf in the original problem. If a prior brings about agreement between the marginal posterior and the pseudoposterior distributions, how should we interpret this?

Suppose we start from the pseudoproblem. It seems that if introducing a new parameter η and using the prior $p(\eta|I)$ makes no difference, then it has conveyed no *information* at all about ζ : that prior must express ‘complete ignorance’ of η in a rather fundamental sense. In all cases yet found the prior $p(\eta|I)$ which does this on an infinite domain is improper; this lends support to that conclusion because, as noted, our common sense should have told us that *any proper prior on an infinite domain is necessarily informative about η* ; it places some finite limits on the range of values that η could reasonably have, whether we interpret ‘reasonably’ as ‘with 99% probability’ or ‘with 99.9% probability’ . . . and so on.

Can this observation be extended to a general technique for constructing uninformative priors beyond the location and scale parameter cases? This is at present an ongoing research project rather than a finished part of probability theory, so we defer it for the future.

15.11 How to mass-produce paradoxes

Having examined a few paradoxes, we can recognize their common feature. Fundamentally, the procedural error was always failure to obey the product and sum rules of probability theory. Usually, the mechanism of this was careless handling of infinite sets and limits, sometimes accompanied by attempts to replace the rules of probability theory by intuitive *ad hoc* devices like B_2 ’s ‘reduction principle’. Indeed, paradoxes caused by careless dealing with infinite sets or limits can be mass-produced by the following simple procedure:

- (1) Start from a mathematically well-defined situation, such as a finite set, a normalized probability distribution, or a convergent integral, where everything is well-behaved and there is no question about what is the correct solution.
- (2) Pass to a limit – infinite magnitude, infinite set, zero measure, improper pdf, or some other kind – without specifying how the limit is approached.
- (3) Ask a question whose answer depends on how the limit was approached.

This is guaranteed to produce a paradox in which a seemingly well-posed question has more than one seemingly right answer, with nothing to choose between them. The insidious thing about it is that, as long as we look only at the limit, and not the limiting process, the source of the error is concealed from view.

Thus, it is not surprising that those who persist in trying to evaluate probabilities directly on infinite sets have been able to study finite additivity and nonconglomerability for decades – and write dozens of papers of impressive scholarly appearance about it. Likewise, those who persist in trying to calculate probabilities conditional on propositions of probability zero, have before them an unlimited field of opportunities for scholarly looking research and publication – without hope of any meaningful or useful results.

In our opening quotation, Gauss had a situation much like this in mind. Whenever we find a belief that such infinite sets possess some kind of ‘existence’ and mathematical properties in their own right, independent of any such limiting process, we can expect to see paradoxes of the above type. But note that this does not in any way prohibit us from using infinite sets to define *propositions*. Thus the proposition

$$G \equiv 1 \leq x \leq 2 \quad (15.94)$$

invokes an uncountable set, but it is still a single discrete proposition, to which we may assign a probability $P(G|I)$ defined on a sample space of a finite number of such propositions without violating our ‘probabilities on finite sets’ policy. We are not assigning any probability directly on an infinite set.

But then if we replace the upper limit 2 by a variable quantity z , we may (and nearly always do) find that this defines a well-behaved function, $f(z) \equiv P(G|zI)$. In calculations, we are then free to make use of whatever analytic properties this function may have, as we noted in Chapter 6. Even if $f(z)$ is not an analytic function, we may be able to define other analytic functions from it, for example by integral transforms. In this way, we are able to deal with any real application that we have been able to imagine, by discrete algebraic or continuum analytical methods, without losing the protection of Cox’s theorems.

15.12 Comments

In this chapter and Chapter 5, we have seen two different kinds of paradox. There are ‘conceptually generated’ ones, such as the Hempel paradox of Chapter 5, which arise from placing faulty intuition above the rules of probability theory, and ‘mathematically generated’ ones, such as nonconglomerability, which arise mostly out of careless use of infinite sets. Marginalization is an elaborate example of a compound paradox, generated by both conceptual errors and mathematical errors which happened to reinforce each other. It seems that nothing in the mathematics can protect us against conceptual errors, but we might ask whether there are better ways of protection against mathematical ones.

Back in Chapter 2 we saw that the rules of probability theory can be derived as necessary conditions for consistency, as expressed by Cox’s functional equations. The proofs applied to finite-sets of propositions, but when the results of a finite-set calculation can be extended to an infinite set by a mathematically well-behaved passage to a limit, we also accept that limit.

It might be thought that it would be possible, and more elegant, to generalize Cox’s proofs so that they would apply directly to infinite sets; and indeed that is what the writer believed and tried to carry out for many years. However, since at least the work of Bertrand (1889), the literature has been turning up paradoxes that result from attempts to apply the rules of probability theory directly and indiscriminately on infinite sets; we have just seen some representative examples and their consequences. Since in recent years there has been a sharp increase in this paradoxing, one must take a more cautious view of infinite sets.

Our conclusion – based on some 40 years of mathematical efforts and experience with real problems – is that, at least in probability theory, an infinite set should be thought of only as the limit of a specific (i.e. unambiguously specified) sequence of finite sets. Likewise, an improper pdf has meaning only as the limit of a well-defined sequence of proper pdfs. The mathematically generated paradoxes have been found only when we tried to depart from this policy by treating an infinite limit as something already accomplished, without regard to any limiting operation. Indeed, experience to date shows that almost any attempt to depart from our recommended ‘finite-sets’ policy has the potentiality for generating a paradox, in which two equally valid methods of reasoning lead us to contradictory results.

The paradoxes studied here stand as counter-examples to any hope that we can ever work with full freedom on infinite sets. Unfortunately, the Borel–Kolmogorov and marginalization paradoxes turn up so seldom as to encourage overconfidence in the inexperienced. As long as one works on problems where they do not cause trouble, the psychological phenomenon: ‘You can’t argue with success!’, noted at the beginning of this Chapter, controls the situation. Our reply to this is, of course, ‘You can and should argue with success that was obtained by fraudulent means’.

Mea culpa

For many years, the present writer was caught in this error just as badly as anybody else, because Bayesian calculations with improper priors continued to give just the reasonable and clearly correct results that common sense demanded. So warnings about improper priors went unheeded; just that psychological phenomenon. Finally, it was the marginalization paradox that forced recognition that we had only been lucky in our choice of problems. If we wish to consider an improper prior, the only correct way of doing it is to approach it as a well-defined limit of a sequence of proper priors. If the correct limiting procedure should yield an improper posterior pdf for some parameter α , then probability theory is telling us that the prior information and data are too meager to permit any inferences about α . Then the only remedy is to seek more data or more prior information; probability theory does not guarantee in advance that it will lead us to a useful answer to every conceivable question.

Generally, the posterior pdf is better behaved than the prior because of the extra information in the likelihood function, and the correct limiting procedure yields a useful posterior pdf that is analytically simpler than any from a proper prior. The most universally useful results of Bayesian analysis obtained in the past are of this type, because they tended to be rather simple problems, in which the data were indeed so much more informative than the prior information that an improper prior gave a reasonable approximation – good enough for all practical purposes – to the strictly correct results (the two results agreed typically to six or more significant figures).

In the future, however, we cannot expect this to continue because the field is turning to more complex problems in which the prior information is essential and the solution is found by computer. In these cases it would be quite wrong to think of passing to an improper prior. That would lead usually to computer crashes; and, even if a crash is avoided,

the conclusions would still be, almost always, quantitatively wrong. But, since likelihood functions are bounded, the analytical solution with proper priors is always guaranteed to converge properly to finite results; therefore it is always possible to write a computer program in such a way (avoid underflow, etc.) that it cannot crash when given proper priors. So, even if the criticisms of improper priors on grounds of marginalization were unjustified, it remains true that in the future we shall be concerned necessarily with proper priors.

Note added

Preliminary versions of this chapter were made available to many interested persons, for comments and suggestions. Several have expressed, both privately and publicly, their appreciation for these clarifications of issues that have long been mysterious and confused, and even some compulsive nitpickers have failed to raise any objections. Only one source has exhibited that psychological phenomenon noted in Chapter 5 in connection with the Hempel paradox; someone asserts a principle that seems to him intuitively right, and when probability analysis reveals the error, instead of taking this opportunity to educate his intuition, he reacts by rejecting the probability analysis. For him, his intuitive *ad hoc* principle takes precedence over the rules of probability theory.

If the issue is only which is to take precedence, there does not seem to be any way to resolve it; if one is not convinced by Cox's theorems and our great deal of experience confirming what they tell us, then we shall just have to agree to disagree. But if the issue is one of mathematically demonstrable fact, then it can be resolved at once – in the minds of everyone except the one who proposed the principle. One can be so deeply committed to his position that mathematical proof to the contrary, and any number of counter-examples, carry no weight for him. That is just what has happened.

In the case of the tumbling tetrahedra problem, we pointed out the error in previous discussions, gave the exact solution (15.26) according to the rules of probability theory, an asymptotic approximation (15.18) to it, and after a little thought could see that the final result (15.34) was really obvious from the start. That should be enough; this is an issue of mathematically demonstrable fact, the mathematics is before us, and every reader can judge it for himself.

The case of the marginalization paradox is very similar. The real purpose of this note is to stress what is the issue here. DSZ (p. 194) purported to have proved that the disagreement of B_1 and B_2 'could not have arisen if B_1 had employed proper prior distributions'. We pointed out that their proof is based on mutually contradictory assumptions, and reinforced this by (1) pointing out that a counter-example to what they claimed to have proved was already present in plain sight in the original DSZ article (their Example #5, where it is evident by inspection that the disagreement is present for all priors, proper or improper), and (2) gave in (15.76) another counter-example, where B_1 uses a proper prior, but B_1 and B_2 still disagree because B_1 's posterior distribution for ζ depends on the datum y and the prior $p(\eta|I_1)$, as common sense tells us it must when B_1 uses a proper – therefore informative – prior for η . In fact, we can leave it as an exercise for the reader to verify that every one of the DSZ

examples is an equally good counter-example, if you look at what happens when B_1 uses a proper prior. Of course, B_1 and B_2 agree when B_1 uses a proper prior in the trivial case where we are concerned with two independent problems; then the sampling distribution factors in the form $p(yz|\eta\zeta) = p(y|\eta)p(z|\zeta)$ and the prior $p(\eta\zeta|I_1)$ also factors.

But if the basic theorem is invalid, then the entire marginalization tale collapses; if one applies the rules of probability theory correctly, as explained long ago by Harold Jeffreys, there is no paradox. Again, this is an issue of mathematically demonstrable fact for which we have given the relevant mathematics, so we see no reason to engage in continuing debate over it; every reader can judge it for himself.