



For our purposes, everything follows from the product rule:

$$p(AB|C) = p(A|BC)p(B|C) = p(B|AC)p(A|C). \quad (14.2)$$

If the propositions  $B$  and  $C$  are not mutually contradictory, this may be rearranged to give the rule of ‘learning by experience’, Bayes’ theorem:

$$p(A|BC) = p(A|C) \frac{p(B|AC)}{p(B|C)} = p(A|B) \frac{p(C|AB)}{p(C|B)}. \quad (14.3)$$

If there are several mutually exclusive and exhaustive propositions  $B_i$ , then, by summing (14.2) over them, we obtain the chain rule

$$p(A|C) = \sum_i p(A|B_iC)p(B_i|C) \quad (14.4)$$

or, in a simple skeleton notation,

$$p(A|C) = \sum_B p(A|BC)p(B|C). \quad (14.5)$$

Now let

$$\begin{aligned} X &= \text{prior knowledge, of any kind whatsoever,} \\ S &= \text{signal,} \\ N &= \text{noise,} \\ V &= V(S, N) = \text{observed voltage,} \\ D &= \text{decision about the nature of the signal.} \end{aligned} \quad (14.6)$$

Thus we have

$$\begin{aligned} p(S|X) &= \text{prior probability for the particular signal } S, \\ p(N|X) &= W(N) = \text{prior probability for the particular sample of noise } N. \end{aligned} \quad (14.7)$$

We understand that the prior information  $X$  is always built into the right-hand side of all our probability symbols, whether or not we write it explicitly. Thus, in a linear system,  $V = S + N$  and

$$p(V|S) \equiv p(V|SX) = W(V - S). \quad (14.8)$$

The reader may be disturbed by the absence of density functions,  $dS$ ,  $dN$ , etc., which might be expected in the case of continuous  $S$ ,  $N$ . Note, however, that our equations are homogeneous in these quantities, so they cancel out anyway. We are trying only to convey the broad ideas, without bothering with fine details which would make the notation very intricate. Thus by  $\sum_A$  we mean ordinary summation over some previously agreed set of possible values if  $A$  is discrete, integration with appropriate density functions if  $A$  is continuous.

A decision rule  $p(D_i|V_j)$ , or for brevity just  $p(D|V)$ , represents the process of drawing inferences about the signal from the observed voltage. If it is always made in a definite way, then  $p(D|V)$  has only the values 0, 1 for any choice of  $D$  and  $V$ ; however, we may also have

a ‘randomized’ decision rule according to which  $p(D|V)$  is a true probability distribution. Maintaining this more general view turns out to be a help in formulating the theory.

The essence of any decision rule, and in particular any one which can be built into automatic equipment, is that the decision must be made on the basis of  $V$  alone;  $V$  is, by definition, the quantity which contains all the information actually used (in addition to the ever-present  $X$ ) in arriving at the decision. Thus, if  $Y \neq D$  is any other proposition, we have

$$p(D|V) = p(D|VY). \quad (14.9)$$

The fact that  $Y$  is to be ignored in the presence of  $V$  might appear a departure from our previous exhortations that the robot is always to take into account all the relevant information it has. However, if we consider that the property (14.9) is a part of the prior information  $X$  there is no difficulty. To put it differently, (14.9) expresses the prior knowledge that there is a direct logical relation by which  $D$  is determined by  $V$  alone. If this relationship was a known law of physics, there would be nothing strange in (14.9). The only difference is that in the present case this relationship does not express any law of Nature, but rather our own design of the apparatus. Then  $Y$  is ignored not because the robot has relaxed its rules, but because our design makes  $Y$  irrelevant.

An equivalent statement is that the probability for reaching a decision  $D$  depends on any proposition  $Y$  only through the intermediate influence of  $V$  on  $V$ :

$$p(D|Y) = \sum_V p(D|V)p(V|Y) \quad (14.10)$$

which is a kind of ‘Huygens principle’ for logic. To see the analogy, think of  $Y$  as a light source which cannot be seen from  $D$ , but it illuminates various points  $V$ . Then the resulting light arriving at  $D$  is the sum of the Huygens wavelets  $p(D|V)$  with amplitudes  $p(V|Y)$ . The almost exact mathematical analogy between conditional information flow and the flow of light according to the Huygens principle of optics appears in statistical mechanics of irreversible processes.

## 14.2 Sufficiency and information

Equation (14.9) has interesting consequences; suppose we wish to judge the plausibility of some proposition  $Y$ , on the basis of knowledge of  $V$  and  $D$ . From the product rule (14.2),

$$p(DY|V) = p(Y|VD)p(D|V) = p(D|VY)p(Y|V) \quad (14.11)$$

and, using (14.9), this reduces to

$$p(Y|VD) = p(Y|V). \quad (14.12)$$

Thus, if  $V$  is known, knowledge of  $D$  is redundant and cannot help us in estimating any other quantity. The reverse is not true, however; we could equally well use (14.9) in

another way:

$$p(VY|D) = p(Y|VD)p(V|D) = p(Y|D)p(V|YD). \quad (14.13)$$

Combining this with (14.12), there results the following theorem.

*Theorem*

Let  $D$  be a possible decision, given  $V$ . Then  $p(V|D) \neq 0$ , and

$$p(Y|V) = p(Y|D) \text{ if and only if } p(V|D) = p(V|YD). \quad (14.14)$$

In words: knowledge of  $D$  is as good as knowledge of  $V$  for judgments about  $Y$  if and only if  $Y$  is irrelevant for judgments about  $V$ , given  $D$ . Stated differently: in the ‘environment’ produced by knowledge of  $D$ , the probabilities for  $Y$  and  $V$  are independent, i.e.

$$p(YV|D) = p(Y|D)p(V|D). \quad (14.15)$$

In this case, in the literature of this field  $D$  is said to be a *sufficient statistic* for judgments about  $Y$ . We shall want to see whether this is in agreement with our earlier definitions of sufficiency, made from a quite different point of view in Chapter 8.

Evidently, a decision rule which makes  $D$  a sufficient statistic for judgments about the signal  $S$  is superior to one without this property, in that it tells us more about the signal. However, such a rule does not necessarily exist. Equation (14.15) is a very restrictive condition, since it must be satisfied for all values of  $Y$ ,  $V$ , and all  $D$  for which  $p(D|V) \neq 0$ .

As you might guess from this, the concept of sufficiency is closely related to that of information. The above definition of sufficiency could be stated equally well as:  $D$  is a sufficient statistic for judgments about  $Y$  if it contains all the information about  $Y$  which  $V$  contains. Since  $D$  is determined from  $V$ , if it is not a sufficient statistic, it necessarily contains *less* information about  $Y$  than does  $V$ . In this statement, the term ‘information’ was used in a loose, intuitive sense; does it remain true if we adopt Shannon’s measure of information?

Imagine that there are several mutually exclusive propositions  $Y_i$ , one of which must be true. For brevity we use, as above, the notation  $\sum_Y f(Y) \equiv \sum_i f(Y_i)$ . With a specific value of  $D$  given, the entropy which measures our information about the propositions  $Y_i$  is

$$H_D(Y) = - \sum_Y p(Y|D) \log[p(Y|D)], \quad (14.16)$$

and its expectation over all values of  $D$  is

$$\overline{H}_D(Y) = \sum_D p(D|X) H_D(Y). \quad (14.17)$$

If

$$\overline{H}_C(Y) < \overline{H}_D(Y), \quad (14.18)$$

we say colloquially that  $C$  contains, ‘on the average’, more information about  $Y$  than does  $D$ . Note, however, that it may be otherwise for specific values of  $C$  and  $D$ .

Acquisition of new information can never increase  $\bar{H}$ ; let  $\{Z_i\}$  be, for the moment, any set of propositions and form the expression

$$\begin{aligned}\bar{H}_V(Z) - \bar{H}_{DV}(Z) &= \sum_{DVZ} p(DV|X)p(Z|DV) \log[p(Z|DV)] \\ &\quad - \sum_{VZ} p(V|X)p(Z|V) \log[p(Z|V)] \\ &= \sum_{DVZ} p(DV|X)p(Z|DV) \log \left[ \frac{p(Z|DV)}{p(Z|V)} \right].\end{aligned}\tag{14.19}$$

Using the fact that on the positive real line  $\log(x) \geq (1 - x^{-1})$ , with equality if and only if  $x = 1$ , this becomes

$$\bar{H}_V(Z) - \bar{H}_{DV}(Z) \geq \sum_{DVZ} p(DV|X)[p(Z|DV) - p(Z|V)] = 0.\tag{14.20}$$

Thus,  $\bar{H}_{DV}(Z) \leq \bar{H}_V(Z)$ , with equality if and only if (14.12) holds for all  $D$ ,  $V$  and  $Z$  for which  $p(DV|X) \neq 0$ .

But now, since (14.20) holds regardless of the meaning of  $D$  and  $V$ , we can conclude equally well that, for all  $D$ ,  $V$ ,  $Z$ ,

$$\bar{H}_D(Y) \geq \bar{H}_{DV}(Z) \leq \bar{H}_V(Z).\tag{14.21}$$

Choosing  $Z = Y$ , we have in consequence of (14.12)  $H_V(Y) = H_{DV}(Y)$ , so that

$$\bar{H}_V(Y) \leq \bar{H}_D(Y),\tag{14.22}$$

with equality if and only if (14.15) holds, i.e. if and only if  $D$  is a sufficient statistic as just defined. Thus, if by ‘information’ we mean minus the expectation of the entropy of  $Y$  over the prior distribution of  $D$  or  $V$ , zero information loss in going from  $V$  to  $D$  is equivalent to sufficiency of  $D$ . Note that inequality (14.20) holds only for the expectations of  $\bar{H}$ , not for the  $H$ . Acquisition of a specific piece of information (that an event previously considered improbable had in fact occurred) may in some cases increase the entropy of  $Y$ . However, this is an improbable situation, and on the average the entropy can only be lowered by additional information. This shows again that the term ‘information’ is not a happy choice of words to describe entropy expressions. In spite of the entropy increases, the situation just described could hardly be called one of less *information* in the colloquial sense of that word; but rather one of less *certainty*.

### 14.3 Loss functions and criteria of optimum performance

In order to say that one decision rule is better than another, we need some specific criterion of what we want our detection system to accomplish. The criterion will vary with the application, and obviously no single decision rule can be best for all purposes. But our

discussion in Chapter 13 will apply, almost unchanged, in this slightly different language. A very general type of criterion is obtained by assigning a *loss function*  $L(D, S)$  which represents our judgment of how serious it is to make decision  $D$  when signal  $S$  is in fact present.

In the case where there are only two possible signals,  $S_0 = 0$  (i.e. no signal), and  $S_1 > 0$ , and consequently two possible decisions  $D_0, D_1$  about the signal, there are two types of error, the false alarm  $A = (D_1, S_0)$  and the false rest  $R = (D_0, S_1)$ . In some applications, one type of error might be much more serious than the other.

Suppose that a false rest is considered ten times as serious as is a false alarm, while a correct decision of either type represents no 'loss'. We could then take  $L(D_0, S_0) = L(D_1, S_1) = 0$ ,  $L(D_0, S_1) = 10$ ,  $L(D_1, S_0) = 1$ . Whenever the possible signals and the possible decisions form discrete sets, the loss function becomes a *loss matrix*. In the above example,

$$L_{ij} = \begin{pmatrix} 0 & 10 \\ 1 & 0 \end{pmatrix}. \quad (14.23)$$

Instead of assigning arbitrarily a certain loss value to each possible type of detection error, we may consider *information loss* by the assignment  $L(D, S) = -\log[p(S|D)]$ . This is somewhat more difficult to manipulate, because now  $L(D, S)$  depends on the decision rule. A decision rule which minimizes information loss is one which makes the decision in some sense as close as possible to being a sufficient statistic for judgments about the signal. In exactly what sense seems never to have been clarified. The *conditional loss*  $L(S)$  is the expected loss incurred when the specific signal  $S$  is present:

$$L(S) = \sum_D L(D, S)p(D|S), \quad (14.24)$$

which may in turn be expressed in terms of the decision rule and the properties of the noise by using (14.10). What is often called colloquially the 'average loss' is the expectation of the conditional loss over all possible signals:

$$\langle L \rangle = \sum_S L(S)p(S|X). \quad (14.25)$$

Two different criteria of optimal performance now suggest themselves:

**The minimax criterion.** For a given decision rule  $p(D|V)$ , consider the conditional loss  $L(S)$  for all possible signals, and let  $[L(S)]_{\max}$  be the maximum value attained by  $L(S)$ . We seek that decision rule for which  $[L(S)]_{\max}$  is as small as possible. As we noted in Chapter 13, this criterion concentrates attention on the worst possible case, regardless of the probability for occurrence of this case, and it is thus in a sense too conservative. However, it gives some psychological comfort that it does not involve the prior probabilities for the different signals  $p(S|X)$ , and therefore it can be applied by persons who, under the handicap of orthodox training, have a mental hangup against prior probabilities.

**The Bayes criterion.** We seek that decision rule for which the expected loss  $\langle L \rangle$  is minimized. In order to apply this, a prior distribution  $p(S|X)$  must be available.

Other criteria were proposed before the days of Wald's decision theory. In the Neyman–Pearson theory, we fix the probability for occurrence of one type of error at some small value  $\epsilon$ , and then minimize the probability  $\delta$  of the other type of error subject to this constraint.<sup>1</sup> Arnold Siebert's 'ideal observer' minimizes the total probability for error ( $\epsilon + \delta$ ).

After having invented many different such *ad hoc* criteria from various viewpoints, and arguing their relative merits on philosophical grounds, the basic mathematical identity of all these criteria came as quite a surprise to the early workers in this field. We shall see below that all of them are special cases of the Bayes criterion, for particular prior probabilities.

Let us find the Bayes solution, as it was rationalized in decision theory. Substituting in succession (14.24), (14.10), and (14.9) into (14.25), we obtain for the expected loss

$$\langle L \rangle = \sum_{DV} \left[ \sum_S L(D, S) p(VS|X) \right] p(D|V). \quad (14.26)$$

If  $L(D, S)$  is a definite function independent of  $p(D|V)$  (this assumption excludes for the moment the information loss function), there is no function  $p(D|V)$  for which this expression is stationary in the sense of the calculus of variations. We then minimize  $\langle L \rangle$  merely by choosing for each possible  $V$  that decision  $D_1(V)$  for which the coefficient in (14.26)

$$K(D, V) \equiv \sum_S L(D_1, S) p(VS|X) \quad (14.27)$$

is a minimum. Thus, we adopt the decision rule

$$p(D|V) = \delta(D, D_1). \quad (14.28)$$

In general, there will be only one such  $D_1$ , and the best decision rule is nonrandom. However, in case of 'degeneracy',  $K(D_1, V) = K(D_2, V)$ , any randomized rule of the form

$$p(D|V) = a\delta(D, D_1) + b\delta(D, D_2), \quad a + b = 1, \quad (14.29)$$

is just as good by the criterion being used. This degeneracy occurs at 'threshold' values of  $V$ , where we change from one decision to another.

#### 14.4 A discrete example

Consider the case already mentioned, where there are two possible signals,  $S_0$  and  $S_1$ , and a loss matrix

$$L_{ij} = \begin{pmatrix} L_{00} & L_{01} \\ L_{10} & L_{11} \end{pmatrix} = \begin{pmatrix} 0 & L_r \\ L_a & 0 \end{pmatrix}, \quad (14.30)$$

<sup>1</sup> For example, we suspect that at an Early Warning Radar Installation, the primary constraint might be that the Commanding Officer shall not be roused out of bed by a false alarm more often than once per month, and, subject to that requirement, we minimize the probability for a false rest.

where  $L_a$ ,  $L_r$  are the losses incurred by a false alarm and a false rest, respectively. Then

$$\begin{aligned} K(D_0, V) &= L_{01}p(VS_1|X) = L_r p(VS_1|X), \\ K(D_1, V) &= L_{10}p(VS_0|X) = L_a p(VS_0|X), \end{aligned} \quad (14.31)$$

and the decision rule that minimizes  $\langle L \rangle$  is

$$\begin{aligned} \text{choose } D_1 & \text{ if } \frac{p(VS_1|X)}{p(VS_0|X)} > \frac{L_a}{L_r}, \\ \text{choose } D_0 & \text{ if } \frac{p(VS_1|X)}{p(VS_0|X)} < \frac{L_a}{L_r}, \\ & \text{choose either at random in case of equality.} \end{aligned} \quad (14.32)$$

If the prior probabilities for a signal and no signal are

$$p(S_1|X) = p, \quad p(S_0|X) = q = 1 - p, \quad (14.33)$$

respectively, the decision rule becomes

$$\text{choose } D_1 \text{ if } \frac{p(V|S_1)}{p(V|S_0)} > \frac{qL_a}{pL_r}, \quad \text{etc.} \quad (14.34)$$

The left-hand side of (14.34) is a likelihood ratio, which depends only on the pdf assigned to the noise, and is the quantity which should be computed by the optimum receiver according to the Bayes criterion.

This same quantity is the essential one regardless of the assumed loss function and regardless of the probability for the occurrence of the signal; these affect only the threshold of detection. Furthermore, if the receiver merely computes this likelihood ratio and delivers it at the output without making any decision, it provides us with all the information we need to make optimum decisions in the Bayes sense. Note the generality of this result, which is important for applications; no assumptions were needed as to the type of signal, linearity of the system, or properties of the noise.

We now work out, for purposes of illustration, the decision rules and their degree of reliability, for several of the above criteria, in the simplest possible problem. We have a linear system in which the voltage is observed at a single instant. We are to decide whether a signal, which can have only amplitude  $S_1$ , is present in noise. We assign a Gaussian pdf for the noise with variance  $\sigma^2$ :

$$W(N) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{N^2}{2\sigma^2} \right\}. \quad (14.35)$$

The likelihood ratio in (14.34) then becomes

$$\frac{p(V|S_1)}{p(V|S_0)} = \frac{W(V - S_1)}{W(V)} = \exp \left\{ \frac{2VS_1 - S_1^2}{2\sigma^2} \right\}, \quad (14.36)$$



and, since this is a monotonic function of  $V$ , the Bayesian decision rule,  $V_b$  can be written as

$$\text{choose } \begin{pmatrix} D_1 \\ D_0 \end{pmatrix} \quad \text{when} \quad V \begin{pmatrix} > \\ < \end{pmatrix} V_b, \quad (14.37)$$

with

$$\frac{V_b}{\sigma} = \frac{1}{2s} \left[ 2 \log \left( \frac{qL_a}{pL_r} \right) + s^2 \right] = v_b, \quad (14.38)$$

in which

$$s \equiv \frac{S_1}{\sigma} \text{ is the voltage signal-to-noise ratio,} \quad (14.39a)$$

and

$$v \equiv \frac{V}{\sigma} \text{ is the normalized voltage.} \quad (14.39b)$$

Now we find the probability for a false rest:

$$\begin{aligned} p(R|X) &= p(D_0 S_1 | X) \\ &= p \sum_V p(D_0 | V) p(V | S_1) \\ &= p \int_{-\infty}^{V_b} dV W(V - S_1) \\ &= p \Phi(v_b - s) \end{aligned} \quad (14.40)$$

and for a false alarm

$$\begin{aligned} p(A|X) &= p(D_1 S_0 | X) \\ &= q \sum_V p(D_1 | V) p(V | S_0) \\ &= q \int_{V_b}^{\infty} dV W(V) \\ &= q [1 - \Phi(v_b)]. \end{aligned} \quad (14.41)$$

Here  $\Phi(x)$  is the cumulative normal distribution function and, as shown in (7.2), it may be computed from an error function:

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x dt \exp\{-t^2/2\} = \frac{1}{2} [1 + \text{erf}(x)]. \quad (14.42)$$

For  $x > 2$ , a good approximation is

$$1 - \Phi(x) \approx \frac{\exp\{-x^2/2\}}{x\sqrt{2\pi}}. \quad (14.43)$$

As a numerical example, if  $L_r = 10L_a$ ,  $q = 10p$ , these expressions reduce to

$$p(A|X) = 10p(R|X) = \frac{10}{11} \left[ 1 - \Phi\left(\frac{s}{2}\right) \right]. \quad (14.44)$$

The probability for a false alarm is less than 0.027, and for a false rest less than 0.0027 for  $s > 4$ . For  $s > 6$ , these numbers become  $1.48 \times 10^{-3}$ ,  $1.48 \times 10^{-4}$ , respectively.

Let us see what the minimax criterion would give in this problem. The conditional losses are

$$\begin{aligned} L(S_0) &= L_a \sum_V p(D_1|V)p(V|S_0) = L_a \int_{-\infty}^{\infty} dV p(D_1|V)W(V), \\ L(S_1) &= L_r \sum_V p(D_0|V)p(V|S_1) = L_r \int_{-\infty}^{\infty} dV p(D_0|V)W(V - S_1). \end{aligned} \quad (14.45)$$

Writing  $f(V) \equiv p(D_1|V) = 1 - p(D_0|V)$ , the only restriction on  $f(V)$  is  $0 \leq f(V) \leq 1$ . Since  $L_a$ ,  $L_r$ , and  $W(V)$  are all positive, a change  $\delta f(V)$  in the neighborhood of any given point  $V$  will always increase one of the quantities in (14.45) and decrease the other. Thus, when the maximum  $L(S)$  has been made as small as possible, we will certainly have  $L(S_0) = L(S_1)$ , and the problem is thus to minimize  $L(S_0)$  subject to this constraint.

Suppose that for some particular  $p(S|X)$  the Bayes decision rule happened to give  $L(S_0) = L(S_1)$ . Then this particular solution must be identical with the minimax solution, for with the above constraint,  $\langle L \rangle = [L(S)]_{\max}$ , and, if the Bayes solution minimizes  $\langle L \rangle$  with respect to all variations  $\delta f(V)$  in the decision rule, it *a fortiori* minimizes it with respect to the smaller class of variations which keep  $L(S_0) = L(S_1)$ . Therefore the decision rule will have the same form as before: there is a minimax threshold  $V_m$  such that

$$f(V) = \begin{cases} 0 & V < V_m \\ 1 & V > V_m. \end{cases} \quad (14.46)$$

Any change in  $V_m$  from the value which makes  $L(S_0) = L(S_1)$  necessarily increases one or the other of these quantities. The equation determining  $V_m$  is therefore

$$L_a \int_{V_m}^{\infty} dV W(V) = L_r \int_{-\infty}^{V_m} dV W(V - S_1), \quad (14.47)$$

or, in terms of normalized quantities,

$$L_a[1 - \Phi(v_m)] = L_r \Phi(v_m - s). \quad (14.48)$$

Note that (14.40) and (14.41) give the conditional probabilities for a false rest and false alarm for any decision rule of type (14.46), regardless of whether the threshold was determined from (14.38) or not; for the arbitrary threshold  $V_0$

$$\begin{aligned} p(R|S_1) &= p(V < V_0|S_1) = \Phi(v_0 - s) \\ p(A|S_0) &= p(V > V_0|S_0) = \frac{1}{2}[1 - \Phi(v_0)]. \end{aligned} \quad (14.49)$$

From (14.38) we see that there is always a particular ratio ( $p/q$ ) which makes the Bayes threshold  $V_b$  equal to the minimax threshold  $V_m$ . For values of ( $p/q$ ) other than this worst value, the Bayes criterion gives a lower expected loss than does the minimax, although one of the conditional losses  $L(S_0)$ ,  $L(S_1)$  will be greater than the minimax value.

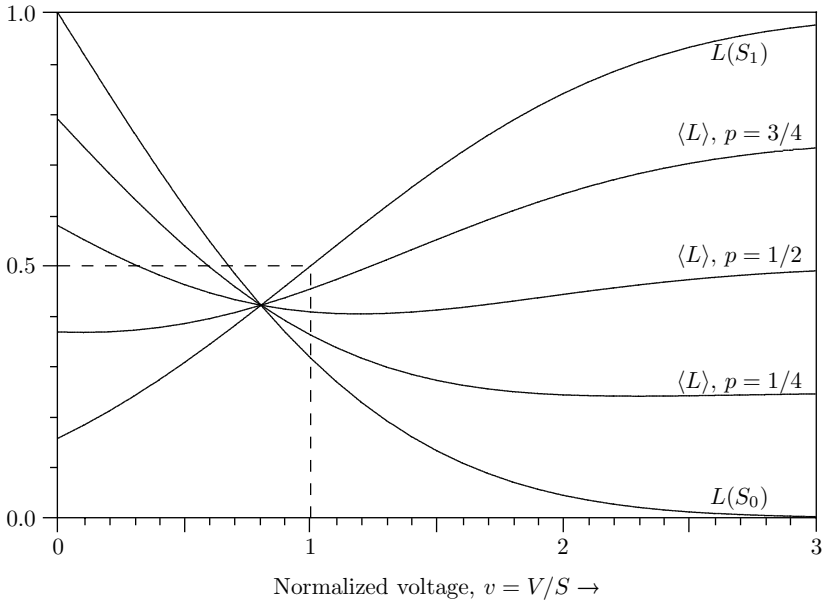


Fig. 14.1. Various risks as a function of voltage for  $L_r = 1$ ,  $L_a = 2$ ,  $p = 1/4, 1/2, 3/4$ .

These relations and several previous remarks are illustrated in Figure 14.1, in which we plot the conditional losses  $L(S_0)$ ,  $L(S_1)$  and the expected loss  $\langle L \rangle$  as functions of the threshold  $V_0$ , for the case  $L_a = (3/2)L_r$ ,  $p = q = 1/2$ . The minimax threshold is at the common crossing point of these curves, while the Bayes threshold occurs at the lowest point of the  $\langle L \rangle$  curve.

One sees how the Bayes threshold moves as the ratio  $(p/q)$  is varied, and in particular that the value of  $(p/q)$  which makes  $V_b = V_m$  also leads to the maximum value of the  $\langle L \rangle_{\min}$  obtained by the Bayes criterion. Thus we could also define a ‘maximin’ criterion; first find the Bayes decision rule which gives minimum  $\langle L \rangle$  for a given  $p(S|X)$ , then vary the prior probability  $p(S|X)$  until the maximum value of  $\langle L \rangle_{\min}$  is attained. The decision rule thus obtained is identical with the one resulting from the minimax criterion; this is the worst possible prior probability, in the sense that the most pessimistic rule is the best that can be done.

The Neyman–Pearson criterion is easily discussed in this example. Suppose the conditional probability for a false alarm  $p(D_1|S_0)$  is held fixed at some value  $\epsilon$ , and we wish to minimize the conditional probability  $p(D_0|S_1)$  of a false rest, subject to this constraint. Now the Bayes criterion minimizes the expected loss

$$\langle L \rangle = pL_r p(D_0|S_1) + qL_a p(D_1|S_0) \quad (14.50)$$

with respect to any variation  $\delta p(D|V)$  in the decision rule. In particular, therefore, it minimizes it with respect to the smaller class of variations which hold  $p(D_1|S_0)$  constant at

the value finally obtained. Thus it minimizes  $p(D_0|S_1)$  with respect to these variations and solves the Neyman–Pearson problem; we need only choose the particular value of the ratio  $(qL_a/pL_r)$  which results in the assumed value of  $\epsilon$  according to (14.38) and (14.41).

We find for the Neyman–Pearson threshold, from (14.49),

$$\Phi(v_{np}) = 1 - \epsilon, \quad (14.51)$$

and the conditional probability for detection is

$$p(D_1|S_1) = 1 - p(D_0|S_1) = \Phi(s - v_{np}). \quad (14.52)$$

If  $\epsilon = 10^{-3}$ , a detection probability of 99% or better is attained for  $s > 6$ .

It is important to note that these numerical examples depend critically on our noise pdf assignment. If we have prior information about the noise beyond its first and second moments, the noise pdf expressing this may not be Gaussian, and the actual situation may be either more or less favorable than indicated by the above relations.

It is well known that in one sense noise with a Gaussian frequency distribution is the worst possible kind; because of its maximum entropy properties, it can obscure a weak signal more completely than can any other noise of the same average power. On the other hand, Gaussian noise is a very favorable kind from which to extract a fairly strong signal, because the probability that the noise will exceed a few times the RMS value  $\sigma = \sqrt{\langle N^2 \rangle}$  becomes vanishingly small. Consequently, the probability for making an incorrect decision on the presence or absence of a signal goes to zero very rapidly as the signal strength is increased. The high reliability of operation found above for  $s > 6$  would not be found for noise possessing a frequency distribution with wider tails.

The type of noise frequency distribution to be expected in any particular case depends, of course, on the physical mechanism which gives rise to the noise. When the noise is the result of a large number of small, independent effects, the Landon derivation of Chapter 7 and the central limit theorem both tell us that a Gaussian frequency distribution for the total noise is by far the most likely to be found, regardless of the nature of the individual sources.

All of these apparently different decision criteria lead to a probability ratio test. In the case of a binary decision, it took the simple form (14.32). Of course, any decision process can be broken down into successive binary decisions, so this case really has the whole story in it. All the different criteria amounted, in the final analysis, only to different philosophies about how you choose the threshold value at which you change your decision.

### 14.5 How would our robot do it?

Now, let's see how this problem appears from the viewpoint of our robot. The rather long arguments we had to go through above (and even they are very highly condensed from the original literature) to obtain the result are due only to the orthodox view which insists on looking at the problem backwards, i.e. on concentrating attention on the final decision rather than on the inference process which logically has to precede it.

To the robot, if our job is to make the best possible decision as to whether the signal is present, the obvious thing we must do is to calculate the *probability* that the signal is present, conditional on all the evidence at hand. If there are only two possibilities,  $S_0$ ,  $S_1$ , to be taken into account, then, after we have seen voltage  $V$ , the posterior odds on  $S_1$  are, from (4.7),

$$O(S_1|VX) = O(S_1|X) \frac{p(V|S_1)}{p(V|S_0)}. \quad (14.53)$$

If we give the robot the loss function (14.31) and ask it to make the decision which minimizes the expected loss, it will evidently use the decision rule

$$\text{choose } D_1 \text{ if } O(S_1|V) = \frac{p(S_1|V)}{p(S_0|V)} > \frac{L_a}{L_r}, \quad (14.54)$$

etc. But from the product rule,  $p(VS_1|X) = p(S_1|V)p(V|X)$ ,  $p(VS_0|X) = p(S_0|V)p(V|X)$ , and (14.54) is identical with (14.32). So, just from looking at this problem the other way around, our robot obtains the same final result in just two lines!

You see that all this discussion of strategies, admissibility, conditional losses, etc. was unnecessary. Except for the introduction of the loss function at the end, there is nothing in the actual functional operation of Wald's decision theory that isn't contained already in basic *probability* theory, if we will only use it in the full generality given to it by Laplace and Jeffreys.

## 14.6 Historical remarks

This comparison shows why the development of decision theory has, more than any other single factor, touched off our 'Bayesian revolution' in statistical thought. For some 50 years, Harold Jeffreys tried valiantly to explain the great advantages of the Laplace methods to statisticians, and his efforts met only with a steady torrent of denials and ridicule. It was then a real irony that the work of one of the most respected of 'orthodox' statisticians (Abraham Wald), which was hailed, very properly, as a great advance in statistical practice, turned out to give, after very long and complicated arguments, exactly the same final results that the Laplace methods give you immediately. Wald showed in great generality what we have just illustrated by one simple example.

The only proper conclusion, as a few recognized at once, is that the supposed distinction between statistical inference and probability theory was entirely artificial – a tragic error of judgment which has wasted perhaps 1000 man-years of our best mathematical talent in the pursuit of false goals.

In the works cited, addressed to electrical engineers, the viewpoint of Middleton and van Meter was that of the Neyman–Pearson and Wald decision theories. At about the same time, Herbert Simon expounded the Neyman–Pearson viewpoint to economists. The writer collaborated with David Middleton for a short time while he was writing his large work, and tried to persuade him of the superiority of the straight Bayesian approach to

decision theory. The success of the effort may be judged by comparing Middleton (1960, Chap. 18) – particularly its length – with our exposition deriving (14.54). It seems that persons with orthodox training had received such strong anti-Bayesian indoctrination that they were locked in an infinite regress situation; although they could not deny the results that Bayesians got on any specific problem, they could never believe that Bayesian methods would work on the next problem until that next solution was also presented to them.

### 14.6.1 The classical matched filter

A funny thing happened in the history of this subject. In the 1930s, electrical engineers knew nothing whatsoever about probability theory; they knew about signal to noise ratios. Receiver input circuits were designed for many years on the basis that signal to noise ratio was maximized by empirical trial and error. Then a general theoretical result was found: if you take the ratio of (peak signal)<sup>2</sup> to mean-square noise, and find, as a variational principle, the design of input stages of the receiver which will maximize it, this turned out to have an analytically neat and useful solution. It is now called the *classical matched filter*, and it has been discovered independently by dozens of people.

To the best of our knowledge, the first person to derive this matched filter solution was the late Professor W. W. Hansen of Stanford University. The writer was working with him, beginning in May 1942, on problems of radar detection. Shortly before then, Hansen had circulated a little memorandum dated 1941, in which he gave this solution for the design of the optimum response curve for the receiver first stage. Years later, I was thinking about an entirely different problem (an optimum antenna pattern for a radar system to maximize the ratio (signal)/(ground clutter response)), and when I finally got the solution, I recognized it as the same result that Bill Hansen had shown me many years before.

Throughout the 1950s, almost every time one opened a journal concerned with these problems, somebody else had a paper announcing the discovery of the same solution. The situation was satirized in a famous editorial by Peter Elias (1958), entitled ‘Two famous papers’. He suggested that it was high time that people stopped rediscovering the easiest solution, and started to think about the many harder problems still in need of solution.

But also in the 1950s people became more sophisticated about the way they handled their detection problems, and they started using this wonderful new tool, statistical decision theory, to see if there were still better ways of handling these design problems. The strange thing happened that in the case of a linear system with Gaussian noise, the optimum solution which decision theory leads you to turns out to be exactly the same old classical matched filter! At first glance, it was very surprising that two approaches so entirely different conceptually should lead to the same solution. But note that our robot represents a viewpoint from which it is obvious that the two lines of argument would have to give the same result.

To our robot, it is obvious that the best analysis we can make of the problem will always be one in which we calculate the *probabilities* that the various signals are present by means of Bayes’ theorem. But let us apply Bayes’ theorem in the logarithmic form of Chapter 4. If we now let  $S_0$  and  $S_1$  stand for numerical values giving the amplitude of two possible

signals as a function of  $V$ , the *evidence* for  $S_1$  is increased by

$$\log \left[ \frac{p(V|S_1)}{p(V|S_0)} \right] = \frac{(V - S_0)^2 - (V - S_1)^2}{2\langle\sigma^2\rangle} = \text{const.} + \frac{(S_1 - S_0)}{\langle\sigma^2\rangle} V. \quad (14.55)$$

In the case of a linear system with Gaussian noise, the observed voltage is itself just a linear function of the posterior probability measured in decibels. So, they are essentially just two different ways of formulating the same problem. Without recognizing it, we had essentially solved this problem already in the Bayesian hypothesis testing discussion of Chapter 4.

In England, P. M. Woodward had perceived much of this correctly in the 1940s – but he was many years ahead of his time. Those with conventional statistical training were unable to see any merit in his work, and simply ignored it. His book (Woodward, 1953) is highly recommended reading; although it does not solve any of our current problems, its thinking is still in advance of some current literature and practice.

We have seen that the other non-Bayesian approaches to the theory all amounted to different philosophies of how you choose the threshold at which you change your decision. Because of the fact that they all lead to the same probability ratio test, they must necessarily all be derivable from Bayes' theorem.

The problem just examined by several different decision criteria is, of course, the simplest possible one. In a more realistic problem, we will observe the voltage  $V(t)$  as a function of time, perhaps several voltages  $V_1(t)$ ,  $V_2(t)$ , ... in several different channels. We may have many different possible signals  $S_a(t)$ ,  $S_b(t)$ , ... to distinguish and correspondingly many possible decisions. We may need to decide not only *whether* a given signal is present, but also to make the best estimates of one or more signal parameters (such as intensity, starting time, frequency, phase, rate of frequency modulation, etc.). Therefore, just as in the problem of quality control discussed in Chapter 4, the details can become arbitrarily complicated. But these extensions are, from the Bayesian viewpoint, straightforward in that they require no new principles beyond those already given, only mathematical generalization.

We shall return to some of these more complicated problems of detection and filtering when we take up frequency/shape estimation; but for now let's look at another elementary kind of decision problem. In the ones just discussed, we needed Bayes' theorem, but not maximum entropy. Now we examine a kind of decision problem where we need maximum entropy, but not Bayes' theorem.

### 14.7 The widget problem

This problem was first propounded at a symposium held at Purdue University in November, 1960 – at which time, however, the full solution was not known. This was worked out later (Jaynes, 1963c), and some numerical approximations were improved in the computer work of Tribus and Fitts (1968).

The widget problem has proved to be interesting in more respects than originally realized. It is a decision problem in which there is no occasion to use Bayes' theorem, because no 'new' information is acquired. Thus it would be termed a 'no data' decision problem in the

sense of Chernoff and Moses (1959). However, at successive stages of the problem we have more and more prior information; and digesting it by maximum entropy leads to a sequence of prior probability assignments, which lead to different decisions. Thus it is an example of the ‘pure’ use of maximum entropy, as in statistical mechanics. It is hard to see how the problem could be formulated mathematically at all without use of maximum entropy, or some other device (such as the method of Darwin and Fowler (Fowler, 1929) in statistical mechanics, or the ‘method of the most probable distribution’ dating back to Boltzmann (1871)) which turns out in the end to be mathematically equivalent to maximum entropy.

The problem is interesting also in that we can see a continuous gradation from decision problems so simple that common sense tells us the answer instantly, with no need for any mathematical theory, through problems more and more involved so that common sense has more and more difficulty in making a decision, until finally we reach a point where nobody has yet claimed to be able to see the right decision intuitively, and we require the mathematics to tell us what to do.

Finally, the widget problem turns out to be very close to an important real problem faced by oil prospectors. The details of the real problem are shrouded in proprietary caution; but it is not giving away any secrets to report that, a few years ago, the writer spent a week at the research laboratories of one of our large oil companies, lecturing for over 20 hours on the widget problem. We went through every part of the calculation in excruciating detail – with a room full of engineers armed with calculators, checking up on every stage of the numerical work.

Here is the problem: Mr A is in charge of a widget factory, which proudly advertises that it can make delivery in 24 hours on any size order. This, of course, is not really true, and Mr A’s job is to protect, as best he can, the advertising manager’s reputation for veracity. This means that each morning he must decide whether the day’s run of 200 widgets will be painted red, yellow or green. (For complex technological reasons, not relevant to the present problem, only one color can be produced per day.) We follow his problem of decision through several stages of increasing knowledge.

### *Stage 1*

When he arrives at work, Mr A checks with the stock room and finds that they now have in stock 100 red widgets, 150 yellow, and 50 green. His ignorance lies in the fact that he does not know how many orders for each type will come in during the day. Clearly, in this state of ignorance, Mr A will attach the highest significance to any tiny scrap of information about orders likely to come in today; and if no such scraps are to be had, we do not envy Mr A his job. Still, if a decision must be made here and now on no more information than this, his common sense will probably tell him that he had better build up that stock of green widgets.

### *Stage 2*

Mr A, feeling the need for more information, calls up the front office and asks, ‘Can you give me some idea of how many orders for red, yellow, and green widgets are likely to come



Table 14.1. *Summary of four stages of the widget problem.*

Stage	R	Y	G	Decision
1. In stock	100	150	50	G
2. Av. daily order total	50	100	10	Y
3. Av. individual order	75	10	20	R
4. Specific order			40	?

in today?’ They reply, ‘Well, we don’t have the breakdown of what has been happening each day, and it would take us a week to compile that information from our files. But we do have a summary of the total sales last year. Over the last year, we sold a total of 13 000 red, 26 000 yellow, and 2600 green. Figuring 260 working days, this means that last year we sold an average of 50 red, 100 yellow, and 10 green each day.’ If Mr A ponders this new information for a few seconds, I think he will change his mind, and decide to make yellow ones today.

*Stage 3*

The man in the front office calls Mr A back and says, ‘It just occurred to me that we do have a little more information that might possibly help you. We have at hand not only the total number of widgets sold last year, but also the total number of orders we processed. Last year we got a total of 173 orders for red, 2600 for yellow, and 130 for green. This means that the customers who use red widgets order, on the average,  $13\,000/173 = 75$  widgets per order, while the average order for yellow and green were  $26\,000/2600 = 10$ , and  $2600/130 = 20$ , respectively.’ These new data do not change the expected daily demand; but, if Mr A is very shrewd and ponders it very hard, I think he may change his mind again, and decide to make red ones today.

*Stage 4*

Mr A is just about to give the order to make red widgets when the front office calls him again to say, ‘We just got word that a messenger is on his way here with an emergency order for 40 green widgets.’ Now, what should he do? Up to this point, Mr A’s decision problem has been simple enough so that reasonably good common sense will tell him what to do. But now he is in trouble; qualitative common sense is just not powerful enough to solve his problem, and he needs a mathematical theory to determine a definite optimum decision.

Let’s summarize all the above data in Table 14.1. In the final column, we give the decision that seemed intuitively to be the best one before we had worked out the mathematics. Do other people agree with this intuitive judgment? Professor Myron Tribus has put this to a test by giving talks about this problem, and taking votes from the audience before the solution is given. We quote his findings as given in Tribus and Fitts

(1968). They use  $D_1, D_2, D_3, D_4$  to stand for the optimum decisions in stages 1, 2, 3, 4, respectively:

Before taking up the formal solution, it may be reported that Jaynes' widget problem has been presented to many gatherings of engineers who have been asked to vote on  $D_1, D_2, D_3, D_4$ . There is almost unanimous agreement about  $D_1$ . There is about 85% agreement on  $D_2$ . There is about 70% agreement on  $D_3$ , and almost no agreement on  $D_4$ . One conclusion stands out from these informal tests; the average engineer has remarkably good intuition in problems of this kind. The majority vote for  $D_1, D_2$ , and  $D_3$  has always been in agreement with the formal mathematical solution. However, there has been almost universal disagreement over how to defend the intuitive solution. That is, while many engineers could agree on the best course of action, they were much less in agreement on *why* that course was the best one.

### 14.7.1 Solution for Stage 2

Now, how are we to set up this problem mathematically? In a real life situation, evidently, the problem would be a little more complicated than indicated so far, because what Mr A does today also affects how serious his problem will be tomorrow. That would get us into the subject of dynamic programming. But for now, just to keep the problem simple, we shall solve only the truncated problem in which he makes decisions on a day to day basis with no thought of tomorrow.

We have just to carry out the steps enumerated in Section 13.11. Since Stage 1 is almost too trivial to work with, consider the problem of Stage 2. Firstly, we define our underlying hypothesis space by enumerating the possible 'states of nature'  $\theta_j$  that we will consider. These correspond to all possible order situations that could arise; if Mr A knew in advance exactly how many red, yellow, and green widgets would be ordered today, his decision problem would be trivial. Let  $n_1 = 0, 1, 2, \dots$  be the number of red widgets that will be ordered today, and similarly  $n_2, n_3$  for yellow and green, respectively. Then, any conceivable order situation is given by specifying three non-negative integers  $\{n_1, n_2, n_3\}$ . Conversely, every ordered triple of non-negative integers represents a conceivable order situation.

Next, we are to assign prior probabilities  $p(\theta_j|X) = p(n_1 n_2 n_3|X)$  to the states of nature, which maximize the entropy of the distribution subject to the constraints of our prior knowledge. We solved this problem in general in Chapter 11, Eqs. (11.39)–(11.50); and so we just have to translate the result into our present notation. The index  $i$  on  $x_i$  in Chapter 11 now corresponds to the three integers  $n_1, n_2, n_3$ ; the function  $f_k(x_i)$  also corresponds to the  $n_i$ , since the prior information at this stage will be used to fix the expectations  $\langle n_1 \rangle, \langle n_2 \rangle, \langle n_3 \rangle$  of orders for red, yellow, and green widgets at 50, 100, 10, respectively. With three constraints we will have three Lagrange multipliers  $\lambda_1, \lambda_2, \lambda_3$ , and the partition function (11.45) becomes

$$\begin{aligned} Z(\lambda_1, \lambda_2, \lambda_3) &= \sum_{n_1=0}^{\infty} \sum_{n_2=0}^{\infty} \sum_{n_3=0}^{\infty} \exp\{-\lambda_1 n_1 - \lambda_2 n_2 - \lambda_3 n_3\} \\ &= \prod_{i=1}^3 (1 - \exp\{-\lambda_i\})^{-1}. \end{aligned} \quad (14.56)$$

The  $\lambda_i$  are determined from (11.46):

$$\langle n_i \rangle = -\frac{\partial \log(Z)}{\partial \lambda_i} = \frac{1}{\exp\{\lambda_i\} - 1}. \quad (14.57)$$

The maximum entropy probability assignment (11.41) for the states of nature  $\theta_j = \{n_1 n_2 n_3\}$  therefore factors:

$$p(n_1 n_2 n_3) = p_1(n_1) p_2(n_2) p_3(n_3) \quad (14.58)$$

with

$$\begin{aligned} p_i(n_i) &= (1 - \exp\{-\lambda_i\}) \exp\{-\lambda_i n_i\}, \quad n_i = 1, 2, 3 \dots \\ &= \frac{1}{\langle n_i \rangle + 1} \left[ \frac{\langle n_i \rangle}{\langle n_i \rangle + 1} \right]^{n_i}. \end{aligned} \quad (14.59)$$

Thus, in Stage 2, Mr A's state of knowledge about today's orders is given by three exponential distributions:

$$p_1(n_1) = \frac{1}{51} \left( \frac{50}{51} \right)^{n_1}, \quad p_2(n_2) = \frac{1}{101} \left( \frac{100}{101} \right)^{n_2}, \quad p_3(n_3) = \frac{1}{11} \left( \frac{10}{11} \right)^{n_3}. \quad (14.60)$$

Applications of Bayes' theorem to digest new evidence  $E$  is absent because there is no new evidence. Therefore, the decision must be made directly from the prior probabilities (14.60), as is always the case in statistical mechanics.

So, we now proceed to enumerate the possible decisions. These are  $D_1 \equiv$  make red ones today,  $D_2 \equiv$  make yellow ones,  $D_3 \equiv$  make green ones, for which we are to introduce a loss function  $L(D_i, \theta_j)$ . Mr A's judgment is that there is no loss if all orders are filled today; otherwise, the loss will be proportional to – and in view of the invariance of the decision rule under proper linear transformations that we noted at the end of Chapter 13, we may as well take it equal to – the total number of unfilled orders.

The present stock of red, yellow, and green widgets is  $S_1 = 100$ ,  $S_2 = 150$ ,  $S_3 = 50$ , respectively. On decision  $D_1$  (make red widgets) the available stock  $S_1$  will be increased by the day's run of 200 widgets, and the loss will be

$$L(D_1; n_1, n_2, n_3) = R(n_1 - S_1 - 200) + R(n_2 - S_2) + R(n_3 - S_3), \quad (14.61)$$

where  $R(x)$  is the ramp function

$$R(x) \equiv \begin{cases} x & x \geq 0 \\ 0 & x \leq 0. \end{cases} \quad (14.62)$$

Likewise, on decisions  $D_2$ ,  $D_3$ , the loss will be

$$L(D_2; n_1, n_2, n_3) = R(n_1 - S_1) + R(n_2 - S_2 - 200) + R(n_3 - S_3), \quad (14.63)$$

$$L(D_3; n_1, n_2, n_3) = R(n_1 - S_1) + R(n_2 - S_2) + R(n_3 - S_3 - 200). \quad (14.64)$$

So, if decision  $D_1$  is made, the expected loss will be

$$\begin{aligned}
 \langle L \rangle_1 &= \sum_{n_i} p(n_1 \ n_2 \ n_3) L(D_1; n_1, n_2, n_3) \\
 &= \sum_{n_1=0}^{\infty} p_1(n_1) R(n_1 - S_1 - 200) \\
 &\quad + \sum_{n_2=0}^{\infty} p_2(n_2) R(n_2 - S_2) \\
 &\quad + \sum_{n_3=0}^{\infty} p_3(n_3) R(n_3 - S_3),
 \end{aligned} \tag{14.65}$$

and similarly for  $D_2, D_3$ . The summations are elementary, giving

$$\begin{aligned}
 \langle L \rangle_1 &= \langle n_1 \rangle \exp\{-\lambda_1(S_1 + 200)\} + \langle n_2 \rangle \exp\{-\lambda_2 S_2\} + \langle n_3 \rangle \exp\{-\lambda_3 S_3\}, \\
 \langle L \rangle_2 &= \langle n_1 \rangle \exp\{-\lambda_1 S_1\} + \langle n_2 \rangle \exp\{-\lambda_2(S_2 + 200)\} + \langle n_3 \rangle \exp\{-\lambda_3 S_3\}, \\
 \langle L \rangle_3 &= \langle n_1 \rangle \exp\{-\lambda_1 S_1\} + \langle n_2 \rangle \exp\{-\lambda_2 S_2\} + \langle n_3 \rangle \exp\{-\lambda_3(S_3 + 200)\},
 \end{aligned} \tag{14.66}$$

or, inserting numerical values,

$$\begin{aligned}
 \langle L \rangle_1 &= 0.131 + 22.48 + 0.085 = 22.70, \\
 \langle L \rangle_2 &= 6.902 + 3.073 + 0.085 = 10.6, \\
 \langle L \rangle_3 &= 6.902 + 22.48 + 4 \times 10^{-10} = 29.38,
 \end{aligned} \tag{14.67}$$

showing a strong preference for decision  $D_2 \equiv$  ‘make yellow ones today’, as common sense had already anticipated.

Physicists will recognize that Stage 2 of Mr A’s decision problem is mathematically the same as the theory of harmonic oscillators in quantum statistical mechanics. There is still another engineering application of the harmonic oscillator equations, in some problems of message encoding, to be noted when we take up communication theory. We are trying to emphasize the generality of this theory, which is mathematically quite old and well-known, but which has been applied in the past only in some specialized problems in thermodynamics. This general applicability can be seen only after we are emancipated from the orthodox view of probability.

### 14.7.2 Solution for Stage 3

In Stage 3 of Mr A’s problem we have some additional pieces of information: the average *individual* orders for red, yellow, and green widgets. To take account of this new information, we need to go down into a deeper hypothesis space; set up a more detailed enumeration of the states of nature in which we take into account not only the total orders for each type, but also the breakdown into individual orders. We could have done this also in Stage 2, but since at that stage there was no information available bearing on this breakdown, it would have added nothing to the problem (the subtle difference that this makes after all will be noted later).

In Stage 3, a possible state of nature can be described as follows. We receive  $u_1$  individual orders for one red widget each,  $u_2$  orders for two red widgets each,  $\dots$ ,  $u_r$  individual orders for  $r$  red widgets each. Also, we receive  $v_y$  orders for  $y$  yellow widgets each, and  $w_g$  orders for  $g$  green widgets each. Thus, a state of nature is specified by an infinite number of non-negative integers

$$\theta = \{u_1 \dots; v_1 \dots; w_1 \dots\}, \quad (14.68)$$

and conversely every such set of integers represents a conceivable state of nature, to which we assign a probability  $p(u_1 \dots; v_1 \dots; w_1 \dots)$ .

Today's total demands for red, yellow, and green widgets are, respectively,

$$n_1 = \sum_{r=1}^{\infty} r u_r, \quad n_2 = \sum_{y=1}^{\infty} y v_y, \quad n_3 = \sum_{g=1}^{\infty} g w_g, \quad (14.69)$$

the expectations of which were given in Stage 2 as  $\langle n_1 \rangle = 50$ ,  $\langle n_2 \rangle = 100$ ,  $\langle n_3 \rangle = 10$ . The total number of individual orders for red, yellow, and green widgets are, respectively,

$$m_1 = \sum_{r=1}^{\infty} u_r, \quad m_2 = \sum_{y=1}^{\infty} v_y, \quad m_3 = \sum_{g=1}^{\infty} w_g, \quad (14.70)$$

and the new feature of Stage 3 is that  $\langle m_1 \rangle$ ,  $\langle m_2 \rangle$ ,  $\langle m_3 \rangle$  are also known. For example, the statement that the average individual order for red widgets is 75 means that  $\langle n_1 \rangle = 75 \langle m_1 \rangle$ .

With six average values given, we will have six Lagrange multipliers  $\{\lambda_1, \mu_1; \lambda_2, \mu_2; \lambda_3, \mu_3\}$ . The maximum entropy probability assignment will have the form

$$p(u_1 \dots; v_1 \dots; w_1 \dots) = \exp \{-\lambda_0 - \lambda_1 n_1 - \mu_1 m_1 - \lambda_2 n_2 - \mu_2 m_2 - \lambda_3 n_3 - \mu_3 m_3\}, \quad (14.71)$$

which factors:

$$p(u_1 \dots; v_1 \dots; w_1 \dots) = p_1(u_1 \dots) p_2(v_1 \dots) p_3(w_1 \dots). \quad (14.72)$$

The partition function also factors:

$$Z = Z_1(\lambda_1 \mu_1) Z_2(\lambda_2 \mu_2) Z_3(\lambda_3 \mu_3), \quad (14.73)$$

with

$$\begin{aligned} Z_1(\lambda_1 \mu_1) &= \sum_{u_1=1}^{\infty} \sum_{u_2=1}^{\infty} \dots \exp \{-\lambda_1(u_1 + 2u_2 + 3u_3 + \dots) - \mu_1(u_1 + u_2 + u_3 + \dots)\} \\ &= \prod_{r=1}^{\infty} \frac{1}{1 - \exp\{-r\lambda_1 - \mu\}} \end{aligned} \quad (14.74)$$

and similar expressions for  $Z_2, Z_3$ . To find  $\lambda_1, \mu_1$  we apply the general rule, Eq. (14.57):

$$\langle n_1 \rangle = \frac{\partial}{\partial \lambda_1} \sum_{r=1}^{\infty} \log(1 - \exp\{-r\lambda_1 - \mu_1\}) = \sum_{r=1}^{\infty} \frac{r}{\exp\{r\lambda_1 + \mu_1\} - 1}, \quad (14.75)$$

$$\langle m_1 \rangle = \frac{\partial}{\partial \mu_1} \sum_{r=1}^{\infty} \log(1 - \exp\{-r\lambda_1 - \mu_1\}) = \sum_{r=1}^{\infty} \frac{1}{\exp\{r\lambda_1 + \mu_1\} - 1}. \quad (14.76)$$

Combining with (14.69) and (14.70), we see that

$$\langle u_r \rangle = \frac{1}{\exp\{r\lambda_1 + \mu_1\} - 1}, \quad (14.77)$$

and now the secret is out – Stage 3 of Mr A's decision problem is just the theory of the ideal Bose–Einstein gas in quantum statistical mechanics!

If we treat the ideal Bose–Einstein gas by the method of the Gibbs grand canonical ensemble, we obtain just these equations, in which the number  $r$  corresponds to the  $r$ th single-particle energy level,  $u_r$  to the number of particles in the  $r$ th state, and  $\lambda_1$  and  $\mu_1$  to the temperature and chemical potential.

In the present problem it is clear that for all  $r$ ,  $\langle u_r \rangle \ll 1$ , and that  $\langle u_r \rangle$  cannot decrease appreciably below  $\langle u_1 \rangle$  until  $r$  is of the order of 75, the average individual order. Therefore,  $\mu_1$  will be numerically large, and  $\lambda_1$  numerically small, compared with unity. This means that the series (14.75), (14.76) converge very slowly and are useless for numerical work unless you write a computer program to do it. However, we can do it analytically if we transform them into rapidly converging sums as follows:

$$\begin{aligned} \sum_{r=1}^{\infty} \frac{1}{\exp\{\lambda r + \mu\} - 1} &= \sum_{r=1}^{\infty} \sum_{n=1}^{\infty} \exp\{-n(\lambda r + \mu)\} \\ &= \sum_{n=1}^{\infty} \frac{\exp\{-n\mu\}}{1 - \exp\{-n\lambda\}}. \end{aligned} \quad (14.78)$$

The first term is already an excellent approximation. Similarly,

$$\sum_{r=1}^{\infty} \frac{r}{\exp\{\lambda r + \mu\} - 1} = \sum_{n=1}^{\infty} \frac{\exp\{-n(\lambda + \mu)\}}{(1 - \exp\{-n\lambda\})^2}, \quad (14.79)$$

and so (14.75) and (14.76) become

$$\langle n_1 \rangle = \frac{\exp\{-\mu_1\}}{\lambda_1^2}, \quad (14.80)$$

$$\langle m_1 \rangle = \frac{\exp\{-\mu_1\}}{\lambda_1}, \quad (14.81)$$

$$\lambda_1 = \frac{\langle m_1 \rangle_1}{\langle n_1 \rangle} = \frac{1}{75} = 0.0133, \quad (14.82)$$

$$\exp\{\mu_1\} = \frac{\langle n_1 \rangle_1}{\langle m_1 \rangle} = 112.5, \quad (14.83)$$

$$\mu_1 = 4.722. \quad (14.84)$$

Tribus and Fitts, evaluating the sums by computer, obtain  $\lambda_1 = 0.0131$ ,  $\mu_1 = 4.727$ ; so our approximations (14.80), (14.81) are very good, at least in the case of red widgets.

The probability that  $u_r$  has a particular value is, from (14.72) or (14.74),

$$p(u_r) = (1 - \exp\{-r\lambda_1 - \mu\}) \exp\{(-r\lambda_1 + \mu_1)u_r\}, \quad (14.85)$$

which has the mean value (14.77) and the variance

$$\text{var}(u_r) = \langle u_r^2 \rangle - \langle u_r \rangle^2 = \frac{\exp\{r\lambda_1 + \mu_1\}}{\exp\{r\lambda_1 + \mu_1\} - 1}. \quad (14.86)$$

The total demand for red widgets

$$n_1 = \sum_{r=1}^{\infty} r u_r \quad (14.87)$$

is expressed as the sum of a large number of independent terms. The pdf for  $n_1$  will have the mean value (14.80) and the variance

$$\text{var}(n_1) = \sum_{r=1}^{\infty} r^2 \text{var}(u_r) = \sum_{r=1}^{\infty} \frac{r^2 \exp\{r\lambda_1 + \mu_1\}}{(\exp\{r\lambda_1 + \mu_1\} - 1)^2}, \quad (14.88)$$

which we convert into the rapidly convergent sum

$$\sum_{r,n=1}^{\infty} n r^2 \exp\{-n(r\lambda + \mu)\} = \sum_{n=1}^{\infty} n \frac{\exp\{-n(\lambda + \mu)\} + \exp\{-n(2\lambda + \mu)\}}{(1 - \exp\{-n\lambda\})^3} \quad (14.89)$$

or, approximately,

$$\text{var}(n_1) = \frac{2 \exp\{-\mu_1\}}{\lambda_1^3} = \frac{2}{\lambda_1} \langle n_1 \rangle. \quad (14.90)$$

At this point we can use some mathematical facts concerning the central limit theorem. Because  $n_1$  is the sum of a large number of small terms to which we have assigned independent probabilities, our probability distribution for  $n_1$  will be very nearly Gaussian:

$$p(n_1) \approx A \exp \left\{ -\frac{\lambda_1 (n_1 - \langle n_1 \rangle)^2}{4 \langle n_1 \rangle} \right\} \quad (14.91)$$

for those values of  $n_1$  which can arise in many different ways. For example, the case  $n = 2$  can arise in only two ways:  $u_1 = 2$ , or  $u_2 = 1$ , all others  $u_k$  being zero. On the other hand, the case  $n_1 = 150$  can arise in an enormous number of different ways, and the ‘smoothing’ mechanism of the central limit theorem can operate. Thus, Eq. (14.91) will be a good approximation for the large values of  $n_1$  of interest to us, but not for small  $n_1$ .

The expected loss on the various decisions is, as we saw in (14.66), the sum of three terms arising from failure to meet orders for red, yellow, or green widgets, respectively. If we do not make red ones today, then the possibility of failing to meet orders for red widgets contributes to the expected loss the amount

$$\begin{aligned} \sum_{n_1=0}^{\infty} p(n_1)R(n_1 - S_1) &\simeq \sqrt{\left[\frac{\lambda_1}{4\pi\langle n_1 \rangle}\right]} \int_{S_1}^{\infty} dn_1 (n_1 - S_1) \exp\left\{-\frac{\lambda_1(n_1 - \langle n_1 \rangle)^2}{4\langle n_1 \rangle}\right\} \\ &= (\langle n_1 \rangle - S_1)\Phi\left[\alpha_1\sqrt{2}(\langle n_1 \rangle - S_1)\right] \\ &\quad + \frac{1}{2\alpha_1\sqrt{\pi}} \exp\left\{-\alpha_1^2(\langle n_1 \rangle - S_1)^2\right\}, \end{aligned} \quad (14.92)$$

where  $\alpha_1^2 = \lambda_1/4\langle n_1 \rangle$ , and  $\Phi(x)$  is the cumulative normal distribution function (14.42).

If we do decide to make red widgets today, the possibility of failing to meet red orders contributes to the expected loss the above expression (14.92) with  $S_1$  replaced by  $(S_1 + 200)$ .

Similar equations hold for yellow and green widgets. Although the approximations we made are not equally good in all cases, let us use (14.92) for the partial losses and apply it three times with the given numerical values

$$\begin{array}{lll} S_1 = 100, & S_2 = 150, & S_3 = 50, \\ \langle n_1 \rangle = 50, & \langle n_2 \rangle = 100, & \langle n_3 \rangle = 10, \\ \alpha_1 = 0.0082, & \alpha_2 = 0.0160, & \alpha_3 = 0.035. \end{array} \quad (14.93)$$

Doing the indicated calculations, we find that on the decisions  $D_1, D_2, D_3$  the expected losses are

$$\begin{aligned} \langle L \rangle_1 &= (0) + 2.86 + 0.18 = 3.04 \text{ unfilled orders} \\ \langle L \rangle_2 &= 14.9 + (0) + 0.18 = 15.1 \text{ unfilled orders} \\ \langle L \rangle_3 &= 14.9 + 2.86 + (0) = 17.8 \text{ unfilled orders} \end{aligned} \quad (14.94)$$

where (0) stands for a term orders of magnitude smaller than the others. The breakdown indicated is to be read as follows. If Decision  $D_1$  (make red widgets) is made, there is negligible loss from the possibility of failing to meet red orders, while the possibility of failure with yellow orders contributes an expected loss of 2.86, and only 0.18 for green.

These results show the great preference for  $D_1$  caused by the additional information about average individual orders, which had the intuitive effect of making the situation with respect to yellow widgets much safer than it seemed in Stage 2.

### 14.7.3 Solution for Stage 4

It is in the passage from Stage 3 to Stage 4 (where the new information consists of a specific order for 40 green widgets) that our common sense first fails us. Now both the red and green situations seem rather precarious, and our common sense lacks the 'resolving power' to tell which is the more serious. Strangely enough, this new knowledge, which makes the



problem so hard for our common sense, causes no difficulty at all in the mathematics. The previous equations still apply, with the sole difference that the stock  $S_3$  of green widgets is reduced from 50 to 10. We now have  $(\langle n_3 \rangle - S_3) = 0$  so that (14.92) reduces to

$$\frac{1}{2\alpha_3\sqrt{\pi}} = 8.08, \quad (14.95)$$

and in place of (14.94) we have

$$\begin{aligned} \langle L \rangle_1 &= (0) + 2.86 + 8.08 = 10.9 \text{ unfilled orders} \\ \langle L \rangle_2 &= 14.9 + (0) + 8.08 = 23.0 \text{ unfilled orders} \\ \langle L \rangle_3 &= 14.9 + 2.86 + (0) = 17.8 \text{ unfilled orders.} \end{aligned} \quad (14.96)$$

So, Mr A should stick to his decision to make red widgets! Our common sense fails just because there is now so little difference between  $\langle L \rangle_1$  and  $\langle L \rangle_3$ .

### 14.8 Comments

We have tried to show that use of probability theory in the sense of Laplace, with prior probabilities determined by the principle of maximum entropy, leads to a reasonable method of treating decision problems and to results in good correspondence with common sense. Mathematically, our equations are nothing but the Gibbs formalism in statistical mechanics, the only new feature being the recognition that the Gibbs methods are of far more general applicability than had been supposed.

The moral of this is simply that questions about ‘interpretation of a formalism’, which the positivist philosophy tends to reject as meaningless and useless, are, on the contrary, of central importance in scientific work. It is, of course, true that, in an application already established, a different interpretation of the equations cannot lead to any new numerical results. But our judgment as to the *range of validity* of a formalism can depend entirely on how we interpret it. The interpretation (probability)  $\equiv$  (frequency) has led to a great and unnecessary restriction on the kinds of problem where probability theory can be applied. Today, the scientist, engineer, and economist face many problems which require the broader Laplace–Jeffreys interpretation.