

## CHAPTER 19

### PHYSICAL MEASUREMENTS

We have seen, in Chapter 7, how the great mathematician Leonhard Euler was unable to solve the problem of estimating 8 orbital parameters from 75 discrepant observations of the past positions of Jupiter and Saturn. Thinking in terms of deductive logic, he could not even conceive of the principles by which such a problem could be solved. But 38 years later Laplace, thinking in terms of probability theory as logic, was in possession of exactly the right principles to resolve the Great Inequality of Jupiter and Saturn. In this Chapter we develop the solution as it would be done today by considering a simpler problem, estimating two parameters from three observations. But our general solution, in matrix notation, will include Laplace's automatically.

#### Reduction of Equations of Condition

Suppose we wish to determine the charge  $e$  and mass  $m$  of the electron. The Millikan oil-drop experiment measures  $e$  directly. The deflection of an electron beam in a known electromagnetic field measures the ratio  $e/m$ . The deflection of an electron toward a metal plate due to attraction of image charges measures  $e^2/m$ .

From the results of any two of these experiments we can calculate values of  $e$  and  $m$ . But all the measurements are subject to error, and the values of  $e$ ,  $m$  obtained from different experiments will not agree. Yet each of the measurements does contain some information relevant to our question, that is not contained in the others. Then how are we to process the data so as to make use of all the information available and get the best estimates of  $e$ ,  $m$ ? What is the probable error remaining? How much would the situation be improved by including still another experiment of given accuracy? Probability theory gives simple and elegant answers to these questions.

More specifically, suppose we have the results of these experiments:

- (1) measures  $e$  with  $\pm 2\%$  accuracy
- (2) measures  $(e/m)$  with  $\pm 1\%$  accuracy
- (3) measures  $(e^2/m)$  with  $\pm 5\%$  accuracy

Supposing the values of  $e$ ,  $m$  approximately known in advance,  $e \approx e_0$ ,  $m \approx m_0$ , the measurements are then linear functions of the corrections. Write the unknown true values of  $e$  and  $m$  as

$$\begin{aligned} e &= e_0(1 + x_1) \\ m &= m_0(1 + x_2) \end{aligned} \tag{19-1}$$

then  $x_1$ ,  $x_2$  are dimensionless corrections, small compared to unity, and our problem is to find the best estimates of  $x_1$ ,  $x_2$ . The results of the three measurements are three numbers  $M_1$ ,  $M_2$ ,  $M_3$  which we write as

$$\begin{aligned} M_1 &= e_0(1 + y_1) \\ M_2 &= \frac{e_0}{m_0}(1 + y_2) \\ M_3 &= \frac{e_0^2}{m_0}(1 + y_3) \end{aligned} \tag{19-2}$$

where the  $y_i$  are also small dimensionless numbers which are defined by (19-2) and are therefore known in terms of the old estimates  $e_0$ ,  $m_0$  and the new measurements  $M_1$ ,  $M_2$ ,  $M_3$ . On the other hand, the true values of  $e$ ,  $e/m$ ,  $e^2/m$  are expressible in terms of the  $x_i$ :

$$\begin{aligned} e &= e_0(1 + x_1) \\ \frac{e}{m} &= \frac{e_0(1 + x_1)}{m_0(1 + x_2)} = \frac{e_0}{m_0}(1 + x_1 - x_2 + \dots) \\ \frac{e^2}{m} &= \frac{e_0^2(1 + x_1)^2}{m_0(1 + x_2)} = \frac{e_0}{m_0}(1 + 2x_1 - x_2 + \dots) \end{aligned} \quad (19-3)$$

where higher order terms are considered negligible. Comparing (19-2) and (19-3) we see that if the measurements were exact we would have

$$\begin{aligned} y_1 &= x_1 \\ y_2 &= x_1 - x_2 \\ y_3 &= 2x_1 - x_2 \end{aligned} \quad (19-4)$$

But taking into account the errors, the known  $y_i$  are related to the unknown  $x_j$  by

$$\begin{aligned} y_1 &= a_{11}x_1 + a_{12}x_2 + \delta_1 \\ y_2 &= a_{21}x_1 + a_{22}x_2 + \delta_2 \\ y_3 &= a_{31}x_1 + a_{32}x_2 + \delta_3 \end{aligned} \quad (19-5)$$

where the coefficients  $a_{ij}$  form a  $(3 \times 2)$  matrix:

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 1 & -1 \\ 2 & -1 \end{pmatrix} \quad (19-6)$$

and the  $\delta_i$  are the unknown fractional errors of the three measurements. For example, the statement that  $\delta_2 = -0.01$  means that the second measurement gave a result one per cent too small.

More generally, we have  $n$  unknown quantities  $\{x_1 \dots x_n\}$  to be estimated from  $N$  imperfect observations  $\{y_1 \dots y_N\}$ , and the  $N$  "equations of condition:"

$$y_i = \sum_{j=1}^n a_{ij}x_j + \delta_i, \quad i = 1, 2, \dots, N \quad (19-7)$$

or, in matrix notation,

$$y = Ax + \delta \quad (19-8)$$

where  $A$  is an  $(N \times n)$  matrix. In the present discussion we suppose the problem "overdetermined" in the sense that  $N > n$ . This condition defeated Euler (1749), who was facing the case  $N = 75$ ,  $n = 8$ . But we keep in mind that the cases  $N = n$  (ostensibly well-posed) and  $N < n$  (underdetermined) can also arise in real problems, and it will be interesting to see what probability theory has to say about those cases.

In the early 19'th Century, it was common to reason as follows. It seems plausible that the best estimate of each  $x_j$  will be some linear combination of all the  $y_i$ , but if  $N > n$  we cannot simply solve equation (19-8) for  $x$ , since  $A$  is not a square matrix and has no inverse. However, we can get a system of equations solvable for  $x$  if we take  $n$  linear combinations of the equations of condition; *i.e.*, if we multiply (19-8) on the left by some  $(n \times N)$  matrix  $B$ . Then the product  $BA$  exists and

is a square ( $n \times n$ ) matrix. Choose  $B$  so that  $(BA)^{-1}$  exists. Then the linear combinations are the  $n$  rows of

$$By = BAx + B\delta, \quad (19-9)$$

which has the unique solution

$$x = (BA)^{-1}B(y - \delta). \quad (19-10)$$

If the probabilities of various fractional errors  $\delta_i$  are symmetric:  $p(\delta_i) = p(-\delta_i)$  so that  $\langle \delta_i \rangle = 0$ , then corresponding to any given matrix  $B$  the “best” estimate of  $x_j$  by almost any reasonable loss function criterion will be the  $j$ ’th row of

$$\hat{x} = (BA)^{-1}By, \quad (19-11)$$

but by making different choices of  $B$  (*i.e.*, taking different linear combinations of the equations of condition) we get different estimates. In Euler’s problem there were billions of possible choices. Which choice of  $B$  is best?

In the above we have merely restated, in modern notation but old language, the problem of “reduction of equations of condition” described in Laplace’s *Essai Philosophique* (1819). A popular criterion for solution was the principle of least squares; find that matrix  $B$  for which the sum of the squares of the errors in  $\hat{x}_j$  is a minimum; or perhaps use a weighted sum. This problem can be solved directly; we shall find the same solution by different reasoning below.

### Reformulation as a Decision Problem

But we really solved this problem in Chapter 13, where we have seen in generality that the best estimate of any parameter, by the criterion of any loss function, is found by applying Bayes’ theorem to find the probability, conditional on the data, that the parameter lies in various intervals, then making that estimate which minimizes the expected loss taken over the posterior probabilities.

Now in the original formulation of the problem, as given above, it was only a plausible conjecture that the best estimate of  $x_j$  is a linear combination of the  $y_i$  as in Equation (19-11). The material in Chapter 13 shows us a much better way of formulating the problem, in which we don’t have to depend on conjecture. Instead of trying to take linear combinations without knowing which combinations to take, we should apply Bayes’ theorem directly to the equations of condition. Then, if the best estimates are indeed of the linear form (19-11), Bayes’ theorem should not only tell us that fact, it will give us automatically the best choice of the matrix  $B$  and also tell us the accuracy of those estimates, which least squares does not give at all.

Let’s do this calculation for the case that we assign independent gaussian probabilities to the errors  $\delta_i$  of the various measurements. From our discussion in Chapter 7 we expect this to be, nearly always, the best error law we can assign from the information we have. But in the orthodox literature one would not see it that way; instead one would argue that in most physical measurements the total error is the sum of contributions from many small, causally independent imperfections, and the central limit theorem would then lead us to a gaussian *frequency distribution* of errors.<sup>†</sup> There is nothing wrong with that argument, except that it has been psychologically misleading to generations of workers, who concluded that if the frequency distribution of errors is not in fact gaussian, then to assign a gaussian probability distribution is to “assume” something that is not true; and this will lead to some horrible kind of error in our final conclusions.

**Sermon on Gaussian Error Distributions.** The considerations of Chapter 7 reassure us that this danger is grossly exaggerated; the point is that in probability theory as logic, the gaussian

---

<sup>†</sup> As noted in Chapter 14, this is subject to an important qualification; that in general the gaussian approximation will be good only for those values of total error  $\delta$  which can arise in many different ways by combination of the individual elementary errors. For unusually wide deviations we do not expect, and hardly ever observe, gaussian frequencies.

probability assignment is not an *assumption* about the frequencies of the errors; it is a *description* of our state of knowledge about the errors. We hardly ever have prior knowledge about the errors beyond the general magnitude to be expected, which we can interpret reasonably as specifying the first two moments of the error distribution. This leads, by the principle of maximum entropy, to a gaussian probability assignment as the one which agrees with that information without assuming anything else. The region  $\Omega$  of reasonably probable noise vectors  $(\delta_1 \cdots \delta_N)$  or the region  $Ax + \Omega$  of reasonably probable data vectors, is then as large as it can be while agreeing with the second moment constraints. The frequency distribution of errors is almost always unknown before seeing the data; but even if it is far from gaussian, the gaussian probability assignment will still lead us to the best inferences possible from the information we have.

The privileged status of a gaussian frequency distribution lies in a more subtle fact: acquisition of new information does not affect our inferences if that new information is only what we would have predicted from our old information. Thus if we assigned gaussian probabilities and then acquired new information that the true frequency distribution of errors is indeed gaussian with the specified variance, *this would not help us* because it is only what we would have predicted. But if we had additional prior information about the specific way in which the error frequencies depart from gaussian, that would be cogent new information constraining the possible error vectors to a smaller domain  $\Omega_1 \subset \Omega$ . This would enable us to improve our parameter estimates over the ones to be obtained below, because data vectors in the complementary set  $\Omega - \Omega_1$ , which were previously dismissed as noise, are now recognized as indicating a real “signal”. Bayes’ theorem does all this for us automatically.

Thus the covenant that we have with Nature is more favorable than supposed in orthodox teaching; for given second moments a nongaussian frequency distribution will not make our inferences worse; but *knowledge* of a nongaussian distribution would make them still better than the results to be found below.

Encouraged by the message of this Sermon, we assign the probability for the errors  $\{\delta_1 \cdots \delta_N\}$  lie in the intervals  $\{d\delta_1 \cdots d\delta_N\}$  respectively, as

$$p(\delta_1 \cdots \delta_N) d\delta_1 \cdots d\delta_N = (\text{const.}) \exp \left[ -\frac{1}{2} \sum_{i=1}^N w_i \delta_i^2 \right] d\delta_1 \cdots d\delta_N \quad (19-12)$$

where the “weight”  $w_i$  is the reciprocal variance of the error of the  $i$ ’th measurement. For example, the crude statement that the first measurement has  $\pm 2$  per cent accuracy, now becomes the more precise statement that the first measurement has weight

$$w_1 = \frac{1}{\langle \delta_1^2 \rangle} = \frac{1}{(0.02)^2} = 2500 \quad (19-13)$$

For the time being we suppose these weights known, as is generally the case with astronomical and other physical data. From (19-7) and (19-12) we have immediately the sampling probability density for obtaining measured values  $\{y_1 \cdots y_N\}$  given the true values  $\{x_1 \cdots x_N\}$ :

$$p(y_1 \cdots y_N \mid x_1 \cdots x_N) = C_1 \exp \left\{ -\frac{1}{2} \sum_{i=1}^N w_i \left[ y_i - \sum_{j=1}^n a_{ij} x_j \right]^2 \right\} \quad (19-14)$$

where  $C_1$  is independent of the  $y_i$ . According to Bayes’ theorem, if we assign uniform prior probabilities to the  $x_j$ , then the posterior probability density for the  $x_j$ , given the actual measurements  $y_i$ , is of the form

$$p(x_1 \cdots x_n \mid y_1 \cdots y_N) = C_2 \exp \left\{ -\frac{1}{2} \sum_{i=1}^N w_i \left[ y_i - \sum_{j=1}^n a_{ij} x_j \right]^2 \right\} \quad (19-15)$$

where now  $C_2$  is independent of the  $x_j$ . Next, as in nearly all gaussian calculations, we need to reorganize this quadratic form to bring out the dependence on the  $x_i$ . Expanding it, we have

$$\begin{aligned} \sum_{i=1}^N w_i \left( y_i - \sum_{j=1}^n a_{ij} x_j \right)^2 &= \sum_{i=1}^N w_i \left\{ y_i^2 - 2y_i \sum_{j=1}^n a_{ij} x_j + \sum_{j,k=1}^n a_{ij} a_{ik} x_j x_k \right\} \\ &= \sum_{j,k=1}^n K_{jk} x_j x_k - 2 \sum_{j=1}^n L_j x_j + \sum_{i=1}^N w_i y_i^2 \end{aligned} \quad (19-16)$$

where

$$K_{jk} = \sum_{i=1}^N w_i a_{ij} a_{ik}, \quad L_j = \sum_{i=1}^N w_i y_i a_{ij} \quad (19-17)$$

or, defining a diagonal “weight” matrix  $W_{ij} = w_i \delta_{ij}$ , we have a matrix  $K$  and a vector  $L$ :

$$K = \tilde{A} W A, \quad L = \tilde{A} W y \quad (19-18)$$

where  $\tilde{A}$  is the transposed matrix. We want to write (19-15) in the form

$$p(x_1 \cdots x_n \mid y_1 \cdots y_N) = C_3 \exp \left\{ -\frac{1}{2} \sum_{j,k=1}^n K_{jk} (x_j - \hat{x}_j)(x_k - \hat{x}_k) \right\} \quad (19-19)$$

whereupon the  $\hat{x}_j$  will be the mean value estimates desired. Comparing (19-16) and (19-19) we see that

$$\sum_{k=1}^n K_{jk} \hat{x}_k = L_j \quad (19-20)$$

so if  $K$  is nonsingular we can solve uniquely for  $\hat{x}$ .

### The Underdetermined Case: $K$ is Singular

If we have fewer observations than parameters,  $N < n$ , then from (19-17),  $K$  is still an  $(n \times n)$  matrix, but it is at most of rank  $N$ , and so is necessarily singular. Then the trouble is not that (19-20) has no solution; but rather that it has an infinite number of them. The maximum likelihood is attained not at a point, but on an entire linear manifold of dimensionality  $n - N$ . Of course, maximum likelihood solutions still exist as is seen from the fact that, although  $(\tilde{A} W A)^{-1}$  does not exist,  $(A \tilde{A})^{-1}$  does, and so the parameter estimate

$$x^* = \tilde{A} (A \tilde{A})^{-1} y \quad (19-21)$$

now makes the quadratic form in (19-15) vanish:  $y = Ax^*$ , achieving the maximum possible likelihood. This is called the canonical inverse solution, and the MAXENT.EXE program described in Appendix H has an option which will calculate it for you. But the canonical inverse is far from unique; for we see from (19-8) that if we add to the estimate (19-21) any solution  $z$  of the

homogeneous equation  $Az = 0$ , we have another estimate  $x^* + z$  with just as high a likelihood; and there is a linear manifold  $\Delta$  of such vectors  $x^* + z$ , of dimensionality  $n - N$ .

**Exercise 19.1.** Show that the canonical inverse solution (19–21) is also a least squares one, making  $\sum (x_i^*)^2$  a minimum on the manifold  $\Delta$ . Unfortunately, there seems to be no compelling reason why one should want the vector of estimates to have minimum length.

For a long time no satisfactory way of dealing with such problems was recognized; yet we are not entirely helpless, for the data do restrict the possible values of the parameters  $\{x_i\}$  to a “feasible set”  $\Delta$  satisfying (19–20). The data alone are incapable of picking out any unique point in this set; but the data may be supplemented with prior information which enables us to make a useful choice in spite of that. These are “generalized inverse” problems, which are of current importance in many applications such as image reconstruction. In fact, in the real world, generalized inverse problems probably make up the great majority, because the real world seldom favors us with all the information needed to make a well-posed problem. Yet useful solutions may be found in many cases by maximum entropy which resolves the ambiguity in a way that is “optimal” by several different criteria, as described in Chapters 11 and 24.

### The Overdetermined Case: $K$ Can be Made Nonsingular

By its definition (19–17),  $K$  is an  $(n \times n)$  matrix, and for all real  $\{q_1 \cdots q_n\}$  such that  $\sum q_i^2 > 0$ ,

$$\sum_{j,k=1}^n K_{jk} q_j q_k = \sum_{i=1}^N w_i \left( \sum_{j=1}^n a_{ij} q_j \right)^2 \geq 0 \quad (19-22)$$

so if  $K$  is of rank  $n$  it is not only nonsingular, but positive definite. If  $N \geq n$  this will be the case unless we have done something foolish in setting up the problem – including a useless observation or an irrelevant parameter.

For, in the first place, we suppose all the weights  $w_i$  to be positive; if any observation  $y_i$  has weight  $w_i = 0$ , then it is useless in our problem; that is, it can convey no information about the parameters and we should not have included it in the data set at all. We can reduce  $N$  by one.

Secondly, if there is a nonzero vector  $q$  for which  $\sum_j a_{ij} q_j$  is zero for all  $i$ , then in (19–7) for all  $c$ , the parameter sets  $\{x_j\}$  and  $\{x_j + c q_j\}$  would lead to identical data, and so could not be distinguished whatever the data. In other words, there is an irrelevant parameter in the problem which has nothing to do with the data; we can reduce  $n$  by one. Mathematically, this means that the columns of the matrix  $A$  are not linearly independent; then if  $q_k \neq 0$ , we can remove the parameter  $x_k$  and the  $k$ 'th column of  $A$  with no essential change in the problem (*i.e.*, no change in the information we get from it).

Removing irrelevant observations and parameters if necessary, if finally the number of cogent observations is at least as great as the number of relevant parameters, then  $K$  is a positive definite matrix and (19–20) has a unique solution

$$\hat{x}_k = \sum_{j=1}^n (K^{-1})_{kj} L_j. \quad (19-23)$$

From (19–18), we can write the result as

$$\hat{x} = (\tilde{A}W A)^{-1} \tilde{A}W y \quad (19-24)$$

and, comparing with (19–11), we see that in the gaussian case with uniform prior probabilities, the best estimates are indeed linear combinations of the measurements, of the form (19–11), and the best choice of the matrix  $B$  is

$$B = \tilde{A}W, \quad (19-25)$$

a result perhaps first found by Gauss, and repeated in Laplace's work referred to. Let us evaluate this solution for our simple problem.

### Numerical Evaluation of the Result

Applying the solution (19-24) to our problem of estimating  $e$  and  $m$ , the measurements of  $e$ ,  $(e/m)$ ,  $(e^2/m)$  were of 2%, 1%, 5% accuracy respectively, and so

$$\begin{aligned} w_2 &= \frac{1}{(0.01)^2} = 10,000 \\ w_3 &= \frac{1}{(0.05)^2} = 400 \end{aligned} \quad (19-26)$$

and we found  $w_1 = 2500$  before. Thus we have

$$B = \tilde{A}W = \begin{pmatrix} 1 & 1 & 2 \\ 0 & -1 & -1 \end{pmatrix} \begin{pmatrix} w_1 & 0 & 0 \\ 0 & w_2 & 0 \\ 0 & 0 & w_3 \end{pmatrix} = \begin{pmatrix} w_1 & w_2 & 2w_3 \\ 0 & -w_2 & -w_3 \end{pmatrix} \quad (19-27)$$

$$K = \tilde{A}WA = \begin{pmatrix} (w_1 + w_2 + 4w_3) & -(w_2 + 2w_3) \\ -(w_2 + 2w_3) & (w_2 + w_3) \end{pmatrix} \quad (19-28)$$

$$K^{-1} = (\tilde{A}WA)^{-1} = \frac{1}{|K|} \begin{pmatrix} (w_2 + w_3) & (w_2 + 2w_3) \\ (w_2 + 2w_3) & (w_1 + w_2 + 4w_3) \end{pmatrix} \quad (19-29)$$

where

$$|K| = \det(K) = w_1w_2 + w_2w_3 + w_3w_1 \quad (19-30)$$

Thus the final result is

$$(\tilde{A}WA)^{-1}\tilde{A}W = \frac{1}{|K|} \begin{pmatrix} w_1(w_2 + w_3) & -w_2w_3 & w_2w_3 \\ w_1(w_2 + 2w_3) & -w_2(w_1 + 2w_3) & w_3(w_2 - w_1) \end{pmatrix} \quad (19-31)$$

and the best point estimates of  $x_1, x_2$  are

$$\begin{aligned} \hat{x}_1 &= \frac{w_1(w_2 + w_3)y_1 + w_2w_3(y_3 - y_2)}{w_1w_2 + w_2w_3 + w_3w_1} \\ \hat{x}_2 &= \frac{w_1w_2(y_1 - y_2) + w_2w_3(y_3 - 2y_2) + w_3w_1(2y_1 - y_3)}{w_1w_2 + w_2w_3 + w_3w_1} \end{aligned} \quad (19-32)$$

Inserting the numerical values of  $w_1, w_2, w_3$ , we have

$$\begin{aligned} \hat{x}_1 &= \frac{13}{15}y_1 + \frac{2}{15}(y_2 - y_3) \\ \hat{x}_2 &= \frac{5}{6}(y_1 - y_2) + \frac{2}{15}(y_3 - 2y_2) + \frac{1}{30}(2y_1 - y_3) \end{aligned} \quad (19-33)$$

which exhibits the best estimates as weighted averages of the estimates taken from all possible pairs of experiments. Thus,  $y_1$  is the estimate of  $x_1$  obtained in the first experiment, which measures  $e$  directly. The second and third experiments combined yield an estimate of  $e$  given by  $(e^2/m)(e/m)^{-1}$ . Since

$$\frac{\frac{e_0^2}{m_0}(1+y_3)}{\frac{e_0}{m_0}(1+y_2)} \approx e_0(1+y_3-y_2) \quad (19-34)$$

$(y_3 - y_2)$  is the estimate of  $x_1$  given by experiments 2 and 3. Equation (19-33) says that these two independent estimates of  $x_1$  should be combined with weights 13/15, 2/15. Likewise,  $\hat{x}_2$  is given as a weighted average of three different (although not independent) estimates of  $x_2$ .

### Accuracy of the Estimates

From (19-19) we find the second central moments of  $p(x_1 \dots x_n \mid y_1 \dots y_N)$ :

$$\langle (x_j - \hat{x}_j)(x_k - \hat{x}_k) \rangle = \langle x_j x_k \rangle - \langle x_j \rangle \langle x_k \rangle = (K^{-1})_{jk} \quad (19-35)$$

Thus from the  $(n \times n)$  inverse matrix

$$K^{-1} = (\tilde{A}WA)^{-1} \quad (19-36)$$

already found in our calculation of  $\hat{x}_j$ , we can also read off the probable errors, or more conveniently, the standard deviations. From (19-29) we can state the results in the form (mean)  $\pm$  (standard deviation) as

$$(x_j)_{est} = \hat{x}_j \pm \sqrt{(K^{-1})_{jj}}. \quad (19-37)$$

Equations (19-24) and (19-37) represent the general solution of the problem, which Euler needed. In the present case this is

$$\begin{aligned} (x_1)_{est} &= \hat{x}_1 \pm \left[ \frac{w_2 + w_3}{w_1 w_2 + w_2 w_3 + w_3 w_1} \right]^{1/2} \\ (x_2)_{est} &= \hat{x}_2 \pm \left[ \frac{w_1 + w_2 + 4w_3}{w_1 w_2 + w_2 w_3 + w_3 w_1} \right]^{1/2} \end{aligned} \quad (19-38)$$

with numerical values

$$\begin{aligned} x_1 &= \hat{x}_1 \pm 0.0186 \\ x_2 &= \hat{x}_2 \pm 0.0216 \end{aligned} \quad (19-39)$$

so that from the three measurements we obtain  $e$  with  $\pm 1.86$  per cent accuracy,  $m$  with  $\pm 2.16$  per cent accuracy.

How much did the rather poor measurement of  $(e^2/m)$ , with only  $\pm 5$  per cent accuracy, help us? To answer this, note that in the absence of this experiment we would have arrived at conclusions given by (19-28), (19-46) and (19-32) in the limit  $w_3 \rightarrow 0$ . The results (also easily verified directly from the statement of the problem) are

$$\begin{aligned} \hat{x}_1 &= y_1 \\ \hat{x}_2 &= y_1 - y_2 \end{aligned} \quad (19-40)$$

$$K^{-1} = \frac{1}{w_1 w_2} \begin{pmatrix} w_2 & w_2 \\ w_2 & (w_1 + w_2) \end{pmatrix} \quad (19-41)$$

or, the (mean)  $\pm$  (standard deviation) values are

$$\begin{aligned} x_1 &= y_1 \pm \frac{1}{w_1} = y_1 \pm 0.020 \\ x_2 &= y_1 - y_2 \pm \left[ \frac{w_1 + w_2}{w_1 w_2} \right]^{1/2} = y_1 - y_2 \pm 0.024 \end{aligned} \quad (19-42)$$



As might have been anticipated by common sense, a low-accuracy measurement can add very little to the results of accurate measurements, and if the  $(e^2/m)$  measurement had been much worse than  $\pm 5$  per cent it would hardly be worthwhile to include it in our calculations. But suppose that an improved technique gives us an  $(e^2/m)$  measurement of  $\pm 2$  per cent accuracy. How much would this help? The answer is given by our previous formulas with  $w_1 = w_3 = 2500$ ,  $w_2 = 10,000$ . We find now that the mean-value estimates give much higher weight to the estimates using the  $(e^2/m)$  measurement:

$$\begin{aligned}\hat{x}_1 &= 0.556y_1 + 0.444(y_3 - y_2) \\ \hat{x}_2 &= 0.444(y_1 - y_2) + 0.444(y_3 - 2y_2) + 0.112(2y_1 - y_3)\end{aligned}\tag{19-43}$$

which is to be compared with (19-33). The standard deviations are given by

$$\begin{aligned}x_1 &= \hat{x}_1 \pm 0.0149 \\ x_2 &= \hat{x}_2 \pm 0.020\end{aligned}\tag{19-44}$$

The accuracy of  $e$  ( $x_1$ ) is improved roughly twice as much as that of  $m$  ( $x_2$ ), since the improved measurement involves  $e^2$ , but only the first power of  $m$ .

**Exercise 19.2.** Write a computer program which solves this problem for general  $N$  and  $n$ , with  $N \geq n$ , and test it on the problem just solved. Estimate how long it would require for the compiled program to solve Euler's problem. A usable routine for matrix inversion is given in the source code file MAXENT.BAS, lines 600-850. The same routine is included in many FORTRAN packages.

### Generalization

In the above we supposed the weights  $w_i$  known from prior information. If this is not the case, there are many different conceivable kinds of partial prior information about them, leading to many different possible prior probability assignments  $p(w_1 \cdots w_n | I)$ . Let us see how this circumstance would change the above solutions.

\*\*\*\*\* MORE TO COME HERE! \*\*\*\*\*

### COMMENTS

**A Paradox.** We can learn many more things from studying this problem. For example, let us note something which you will find astonishing at first. If you study Equation (19-32), which gives the best estimate of  $m$  from the three measurements, you will see that  $y_3$ , the result of the  $(e^2/m)$  measurement, enters into the formula in a different way than  $y_1$  and  $y_2$ . It appears once with a positive coefficient, and once with a negative one. If  $w_1 = w_2$ , these coefficients are equal and (19-32) collapses to

$$\hat{x}_2 = y_1 - y_2\tag{19-45}$$

Now, realize the full implications of this: *it says that the only reason we make use of the  $(e^2/m)$  measurement in estimating  $m$  is that the  $(e)$  measurement and the  $(e/m)$  measurement have different accuracy.* No matter how accurately we know  $(e^2/m)$ , if the  $(e)$  and  $(e/m)$  measurements happen to have the same accuracy, however poor, then we should ignore the good measurement and base our estimate of  $m$  only on the  $(e)$  and  $(e/m)$  measurements!

We think that, on first hearing, your intuition will revolt against this conclusion, and your first reaction will be that there must be an error in Equation (19-32). So, check the derivation at your

leisure. This is a perfect example of the kind of result which probability theory gives us almost without effort, but which our unaided common sense might not notice in years of thinking about the problem. We won't deprive you of the pleasure of resolving this "paradox" for yourself, and explaining to your friends how it can happen that consistent inductive reasoning may demand that you throw away your best measurement.

In Chapter 17, we complained about the fact that orthodox statisticians sometimes throw away relevant data in order to fit a problem to their preconceived model of "independent random errors." Are we now guilty of the same offense? No doubt, it looks very much that way! Yet we plead innocence; the numerical value of  $(e^2/m)$  is in fact *irrelevant* to inference about  $m$ , if we already have measurements of  $e$  and  $e/m$  of equal accuracy. To see this, suppose that we knew  $(e^2/m)$  exactly from the start. How would you make use of that information in this problem? If you try to do this, you will soon see why  $(e^2/m)$  is irrelevant. But to clinch matters, try the following exercise:

**Exercise 19.3.** Consider a specific case:  $w_1 = w_2 = 1$ ,  $w_3 = 100$ ; the third measurement is ten times more accurate than the first two. If the third measurement cancels out when we try to use all three as in (19-22), then it seems that the only way we could use the third measurement is by discarding either the first or second. Show that, nevertheless, the estimates made by (19-32) using only the first and second measurements are more accurate than those made by using the first and third; or the second and third. Now explain intuitively why the paradox is not real.

As another example, it is important that we understand the way our conclusions depend on our choice of loss functions and probability distributions for the errors  $\delta_i$ . If we use instead of the gaussian distribution (19-12) one with wider tails, such as the Cauchy distribution  $p(\delta) \propto (1 + w\delta^2/2)^{-1}$ , the posterior distribution  $p(x_1x_2 | y_1y_2y_3)$  may have more than one peak in the  $(x_1, x_2)$ -plane. Then a quadratic loss function, or more generally any concave loss function (*i.e.*, doubling the error more than doubles the loss) will lead one to make estimates of  $x_1$  and  $x_2$  which lie between the peaks, and are known to be very unlikely. With a convex loss function a different "paradox" appears, in that the basic equation (19-46) for constructing the best estimator may have more than one solution, with nothing to tell us which one to use.

The appearance of these situations is the robot's way of telling us this: our state of knowledge about  $x_1$  and  $x_2$  is too complicated to be described adequately simply by giving best estimates and probable errors. The only honest way of describing what we know is to give the actual distribution  $p(x_1x_2|y_1y_2y_3)$ . This is one of the limitations of decision theory, which we need to understand in order to use it properly.