

would have to coincide with the data perfectly. Thus the quantity,  $\chi^2$  serves as a measure of how far the theory is from the experiment. To get a good "fit" to the data we want the smallest  $\chi^2$  possible. However, it is unreasonable to expect a value of  $\chi^2$  smaller than the number of data points because of the (assumed) random errors on the data. If it is much less than the number of data points some problem is to be suspected, the errors in the data are not uncorrelated, the errors are overestimated, etc.

We will suppose that the theoretical values,  $T_i$ , to be compared with data can be written as:

$$T_i = \sum_{j=1}^n b_j f_{ji} \quad (5.3)$$

where the  $f_{ji}$  are a series of functions, labeled by  $j$ , evaluated for the experimental conditions  $i$ , and the  $b_j$  are the parameters to be fit. For example, if the data were taken at series of points  $t_i$ , then the functions,  $f_{ji}$  might be monomials in  $t_i$  with  $j$  denoting the power of  $t_i$ . The theoretical values would then be polynomials in the independent variable  $t$ . For the example of constant acceleration given above we would have

$$f_{1i} = t_i^0 \equiv 1; \quad f_{2i} = t_i \quad \text{and} \quad f_{3i} = t_i^2 \quad (5.4)$$

and

$$b_1 = x_0; \quad b_2 = v_0 \quad \text{and} \quad b_3 = \frac{1}{2}a. \quad (5.5)$$

In the following discussion we assume that the theoretical function is linear in the parameters  $b_j$  as expressed by Eq. 5.3. If it is not, then one may still apply the following analysis if we suppose that the  $b_j$  represent the coefficients in a Taylor's series expansion about some point. We wish to formulate the problem of finding the "best" fit of the theory to the data by minimizing the function

$$\chi^2 = \sum_{i=1}^N \frac{(\sum_{j=1}^n b_j f_{ji} - d_i)^2}{e_i^2}. \quad (5.6)$$

Differentiating with respect to each parameter in turn,

$$\frac{d\chi^2}{db_k} = 2 \sum_{i=1}^N \frac{(\sum_{j=1}^n b_j f_{ji} - d_i) f_{ki}}{e_i^2} = 0; \quad k = 1, 2, 3, \dots, n \quad (5.7)$$

leading to the system of equations:

$$\sum_{j=1}^n a_{kj} b_j = c_k; \quad k = 1, 2, 3, \dots, n \quad (5.8)$$

where

$$c_k = \sum_{i=1}^N \frac{d_i f_{ki}}{e_i^2}; \quad \text{and} \quad a_{kj} = \sum_{i=1}^N \frac{f_{ki} f_{ji}}{e_i^2}.$$

## 5.2. SOLUTION OF LINEAR EQUATIONS

The solution of these equations will give the desired parameters. Note that the number of equations to be solved is given by the number of parameters and has nothing to do with number of data points, except that if there are not (at least) as many independent data points as parameters the system will not have a unique solution.

## 5.2 Solution of Linear Equations

Let us consider a set of equations of the form:

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n &= y_1 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 + \dots + a_{2n}x_n &= y_2 \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 + \dots + a_{3n}x_n &= y_3 \\ \vdots &\vdots \\ a_{n1}x_1 + a_{n2}x_2 + a_{n3}x_3 + \dots + a_{nn}x_n &= y_n \end{aligned} \quad (5.9)$$

which we can write in abbreviated form as an augmented matrix

$$\left( \begin{array}{cccc|c} a_{11} & a_{12} & a_{13} & \dots & a_{1n} & y_1 \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} & y_2 \\ a_{31} & a_{32} & a_{33} & \dots & a_{3n} & y_3 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \dots & a_{nn} & y_n \end{array} \right). \quad (5.10)$$

In matrix notation Eq. 5.9 is written as:

$$Ax = y \quad (5.11)$$

where we ask for the solution of  $x$  for a given  $y$ .

We may multiply any of the equations in 5.9 by a constant and subtract it from another without changing the value of the system. Thus we can multiply the first equation by  $\frac{a_{21}}{a_{11}}$  and subtract it from the second equation. This leads to:

$$\left( \begin{array}{cccc|c} a_{11} & a_{12} & a_{13} & \dots & a_{1n} & y_1 \\ 0 & a_{22} - m_{21}a_{12} & a_{23} - m_{21}a_{13} & \dots & a_{2n} - m_{21}a_{1n} & y_2 - m_{21}y_1 \\ a_{31} & a_{32} & a_{33} & \dots & a_{3n} & y_3 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \dots & a_{nn} & y_n \end{array} \right). \quad (5.12)$$

Next we can multiply the first equation by  $m_{31} \equiv \frac{a_{31}}{a_{11}}$  and subtract it from the third equation to reduce the entry in the third row, first column to zero. This process



where  $L$  is a matrix with zeros above the diagonal and  $U$  is a matrix with zeros below the diagonal. The matrix  $U$  for the Gaussian elimination presented in the previous section is the one which results after the elimination, i.e. the  $b_{ij}$  in Eq. 5.13. It can be shown (see e.g. Ref. 1, p. 56) that the lower matrix for this case is:

$$\begin{pmatrix} 1 & 0 & 0 & \dots & 0 & 0 \\ m_{21} & 1 & 0 & \dots & 0 & 0 \\ m_{31} & m_{32} & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ m_{n1} & m_{n2} & m_{n3} & \dots & m_{n,n-1} & 1 \end{pmatrix} \quad (5.18)$$

where the  $m_{i,j}$  are just the ratios needed to reduce the  $ij$  element to zero.

Thus the system,

$$Ax = y, \quad \text{or} \quad LUx = y \quad (5.19)$$

may be considered as two systems;

$$Ux = z, \quad \text{and} \quad Lz = y. \quad (5.20)$$

Since both systems are triangular they only require back substitution to solve them and one can first solve for  $z$  followed by  $x$  with two “ $n^2$ ” procedures. The back substitution for  $L$  in the Gaussian elimination method is done at the same time as the elimination by operating on the  $y$  vector simultaneously. Since we know that there are only zeros below the diagonal when we finish the elimination, and that the diagonal values of  $L$  are unity, there is no reason to store either the “zeros” of  $U$  or the “ones” of  $L$  and we have just enough room to store the non-trivial elements of  $L$  in the locations that would be occupied by the zeros of  $U$ . Thus, with (essentially) no cost we obtain an LU decomposition of the matrix at the same time that a solution is obtained. Many equation solvers operate in two steps. They first perform an LU decomposition with one subroutine and then perform the back substitution with another.

Following is a subroutine which eliminates the first column in a matrix from the second element on. As input it has the preceding row, normalized to have the first element equal to -1. Note that it operates on  $m$  columns, independent of the size of the real matrix so it can perform the elimination for several input vectors or the unit matrix in the case of a matrix inversion. The first step in the “do 3” loop is to zero the first element in the row. After the loop is finished the next statement fills in that element with the value of “ $m_{i,j}$ ” to create the  $L$  matrix

```
subroutine elim(a,t,nn,m,ldim,con1) ! eliminates 1 column
dimension a(ldim,1),t(1)
do 4 i=2,nn+1 ! count down the column
```

## 5.2. SOLUTION OF LINEAR EQUATIONS

```
con=a(i,1)
do 3 j=1,m
3 a(i,j)=a(i,j)+con*t(j) ! SAXPY
a(i,1)=con*con ! fill in L matrix
4 continue
return
end
```

The following routine calls “elim”, column by column, to perform a complete elimination (or LU decomposition). The two routines were designed to work together and could have been constructed as a single unit. The reason for separating them is the the inner code might well be replaced by one written in machine language for greater speed.

```
subroutine gauss(a,n,m,ldim)
dimension t(1000),a(ldim,1)
do 10 i=1,n-1
con=1./a(i,i)
do 5 k=1,m+1-i
5 t(k)=-con*a(i,k+i-1)
call elim(a(i,i),t,n-i,m+1-i,ldim,con)
10 continue
return
end
```

The LU decomposition just studied is by no means unique nor is the way of obtaining it. Doolittle’s method finds the same LU decomposition as Gauss’ method by a different technique and Crout’s method chooses the diagonal elements of  $U$  to be unity instead of those of  $L$ . The principal advantage of these methods over Gauss’ method is that the accumulation of errors is better controlled. On a machine with a large word length there is little to choose among them and we will primarily consider the Gauss method.

For studies of the matrices it is useful to have a practice system. Below is such a matrix to be used in the problems.

$$\begin{pmatrix} 0.546 & 0.447 & 0.242 & 0.194 & 0.795 \\ 0.380 & 0.276 & 0.581 & 0.108 & 0.416 \\ 0.721 & 0.022 & 0.853 & 0.068 & 0.312 \\ 0.151 & 0.759 & 0.186 & 0.597 & 0.757 \\ 0.192 & 0.509 & 0.041 & 0.411 & 0.632 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} = \begin{pmatrix} 0.468 \\ 0.695 \\ 0.398 \\ 0.913 \\ 0.483 \end{pmatrix} \quad (5.21)$$

### 5.2.2 The Gauss-Seidel Method

For those systems in which the diagonal elements of the matrix  $A$  are dominant iterative methods for the solution of linear systems are often useful. Let us consider the solution to Eq. 5.9:

$$\begin{array}{ccccccc} a_{11}x_1 & +a_{12}x_2 & +a_{13}x_3 & \dots & +a_{1n}x_n & = & y_1 \\ a_{21}x_1 & +a_{22}x_2 & +a_{23}x_3 & \dots & +a_{2n}x_n & = & y_2 \\ a_{31}x_1 & +a_{32}x_2 & +a_{33}x_3 & \dots & +a_{3n}x_n & = & y_3 \\ \vdots & \vdots & \vdots & \vdots & \vdots & & \vdots \\ a_{n1}x_1 & +a_{n2}x_2 & +a_{n3}x_3 & \dots & +a_{nn}x_n & = & y_n \end{array} \quad (5.22)$$

in a situation in which the matrix  $A$  is dominated by the diagonal elements. In this case we might believe that a first crude approximation would be given by the solution that one would have if only the diagonal elements were non-zero i.e.,

$$\begin{array}{l} x_1 = y_1/a_{11} \\ x_2 = y_2/a_{22} \\ x_3 = y_3/a_{33} \\ \vdots \\ x_n = y_n/a_{nn}. \end{array} \quad (5.23)$$

We can imagine an approximation scheme in which we refine this first approximation by computing successive approximations to the  $x_i$  as,

$$\begin{array}{l} x'_1 = \frac{y_1 - \sum_{j=2}^n a_{1j}x'_j}{a_{11}} \\ x'_2 = \frac{y_2 - a_{21}x'_1 - \sum_{j=3}^n a_{2j}x'_j}{a_{22}} \\ x'_3 = \frac{y_3 - \sum_{j=1}^2 a_{3j}x'_j - \sum_{j=4}^n a_{3j}x'_j}{a_{33}} \\ \vdots \\ x'_i = \frac{y_i - \sum_{j=1}^{i-1} a_{ij}x'_j - \sum_{j=i+1}^n a_{ij}x'_j}{a_{ii}} \\ \vdots \\ x'_n = \frac{y_n - \sum_{j=1}^{n-1} a_{nj}x'_j}{a_{nn}} \end{array} \quad (5.24)$$

### 5.2. SOLUTION OF LINEAR EQUATIONS

where we have used the new values of  $x_i$  ( $x'_i$ ) if they are defined and the old values if they are not. We may attempt to improve on this algorithm by defining the "new" values as a linear combination of the result of Eqs. 5.24 and the old values.

$$\begin{array}{l} x''_1 = \frac{y_1 - \sum_{j=2}^n a_{1j}x_j}{a_{11}} \\ x''_2 = \frac{y_2 - a_{21}x'_1 - \sum_{j=2}^n a_{2j}x_j}{a_{22}} \\ x''_3 = \frac{y_3 - \sum_{j=1}^2 a_{3j}x'_j - \sum_{j=4}^n a_{3j}x_j}{a_{33}} \\ \vdots \\ x''_i = \frac{y_i - \sum_{j=1}^{i-1} a_{ij}x'_j - \sum_{j=i+1}^n a_{ij}x_j}{a_{ii}} \end{array}$$

where now

$$x''_n = \frac{y_n - \sum_{j=1}^{n-1} a_{nj}x'_j}{a_{nn}} \quad (5.25)$$

$$x'_i = (1-w)x_i + wx''_i. \quad (5.26)$$

The parameter  $w$  is allowed to vary between 1 and 2. This last method is referred to as the Gauss-Seidel algorithm with successive overrelaxation (SOR). Of course, if  $w = 1$ , it is the same as the original Gauss-Seidel method given in Eqs. 5.24. By varying the parameter,  $w$ , the speed of convergence of the method may be increased by as much as a factor of 10. Of course, there is no guarantee that this process converges. One condition that assures convergence is

$$|a_{ii}| > \sum_{i \neq j} |a_{ij}| \quad \text{for all } i. \quad (5.27)$$

For matrices which are nearly diagonal this technique may be much faster than the elimination method. Of course, if most of the off-diagonal elements are zero it is particularly efficient since those elements do not need to be included in the sums.

### 5.2.3 The Householder Transformation

There is one other method, introduced by Householder in 1958, which is useful not only for the solution of equations but for the solution of eigenvalue problems considered in the next section. Let us treat a slightly more general decomposition than the LU form, called the QR decomposition. The  $R$  denotes the same type of matrix as

U before, an upper (or right) triangular matrix, but Q can represent any matrix. We now have (as before) the two systems,

$$Qz = y, \text{ and } Rx = z. \quad (5.28)$$

Naturally this decomposition is only useful if the system  $Qz = y$  can be easily solved, as was the case for the lower triangular matrix. Another type of matrix which is easily inverted is the orthogonal matrix which has the property that

$$P\bar{P} = I \quad (5.29)$$

or

$$\sum_j P_{ij} P_{kj} = \delta_{ik} \quad (5.30)$$

i.e. its transpose is its inverse. The problem of the solution of linear equations can now be restated:

Find a matrix (an orthogonal matrix),  $Q^{-1}$ , such that

$$Q^{-1}A = R. \quad (5.31)$$

That is, find an orthogonal matrix such that, when multiplied by A the result is an upper triangular matrix. The (orthogonal) matrix,  $Q^{-1}$  can be broken down into a product of several (orthogonal) matrices, each one producing some zeros in the lower triangle (without destroying the ones that were produced by the previously acting matrices).

$$Q^{-1} \equiv P^{(n)} P^{(n-1)} \dots P^{(1)} P^{(0)}. \quad (5.32)$$

Householder solved this problem in the following way. Consider a matrix of the form

$$P^{(r)} = I - Gh^{(r)}\bar{h}^{(r)} \quad (5.33)$$

i.e., in terms of components,

$$P_{ij}^{(r)} = \delta_{ij} - Gh_i^{(r)}h_j^{(r)} \quad (5.34)$$

where

$$h_i^{(r)} = 0 \quad \text{for } i = 1, 2, \dots, r \quad (5.35)$$

but the values of  $h_i$  for  $i = r + 1, r + 2, \dots, n$  are yet to be determined and will be chosen as a function of the elements of the matrix to be put in upper triangular

## 5.2. SOLUTION OF LINEAR EQUATIONS

form. Thus  $P^{(r)}$  has the form:

$$\begin{pmatrix} 1 & 0 & 0 & \dots & 0 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots & 0 & 0 & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & 0 & 0 & \dots \\ 0 & 0 & 0 & \dots & 0 & 1 - Gh_{r+1}^2 & -Gh_{r+1}h_{r+2} & -Gh_{r+1}h_{r+3} \\ 0 & 0 & 0 & \dots & 0 & -Gh_{r+1}h_{r+2} & 1 - Gh_{r+2}^2 & -Gh_{r+2}h_{r+3} \\ 0 & 0 & 0 & \dots & 0 & -Gh_{r+1}h_{r+3} & -Gh_{r+2}h_{r+3} & 1 - Gh_{r+3}^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{pmatrix}. \quad (5.36)$$

Since we are considering the case in which  $P^{(r)}$  is orthogonal

$$\sum_{j=1}^n P_{ij} P_{jk} = \delta_{ik} \quad (5.37)$$

$$= \sum_j \delta_{ij} \delta_{jk} - G \sum_j \delta_{ij} h_j h_k - G \sum_j h_i h_j \delta_{jk} + G^2 \sum_j h_i h_j^2 h_k \quad (5.38)$$

so

$$G^2 \sum h_j^2 = 2G$$

or

$$G = \frac{2}{\sum h_j^2}. \quad (5.39)$$

Thus we are free to choose the  $h_i$  to be anything we want and the matrix  $P^{(r)}$  will be orthogonal provided we chose  $G$  according to this equation. Since the goal is to use these matrices to bring another matrix, A, into triangular form column by column consider the action of the Householder matrix on a single column vector x, where is to be thought of as one of the columns of the matrix A.

$$b = P^{(r)}x = x - Gh^{(r)}\bar{h}^{(r)}x \quad (5.40)$$

or

$$b_i = x_i - Gh_i^{(r)} \sum_j h_j^{(r)} x_j. \quad (5.41)$$

Since the first  $r$  values of  $h$  are zero, the first  $r$  components of  $b$  are the same as the of  $x$ . We see that there is a good chance to reduce the remainder of the component to zero by choosing  $h_i^{(r)}$  to be proportional to  $x_i$ . In fact we cannot make all of the zero because we have Eq. 5.39 to satisfy, but by an appropriate choice of  $h_i^{(r)}$  we can annihilate all but one entry. Let us pick:

$$h_{r+1}^{(r)} = x_{r+1} + z \quad \text{and} \quad h_m^{(r)} = x_m; \quad m = r + 2, r + 3, \dots, n \quad (5.42)$$

where  $z$  is a constant to be determined. In this case the sum in Eq. 5.41 is

$$\sum_j h_j^{(r)} x_j = (x_{r+1} + z)x_{r+1} + \sum_{j=r+2}^n x_j^2 \quad (5.43)$$

$$= zx_{r+1} + S^2; \quad S^2 \equiv \sum_{j=r+1}^n x_j^2. \quad (5.44)$$

In order to make the higher components of  $\mathbf{b}$  zero from Eq. 5.41 we must have

$$G(zx_{r+1} + S^2) = 1 \quad (5.45)$$

giving

$$b_{r+1} = -z. \quad (5.46)$$

and

$$b_i = 0 \quad i = r+2, r+3, \dots, n \quad (5.47)$$

To determine  $z$  we combine Eqs. 5.39 and 5.45:

$$2(zx_{r+1} + S^2) = \sum h_j^2 = z^2 + 2zx_{r+1} + S^2$$

or

$$z = \pm\sqrt{S^2}. \quad (5.48)$$

By operating successively with the matrices  $P^{(r)}$ , a general matrix can be reduced to triangular form.

$$\mathbf{R} = \mathbf{P}^{(n)}\mathbf{P}^{(n-1)} \dots \mathbf{P}^{(1)}\mathbf{P}^{(0)}\mathbf{A}. \quad (5.49)$$

Here  $P^{(0)}$  is chosen to wipe out all elements of the first column except the first one,  $P^{(1)}$  zeros all elements of the second column except the first two, etc. Note that the application of  $P^{(1)}$  does not change the first column of  $P^{(0)}\mathbf{A}$  so that the zeros created there remain. Recalling Eq. 5.32

$$\mathbf{Q}^{-1} \equiv \mathbf{P}^{(n)} \dots \mathbf{P}^{(1)}\mathbf{P}^{(0)} \quad (5.50)$$

then  $\mathbf{A} = \mathbf{Q}\mathbf{R}$  and the equation  $\mathbf{A}\mathbf{x} = \mathbf{y}$  can be solved by back substitution on

$$\mathbf{Q}^{-1}\mathbf{A}\mathbf{x} = (\mathbf{Q}^{-1}\mathbf{A})\mathbf{x} = \mathbf{R}\mathbf{x} = \mathbf{Q}^{-1}\mathbf{y} = \mathbf{P}^{(n)} \dots \mathbf{P}^{(1)}\mathbf{P}^{(0)}\mathbf{y}. \quad (5.51)$$

It is important to realize that the number of steps for multiplication of an  $n \times n$  matrix by a Householder matrix is of order  $n^2$  and not  $n^3$  since

$$(PA)_{ij} = A_{ij} - Gh_i \sum_k h_k A_{kj}$$

so that multiplication can be done one row at a time, first calculating

$$C_j = \sum_k h_k A_{kj} \quad n \text{ operations}$$

then for each  $j$ .

$$(PA)_{ij} = A_{ij} - Gh_i C_j.$$

### 5.3. THE EIGENVALUE PROBLEM

## 5.3 The Eigenvalue Problem

Eigenvalues of matrix (or differential) systems form the basis of much of modern physics. They represent the frequency of the principal modes of motion of systems, energy levels of atoms, nuclei (or any quantum system) and indeed, the masses of the particles making up the universe. Hence the study of numerical techniques for finding the eigenvalues of a mathematical system are of fundamental importance.

### 5.3.1 Coupled Oscillators

Consider, as a model physics problem, a system of equal masses constrained to move in one dimension. Mass 1 is connected to a fixed wall by a spring with constant  $k_{01}$ , mass 2 is connected to mass 1 with a spring with constant  $k_{12}$  and mass 3 with  $k_{23}$ . In addition there is a spring connecting mass 1 and mass 3 with constant  $k_{13}$ . If  $u_1$ ,  $u_2$ , and  $u_3$  denote the displacement of the masses from the unstretched configuration then Newton's equations become:

$$\begin{aligned} m\ddot{u}_1 &= -(k_{01} + k_{12} + k_{13})u_1 + k_{12}u_2 + k_{13}u_3 \\ m\ddot{u}_2 &= k_{12}u_1 - (k_{12} + k_{23})u_2 + k_{23}u_3 \\ m\ddot{u}_3 &= k_{13}u_1 + k_{23}u_2 - (k_{23} + k_{13})u_3 \end{aligned} \quad (5.52)$$

In matrix form:

$$m \begin{pmatrix} \ddot{u}_1 \\ \ddot{u}_2 \\ \ddot{u}_3 \end{pmatrix} = \begin{pmatrix} -k_{01} - k_{12} - k_{13} & k_{12} & k_{13} \\ k_{12} & -k_{12} - k_{13} & k_{23} \\ k_{13} & k_{23} & -k_{23} - k_{13} \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} \quad (5.53)$$

or

$$m\ddot{\mathbf{u}} = \mathbf{A}\mathbf{u} \quad (5.54)$$

If we suppose a solution of the form:

$$\mathbf{u} = \sum_j \alpha_j \mathbf{x}_j e^{i\omega_j t}, \quad (5.55)$$

then

$$-m \sum_j \alpha_j \mathbf{x}_j \omega_j^2 e^{i\omega_j t} = \sum_j \mathbf{A} \alpha_j \mathbf{x}_j e^{i\omega_j t}. \quad (5.56)$$

Since each function  $e^{i\omega_j t}$  is linearly independent of the others (assuming non-degenerate  $\omega_j$ ) each term will individually satisfy a matrix equation of the form: