ILL-CONDITIONED EQUATIONS

by

DONALD LEE MYERS

B. A., Washburn University of Topeka, 1962

A MASTER'S REPORT

submitted in partial fulfillment of the

requirements for the degree

MASTER OF ARTS

Department of Mathematics

KANSAS STATE UNIVERSITY Manhattan, Kansas 1964

Approved By:

J. Flomas Parker

Major Professor



TABLE OF CONTENTS

C 2												
INTRODUCTION	•••	• •	• •	• •	• •	•••	•	•	•	•	•	1
CLASSIFICATION TECH	NIQU	ES.	•••	• •	• •	•••	•	٠	•	•	•	3
MODIFICATION OF AN I	LL-C	CON	DITI	ONE	D M.	ATR	IX	•	•	•	•	11
AN ITERATIVE TECHNI	QUE	• •	•••	• •	•••	•••	•	•	•	•	•	19
APPENDIX A	• •	•••			• •	•••	•	•	•	•	•	24
BIBLIOGRAPHY	• •				•••	•••	•	•	•		•	27
ACKNOWLEDGMENT .												29

INTRODUCTION

The purposes of this report are the examination of existing techniques of classifying a system of linear equations, and, if the system is classified as ill-conditioned, the examination of iterative methods available by use of which it is possible to obtain an accurate representation of the solution of the system.

A system of equations is said to be "ill-conditioned" if it has a solution that is extremely sensitive to slight changes in the coefficients of the variables. Such sensitivity of the coefficients impairs the accuracy and dependability of the solution obtained in the usual iterative procedures, especially those involving matrix techniques, because of round-off error which is inherent in these techniques.

Many systems of equations occur as a result of experimental work. The coefficients are subject to error resulting from measurement techniques, etc. Hence it is necessary that the stability of the system be known so that appropriate precautions may be taken to insure a reasonably accurate solution of the system.

The solution of a two-dimensional system is, in a geometric sense, represented by the intersection of the graphs of the functions. As illustration, consider the set of linear equations:





It is obvious from the graph that a slight change in the coefficients in this system may result in a relatively much greater change in the coordinates of the point of intersection of the lines.

CLASSIFICATION TECHNIQUES

Consider the equations

$$a_{11}x + a_{12}y = d_1$$

 $a_{21}x + a_{22}y = d_2$

whereas the lines corresponding to these equations are nearly parallel. The angle(s) between the lines must be nearly zero. To utilize this fact, one classification technique is the consideration of the angle between the lines. One can calculate $\cos \theta$ or $\cos^2 \theta$ and arbitrarily decide how near unity must $\cos \theta$ or $\cos^2 \theta$ be to term the system of equations ill-conditioned. Stanton (11) infers that as a reference point a criterion number could be set in the range $0.90^{\pm}.05$, and that ill-conditionedness may result if $\cos^2 \theta$ is greater than this number. This would make θ slightly less than 18 degrees. To calculate $\cos \theta$ recall,

$$\cos \theta = \frac{\begin{vmatrix} a_{11}a_{21}+a_{12}a_{22} \\ a_{11}^2+a_{12}^2 \\ a_{21}^2+a_{22}^2 \end{vmatrix}}{\sqrt{a_{21}^2+a_{22}^2}} \qquad \text{where } \theta \text{ is the angle}$$

between the lines with direction numbers $a_{11}^{}$, $a_{12}^{}$ and $a_{21}^{}$, $a_{22}^{}$. (Recall that these are also the coefficients in the two equations.) Generalized to n dimensions this becomes:

$$\cos \theta_{ij} = \frac{\left| \sum_{a_{ik}} a_{jk} \right|}{\sqrt{\sum_{a_{ik}}^2 \sqrt{\sum_{a_{jk}}^2}}}$$

This criterion is not wholly satisfactory because of the arbitrariness involved in choosing a critical value for $\cos^2 \theta$.

One of the most widely used methods of classifying a system of equations is in terms of the value of the determinant of coefficients. If the determinant is near zero, or relatively small as compared with the determinant formed by replacing the coefficients of one of the variables by the constants, the system probably is ill-conditioned. Another way of expressing this fact involves the determinant of the normalized coefficient matrix. If the absolute value of this determinant is very small as compared with unity (i.e. $|A_N| <<1$) then the system is probably ill-conditioned.

Neither of the two preceding methods is entirely satisfactory by itself. As illustration, consider the set of linear equations already discussed:

> 4.001x₁+4.012x₂=0.001 4.012x₁+4.014x₂=0.002.

The true solution of this system is $x_1 = -1$ and $x_2 = +1$. The corresponding solution obtained by method of elimination with round-off in

the twelfth position is $x_1 = -1.00000001819$ and $x_2 = +1.00000001819$. This solution agrees quite well with the true solution considering the fact that the normalized determinant of the coefficient matrix is computed to be 0.000705 which is quite small relative to unity.

A. M. Turning (12), while considering the solution of the system AX=B, suggests that instead of using the coefficient matrix A, the matrix A-S be used where S is a small predetermined matrix which he calls an average or typical matrix. The solution of this new system will be X₀+A⁻¹SX₀ where X₀ is the true solution. By averaging the effect of the "average" transformation he arrives at two classification schemes, one termed the "N-condition number of A" and the other is the "M-condition number of A." The "N-condition number is equal to $(1/n)N(A)N(A^{-1})$ where n is the dimension of A and N(A) is the norm of A. The second condition number is used in conjunction with the first. The "M-condition number" is equal to $n(M(A)M(A^{-1}))$ where n is again the dimension of A and M(A) is $\max_{ij} |a_{ij}|$. The coefficients of a well-conditioned matrix give the "N-condition number" of order $n^{1/2}$ and the "M-condition number" approximately ln(n) times larger. The larger the condition numbers, the more ill-conditioned the system. However, as in the other classification procedures, this scheme has its disadvantages. For one, it is often quite tedious to compute A⁻¹. For another, it is somewhat arbitrary as to how large

the condition numbers need be, before the system is considered illconditioned.

Among the most popular measures of the ill-conditionedness of a system are those which employ knowledge of eigenvalues and eigenvectors of the coefficient matrix, determination of which in itself is often a challenging problem. The usual criterion is the quotient $\frac{|\lambda|}{|\lambda|} \max$ The larger the resulting quotient, the more ill-conditioned

the system is. The disadvantage of this method, as previously mentioned, lies in the computation of the eigenvalues and eigenvectors.

The last classification scheme being considered in this report is attributed to Riley (9). When attempting to solve a system AX=B where A has real eigenvalues, it consists, instead, of considering the solution of the system CY=B where C=A+kI. The matrix C is "better conditioned" than is the matrix A. This will be proved by showing that if λ_1 , λ_2 , \cdots , λ_n are the eigenvalues of A then λ_1 +k, λ_2 +k, \cdots , λ_n +k are the eigenvalues of C.

To show that this is so, write the characteristic equation for C:

$$|\mathbf{r}I-C| = 0.$$

It follows, successively, that

$$|rI-(A+kI)| = 0,$$

 $|rI-kI-A| = 0,$
 $|(r-k)I-A| = 0.$

If the eigenvalues of C are λ_i , then the eigenvalues of A are λ_i -k.

Now, without loss of generality, it is possible to assume that the eigenvalue of least numerical value is positive.

Note, $\lambda_j^{+k} > \lambda_j$ where λ_j^{+k} and λ_j^{-} are the minimum eigenvalues of C and A respectively. Thus $\frac{\lambda_j^{+k}}{\lambda_i^- \lambda_j} > \frac{\lambda_j}{\lambda_i^- \lambda_j}$ where $\lambda_i^- \lambda_j^{-}$ is

positive. Hence, it follows, successively, that

$$\begin{split} & \frac{\lambda_{i} - \lambda_{j}}{\lambda_{j}^{+k}} \langle \frac{\lambda_{i} - \lambda_{j}}{\lambda_{j}}, \\ & \frac{\lambda_{i} - \lambda_{j}}{\lambda_{j}^{+k}} - \frac{\lambda_{i} - \lambda_{j}}{\lambda_{j}} < 0, \\ & \frac{\lambda_{i}}{\lambda_{j}^{+k}} - \frac{\lambda_{j}}{\lambda_{j}^{+k}} + 1 - \frac{\lambda_{i}}{\lambda_{j}} < 0, \\ & \frac{\lambda_{i}}{\lambda_{j}^{+k}} - \frac{\lambda_{j}}{\lambda_{j}^{+k}} + \frac{\lambda_{j}^{+k}}{\lambda_{j}^{+k}} - \frac{\lambda_{i}}{\lambda_{j}} < 0, \\ & \frac{\lambda_{i}}{\lambda_{j}^{+k}} + \frac{k}{\lambda_{j}^{+k}} - \frac{\lambda_{i}}{\lambda_{j}} < 0, \\ & \frac{\lambda_{i}^{+k}}{\lambda_{j}^{+k}} < \frac{\lambda_{i}}{\lambda_{j}}. \end{split}$$

Therefore

$$\frac{|\lambda_{i^{+k}}|}{|\lambda_{j^{+k}}|} < \frac{|\lambda_{i}|}{|\lambda_{j}|} .$$

Thus the matrix C is better conditioned than the matrix A.

Riley notes that his experience suggests that a reasonable choice is for k, a small positive constant, to be somewhere between 10^{2-s} and 10^{3-s} where s is the number of decimals being carried. X is computed from a series expansion in a matrix Y as follows:

 $\begin{array}{l} C=A+kI \mbox{ implies } A=C-kI. \mbox{ Thus one has} \\ AA^{-1}=(C-kI)A^{-1}. \\ \\ \mbox{ Hence it follows successively, that} \\ I=CA^{-1}-kA^{-1}, \\ C^{-1}=C^{-1}CA^{-1}-C^{-1}kA^{-1}, \\ A^{-1}=C^{-1}+kC^{-1}A^{-1}, \\ A^{-1}=C^{-1}+kC^{-1}(C^{-1}+kC^{-1}A^{-1})=C^{-1}+kC^{-2}+k^{2}C^{-2}A^{-1}, \\ A^{-1}=C^{-1}+kC^{-2}+\cdots+k^{n}C^{-(n+1)}+\cdots \mbox{ Thus} \\ X=A^{-1}B=C^{-1}B+kC^{-2}B+\cdots+k^{n}C^{-(n+1)}B+\cdots . \end{array}$

Hence,

Now, letting $Y=C^{-1}B$, it follows that

 $X=Y+kC^{-1}Y+\cdots + (kC^{-1})^{n}Y+\cdots$

To obtain an indication of the condition of the system, use a value of k of one in the next to last decimal place (k=0.000 \cdot \cdot 010). Y is then computed by an inversion technique such as the Cholesky or Square Root method. Then using Y as the right member, one computes the second term in the series expansion of X; this is then used to obtain the next member of the series. If it appears that additional computation would contribute nothing, it is likely that the system is well-conditioned, providing k was appropriately chosen. On the other hand, if it appears that additional computation would contribute to the sum of the series then the system is ill-conditioned and thus the original system must have been quite ill-conditioned. The disadvantages of the method lie in the choosing of suitable k and in the series computation of X. It should be noted that in order for the condition of C to be better than that of A, k must be approximately as large as λ_i [min. However it must not be too large, lest convergence of the series be slowed down. The rate of convergence depends on k times the eigenvalues of C^{-1} .

This technique presents a method to improve an approximate solution of the system. There exist matrices A such that A^{-1} cannot be expressed exactly carrying any specified number of decimals. However, for such matrices the inverse of the corresponding C can be expressed exactly with appropriate choice of k and can be used to improve the approximate solution. For example:

Let X_0 be an approximate solution of AX=B. Consider "CZ₀=B-AX₀".

Now

$$C(X-(X_0+Z_0)) = CX-CX_0-CZ_0$$

= AX+kX-AX_0-kX_0-CZ_0
= B-AX_0-CZ_0+k(X-X_0).

Thus $X-(X_0+Z_0)=kC^{-1}(X-X_0)$. Now letting $X_1=X_0+Z_0$ and $Z_1=C^{-1}(B-AX_1)$, repeat the process. Continuing, one notes $Z_n=C^{-1}(B-AX_n)$. Thus if $Z_n \rightarrow 0$, one has $X \rightarrow X_n$. Convergence of Z_n is highly dependent upon the appropriate choice of k, as is seen from the following development.

If one has

$$\begin{split} & \operatorname{CZ}_{n} = \operatorname{B-AX}_{n}, \\ & \operatorname{Z}_{0} = \operatorname{C}^{-1}(\operatorname{B-AX}_{0}), \\ & \operatorname{Z}_{1} = \operatorname{C}^{-1}(\operatorname{B-AX}_{1}) = \operatorname{C}^{-1}(\operatorname{B-A}(\operatorname{X}_{0} + \operatorname{Z}_{0})), \\ & \operatorname{Z}_{2} = \operatorname{C}^{-1}(\operatorname{B-AX}_{2}) = \operatorname{C}^{-1}(\operatorname{B-A}(\operatorname{X}_{1} + \operatorname{Z}_{1})), \\ & \operatorname{Z}_{3} = \operatorname{C}^{-1}(\operatorname{B-AX}_{3}) = \operatorname{C}^{-1}(\operatorname{B-A}(\operatorname{X}_{2} + \operatorname{Z}_{2})), \end{split}$$

 $Z_n = C^{-1}(B-AX_n) = C^{-1}(B-A(X_{n-1}+Z_{n-1})).$

then

Consider:

$$\begin{split} & Z_{3} = C^{-1}(B - A(X_{2} + Z_{2})) \\ & = Z_{2} - C^{-1}AZ_{2} \\ & = Z_{2} - C^{-1}A(Z_{1} - C^{-1}AZ_{1}) \\ & = Z_{2} - C^{-1}AZ_{1} + (C^{-1}A)^{2}Z_{1} \\ & = Z_{2} - C^{-1}AZ_{1} + (C^{-1}A)^{2}(Z_{0} - C^{-1}AZ_{0}) \\ & = Z_{2} - C^{-1}AZ_{1} + (C^{-1}A)^{2}Z_{0} - (C^{-1}A)^{3}Z_{0} \\ & = Z_{2} - C^{-1}AZ_{1} + (C^{-1}A)^{2}Z_{0} - (C^{-1}A)^{3}(C^{-1}(B - AX_{0})) . \end{split}$$

In general, for Z_n, one has

$$Z_n = Z_{n-1} - C^{-1}AZ_{n-2} + \cdots + (-1C^{-1}A)^{n-1}Z_0 + (-1C^{-1}A)^n (C^{-1}(B-AX_0)).$$

Now the eigenvalues of C^{-1} (k/ λ_{i}^{+k}) are numerically less than one. Thus since raising C^{-1} to the nth power also raises the eigenvalues of C^{-1} to the nth power, the eigenvalues will approach zero. Recall, if the eigenvalues of a matrix are zero, then the matrix is the zero matrix. Thus an appropriate choice of k which will allow the eigenvalues of C^{-1} to approach zero rapidly, will insure convergence of the series.

MODIFICATION OF AN ILL-CONDITIONED MATRIX

The following technique is one which presents a method of improving the condition of the coefficient matrix of a system of equations. The procedure is concerned with the eigenvalues of the coefficient matrix,

principally those which are very small. Thus, since most methods for extraction of eigenvalues obtain the dominant eigenvalue first, a technique for obtaining the smallest eigenvalue is presented in Appendix A.

As illustration, it is convenient to apply the method to a four by four system where the coefficient matrix is real and symmetric. Further, for the present it is assumed that the eigenvalues are distinct and are such that $|\lambda_1| > |\lambda_2| > |\lambda_3| > |\lambda_4|$. A method attributed to J. W. Head and G. M. Oulton (2) is as follows.

Consider the ill-conditioned system AX=D where

A=	[a11	a ₁₂	a ₁₃	a14	with a j = a ji,	X= x	and $D = \begin{bmatrix} d_1 \end{bmatrix}$
	^a 21	a ₂₂	^a 23	a ₂₄		У	d ₂
	a ₃₁	^a 32	a ₃₃	^a 34		z	d ₃
	a41	^a 42	^a 43	^a 44		[t_]	d ₄

The equation $\begin{bmatrix} a_{11} - \lambda a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} - \lambda a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} - \lambda a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} - \lambda \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ t \end{bmatrix} = 0$

is satisfied by the values λ_1 , λ_2 , λ_3 , λ_4 which allow the column vector X to have a non-trival solution. If to the non-trival solution corresponding to λ_1 the additional condition that $x_1 > 0$; $x_1^2 + y_1^2 + z_1^2 + t_1^2 = 1$

be imposed, the resulting $X_1 = \begin{bmatrix} x_1 \\ y_1 \\ z_1 \\ t_1 \end{bmatrix}$ will be a normalized eigenvector

corresponding to λ_1 . This vector furnishes the "mode", $U_1 = xx_1 + yy_1 + zz_1 + tt_1$, associated with λ_1 . In a like manner the modes corresponding to the other eigenvalues can be represented.

Since these eigenvectors X_i correspond to distinct eigenvalues of a real symmetric matrix they are orthogonal. Hence the matrix $P = \begin{bmatrix} X_1 & X_2 & X_3 & X_4 \end{bmatrix}$ is normal, orthogonal, and unitary. Thus these modes satisfy the orthornormal properties:

> 1.) $x_{r}^{2} + y_{r}^{2} + z_{r}^{2} + t_{r}^{2} = 1$ 2.) $x_{r}^{x} + y_{r}^{y} + z_{r}^{z} z_{s}^{s} + t_{r}^{t} z_{s}^{=0}$ for all $r \neq s$ 3.) $\sum_{r}^{x} z_{r}^{2} = 1$, $\sum_{r}^{y} z_{r}^{2} = 1$, $\sum_{r}^{z} z_{r}^{2} = 1$, $\sum_{r}^{t} z_{r}^{2} = 1$ 4.) $\sum_{r}^{x} x_{r}^{y} - \sum_{r}^{x} x_{r}^{z} - \sum_{r}^{x} x_{r}^{t} - \sum_{r}^{y} y_{r}^{z} - \sum_{r}^{y} y_{r}^{t} - \sum_{r}^{z} y_{r}^{t} + \sum_{r}^{z} y_{r}^{t} +$

Then the solution of the original system is:



*The limit of summation on all sums ranges from one to four.

$$\begin{array}{l} x = d_1 \sum_{\mathbf{x}_r \mathbf{x}_r}^{\mathbf{x}_r \mathbf{x}_r} + d_2 \sum_{\mathbf{x}_r \mathbf{x}_r}^{\mathbf{y}_r \mathbf{x}_r} + d_3 \sum_{\mathbf{x}_r}^{\mathbf{z}_r^2} + d_4 \sum_{\mathbf{x}_r}^{\mathbf{z}_r \mathbf{t}_r} \\ t = d_1 \sum_{\mathbf{x}_r}^{\mathbf{x}_r \mathbf{t}_r} + d_2 \sum_{\mathbf{x}_r}^{\mathbf{y}_r \mathbf{t}_r} + d_3 \sum_{\mathbf{x}_r}^{\mathbf{z}_r \mathbf{t}_r} + d_4 \sum_{\mathbf{x}_r}^{\mathbf{z}_r^2} \\ \end{array}$$

To show that this is so, consider the original system AX=D. Let $X = \sum_{r}^{a} X_{r}$ where X_{r} is the eigenvector corresponding to λ_{r} . Also let $D = \sum_{r}^{b} X_{r}$ where $b_{j} = X_{j}^{t}D$. Thus AX=D becomes $A(\sum_{r}^{a} A_{r}^{X}) = \sum_{r}^{b} X_{r}^{X}$. This is a finite sum, hence one has $\sum_{r}^{a} AX_{r}^{x} = \sum_{r}^{b} X_{r}^{X}$. Thus $\sum_{r}^{a} \lambda_{r} X_{r} = \sum_{r}^{b} X_{r}^{X}$ which implies $a_{i} = b_{i}/\lambda_{i}$. Hence $X = \sum_{r}^{b} \frac{X_{r}}{\lambda_{r}} = \sum_{r}^{T} \frac{X_{r}b_{r}}{\lambda_{r}} = \left[\frac{X_{1}}{\lambda_{1}} \frac{X_{2}}{\lambda_{2}} \frac{X_{3}}{\lambda_{3}} \frac{X_{4}}{\lambda_{4}}\right] \begin{bmatrix} b_{1} \\ b_{2} \\ b_{3} \\ b_{4} \end{bmatrix}$.

Now let
$$\Lambda = \begin{bmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ 0 & 0 & \lambda_3 & 0 \\ 0 & 0 & 0 & \lambda_4 \end{bmatrix}$$
.
Then $X = \begin{bmatrix} X_1 & X_2 & X_3 & X_4 \end{bmatrix} \Lambda^{-1} \begin{bmatrix} X_1 & X_2 & X_3 & X_4 \end{bmatrix} D$.
Thus $\begin{bmatrix} x \\ y \\ z \\ t \end{bmatrix} = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 \\ y_1 & y_2 & y_3 & y_4 \\ z_1 & z_2 & z_3 & z_4 \\ t_1 & t_2 & t_3 & t_4 \end{bmatrix} \begin{bmatrix} x_1 & y_1 & z_1 & t_1 \\ x_2 & y_2 & z_2 & t_2 \\ x_3 & y_3 & z_3 & t_3 \\ x_4 & y_4 & z_4 & t_4 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ d_3 \\ d_4 \end{bmatrix}$

From this expansion it is clear that large coefficients of d_1 are contributed mainly by terms having λ_4 in the denominator. The terms involving this eigenvalue can be isolated by multiplying the expression for x by x_4 , the expression for y by y_4 , the expression for z by z_4 , the expression for t by t_4 and adding. The result is U_4 = $xx_4+yy_4+z_4+tt_4$, $U_4=(d_1x_4+d_2y_4+d_3z_4+d_4t_4)/\lambda_4$. Head and Oulton infer that all uncertainty in the solution of the system can be attributed to the fact that U_4 is uncertain. This is a consequence of the fact that slight errors in the d_1 are greatly magnified when the expression is divided by λ_4 .

The procedure in the next step is to replace one of the equations of the original system by the expression for U_4 . As for which equation to replace, it is suggested that if the left side of any equation is approximately proportional to the expression for U_4 , then that equation is the one to be replaced. If no equation appears proportional to U_4 then the equation to replace to improve the condition of the system is arbitrary. It may be necessary to try omitting several of the original equations before the condition of the system is improved. The condition will be improved when the determinant of coefficients of the new system is greater than that of the original system. If, in the example the fourth equation is replaced by $xx_4+yy_4+zz_4+tt_4=U_4$, the

determinant of the coefficient matrix will be
$$\lambda_1 \lambda_2 \lambda_3 \begin{vmatrix} x_1 & x_2 & x_3 \\ y_1 & y_2 & y_3 \\ z_1 & z_2 & z_3 \end{vmatrix}$$

To show that this expression is the determinant of the new system,

let
$$B = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ x_4 & y_4 & z_4 & t_4 \end{bmatrix}$$
. Consider

the matrix product B $\begin{bmatrix} X_1 & X_2 & X_3 & X_4 \end{bmatrix}$. Note, since $\begin{bmatrix} X_1 & X_2 & X_3 & X_4 \end{bmatrix}$ is unitary, its determinant is one.

Now

$$BX_{1} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ x_{4} & y_{4} & z_{4} & t_{4} \end{bmatrix} \begin{bmatrix} x_{1} \\ y_{1} \\ z_{1} \\ t_{1} \end{bmatrix}$$
$$= \begin{bmatrix} a_{11}x_{1} + a_{12}y_{1} + a_{13}z_{1} + a_{14}t_{1} \\ a_{21}x_{1} + a_{22}y_{1} + a_{23}z_{1} + a_{24}t_{1} \\ a_{31}x_{1} + a_{32}y_{1} + a_{33}z_{1} + a_{34}t_{1} \end{bmatrix}$$
$$BX_{1} = \lambda_{1} \begin{bmatrix} x_{1} \\ y_{1} \\ z_{1} \end{bmatrix} \quad A \text{ similar procedure}$$

Thus

confirms that
$$BX_2 = \lambda_2 \begin{bmatrix} x_2 \\ y_2 \\ z_2 \\ 0 \end{bmatrix}$$
 and $BX_3 = \lambda_3 \begin{bmatrix} x_3 \\ y_3 \\ z_3 \\ 0 \end{bmatrix}$

However,
$$BX_4 = \begin{bmatrix} a_{11}x_4 + a_{12}y_4 + a_{13}z_4 + a_{14}t_4 \\ a_{21}x_4 + a_{22}y_4 + a_{23}z_4 + a_{24}t_4 \\ a_{31}x_4 + a_{32}y_4 + a_{33}z_4 + a_{34}t_4 \\ x_4x_4 + y_4y_4 + z_4z_4 + t_4t_4 = 1 \end{bmatrix}$$

Hence if the determinant of $B\begin{bmatrix} x_1 & x_2 & x_3 & x_4 \end{bmatrix}$ is expanded in terms of the elements of the fourth row the indicated expression is obtained. Similar expressions are obtained when one of the other equations is replaced by the expression for U_4 . These determinants have the property that

since P is orthogonal and unitary. Thus not all of the determinants can be small. Consequently, the condition of the system will be improved when the equation which corresponds to one of these minors is replaced by $U_{\mathcal{A}}$.

It should be noted that if additional data cannot be obtained to fix the value of U_4 very precisely, then useful information can be gained

with $U_{4}=0$.

If there exist two small eigen values then the homogeneous system of equations from which the eigenvectors are derived reduces to two equations with four unknowns. Thus a non-trival solution must be obtained in terms of two of the unknowns. The procedure is to find the non-trival solution with the condition $x^2+y^2+z^2+t^2=1$ and the additional condition that one of the unknowns be zero. Let this solution be denoted as x_3 , y_3 , z_3 , and t_3 . Then another non-trival solution is obtained with the condition $x^2+y^2+z^2+t^2=1$ and $xx_3+yy_3+zz_3+tt_3=0$. These two non-trival solutions yield the corresponding "modes":

$$\begin{array}{l} {}_{U_3=xx_3+yy_3+zz_3+tt_3=(x_3d_1+y_3d_2+z_3d_3+t_3d_4)/} \lambda_3 \\ {}_{U_4=xx_4+yy_4+zz_4+tt_4=(x_4d_1+y_4d_2+z_4d_3+t_4d_4)/} \lambda_4 \end{array}$$

The procedure now is the same as that for the distinct eigenvalue case with the exception that two equations of the original system are replaced by the expressions for U_3 and U_4 .

After employing the preceding technique for improving the condition of the coefficient matrix one may proceed to solve the resulting system of linear equations by the same method used to solve a system of well-conditioned equations.

AN ITERATIVE TECHNIQUE

This technique employs a modification of the coefficient matrix similar to that used by Riley in his classification technique. The procedure is, when contemplating the solution of the ill-conditioned system AX=D, to consider instead the system (A+K)X=D+KX where K=kI and k is a small positive constant. These two systems are mathematically equivalent, but the matrix A+K is better conditioned, as was previously shown.

The method is to delete the term KX and solve the resulting system (A+K)Y=D. The solution is denoted as $Y^{(1)}$. Thus if X is the true solution then X-Y⁽¹⁾=E⁽¹⁾ is the error after the first iteration. Now if one subtracts (A+K)Y=D from AX=D, the result is AX-(A+K)Y=0. Thus A(X-Y)=KY. Therefore, $AE^{(1)}=KY^{(1)}$. Hence one has obtained a new system with the same coefficient matrix A. Modifying the coefficient matrix in the manner above one then solves the system (A+K)Y⁽²⁾=KY⁽¹⁾. Then the error after the second iteration is given by $E^{(2)}=E^{(1)}-Y^{(2)}=X-Y^{(1)}-Y^{(2)}=X-Y^{(1)}-Y^{(2)}$. After m iterations, $E^{(m)}=E^{(m-1)}-Y^{(m)}=X-\sum_{i=1}^{m}Y^{(i)}$. Thus $X=\sum_{i=1}^{m}Y^{(i)}+E^{(m)}$. The formula for the determination of the $Y^{(j)}$ is (A+K)Y^(j)=D if j=1 =KY^(j-1) if j>1.

If the method converges, the solution of the original system is X= $\sum_{j=0}^{\infty} Y^{(j)}$.

To obtain conditions for convergence of the series, consider:

$$\begin{array}{c} Y^{(1)}{=}(A{+}K)^{-1}D\\ Y^{(2)}{=}(A{+}K)^{-1}KY^{(1)}\\ \vdots\\ Y^{(m)}{=}(A{+}K)^{-1}KY^{(m-1)}. \end{array}$$

Therefore, let (A+K)⁻¹D=C and (A+K)⁻¹K=B. Thus

$$y^{(m)}_{=BY}^{(m-1)}$$

 $y^{(m-1)}_{=BY}^{(m-2)}$

$$Y^{(3)} = BY^{(2)}$$

 $Y^{(2)} = BY^{(1)}$

Hence Y^(m)=B^{m-1}C. Also

If $\sum_{j=0}^{\infty}$

$$\sum_{Y} (j) = (I+B+B^{2}+\cdots+B^{m-1})C$$

$$B^{j} \text{ converges then } (A+K)X=D+KX \text{ implies}$$

$$X=(A+K)^{-1}D+(A+K)^{-1}KX$$

=C+BX.

Thus X-BX=C and X=(I-B)⁻¹C=
$$\sum_{j=1}^{m} B^{j}C$$
.
Hence, X- $\sum_{j=1}^{m} Y^{(j)} = (\sum_{j=0}^{\infty} B^{j} - \sum_{j=0}^{m-1} B^{j})C$
 $= \sum_{m}^{j=0} B^{j}C=(I-B)^{-1}B^{m}C$
 $= B^{m}(I-B)^{-1}C=E^{(m)}$.

Thus if $\sum_{j=0}^{\infty} B^j$ converges then $\lim_{m \to \infty} B^m = 0$. Therefore $E^{(m)} = 0$ and $X = \sum_{j=0}^{\infty} Y^{(j)}$.

Since the choice of k, which will insure the convergence of $\sum_{j=0}^{\infty} E^j$, is a difficult problem and then, the computation of $E^{(m)}$ can become tedious, consider an expression for the upper bound of error. This expression is developed as follows.

Define
$$(A+K)^{-1}=D=\begin{bmatrix} d_{ij} \end{bmatrix}$$
. Recall $Y^{(m)}=DKY^{(m-1)}$.

Thus
$$\mathbf{Y}^{(m)} = \begin{bmatrix} \mathbf{y}_{1}^{(m)} \\ \mathbf{y}_{2}^{(m)} \\ \vdots \\ \mathbf{y}_{n}^{(m)} \end{bmatrix} = \begin{bmatrix} \mathbf{d}_{11} & \mathbf{d}_{12} \cdot \cdot \cdot & \mathbf{d}_{1n} \\ \mathbf{d}_{21} & \mathbf{d}_{22} \cdot \cdot \cdot & \mathbf{d}_{2n} \\ \vdots & & \vdots \\ \mathbf{d}_{n1} & \mathbf{d}_{n2} \cdot \cdot \cdot & \mathbf{d}_{nn} \end{bmatrix} \begin{bmatrix} \mathbf{k}\mathbf{y}_{1}^{(m-1)} \\ \mathbf{k}\mathbf{y}_{n}^{(m-1)} \\ \mathbf{k}\mathbf{y}_{n}^{(m-1)} \end{bmatrix}$$

Hence, $y_r^{(m)} = d_{r1} k y_1^{(m-1)} + d_{r2} k y_2^{(m-1)} + \cdots + d_{rn} k y_n^{(m-1)}$; r=1,...n. Thus one has

$$\begin{split} \left| y_r^{(m)} \right| &\leq \left| \mathbf{d}_{r1} \mathbf{k} \mathbf{y}_1^{(m-1)} \right| + \left| \mathbf{d}_{r2} \mathbf{k} \mathbf{y}_2^{(m-1)} \right| + \cdots + \left| \mathbf{d}_{rn} \mathbf{k} \mathbf{y}_n^{(m-1)} \right| \\ &\leq \left| \mathbf{k} \right| \left| \sum_{i=1}^n \mathbf{d}_{ri} \right| \quad \left| \mathbf{y}_m^{(m-1)} \right| \quad \text{where } \mathbf{y}_m^{(m-1)} = \max \mathbf{y}_i^{(m-1)}. \end{split}$$

Since this holds for all r, it will hold if $|y_r^{(m)}|$ is replaced by $|y_M^{(m)}| = \left\{ |y_2^{(m)}|, \cdots, |y_n^{(m)}| \right\}$.

Therefore

Let d

$$\begin{split} \left| \mathbf{y}_{M}^{(m)} \right| &\leq \left| \mathbf{k} \right| \left\{ \sum_{i} \left| \mathbf{d}_{ri} \right| \right\} \left| \mathbf{y}_{m}^{(m-1)} \right| : \mathbf{r} = \mathbf{i}, \ \mathbf{2}, \ \dots \mathbf{n}. \end{split} \tag{c}$$

Then define a convergence constant H and restrict H so that $|\mathbf{k}| d_{M}^{=}H < 1$. This restriction can always be satisfied by choosing k sufficiently small.

From (c), for m>1, one has

$$\frac{y_{M}^{(m)}}{y_{m}^{(m-1)}} \leq H \text{ provided } y_{m}^{(m-1)} \neq 0.$$
 (d)

Hence the series $\sum_{m=1}^{\infty} \ y_M^{(m)}$ is absolutely convergent since H<1.

Considering the expressions for the individual components of the matrix X one observes

$$x_{i}^{*} = \sum_{0}^{\infty} y_{i}^{(j)} = \sum_{j=1}^{m} y_{i}^{(j)} + \sum_{j=m+1}^{\infty} y_{i}^{(j)} : i=1, 2, \dots, n.$$

Thus the error after m iterations for each x, is

$$\begin{split} \mathbf{e}_{i}^{(m)} &= \mathbf{x}_{i}^{-} \sum_{l}^{m} \mathbf{y}_{i}^{(j)} = \sum_{m+1}^{\infty} \mathbf{y}_{i}^{(j)} : i=1, 2, \dots, n. \\ \\ \text{Also } \left| \mathbf{e}_{i}^{(m)} \right| &\leq \sum_{m+1}^{i} \left| \mathbf{y}_{i}^{(j)} \right| : k=1, 2, \dots, n. \\ \\ \text{Now since } \left| \mathbf{y}_{i}^{(m)} \right| &\leq \left| \mathbf{y}_{M}^{(m)} \right| \text{ , from (d) one obtains} \end{split}$$

$$|e_{i}^{(m)}| \leq H |y_{M}^{(m)}| (1+H+H^{2}+\cdots) = |y_{M}^{(m)}| \frac{H}{1-H}$$

This technique is highly adaptable to machine computation as the coefficient matrix used when obtaining a solution of the modified system is the same for each of the $Y^{(i)}$. It is presented here merely as a supplement to the usual methods for solution of systems of equations.

APPENDIX A

In the determination of small eigenvalues only those which are less than 0.2 in absolute value are being considered. This is so since if the eigenvalue exceeds 0.2, then the maximum contribution possible to any coefficient is numerically less than 5.0.

In this procedure it is assumed that the evaluation of a determinant or the solution of a polynomial in with numerical coefficients presents no great problem.

Since attention is focused on those values of less than 0.2, the technique is to evaluate the determinant of the characteristic matrix with

= -0.2+0.4r/n, r=0, 1, ...n, where is possibly negative. The evaluation of these determinants gives rough idea of the number and position of the roots.

To illustrate the procedure, consider the following characteristic matrix:

2.557-2	2.624	2.468	2.361
2.624	3.493-X	2,351	1.532
2.468	2,351	2.557-2	2.624
2.351	1,532	2.624	3.493-2.

Since all values of λ are positive for this example, the determinant, $\Delta(\lambda)$, is evaluated with $\lambda = 0$, $\lambda = 0.04$, $\lambda = 0.08$, $\lambda = 0.12$, $\lambda = 0.16$. The corresponding values for $\Delta(\lambda)$ are

$$\Delta(0) = 0.05$$

$$\Delta(0.04) = 0.0019526$$

$$\Delta(0.08) = 0.01761476$$

$$\Delta(0.12) = 0.09101456$$

$$\Delta(0.16) = 0.2190776.$$

 $\mathrm{Now} \Delta(\lambda) \text{ for a system of four equations is a polynomial of} \\ \texttt{degree four in } \lambda \text{ . Thus } f(\lambda) \text{, defined by} \\$

$$f(\lambda) = \Delta(0) \qquad (\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.12)(\lambda - 0.16) \\ - \Delta(0.04) \frac{\lambda(\lambda - 0.08)(\lambda - 0.12)(\lambda - 0.16)}{0.04 - 0.04} \\ + \Delta(0.08) \frac{\lambda(\lambda - 0.04)(\lambda - 0.12)(\lambda - 0.16)}{0.08 - 0.04} \\ - \Delta(0.12) \frac{\lambda(\lambda - 0.04)(\lambda - 0.12)(\lambda - 0.16)}{0.12 - 0.08} \\ - \Delta(0.12) \frac{\lambda(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.16)}{0.12 - 0.08} \\ + \Delta(0.16) \frac{\lambda(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.12)}{0.16 - 0.12} \\ - \lambda(0.16) \frac{\lambda(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.12)}{0.08 - 0.04} \\ + \Delta(0.16) \frac{\lambda(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.12)}{0.16 - 0.12} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.12)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ + \lambda(0.04) \frac{\lambda(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)}{0.04} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.08)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.04)(\lambda - 0.04)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.04)(\lambda - 0.04)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.04)(\lambda - 0.04)(\lambda - 0.04)(\lambda - 0.04)} \\ - \lambda(0.04)(\lambda - 0.04)(\lambda - 0.04$$

is identical to $\triangle(\lambda)$ when λ has these five values. Hence $f(\lambda)$ will represent $\triangle(\lambda)$ very accurately for small values of λ . Thus

$$f(\lambda) = 0.05 \frac{(\lambda^4 - 0.4 \lambda^3 + 0.056 \lambda^2 - 0.0032 \lambda + 0.0006144)}{0.00006144}$$

$$- 0.0019526 \frac{(\lambda^4 - 0.36 \lambda^3 + 0.0416 \lambda^2 - 0.001536 \lambda)}{0.00001536}$$

$$+ 0.0171476 \frac{(\lambda^4 - 0.32 \lambda^3 + 0.0304 \lambda^2 - 0.000768 \lambda)}{0.00001024}$$

$$- 0.09101456 \frac{(\lambda^4 - 0.28 \lambda^3 + 0.0224 \lambda^2 - 0.000512 \lambda)}{0.00001546}$$

$$+ 0.2190776 \frac{(\lambda^4 - 0.24 \lambda^3 + 0.0176 \lambda^2 - 0.000384 \lambda)}{0.00006144}.$$

This simplifies to

$$f(\lambda) = 9449.49115\lambda^4 - 835.28266\lambda^3 + 74.19092\lambda^2$$

- 3.8285 \lambda + 0.05.

The technique at this point is to equate $f(\lambda)$ to zero and solve by an iterative measure for λ .

In this particular example the value of λ obtained is 0.05.

BIBLIOGRAPHY

- Frazer, R. A. <u>Elementary Matrices</u>, Cambridge University Press, 1938.
- Head, J. W. and Oulton, G. M. "The Solution of 'Ill-Conditioned' Linear Simultaneous Equations", Aircraft Engineering, Volume 30 (October, 1958), London: Dunhill Publications.
- Hildebrand, F. B. <u>Introduction to Numerical Analysis</u>, New York: <u>McGraw-Hill Book Company</u>, Inc., 1956.
- Householder, A. S. <u>Principles of Numerical Analysis</u>, New York: <u>McGraw-Hill Book Company</u>, Inc., 1953.
- 5. Kunz, K. S. <u>Numerical Analysis</u>, New York: McGraw-Hill Book Company, Inc., 1957.
- Lanczos, Cornelius <u>Applied Analysis</u>, Englewood Cliffs, New Jersey: Prentice-Hall, Inc., 1956.
- Nielson, K. L. <u>Methods in Numerical Analysis</u>, New York: <u>The Macmillan Company</u>, 1956.
- Perlis, Sam <u>Theory of Matrices</u>, Reading, Massachusetts: <u>Addison-Wesley Publishing Company</u>, Inc., 1958.
- 9. Riley, J. D.

"Solving Systems of Linear Equations with a Positive Definite, Symmetric, but Possibly Ill-Conditioned Matrix", <u>Math Tables and Other</u> <u>Aids to Computation</u>, Washington, D.C.: The National Research Council, July, 1955.

- Scott, E. J., Li, H. L., and Chao, B. T. On the Solution of Ill-Conditioned Simultaneous, Linear, Algebraic Equations by Machine Computation, Urbana, Illinois: University of Illinois Bulletin, 1961.
- 11. Stanton, Ralph G.

Numerical Methods for Science and Engineering, Englewood Cliffs, New Jersey: Prentice-Hall, Inc., 1961.

12. Turning, A. M.

"Rounding Off Errors in Matrix Process," Quarterly Journal of Mechanics and Applied Mathematics, Volume 1, 1948.

ACKNOWLEDGMENT

The author wishes to express his sincere thanks and appreciation to Professor S. Thomas Parker for his helpful suggestions and assistance with the preparation of this report.

ILL-CONDITIONED EQUATIONS

by

DONALD LEE MYERS

B. A., Washburn University of Topeka, 1962

AN ABSTRACT OF A MASTER'S REPORT

submitted in partial fulfillment of the

requirements for the degree

MASTER OF ARTS

Department of Mathematics

KANSAS STATE UNIVERSITY Manhattan, Kansas 1964 The purposes of this report are to define "Ill-Conditioned Equations," and examine existing methods of classifying a system of equations, develop a procedure to improve the condition of the system of equations, and present an iterative technique for use in the solution of a system of "Ill-Conditioned Equations."

A set of "Ill-Conditioned Equations" is a system of equations which has a solution that is extremely sensitive to slight changes in the coefficients of the variables. As a result, the accuracy and dependability of the solution obtained in the usual iterative procedures is impaired.

The classification techniques considered are:

- The (generalized) angle(s) between the "lines" which represent the system of equations in a geometric sense.
- The ratio of the determinant of the coefficient matrix to the determinant formed from it by replacing the coefficient of one of the variables by the constant.
- The consideration of the determinant of the normalized coefficient matrix as compared with unity.
- The use of "M-condition numbers" and "N-condition numbers."
- 5. The largest ratio of the absolute values of two eigenvalues.

 A consideration of the convergence of a series expansion of the system of equations.

A procedure developed to improve the condition of the original system is to replace one of the original equations with a "mode." These modes are derived from eigenvalues and their corresponding eigenvectors. A technique for determination of eigenvalues less in absolute value than 0.2 is presented in Appendix A.

The remainder of the report deals with an iterative technique which is highly adaptable to electronic computation. An expression for the determination of the error after m iterations is also presented.