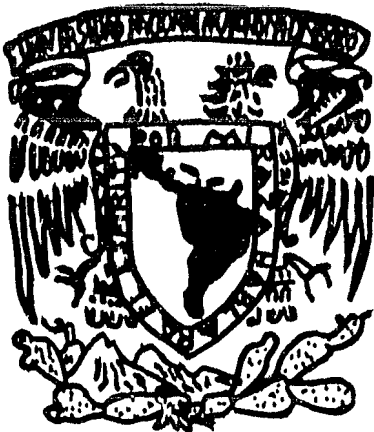


2ej. 52



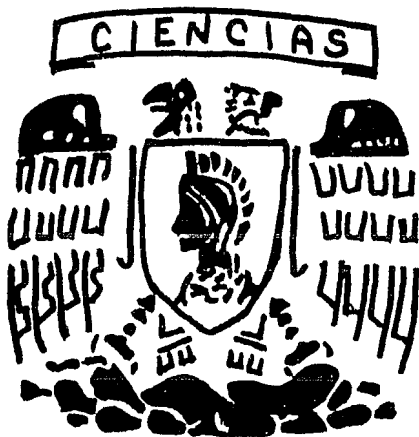
UNIVERSIDAD NACIONAL AUTONOMA DE MEXICO

FACULTAD DE CIENCIAS

ASPECTOS DE CALCULO (

DEL

ALGEBRA LINEAL



T E S I S
QUE PARA OBTENER EL TITULO DE :
A C T U A R I A
P R E S E N T A
BEATRIZ VALADEZ BAUTISTA

MEXICO, D.F.

1988



Universidad Nacional
Autónoma de México



UNAM – Dirección General de Bibliotecas Tesis Digitales Restricciones de uso

DERECHOS RESERVADOS © PROHIBIDA SU REPRODUCCIÓN TOTAL O PARCIAL

Todo el material contenido en esta tesis está protegido por la Ley Federal del Derecho de Autor (LFDA) de los Estados Unidos Mexicanos (México).

El uso de imágenes, fragmentos de videos, y demás material que sea objeto de protección de los derechos de autor, será exclusivamente para fines educativos e informativos y deberá citar la fuente donde la obtuvo mencionando el autor o autores. Cualquier uso distinto como el lucro, reproducción, edición o modificación, será perseguido y sancionado por el respectivo titular de los Derechos de Autor.

I N D I C E

INTRODUCCION	I
CAPITULO I MATERIAL ELEMENTAL DE ALGEBRA	
LINEAL	
1.1 Matrices	1
1.2 Multiplicación de una matriz por un número; la adición de matrices	3
1.3 La multiplicación de matrices	4
1.4 Asociatividad de matrices	7
1.5 Partición de matrices	8
1.6 Matrices inversas y adjuntas	12
1.7 El polinomio característico	15
1.8 Matrices semejantes	18
1.9 Transformaciones elementales	19
1.10 Descomposición de matrices en el producto de dos matrices triangulares	22
1.11 Notación matricial para un sistema de ecuaciones lineales	24
1.12 Espacio vectorial n-dimensiones	26
1.13 Dependencia lineal	28
1.14 Bases	29

Capitulo III UN ENFOQUE MAS SUTIL DE
DE LOS PROBLEMAS

3.1 Un enfoque más sutil de los
de los problemas 64

Capitulo IV ESTADO DE LA COMPUTACION DESDE
1946 HASTA AHORA Y LA NATURALEZA
DEL HARDWARE Y SOFTWARE DE LA
COMPUTACION

4.1 Estado de la computación desde 1946
hasta ahora y la naturaleza del hardware
y software de la computación 74

Capitulo V PROBLEMAS DE LA ECUACIONES
LINEALES

5.1 Problemas de las ecuaciones lineales 83

Capitulo VI INEXACTITUD INHERENTE
EN SISTEMAS LINEALES
(SOLUCIONES)

6.1 Inexactitud inherente en sistemas
lineales (soluciones) 111

Capítulo VII	PRECISION REALIZABLE EN ELIMINACION GAUSSIANA	
7.1	Presicion realizable en eliminación Gaussiana	118
Capítulo VIII	MAS SOLUCIONES EXACTAS	
8.1	Más soluciones exactas	124
Capítulo IX	ESCALONAMIENTO OPTIMO DE MATRICES	
9.1	Escalonamiento optimo de matrices	133
Capítulo X	ANALISIS DE ERROR DE REDONDEO	
10.1	Análisis de error de redondeo	143
Capítulo XI	VALORES CARACTERISTICOS DE MATRICES SIMETRICAS	
11.1	Valores caracteristicos de matrices simétricas	158
Capítulo XII	VALORES CARACTERISTICOS DE MATRICES NO SIMETRICAS	
12.1	Valores característicos de matrices no simétricas	174
	Conclusión	185

I N T R O D U C C I O N

Este estudio de aspectos del Cálculo en Algebra Lineal está dirigido a personas que sin ser especialistas en Análisis Numérico, con frecuencia se enfrentan a este tipo de cálculo en la práctica. El objetivo de este trabajo será conducir a este tipo de personas a aspectos generales del Algebra Lineal referentes a los cálculos matriciales, incluyendo normas para conocer el pocedimiento de la Eliminación Gaussiana, así mismo se darán exposiciones intrínsecas de los métodos de cálculo del Algebra Lineal.

I N T R O D U C C I O N

Este estudio de aspectos del Cálculo en Algebra Lineal está dirigido a personas que sin ser especialistas en Análisis Numérico, con frecuencia se enfrentan a este tipo de cálculo en la práctica. El objetivo de este trabajo será conducir a este tipo de personas a aspectos generales del Algebra Lineal referentes a los cálculos matriciales, incluyendo normas para conocer el pocedimiento de la Eliminación Gaussiana, así mismo se darán exposiciones intrínsecas de los métodos de cálculo del Algebra Lineal.

I N T R O D U C C I O N

C A P I T U L O

I.

M A T E R I A L E L E M E N T A L

D E

A L G E B R A L I N E A L

I) MATERIAL ELEMENTAL DE ALGEBRA LINEAL

1.1.- MATRICES.

Un conjunto de números arreglados en forma rectangular, es llamada una matriz rectangular. Este arreglo puede tener m renglones y n columnas, y se puede representar en la forma.

$$1) \quad A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix},$$

donde el subíndice m describe los renglones y el subíndice n describe las columnas; así un punto ó elemento de la matriz es localizado por renglón y columna.

Esto puede ser abreviado de la forma :

$$A = (a_{ij}) \quad (i = 1,2,\dots,m \quad ; \quad j = 1,2,\dots,n)$$

Dos matrices son iguales si sus correspondientes elementos son iguales.

Las matrices compuestas solo de un renglón son llamadas Vector Renglón.

Las matrices compuestas solo de una columna son llamadas Vector Columna. Si el número de renglones de una matriz es igual al número de columnas entonces se dice que es una matriz cuadrada.

Entre las matrices cuadradas una regla importante son las Matrices Diagonales, es decir, matrices en las cuales todos los elementos son cero a excepción de los elementos principales de la diagonal, los cuales son diferentes de cero.

$$2) \begin{vmatrix} a & 0 & \dots & 0 \\ 1 & & & \\ 0 & a & \dots & 0 \\ \cdot & 2 & & \\ \cdot & & & \\ 0 & 0 & \dots & a \\ & & & n \end{vmatrix} = \begin{pmatrix} a & & & \\ & a & & \\ & & \dots & \\ & & & a \\ & & & & a \end{pmatrix}$$

Si todos los números componentes de la diagonal son iguales, entonces se dice que la Matriz es Escalar.

$$3) \begin{vmatrix} a & 0 & \dots & 0 \\ 0 & a & \dots & 0 \\ \cdot & & & \\ \cdot & & & \\ 0 & 0 & \dots & a \end{vmatrix} = a$$

Y si a es igual a 1, se dice que es una Matriz Unitaria.

$$4) \begin{vmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \cdot & & & \\ \cdot & & & \\ 0 & 0 & \dots & 1 \end{vmatrix} = I.$$

Por último una matriz en la cual todos sus elementos son cero se le llama Matriz Nula, y la designaremos por el número cero.

El determinante cuyos elementos son los elementos de una Matriz Cuadrada, se dice que es el Determinante de esta Matriz, y escribimos el determinante de la Matriz, A como, $|A|$, o bien $d(A)$.

MULTIPLICACION DE UNA MATRIZ POR UN NUMERO; LA ADICION DE MATRICES.

Una matriz cuyos elementos son obtenidos por multiplicar todos los elementos de una Matriz A por una escalar a es llamado producto del número a y la Matriz A.

$$5) \quad ca = \begin{vmatrix} ca & ca & \dots & ca \\ 11 & 12 & & 1n \\ ca & ca & \dots & ca \\ 21 & 22 & & 2n \\ \cdot & & & \\ \cdot & & & \\ ca & ca & \dots & ca \\ m1 & m2 & & mn \end{vmatrix}$$

Una matriz C, cuyos elementos son las sumas de los elementos correspondientes de A y B, matrices que tienen igual número de columnas y renglones se dice que es la suma de A y B :

$$6) \quad A + B = \begin{vmatrix} a + b & & & a + b \\ 11 & 11 & & 1n & 1n \\ a + b & & & a + b \\ 21 & 21 & & 2n & 2n \\ \cdot & & & \\ \cdot & & & \\ a + b & & & a + b \\ m1 & m1 & & m2 & m2 & & mn & mn \end{vmatrix}$$

Las operaciones que se han introducido arriba tiene las siguientes propiedades:

$$1.- A + (B + C) = (A + B) + C$$

$$2.- A + B = B + A$$

$$3.- A + 0 = A$$

$$4.- (r + s) A = rA + sA \text{ donde } r, s \text{ son números}$$

$$5.- r(A + B) = rA + rB$$

Aquí, A, B, C, son matrices, r, s son números, generalmente.

1.3.- LA MULTIPLICACION DE MATRICES.

La multiplicación de dos matrices A y B, es definida siempre que el número de columnas de A sea igual al número de renglones de B. Así los elementos del producto $C = AB$, están definidos de la siguiente manera; el elemento en el i-ésimo renglón y la j-ésima columna de la matriz C es igual a la suma de los productos de los elementos del i-ésimo renglón de la matriz A por los elementos correspondientes de la j-ésima columna de la matriz B, así.

$$7) AB \begin{array}{c} \left| \begin{array}{cccc} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{array} \right| \left| \begin{array}{ccc} b_{11} & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ b_{m1} & b_{m2} & \dots & b_{mn} \end{array} \right| \end{array}$$

$$c_{ij} = \begin{vmatrix} c_{11} & c_{12} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2p} \\ \dots & \dots & \dots & \dots \\ c_{m1} & c_{m2} & \dots & c_{mp} \end{vmatrix} = C, \text{ donde}$$

$$(8) \quad c_{ij} = a_{i1} b_{1j} + a_{i2} b_{2j} + \dots + a_{in} b_{nj} = \sum_{ik} a_{ik} b_{kj}$$

$$(i = 1, 2, \dots, m; \quad j = 1, 2, \dots, p).$$

Notese que el producto de dos matrices rectangulares es una Matriz Rectangular, cuando no es una matriz rectangular, entonces nos da una matriz en la cual el número de renglones, es igual al número de renglones de la primera matriz (que se multiplica) y el número de columnas es igual al número de columnas de la segunda matriz es decir :

$$\begin{matrix} A & \cdot & B & = & C \\ m \cdot n & & n \cdot p & & m \cdot p. \end{matrix}$$

Asi tenemos que el producto de una matriz cuadrada por una matriz compuesta por una sola columna, nos da una matriz de una columna.

El producto de las matrices AB y BA se puede hacer solo si el número de renglones de la primera matriz es igual al número de columnas de la segunda, y el número de columnas de la primera es igual al número de renglones de la segunda, dada ésta última condición, se sigue que las matrices deben de ser cuadradas, así el producto AB y BA, es posible siempre que A y B sean matrices cuadradas, en caso contrario y en general tenemos que $AB \neq BA$.

En casos particulares la multiplicación puede ser conmutativa, y por ejemplo, las matrices conmutan con cualquier matriz cuadrada del mismo orden, así.

$$\begin{aligned}
 & \begin{vmatrix} a & 0 & \dots & 0 \\ 0 & a & \dots & 0 \\ \cdot & & & \\ 0 & 0 & \dots & a \end{vmatrix} \quad \begin{vmatrix} a & a & \dots & a \\ 11 & 12 & & 1n \\ a & a & \dots & a \\ 21 & 22 & & 2n \\ \cdot & & & \\ a & a & \dots & a \\ n1 & n2 & & nn \end{vmatrix} \\
 = & \begin{vmatrix} a & a & \dots & a \\ 11 & 12 & & 1n \\ a & a & \dots & a \\ 21 & 22 & & 2n \\ \cdot & \cdot & \cdot & \cdot \\ a & a & \dots & a \\ n1 & n2 & & nn \end{vmatrix} \quad \begin{vmatrix} a & 0 & \dots & 0 \\ 0 & a & \dots & 0 \\ \cdot & & & \\ 0 & 0 & \dots & a \end{vmatrix} \\
 = & \begin{vmatrix} aa & aa & \dots & aa \\ 11 & 12 & & 1n \\ aa & aa & \dots & aa \\ 21 & 22 & & 2n \\ \cdot & \cdot & \cdot & \cdot \\ aa & aa & \dots & aa \\ n1 & n2 & & nn \end{vmatrix}
 \end{aligned}$$

Así se sigue que para una Matriz unitaria la multiplicación por otra Matriz (cuadrada) es conmutativa, esto es,

$$AI = IA = A .$$

ASOCIATIVIDAD DE MATRICES.- La multiplicación de Matrices es asociativa, si AB y $(AB)C$, es decir, si multiplicamos primero AB y luego efectuamos la multiplicación por C , también se puede hacer primero la operación BC y después hacer la multiplicación $A(BC)$.

entonces,

$$i) \quad A(BC) = (AB)C$$

El producto de matrices tiene la siguientes propiedades :

$$ii) \quad a(AB) = (aA)B = A(aB)$$

$$iii) \quad (A+B)C = AC + BC$$

$$iv) \quad C(A+B) = CA + CB$$

Donde A, B, C son matrices y a un número.

TRANSPUESTA.- La transpuesta de una Matriz se obtiene intercambiando los renglones y columnas de una matriz A .

$$A = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{pmatrix} = (a_{ij})$$

entonces obtenemos la transpuesta de la matriz como :

$$A^t = A^t = (a_{ji}) = \begin{pmatrix} a_{11} & a_{21} & \dots & a_{m1} \\ a_{12} & a_{22} & & a_{m2} \\ \vdots & \vdots & & \vdots \\ a_{1n} & a_{2n} & & a_{mn} \end{pmatrix}$$

Así, tenemos la siguiente regla para las matrices transpuestas,

$$(AB)'' = B' A'$$

En conclusión podemos remarcar que el determinante de un producto de matrices cuadradas es igual al producto de los determinantes de las matrices multiplicadas :

$$|AB| = |A||B|,$$

y este resultado es tomado de la teoría de determinantes.

PARTICION DE MATRICES.- Al tomar matrices de orden grande, requiere como regla un gran número de operaciones, y es conveniente reducir los cálculos de matrices de orden grande a cálculos con matrices de menor orden. Tal reducción puede ser efectuada por particionar las matrices dadas; cada matriz puede ser concebida como compuesta de varias matrices de menor orden, y esta subdivisión puede ser llevada a cabo de varias maneras, por ejemplo,

$$\begin{vmatrix}
 a & a & a & a \\
 11 & 12 & 13 & 14 \\
 a & a & a & a \\
 21 & 22 & 23 & 24 \\
 a & a & a & a \\
 31 & 32 & 33 & 34
 \end{vmatrix}
 =
 \begin{vmatrix}
 a & . & a & a & a \\
 11 & . & 12 & 13 & 14 \\
 . & . & . & . & . \\
 . & . & . & . & . \\
 a & . & a & a & a \\
 21 & . & 22 & 23 & 24 \\
 . & . & . & . & . \\
 a & . & a & a & a \\
 31 & . & 32 & 33 & 34 \\
 . & . & . & . & .
 \end{vmatrix}
 =
 \begin{vmatrix}
 a & a & . & a & a \\
 11 & 12 & . & 13 & 14 \\
 . & . & . & . & . \\
 . & . & . & . & . \\
 a & a & . & a & a \\
 21 & 22 & . & 23 & 24 \\
 . & . & . & . & . \\
 a & a & . & a & a \\
 31 & 32 & . & 33 & 34
 \end{vmatrix}$$

Las matrices es las cuales se subdivide la matriz dada son llamadas submatrices.

Las operaciones básicas con matrices particionadas en las cuales las matrices diagonales son del mismo orden, se relacionan de manera natural con las matrices, entonces si tenemos.

$$A = \begin{vmatrix} A_{11} & A_{12} & \dots & A_{1K} \\ A_{21} & A_{22} & \dots & A_{2K} \\ \cdot & \cdot & \cdot & \cdot \\ A_{K1} & A_{K2} & \dots & A_{KK} \end{vmatrix}$$

y

$$B = \begin{vmatrix} B_{11} & B_{12} & \dots & B_{1K} \\ B_{21} & B_{22} & \dots & B_{2K} \\ \cdot & \cdot & \cdot & \cdot \\ B_{K1} & B_{K2} & \dots & B_{KK} \end{vmatrix}$$

Donde A_{ii} y B_{ii} son matrices cuadradas del mismo orden, entonces

$$A + B = \begin{vmatrix} A_{11} + B_{11} & A_{12} + B_{12} & \dots & A_{1K} + B_{1K} \\ A_{21} + B_{21} & A_{22} + B_{22} & \dots & A_{2K} + B_{2K} \\ \cdot & \cdot & \cdot & \cdot \\ A_{K1} + B_{K1} & A_{K2} + B_{K2} & \dots & A_{KK} + B_{KK} \end{vmatrix}$$

Un caso especial muy importante de una matriz particionada, es la matriz bordeada. Teniendo una matriz cuadrada de orden $n-1$:

$$A_{n-1} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1,n-1} \\ a_{21} & a_{22} & \dots & a_{2,n-1} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n-1,1} & a_{n-1,2} & \dots & a_{n-1,n-1} \end{bmatrix}$$

formamos una matriz cuadrada de orden n , A_n , al agregar a la matriz A_{n-1} un renglón : $v = (a_{n1}, \dots, a_{n,n-1})$, una columna : $u = (a_{1n}, \dots, a_{n-1,n})$, y un número a_{nn} :

$$A_n = \begin{bmatrix} A_{n-1} & u_{n-1} \\ v_{n-1} & a_{nn} \end{bmatrix} = \begin{bmatrix} A_{n-1} & u_{n-1} \\ v_{n-1} & a_{nn} \end{bmatrix}$$

Diremos que la matriz A_n ha sido obtenida por bordear la matriz A_{n-1} y la matriz A_n es particionable naturalmente.

Las operaciones bajo matrices bordeadas son realizadas de acuerdo con las reglas generales de operaciones con matrices particionadas.

sean

$$A = \begin{vmatrix} M & u \\ v & a \end{vmatrix} \quad \text{y} \quad B = \begin{vmatrix} P & y \\ x & b \end{vmatrix}$$

dos matrices bordeadas de orden n , entonces los siguientes exposiciones son válidas;

$$i) aA = \begin{vmatrix} aM & au \\ av & aa \end{vmatrix}$$

$$ii) A + B = \begin{vmatrix} M + P & u + y \\ v + x & a + b \end{vmatrix}$$

$$iii) AB = \begin{vmatrix} MP + ux & My + ub \\ vP + ax & vy + ab \end{vmatrix}$$

aquí MP y ux son matrices de orden $(n-1)$; My , ub , son columnas compuestas de $n-1$ elementos; vP y ax son renglones análogos y, por último $vy + ab$, es un número.

Matriz Quasi-diagonal, son de la forma :

$$\begin{bmatrix} A & \cdot & 0 & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & B & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & 0 & \cdot & C \end{bmatrix}$$

Donde las celdas de la matriz son obviamente,

$$A = \begin{bmatrix} a & a \\ 11 & 12 \\ a & a \\ 21 & 22 \end{bmatrix} ; B = \begin{bmatrix} b & b & b \\ 11 & 12 & 13 \\ b & b & b \\ 21 & 22 & 23 \\ b & b & b \\ 31 & 32 & 33 \end{bmatrix} ; ; C = \begin{bmatrix} c & c \\ 11 & 12 \\ c & c \\ 21 & 22 \end{bmatrix}$$

y las otras seis celdas son nulas.

El determinante de una matriz Quasi diagonal es el producto de los determinantes de las celdas diagonales.

El producto de dos matrices Quasi-diagonales es una matriz quasi-diagonal

MATRICES INVERSAS Y ADJUNTA.- Si una matriz cuadrada $A = (a_{ij})$ se dice que es no singular si su determinante es distinto de cero, en caso contrario la matriz A es singular.

Un concepto muy importante de una matriz inversa es ahora introducido.

Una matriz B es llamada la inversa ó reciproca de una matriz A si

$$AB = I.$$

Una condición necesaria y suficiente para la existencia de la matriz inversa es la no singularidad de la matriz A.

entonces, $AB = I, |A| |B| = |I|$, y consecuentemente $|A| \neq 0$

Supongamos ahora que $A \neq 0$, para construir la inversa de la matriz primero definamos la adjunta de una matriz. i.e., la matriz.

$$C = \begin{bmatrix} A_{11} & A_{21} & \dots & A_{n1} \\ A_{12} & A_{22} & \dots & A_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ A_{1n} & A_{2n} & \dots & A_{nn} \end{bmatrix}$$

Aquí A_{ij} es el complemento algebraico (cofactor de los elementos a_{ij}) en el determinante de la matriz A , i.e., es designado por el menor determinante de los elementos a_{ij} notese que esta puesta en posición transpuesta.

y tiene la siguiente propiedad (la matriz adjunta)

$$AC = A I$$

De esta igualdad se sigue, $B = \frac{1}{A} C$

como A es no singular entonces, se tiene ,

$$AB = A \frac{1}{A} C = C = \frac{1}{A} AC = I$$

La matriz así construida tiene también la propiedad.

$$BA = I$$

La matriz inversa es única, si esta existe.

Polinomios en una matriz.- Definiremos la potencia entera positiva de una matriz cuadrada, haciendo.

$$A.A.A.A.A. \dots A = A^n$$

n veces

en vista de la ley asociativa, es evidente de la definicion que,

$$A^{n \cdot m} = A^{n \cdot m}$$

$$(A^m)^n = A^{nm}$$

de aqui se sigue que las potencias de una misma matriz son conmutativas.

Y por definición tenemos que

$$A^0 = I$$

Una expresion de la forma

$$\alpha_0 A^n + \alpha_1 A^{n-1} + \dots + \alpha_n I$$

donde $\alpha_0, \alpha_1, \dots, \alpha_n$, son números complejos, es llamada un Polinomio en una matriz ó bien Matriz Polinomial. Esta matriz polinomial puede ser tomado, como el resultado de retomar la variable γ en un polinomio algebraico.

$$p(\gamma) = \alpha_0 \gamma^n + \alpha_1 \gamma^{n-1} + \dots + \alpha_n$$

Por la matriz A.

Es importante notar que las reglas de operación para Matrices Polinomiales no difieren de las reglas de operación para polinomios algebraicos, :

$$\begin{aligned} \text{Dado } p(\gamma) &= \#(\gamma) \pm X(\gamma) \\ w(\gamma) &= \#(\gamma)X(\gamma) \end{aligned}$$

$$\begin{aligned} \text{entonces } p(A) &= \#(A) \pm X(A) \\ w(A) &= \#(A)X(A) \end{aligned}$$

Esto se sigue de la conmutatividad de las potencias de una matriz.

1.7.- EL POLINOMIO CARACTERISTICO.

El teorema de Cay Ley - HAMILTON. El polinomio mínimo La ecuación.

$$\begin{vmatrix} a_{11} - \gamma & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} - \gamma & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} - \gamma \end{vmatrix} = 0$$

Es llamada la ecuación característica, ó singular de la matriz $A = (a_{ij})$

El miembro izquierdo de esta ecuación, el cual puede ser escrito en la forma abreviada $(A - \gamma I)$, recibe el nombre de polinomio característico la función característica de la matriz. Las ecuaciones características son frecuentemente encontradas en matematicas aplicadas.

El cálculo directo de la función característica presenta dificultades técnicas considerables.

$$P(\gamma) = \det(A - \gamma I) = (-1)^n (\gamma^n - p_1 \gamma^{n-1} - p_2 \gamma^{n-2} - \dots - p_n)$$

Entonces,

$$p_1 = a_{11} + a_{22} + \dots + a_{nn}$$

$$p_n = (-1)^n \det A$$

Y los coeficientes restantes p_k son las sumas tomadas con el signo $(-1)^k$ de todos los menores principales del

determinante de la matriz A de orden K es decir, de los menores involucrados en la diagonal principal. El número de tales menores es igual al número de combinaciones de n cosas tomadas de K en K.

Las raíces de la ecuación característica, son llamados los números propios (de raíces) características valores propios, ó eigen valores de la matriz A, del bien conocido teorema de Vieta dando la relación entre las raíces de una ecuación y sus coeficientes tenemos.

$$\gamma_1 + \gamma_2 + \dots + \gamma_n = p_1 = a_{11} + a_{22} + \dots + a_{nn}$$

$$\gamma_1 \gamma_2 \dots \gamma_n = (-1)^{n-1} p_n = A_n$$

La cantidad $p_1 = a_{11} + a_{22} + \dots + a_{nn}$ es llamada la traza

de la matriz A, y se denota con $\text{tr}A$.

Para cualquier matriz cuadrada la siguiente relación conocida como el teorema CAYLEY - HAMILTON, obtenemos que si $\mu(\gamma)$ es el polinomio característico de la matriz A, entonces $\mu(A) = 0$, que es en un sentido que la matriz es una raíz de su propia ecuación característica.

Para probarlo consideremos la matriz B, la adjunta de la matriz $A - \gamma I$, dado que cada cofactor en el determinante $|A - \gamma I|$ es un polinomio en γ de grado que no exceda $n-1$, la matriz adjunta puede ser representada como un polinomio algebraico, con coeficientes matriciales, es decir, en la forma.

$$B = B_{n-1} + B_{n-2} \gamma + \dots + B_0 \gamma^{n-1}$$

donde B_{n-1}, \dots, B_0 , son ciertas matrices no dependientes de γ .

En lo fuerte de la propiedad fundamental de la matriz adjunta tenemos

$$\begin{aligned}
 (B_{n-1} + B_{n-2} \gamma + \dots + B_0 \gamma^{n-1}) (A - \gamma I) &= |A - \gamma I| I \\
 &= (-1)^n (\gamma^n - p_1 \gamma^{n-1} - \dots - p_n) I
 \end{aligned}$$

esta ecuación es equivalente al sistema de ecuaciones

$$\begin{array}{l}
 \left[\begin{array}{l}
 B_{n-1} A = (-1)^{n+1} p_n I \\
 B_{n-2} A - B_{n-1} = (-1)^{n+1} p_{n-1} I \\
 \dots \\
 B_0 A - B_{n-1} = (-1)^{n+1} p_1 I \\
 -B_0 = (-1)^n I
 \end{array} \right.
 \end{array}$$

multiplicando estas ecuaciones por la derecha por $I, A, A, \dots, A^{n-1}, A^n$, respectivamente, y sumando obtenemos una matriz nula en el lado izquierdo y en la derecha.

$$(-1)^n (-p_n I - p_{n-1} A - p_{n-2} A^2 - \dots - p_1 A^n) = \mu(A) I$$

Así $\mu(A) = 0$ la cual es lo que se quería probar.

La relación CAYLEY-HAMILTON muestra que para una matriz cuadrada dada un polinomio existe, para la cual es una

raíz. Evidentemente tal polinomio no es Único si $\phi(\gamma)$, tiene tal propiedad, tiene cualquier polinomio divisible por $\phi(\gamma)$ el polinomio de menor grado teniendo la propiedad de que la matriz A es un cero del polinomio, y es llamado el polinomio mínimo de A.

En seguida se probará que el polinomio característico es divisible por el polinomio mínimo.

Sean $q(\gamma)$ $r(\gamma)$, el coeiciente y residuo obtenidos al dividir el polinomio característico $\mu(\gamma)$ por el polinomio mínimo $\phi(\gamma)$:

$$\mu(\gamma) = \phi(\gamma) q(\gamma) + r(\gamma),$$

El grado $r(\gamma)$ siendo por supuesto menor que el grado de $\phi(\gamma)$, sustituyendo A por γ en esta ecuación tenemos

$$r(A) = \mu(A) - \phi(A) q(A) = 0$$

Así la matriz A prueba que es un cero del polinomio $r(\gamma)$ de aquí se sigue que $r(\gamma) = 0$ dado que de otro modo $\phi(\gamma)$ no sería el polinomio mínimo, consecuente mente $\phi(\gamma)$ divide a $\mu(\gamma)$.

1.8 Matrices Semejantes. Decimos que la matriz B, es semejante a la matriz A si existe una matriz C no singular tal

$$\text{que } B = C^{-1} A C$$

La matriz B se puede obtener de la matriz A por una transformación semejante.

La transformación semejanza tiene las siguientes propiedades.

$$1.- C^{-1} A C + C^{-1} A C + \dots + C^{-1} A C = C^{-1} (A + A + \dots + A) C.$$

$$2.- C^{-1} A C C^{-1} A C \dots C^{-1} A C = C^{-1} (A A \dots A) C.$$

En particular, $C^{-1} A C = C^{-1} A C$.

De aquí,

3.- $f(C^{-1} A C) = C^{-1} f(A) C$, para cualquier polinomio $f(\lambda)$

de la última propiedad se sigue directamente que las matrices semejantes tienen la misma función mínima.

Mostraremos que las matrices semejantes tienen también la misma función característica.
tenemos

$$\begin{aligned} |B - \lambda I| &= |C^{-1} A C - \lambda I| = |C^{-1} A C - \lambda C^{-1} I C| \\ &= \begin{bmatrix} C^{-1} \\ C \end{bmatrix}^{-1} \begin{bmatrix} A - \lambda I \end{bmatrix} \begin{bmatrix} C \end{bmatrix} = \begin{bmatrix} A - \lambda I \end{bmatrix}. \end{aligned}$$

1.9.- Transformaciones Elementales.

Frecuentemente es necesario efectuar las siguientes operaciones con matrices.

a).- Multiplicación de los elementos de algunos renglones por un número.,

b).- Sumando a los elementos de algún renglón, números proporcionales a los elementos de algún renglón siguiente.

b").- Sumando a los elementos de algún renglón números proporcionales a los elementos de algún renglón precedente.

Algunas veces estas transformaciones pueden ser hechas en las columnas. Las transformaciones del tipo indicado son llamadas transformaciones elementales de la matriz.

Cualquier transformación elemental de los renglones es equivalente a una Premultiplicación de la matriz (Multiplicación por la izquierda) por una matriz no singular de una forma especial.

(33) ahora efectuando las operaciones B en las columnas y
 (34) efectuando B.

Ejemplos de operaciones por renglones.

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & \alpha & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & b & c \\ x & y & z \\ u & v & w \end{bmatrix} = \begin{bmatrix} a & b & c \\ \alpha x & \alpha y & \alpha z \\ u & v & w \end{bmatrix} ;$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & \alpha & 0 \end{bmatrix} \begin{bmatrix} a & b & c \\ x & y & z \\ u & v & w \end{bmatrix} = \begin{bmatrix} a & b & c \\ x & y & z \\ u + \alpha x & v + \alpha y & w + \alpha z \end{bmatrix} ;$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & \alpha \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} a & b & c \\ x & y & z \\ u & v & w \end{bmatrix} = \begin{bmatrix} a & b & c \\ x + \alpha u & y + \alpha v & z + \alpha w \\ u & v & w \end{bmatrix}$$

En el trabajo futuro tendremos que ejecutar transformaciones de los tipos a y b' sobre matrices.

El resultado de varias transformaciones de ese tipo es equivalente a la premultiplicación de la matriz por alguna matriz triangular, es decir, una de la forma

$$\begin{bmatrix} \# & 0 & \dots & 0 \\ & 1 & & \\ \# & \# & \dots & 0 \\ & 21 & 22 & \\ \# & \# & \dots & \# \\ & n1 & n2 & \dots & nn \end{bmatrix}$$

Con elementos diagonales $\# \neq 0$, es una matriz triangular inferior

1.10.- Descomposición de matrices en el producto de dos matrices triangulares.

Las matrices triangulares, esto es, matrices de la forma.

$$\begin{bmatrix} c_{11} & 0 & \dots & 0 \\ c_{21} & c_{22} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ c_{n1} & c_{n2} & \dots & c_{nn} \end{bmatrix} \quad \text{y} \quad \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ 0 & b_{22} & \dots & b_{2n} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & b_{nn} \end{bmatrix}$$

con $c_{ii} \neq 0, b_{ii} \neq 0$

El producto de dos matrices triangulares de la misma estructura nos da otra vez una matriz triangular de la misma forma, una matriz no singular triangular es invertible facilmente y su inversa es de la misma forma.

Los siguientes teoremas son por consiguiente de interes.

TEOREMA. Una condición para que las submatrices resultantes de la matriz A sean no singulares es decir que, A pueda ser representada como el producto de una matriz triangular superior y una matriz triangular inferior,

Demostración :

Tenemos por inducción que para $n = 1$ es obvio que $\begin{pmatrix} a_{11} \end{pmatrix} = \begin{pmatrix} b_{11} \end{pmatrix} \begin{pmatrix} c_{11} \end{pmatrix}$, y uno de los factores puede ser tomado arbitrariamente.

Luego para $(n - 1)$ el teorema es verdadero, para una matriz de orden $(n - 1)$.

Ahora se mostrará que para una matriz de orden n la partición de la matriz A es una matriz bordeada.

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ a_{21} & \dots & a_{2n} \\ \dots & \dots & \dots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} = \begin{bmatrix} & & & a_{1n} \\ & A_{n-1} & & \\ & & & \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} = \begin{bmatrix} A_{n-1} & u \\ v & a_{nn} \end{bmatrix}$$

Buscamos una descomposición $A = CB$ de la matriz A en el producto de dos matrices B y C de las formas requeridas, primero tomando estas matrices particionadas en la forma bordeada como la de A :

$$C = \begin{bmatrix} C_{n-1} & 0 \\ x & c_{nn} \end{bmatrix} \quad B = \begin{bmatrix} B_{n-1} & y \\ 0 & b_{nn} \end{bmatrix}$$

Por regla de multiplicación para matrices particionadas.

$$CB = \begin{bmatrix} C_{n-1} & 0 \\ x & c_{nn} \end{bmatrix} \begin{bmatrix} B_{n-1} & y \\ 0 & b_{nn} \end{bmatrix} = \begin{bmatrix} C_{n-1} B_{n-1} & C_{n-1} y \\ x B_{n-1} & x y + c_{nn} b_{nn} \end{bmatrix} = A$$

De aquí tenemos:

$$\begin{bmatrix} C_{n-1} & B_{n-1} \end{bmatrix} = A_{n-1}$$

Ahora bien, tales matrices triangulares, C_{n-1} y B_{n-1} ,

Existen por hipótesis de inducción.

Además la suposición de que $|A_{n-1}| \neq 0$ se sigue que

$$|C_{n-1}| \neq 0 \quad \text{y} \quad |B_{n-1}| \neq 0$$

Utilizando la definición de multiplicación de matrices, el sistema puede ser escrito como una simple ecuación en notación matricial.

$$(2) \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \dots \\ b_n \end{bmatrix}$$

ó simplemente como $Ax = b$

Donde A es la matriz de coeficientes del sistema, b la columna de términos constantes, y x la columna cuyos elementos son las incógnitas.

Si la matriz del sistema es no singular, obtenemos la solución del sistema (1) por premultiplicar (2) por

$$A^{-1}$$

$$x = A^{-1} b = \frac{1}{|A|} B b$$

Donde B es la matriz adjunta de A y esta última expresión es la notación matricial para la Regla de Cramer.

$$x_i = \frac{|A_i|}{|A|}$$

Donde A_i es la matriz que se obtiene de A por sustituir los elementos a_{ij} de la i -ésima columna por las componentes b_i de b .

Realmente la ecuación matricial.

$$x = A^{-1} b = \frac{1}{|A|} Bb$$

es equivalente a las n-ecuaciones.

$$x_i = \frac{A_{i1} b_1 + A_{i2} b_2 + \dots + A_{in} b_n}{A} \quad \text{con } i=1, \dots, n$$

Dado que los A_{ki} son los cofactores de los elementos a_{ki} en el determinante de la matriz A. Tenemos obviamente .

$$A_{i1} b_1 + A_{i2} b_2 + \dots + A_{in} b_n = |A|_i$$

y con esto queda probada nuestra exposición.

1.12.- Espacio Vectorial n-dimensional.

Un punto X de éste espacio es un arreglo de n números reales en un orden definido.

$$X = (x_1, x_2, \dots, x_n)$$

X es también llamado vector n-dimensional y los números x_1, x_2, \dots, x_n son las componentes del vector y n es la dimensión del espacio.

Dos vectores son iguales solo si sus componentes correspondientes son iguales. Las operaciones fundamentales con vectores están definidas a continuación.

$$\text{si } X = (x_1, x_2, \dots, x_n) \quad \text{y } Y = (y_1, y_2, \dots, y_n)$$

Son dos vectores n-dimensionales y a es un número complejo arbitrario entonces tenemos por definición.

$$X + Y = (x_1 + y_1, x_2 + y_2, \dots, x_n + y_n)$$

$$aX = (ax_1, ax_2, \dots, ax_n)$$

La suma de vectores sastiface las leyes asociativa y conmutativa .

$$X + Y = Y + X$$

$$(X + Y) + Z = X + (Y + Z)$$

En el espacio n-dimensional, el producto escalar está definido de acuerdo con la fórmula.

$$(X, Y) = \sum_{k=1}^n X_k \bar{Y}_k$$

Donde \bar{Y}_k designa el conjugado complejo de Y_k

y el producto escalar tiene las siguientes propiedades.

1).- $(X, X) \neq 0$ si, $X \neq 0$; $(X, X) = 0$, si $X = 0$

2).- $(X, Y) = \overline{(Y, X)}$

3).- $(X_1 + X_2, Y) = (X_1, Y) + (X_2, Y)$

4).- $(aX, Y) = a (X, Y)$

5).- $(X, Y_1 + Y_2) = (X, Y_1) + (X, Y_2)$

6).- $(X, aY) = \bar{a} (X, Y)$

En resumen $\sqrt{(x, x)}$ es llamada la longitud del vector, en lo que sigue lo designaremos como $\|x\|$. También el espacio n-dimensional complejo, introducido arriba, es usado para considerar el espacio real n-dimensional, es decir arreglos de vectores con componentes reales.

En el espacio real, el producto escalar es igual a la suma de los productos de las componentes correspondientes de los vectores.

La longitud de un vector es igual a la raíz cuadrada de la suma de los cuadrados de sus componentes. A menudo tendremos que trabajar con el espacio real n -dimensional, mientras que el espacio complejo solo en algunas ocasiones se requiere.

1.13.- Dependencia Lineal.

Un vector $Y = c_1 X_1 + c_2 X_2 + \dots + c_m X_m$, decimos que es una combinación lineal de los vectores X_1, X_2, \dots, X_m .

Es fácil ver que si los vectores Y_1, \dots, Y_k , son combinaciones lineales de los vectores X_1, \dots, X_m , cualquier combinación lineal $\mu_1 Y_1 + \dots + \mu_k Y_k$ será también una combinación de los vectores X_1, X_2, \dots, X_m .

Los vectores X_1, X_2, \dots, X_m , se dice que son linealmente dependientes, si existen constantes c_1, c_2, \dots, c_m no todas cero, tales que la ecuación

$$c_1 X_1 + c_2 X_2 + \dots + c_m X_m = 0 \quad (5)$$

se cumple. Sin embargo, si esta ecuación se cumple solo cuando todas las constantes c_i son iguales a cero los vectores X_1, X_2, \dots, X_m se dice que son linealmente independientes.

Si los vectores X_1, \dots, X_m son linealmente dependientes, entonces al menos uno de ellos será una combinación lineal del resto.

Por ejemplo, $c \neq 0$, encontramos de (5)

$$X_m = \frac{c_1}{c_m} X_1 - \dots - \frac{c_{m-1}}{c_m} X_{m-1}$$

1.14.- Bases.

Un sistema de vectores linealmente independiente constituye una base para un espacio si cualquier vector del espacio es una combinación lineal de los vectores del sistema.

Un ejemplo de una base, es el conjunto de vectores.

$$\left\langle \begin{array}{l} e_1 = (1, 0, \dots, 0) \\ e_2 = (0, 1, \dots, 0) \\ \dots \\ e_n = (0, 0, \dots, 1) \end{array} \right\rangle$$

Y es obvio que para cualquier vector $X = (x_1, x_2, \dots, x_n)$ tenemos que

$$X = x_1 e_1 + x_2 e_2 + \dots + x_n e_n$$

Y a esta expresión la llamaremos, la base inicial del espacio ó base canónica. Y dicha base no es la única posible, sino todo lo contrario. El número de vectores que forma una base no depende de la selección de la misma. La prueba de esto es. Sean Y_1, \dots, Y_k y Z_1, \dots, Z_m , dos

bases, donde $k > m$.

Los vectores Y_1, \dots, Y_k , son combinación de los lineal de los vectores Z_1, \dots, Z_m , entonces Y_1, \dots, Y_k , son

linealmente dependientes, lo cual contradice la definición de base, dado que $k = n$.

Adicionalmente, dado que la base inicial está constituida por n vectores, cualquier otra base consistirá también de n vectores. El número de vectores que forma una base, coincide con la dimensión del espacio.

Sea U_1, \dots, U_n , los vectores que forman una base para el espacio. Cualquier vector X será entonces una combinación lineal de U_1, \dots, U_n :

$$X = \epsilon_1 U_1 + \epsilon_2 U_2 + \dots + \epsilon_n U_n.$$

Los coeficientes en esta expresión únicamente definen al vector X , para

$$X = \epsilon_1 U_1 + \dots + \epsilon_n U_n = \epsilon'_1 U'_1 + \dots + \epsilon'_n U'_n,$$

$$\text{entonces } (\epsilon_1 - \epsilon'_1)U_1 + \dots + (\epsilon_n - \epsilon'_n)U_n = 0$$

y por lo tanto

$$\epsilon_1 - \epsilon'_1 = 0, \dots, 0, \epsilon_n - \epsilon'_n = 0$$

En vista de la independencia lineal de los vectores U_1, \dots, U_n los coeficientes $\epsilon_1, \dots, \epsilon_n$, son llamados las coordenadas del vector X con respecto a la base U_1, \dots, U_n . Nótese que las componentes de un vector x_1, \dots, x_n son las coordenadas del vector X con respecto a la base inicial.

1.15.- Sistemas Ortogonales de Vectores.

Los vectores distintos de cero de un espacio se dice que son ortogonales si su producto escalar es igual a cero. Un sistema de vectores se dice que es ortogonal si cualquiera dos vectores son ortogonales mutuamente. Al referirnos a un sistema ortogonal, debemos asumir que todos los vectores de este sistema son diferentes de cero.

Teorema.- Los vectores que forman un sistema ortogonal son linealmente independientes.

Teorema.- Sean X_1, \dots, X_n , linealmente independientes.

Un sistema ortogonal de vectores X'_1, \dots, X'_n puede ser construido si este es relacionado con el conjunto original por las relaciones :

$$X'_1 = X_1$$

$$X'_2 = X_2 + \sigma_{21} X_1$$

.....

$$X'_k = X_k + \alpha_{k1} X_1 + \dots + \alpha_{k,k-1} X_{k-1}$$

La prueba se realiza por inducción

De manera que se puede pasar de cualquier sistema ortogonal de vectores al sistema ortonormal correspondiente, dividiendo cada vector por su longitud (norma)

El proceso descrito permite una gran elasticidad en la elección de una base ortonormal, de una se puede pasar de cualquier base a una ortonormal por ortogonalización y normalización.

El producto escalar de dos vectores es muy fácil de expresar en terminos de las coordenadas de esos vectores con respecto a cualquier base ortonormal, si para U_1, \dots, U_n

una base ortonormal, y

$$X = \varepsilon_1 U_1 + \dots + \varepsilon_n U_n, \quad Y = n_1 U_1 + \dots + n_n U_n$$

entonces

$$\begin{aligned} (X, Y) &= (\varepsilon_1 U_1 + \dots + \varepsilon_n U_n, n_1 U_1 + \dots + n_n U_n) \\ &= \sum_{i=1}^n \sum_{j=1}^n (\varepsilon_i U_i, n_j U_j) = \sum_{i=1}^n \sum_{j=1}^n \varepsilon_i n_j (U_i, U_j) = \sum_{i=1}^n \varepsilon_i n_i \end{aligned}$$

Así la expresión del producto escalar en términos de vectores (de sus coordenadas) con respecto a cualquier base ortonormal coincide con su expresión en términos de las componentes de los vectores, es decir, en términos de las coordenadas con respecto a la base inicial.

1.16.- Transformación de Coordenadas.

Veamos ahora el cambio en las coordenadas de un vector que acompaña un cambio de base.

Sean e_1, e_2, \dots, e_n , y e'_1, e'_2, \dots, e'_n dos bases, y sean

$$e'_1 = a_{11} e_1 + a_{21} e_2 + \dots + a_{n1} e_n$$

$$e'_2 = a_{12} e_1 + a_{22} e_2 + \dots + a_{n2} e_n$$

.....

$$e'_n = a_{1n} e_1 + a_{2n} e_2 + \dots + a_{nn} e_n$$

La relación con la transformación de coordenadas en la matriz A , donde las columnas son las coordenadas de los vectores e'_1, e'_2, \dots, e'_n , con respecto a la base e_1, e_2, \dots, e_n , es decir la matriz.

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}$$

Donde la matriz A es no singular, para que tenga inversa, por medio de la cual los vectores e_1, e_2, \dots, e_n , son expresados en términos de los vectores e'_1, e'_2, \dots, e'_n . Ahora designamos por x_1, x_2, \dots, x_n , las coordenadas del vector X con respecto a la base e_1, e_2, \dots, e_n y por x'_1, x'_2, \dots, x'_n , sus coordenadas con respecto a la base e'_1, e'_2, \dots, e'_n . Determinemos la relación de dependencia entre las coordenadas anteriores y las nuevas. Donde por la independencia lineal de los vectores e_1, e_2, \dots, e_n :

$$x_1 = a_{11}x'_1 + a_{12}x'_2 + \dots + a_{1n}x'_n$$

$$x_2 = a_{21}x'_1 + a_{22}x'_2 + \dots + a_{2n}x'_n$$

.....

$$x_n = a_{n1}x'_1 + a_{n2}x'_2 + \dots + a_{nn}x'_n$$

Y esta ecuación puede ser escrita en la forma,

$$x = Ax'$$

donde

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad y \quad x' = \begin{bmatrix} x'_1 \\ x'_2 \\ \vdots \\ x'_n \end{bmatrix}$$

Son las coordenadas de las columnas del vector X con respecto a las bases e_1, e_2, \dots, e_n y e'_1, e'_2, \dots, e'_n respectivamente.

1.17.- Subespacios.

Un subespacio de un espacio vectorial es un subconjunto que satisface dos condiciones.

i).- Si sumamos dos vectores cualesquiera x e y del subespacio su suma $x+y$ permanece en el subespacio.

ii).- Si multiplicamos cualquier vector x del subespacio por cualquier escalar c , el múltiplo cx permanece en el subespacio.

En otras palabras, un subespacio es un conjunto "cerrado" bajo la suma y la multiplicación por un escalar; estas operaciones nos mantienen dentro del espacio, sin sacarnos. No es necesario verificar las ocho propiedades del Espacio Vectorial ya que estas se satisfacen en el espacio mayor y por ende se satisfacen en cualquier subespacio. Nótese que en particular el vector cero pertenece a cada subespacio debido a la propiedad ii) : basta elegir el escalar $c = 0$

1.18.- Rango de una Matriz.

Dado que una base de un espacio vectorial es un conjunto de vectores con dos propiedades simultáneas :

1).- Es linealmente independiente

2).- Genera espacio

Entonces dos bases cualesquiera de un espacio vectorial V contienen el mismo número de vectores. Este número, compartido por todas las bases, y que expresa el número de " grados de libertad " del espacio, se llama dimensión de V .

Ahora bien, supongamos que reducimos la matriz A de $m \times n$, mediante operaciones elementales e intercambio de filas a una matriz triangular escalonada. Habrá entonces r elementos distintos de cero; las últimas $m-r$ filas de la matriz escalonada son cero. Entonces habrá r variables básicas y $n-r$ variables libres, correspondientes a las columnas de la matriz escalonada. El espacio nulo formado por las soluciones de $Ax = 0$ tiene las $n-r$ variables libres como parámetros independientes.

Si $n = r$, no hay variables libres y el espacio nulo sólo contiene $ax = 0$.

Existen soluciones para cada lado derecho b si y solo si $r = m$ con este número de pivotes la matriz escalonada no tiene filas iguales a cero y puede resolverse mediante sustitución regresiva. En el caso de que $r > m$, la matriz escalonada tendrá $m-r$ filas iguales a cero y habrá $m-r$ restricciones para b de modo que $Ax = b$ se puede resolver. Si existe una solución particular, entonces todas las demás soluciones difieren de la particular en un vector del espacio nulo de A .

En otras palabras el rango de una matriz es un número r tal que entre los menores de la matriz existe un menor distinto de cero de orden r , pero todos los menores de orden $r + 1$ y más grande son iguales a cero.

Y es válido el siguiente teorema importante.

Teorema :

El número máximo de renglones linealmente independientes de una matriz, como también el número máximo de columnas linealmente independientes, coinciden con el rango de una matriz.

1.19.- Transformaciones Lineales.

Asociaremos a cada vector X de un espacio un cierto vector Y del mismo espacio. Tal asociación la llamaremos una transformación del espacio. Designaremos el resultado de la aplicación de transformación A al vector X por AX . Y llamaremos la transformación A lineal si

1).- $A(\gamma X) = \gamma AX$ para cualquier número complejo γ ;

2).- $A(X_1 + X_2) = AX_1 + AX_2$.

Ahora definiremos las operaciones sobre las transformaciones lineales.

El producto de las transformaciones lineales A y B , $AB = C$, será una transformación constituida por las transformaciones B y A por turnos, primero se completará la transformación B y entonces se llevará a cabo A .

El producto de transformaciones lineales es una transformación lineal.

La suma de transformaciones lineales A y B será una transformación C la cual asocia el vector X con el vector $AX + BX$, y la suma de estas es asimismo, una transformación lineal.

Representación de una Transformación Lineal por una Matriz.
Sea una transformación lineal.

$$L : K \xrightarrow{n} K$$

entonces existe un vector A en K^n , tal que $L = L_A$, es decir tal que para todo X tenemos.

$$L(X) = AX$$

Sean E_1, \dots, E_n , los vectores unitarios en K^n .

Si $X = x_1 E_1 + \dots + x_n E_n$ es cualquier vector, entonces

$$\begin{aligned} L(X) &= L(x_1 E_1 + \dots + x_n E_n) \\ &= x_1 L(E_1) + \dots + x_n L(E_n) \end{aligned}$$

Si hacemos $a_i = L(E_i)$

vemos que $L(X) = x_1 a_1 + \dots + x_n a_n = XA$

De modo que las componentes de A son, precisamente los valores: $L(E_1), \dots, L(E_n)$, donde E_i , donde $i = 1, \dots, n$ son los vectores unitarios de K .
En general, tenemos en términos de los vectores columna.

$$L(E_1) = \begin{bmatrix} a_{11} \\ \vdots \\ a_{m1} \end{bmatrix}, \dots, L(E_n) = \begin{bmatrix} a_{1n} \\ \vdots \\ a_{mn} \end{bmatrix}$$

Por tanto.

$$\begin{aligned} L(X) &= x_1 (a_{11} e_1 + \dots + a_{m1} e_m) + \dots + x_n (a_{1n} e_1 + \dots + a_{mn} e_m) \\ &= (a_{11} x_1 + \dots + a_{1n} x_n) e_1 + \dots + (a_{m1} x_1 + \dots + a_{mn} x_n) e_m \end{aligned}$$

En consecuencia si hacemos $A = (a_{ij})$, entonces se ve que

$$L(X) = AX$$

Desarrollando completamente, ésta expresión tiene la siguiente forma.

$$\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + \dots + a_{1n}x_n \\ \vdots \\ a_{m1}x_1 + \dots + a_{mn}x_n \end{bmatrix}$$

Así L es la aplicación lineal asociada con la matriz A .

También designamos a A como la matriz asociada a la transformación lineal L , y se sabe que esta matriz está determinada de manera única.

1.20.- El Rango de Una Transformación Lineal.

Sea A una cierta transformación lineal, el conjunto de vectores AX constituirá obviamente un subespacio, el cual denotaremos por AR_n .

La dimensión de este espacio se dice que es el rango de la transformación A .

Mostraremos ahora que el rango de una transformación es igual al rango de la matriz correspondiente a esta transformación en cualquier base, e_1, e_2, \dots, e_n .

Obviamente el subespacio AR_n está constituido por los vectores

Ae_1, Ae_2, \dots, Ae_n , es decir, la dimensión de AR_n es igual al rango de la matriz cuyas columnas están compuestas de las coordenadas de los vectores Ae_1, Ae_2, \dots, Ae_n , esto es, el rango de la matriz correspondiente a la transformación.

Dado que la dimensión de un subespacio no depende de la elección de la base, se sigue que los rangos de matrices semejantes son iguales.

1.21.- Los Vectores Propios de una Transformación Lineal.

Por un vector propio (vector característico, ó eigenvector) de una transformación lineal A por medio de un vector X diferente de cero tal que.

$$AX = \gamma X$$

Donde γ es un número.

El número γ es llamado el número propio (valor propio, ó eigenvalor) de la transformación.

Los valores propios y vectores propios de la transformación pueden ser determinados de la siguiente manera. Sea A la transformación relacionada con la matriz $A = (a_{ik})$ con

respecto a alguna base; sea X el vector propio de coordenadas, con respecto a esta base x_1, \dots, x_n .

Este sistema de ecuaciones homogéneo en x_1, \dots, x_n , tendrá una solución diferente de cero en el caso solo.

$$\begin{bmatrix} a_{11} - \gamma & a_{12} & a_{1n} \\ a_{21} & a_{22} - \gamma & a_{2n} \\ \dots & \dots & \dots \\ a_{n1} & a_{n2} & a_{nn} - \gamma \end{bmatrix} = 0$$

es decir, si γ es un cero del polinomio característico de la matriz.

Así lo siguiente es válido.

TEOREMA.

Los valores propios de una transformación coinciden con los ceros del polinomio característico de una matriz que está relacionado con ésta transformación con respecto a una base arbitraria.

Ejemplo.

Consideremos la transformación lineal con la matriz.

$$A = \begin{bmatrix} 3 & 1 \\ 0 & 3 \end{bmatrix}$$

Entonces $|A - \gamma I| = (\gamma - 3)^2$, y así $\gamma = 3$ es una raíz doble del polinomio característico.

El sistema de ecuaciones para determinar las coordenadas del vector propio de la transformación A, será :

$$\begin{matrix} 3x_1 + x_2 = 3x_1 \\ 1 & 2 & 1 \end{matrix}$$

$$\begin{matrix} 3x_2 = 3x_2 \\ 2 & 2 \end{matrix}$$

Donde $x_2 = 0$, y así todos los vectores propios de la transformación en cuestión serán $(x_1, 0) = x_1(1, 0)$.

Propiedades de los valores propios y los vectores propios de una matriz.

- 1).- Una matriz y su transpuesta tienen polinomios característicos idénticos.
- 2).- Para una matriz real simétrica, las propiedades de ortogonalidad son simples gracias al hecho de que todos sus valores propios son reales.
- 3).- Los vectores propios que pertenecen a valores propios distintos son ortogonales, y viceversa.

1.22.- Valores Propios de una Forma Cuadrática Positiva Definida.

Un polinomio homogéneo de segundo grado en varias variables x_1, \dots, x_n es llamado una forma cuadrática. Consideremos solamente aquellos con coeficientes reales. Cualquier forma cuadrática puede ser escrita como.

$$\theta(x_1, x_2, \dots, x_n) = \sum_{i,k=1}^n a_{ik} x_i x_k,$$

Donde $a_{ik} = a_{ki}$

Una forma cuadrática se dice que es positiva definida si sus valores son positivos para cualquier valor x_1, \dots, x_n , no todos son cero simultáneamente.

Es evidente que los coeficientes diagonales de una forma positiva definida son positivos, para

$$a_{11} = \theta(1, 0, \dots, 0), \quad a_{22} = \theta(0, 1, \dots, 0), \dots$$

$$a_{nn} = \theta(0, 0, \dots, 1)$$

Denotado por X el vector con componentes (x_1, \dots, x_n) ,

podemos escribir una forma cuadrática como (AX, X) , A siendo la matriz de los coeficientes de la forma. Esta matriz es simétrica por la definición. Los valores propios de la matriz son llamados los valores propios de la forma cuadrática. En vista de los resultados anteriores, todos los valores propios de una forma cuadrática son reales. Y, si una forma cuadrática es positiva definida, sus valores propios son positivos.

1.23.- Los Valores Propios y Vectores Propios de Matrices Semejantes.

Hemos establecido previamente, que las matrices semejantes tienen polinomios característicos idénticos, y

consecuentemente idéntico espectro de valores propios. Por consiguiente los vectores propios de matrices semejantes son las columnas de las coordenadas de los vectores propios de la transformación bajo consideración, con respecto a bases diferentes, y así están relacionados por la relación

$x' = C^{-1} x$, donde C es la matriz de transformación de coordenadas. Esta circunstancia puede ser verificada formalmente: Si $Ax = \gamma x$,

$$(C^{-1} A C)^{-1} (C^{-1} x) = \gamma (C^{-1} x)$$

1.24.- Los Valores Propios de un Polinomio en una Matriz.

Sea A una matriz con valores propios $\gamma_1, \dots, \gamma_n$ y sea

$$\mu(x) = a_0 + a_1 x + \dots + a_m x^m \text{ el polinomio dado.}$$

Entonces, los valores propios de la matriz $\mu(A)$ serán: $\mu(\gamma_1), \mu(\gamma_2), \dots, \mu(\gamma_n)$.

Esto está establecido para una matriz cuyos valores propios son distintos (todos). Realmente, tal matriz puede ser reducida a la forma diagonal por una transformación semejante.

$$A = C^{-1} \begin{bmatrix} \gamma_1 & & \\ & \gamma_2 & \\ & & \dots \\ & & & \gamma_n \end{bmatrix} C$$

$$\text{De acuerdo a : } \mu(A) = C^{-1} \begin{bmatrix} \mu(\gamma_1) & & \\ & \mu(\gamma_2) & \\ & & \dots \\ & & & \mu(\gamma_n) \end{bmatrix} C$$

$$\text{Pero } \mu \left(\begin{bmatrix} \gamma_1 & & \\ & \gamma_2 & \\ & & \dots \\ & & & \gamma_n \end{bmatrix} \right) = \begin{bmatrix} \mu(\gamma_1) & & \\ & \mu(\gamma_2) & \\ & & \dots \\ & & & \mu(\gamma_n) \end{bmatrix}$$

lo cual se sigue del hecho de que,

$$\begin{bmatrix} \gamma_1 & & & \\ & \dots & & \\ & & \gamma_{n-1} & \\ & & & \gamma_n \end{bmatrix}^k = \begin{bmatrix} \gamma_1^k & & & \\ & \dots & & \\ & & \gamma_{n-1}^k & \\ & & & \gamma_n^k \end{bmatrix}$$

Consecuentemente

$$\begin{aligned} \mu \left(\begin{bmatrix} \gamma_1 & & & \\ & \dots & & \\ & & \gamma_{n-1} & \\ & & & \gamma_n \end{bmatrix} \right) &= \sum_{k=0}^m a_k \begin{bmatrix} \gamma_1^k & & & \\ & \dots & & \\ & & \gamma_{n-1}^k & \\ & & & \gamma_n^k \end{bmatrix} \\ &= \begin{bmatrix} \sum_{k=0}^m a_k \gamma_1^k & & & \\ & \dots & & \\ & & \sum_{k=0}^m a_k \gamma_{n-1}^k & \\ & & & \sum_{k=0}^m a_k \gamma_n^k \end{bmatrix} \\ &= \begin{bmatrix} \mu(\gamma_1) & & & \\ & \mu(\gamma_2) & & \\ & & \dots & \\ & & & \mu(\gamma_n) \end{bmatrix}. \end{aligned}$$

Así la matriz $\mu(A)$ es similar a la matriz $\begin{bmatrix} \mu(\gamma_1) & & & \\ & \dots & & \\ & & \mu(\gamma_{n-1}) & \\ & & & \mu(\gamma_n) \end{bmatrix}$,

y de acuerdo a sus valores propios son $\mu(\gamma_1), \mu(\gamma_2), \dots, \mu(\gamma_n)$

q.e.d

Este resultado es verdadero para cualquier matriz.

Notese en particular que los valores de la matriz A son $\gamma_1, \gamma_2, \dots, \gamma_n$

1.25.- La Normalización de los Vectores Propios de una Matriz. Segundo Grupo de Relaciones de Ortogonalidad.

Sea A una matriz real cuyos valores propios $\gamma_1, \dots, \gamma_n$ son distintos, y sean X_1, \dots, X_n los vectores propios correspondientes a ellos.

Como sabemos, la matriz transpuesta A' , tiene los mismos vectores propios. Sean X'_1, \dots, X'_n , los vectores propios de la matriz A' , y su enumeración elegida tal que X_i y X'_i corresponden al conjugado complejo de valores propios, y por consiguiente, se sigue la relación de ortogonalidad :

$$(X'_i, X'_j) = 0 \text{ para } i \neq j$$

Mostraremos ahora que habiendo elegido X'_1, \dots, X'_n de cualquier manera (están determinados pero, para un multiplicador numérico), podemos normalizar los vectores X'_1, \dots, X'_n , tal que $(X'_i, X'_i) = 1$

Demostración.

Sean X_1, \dots, X_n linealmente independientes.

resolviendo X'_i , en terminos de está base :

$$X'_i = \mu_{i1} X_1 + \mu_{i2} X_2 + \dots + \mu_{in} X_n$$

Formando el producto escalar (X'_i, X'_i) , obtenemos

$$\begin{aligned} (X'_i, X'_i) &= \mu_{i1} (X'_i, X_1) + \dots + \mu_{in} (X'_i, X_n) \\ &= \mu_{ii} (X'_i, X_i) \end{aligned}$$

De aquí concluimos que :

$$(X'_i, X_i) = \alpha_i \neq 0, \text{ para } (X'_i, X'_i) > 0$$

Tomando en lugar de los vectores X_1, \dots, X_n , los vectores X'_1, \dots, X'_n

$1/\alpha X_{11}, \dots, 1/\alpha X_{nn}$, llegamos a la normalización requerida

$$(X'_i, 1/\alpha X_{ii}) = 1/\alpha (X'_i, X_{ii}) = 1$$

1.26.- La forma Canónica de Jordan.

Hemos visto que si una matriz tiene n valores propios, puede trasladarse a su forma diagonal por transformaciones semejantes.

Dada la presencia de raíces múltiples, sin embargo, no siempre la transformación puede ser posible. Así el problema de reducir una matriz a una forma tan simple como sea posible por medio de una transformación semejante hecha.

El problema es equivalente al encontrar una base con respecto a la cual la transformación lineal relacionada con la matriz dada tuviera una matriz de la forma más simple y esto último es la forma canónica de Jordan.

Omitiremos la demostración y nos concretaremos a la descripción de ésta forma canónica.

La forma (clásica) canónica es una matriz cuasi-diagonal compuesta de bloques canónicos.

es el orden de la i -ésima caja ó bloque.

La determinación de los bloques de Jordan para una matriz dada A presenta ciertas dificultades. La característica

polinomial $\pi_i(\gamma - \gamma)$ coincide con la característica

polinomial de la matriz original, y consecuentemente es posible encontrarla sin conocer la matriz canónica en sí misma. Pero el conocimiento de la característica polinomial no hace posible la determinación completa de la forma canónica para un valor propio γ de multiplicidad

k estos pueden corresponder a varios bloques de Jordan

conteniendo este valor como un elemento diagonal, y considerandolos solamente la suma de sus órdenes, no el orden de cada bloque en particular. Si la forma canónica debe ser determinada, el conocimiento de los "divisores elementales" de la matriz debe tenerse.

Designamos por $D_i(\gamma)$ el común divisor más grande de todos los menores del i -ésimo orden del determinante $|A - \gamma I|$. En particular, $D_n(\gamma)$ coincide con la característica polinomial.

Puede ser probado que todo $D_i(\gamma)$, como $D_n(\gamma)$, son comunes a la clase de matrices similares, y puede demostrarse que $D_{i-1}(\gamma)$ divide a $D_i(\gamma)$.

Pero

$$\frac{D_i(\gamma)}{D_{i-1}(\gamma)} = E_i(\gamma);$$

Obviamente

$$D_n(\gamma) = \prod_{i=1}^n E_i(\gamma)$$

Se sigue que

$$E(\gamma) = \frac{D_n(\gamma)}{D_{n-1}(\gamma)}$$

es el polinomio mínimo de la matriz

Resolviendo $E(\gamma)$ en factores lineales, entonces

$$E(\gamma) = \prod_{j=1}^s (\gamma - \gamma_j)^{m_{ij}}$$

Aquí s denota el número de valores propios distintos,

$$\sum_{i=1}^n m_{ij} = k_j$$

Es obvio que entre los exponentes m_{ij} solamente algunos serán diferentes de cero.

Consideremos la matriz que transforma una matriz A dada en la forma canónica. Entonces las columnas de la matriz transformada serán las componentes de los vectores de esas bases en términos de lo cual la transformación en cuestión se describe como una matriz canónica.

Esta matriz canónica tiene la forma .

$$\begin{bmatrix} \wedge_1 & & & & \\ & \wedge_2 & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & \wedge_s \end{bmatrix}$$

Donde

$$\wedge_r = \begin{bmatrix} \gamma_r & 0 & \dots & 0 \\ 0 & \gamma_r & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \gamma_r \end{bmatrix}$$

Y el orden de \wedge_r es igual a m ; sean

$U_1^{(1)}, \dots, U_{m_1}^{(1)}, \dots, U_1^{(s)}, \dots, U_{m_s}^{(s)}$ la base correspondiente.

Entonces las siguientes formulas para la transformación son válidas.

$$AU_1^{(r)} = \gamma U_{r-1}^{(r)} + U_2^{(r)}$$

$$AU_2^{(r)} = \gamma U_{r-2}^{(r)} + U_3^{(r)}$$

.....

$$AU_{mr-1}^{(r)} = \gamma U_{r-mr-1}^{(r)} + U_{mr}^{(r)}$$

$$AU_{mr}^{(r)} = \gamma U_{mr}^{(r)}$$

Para todo $r = 1, \dots, s$.

Veamos que entre los vectores de la base canónica son los vectores propios de la matriz dada, uno por block. Puede ser probado que con todos estos vectores propios linealmente independientes de la matriz dada, son iguales al número de blocks canónicos en su forma canónica. En particular, el número de vectores propios linealmente independientes corresponden a un valor propio que es igual al número de blocks canónicos que contienen este número. No es de mayor multiplicidad que el valor propio, y es igual a esta multiplicidad en caso, y solo en caso de que, todos los bloques contengan el valor propio dado y sean de orden 1, es decir, cuando los divisores elementales correspondientes sean lineales.

C A P I T U L O

II.

PROBLEMAS DE CALCULO
DEL
ALGEBRA LINEAL

PROBLEMAS DE CALCULO DE ALGEBRA LINEAL

Los problemas de cálculo más comunes del Algebra Lineal son los concernientes a las matrices de números reales. Este capítulo nos llevará unos pasos más en la resolución de estos problemas y se tratan los principales problemas y los más antiguos del tema $Ax = b$ y $Ax = \gamma x$, sin embargo podemos decir que cada uno de ellos ha sido considerado por la nueva generación de matemáticos, y los dividimos en cinco grupos.

1).- Sea A una matriz de números reales con n renglones y n columnas, es decir una matriz cuadrada de $n \times n$, y b un vector columna de n renglones. El problema que frecuentemente se presenta en las ecuaciones lineales es encontrar un vector columna " x " de n renglones, tal que.

$$Ax = b$$

Por regla general se supone que A es una matriz no singular puesto que, solo si esto cumple entonces existe una solución única para toda b . En caso contrario, es decir si A es una matriz singular, entonces se tendrá una infinidad de soluciones para b .

2).- Otro problema tradicional, si tenemos una matriz de números reales, con n renglones y n columnas es encontrar

-1

la matriz inversa, esto es A^{-1} , y la condición para que está exista es que su determinante sea distinto de cero.

3).- El tercer problema que se encuentra con frecuencia, es que si A es una matriz de n renglones y n columnas de números reales, la cual es simétrica, esto es que cumpla que

sea cuadrada e igual a su transpuesta, es decir, si $A^T = A$, es simétrica.

Ejemplo:

$$\text{Sea } A = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 0 \end{bmatrix} \quad \text{entonces, } A^T = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$

y , $A = A^T$, por lo tanto A es una matriz simétrica.

Dado que $A = A^T = (a_{ij}) = (a_{ji})$.

Aquí el problema consiste en encontrar alguno de todos (necesariamente real) los eigen valores de A. Recordando que un eigen valor de A es un número γ para el cual existe un vector columna v, tal que ;

$$Av = \gamma v$$

Así un vector (columna) v es llamado eigen vector de A correspondiente a γ , y con frecuencia el problema computacional encierra el encontrar un v correspondiente a cada eigen valor computado. Existen n eigen vectores ortonormales de A donde uno corresponde a cada eigen valor de A.

4).- Sea A una matriz no simétrica de n renglones y n columnas de números reales. Otro problema tradicional del Algebra Lineal, es de encontrar alguno de todos sus eigen valores y algunas veces tambien su eigen vector columna correspondiente. Recordando que un eigen vector renglón corresponde a γ y un vector renglón de n columnas v, es tal que ;

$$Av = \gamma v$$

Cuando A es no-simétrica el problema se complica de muchas maneras:

PRIMERO.- Algunos de los eigen valores de A son con frecuencia números complejos.

SEGUNDO.- Puede ser que no sea un eigen vector de n columnas linealmente independientes, y estos, no son por lo general ortogonales. Realmente es probable que se acerquen a la dependencia lineal, y lo mismo se tiene para los eigen vectores renglón.

TERCERO.- Si un eigen valor γ es una raíz de multiplicidad $K > 1$ de la ecuación característica $\det (A - \gamma I) = 0$,

entonces puede existir en cualquier parte de 1 a k, de la columna linealmente independiente de los ingen -vectores correspondientes a γ_i . (Si A fuera simétrica, habría siempre un k) si el número es menor que k, corresponde a uno ó más blocks no diagonales en la forma canónica de Jordan de A, es decir.

$$\begin{bmatrix} \gamma_i & 0 & 0 & \dots & 0 & 0 \\ 1 & \gamma_i & 0 & \dots & 0 & 0 \\ 0 & 1 & \gamma_i & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 1 & \gamma_i \end{bmatrix}$$

Donde en la diagonal principal se encuentra γ_i ; y directamente bajo la diagonal (en la subdiagonal) están dispuestos elementos los cuales son todos la unidad; todo el resto de los elementos son cero.

O equivalentemente los tan llamados divisores elementales no lineales de A.

CUARTO.- Múltiples o cercanamente múltiples eigen-valores de A son probables de ser funciones de cambio muy rapido de los elementos $a_{i,j}$ de A, de modo que los cálculos son de lo más complicado.

5).- Para cualquier vector columna y definimos la norma potencia p - ésima de y, como;

$$(1) \quad || y ||_p = \left(\sum_{i=1}^n | y_i |^p \right)^{1/p}$$

Donde p es un número real con $1 \leq p < \infty$, $y = y_1, \dots, y_n$

son las componentes de y , en un sistema coordenado dado. Definiremos la norma máxima como en el caso del límite de p cuando $p \rightarrow \infty$ de (1).

$$\|y\|_{\infty} = \max_{1 \leq i \leq n} |y_i|$$

Las normas más usadas en análisis numérico son $p = 1, 2, \dots, \infty$, pero las estadísticas le están prestando atención a valores de p , entre 1 y 2.

Sea una matriz de n renglones y k columnas de números reales y sea b el vector columna de n renglones. Dado algún p , un problema computacional más reciente es el de encontrar un vector columna x de k renglones tal que,

$$\|Ax - b\|_p$$

Sea minimizado cuando $p = 2$, el caso común, este es el caso del problema de mínimos cuadrados lineales, dado que tenemos una matriz de n renglones y k columnas, entonces podemos imaginar que el número de observaciones es mayor que el número de incógnitas, así que se debe expresar que el sistema $Ax = b$ sea inconsistente. Probablemente no exista una selección de x que coincida perfectamente con los datos b ; o, en otras palabras probablemente el vector b no sea una combinación de las columnas de A .

Y el problema es seleccionar x que minimice el error, y de nuevo esta minimización se hará en el sentido de los mínimos cuadrados. El error es $E = \|Ax - b\|$, que es exactamente la distancia de b al punto Ax en el espacio columna de A , por lo tanto lo que hacemos al usar mínimos

cuadrados es localizar el punto $p = \bar{Ax}$ que está más cerca de b que cualquier otro punto en el espacio columna.

Así la solución en términos de mínimos cuadrados para un sistema inconsistente $Ax = b$ de n ecuaciones con k incógnitas satisface.

$$A^T Ax = A^T b$$

Ahora bien para $p = 2$ la esfera unitaria es muy pareja y los métodos de trabajo de análisis buenos.

Sin embargo para $p = 1$ ó ∞ , la esfera unitaria tiene muchos ángulos y métodos para minimizar $\|Ax - b\|_p$, haciéndose combinatorio ó discreto.

6).- Para dos vectores columna de n renglones x y y y definimos $x \geq y$, para hacer notar que $x_i \geq y_i$, para todas las componentes de x y y

Sea A y b , una matriz de n renglones y q columnas de números reales y un vector columna de n renglones respectivamente, entonces un problema de cálculo importante es el describir el conjunto " S " de vectores columna x de n renglones tal que ;

$$Ax \geq b .$$

Algunas veces, como en los problemas de programación lineal uno parecido para vectores x en S tal que, $C^t x$ es un mínimo, donde C es un vector columna dado, con k renglones.

Hasta este momento hemos hablado extensamente solo de matrices de números reales, problemas similares se presentan ocasionalmente para matrices de números complejos de la forma.

$$A = \begin{bmatrix} 0 & i \\ -i & 0 \end{bmatrix}$$

Mucho de los problemas pueden ser formulados también para matrices en las cuales sus elementos son expresiones incógnitas ó literales. Como método para simbolizar la manipulación en computadoras digitales, para hacerlo más accesible a los usuarios de las computadoras. Los problemas del Algebra Lineal con matrices literales serán posteriormente más estudiadas. La manipulación de simbolos prácticos probablemente será de más interes para los matemáticos en computación que para cualquier otro.

La presente discusión está limitada a las matrices de números ó mejor dicho a los problemas :

1).- Donde A es una matriz de $n \times n$ de números reales y b un vector columna de $n \times 1$ y se desea encontrar x tal que,

$$Ax = b$$

2).- Donde A es una matriz de $n \times n$ de números reales, y se requiere,

$$A^{-1}$$

3).- Donde A es una matriz simétrica de $n \times n$ de números reales en la cual el problema es encontrar alguno de todos los eigen-valores de A, es decir,

$$Av = \lambda v$$

4).- Donde A es una matriz no simétrica de $n \times n$ de números reales, en la cual el problema es análogo al anterior.

$$Av = \lambda v$$

Porque los problemas lineales anteriores, surgen con frecuencia ?. Por que son importantes ?. La respuesta es que los operadores lineales son los más simples en matemáticas y solo los operadores lineales son plenamente estudiados, por consiguiente son un modelo natural para una matemática aplicada, para ser usada al atacar un problema.

Aun cuando los operadores lineales en espacios de dimensión infinita ocurrieran en análisis ó ecuaciones diferenciales. (por ejemplo), la realidad del significado de la computación en el cual solo espacios de dimensión finita pueden ser manejados con computadoras digitales.

Modelos más realistas de matemáticas aplicadas son generalmente no lineales. Pero cuando algunos operadores no lineales son usados, la solución actual de ecuaciones, funcionales siempre involucra la aproximación de operadores no lineales por algunos lineales. Un ejemplo típico de esto es el uso del método de Newton para resolver un sistema de ecuaciones no lineales en el cual a todo paso un sistema de ecuaciones localmente más apropiado debe ser resuelto. Los problemas no lineales generalmente son muy difíciles. Al atacarlos por métodos lineales, es esencial que nuestros conceptos lineales sean muy claros, de modo que ellos

puedan formalizarse para trabajar sin falla. Sólo de esta manera puede el análisis concentrarse en las dificultades reales del aspecto no lineal.

Este punto de vista no solo enfatiza la importancia de poder resolver problemas lineales, sino también la necesidad de resolver sistemas lineales con métodos extremadamente formales.

Del sistema de ecuaciones lineales 1) es decir, $Ax = b$ surgen dos opciones.

- Una de aproximación a ecuaciones lineales funcionales generalmente ecuaciones diferenciales parciales u ordinarias.

- La otra opción es un problema en el que dado un conjunto de datos, implica que el método de interpolación ó aproximación por familias de funciones lineales, tendran que ser aplicados.

El problema de los eigen-valores generalmente surge de estudios de vibración ó estabilidad ó resonancia de sistemas físicos lineales (aeronaves, reactore, etc) ó problemas de análisis de factor, valiendose del texto de Noble tomamos los siguientes ejemplos físicos de problemas del cálculo de matrices.

EJEMPLO No 1

Primero, retomando la ecuación diferencial para el movimiento armónico simple.

Si una partícula de masa m se mueve en línea recta y es atraída a un punto en la línea por una fuerza la cual es proporcional a la distancia X de la partícula desde el punto, entonces la ecuación de movimiento de la partícula está dada.

$$m \frac{d^2 x}{dt^2} + kx = 0 \quad (a)$$

donde P es una constante positiva de proporcionalidad, para resolver esta ecuación agrupamos, en la forma acostumbrada,

$$x'' = -\frac{P}{m} x e^{i\omega t}$$

Donde ω es la frecuencia angular de vibración y x_0 es una constante que representa la amplitud de vibración, la ecuación (a) da.

$$(-m\omega^2 + P) x_0 e^{i\omega t} = 0$$

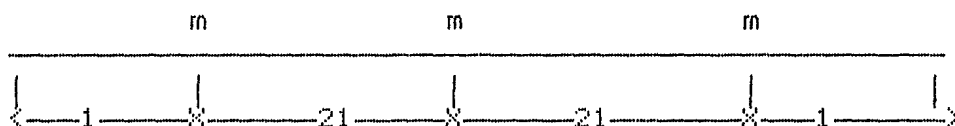
Dado que la partícula está en movimiento, $x \neq 0$. También $e^{i\omega t}$

es distinto de cero para toda t . De aquí la ecuación anterior implica que.

$$m\omega^2 = P \quad \text{ó} \quad \omega = \frac{\sqrt{P}}{m}$$

Esto significa físicamente que hay solo una frecuencia de vibración libre.

Ahora veamos un ejemplo más complicado. Consideremos tres partículas, cada una de masa m , tomadas en las posiciones 1, 31, 51, respectivamente, a lo largo de un alambre elástico de longitud 61, como lo muestra la figura a continuación.



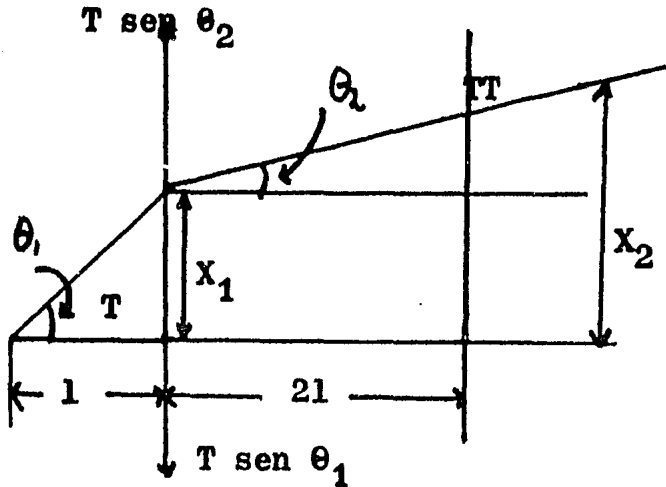
Tres masas en un alambre.

El alambre es fijo por ambos lados y está bajo una tensión T .

Supongamos que las partículas ejecutan pequeñas vibraciones transversales, bajo fuerzas no externas, los desplazamientos de estas tres partículas en una dirección perpendicular a la línea de equilibrio del alambre siendo.

x_1, x_2, x_3 , respectivamente.

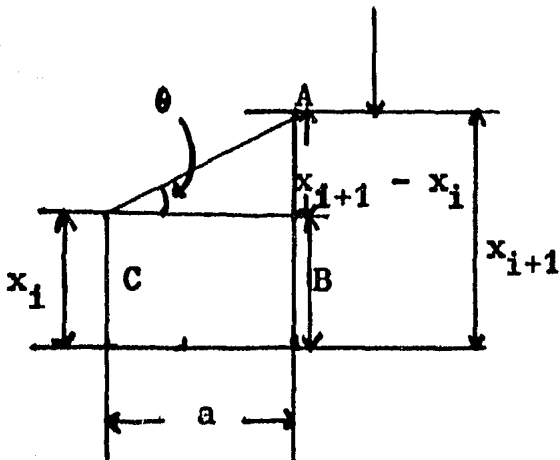
Las fuerzas en la cama 1 se muestra en la figura siguiente.



La fuerza resultante perpendicular a la línea de equilibrio es.

$$T (\sin \theta_2 - \sin \theta_1)$$

Tomando en cuenta que los desplazamientos son pequeños, tal que la tensión T puede llegar a ser constante, y el $\sin \theta_2 - \sin \theta_1$ pueden ser aproximados por $\tan \theta_2 - \tan \theta_1$ como se muestra en la figura a continuación.



La aproximación para $\sin \theta$

$$\sin \theta = \frac{AB}{AC} \approx \frac{AB}{BC} = \frac{x_{i+1} - x_i}{a}$$

Vibración de camas perpendicular a un alambre.

Por la segunda ley de movimiento de Newton (los tiempos de aceleración de la masa es igual a la fuerza en la partícula) da para la cama 1, al usar los resultados anteriores.

$$m \frac{d^2 X_1}{dt^2} = - \frac{T X_1}{1} + \frac{T(X_2 - X_1)}{21}$$

Similarmente, para el movimiento en las camas 2,3, obtenemos

$$m \frac{d^2 X_2}{dt^2} = - \frac{T(X_2 - X_1)}{21} + \frac{T(X_3 - X_2)}{21}$$

$$m \frac{d^2 X_3}{dt^2} = - \frac{T(X_3 - X_2)}{21} - \frac{TX_3}{1}$$

Ahora asumimos que todas las cantidades varían sinusoidalmente con el tiempo y con la misma frecuencia, y agrupando.

$$X_r = x_r e^{i\omega t} \quad \text{donde } r = 1,2,3,$$

Entonces las ecuaciones anteriores se convierten en :

$$\begin{aligned} (3 - \gamma) x_1 - x_2 &= 0 \\ -x_1 + (2 - \gamma) x_2 - x_3 &= 0 \\ -x_2 + (3 - \gamma) x_3 &= 0, \end{aligned} \quad (b)$$

Donde $\gamma = 2w \text{ ml/T}$; esto establece que γ debe ser eigen valor y $x = (x_1, x_2, x_3)^T$ su eigen vector asociado para la matriz

$$\begin{bmatrix} 3 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 3 \end{bmatrix}$$

Como vimos antes, las soluciones distintas de cero x existiran precisamente cuando el coeficiente del determinante en la ecuación (b) sea cero, que es cuando,

$$\det \begin{bmatrix} 3 - \gamma & -1 & 0 \\ -1 & 2 - \gamma & -1 \\ 0 & -1 & 3 - \gamma \end{bmatrix} = 0 .$$

Este determinante es fácil de evaluar, tal que de la expresión inmediatamente anterior se tiene que

$$-\gamma^3 + 8\gamma^2 - 19\gamma + 12 = 0,$$

cuyas raices son

$$\gamma = 1, \quad \gamma = 3, \quad \gamma = 4 .$$

Para $\gamma = 1$ (b) se convierte en

$$2x_1 - x_2 = 0$$

$$-x_1 + x_2 - x_3 = 0$$

$$-x_2 + 2x_3 = 0$$

la cual se resuelve para obtener ,

$$\gamma = 1, x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \alpha \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad \text{para una } \alpha \text{ arbitraria } \neq 0$$

Para $\gamma = 3$ (b) es ,

$$\gamma = 3, x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \beta \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \quad \text{para una } \beta \text{ arbitraria } \neq 0$$

Para $\gamma = 4$ (b) es ,

$$\gamma = 4, x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \Gamma \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \quad \text{para una } \Gamma \text{ arbitraria } \neq 0$$

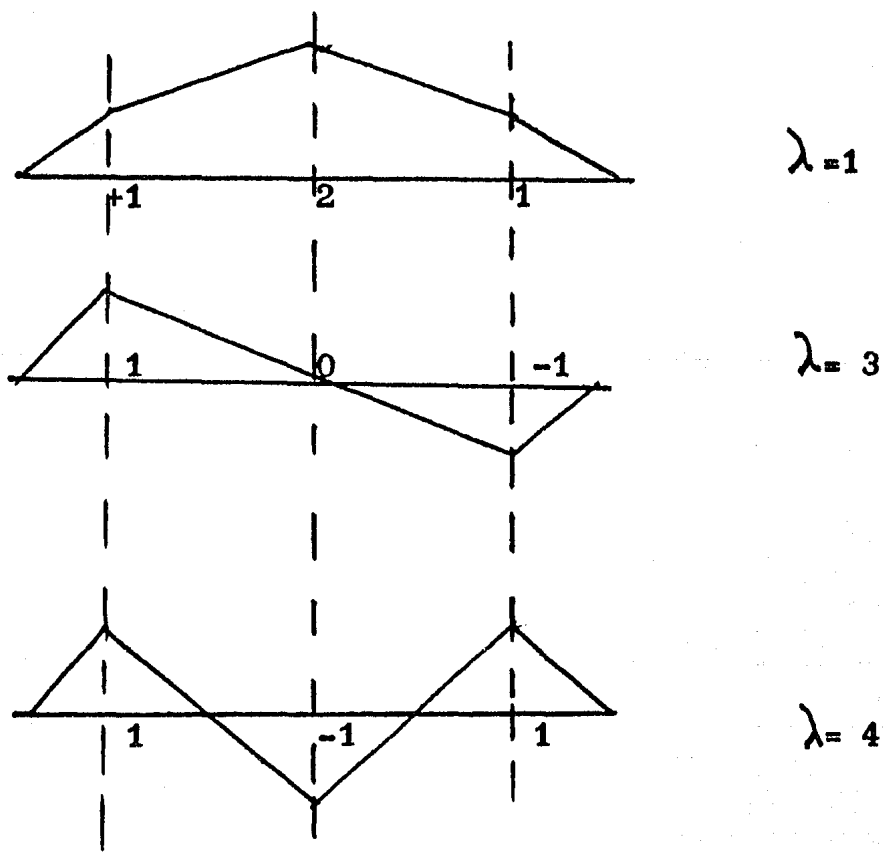
El significado físico de los resultados matemáticos anteriores en eigen sistemas es el siguiente. Correspondiendo a $\lambda = 1$, por ejemplo, hay una vibración

libre con frecuencia angular ω dado por $\omega = \sqrt{\frac{g}{l}}$, y correspondiendo a esta frecuencia de vibración hay un modo de oscilación dado por $X = x e^{i\omega t}$ y (b), tal que las

proporciones $x_1 : x_2 : x_3$, son $1 : 2 : 1$.

Similarmente para $\lambda = 3$ y $\lambda = 4$.

Los modos de vibración están ilustrados gráficamente en la figura siguiente. Estas tres frecuencias y modos de vibración son solo algunos que pueden existir.



Modos de vibración.

C A P I T U L O

III.

UN ENFOQUE MAS SUTIL
DE LOS
PROBLEMAS

UN ENFOQUE MAS SUTIL DE LOS PROBLEMAS

Dado que las computadoras actuales tienen capacidad de almacenaje finito y desde que tienen precisión finita, necesitamos estudiar más de cerca la naturaleza de las matrices A y sus problemas de cálculo.

Es A una matriz densa (tiene más elementos $a_{ij} \neq 0$, ó es una matriz hueca (tiene más elementos $a_{ij} = 0$) ?

$$A = \begin{bmatrix} 1 & 0 & -2 \\ 0 & 1 & 3 \\ 4 & 6 & 4 \end{bmatrix}$$

Matriz densa

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 0 & 0 \\ 6 & 0 & 0 \end{bmatrix}$$

Matriz hueca.

Si A es hueca, los elementos distintos de cero forman un paso significativo ? Por ejemplo es A triangular (con $a_{ij} = 0$ para $i > j$ ó para $i < j$) ?

$$A = \begin{bmatrix} c_{11} & 0 & \dots & 0 \\ c_{21} & c_{22} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ c_{n1} & c_{n2} & \dots & c_{nn} \end{bmatrix}$$

Matriz triangular inferior

$$A = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ 0 & b_{22} & \dots & b_{2n} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & b_{nn} \end{bmatrix}$$

Matriz triangular superior

O es de la forma de Hessenberg. ($a_{ij} = 0$ para $i > j + 1$ ó $j = i + 1$) ?

Es una matriz banda ($a_{ij} = 0$ para $|i - j| > m$, donde $m \ll n$)?

$$A = \begin{bmatrix} a_{11} & 0 & \dots & 0 & 0 \\ a_{21} & a_{22} & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & a_{n,n-1} & a_{nn} \end{bmatrix}$$

Es una matriz tridiagonal (es decir una matriz banda con $m = 1$) ?

$$A = \begin{bmatrix} a_{11} & 0 & \dots & 0 \\ 0 & a_{22} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & a_{nn} \end{bmatrix}$$

Todas estas formas ocurren con frecuencia y pueden ser dadas con consideración especial. Si A es simétrica (es decir $A = A^T$). Positivamente definida. Esto es que cada uno de los siguientes criterios es una condición necesaria y suficiente para que una matriz simétrica real A sea positivamente definida :

I .- $x^T A x = 0$ para todos los vectores x distintos de cero.

II.- Todos los valores propios de A satisfacen $\lambda_i > 0$

III.- Todas las submatrices A_k tienen determinantes positivos

IV.- Todos los pivotes (sin intercambio de filas) satisfacen $d_i > 0$.

Ahora bien si es hueca, es el paso asociado con la matriz adyacente de alguna gráfica. Frecuentemente las matrices asociadas con estructuras ó con ecuaciones diferenciales parciales son mejor entendidas en términos de una gráfica asociada. Por ejemplo :

En los elipsoides de n dimensiones, para entender la geometría de una matriz positivamente definida es $x^T A x = 1$. Los vectores x que satisfacen esta ecuación se localizan en una superficie en el espacio n - dimensional. Y se quiere mostrar que es una elipsoide n - dimensional con centro en el origen.

$$\text{Sea } A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \text{ y } x^T A x = 2u^2 - 2uv + 2v^2 = 1$$

Reescribiendo esta ecuación como una suma de cuadrados con $\gamma_1 = 1$ y $\gamma_2 = 3$ como coeficientes. Se tiene ,

$$2u^2 - 2uv + 2v^2 = 1 \left(\frac{u}{2} + \frac{v}{2} \right)^2 + 3 \left(\frac{u}{2} - \frac{v}{2} \right)^2 = 1.$$

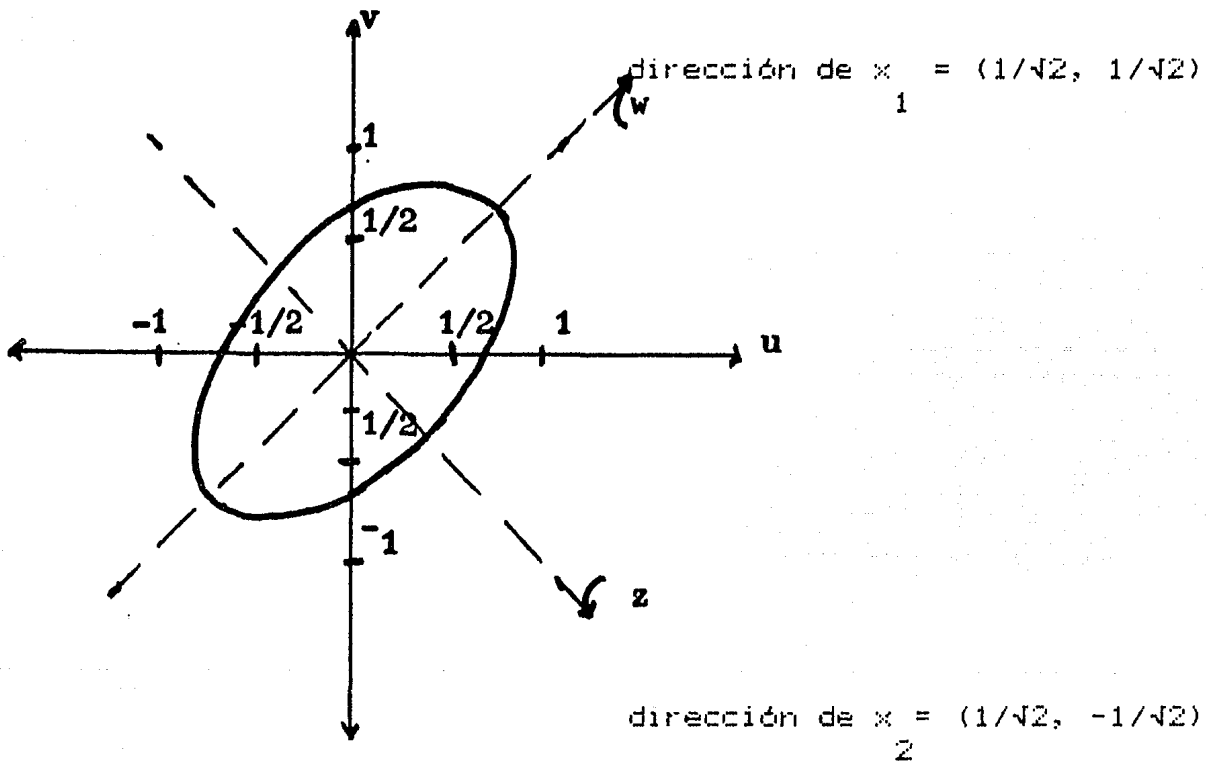
$$\text{Haciendo } \left(\frac{u}{2} + \frac{v}{2} \right) = w \text{ , y , } \left(\frac{u}{2} - \frac{v}{2} \right) = z$$

Resulta la ecuación $w^2 + 3z^2 = 1$, la cual corresponde a una elipse. Cuyo eje más largo termina en el punto $w = 1, z = 0$ y el eje más corto termina en $w = 0, z = 1/\sqrt{3}$. Así la

longitud del eje mayor es $2 \frac{1}{\sqrt{2}}$, y la longitud del eje menor es $2 \frac{1}{\sqrt{6}}$. Además, estos ejes están dirigidos hacia

los vectores propios $x_1 = (1/\sqrt{2}, 1/\sqrt{2})$ y $x_2 = (1/\sqrt{2}, -1/\sqrt{2})$,
 por lo tanto, la geometría está completamente vinculada con
 los valores y vectores propios.

Gráfica de la elipse $2u^2 - 2uv + 2v^2 = 1$



Otro ejemplo ilustrativo es el de un Problema de Programación Lineal. En el problema del mínimo, se tiene que deseamos minimizar Cx sujeto a, $Ax = b$ con $x \geq 0$ es decir.

En forma matricial podemos escribir.

$$\text{Min } Cx$$

s.a.

$$Ax \geq b$$

$$x \geq 0$$

En un ejemplo numérico sea

$$\text{Mín } x + y$$

s.a.

$$x + 2y \geq 6$$

$$2x + y \geq 6$$

$$x \geq 0$$

$$y \geq 0$$

Matricialmente se puede escribir

$$\text{Mín } C = (1, 1) \begin{bmatrix} x \\ y \end{bmatrix}$$

$$A = \begin{vmatrix} 1 & 2 \\ 2 & 1 \end{vmatrix} \begin{vmatrix} x \\ y \end{vmatrix} = \begin{vmatrix} 6 \\ 6 \end{vmatrix} > b$$

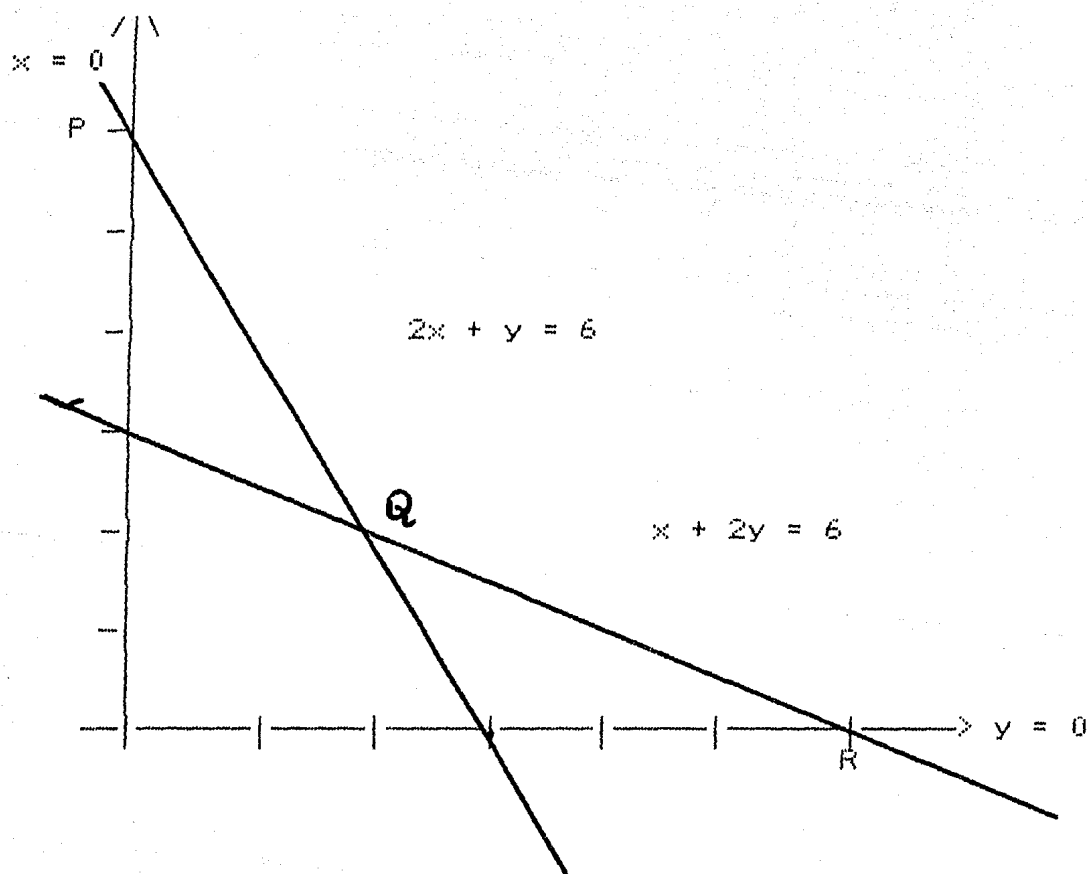
$$x \geq 0 \quad y \geq 0$$

Interpretación geométrica, la primera condición $x \geq 0$ restringe al vector al cuadrante positivo, que en dos dimensiones es un cuarto del plano. Las otras m restricciones producen m semiespacios adicionales y los vectores factibles son todos aquellos que cumplen las $m + n$ condiciones; están en el cuadrante $x \geq 0$ y al mismo tiempo satisfacen $Ax \geq b$. Dicho de otra manera, el conjunto factible es la intersección de $m + n$ semiespacios; puede ser acotado, puede no estar acotado o puede estar vacío. Ahora bien, la función de costo Cx , introduce al problema una familia de planos paralelos, que dan los costos posibles, cuando varían, estos planos barren todo el espacio y el vector óptimo x^* está en el punto donde tocan primero al conjunto factible.

En nuestro ejemplo numérico se tiene $n=2$ variables y $m = 2$ restricciones; entonces hay seis intersecciones posibles, como se muestra en la gráfica, donde tres de ellas son realmente esquinas del conjunto factible y son los puntos

marcados con las letras A, B, C. respectivamente $(0,6)$, $(2,2)$ y $(6,0)$ y uno de ellos debe ser el vector óptimo.

Gráfica.



Ahora, volviendo a lo anterior, están los elementos a_{ij} , almacenados en la memoria de la computadora, ó son generados de algún algoritmo conforme se necesite ? ó podrían definir el contenido informacional de una matriz como el número de espacios necesitados (en una cierta computadora) para almacenar datos y programar para obtener todos los a_{ij} ?

Cuál es el tamaño de la matriz A, relativo al tamaño de la memoria y rapidez de una computadora dada ? . Si nosotros estamos resolviendo un sistema de ecuaciones lineales $Ax=b$ tenemos muchas soluciones para b ó solo una ? Tenemos muchas matrices A diferentes que están juntas o sólo una A? Son los elementos de A precisamente números matemáticos (por ejemplo enteros) ó son números físicos (aproximaciones) ?

Los mejores propósitos de un sistema de ecuaciones lineales no han considerado esta cuestión y miran al análisis numérico para explicar las posibilidades y seleccionar opciones.

Si para un propósito se necesita la matriz inversa A^{-1} es generalmente el descubrir por qué. Frecuentemente se desea un camino lo más práctico para resolver $Ax = b$, para un

vector arbitrario b. Teniendo A^{-1} almacenada, el propósito que se espera es obtener la solución " x " en la forma

$A^{-1} b$, para algún b nuevo que venga. Y será mostrado que hay

otros métodos para obtener $A^{-1} b$, para nuevos vectores b, los métodos que no requieren de más almacenaje y no toman el camino más largo para la misma precisión, que la

multiplicación de A^{-1} por b. Porque de estos factores, la

computación de A^{-1} puede frecuentemente ser omitida. Sin embargo ciertas aplicaciones estadísticas realmente requieren de conocimientos de un número de elementos

diagonales de A.

El problema de eigen valor, es decir al tener una matriz de $n \times n$, A la cual es simétrica, el encontrar alguno de todos los eigen valores de A , $Av = \lambda v$, puede requerir de encontrar todos los eigen valores ó solo algunos, dependiendo de los requerimientos que se tengan, ya sea que los correspondientes eigen vectores sean necesitados ó no?

Si un conjunto completo de eigen-vectores es necesitado, es importante que sean mutuamente ortogonales?. Consiguiendo eigen-vectores ortogonales correspondientes a eigen-valores múltiples, es más largo y más difícil que encontrar eigen-valores?.

En el problema del valor característico, si las respuestas a las preguntas anteriores son afirmativas, entonces, para matrices no-simétricas A , uno tiene similares selecciones? Queremos todos los valores característicos ó solo algunos? Queremos eigen-vectores columna? Queremos eigen-vectores renglon ó ambos?. Pero entonces viene una nueva elección. Si algunos eigen valores son multiplicados y corresponden a divisores elementales no lineales, que vectores es nuestro propósito encontrar? En algunos textos de álgebra se aprende sobre cadenas de vectores principales que con el vector característico forman una base para el espacio nulo

N (Kernel) de $(A - \lambda I)$, donde λ es un eigen-vector de multiplicidad k , estos vectores principales están asociados con la forma canónica de Jordan de A . Un propósito que tenga una buena base en álgebra, querrá ver un conjunto de vectores principales (no son únicos). Pero estos vectores principales son extremadamente difíciles de evaluar, en parte por que son funciones discontinuas de los datos. Es como una base ortogonal para el espacio nulo N que sería un conjunto más usual de vectores.

El asunto parece ser poco entendido, para propósitos de problemas y análisis numérico.

Las matrices con eigen-valores múltiples actualmente son muy raras, y una pequeña perturbación computacional de estos destruye la igualdad de eigen-valores.

Por consiguiente se da por hecho que no interese en la práctica que hacer con ellos. Pero en efecto el comportamiento de los divisores no lineales arrastrados en la práctica, para un largo y sorprendente conjunto de

matrices cercanas. Estas matrices cercanas tienen distintos eigen-valores, pero la columna k de eigen-vectores son tan cercanos a la dependencia lineal que no pueden ser separados en una computación normal.

Un aspecto ahora en el problema de mínimos cuadrados, se dice que es un análisis para x, para lograr minimizar.

$$f(x) = \left\| Ax - b \right\|_2$$

Y el propósito realmente es que se desea un mínimo para $f(x)$ ó realmente es que se desea una x que de un valor de $f(x)$ lo más aproximado al mínimo?. En un problema de curva propia por ejemplo, uno puede muchas veces conseguir una cantidad sorprendentemente buena para un polinomio con coeficientes muy diferentes de los polinomios minimizados.

En todos los problemas de los cálculos anteriores es importante verificar cual de los siguientes tipos de respuesta al problema propuesto se está buscando :

- a).- Una supuesta respuesta sin estimadores de su exactitud.
- b).- Alguna respuesta, junto con algunas afirmaciones probabilísticas acerca de su exactitud.
- c).- Alguna respuesta, junto con límites matemáticos probables para su error.

Normalmente es más largo para obtener, alguna respuesta junto con algunas afirmaciones probabilísticas acerca de su exactitud y, una supuesta respuesta sin estimadores de su exactitud, y hasta más largo de obtener alguna respuesta, junto con límites matemáticos probables para su error.

No es muy fácil ver cual de los tipos de respuesta anteriores al problema propuesto querriamos. Frecuentemente el desear una supuesta respuesta sin extimadores de su exactitud es completamente sastifactoria, los Científicos, Físicos e Ingenieros, frecuentemente hacen sus propios chequeos en cuanto a la validez de una respuesta y pueden no necesitar ninguna ni desear los métodos rigurosos de los

Matemáticos. Ellos pueden reconocer por ejemplo, que el modelo matemático es como una ligera aproximación a la realidad y que los métodos matemáticos serían solo utópicos. Cuando los matemáticos entran al mundo práctico de ingeniería, las reglas con las cuales los matemáticos están jugando frecuentemente tienen poca relevancia. El analista numérico frecuentemente tiene el problema de decidir cuando debe de jugar el juego de acuerdo a la reglas matemáticas y cuando jugar el juego de acuerdo a los ingenieros.

Es por supuesto, extremadamente favorable para encontrar esos ejemplos ocasionales, donde los métodos matemáticos probables pueden ser encontrados, estos son justamente tan exactos y fáciles de encontrar como las supuestas respuestas.

Y nunca cesará uno de buscar tales espejismos, porque estos espejismos ocurren .

C A P I T U L O
I V.

EL ESTADO DE LA COMPUTACION
DESDE 1946 HASTA AHORA

Y

LA NATURALEZA DEL HARDWARE
Y SOFTWARE DE LA COMPUTADORA

En las tres últimas décadas, se han logrado avances asombrosos en la construcción y el uso de los sistemas de cómputo. Cada avance hacía obsoletos tecnológicamente los modelos anteriores.

En el año de 1946 vio el nacimiento de la primera generación de computadoras con la computadora ENIAC (Electronic Numerical Integrator and Calculator) una de sus primeras utilidades fue para la construcción de tablas para el cálculo de trayectoria de proyectiles, es decir, este avance vino, como en muchas otras ciencias, por las necesidades militares que surgieron con la segunda guerra mundial.

Después en 1951 surge la computadora UNIVERSAL I, y esto también marcó el inicio del desarrollo de los programas y de los lenguajes de programación. Muchas de las computadoras de la Primera Generación tenían que programarse en un lenguaje de máquina, que consistía en una serie de ceros y unos, lenguaje que era muy difícil de aprender y de utilizar. Como resultado de lo anterior, fueron desarrollados los lenguajes simbólicos y ensamblador. Estos nuevos lenguajes utilizaban símbolos como D para dividir y M para multiplicar, pero tenía que existir un programa que tradujera el lenguaje simbólico ó ensamblador al lenguaje máquina. En 1952 se desarrolló uno de los primeros traductores de programas y este hecho determinó el inicio del desarrollo en los lenguajes de programación.

Durante esta primera generación, los sistemas de computación se utilizaron principalmente con fines científicos.

Mientras que la primera generación de computadoras definitivamente significó una mayoría sobre las máquinas antecesoras, aún dejaba mucho que desear la utilización de bulbos al vacío. Estos bulbos producían altas temperaturas y, en consecuencia, se quemaban con mucha facilidad y obligaba a utilizar costosos sistemas de refrigeración. Otra característica de las computadoras de esta generación era la escasa fiabilidad; por ejemplo, el tiempo medio entre averías de una unidad central era inferior a la hora, esto implicaba que para garantizar el buen funcionamiento de un equipo se necesitaba la total dedicación de un grupo

de personas encargadas del mantenimiento. La forma de ejecutar los trabajos era estrictamente secuencial; el programa que previamente se había perforado en tarjetas, se cargaba en la memoria de la computadora y, a continuación se ejecutaba, procesando las instrucciones de cálculo y las de salidas de información. Estos problemas y el deseo de producir mejores sistemas de cómputo alentaron el desarrollo de transistores, lo cual marcó el inicio de la segunda generación de computadoras.

LA SEGUNDA GENERACION de computadoras utilizó transistores en lugar de bulbos. La razón de ello es que los transistores ofrecían mayores ventajas, tales como: Eran más pequeños, más rápidos, más confiables y producían menos calor durante su operación. Esta segunda generación de computadoras se ubica aproximadamente de 1959 a 1965.

Algunas mejoras en los equipos de computación acompañaron a esta segunda generación de computadoras. Disminuyó el empleo de tarjetas como instrumento de almacenamiento y se introdujo el uso de dispositivos de cinta magnética y discos. Por primera vez, un sistema de cómputo tenía acceso directo a los datos almacenados de manera permanente por medio del disco.

Los dispositivos de entrada y salida fueron perfeccionados y también se desarrollaron mejores impresoras, terminales, lectoras de tarjetas, y además, se inició la utilización de núcleos de memoria (almacenamiento temporal). Los dispositivos de equipo también fueron modulares, una innovación que facilitaba la localización de problemas de equipo y mantenimiento.

También hubo algunas mejoras en los sistemas de programación. Los sistemas operativos ya estaban reemplazando muchas de las funciones que estaban siendo desempeñadas por la persona que operaba el sistema de cómputo. Se desarrollaron los compiladores para traducir lenguajes de programación de alto nivel a lenguaje de máquina y los programas de utilería se desarrollaron para intercalar, ordenar, transferir archivos de datos de una localidad de almacenamiento a otra.

Esta segunda generación fué testigo del cambio en el énfasis de equipo a programas. En muchas instalaciones, el costo de los programas ya se aproximaba al costo del equipo.

Más aún, algunas empresas empezaron a utilizar la computadora para suministrarle a sus ejecutivos la información que podía ayudarles a incrementar las utilidades y a mejor administrar las empresas. Esto fué el inicio de los sistemas computarizados de información administrativa.

Al mismo tiempo empezaron a surgir emopresas de procesamiento de datos. Agunas de éstas se fueron especializando en diversas funciones específicas dentro de la industria del procesamiento de datos, en términos de organizaciones avocadas a la programación ó bien a un servicio de procesamiento de datos. Al principio de la década de los sesentas, muchas de las primeras organizaciones de procesaminto de datos se empezaron a ocupar en la capacidad del personal encargado del procesamiento de datos.

Aún cuando el uso de transistores mejoró considerablemente las capacidades de las computadoras de la segunda generación, científicos e ingenieros seguían trabajando con ahínco en areas de una nueva y más depurada tecnología. Estos esfuerzos llevaron al desarrollo de los circuitos miniaturizados y a la tercera generación de computadoras.

LA TERCERA GENERACION de computadoras se inició con la introducción de la IBM 360, lo cual aconteció en 1965. Los transistores dieron paso a los circuitos integrados y miniaturizados. De 1965 a la fecha, el desarrollo de las computadoras ha sido muy rápido y se han dado muchísimas mejoras en todos los aspectos del sistema de cómputo. Estas mejoras incluyen :

- 1).- Mayor velocidades de precesamiento.
- 2).- Mayor exactitud.
- 3).- La integración del equipo y los programas.

- 4).- La capacidad de desempeñar simultáneamente varias operaciones.
- 5).- Avances en materia de comunicación de datos.
- 6).- Una mejoría en la proporción entre desempeño y precio.

Los programas han mejorado significativamente, al grado de que, hoy en día, el mayor costo de un sistema de cómputo le corresponde a los programas. En un futuro cercano, los programas pudieran significar tanto como un 75% del costo total de un sistema de cómputo. Los lenguajes de programación de alto nivel se usan casi exclusivamente para el desarrollo de programas de aplicación, y ya existe un buen número de nuevos lenguajes de programación de alto nivel, aclarando que en este sentido, la tendencia es hacer más fácil la utilización de estos lenguajes. También existen nuevas formas de diseñar y desarrollar programas de aplicación, y un ejemplo de ello lo constituye el diseño estructurado. Así mismo, han existido unas mejorías realmente dramáticas en lo tocante a programas del sistema y lo mismo puede decirse de los sistemas operativos, los compiladores, ensambladores y los programas de utilería.

La tercera generación también ha permitido hacercar la informática a los usuarios finales, tanto a los profesionales infomáticos como a los de otras especialidades, a través del teleprocesamiento, de los sistemas conversacionales y, sobre todo a través de las computadoras personales. El teleprocesamiento permite realizar la entrada de datos desde terminales remotas, y recibir los resultados en el mismo lugar. Los sistemas conversacionales permiten a los usuarios no sólo enviar y recibir datos desde sus terminales, sino también seguir e intervenir en el desarrollo de sus programas ha través de una " conversación " con el sistema. Por último, los computadores personales han conseguido popularizar la informática y es sorprendente comprobar cómo las prestaciones de una computadora personal, micro-computadora ó minicomputadora, son cada vez más parecidas a las ofrecidas por los grandes equipos tradicionales.

LA CUARTA GENERACION de computadoras ha surgido con los circuitos integrados de alta escala de integración. Algunos autores cuestionan su existencia por no haber sufrido

modificación el método de explotación de los grandes equipos. Sin embargo, el importante auge de las mini-computadoras y microcomputadoras, ha supuesto importantes modificaciones en el método de explotación, por lo que parece incuestionable que nos encontremos en la cuarta generación de computadoras. Algunos discrepan de esta opinión, dado que los procedimientos de explotación no han variado sustancialmente. En cualquier caso, lo que si se puede afirmar es que, en la década de los noventa superaremos e inauguraremos la quinta generación de sistemas informáticos.

Por otra parte han continuado su crecimiento las organizaciones de procesamiento de datos y estas mismas organizaciones están ahora publicando revistas especializadas sobre computadoras y procesamiento de datos. A medida que se constituyen nuevas organizaciones, son reconocidos los logros de sus predecesores.

Además es conveniente decir que la computación de matrices ha pasado por diversas etapas, así como la computación en general. Ciertamente no se puede aprender álgebra lineal teórica un lenguaje de programación algebraico ó algún otro, y enpezar a escribir programas los cuales se ejecutarán aceptablemente por modelos actuales. Ya se ha ganado mucha experiencia en los mejores algoritmos, y esto es mencionado en los textos matemáticos de álgebra lineal.

El monto de literatura en computación de matrices está en etapas. En 1958 Fadeev y Fadeeva en su libro registraron una reseña muy completa de métodos computacionales. Wilkinson da más acerca del conocimiento sobre la computación de eigen-valores y de matrices densas almacenadas y escasas (ambas, simétricas y no simétricas), con límites de error para muchos algoritmos. Hay una diferencia muy pequeña entre estos libros, porque Wilkinson y unos pocos contemporaneos crearon más del material de su libro en los años posteriores a 1958. Posiblemente nadie podría enpezar la búsqueda en las matemáticas numéricas del álgebra lineal sin un completo conocimiento del material relevante en estos libros.

Lo más conocido para los principiantes, para matrices densas ó huecas, la eliminación gaussiana, de supuesto. Sabemos que algunas veces produce pobres resultados. No siempre estamos seguros por que, debatimos eternamente

sobre como opera el pivote para la eliminación, sin arreglarlo. El debate aún continúa, pero ahora principalmente entre personas quienes no entienden que las líneas principales de la respuesta han sido arregladas. Porque de no haber comprendido las dificultades de la eliminación gaussiana, buscamos otros métodos los cuales podrían hacerlo mejor. El método de gradiente conjugado ha sido ideado para matrices densas ó huecas por Lanczos.

Von Neumann y Goldstine, evitaron el problema del pivoteo para reducir cualquier sistema de ecuaciones lineales $Ax=b$ a un sistema positivo definido $A^T A x = A^T b$. Sabemos que está normalización del problema fué costoso en tiempo y de empeorarse la condición del problema Von Neumann y Goldstine presentaron límites de error para la solución; actualmente los errores observados fueron encontrados para ser acaso 100 veces más pequeños en casos razonables.

La forma de error de análisis fué una comparación directa de la máquina aritmética con operaciones exactas. La no-asociatividad y no distributividad de la máquina aritmética hizo el análisis extremadamente difícil. En algún caso, podría solo manejarse encerrado en la aritmética de punto fijo. Porque del tamaño de los límites de los errores, Von Neumann y Goldstine eran incesantemente pesimistas sobre la posibilidad de invertir matrices en general de ordenes sobre 15. en máquinas con 27 bits de precisión.

Para los problemas de eigen valores, el estado de cosas estaba mucho peor, se tenía el método potencia con deflación de matriz. Mientras era razonablemente satisfactorio para unos cuantos renglones dominantes, su aplicación general requiere de intuición y suerte, y supone una algoritmización completa. Para matrices densas y matrices simétricas almacenadas teníamos el método de Jacobi y era muy satisfactorio.

Para matrices no simétricas, las cosas eran horribles. Si el método potencia no trabajara no tendríamos practicamente alternativas. Podríamos examinar para ceros del $\det(A - zI)$ de alguna manera u otra. Se han tratado métodos para determinar la característica polinomial, tan descrito en Faddeeva y encontrandolas irremediables. Fué también increíble tanto como malos los métodos estandar

para $n = 4$, y ejecutarlos para $n = 10$. No obstante el método original de Lanczos necesitó cuidadoso manejo porque el nuevo resultó que con frecuencia era pobre.

El carácter de soluciones realizables a los problemas computacionales de álgebra lineal, está gratamente influenciada por la naturaleza de los sistemas computacionales aprovechables para nosotros. Es acostumbrado separar los sistemas computacionales en las siguientes líneas.

- a).- **HARDWARE.**- Es la naturaleza de los circuitos electrónicos de una computadora, con la palabra Hardware se hace referencia en informática a las diferentes máquinas llamadas computadoras.

Analícemos las características y peculiaridades de las microcomputadoras (computadoras cuya unidad central de proceso - el cerebro ejecutor - es un microprocesador) de mayor vigencia y difusión. Para poder evaluar las diferencias y similitudes que existen entre los diversos sistemas, es preciso aplicar a todos ellos un cuestionario análogo que evidencie claramente la potencia y posibilidades de cada computadora. El avance de la informática y su definitiva penetración en la sociedad actual, se apoya sobre la base de los microprocesadores, la amplia variedad y diversificación de este tipo de sistemas hacen que su campo de aplicación sea casi ilimitado.

La entrada definitiva de IBM en el mundo de la microinformática se ha realizado con la llegada del IBM-PC y el microprocesador 8088, cuya arquitectura interna es de 16 bits, un procesador aritmético destinado al cálculo rápido de operaciones matemáticas en punto flotante, con este se consigue mayor versatilidad y rapidez. El rango de memoria disponible en la IBM-PC oscila entre 64 y 544 kbytes.

- b).- **LENGUAJES DE COMPUTACION.**- Los lenguajes en el cual son descritos los algoritmos para la solución de un problema dado en una computadora.

Con el fin de facilitar el trabajo del programador, surge la necesidad de que la computadora entienda un lenguaje diferente al suyo propio. Entramos en la etapa de la simbolización. Ya el programador no necesita conocer realmente donde ubica sus datos, le basta con referirse a direcciones simbólicas. Así nacen los lenguajes de programación del tipo ensamblador y, consecuentemente, nace el soft-ware traductor ó conjunto de programas que permiten convertir los programas escritos en lenguaje del programador al lenguaje que entiende la máquina. Por esta vía se avanza más y se llega a un nuevo paso que permite al hombre dar al ordenador las fórmulas o notaciones que normalmente usa en su trabajo. Aparecen los lenguajes de programación de alto nivel y los compiladores programas cuya misión es traducirlos al lenguaje de la computadora.

c).- SOFT-WARE.- Son los programas los cuales hacen posible para una computadora actualmente ejecutar los algoritmos descritos en el lenguaje de la computadora.

La característica fundamental de las computadoras es la posibilidad de almacenar en su memoria programas que pueden modificarse y ejecutarse automáticamente. Pero esa ejecución automática implica que existan otros programas que permiten a la computadora coordinar el que varios programas puedan ejecutarse, bien secuencialmente. También es evidente que existe un gran número de operaciones repetitivas, tales como clasificar un grupo de datos en un orden determinado, que podrían prepararse de forma que el usuario del ordenador no tenga que preocuparse de ello. Esto es el soft-ware, también denominado " el componente lógico del sistema informático ".

En un sentido más estricto el soft-ware es el conjunto de programas que se utilizan en una computadora. Es el conjunto de programas y otras ayudas generalmente proporcionados por el fabricante de la computadora que facilitan al usuario una operación más eficiente del equipo.

Tipos de soft-ware. Podemos clasificar los programas en cuatro grandes grupos.

Soft-ware específico, ó programas de aplicación que pueden ser escritos tanto por el constructor, como por el usuario, todos los problemas técnicos y científicos pertenecen a este grupo.

Soft-ware traductor ó programas de ayuda para escribir nuevos programas. Constituido por los programas que permiten que los programas escritos por los usuarios en un lenguaje distinto al de la máquina se conviertan en programas con instrucciones en código de lenguaje de máquina. Son escritos y proporcionados por los fabricantes.

Soft-ware funcional ó programas de ayuda para ejecutar otros programas. Más comunmente conocido como sistema operativo. Es un conjunto de programas que facilitan una explotación más racional de las computadoras, guiando todas las tareas y ayudando a los programas en ciertas funciones. Son desarrollados por los constructores.

Soft-ware ó rutinas de utilidad. Programas que permiten la realización de funciones de uso frecuente y que generalmente son escritos por los fabricantes; aunque también pueden desarrollarlos los propios usuarios.

Al buscar en el Hardware de la computadora, para cálculos de álgebra lineal, se quiere saber que precisión es aprovechable para computación, cuantos dígitos están en el significante de los operadores de punto flotante, y en que base.

También se está interesado en el costo y rapidez de las operaciones de doble precisión. En el trabajo de álgebra de matrices la operación crítica es frecuentemente la computación de la aproximación del redondeo de precisión simple al de doble precisión, el producto interno de dos vectores en los cuales sus componentes tienen números de punto flotante con precisión simple.

C A P I T U L O

V

PROBLEMAS DE LAS
ECUACIONES LINEALES

Para resolver el problema de las ecuaciones lineales de la forma $Ax = b$. La eliminación gaussiana es un algoritmo perfecto, excepto si un problema en particular tiene algunas propiedades especiales, como las tienen casi todos los problemas.

Para matrices grandes y huecas como aquellas que aparecen en aproximaciones de diferencias finitas en ecuaciones diferenciales parciales. Los métodos vistos dependen para su ocurrencia de sus éxitos en la naturaleza de los problemas continuos siendo aproximados. Porque si las matrices son huecas, los métodos que prevalecen son los iterativos. Omitiremos la discusión adicional de estas y daremos atención a las matrices densas almacenadas.

El análisis simple de multiplicaciones anidadas (Capitulo XII) produjo grandes resultados; muestra que un detalle de análisis de error de redondeo de muchos posibles algoritmos es completamente innecesario. Son excluidos a-priori por simples consideraciones de las limitaciones inherentes.

Para la solución de sistemas lineales la situación no es tan extrema; sin embargo un análisis de las limitaciones inherentes es otra vez muy ilustrativa. Y ha precedido al intento de llevar a cabo un detallado análisis de error de la eliminación gaussiana muchas de las confusiones que han sido evitadas.

Antes de esto examinemos las propiedades de la solución exacta del sistema :

$$Ax = b \quad (1)$$

Para cualquier matriz consistente y un vector normado tenemos de.

$$Ax = b$$

$$\text{que } \| A \| \| x \| \leq \| b \| \quad (2)$$

La solución de

$$Ax = b$$

-1

$$\text{Está dada por : } x = A^{-1} b \quad (3)$$

Y dado que

$$\|x\| \leq \|A^{-1}\| \|b\| \quad (4)$$

muestra que

$$\|b\| / \|A\| \leq x \leq \|A^{-1}\| \|b\| \quad (5)$$

Si A está definida por

$$\|A\| = \min. \lim. \sup. \|Ax\| / \|x\| \quad (6)$$

Entonces ambos límites son alcanzables en

$$\|b\| / \|A\| \leq x \leq \|A^{-1}\| \|b\|$$

de aquí que el radio del valor más grande de x, al más pequeño correspondiente a un valor dado de \|b\|, es

$$\|A\| \|A^{-1}\|.$$

y es escribimos :

$$k(A) = \|A\| \|A^{-1}\| \quad (7)$$

donde k(A) es llamada la condición de A con respecto a la inversión en la norma dada; la A en k(A) será de nuestro interés principalmente con las normas L1, L2 y L ∞ y la

norma Frobenius. De la relación AA⁻¹ = I, tenemos,

$$k(A) = \|A\| \|A^{-1}\| \geq 1 \quad (8)$$

Cuando k(A) es grande se dice que la matriz está mal condicionada con respecto a la inversión, por algunas razones, las cuales serán discutidas después.

En el caso de la norma L2, entra en la variación de x para una b fija, está estipulada por el valor singular de descomposición.

$$A = U \text{diag}(\sigma_i) V^T \quad (9)$$

De aquí U y V son ortogonales y la σ_i , la cual puede ser ordenada de tal que,

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0 \quad (10)$$

Son los valores singulares de A. La condición $\sigma_n \geq 0$ está basada en la suposición de que A es no singular. De las relaciones,

$$\|A\|_2 = \sigma_1, \quad \|A^{-1}\|_2 = 1 / \sigma_n \quad (11)$$

Tenemos

$$k(A) = \sigma_1 / \sigma_n \quad (12)$$

Y las matrices mal condicionadas se caracterizan por la propiedad $\sigma_1 \gg \sigma_n$. Supongamos ahora que b es un vector por la derecha. Podemos escribir, entonces

$$b = \|b\|_2 \sum \alpha_i u_i, \text{ donde } \sum \alpha_i^2 = 1 \quad (13)$$

Ahora de la relación,

$$A = U \text{diag}(\sigma_i) V^T$$

implica que :

$$AV = U \text{diag}(\sigma_i), \text{ es decir, que } AV = \sigma_i u_i \quad (14)$$

De aquí la solución de $Ax = b$, está dada por :

$$x = \| \| b \| \|_2 \sum (\alpha_i / \sigma_i) v_i \quad (15)$$

Dando
$$\|x\|_2 / \|b\|_2 = \left[\sum_i (\alpha_i / \sigma_i)^2 \right]^{1/2}$$

donde
$$\alpha_i / \sigma_i = \alpha_i \|A\|_2^{-1} \quad (16)$$

A no ser que b pase a ser patologicamente deficiente en su componente α_i de u_i , la relación

$$\|x\|_2 / \|b\|_2 = \left[\sum_i \alpha_i^2 / \sigma_i^2 \right]^{1/2}, \quad \alpha_i / \sigma_i = \alpha_i \|A\|_2^{-1}$$

muestra que $\|x\|_2 / \|b\|_2$ es del orden de magnitud de $\|A\|_2^{-1}$.

Dado que los u_i son ortogonales, cuando A está mal condicionada, la probabilidad de que un rango b de una x , la cual es del orden de magnitud del límite superior en la siguiente relación es grande, más bien que del límite inferior.

$$\|b\|_2 / \|A\|_2 \leq x \leq \|A\|_2^{-1} \|b\|_2$$

Si varias de las σ_i son comparables con σ_n , entonces

$\|x\|_2 / \|b\|_2$ serán del orden de magnitud de

$\|A\|_2^{-1}$, a menos que todos los α_i correspondientes sean

pequeños y la probabilidad de que ésto suceda es totalmente mínima.

Si en el otro lado tomamos un rango x , y escribimos :

$$x = \|x\|_2 \sum_i \beta_i v_i \quad \text{donde} \quad \sum_i \beta_i^2 = 1 \quad (17)$$

entonces

$$b = Ax = \left\| \left\| x \right\| \right\| \Sigma \beta_i \sigma_i u_i$$

Dado que $\sigma_n \ll \sigma_1$, para una A mal condicionada la probabilidad es grande de que una b construida de esta manera, sea particularmente débil en su componente u_n .

Realmente esto será verdad, a menos que una x pase a ser excepcionalmente deficiente en todos estos u_i , para los cuales $\sigma_i \gg \sigma_n$. Este resultado es importante, porque se ha experimentado con frecuencia, y construye del lado derecho vectores de rango x, para tener sistemas de los cuales se saben las soluciones. Tales lados derechos constituyen una inclinación muy simple.

Teoría de Perturbación para Sistemas Lineales.

El objetivo principal es que es casi inconcebible que pueda haber un algoritmo general propuesto para resolver $Ax = b$, el cual podría garantizar hacer mejor el trabajo que para dar la solución exacta de algún sistema perturbado.

$$(A + \delta A) x = b + \delta b$$

con

$$\delta a_{ij} \leq k \beta_{ij}^{1-t} |a_{ij}|, \quad |\delta b_i| \leq k \beta_i^{1-t} |b_i|$$

Y en cuanto es probable para involucrar $O(n^3)$ operaciones, los límites grandes para δA y δb son virtualmente inevitables.

Por lo tanto discutiremos aquí el efecto de las perturbaciones que satisface.

$$|\delta a_{ij}| \leq \epsilon |a_{ij}|, \quad |\delta b_i| \leq \epsilon |b_i|$$

Dado que tenemos límites para las perturbaciones en elementos individuales, no usaremos ninguno de estos; en efecto nuestro análisis involucra solo los requerimientos.

$\| \delta A \| \leq \epsilon \| A \|$, $\| \delta b \| \leq \epsilon \| b \|$, los cuales son consecuencia de .

$$| \delta a_{ij} | \leq \epsilon | a_{ij} | \quad , \quad | \delta b_i | \leq \epsilon | b_i |$$

Para las normas L1, L ∞ . Para la norma L2

$$\delta a_{ij} \leq k \beta^{1-t} | a_{ij} | \quad , \quad | \delta b_i | \leq k \beta^{1-t} | b_i |$$

implica solo que :

$$\| \delta A \|_2 \leq n^{\frac{1}{2}} \epsilon \| A \|$$

pero ignoraremos esto por el momento. Observamos primero que podemos dejar la perturbación b. En efecto

$$(A + \delta A) x = b + \delta b$$

puede ser expresado de la forma.

$$(A + \delta A - \delta b x^T / x^T x) x = b$$

y tenemos.

$$\| \delta b x^T / x^T x \|_2 = \| \delta b \|_2 / \| x \|_2 \leq \| b \|_2 / \| x \|_2$$

y de $(A + \delta A) x = b + \delta b$

resulta

$$(\| A \|_2 + \| \delta A \|_2) \| x \|_2 \geq \| b \|_2 - \| \delta b \|_2$$

dando

$$\| A \|_2 (1 + \epsilon) \| x \|_2 \geq \| b \|_2 (1 - \epsilon)$$

de aquí

$$\| \delta b \|^T / \| x \|^T \leq \epsilon \| A \|^T (1 + \epsilon) / (1 - \epsilon) = \epsilon \| A \|^T$$

Si $\epsilon \ll 1$, y si estamos preparados para doblar el límite en δA , podemos distribuir con δb . Nótese que cuando A está mal condicionada.

$$\| A \|^T (1 + \epsilon) \| x \|^T \geq \| b \|^T (1 - \epsilon)$$

será muy débil y para un rango b la probabilidad es grande de que.

$$\| b \|^T / \| x \|^T \ll \epsilon \| A \|^T (1 + \epsilon) / (1 - \epsilon)$$

y de aquí δb puede ser retomada por una perturbación adicional, y A es mucho más pequeña que $\epsilon \| A \|^T$.

En cambio no es posible retomar la δA con el modesto argumento b , tenemos en efecto.

$$Ax = b + \delta b - \delta Ax$$

Para un rango b vimos que $\| x \|^T$ es probable de ser del orden de magnitud de $\| A \|^T \| b \|^T$ y dado que

$$\| \delta Ax \|^T \leq \| \delta A \|^T \| x \|^T = 0$$

$$(\epsilon \| A \|^T \| A \|^T \| b \|^T) = 0 \approx k \| b \|^T$$

si ahora escribimos

$$(A + \delta A) y = b, \quad Ax = b, \quad \| \delta A \|^T \leq \epsilon \| A \|^T$$

Tenemos

$$(I + A^{-1} \delta A) y = x$$

Observamos que, a menos que

$$\| A^{-1} \delta A \| < 1$$

No podemos garantizar que $A + \delta A$ no sea singular, tenemos

$$\| A + \delta A \|^{-1} \leq \| A \|^{-1} \| \delta A \| \leq \epsilon \| A \| \| A \|^{-1} = \epsilon k$$

En general no podemos esperar que y sea en algún sentido una aproximación de x a no ser que ϵk sea significativamente más pequeña que la unidad y esto se puede asumir desde ahora.

De

$$(I + A^{-1} \delta A) y = x$$

tenemos

$$y - x = -A^{-1} \delta A y$$

$$\| y - x \| \leq \epsilon k \| y \|$$

y

$$\| x \| \geq \| y \| (1 - \epsilon k)$$

Dado

$$\| y - x \| / \| x \| \leq \epsilon k / (1 - \epsilon k)$$

La expresión de la izquierda es el error relativo en y y vemos que está efectivamente limitado por ϵk . Para casi todos los x y δA ésta es una estimación realista para el error relativo. La presencia del factor k en el límite de error para cualquier algoritmo, para resolver sistemas lineales es por consiguiente natural y realmente inevitable.

Cuando A es positiva definida se cumple que

$$\sigma = \gamma_1 / \gamma_n$$

$$\sigma = \gamma_1 / \gamma_n$$

y

$$k = \gamma_1 / \gamma_n$$

Supongamos ahora que se intenta asignar un aproximado vía el tamaño del correspondiente residual, tenemos

$$r = b - A y = \delta A y$$

$$\| r \| \leq \| \delta A \| \| y \| \leq \epsilon \| A \| \| x \| / (1 - \epsilon k)$$

este resultado es trivial y mucho más remarcable que el visto.

Consideremos una solución aproximada z teniendo un error del mismo orden de magnitud, como aquel en y pero arbitrario.

$$\text{De } \| y - x \| / \| x \| \leq \epsilon k / (1 - \epsilon k)$$

podemos escribir

$$z = x + w \quad \text{donde } \| w \| \leq \epsilon k \| x \| / (1 - \epsilon k)$$

y de aquí

$$b - Az = b - Ax - Aw = -Aw$$

Dando

$$\| b - Ax \| \leq \| A \| \| w \| \leq \epsilon k \| A \| \| x \| / (1 - \epsilon k)$$

Este límite es más grande por el factor k que el que se tiene en y . De este la totalidad de las soluciones, el tener errores del orden de magnitud $\epsilon k \| x \| / (1 - \epsilon k)$ casi todos ellos dan muy grandes residuales que los otros, las cuales son soluciones exactas correspondientes a una matriz $A + \delta A$ con una perturbación relativa de orden ϵ . Hay mucho más aspectos remarcables de esta propiedad de y . Consideremos una solución aproximada $x + w$ con

$$\| w \| \leq \epsilon \| x \|$$

Es decir, las cuales realmente tienen un error relativo bajo ahora tenemos

$$\begin{aligned} & \| b - A(x + w) \| = \\ & = \| -Aw \| \leq \| A \| \| w \| \leq \epsilon \| A \| \| x \| \end{aligned}$$

ignorando el factor $1/(1 - \epsilon k)$, el límite en el residual para una solución muy precisa no es más pequeña que para y . Un ejemplo particular de como una $x + w$ es estipulada por la correcta versión de redondeo de x , en el caso cuando los elementos de x no son digitales. Para esto tenemos,

$$\| w \| \leq k \beta^{1-t} \| x \|$$

Veamos que en toda esta cantidad de virtudes entre soluciones aproximadas, no cabe esperar que de un residual más pequeño que la solución exacta de y .

$$(A + \delta A) y = b \quad \text{con} \quad \| \delta A \| \leq k \beta^{1-t} \| A \|$$

Nos lleva al hecho de que

$$y = x - A^{-1} \delta Ax + A^{-1} \delta AA^{-1} \delta Ax$$

y de aquí

$$\begin{aligned} Ay &= Ax - AA^{-1} \delta Ax + \text{terminos más chicos} \\ &= b - \delta Ax + \dots \end{aligned}$$

El principal error en y es $A^{-1} \delta Ax$, y cuando esto se multiplica por A , la A y A^{-1} se cancelan si A es sustituida por cualquier otra matriz B con la misma norma, entonces podemos decir solo que,

$$\| AB \| \leq \| A \| \| B \| = \| A \| \| A^{-1} \| = k(A)$$

De aquí el factor extra k .

Notese que a no ser que la solución exacta tenga elementos digitales entonces, no habrá en general alguna solución digital aproximada, dando un residual más pequeño que

$$k \beta^{1-t} \| A \| \| x \| \quad \text{el cual corresponda a un residual relativo.}$$

$$k \beta^{1-t} \|A\| \|x\| / \|b\| \leq k \beta^{-t} \|A\| \|A^{-1}\|$$

y para una b aleatoria hemos visto que la x correspondiente es tal que este es un límite realista.

Hemos propuesto consideraciones detalladas de una perturbación en b porque en general esto es menos importante.

Supongamos ahora que

$$Ay = b + \delta b \quad \|\delta b\| \leq \epsilon \|b\|$$

entonces

$$y = A^{-1} b + A^{-1} \delta b = x + A^{-1} \delta b$$

Dando

$$\|y - x\| \leq \epsilon \|A^{-1}\| \|b\| \leq \epsilon \|A^{-1}\| \|A\| \|x\|$$

ó

$$\|y - x\| / \|x\| \leq \epsilon k$$

Aparte del factor $1 / (1 - \epsilon k)$, este es el mismo límite para el error relativo en y, como está dado en

$$\|y - x\| / \|x\| \leq \epsilon k / (1 - \epsilon k)$$

Para la perturbación δA en A. No obstante éste está perdido en

$$\|y - x\| \leq \epsilon \|A^{-1}\| \|b\| \leq \epsilon \|A^{-1}\| \|A\| \|x\|$$

b es sustituido por $\|A\| \|x\|$ en la última desigualdad y para la mayoría de b esta en una desigualdad extremadamente débil.

Para la mayoría de b hemos visto que ,

$$\|x\| = 0 \|A^{-1}\| \|b\|$$

y cuando esto es verdad tenemos

$$\| y - x \| \leq \epsilon \| A \|^{-1} \| b \| = 0 \epsilon \| x \|$$

Tal que para la mayoría de b, cualquier perturbación δb tal que $\| \delta b \| \leq \epsilon \| b \|$ da solo un error relativo pequeño en la y correspondiente.

Quizá en este punto debemos enfatizar dos factores.

- 1).- Hay siempre una aproximación digital a x teniendo un error relativo bajo. En efecto si x es el vector R derivado de x por redondear cada elemento a la aproximación digital correspondiente será.

$$\| x - x_R \| / \| x \| \leq \frac{1}{2} \beta^{1-t}$$

- 2).- Cuando A está mal condicionada entonces, para la mayoría de b no habrá aproximación digital para x dando un residual relativo bajo. Hemos visto que toda x (la cual R puede ser considerada como la mejor aproximación digital)

dará un residual relativo del orden de magnitud de $k \beta^{1-t}$ para la mayoría de b . En general es solo cuando b es tal que $x = 0$ ($\| b \| / \| A \|$) y esto hará aproximaciones digitales para x dando un residual relativo bajo.

Solución Práctica de Sistemas de Ecuaciones Lineales

Aquí discutiremos brevemente algunas concideraciones computacionales prácticas, concernientes a la eliminación gaussiana.

Muchos problemas en las ciencias aplicadas involucran la inversión de matrices de $m \times n$ ó la solución de n ecuaciones con n incógnitas. Si n es grande (más grande que 3 o 4) ó las entradas en la matriz no son enteros con solo pocos dígitos la solución de tales problemas a mano son tediosos e ingratos, en la práctica estos problemas son resueltos por lo general por programas de cómputo, en computadoras digitales electrónicas, pero con la guía

humana para simplificar el problema del comienzo. El procedimiento sistemático da una excelente base para determinar la naturaleza del conjunto de soluciones y para calcular las soluciones, podemos hacer de este procedimiento las bases para el método computacional elegido, si se toman las precauciones propias. Antes debemos considerar como hacerlo, haciendo algunos comentarios acerca de la naturaleza del sistema de ecuaciones, el cual surge en la práctica, tal que el lector pueda ver el criterio usado para medir las eventualidades de cualquier procedimiento propuesto.

Consideremos el sistema de ecuaciones lineales.

1).- El problema real que queremos resolver ha sido modelado de alguna manera matemática.

Supongamos que un paso en la solución del modelo matemático requiere la solución del conjunto de ecuaciones

$$A^* x^* = b^*$$

Donde A^* y b^* , son matrices, cuyos elementos están precisamente definidos, pero generalmente no necesariamente conocidos.

2).- Dado el sistema ideal $A^* x^* = b^*$, tenemos actualmente a la mano un conjunto de ecuaciones

$$A x = b$$

Donde tenemos que $A \neq A^*$ y $b \neq b^*$. Los errores en los datos son a menudo también al menos en parte debido a errores de redondeo y para una conversión binaria cuando los números son almacenados en la computadora; comúnmente los datos también llevan un error experimental causado por nuestra poca habilidad para medir A^* , b^* precisamente.

El sistema

$$A x = b$$

por supuesto tiene una solución exacta, que es, una para la cual $A x - b$ es exactamente igual a cero, con tal de que las computadoras se involucren en evaluar $A x - b$ exactamente, pero lo mejor que podemos esperar hacer en una computadora digital ó en una calculadora de escritorio ó

probablemente a mano, es obtener una solución aproximada x' , no igual a la solución exacta x generalmente, para

$$A x = b$$

Nuestros métodos prácticos para resolver las ecuaciones podría ser, por supuesto, solamente con el dato aproximado A, b en la ecuación $Ax = b$. Cuando tratamos de evaluar lo que hacen nuestros métodos, esto es, como se produce la precisión y aproximación x' , podemos solamente medir la precisión con respecto a la solución exacta x para los datos que nos dan en $Ax = b$, trataremos de garantizar que nuestros procedimientos hacen el error $x - x'$ "pequeño".

Los científicos aplicados, por supuesto que solo quieren la respuesta x' , y dado que los errores $x - x'$ no nos conciernen de momento, entramos al procedimiento de eliminación gaussiana para obtener alguna idea sobre el monto del trabajo involucrado en la eliminación gaussiana, ambos para resolver ecuaciones y para invertir matrices, y enseguida lo concerniente a ejemplos concretos con un tratamiento más teórico de la solución de sistemas lineales de la forma $A x = b$.

Por simplicidad consideremos las adiciones y las sustracciones juntas como "adiciones" multiplicaciones y divisiones son consideradas juntas como multiplicaciones. Consideremos la solución de k conjuntos de ecuaciones como en

$$A x = b$$

entonces

$$A_j x_j = b_j \quad \text{con } j = 1, 2, \dots, k$$

El cálculo de los elementos de

$$A^{(r+1)}, b^{(r+1)}, \text{ de } A^{(r)}, b^{(r)}$$

involucran lo siguiente,

1).- Cálculo de $m_{ir}^{(r)}$, es decir.

$$m_{ir}^{(r)} = \frac{a_{ir}^{(r)}}{a_{rr}^{(r)}}, \quad \text{con } i = r+1, \dots, n$$

involucrando $n-r$ divisiones.

2) El cálculo de $a_{ij}^{(r+1)}$, de

$$a_{ij}^{(r+1)} = a_{ij}^{(r)} - m_{ir}^{(r)} a_{rj}^{(r)} \quad \text{con } i, j = r+1, \dots, n$$

involucrando $(n-r)^2$ multiplicaciones y sumas.

3) El cálculo de k conjuntos $b_i^{(r+1)}$ involucrando $k(n-r)$ multiplicaciones y adiciones.

$$b_i^{(r+1)} = b_i^{(r)} - m_{ir}^{(r)} b_r^{(r)} \quad \text{con } i=r+1, \dots, n$$

donde para $r = 1$, debemos tomar

$$a_{ij}^{(1)} = a_{ij}$$

$$b_i^{(1)} = b_i$$

La eliminación progresiva para obtener.

$$a_{11} x_1 + a_{12} x_2 + \dots + a_{1n} x_n = b_1 \quad (1)$$

$$a_{22} x_2 + \dots + a_{2n} x_n = b_2 \quad (2)$$

$$a_{nn} x_n = b_n \quad (n)$$

Por consiguiente se requiere de

$$\sum_{r=1}^{n-1} \{ (n-1)^2 + (k+1)(n-r) \} = n \{ (1/3 n^2 - 1/3) + 1/2 k(n-1) \}$$

multiplicaciones

$$\sum_{r=1}^{n-1} \{ (n-1) + k(n-r) \} = n \{ (1/3 n^2 - 1/2 n + 1/6) + 1/2 k(n-1) \}$$

sumas ó adiciones

donde hemos usado los resultados

$$\sum_{s=1}^m s = 1/2 m(m+1)$$

$$\sum_{s=1}^m s^2 = 1/6 m(m+1)(2m+1)$$

En la solución de (1) por sustitución hacia atrás de x , involucra $n-r$ multiplicaciones y adiciones y una división. Dejando la división como equivalente a una multiplicación, las k sustituciones hacia atrás involucran.

$$k(1 + 2 + \dots + n) = 1/2 k n(n+1) \text{ multiplicaciones}$$

y

$$k(1 + 2 + \dots + n-1) = 1/2 k n(n-1) \text{ adiciones.}$$

Los totales para la solución completa de k conjuntos de ecuaciones son ;

$$n \left(\frac{1}{3} n^2 - \frac{1}{3} n + k n \right) \text{ multiplicaciones}$$

y

$$n \left(\frac{1}{3} n^2 - \frac{1}{2} n + \frac{1}{6} + k (n - 1) \right) \text{ adiciones}$$

Cuando invertimos A , empezamos con n vectores unitarios e_j en vez de b_j a la derecha de

$$A x_j = b_j$$

Al recordar con cuidado donde hay ceros y unos a la derecha y entonces, evitando multiplicaciones por unos y ceros y sumas de ceros, podemos facilmente reducir el trabajo total para la eliminación progresiva a,

$$\frac{1}{2} n^2 (n - 1) \text{ multiplicacion}$$

y

$$\frac{1}{2} n^2 (n - 1) \text{ adiciones}$$

Y también siendo cautelosos en la sustitución hacia atrás, finalmente llegamos a los requerimientos totales para la inversión entera.

$$n^3 \text{ multiplicaciones}$$

y

$$n^3 - 2n^2 + n \text{ adiciones}$$

Las conclusiones importantes de estos resultados para la eliminación gaussiana son .

- 1).- El número de multiplicaciones y adiciones requeridas para resolver un conjunto simple de ecuaciones son del orden de $\frac{1}{3} n^3$

2).- El número de multiplicaciones y adiciones requeridas para invertir una matriz son del orden de n^3

Esto no dice que :

1).- El monto de trabajo es proporcional a n^3 tal que si duplicamos el número de ecuaciones tenemos que hacer 8 veces el monto del trabajo involucrado, y rápidamente se convierte en sustancial, a medida que el sistema se incrementa.

2).- Alrededor de 3 veces del monto de trabajo es involucrado cuando invertimos una matriz, comparado con resolver un conjunto de ecuaciones, no n veces, como podríamos esperar. El ratio es esencialmente independiente de la medida del sistema.

Para una n muy grande estas operaciones contadas, son

consideradas un tanto impropias, dado que n^3 será extremadamente grande. Sin embargo, se tiende a olvidar que

multiplicar dos matrices de orden n juntas requieren de n^2 multiplicaciones y adiciones. De aquí el monto de trabajo

requerido para formar A^{-1} ó $A^{-1}A$. Además se sabe que en ningún método las operaciones son enteramente por renglón y columna, y puede requerir menos que este monto de trabajo.

En suma, justamente olvidamos que las computadoras pueden operar rápidamente, para implementar la eliminación gaussiana. Pues a mano es laborioso resolver 5 ecuaciones con 5 incógnitas. En el presente las computadoras es lo más práctico para resolver sistemas de 100 ecuaciones con 100

incógnitas, con todos los 100 coeficientes distintos de cero. Las operaciones aritméticas toman más o menos 10 segundos en hacerlo y el costo es menor que un dolar. Este es el porque las computadoras han revolucionado el camino en el cual formulamos los problemas en las ciencias aplicadas. Sin embargo todavía titubearíamos antes de tomar un sistema con $n = 100$, donde la mayoría de los coeficientes son cero.

El Método de Gauss

En las páginas anteriores se ilustraron algunas maneras en las cuales surgen las matrices en las ciencias aplicadas. En lo que sigue veremos las matrices desde un punto de vista más teórico, para obtener un entendimiento más amplio de las potencialidades y las limitaciones de los métodos matriciales.

Como se ha ilustrado en los ejemplos considerados, la utilidad de las matrices en las ciencias aplicadas, está relacionada a menudo con el hecho de que se lleve a cabo un método conveniente para formular problemas físicos en términos de un conjunto de ecuaciones lineales algebraicas. Sin embargo es importante entender teóricamente las diversas situaciones que pueden surgir cuando se resuelven conjuntos de ecuaciones lineales. Este es el objetivo de lo siguiente.

Pasando a la eliminación gaussiana. El punto inicial es el sistema.

$$A x = b$$

comenzaremos con un ejemplo en tres dimensiones, con el sistema.

$$\begin{aligned} 2x + y + z &= 1 \\ 4x + y &= -2 \\ -2x + 2y + z &= 7 \end{aligned}$$

Nuestro problema es encontrar el valor de las incógnitas x, y, z , y aplicando la eliminación gaussiana tenemos.

- a).- Sustrayendo de la segunda ecuación, la primera multiplicada por 2;
- b).- Sustrayendo de la tercera ecuación, la primera multiplicada por -1.

el resultado es un sistema equivalente de ecuaciones.

$$\begin{aligned} 2x + y + z &= 1 \\ -y - 2z &= -4 \\ 3y + 2z &= 8 \end{aligned}$$

c).- Al sustraer de la tercera ecuación la segunda multiplicada por -3 nos queda.

De tal forma que hemos obtenido de la ecuación original.

$$Ax = \begin{bmatrix} 2 & 1 & 1 \\ 4 & 1 & 0 \\ -2 & 2 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 1 \\ -2 \\ 7 \end{bmatrix}$$

Un sistema equivalente pero más simple, con una matriz de coeficientes nueva que denotaremos por U :

$$Ux = \begin{bmatrix} 2 & 1 & 1 \\ 0 & -1 & -2 \\ 0 & 0 & 4 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 1 \\ -4 \\ -4 \end{bmatrix}$$

Ahora la matriz de coeficientes es triangular superior : son cero todas las entradas debajo de la diagonal principal.

El lado derecho, que es un nuevo vector c, se derivó del vector original b mediante los mismos pasos que cambiaron a A en U, de modo que la eliminación gaussiana consiste en: Comenzar con A y b;

Aplicar los pasos a), b) y c), en ese orden;

Terminar con U y c.

La última etapa de la solución $Ux = c$, mediante la sustitución hacia atrás.

La matriz A puede escribirse, mientras que ninguno de los

pivotes sea cero, como un producto LU de una matriz triangular inferior L y una matriz triangular superior U, las entradas de la diagonal principal de L son unos; debajo de la diagonal están los multiplicadores l_{ij} de la fila j

que fueron sustraídos de la fila i en el proceso de eliminación. U es la matriz de coeficientes que obtenemos después de la eliminación y antes de la sustitución

regresiva; las entradas de su diagonal son los pivotes.

La Existencia de Soluciones para un conjunto de ecuaciones.
Algunos ejemplos.

Hemos visto hasta aquí los métodos para encontrar la solución numérica de ecuaciones lineales simultáneas, pero no hemos puesto mucha atención a la cuestión de que tales soluciones existan, ó que la solución sea única. Este es el principal problema que ahora trataremos.

Consideremos ahora la ecuación simple $ax = b$, donde a, x, b , son escalares. Y decimos inmediatamente que la solución de esta ecuación es $x = b / a$, pero en esto hay tres posibilidades.

1).- Si $a \neq 0$, entonces $x = b/a$, y ésta es la única solución de la ecuación, para cualquier valor de b .
(El cual puede ser cero, en tal caso la solución es $x = 0$)

2).- Si $a = 0$, hay dos posibilidades, dependiendo del valor de b .

a).- si $b \neq 0$, entonces la ecuación es $0 x = b \neq 0$ y no existe una solución finita. La solución " $x = \infty$ ", no puede ser considerada como una solución posible. Decimos que "la solución no existe", ó, alternativamente, que "la ecuación es inconsistente", implica que $0 = b \neq 0$, lo cual es una contradicción.

b).- si $b = 0$, entonces cualquier número es una solución de la ecuación, para $0 x = 0$, siempre que el valor sea dado a x . Se excluye otra solución infinita, dado que 0∞ no está definido. Es un notable resultado que precisamente las mismas posibilidades existen en el caso de dos ecuaciones con dos incógnitas. Así,

$$\begin{matrix} x & + & x & = & 2 \\ 1 & & 2 & & \end{matrix}$$

$$\begin{matrix} x & - & x & = & 0 \\ 1 & & 2 & & \end{matrix}$$

tienen una solución única.

$$x_1 + x_2 = 2$$

$$x_1 + x_2 = 1$$

son inconsistentes, y

$$x_1 + x_2 = 2$$

$$2x_1 + 2x_2 = 4$$

tienen una infinidad de soluciones,

sean $x_1 = k$, $x_2 = 2 - k$ para todo k .

Y un resultado más notable es que precisamente las mismas posibilidades existen en el caso general de n ecuaciones con n incógnitas, como veremos. Para algún ejemplo específico es posible descubrir que la situación es por solución directa de la ecuaciones ó, equivalentemente, por usar la reducción triangular de la matriz aumentada, así consideremos las ecuaciones.

$$x_1 + 2x_2 - 5x_3 = 2$$

$$2x_1 - 3x_2 + 4x_3 = 4 \quad (1)$$

$$4x_1 + x_2 - 6x_3 = 8$$

La matriz aumentada es,

$$\begin{bmatrix} 1 & 2 & -5 & 2 \\ 2 & -3 & 4 & 4 \\ 4 & 1 & -6 & 6 \end{bmatrix}$$

Las operaciones por renglón reducen esto a,

$$\begin{bmatrix} 1 & 2 & -5 & 2 \\ 0 & 1 & -2 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Nótese que todos los elementos en el tercer renglón son cero

Esto significa que la tercera ecuación se ha reducido a $0x = 0$ tal que una solución de la ecuación original

existe para la cual $x_3 = k$, donde k es cualquier número

Por la sustitución hacia atrás, nos da.

$$x_2 = 2k$$

$$x_1 = 2 + k$$

Y ésta solución puede ser expresada de la forma

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \\ 0 \end{bmatrix} + k \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$$

Si de otra manera, la tercera ecuación en (1) hubiera sido

$$4x_1 + x_2 - 6x_3 = 0$$

encontraríamos en lugar de

$$\begin{bmatrix} 1 & 2 & -5 & 2 \\ 0 & 1 & -2 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

lo siguiente

$$\begin{bmatrix} 1 & 2 & -5 & 2 \\ 0 & 1 & -2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

La tercera ecuación es ahora $0 \cdot x_3 = 1$, tal que este conjunto de ecuaciones es inconsistente.

En los ejemplos consideramos, que el número de ecuaciones ha sido exactamente igual al número de incógnitas, pero los mismos métodos (llamados reducción de una matriz aumentada

a la forma triangular ó, equivalentemente, la solución directa de las ecuaciones), pueden ser utilizados para checar un conjunto dado de m ecuaciones con n incógnitas, para algún m, n , que posea ó no solución, una solución única, ó una infinidad de soluciones. Así podemos verificar los siguientes ejemplos siguiendo el método descrito anteriormente.

Por ejemplo, si la ecuación

$$x_1 + x_2 + x_3 = 6$$

es sumada al conjunto de ecuaciones

$$\begin{array}{r} x_1 + 2x_2 - 5x_3 = 2 \\ 2x_1 - 3x_2 + 4x_3 = 4 \\ 4x_1 + x_2 - 6x_3 = 8 \end{array}$$

como una cuarta ecuación; el sistema resultante de 4 ecuaciones con 3 incógnitas, tiene solución única.

$$x_1 = 3$$

$$x_2 = 2$$

$$x_3 = 1$$

Si la ecuación

$$3x_1 - x_2 - x_3 = 6$$

es sumada al conjunto de ecuaciones

$$\begin{array}{r} x_1 + 2x_2 - 5x_3 = 2 \\ 2x_1 - 3x_2 + 4x_3 = 4 \\ 4x_1 + x_2 - 6x_3 = 8 \\ 3x_1 - x_2 - x_3 = 6 \end{array}$$

el sistema resulta de 4 ecuaciones con 3 incógnitas, tiene el mismo conjunto de soluciones, del sistema original. El cual es

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \\ 0 \end{bmatrix} + k \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$$

Y el sistema de dos ecuaciones con 3 incógnitas, dado por las primeras dos ecuaciones en el sistema anterior, es decir

$$x_1 + 2x_2 - 5x_3 = 2$$

$$2x_1 - 3x_2 + 4x_3 = 4$$

tiene el mismo conjunto de soluciones como, el sistema original anteriormente descrito en (1)

La moral de estos ejemplo es que en general, no es posible decir si un conjunto de ecuaciones, no tiene solución, una solución única, ó una infinidad de soluciones, meramente del conocimiento del número de ecuaciones y del número de incógnitas.- es decir, m y n. Por tanto 10 ecuaciones con 2 incógnitas pueden tener una solución única, y 2 ecuaciones con 10 incógnitas pueden ser inconsistentes.

Consideremos el siguiente ejemplo.

$$x_1 + 2x_2 - x_3 = 2$$

$$2x_1 + 4x_2 + x_3 = 7$$

$$3x_1 + 6x_2 - 2x_3 = 7$$

Usando la primera ecuación para eliminar x de las 2 restantes obtenemos.

$$\begin{array}{rcl} x_1 + 2x_2 - x_3 & = & 2 \\ & & 3x_3 = 3 \\ & & x_3 = 1 \end{array} \quad \begin{bmatrix} 1 & 2 & -1 & 2 \\ 0 & 0 & 3 & 3 \\ 0 & 0 & 1 & 1 \end{bmatrix}$$

La incógnita x no aparece en la segunda y tercera ecuaciones.

Hemos reducido el coeficiente de x en la segunda ecuación a la unidad por dividir ésta ecuación por 3.

Podemos, entonces, omitir la tercera ecuación, porque es simplemente un múltiplo de la segunda. (en terminos de la matriz, los elementos (2,2) y (3,2) son cero. Y consideremos la tercera columna, reduce los elementos (2,3) a la unidad por multiplicar el renglón siguiente por 1/3, y restando el renglón resultante del último renglón, tal que el último renglón es reducido a un renglón de ceros).

$$\begin{array}{rcl} x_1 + 2x_2 - x_3 & = & 2 \\ & & x_3 = 3 \end{array} \quad \begin{bmatrix} 1 & 2 & -1 & 2 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (2')$$

La solución general del conjunto original de ecuaciones está dado por

$$x_3 = 1, \quad x_2 = p, \quad x_1 = 3 - 2p \quad (2)$$

Donde tenemos el conjunto de $x_2 = p$, una constante arbitraria. Dado que de conjuntar x_2 igual a una constante x_2

arbitraria, podríamos tener el conjunto de x igual a una constante arbitraria, y entonces la solución general es

$$\begin{matrix} x_3 \\ x_1 \\ x_2 \end{matrix} = \begin{matrix} 1 \\ q \\ 1/2(3 - q) \end{matrix} \quad (3)$$

Es claro que las soluciones anteriores son esencialmente las mismas dado que, si asignamos cualquier valor a p en la primera solución general (2), la misma solución será dada al conjuntar $q = 3 - 2p$ en (3). Sin embargo, en casos más complicados la equivalencia de formas diferentes de soluciones puede no ser obvia. Debería notarse que en el ejemplo de arriba no es posible asignar un valor arbitrario a x_3 , dado que una ecuación es $x_3 = 1$

3

El procedimiento dado en la ecuación (2') corresponde a lo que hemos llamado previamente "reducción a la forma triangular", con modificaciones apropiadas para el hecho de que los elementos cero en matrices sucesivas significa que no podemos tener una forma triangular estricta. Para ver la estructura de la solución de (2') más claramente, es conveniente usar la sustitución hacia atrás en las ecuaciones (en la matriz; usando el procedimiento análogo de Gauss - Jordann para reducir el elemento (1,3) a cero. Esto da

$$\begin{matrix} x_1 + 2x_2 - x_3 = 3 \\ x_3 = 1 \end{matrix} \quad \begin{bmatrix} 1 & 2 & 0 & 3 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Esto nos dice inmediatamente que a cualquiera de x_1 ó x_2 (pero no ambos) puede asignarsele valores arbitrarios, pero a x_3 no se le puede dar un valor arbitrario.

3

C A P I T U L O

V. I.

INEXACTITUD INHERENTE
EN SISTEMAS LINEALES
(SOLUCIONES)

Dada una matriz no singular A y un vector b distinto de cero, sea \bar{x} la solución del sistema.

$$A x = b$$

Supongamos que A y b son sujetos de incertidumbre. Cual es la incertidumbre resultante en la solución x ? Para cualquier vector columna y de orden n definimos $\|y\|$ como la longitud euclidiana de y ;

$$\|y\| = \sqrt{y_1^2 + y_2^2 + \dots + y_n^2}$$

Para una matriz cuadrada A de n x n, definimos la norma espectral de $\|A\|$, como.

$$\|A\| = \max_{\|x\|=1} \|Ax\|$$

Estas funciones dan por lo general, la medida ó tamaño de vectores y matrices respectivamente.

Para una matriz no singular A. Definimos la condición de A como $\text{cond}(A)$, por la relación.

$$\text{cond}(A) = \|A\| \cdot \|A^{-1}\|$$

El concepto de condición de una matriz, intenta medir la sensibilidad ó la "vulnerabilidad" de una solución : si cambiamos ligeramente A y b, cuán grande es el efecto en

$x = A^{-1} b$? antes de comenzar con esta cuestión queremos señalar un obstaculo (que podemos salvar fácilmente). Tiene que haber una manera de medir el cambio δA y de calcular el tamaño de A misma.

Ya definimos la longitud de un vector; necesitamos ahora la norma de una matriz. Entonces el número de condición, y con él la sensibilidad de A, se podrán obtener fácilmente

directamente de las normas A y A^{-1}

Uno de los principales usos del concepto de condición yace en respuesta, por ahora retornamos a la cuestión hecha al principio de este capítulo, Qué cambio produce en la solución x un cambio pequeño en b ó en A ?

Comenzaremos con un cambio del lado derecho, de b a $b + \delta b$. Este error puede provenir de los datos experimentales ó de redondeo; podemos suponer que δb es pequeño pero que su dirección está fuera de nuestro control. La solución correspondiente cambia de x a $x + \delta x$:

$$A (x + \delta x) = b + \delta b ,$$

$$x + \delta x = A^{-1} b + A^{-1} \delta b$$

$$\delta x = A^{-1} \delta b$$

$$\| \delta x \| \leq \| A^{-1} \| \cdot \| \delta b \|$$

Este es un caso particularmente sencillo; consideremos todas las perturbaciones δb y calculamos la perturbación.

resultante $\delta x = A^{-1} \delta b$. Habrá grandes cambios en la solución

cuando A^{-1} sea grande (A es casi singular) y será particularmente grande cuando δb vaya en la dirección donde

sea amplificado por A^{-1} .

y puesto que $A x = b$, tenemos ,

$$\| b \| = \| A x \| \leq \| A \| \| x \|$$

dividiendo

$$\| \delta x \| \leq \| A^{-1} \| \cdot \| \delta b \|$$

por

$$\| b \| = \| A x \| \leq \| A \| \cdot \| x \|$$

tenemos

$$\frac{\| \delta x \|}{\| x \|} \leq \| A^{-1} \| \cdot \| A \| \cdot \frac{\| \delta b \|}{\| b \|}$$

ó

$$\frac{\| \delta x \|}{\| x \|} \leq \text{cond}(A) \frac{\| \delta b \|}{\| b \|}$$

La última desigualdad muestra que la incertidumbre relativa a x es limitada por la $\text{cond}(A)$. El límite en esta expresión es alcanzable, para cualquier A no singular y b distinto de cero. Esto es fácil de ver, si llevamos a cabo un cambio de coordenadas en el cual A tome la forma diagonal.

Pero hay un serio inconveniente para medir así la sensibilidad. Supongamos que multiplicamos todas las entradas de A por 1,000; entonces γ (el error del tamaño

γ se amplifica por el factor $1/\gamma$ que no es más que el

valor propio mayor de A . La simplificación es mayor cuando γ está cerca de cero, así que las matrices casi singulares

son las más sensibles) se multiplicará por 1,000 y la matriz será mucho menos singular. Pero se supone que estamos jugando limpio; no es posible que mediante un cambio de escala podamos convertir la matriz mal dispuesta en una bien colocada. Es cierto que δx será menor en un

factor de 1,000, pero también lo será la solución $x = A^{-1} b$ y el error relativo $\| \delta x \| / \| x \|$ será el mismo. El factor x en el denominador normaliza el problema contra pequeños cambios de escala. Al mismo tiempo hay una normalización corespondiente para δb ; nuestro problema es comparar el cambio relativo $\| \delta b \| / \| b \|$ con el error relativo $\| \delta x \| / \| x \|$.

Lo peor sucede cuando el número δx es grande (las perturbaciones están en la dirección del vector propio x)

y cuando el denominador x es pequeño. La solución sin perturbar x debería ser la más pequeña posible comparada con la b sin perturbar.

Ahora bien llevando a cabo un cambio de coordenadas, como una transformación lineal, a toma vectores x en vector b . Un teorema de fundamental importancia pero poco conocido establece que, por un cierto cambio ortogonal de coordenadas en el espacio x y por otro cambio ortogonal de coordenadas en el espacio de b , la matriz A puede ser puesta en la forma diagonal.

$$A = \begin{bmatrix} \mu_1 & & & 0 \\ & \mu_2 & & \\ & & \ddots & \\ 0 & & & \mu_n \end{bmatrix}$$

Aquí los números positivos $\mu_1, \mu_2, \dots, \mu_n$, son llamados los valores singulares de A además.

$$\|A\| = \mu_1 \quad \|A^{-1}\| = \mu_n^{-1}$$

Finalmente, las transformaciones ortogonales no cambian las normas de x y b, tenemos.

$$A^{-1} = \begin{bmatrix} -1 & & & \\ \mu_1 & & & \\ & -1 & & \\ & \mu_2 & & \\ & & \ddots & \\ & & & -1 \\ & & & \mu_n \end{bmatrix}$$

$$\text{si } b = \begin{bmatrix} 1 \\ 0 \\ \cdot \\ 0 \end{bmatrix}, \quad \text{y} \quad \delta b = \begin{bmatrix} 0 \\ 0 \\ \cdot \\ \epsilon \end{bmatrix}$$

$$x = A^{-1} b = \begin{bmatrix} -1 \\ \mu \\ 0 \\ \cdot \\ 0 \end{bmatrix}, \quad y \quad x = A^{-1} \delta b = \begin{bmatrix} 0 \\ \cdot \\ 0 \\ \epsilon \mu \\ n \end{bmatrix}$$

Para estos vectores

$$\frac{\| \delta x \|}{\| x \|} = \frac{\mu_1}{\mu_n} = \frac{\| \delta b \|}{\| b \|} \frac{\mu_1}{\mu_n} = \text{cond}(A) \frac{\| \delta b \|}{\| b \|}$$

Y esta última expresión muestra que

$$\frac{\| \delta x \|}{\| x \|} \leq \text{cond}(A) \frac{\| \delta b \|}{\| b \|}$$

es una igualdad, como prometimos probar aunque.

$$\frac{\| \delta x \|}{\| x \|} = \frac{\mu_1}{\mu_n} = \frac{\| \delta b \|}{\| b \|} \frac{\mu_1}{\mu_n} = \text{cond}(A) \frac{\| \delta b \|}{\| b \|}$$

es solo una igualdad exacta, bajo condiciones excepcionales es mejor cerrar la igualdad y por consiguiente damos por

supuesto la igualdad aproximada. Si $\text{cond}(A) = 10^p$, y si b es correctamente conocida solo para 10 decimales, ahora x puede ser conocida para 10^{-p} decimales, entonces p puede

tomar cualquier rango de 0 a x

El solo esperar tener alguna significancia para x , el sistema de computadora de 10 decimales, es que

$$\text{cond}(A)^{-10} < 1/2$$

En una computadora de base β con t dígitos significativos necesitamos

$$\text{cond}(A)\beta^{-t} < 1/2$$

Para tener alguna significancia para una solución. Recordamos que todas las exposiciones en esta sección son independientes de cualquier método de resolver un sistema $Ax = b$. Son exposiciones acerca de los errores en x , los cuales son inherentes a la incertidumbre en los datos.

Si A está sujeta a un cambio δA y b es conocido exactamente, entonces una desigualdad análoga a,

$$\frac{\| \delta x \|}{\| x \|} \leq \text{cond}(A) \frac{\| \delta b \|}{\| b \|}$$

es la siguiente

$$\frac{\| \delta x \|}{\| x + \delta x \|} = \text{cond}(A) \frac{\| \delta b \|}{\| A \|}$$

Si x es pequeña comparado con $\| x + \delta x \|$, entonces podemos, considerar sin contratiempo alguno, el lado izquierdo de ésta última desigualdad, como un error relativo a x .

Ejemplo.

Los valores propios de

$$A = \begin{bmatrix} 1 & 1 \\ 1 & 1.0001 \end{bmatrix}$$

-4

Son aproximadamente $\gamma_1 = 10^4 / 2$ y $\gamma_2 = 2$. Por lo tanto, su número de condición es alrededor de $c = 4 \cdot 10^4$ y debemos esperar un cambio violento en la solución debido a cambios nada extraños en los datos. En la solución.

$$u + v = 2 \quad u + v = 2$$

$$u + 1.0001 v = 2 ; \quad u + 1.0001 v = 2.0001.$$

Los lados derechos se alteran solamente en $\delta b = 10^{-4}$. Al mismo tiempo, las soluciones cambiaron de $u = 2, v = 0$ a $u = v = 1$. Este es un error relativo de

$$\frac{\| \delta x \|}{\| x \|} = \frac{\| (-1, 1) \|}{\| (2, 0) \|} = \frac{2}{2}$$

que es igual a

$$2 \cdot 10^4 \frac{\| \delta b \|}{\| b \|}$$

Aun sin seleccionar alguna perturbación en especial (nuestra x y δb forma un ángulo de 45° en el peor de los casos, lo cual equivale a la pérdida del 2 en nuestro factor $2 \cdot 10^4$ comparado con ; a posibilidad extrema

$c = 4 \cdot 10^4$) tuvimos un cambio tremendo en la solución.

Notese que el número de condición c , no se afecta directamente por el tamaño de la matriz; si $A = I$ ó aún si $A = I / 10$ el número de condición es $c = \gamma_{\text{máx}} / \gamma_{\text{mín}} = 1$

En comparación, el determinante es un criterio malísimo para probar la mala colocación de una matriz, ya que depende no sólo de la escala considerada, sino también del orden n ; si $A = I/10$ entonces el determinante de A es 10^{-n}

De hecho, esta matriz " casi singular " está lo mejor colocada posible.

C A P I T U L O

V I I .

PRECISION REALIZABLE
EN
ELIMINACION GAUSSIANA

En este capítulo daremos por hecho que el lector conoce la eliminación gaussiana, que es un método para resolver sistemas de ecuaciones lineales. La principal decisión estratégica que contempla el diseñador del algoritmo es la elección de un elemento pivote único para cada uno de los $n-1$ pasos, en el cual una variable es eliminada de las ecuaciones restantes.

Hay dos tipos de estrategias importantes :

- i) Pivoteo completo.- En el cual cada paso uno selecciona como pivote algún elemento a_{ij} de valor absoluto máximo entre todos los elementos restantes de la matriz.
- ii) Pivoteo parcial.- En el cual a cada paso uno selecciona como pivote algún elemento a_{ij} de valor absoluto máximo entre la primera columna de los elementos restantes de la matriz.

Así en el primer paso el pivoteo completo examinaría la sola matriz A , para un elemento maximal en valor absoluto, mientras que el pivoteo parcial examinaría solo la primera columna.

Algunas clases especiales de matrices permiten la eliminación para proceder sucesivamente sin algún pivote - por ejemplo, las matrices simétricas positivas definidas - pero generalmente el examen de pivote es esencial para garantizar el éxito. El siguiente ejemplo simple ilustra el desastre que produce al no examinar un pivote.

Consideremos una máquina de tres dígitos decimales en punto flotante.

El sistema es,

$$.0001 x + 1.000 y = 1.00$$

$$1.00 x + 1.00 y = 2.00$$

La verdadera solución, redondeando a cinco decimales es,

$$x = 1.00010 \quad y = .99990$$

Si aceptamos el elemento .0001 como pivote, la eliminación de x en la segunda ecuación da la ecuación.

$$- 10000 y = - 10000$$

Resolviendo por sustitución hacia atrás, encontramos que $y = 1.00$, mientras que $x = 0.00$, un claro desastre. Por otra parte el pivoteo parcial seleccionaría el elemento pivote $a_{2,1} = 1.00$. La eliminación de x de la primera ecuación da la ecuación.

$$1.00 y = 1.00$$

Resolviendo a la inversa encontramos que

$$y = 1.00$$

y entonces

$$x = 1.00 \quad \text{con éxito obvio.}$$

Ahora diremos que estamos tratando con una computadora de punto flotante de base 2 y t dígitos. Más bien que discutir la solución de un sistema lineal consideremos la

computación de la inversa A^{-1} de una matriz dada. Deseamos establecer los límites de error de redondeo que habian sido probados por la eliminación gaussiana.

Asumiendo una estrategia de pivoteo completo y que la matriz A este razonablemente escalonada, al principio y en todos los pasos intermedios (ver capítulo de Escalonamiento de Matrices) entonces, si todo $a_{ij} > 1$, un cierto algoritmo gaussiano da una matriz X tal que.

$$\frac{\| X - A^{-1} \|}{\| A^{-1} \|} \leq (0.8)^{-t} \frac{7}{2} g(n) \| A^{-1} \|$$

Aquí $g(n)$ es el máximo de todos los elementos de las matrices sucesivas encontradas durante la eliminación. Para expresar el resultado anterior, en una forma para poder compararla con algunos resultados del capítulo

anterior notemos que $1 < A < n$, tal que esperamos que

$$A = n^{\frac{1}{2}}$$

Entonces tenemos.

$$\frac{\|X - A^{-1}\|}{\|A^{-1}\|} = n^{\frac{3-t}{2}} \text{cond}(A) g(n)$$

Qué clase de límite podemos dar para $g(n)$?.

Esto arroja una pregunta abierta, y lo que se sabe que resulta mejor es aproximadamente.

$$g(n) \leq 1.8 n^{(1/2) \log n}$$

Por otro lado para todas las matrices reales que han sido examinadas, siempre se ha observado que.

$$g(n) \leq n$$

el último límite es alcanzado para una n grande ilimitadamente por matrices relacionadas a las matrices de Hadamard.

Para las demás matrices alguna vez se observa que.

$$g(n) \leq 8$$

Y se han encontrado matrices complejas A de una n ilimitadamente grande para lo cual.

$$g(n) = 3.1 n$$

Debería ser más deseable tener un buen límite para $g(n)$, tal que.

$$\frac{\|X - A^{-1}\|}{\|A^{-1}\|} (0.8)^{\frac{-t}{2}} n^{\frac{7}{2}} g(n) A^{-1}$$

Podría ser cambiado en un buen límite de error a - priori

para el cálculo de A^{-1} . Naturalmente para cualquier matriz A en particular, $g(n)$ es fácilmente observado en el curso de la eliminación, tal que en algún evento ésta relación

$$\frac{\|X - A^{-1}\|}{\|A^{-1}\|} \leq (0.8)^{n-1} 2^{7/2} g(n) \|A^{-1}\|^{-1}$$

se transforma en un límite de error a - posteriori.

Sin embargo todavía límites de error mejores pueden estar dados a - posteriori, como será mostrado en el capítulo "Más Soluciones Exactas".

El factor $2^{7/2}$ es esencialmente el nivel de insertidumbre inherente a los datos, y sería igualado a

$$\frac{\|SA\|}{\|A\|}$$

Entonces el límite en

$$\frac{\|X - A^{-1}\|}{\|A^{-1}\|} = n^{-3} 2^{-t} \text{cond}(A) g(n)$$

es más grande que en

$$\frac{\|\delta x\|}{\|x + \delta x\|} \leq \text{cond}(A) \frac{\|\delta A\|}{\|A\|}$$

Por el factor n^{-3} . Tomando en cuenta el resultado empirico

$$g(n) \leq 8$$

que para más matrices reales, entonces interpretamos

$$\frac{\|X - A^{-1}\|}{\|A^{-1}\|} \leq \frac{3}{n^2} \text{cond}(A) g(n)^{-t}$$

como diciendo que la matriz X calculada, generalmente difiere de la verdadera inversa A^{-1} , en terminos relativos por no más de $\frac{3}{n^2}$ veces el error inherente en el problema.

Así la eliminación gaussiana simple es razonablemente buena al tomar el límite, con el error de redondeo bajo control, para valores modestos de n . Mejores resultados pueden obtenerse con algunas ideas que mencionaremos en el siguiente capítulo.

El límite dado a.

$$\frac{\|X - A^{-1}\|}{\|A^{-1}\|} \leq (0.8)^{-t} \frac{7}{2n} g(n) A^{-1}$$

fué

$$\frac{\|X - A^{-1}\|}{\|A^{-1}\|} \leq (5.3 + 14.6 A^{-1})^2 n^{-t} \|A^{-1}\|^2 \approx 15n^2 (\text{cond}(A))^2$$

El factor $(\text{cond}(A))^2$ surge de resolver

$$A^t A^t x = A^t b, \text{ más bien que } Ax = b$$

La prueba de

$$\frac{\|X - A^{-1}\|}{\|A^{-1}\|} \leq (5.3 + 14.6 A^{-t}) \frac{\|A^{-1}\|^2}{15n^{2-t} (\text{cond}(A))^2}$$

es mucho más difícil y tediosa que la prueba de

$$\frac{\|X - A^{-1}\|}{\|A^{-1}\|} \leq (0.8)^{t/2} \frac{1}{n^{7/2}} \frac{1}{g(n) A^{-1}}$$

C A P I T U L O

V I I I .

M A S S O L U C I O N E S
E X A C T A S

Supongamos que la matriz A está dada como un dato en precisión simple y que deseamos conseguir soluciones que garanticen ser más precisas, que los límites,

$$\frac{\| \| X - TA \| \|^{-1}}{\| \| A^{-1} \| \|} = (5.3 + 14.6 \frac{2}{A}) \frac{-t}{n} \frac{2}{\| \| A^{-1} \| \|}^2$$

$$\approx 15n \frac{2}{2} \frac{-t}{2} (\text{cond}(A))^2$$

donde el condicional $\frac{2}{(\text{cond}(A))}$

surge de resolver $A^t Ax = A^t b$

más bien que de $Ax = b$

Ahora la cuestión es, Cómo procederíamos ?. La elección más obvia es ejecutar los cálculos en doble precisión entonces t es sustituido por 2t en los límites de arriba y hasta $\frac{-2t}{-t}$

$\frac{2}{2}$ es mucho más pequeño que $\frac{2}{2}$, con esto ganamos muchos órdenes de magnitud en precisión. El costo en tiempo de computación varía entre diferentes máquinas, pero es sólo un factor de cuatro en la IBM.

El costo del almacenaje es grande, hasta debemos duplicar el almacenaje reservado para la matriz desarrollada.

Quando el costo del tiempo y almacenaje son muy grandes para justificar el completar la doble precisión, es posible aprovecharlo sustancialmente, por un uso más limitado de doble precisión. Muchas de las operaciones en la eliminación gaussiana pueden ser expresadas como productos internos de vectores de números de precisión simple. En muchas máquinas es posible acumular como un producto interno en doble precisión y entonces redondearlo para precisión simple, antes de almacenar. El resultado de esta acumulación es reducir al máximo el error de redondeo de un producto interno por un factor de n.

El efecto total arroja el reducir el límite del error de redondeo por un factor $n^{5/2}$. Así, en lugar del resultado,

$$\frac{\left\| X - A^{-1} \right\|}{\left\| A^{-1} \right\|} = (0.8)^{2n} \frac{t^{7/2}}{g(n) A^{-1}}$$

Una eliminación con pivoteo y acumulación produce una inversa aproximada X tal que

$$\frac{\left\| X - A^{-1} \right\|}{\left\| A^{-1} \right\|} = 3.3 n^{2-t} A^{-1}$$

Bajo ciertas hipótesis adicionales, el provecho del factor $n^{5/2}$

es muy sustancial, aunque la experiencia muestra que los errores actuales en computación como precisión simple son usualmente un tanto más pequeños que los límites. Una desventaja teórica de la estrategia de pivoteo completo es que no se mezcla bien con la acumulación de productos internos. Cuando los productos son acumulados, casi siempre se usa una estrategia de pivoteo parcial y acepta la posibilidad teórica de que estos pivoteos puedan crecer mucho.

Un tercero y el aprovechamiento más afortunado para incrementar la exactitud en soluciones de sistemas lineales densos almacenados, es el tan llamado método de refinamiento iterativo. Describámoslo brevemente.

Cuando un sistema lineal $Ax = b$ ha sido resuelto usando una factorización estable de A tal como el dado en la eliminación gaussiana con pivoteo ó una triangularización ortogonal entonces una solución más precisa puede encontrarse usando esta factorización dando el residual correspondiente a la solución aproximada, es calculada usando gran precisión. Este proceso puede ser usado de una manera interativa para dar una solución mejorada progresivamente. Digamos que el método directo es tal, que

corresponde a algún c del lado derecho que da la solución exacta de algunos

$$(A + E_c) x = c \quad \text{donde} \quad \left\| \begin{matrix} E \\ c \end{matrix} \right\| \leq \left\| A \right\|$$

Donde E_c es una función de c , pero el límite en E_c es independiente de c . Entonces si el residual pudiera ser exactamente un proceso iterativo, podría ser definido por.

$$x_0 = 0, \quad r_k = b - Ax_k, \quad (A + E_k) s_{k+1} = r_k$$

$$x_{k+1} = x_k + s_k$$

Aquí damos por hecho que, no solamente cada una de las r_k está computada exactamente, pero también que este residual puede ser usado por el lado derecho y que la corrección calculada s_k puede ser agregada a la aproximación x_k sin error. Los únicos errores son aquellos en resolver los sistemas lineales.

si escribimos

$$e_k = x - x_k$$

entonces de las ecuaciones

$$x_0 = 0$$

$$r_k = b - Ax_k$$

$$(A + E_k) s_{k+1} = r_k$$

y

$$e_k = x - x_k$$

Resulta,

$$e_{k+1} = (I - (A + E_k)^{-1} A) e_k = P_k e_k$$

Donde

$$P_k = I - (A + E_k)^{-1} A$$

y,

$$\|P_k\| \leq \epsilon_k / (1 - \epsilon_k)$$

De aquí si

$$\epsilon_k / (1 - \epsilon_k) \leq \beta^{-p}$$

La solución gana al final p dígitos de precisión para la iteración.

En la práctica el residual r_k no puede ser calculado

exactamente, pero si se puede acumular productos internos en doble precisión, el error puede ser mantenido a un nivel muy bajo. El residual entonces, puede ser redondeado en precisión simple y usado como el próximo lado derecho y δ_k puede ser obtenido usando la factorización de A .

Esta precisión simple δ_{k+1} , es entonces sumada a x_k y el vector resultante redondeado a precisión simple de la x_{k+1} calculada. Dado que cada x_k es un vector en precisión simple, la precisión final obtenible está obviamente limitada a precisión simple.

Por este método, si la matriz A no esta demasiado mal condicionada, en la práctica se consiguen las soluciones las cuales son las aproximaciones correctas redondeadas a la verdadera respuesta. Describiendo este desarrollo en otras palabras tenemos : Supongamos que por la eliminación gaussiana se ha encontrado una primera solución aproximada x_0 del sistema lineal $Ax = b$. El próximo paso es formar el

vector residual $r = b - Ax_0$. Si x_0 era la

solución exacta del sistema podríamos tener $r = 0$, el vector nulo.

Si no resolvemos el nuevo sistema lineal $Ay = r$, para obtener un vector y , sea

$$x_1 = x_0 + y_1.$$

El proceso es repetido iterativamente, es decir para $k = 0, 1, 1, \dots$, formamos el residual.

$$r_k = b - Ax_k$$

resolver el sistema $Ay = r_k$, para obtener un vector y_{k+1} , y entonces formar

$$x_{k+1} = x_k + y_{k+1}$$

bajo hipótesis apropiadas para ser especificado bajo la serie x_k que converge a la verdadera solución A^{-1} del sistema $Ax = b$.

Varias cuestiones se necesitan para esclarecer esta algoritmo.

PRIMERO.- Parece involucrar una gran cantidad de trabajo para resolver sistemas de la forma $Ay = r_k$,

para muchos valores de k . De hecho esto no es suficiente.

La eliminación gaussiana para resolver un sistema $Ax = b$, involucra tres pasos principales;

- i) Triangularización de la matriz A por transformaciones elementales de renglón.
- ii) Aplicaciones de las mismas transformaciones por renglón al lado derecho b ;

iii) La solución del sistema triangular por sustitución hacia atrás

3/3

El producir el paso i) requiere aproximadamente de n adiciones y multiplicaciones, pero los pasos ii) y iii)

2

requieren solo de n multiplicaciones y adiciones aproximadamente. El paso i) solo es necesario hacerlo una vez para todos los sistemas $Ay = r$. Si los

k

multiplicadores definen que las transformaciones por renglón están listas, los pasos ii) y iii), entonces pueden ser realizados rápidamente para cada nuevo sistema $Ay = r$

k

en turno. Como un resultado, se ve que una serie suficientemente larga de vectores x puede por lo general,

k

ser calculada en más o menos 20% más tiempo que el del cálculo en la primera solución x .

0

Es absolutamente indispensable que cada vector residual r

k

sea calculado, para gran precisión. Esto es un hecho normal para la acumulación de productos internos en doble precisión, seguido por el redondeo de la respuesta en precisión simple en la forma de punto flotante. Si r es

k

calculado con un producto interior en precisión simple, tendrá errores de redondeo de varias unidades en los mínimos dígitos significativos de x . Entonces la desigualdad.

k

$$\frac{\| \delta x \|}{\| x \|} \leq \text{cond}(A) \frac{\| \delta b \|}{\| b \|}$$

será una igualdad aproximada en la práctica, digamos que x

k

será fuerte en algunos casos, cuando el $\text{cond}(A)$ sea el mínimo dígito significativo. Puesto que el $\text{cond}(A)$ puede

4 5

bien ser 10^4 ó 10^5 . La precisión resultante en x es

k

muy bajo y, en efecto, x_k es también tan precisa como cualquier x_0 .

El siguiente teorema de las bases del método expuesto anteriormente del refinamiento iterativo.

TEOREMA.

Sea la matriz A , la cual tiene la propiedad de

$$(0.8) \quad \left\| \begin{matrix} 2^{-t} & 7/2 \\ n & g(n) \end{matrix} \right\| \left\| A^{-1} \right\| < 1/2$$

Sea el algoritmo de arriba, llevándose a cabo, con cada sistema $Ay = r_k$, siendo resuelto en simple precisión en aritmética de punto flotante con base 2, pero con computaciones de,

$$r_k = b - Ax_k$$

$$y, \quad x_{k+1} = x_k + y_{k+1}$$

ejecutando sin error de redondeo, entonces,

$$\left\| x_k - A^{-1} b \right\| \longrightarrow 0$$

cuando $k \longrightarrow \infty$

Si la solución de los sistemas $Ay = r_k$, están hechos con acumulaciones de producto internos en doble precisión, entonces el lado izquierdo de

$$(0.8) \quad \left\| \begin{matrix} 2^{-t} & 7/2 \\ n & g(n) \end{matrix} \right\| \left\| A^{-1} \right\| < 1/2$$

puede ser reemplazado por el lado derecho de

$$\frac{\left\| \begin{matrix} X - A^{-1} \\ \end{matrix} \right\|}{\left\| \begin{matrix} -1 \\ A \end{matrix} \right\|} \leq 3.3 n^{2-t} \left\| \begin{matrix} -1 \\ A \end{matrix} \right\| g(n)$$

En la práctica, por supuesto, r es computado por una
 x_{k+1} acumulación de productos internos en doble precisión y y_{k+1} es computado como la suma en punto flotante x_k y y_k .

Como un resultado la serie x_k no converge a $A^{-1}b$ en el sentido matemático.

En lugar de ello, x_k es observada para convertirse constantemente a un valor, el cual es normalmente el correcto redondeo en precisión simple. (de la aproximación $A^{-1}b$). En el uso actual del refinamiento iterativo, generalmente no se conoce el avance, si la hipótesis.

$$(0.8)^{2-t} n^{7/2} g(n) \left\| \begin{matrix} -1 \\ A \end{matrix} \right\| < 1/2$$

es ó no sastisfecha y puede no ser concluida después de uno u otro. En la práctica es por consiguiente normal contar con el siguiente resultado.

Sea el algoritmo de arriba, ejecutado en cada sistema $Ay = r_k$, siendo resuelto por la misma versión de eliminación gaussiana, con cada r_k computado por una acumulación de productos internos en doble precisión, y con x_{k+1} computado con la suma de punto flotante de x_k y y_k .

Sí para $k \geq k_0$, todos los vectores x_k son iguales para

algún vector x^* , en doble precisión, entonces x^* es la aproximación correctamente redondeada en precisión simple a

-1
A b.

Este resultado, no puede ser probado, y realmente la mayoría de las computadoras ocuparían todo su tiempo en la aplicación del algoritmo de arriba, en algún ejemplo práctico.

Normalmente cuando el $\text{cond}(A)$, consigue acercarse a 2^t los vectores x obviamente divergen, entonces no hay remedio

excepto para incrementar la precisión con lo cual la eliminación es llevada a cabo, a menos que escalar A ayudara, el valor de k es generalmente 3 ó 4.

C A P I T U L O

IX.

E S C A L O N A M I E N T O
O P T I M O D E
M A T R I C E S

Un aspecto que fue mencionado en el capítulo referente a Sistemas de Ecuaciones Lineales, fué el escalar una matriz A, antes de resolver un sistema $Ax = b$, los términos alternativos para escalar una matriz, están precondicionados y equilibrados.

De nuevo veamos el condicional. Una definición de condición de una matriz, de manera más natural, es el radio de su mínimo límite superior y de su máximo límite inferior, con respecto a un par de normas.

Algunas veces, los renglones ó columnas de una matriz son un poco grandes para algunos factores; en cualquier caso escalar una matriz es trivial, computacionalmente hablando, por consiguiente, el problema de minimizar la condición por escalinamiento tiene varias aplicaciones y deseamos determinar el infimo ó mínimo para derivar los límites inferiores y superiores. Para esto queremos conocer la dirección del escalonamiento al mínimo ó al último que se le acerque y queremos mostrar condiciones para que una matriz esté optimamente escalonada.

Así el problema puede ser completamente resuelto por la condición subordinada a un par de normas máximas; y así

-1

cada una de las matrices A y A^{-1} , tienen una distribución en bloque cuadrada. Para derivar los límites superiores, se debe determinar un par de vectores mutuamente duales.

Sean B_1 y B_2 dos espacios vectoriales normados, con

$$\|x\|_i \quad \text{para } x \in B_i$$

Para un mapeo A de B_1 en B_2 , el mínimo límite superior (mín. lím. sup.)-el cual es una norma-y el máximo límite inferior (máx. lím. inf.) están definidos por.

$$\text{mín.lím.sup. (A)} = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|}$$

$$\text{máx.lím.inf. (A)} = \inf_{x \neq 0} \frac{\| Ax \|}{\| x \|}$$

Entonces el condicional subordinado a las normas $\| Ax \|$ y $\| x \|$ se puede definir por:

$$\text{cond}(A) = \text{mím. lím. sup (A)} / \text{máx. lím. inf (A)}$$

si tenemos que

$$\text{máx. lím. inf (A)} > 0$$

y

$$\text{máx. lím. inf (A)} = 0$$

podemos decir que $\text{cond}(A) = \infty$.

El condicional puede ser totalmente finito para una A rectangular, sin embargo.

$$\text{cond}(A) = \text{mím. lím. sup (A)} \text{ mín. lím. sup (A)}^{-1}$$

-1

si A existe.

En lo sucesivo diremos que A es un mapeo regular. Mencionaremos algunas propiedades básicas de la función $\text{cond}(A)$ las cuales son válidas independientemente de las normas fundamentales.

$$1) \text{ cond}(A) \geq 1$$

$$2) \text{ cond}(A) = 1, \text{ implica que } \| Ax \| = \delta \| x \| \text{ para una cierta } \delta > 0, \text{ independiente de } x.$$

el $\text{cond}(A) = 1$, se ve inmediatamente de la definición

$$\text{cond}(A) = \text{mím. lím. sup (A)} / \text{máx. lím. inf (A)}$$

dado que el infimo no excede el supremo y significa igualmente que.

$$\text{mím. lím. sup (A)} \text{ máx. lím. inf (A)} = \delta$$

ó

$$\| Ax \| \leq \delta \| x \| \text{ y } \| Ax \| = \delta \| x \|$$

$$3) \text{cond}(A^{-1}) = \text{cond}(A)$$

$$4) \text{cond}(A B) \leq \text{cond}(A) \text{cond}(B)$$

Se ve de la definición

$$\text{cond}(A) = \frac{\text{mín. lím. sup.}(A)}{\text{mín. lím. sup.}(A^{-1})}$$

sí A^{-1} existe que,

$$\text{cond}(A B) \leq \text{cond}(A) \text{cond}(B)$$

se cumple dado que las normas mín. lím. sup. son multiplicativas.

La importancia de los números condición está basada en el hecho de que dan los mejores límites posibles en las desigualdades.

$$\frac{1}{\text{cond}(A)} \frac{\|x\|}{\|y\|} \leq \frac{\|Ax\|}{\|Ay\|} \leq \text{cond}(A) \frac{\|x\|}{\|y\|}$$

Por consiguiente, las condiciones aparecen, por ejemplo en la estimación de eigen valores para matrices diagonalizables, donde P , puede ser un sistema de eigenvectores diagonales para A , entonces,

$$P^{-1} A P = \text{diag}(\gamma_i) \text{ y por consiguiente}$$

$$\text{mín. lím. sup.}(A) = \text{cond}(A) \max_i |\gamma_i|$$

Para cualquier norma absoluta $\|x\|$, donde:

$$\text{mín. lím. sup.}(\text{diag}(\gamma_i)) = \max_i |\gamma_i|$$

finalmente, se sabe que el análisis de error de métodos de eliminación para matrices positivas definidas, que la condición (por ejemplo el de uno subordinado a la norma

euclidiana) está caracterizada por la amplificación posible de errores de redondeo simple.

Si A es la matriz de un sistema de ecuaciones lineales $Ax = b$, el escalonar los renglones de un sistema, suma a una transformación de A por una matriz diagonal D_1 por la izquierda y el escalonar las sumas no conocidas a la transformación de A por una matriz diagonal D_2 por la derecha. Si A es hermitiana, ésta propiedad se preserva si $D_1 = D_2$. Dado que las transformaciones diagonales de una matriz son modificaciones triviales y dado que de otro modo la calidad de las computaciones numéricas son generalmente mejoradas si la condición de la matriz concerniente es disminuida. Ha llavado a investigar los cuatro problemas de minimización, para conseguir el pré-escalonamiento de una matriz.

$$i) \inf \text{ cond } D_1 D_2 (D_1 A D_2)$$

$$ii) \inf \text{ cond } D_2 (A D_2)$$

$$iii) \inf \text{ cond } D_1 (D_1 A)$$

$$iv) \inf D \text{ cond } (D A D) \text{ para } A \text{ hermitiana.}$$

Aquí podemos asumir que D_1, D_2, D , son no singulares para el problema.

$$\inf. \text{ cond } D_1 D_2 (D_1 A D_2)$$

y en el caso hermitiano el problema $\inf D \text{ cond } (D A D)$ para una A hermitiana, son presentados por ejemplo en el análisis de error de métodos directos para la solución de ecuaciones lineales.

El problema del $\inf. \text{ cond } D_1 D_2 (D_1 A D_2)$, es

importante para obtener los mejores límites posibles para la inclusión de teoremas, finalmente la cantidad $\inf D^2 \text{ cond}(A D^2)$ está también una medida natural para la independencia lineal de los vectores columna los cuales forman A.

En esta relación, A puede ser una matriz rectangular.

Ahora supongamos un ejemplo numérico de un sistema de 2×2 alterado por multiplicar la primera ecuación por 10^5 . Entonces el sistema sería:

$$\begin{aligned} 10.0 x + 100,000 y &= 100,000. \\ 1.0 x + 1.0 y &= 2.00 \end{aligned}$$

El efecto de escalonar es hacer 10.0 el mayor pivote es la primera columna. Entonces la eliminación de x en la segunda ecuación del sistema de arriba en aritmética de punto flotante en tres dígitos, resultará en una nueva segunda ecuación.

$$- 10,000 y = - 10,000.$$

La solución al revés lleva a $y = 1.00$ y el resultado tremendo de $x = 0.00$.

Veamos que el escalonar pobremente con una buena estrategia de pivoteo nos fuerza a caer en el mismo error de redondeo enorme que obtenemos del conjunto original de ecuaciones y una mala estrategia de pivoteo.

La conclusión de esto es que una buena estrategia de pivoteo es solo buena cuando la matriz está propiamente escalonada. Sin embargo puede ser omitido el método que no conocemos tanto por algoritmos garantizados para matrices bien escalonadas.

Es normal escalonar matrices por multiplicación simplemente de renglones y columnas por factores. En efecto, se eligen matrices diagonales no singulares D_1 y D_2 , y entonces se escala A por la transformación. $D_1 A D_2$

$$A \longrightarrow D^{-1} A D$$

Porque el $\text{cond}(A)$ es un factor de todos nuestros límites de error y teoremas de convergencia.

Es natural el desear seleccionar D^{-1} y D , tanto para reducir el $\text{cond}(D^{-1} A D)$ como para reducir un valor tan razonablemente como sea posible. Generalmente se usan potencias de la base de punto flotante para escalar factores, para evitar la introducción de errores de redondeo en el escalonamiento. O alternativamente, se puede usar el escalonamiento sólo implícitamente, sin alterar actualmente los elementos de A .

TEOREMA (B. L. Bauer)

Si el orden del conjunto de elementos pivoteados es seleccionado para avanzar en el escalonamiento de una matriz A , por potencias de punto flotante, no cambia un simple dígito el significado de algún número final intermedio en la solución $Ax = b$, por eliminación gaussiana.

Este teorema fué presentado por Bauer. Así el único efecto posible del escalonamiento de A , en los errores de redondeo, debe ocurrir a través de cambiar el orden de los pivotes.

El ejemplo mostró que el cambio en los pivotes puede hacer una gran diferencia en el resultado.

Algunas veces se advierte al escoger D^{-1} y D , tal que la matriz resultante $D^{-1} A D$, tiene máximo elemento en cada renglón y cada columna (en valor absoluto) en el intervalo $(.1, 1)$, en cuanto al número base que se este usando.

Sin embargo, se ha mostrado que esto no siempre conduce a un buen escalonamiento, si.

$$A = \begin{bmatrix} 1 & 1 & 2 \times 10^9 \\ 2 & -1 & 10^9 \\ 1 & 2 & 0 \end{bmatrix}$$

Entonces ambas, de las siguientes matrices están escalonadas decimalmente y son equivalentes a A.

$$A_C = \begin{bmatrix} .1 & .1 & .2 \\ .2 & -.1 & .1 \\ .1 & .2 & 0 \end{bmatrix}$$

$$A_R = \begin{bmatrix} & -10 & -10 & \\ & 10 & 10 & .2 \\ & -10 & -10 & \\ 2 \times 10 & -10 & -10 & .1 \\ .1 & .2 & & 0 \end{bmatrix}$$

Sin embargo, A_C es una matriz bien condicionada que no ofrece dificultades en la solución de un sistema de ecuaciones mientras que A_R está mal condicionada y da muchos problemas para la eliminación.

Se ha estudiado el problema de encontrar D_1 y D_2 para minimizar el condicional $(D_1^{-1} A D_2)$.

De esto se deduce que la solución depende de ciertas propiedades de las matrices no negativas $|A|^{-1}$ y $|A|^{-1}$.

Claramente podemos calcular A^{-1} para encontrar un escalonamiento razonable, tal que podemos calcular A^{-1} , y está es una cuestión abierta, como encontrar un algoritmo para escalar, demostrablemente bueno y conveniente. Los algoritmos existentes, son unos y otros muy superficiales ó potencialmente muy bajos. El lado bueno de ésta cuestión del escalonamiento, es que está visto que es para tratar con una matriz rara, la cual despues de un buen escalonamiento la cambia de intratable a tratable.

Así el escalar optimamente para la norma máxima $(p = p = \infty)$ está caracterizada por $D^{-1}A$ ó $(A^{-1}D)$ respectivamente, teniendo sumas de renglones iguales a valores absolutos.

Ejemplo.

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 3 & 6 \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} 3 & -3 & 1 \\ -3 & 5 & -2 \\ 1 & -2 & 1 \end{bmatrix}$$

Escalonar optimamente para la norma máxima.

$$D_1 A D_2 = \begin{bmatrix} 3/5 \sqrt{10} & 3 & 2/5 \sqrt{10} \\ 1 & \sqrt{10} & 2 \\ 1/5 \sqrt{10} & 3 & 4/5 \sqrt{10} \end{bmatrix}$$

$$D_2^{-1} A D_1^{-1} = \begin{bmatrix} 1/2 \sqrt{10} & -3 & 1/2 \sqrt{10} \\ -1 & \sqrt{10} & -2 \\ 1/4 \sqrt{10} & -3 & 3/4 \sqrt{10} \end{bmatrix}$$

Donde el condicional $(D_1 A D_2)_2 = (3 + \sqrt{10})^2 = 37.974$

$$A D_2 = \begin{bmatrix} 7 & 10 & 4 \\ 7 & 20 & 12 \\ 7 & 30 & 24 \end{bmatrix}$$

$$D_2^{-1} A = \begin{bmatrix} 3/7 & -3/7 & 1/7 \\ -3/10 & 1/2 & -1/5 \\ 1/4 & -1/2 & 1/4 \end{bmatrix}$$

Donde el cond $(A D_2)_2 = 61$

$$D_1 A = \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/6 & 1/3 & 1/2 \\ 1/10 & 3/10 & 3/5 \end{bmatrix}$$

$$A^{-1} D_1^{-1} = \begin{bmatrix} 9 & -18 & 10 \\ -9 & 30 & -20 \\ 3 & -12 & 10 \end{bmatrix}$$

$$\text{Cond} (D_1 A) = 59$$

Ahora el escalonamiento óptimo para la norma euclidiana

(Ambas A y A⁻¹ tienen distribución de signo cuadrangular)

$$D_1 A D_2 = \begin{bmatrix} 3/5 \sqrt{10} & 3 & 2/5 \sqrt{5} \\ 3 & \sqrt{10} & 6 \\ 2/5 \sqrt{5} & 6 & 4/5 \sqrt{10} \end{bmatrix}$$

$$D_2^{-1} A^{-1} D_1^{-1} = \begin{bmatrix} 1/2 \sqrt{10} & -\sqrt{3} & 1/2 \sqrt{5} \\ -\sqrt{3} & \sqrt{10} & -\sqrt{6} \\ 1/2 \sqrt{5} & -\sqrt{6} & 3/4 \sqrt{10} \end{bmatrix}$$

Donde el $\text{cond} (D_1 A D_2) = (3 + \sqrt{10}) = 37.974$

Que es el mismo del escalonamiento para la norma máxima.

C A P I T U L O

X.

ANALISIS DE ERROR
DE
REDONDEO

El moderno análisis de error de redondeo originado en 1940, centró su atención en el algoritmo de Eliminación Gaussiana para resolver sistemas lineales al mismo tiempo que los analistas numéricos contemplaban el inminente arribo de las maquinas computadoras electrónicas. Encarando la posibilidad de poder llevar a cabo computaciones que involucraban números muy grandes de operaciones aritmeticas y estaban muy interesados en los efectos acumulativos, que los errores de redondeo involucraban .Y centraron su atención en el algoritmo de Gauss, porque este es el problema computacional mas importante, y ademas, porque es matemáticamente muy simple y por consiguiente un punto de partida prometedor.

Muy pronto ensayaron para obtener un límite para el error en la solución computada abandonando las más pesimistas conclusiones. Esto fué un estímulo para tratar de investigar métodos alternativos en la espera de que ellos no sufrieran de la inestabilidad numérica de la eliminación gaussiana. Alguno de estos métodos alternativos han probado ser de gran importancia, aunque mucho menos estables que la eliminación gaussiana.

Desde el principio esos trabajos en análisis de error se encontraron a sí mismos involucrados en las complejidades formales del tan llamado " acumulación de errores de redondeo " en la eliminación Gaussiana y esto tuvo un efecto adverso en la presentación del objeto de dicho análisis. Retomando la situación, dejemos por el momento el análisis de error detallado de algoritmos específicos y nos concentraremos en lugar de ello en las limitaciones inherentes al computo en presencia de errores de redondeo. Esto lleva a conclusiones generales las cuales son muy ilustrativas cuando se aplican a la solución de sistemas lineales y por consiguiente seguiremos ese curso en este capítulo.

Las Operaciones Aritméticas Básicas.

Muchos de los primeros análisis de error estaban formulados en términos de la aritmética de punto flotante, porque ninguna de las computadoras estaban bien preparadas para tener hardware de punto flotante. Esto de nuevo fué desafortunado, dado que los algoritmos en punto flotante casi invariablemente involucran un proceso continuo de

balanceo adecuado y esto complica la descripción del algoritmo computado por sí mismo. Por extraño que parezca hubo una concepción popular en los años tempranos en que sería difícil poner el error de análisis en punto flotante en bases rigurosas. De hecho el análisis de error riguroso de computación en punto flotante es generalmente más simple que el de la computación de punto fijo. Dado que casi todo el software matemático es ahora expresado en términos de la aritmética de punto flotante y nos restringiremos a esto.

Asumiremos que la computación es llevada a cabo en la base β con t dígitos en la mantisa y que la unidad aritmética en la computadora acepta números normalizados y produce números normalizados como la salida de las operaciones aritméticas.

Nos referimos a este conjunto de números que puede ser exactamente representados en la computadora con números digitales. Después será asunto nuestro el uso de números con $2t$ dígitos (números digitales en doble precisión), pero por el momento ignoramos esta posibilidad. Aunque los problemas de overflow y underflow son importantes en consideración al robustecer algoritmos, los ignoraremos aquí, y damos por hecho que el exponente es el adecuado.

Usaremos la notación $fl(a * b)$, donde $*$ denota cualquier operación aritmética básica, $+$, $-$, \times , \div , para enotar el valor de $a * b$. Está implícito en el uso de esta notación que a y b son números y $fl(a * b)$ es, por definición un número digital. Es el resultado establecido por la computadora después de cualquier redondeo, corte, etc. y la renormalización puede ser necesaria. Aunque se le da gran importancia a los intentos que están hechos para asegurar que el futuro de las computadoras tendrá unidades aritméticas realmente satisfactorias, la naturaleza de los procesos de redondeo no serán de fundamental importancia en este capítulo. Meramente se asumirá por el momento que cada una de las operaciones básicas

$$fl(a * b) = (a * b)(1 + \epsilon) \quad , \quad |\epsilon| < k \beta^{-t} \quad (1)$$

Dónde el orden de k es la unidad y es independiente de a y b .

En algunas computadoras la relación

$$fl(a * b) = (a * b)(1 + \epsilon)$$

No siempre se toma como adición y sustracción y debemos conformarnos con resultado de la forma.

$$fl(a + b) = a(1 + \epsilon_1) + b(1 + \epsilon_2) \quad , \quad |\epsilon_i| < k \beta^{-t} \quad (2)$$

Donde en general no es posible tomar $\epsilon_1 = \epsilon_2$. Por consiguiente este RESULTADO MAS DEBIL no tendría serias repercusiones para los resultados dados.

La relación

$$fl(a + b) = a(1 + \epsilon_1) + b(1 + \epsilon_2)$$

implica que el resultado computado de cada operación individual tiene un error relativo bajo. Sin embargo hay un medio alternativo de interpretar la relación

$$fl(a * b) = (a * b)(1 + \epsilon)$$

la cual por consiguiente parece trivial, tiene sorprendentes implicaciones, las cuales nos pueden ilustrar en el caso de la multiplicación. La relación

$$fl(a * b) = (a * b)(1 + \epsilon)$$

establece que el producto computado es el producto exacto de $a(1 + \epsilon_1)$ y b ó de a y, $b(1 + \epsilon_2)$, ó más simétricamente,

de $a(1 + \epsilon_1)$ y $b(1 + \epsilon_2)$. dado que ϵ es pequeña podemos decir que el resultado computado es exacto para dos operandos los cuales difieren de los dados por errores relativos chicos; al presentar el resultado en esta forma estamos retomando los errores computados por errores equivalentes en los datos. Por supuesto hay un infinito de caminos en los cuales esto podría ser hecho pero muchos de ellos no son de interés; así por ejemplo podríamos tomar $a/2$ y $2b(1 + \epsilon)$ pero esto no involucra perturbaciones

relativas grandes en los datos.

Análisis de Error Inverso

La interpretación que hemos dado de la expresión,

$$fl(a * b) = (a * b)(1 + \epsilon) \quad \text{donde } |\epsilon| < k\beta^{-t}$$

Usualmente es referido como un error de análisis inverso y en general puede ser descrito en los siguientes términos.

Consideremos cualquier algoritmo computado con elementos x_i , con $i = 1, p$, y los elementos solución y_j con $j = 1..q$

Cuando los errores de redondeo son hechos en la ejecución de los algoritmos de elementos computados y_j no serán exactos, pero en general, será un número infinito del conjunto de elementos $x_i(1 + \epsilon)$ para los cuales los y_j computados son la solución exacta. En el análisis de error inverso se trata de probar la existencia de tales conjuntos de datos para los cuales todos los ϵ_i son pequeños.

Encontramos que para las operaciones aritméticas básicas los valores naturales de ϵ_i son sugeridos por las relaciones fundamentales satisfechas por los valores computados; con muchos algoritmos de matrices las aplicaciones de la relación,

$$fl(a * b) = (a * b)(1 + \epsilon) \quad \text{donde } |\epsilon| < k\beta^{-t}$$

para todas las operaciones involucradas también lleva de una manera equitativamente natural a un conjunto de datos apropiados.

$$x_i(1 + \epsilon_i)$$

Ahora para todas las operaciones básicas en sí mismas, estamos forzando para aceptar las perturbaciones relativas en el conjunto de datos (aquí son los operandos a, b) de

-t

arriba para $k\beta$, donde k es la constante en la expresión (1). Para un algoritmo muy complicado sería razonable esperar el poder garantizar guardar el ϵ bajo este nivel.

Realmente si el número de operaciones involucradas en el algoritmo es $m\mu$, es decir, un promedio de m operaciones por elemento de datos, un algoritmo para el cual podríamos

establecer límites a-priori de $m\mu k\beta$, para la ϵ podría ser considerado como extremadamente estable con respecto a errores de redondeo.

Un algoritmo para el cual los límites existen satisfactoriamente para ϵ , se dice que es inversamente estable.

Debemos distinguir entre la existencia de estabilidad inversa y nuestra habilidad para establecerla. Adicionalmente si tomamos dos algoritmos A y B para resolver el mismo problema el establecimiento riguroso de los límites más pequeños para las ϵ correspondiendo a A

que aquel correspondiente a B, no necesariamente significa que sea más perceptivo para A que para B.

Sin embargo, ahora tenemos una experiencia suficientemente grande del análisis de error de algoritmos de matrices que en muchas situaciones podemos estar razonablemente confiados de que nuestros límites dan una estimación realista de ejecución cuando hemos hecho una apropiada concesión para la distribución estadística de los errores de redondeo. En la práctica será de nuestro interés la precisión de la solución computada es decir, la distancia entre la y computada y la y exacta correspondiente a la

x dada. La estabilidad inversa no garantiza que la solución computada sea de gran precisión pero hasta ahora como no se hace esto es obviamente debido a la gran sensibilidad de la solución para perturbaciones en los datos.

Esta sensibilidad es una propiedad inherente del problema

computado y no tiene nada que ver con la efectividad del algoritmo.

Hay tres características importantes del análisis de error inverso.

- a).- Pone los errores hechos durante el curso de la solución en la misma posición, como errores en los datos. En la práctica los datos rara vez serán exactos; comúnmente involucran errores de observación, si los errores relativos probables en las a_i dadas son más grandes que ϵ_i dada, por un análisis de error inverso entonces, los errores de redondeo, efectuados durante el curso del algoritmo son menos importantes que los errores iniciales en los datos. Cuando todas las a_i son conocidas exactamente pueden no ser números digitales; los errores serán entonces involucrados en las representaciones computacionales de a_i aunque los errores relativos introducidos de esta manera estarían limitados por $k\beta^{1-t}$.
- b).- Para estimar los errores en la solución computada podemos usar la teoría de las perturbaciones y podemos por consiguiente sacar una fuente muy rica de información.
- c).- La tercera característica es la que parece brotar de la gran naturaleza del análisis de error inverso. Para muchos algoritmos matemáticos algebraicos produce una simplicidad formal inesperada y muchas veces el efecto de los errores es puramente aditivo, cuando esto sucede se puede hacer un análisis de error para cubrir los grandes desatinos aritméticos, tanto como los meros errores de redondeo.

Análisis de Error de Multiplicaciones Anidadas.

Una ilustración de las tres características que hemos mencionado es estipulando en un contexto de gran simplicidad formal, por la evaluación de un polinomio por multiplicaciones anidadas.

Sea,

$$f(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_0 \quad (3)$$

Entonces $f(x)$ puede ser evaluado por el algoritmo.

$$s_n = a_n$$

$$s_r = x s_{r+1} + a_r \quad \text{donde } r = n-1, \dots, 2, 1 \quad (4)$$

$$f(x) = s_0$$

Podemos llevar a cabo el análisis de error inverso de este algoritmo en el cual, incluimos el efecto de cualquier error efectuado en las operaciones aritméticas incluidas, posiblemente grandes desatinos. Denotemos las cantidades

computadas por \bar{s}_r , y tenemos.

$$\bar{s}_r = (x \bar{s}_{r+1} + a_r) = (x \bar{s}_{r+1} + e_r) + a_r + f_r \quad (5)$$

computada

Donde e_r denota el error efectuado al calcular $x \bar{s}_{r+1}$ (esto podría ser meramente un error de redondeo ó un error sustancial debido a un mal funcionamiento de alguna clase) De aquí,

$$\bar{s}_r = x \bar{s}_{r+1} + (a_r + e_r + f_r) \quad (6)$$

Muestra que las \bar{s}_r computadas son exactas para el polinomio con coeficiente $a_r + e_r + f_r$. Nótese la gran simplicidad formal de este resultado; sin intersecciones entre los errores sucesivos que han sido considerados. Supóngase, por ejemplo, que estamos trabajando con a entera y una x

entera los únicos errores cometidos son desatinos. Supóngase que solo dos errores son cometidos, y

que estos ocurrieron cuando se calculaban \bar{s}_p y \bar{s}_q .

Entonces \bar{s}_0 (el valor computado de $f(x)$) sería el valor exacto de un polinomio, el cual difiere del polinomio dado solo en los elementos a_p y a_q , los cuales serían argumentos precisamente por la magnitud de los errores en los pasos p -ésimos y q -ésimos respectivamente.

Como un ejemplo sencillo consideremos la evaluación del polinomio.

$$f(x) = 3x^4 + 4x^3 + 2x^2 + x + 5 \quad \text{con } x = 3. \quad (5)$$

La evaluación correcta está descrita obviamente por la tabla

3	4	2	1	5	
	9	39	123	372	
3	11	41	124	377	$f(3) = 377$

Ahora consideremos la evaluación incorrecta descrita por.

3	4	2	1	5	
	10	45	142	435	
3	15	47	145	440	

Observemos que aquí los errores fueron cometidos en la segunda columna; un error de +1 en la multiplicación, un +1 en la adición. Estos son contados según por una

perturbación de +2 en el coeficiente de x^3 . Los errores también son cometidos en la cuarta columna; un error de +1 en la multiplicación y un error de +2 en la adición. Estos son contados según por una perturbación de +3 en el coeficiente de x . De aquí el valor computado es el valor exacto de $x = 3$ de,

$$3x^4 + (4 + 1 + 1 + x^3 + 2x^2 + (1 + 1 + 2)x + 5$$

es decir,
$$3x^4 + 5x^3 + 2x^2 + 4x + 5 \tag{6}$$

Nótese que la diferencia entre los elementos correspondientes en las columnas en la computación exacta y errónea son irrelevantes y no viene al caso el análisis de error.

Retornando ahora al análisis de error de redondeo es conveniente, expresar el efecto que hace el error de redondeo en la adición y sustracción de la forma,

$$fl(a \pm b) = (a \pm b) / (1 + \epsilon), \text{ donde } |\epsilon| \leq k\beta^{-t} \tag{7}$$

Mejor que la expresión

$$fl(a * b) = (a * b) (1 + \epsilon), \quad |\epsilon| \leq k\beta^{-t}$$

La computación de \bar{s}_r toma lugar en los pasos

$$b_r = fl(x_{r+1} \bar{s}_{r+1}) = x_{r+1} \bar{s}_{r+1} (1 + \epsilon_r), \quad |\epsilon_r| \leq k\beta^{-t} \tag{8}$$

$$\bar{s}_r = fl(b_r + a_{r+1}) = (b_r + a_{r+1}) / (1 + \eta_r), \quad \left| \frac{\eta_r}{r} \right| \leq k\beta^{-t} \tag{9}$$

dando

$$\bar{s}_r = b_r + a_{r+1} - \bar{s}_r \eta_r$$

$$= x_{r+1}^{\bar{s}} (1 + \epsilon_r) + a_{r+1} - \bar{s}_r \Phi$$

$$= x_{r+1}^{\bar{s}} (a_{r+1} + x_{r+1}^{\bar{s}} \epsilon_r - \bar{s}_r \Phi) \quad (10)$$

De aquí la $f(x)$ computada es exacta para el polinomio con coeficientes $a_{r+1} + x_{r+1}^{\bar{s}} \epsilon_r - \bar{s}_r \Phi$. El análisis descrito aquí es mencionado como un análisis de error inverso corrido; para las perturbaciones es el conjunto de datos en términos de los valores computados intermedios actuales (aquí la \bar{s}_r) tal que, cuando la computación ha sido completada tenemos a la mano límites para las perturbaciones equivalentes en el conjunto de datos. Nótese que si ni $\left| x_{r+1}^{\bar{s}} \right|$ ni $\left| \bar{s}_r \right|$ es significativamente más grande que $\left| a_{r+1} \right|$ la perturbación $^{-t}$ relativa equivalente es por sí misma el nivel de β y Φ , el orden del polinomio no está involucrado. Y ésta es una situación muy común.

Para este algoritmo simple podemos dar un error de análisis muy satisfactorio a-priori. De las relaciones,

$$b_r = f_l(x_{r+1}^{\bar{s}}) x_{r+1}^{\bar{s}} (1 + \epsilon_r), \quad \left| \frac{\epsilon_r}{r} \right| \leq k\beta^{-t}$$

y,

$$\bar{s}_r = f_l(b_r + a_{r+1}) = (b_r + a_{r+1}) \div (1 + \Phi_r), \quad \left| \frac{\Phi_r}{r} \right| \leq k\beta^{-t}$$

Tenemos sucesivamente

$$\bar{s}_n \approx a_n$$

$$\frac{\epsilon}{n} \approx \frac{a}{n} \quad (11)$$

$$\frac{\epsilon}{n-1} = \frac{\frac{a}{n} (1 + \epsilon) \times}{(1 + \epsilon) \frac{n-1}{n-1}} + \frac{\frac{a}{n-1}}{(1 + \epsilon) \frac{n-1}{n-1}} \quad (12)$$

$$\begin{aligned} \frac{\epsilon}{n-1} = & \frac{\frac{a}{n} (1 + \epsilon) (1 + \epsilon) \times^2}{(1 + \epsilon) \frac{n-1}{n-1} (1 + \epsilon) \frac{n-2}{n-2}} + \frac{\frac{a}{n-1} (1 + \epsilon)}{(1 + \epsilon) \frac{n-1}{n-1} (1 + \epsilon) \frac{n-2}{n-2}} + \\ & + \frac{\frac{a}{n-2}}{(1 + \epsilon) \frac{n-2}{n-2}} \end{aligned} \quad (13)$$

Y se deduce que, (14)

$$s = a (1 + E) \times \frac{n}{n} + a (1 + E) \times \frac{n-1}{n-1} + \dots + a (1 + E) \frac{0}{0}$$

donde

$$1 + E = \frac{(1 + \epsilon) (1 + \epsilon) \dots (1 + \epsilon)}{(1 + \epsilon) (1 + \epsilon) \dots (1 + \epsilon)} \quad (15)$$

$$1 + E = \frac{(1 + \epsilon) (1 + \epsilon) \dots (1 + \epsilon)}{(1 + \epsilon) (1 + \epsilon) \dots (1 + \epsilon)} \quad (16)$$

donde $r = n - 1, \dots, 1, 0$

De los límites en las ϵ_i y Φ_i , tenemos.

$$\frac{(1 - \alpha)^n}{(1 + \alpha)^n} \leq 1 + \frac{E}{n} \leq \frac{(1 + \alpha)^n}{(1 - \alpha)^n},$$

$$\frac{(1 - \alpha)^r}{(1 + \alpha)^{r+1}} \leq 1 + \frac{E}{r} \leq \frac{(1 + \alpha)^r}{(1 - \alpha)^{r-1}},$$

con $r = n-1, \dots, 0$.

donde $\alpha = k\beta^{-t}$

Igual en una computadora para la cual $k\beta^{-7} = 10$ (que es un tanto baja precisión) esto representa un error relativo bajo a no ser que n sea ordinariamente larga. Nótese que tenemos un error relativo pequeño en cada coeficiente individual y esto es verdadero, sin embargo muchas cancelaciones pueden tomar lugar en la computación. Hay una diferencia interesante entre este análisis a-priori y el análisis corrido. En el análisis corrido encontramos expresiones para perturbaciones en los datos y este

conjunto de datos perturbados da todos los \bar{s}_r computados.

En el análisis de error a-priori los datos perturbados dan la \bar{s}_0 computado pero diferentes perturbaciones (en general más pequeño) son requeridos para cada una de las \bar{s}_r . Esto es evidente de las relaciones,

$$\bar{s}_r = x_{r+1} \bar{s}_{r+1} + (a_{r+1} + x_{r+1} \epsilon - \bar{s}_r \Phi_r)$$

$$\bar{s} \approx a_n$$

$$\bar{s}_{n-1} = \frac{a_n (1 + \epsilon_{n-1}) \times}{(1 + \epsilon_{n-1})} + \frac{a_{n-1}}{(1 + \epsilon_{n-1})}$$

Y en sus sucesores.

El análisis muestra que la multiplicación anidada es muy estable inversamente. En el caso cuando las \bar{x}_{r+1} y \bar{s}_r son comparables con a_r (una situación extremadamente común)

el análisis de error corrido muestra que la multiplicación anidada es más ó menos tan estable como posiblemente se espera.

Además cuando los coeficientes no son números digitales, como que las sola representación de ellos en la computadora involucra un redondeo en cada una, el probable efecto de esto es casi tan grande como el efecto de los errores de redondeo cometidos en la computación.

Ahora consideremos los efectos de estas pequeñas perturbaciones relativas en los coeficientes de los ceros del polinomio. Si asumimos que las perturbaciones en los coeficientes y en los ceros podemos considerar la perturbación en cada coeficiente por separado. Si γ_i , con $i = 1, \dots, n$, son los ceros del polinomio dado (dando por supuesto que es simple) entonces,

$$\frac{\delta \gamma_i}{\gamma_i} = \frac{-\gamma_i^r}{a_n \prod_{j \neq i} (\gamma_i - \gamma_j)} \quad (18)$$

Consideremos el polinomio de orden 20 con $y_i = i$,
 con $i = 1, \dots, 20$ y se verificará que,

$$\frac{\delta y_i}{\delta a_r}$$

es extremadamente grande para algunos valores de i y r .
 De hecho,

$$\frac{\left| \begin{array}{c} \delta y_{16} \\ \delta a_{15} \end{array} \right|}{\left| \begin{array}{c} \delta y_{15} \\ \delta a_{15} \end{array} \right|} = 0.4 \times 10^5$$

Y dado que $a_{15} = 10^{10}$, una perturbación en a_{15} , da una
 perturbación.

$$\delta y_{16} = (0.4 \times 10^{15}) \epsilon$$

De aquí cuando $\epsilon = 0$ (β^{-15}) vemos que, a no ser que $\beta^{-t} \ll \beta^{-15}$

la perturbación inducida en y_{16} no será en ningún
 sentido pequeña. En una computadora de diez dígitos
 decimales algunos de los ceros del polinomio perturbado no
 llevará relación con esos del polinomio verdadero.

Sí bien los resultados de arriba son triviales sus
 implicaciones para el analista numérico son totalmente
 fundamentales los eigen-valores de una matriz de orden n
 son los ceros del polinomio característico. Es atractivo
 superficialmente para determinar los eigen-valores. Nuestro
 análisis muestra que cualquier algoritmo basado en tal
 técnica esta predispuesto a fracasar como un algoritmo
 general. La probabilidad de que el polinomio explícito
 tenga números digitales en sus coeficientes es pequeño y de
 aquí los errores cometidos en la representación del
 polinomio será muchas veces fatal.

Esta simple observación escapó al analista numérico por muchos años y la mayoría de los algoritmos propuestos para resolver el problema de los eigen-valores para matrices no normales estaban basados en el cálculo del polinomio explícito.

C A P I T U L O
XI.

VALORES CARACTERISTICOS
DE
MATICES SIMETRICAS

El espacio no permite un tratamiento tan extenso del problema de los valores característicos, como el dado para el problema de las ecuaciones lineales, podemos solamente mencionar algunos métodos.

Continuaremos resolviendo los nuevos problemas mediante la simplificación de una matriz, convirtiéndola en diagonal ó en triangular superior, pero el paso fundamental ya no será sustraer de una fila el múltiplo de otra. Ya no nos interesará conservar el espacio fila de una matriz, sino preservar sus valores propios. Las operaciones elementales en la filas no lo hacen.

El determinante conduce a una solución formal; a la regla de Cramer en el caso $Ax = b$ y al polinomio $\det(A - \gamma I)$ cuyas raíces serán los valores propios. Como siempre, si $n = 2$ ó 3 , puede realmente usarse el determinante para resolver el problema; el cálculo de los valores propios para una n grande es una tarea más larga y difícil que resolver $Ax = b$, e incluso el mismo Gauss no aportó mucho.

El primer paso es entender que son los valores propios y como pueden ser útiles. Por tanto estamos interesados sólo en aquellos valores particulares y para los cuales existe algún vector propio x distinto de cero. Para poder ser útil el espacio nulo $A - \gamma I$ debe contener algún vector diferente de cero. Esto es que $A - \gamma I$ debe ser singular. Para esto el determinante nos da un criterio definitivo.

El número γ será un valor propio de A , con un vector propio corespondiente distinto de cero, si y sólo si

$$\det(A - \gamma I) = 0$$

esta es la ecuación característica de la matriz A .

En la ecuación $Ax = \gamma x$, la mayoría de los vectores x no satisfacen esta ecuación, ya sea γ un valor propio ó no. Una x típica cambia de dirección cuando se multiplica por A , así que Ax no es múltiplo de x . Esto significa que sólo ciertos números especiales y son valores propios y que sólo ciertos vectores especiales son vectores propios. Claro que si A fuera un múltiplo de la matriz identidad, entonces ningún vector cambiaría de dirección y todos los vectores serían vectores propios. Pero en general los vectores

Propios son pocos y separados.

Sin embargo la sustancia de este capítulo está en otro lado. La cuestión más importante es explicar cómo se descompone un sistema de ecuaciones al ir hallando los vectores propios. Estos vectores propios son los modos normales del sistema y actúan independientemente. Podemos observar el comportamiento de cada vector propio por separado y después combinar estos modos normales para encontrar la solución. Para decir lo mismo en otras palabras, la matriz subyacente ha sido diagonalizada.

Para concluir resumiremos algunos de los hechos básicos acerca de los valores propios y los vectores propios; $Ax = \gamma x$. La matriz A de $n \times n$, está dada y el problema es encontrar aquellos vectores especiales x para los cuales A actúa como una multiplicación simple; Ax apunta en la misma dirección que x . La cuestión es encontrar primero los valores propios, reescribiendo γx como $\gamma I x$ y pasando este término a lado izquierdo: $(A - \gamma I) x = 0$. Un vector propio de A es un vector nulo de $A - \gamma I$. En otras palabras, la cuestión clave es la siguiente: si modificamos A por varios múltiplos de la matriz identidad, ¿qué cambio la hará singular? Estas modificaciones revelarán los valores propios y entonces podremos calcular los vectores propios.

Si A es singular, entonces una posibilidad es no modificarla. Uno de los valores propios de una matriz singular es $\gamma = 0$ y el espacio nulo contiene los vectores propios correspondientes. Pero el que A sea singular no significa nada en especial; significa solo que $\gamma = 0$ es un valor propio. Ya sea que A sea singular ó no, todos sus valores propios están basados en lo mismo: la combinación $A - \gamma I$ es singular y el espacio nulo de esa combinación es el espacio de los vectores propios correspondientes a γ .

Para saber cuando $A - \gamma I$ es singular, calculamos su determinante. Este determinante es un polinomio en γ de grado n y se conoce como polinomio característico de A . La ecuación $\det(A - \gamma I)$ es la ecuación característica y sus raíces $\gamma_1, \gamma_2, \dots, \gamma_n$ (que pueden ser números reales ó no, y

pueden incluir ó no algunas repeticiones de la misma γ) son los valores propios de A . Esto es:

Cada una de las siguientes condiciones es necesaria y suficiente para que el número γ sea un valor propio de A :

- 1).- Hay algún vector x diferente de cero tal que $Ax = \gamma x$
- 2).- La matriz $A - \gamma I$ es singular
- 3).- $\det.(A - \gamma I) = 0$

Como un ejemplo ilustrativo, consideremos la matriz simétrica.

$$A = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{bmatrix}$$

Su polinomio característico es

$$\det(A - \gamma I) = \begin{vmatrix} 1-\gamma & -1 & 0 \\ -1 & 2-\gamma & -1 \\ 0 & -1 & 1-\gamma \end{vmatrix}$$

$$= -\gamma^3 + 4\gamma^2 - 3\gamma$$

Por lo tanto, la ecuación característica es

$$-\gamma^3 + 4\gamma^2 - 3\gamma = -\gamma(\gamma - 1)(\gamma - 3) = 0$$

y los valores propios son reales y distintos : $\gamma_1 = 0$,

$\gamma_2 = 1$, $\gamma_3 = 3$. Para cada valor propio γ_i busquemos un vector propio correspondiente x_i .

$$\text{Así } \gamma_1 = 0$$

$$(A - 0I) \times_1 = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{bmatrix} \times_1 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

y,

$$\times_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$\gamma_2 = 1 :$$

$$(A - I) \times_2 = \begin{bmatrix} 0 & -1 & 0 \\ -1 & 1 & -1 \\ 0 & -1 & 0 \end{bmatrix} \times_2 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

y,

$$\times_2 = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix}$$

$$\gamma_3 = 3 :$$

$$(A - 3I) x_3 = \begin{bmatrix} -2 & -1 & 0 \\ -1 & -1 & -1 \\ 0 & -1 & -2 \end{bmatrix} x_3 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

y,

$$x_3 = \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}$$

Esta cantidad de cálculos es inevitable y desde luego la mayoría de los polinomios cúbicos no se factorizarán tan fácilmente como sucedió con este. No cabe duda que algebraicamente y desde el punto de vista de los cálculos, el problema de los valores propios es mucho más difícil que $Ax = b$. Para un sistema lineal, un número finito de pasos de eliminación produjo la respuesta exacta en un tiempo finito (o similarmente la regla de Cramer nos dió una fórmula exacta para la solución) en el caso de los valores propios no pueden existir tales pasos ni tal fórmula.

Ahora bien, la suma de los n valores propios es igual a la suma de las n entradas de la diagonal de A :

$$\gamma_1 + \dots + \gamma_n = a_{11} + \dots + a_{nn}$$

Esta suma es la traza de A . Además, el producto de los n valores propios es igual al determinante de A .

Esto se confirma con el ejemplo anterior, donde la traza de A es $0 + 1 + 3 = 1 + 2 + 1 = 4$ y el determinante es $0 \cdot 1 \cdot 3 = 0$.

No deben confundirse los valores propios de una matriz y sus entradas diagonales. Normalmente son completamente diferentes. Sin embargo, señalaremos una situación donde esos números coinciden.

Si la matriz A es triangular (puede ser superior ó inferior y en particular diagonal) entonces los valores propios $\gamma_1, \dots, \gamma_n$ son exactamente los mismos que las entradas de la diagonal a_{11}, \dots, a_{nn} .

Por ejemplo si.

$$A = \begin{bmatrix} 1 & 1/4 & 0 \\ 0 & 3/4 & 1/2 \\ 0 & 0 & 1/2 \end{bmatrix}$$

entonces su polinomo característico es

$$\det \begin{bmatrix} 1-\gamma & 1/4 & 0 \\ 0 & 3/4-\gamma & 1/2 \\ 0 & 0 & 1/2-\gamma \end{bmatrix}$$

$$= (1 - \gamma) (3/4 - \gamma) (1/2 - \gamma)$$

El determinante es el producto de las entradas diagonales. Obviamente las raíces son $\gamma = 1$, $\gamma = 3/4$ y $\gamma = 1/2$; los valores propios ya estaban colocados a lo largo de la diagonal principal.

Hay otra situación en que los cálculos son fáciles. Supongamos que hemos encontrado ya los valores propios de A y también los vectores propios. Entonces los valores propios de A^2 son exactamente $\gamma_1^2, \dots, \gamma_n^2$, y cada vector propio de A^2 es también un vector propio de A . La demostración es típica en matemáticas, si comenzamos con

$Ax = \lambda x$, es obvio. Multiplicando nuevamente por A ,

$$A^2 x = A \lambda x = \lambda A x = \lambda^2 x.$$

Entonces λ es un valor propio de A^2 con el mismo vector propio x . Si la primera multiplicación por A deja intacta la dirección de x , lo mismo hará la segunda.

Este razonamiento no se aplica cuando estén involucradas dos matrices diferentes. Supongamos que λ es un valor propio de A y que μ es un valor propio de B . Entonces, en general, $\lambda\mu$ no es valor propio de AB . Dado que A y B no comparten el mismo vector propio.

En relación al problema de las ecuaciones lineales, el cálculo de los valores característicos de matrices se divide en dos casos, de acuerdo a la naturaleza de las matrices. Para matrices escasas ó huecas grandes, los métodos son para mejorar iteraciones infinitas y no serán concideradas aquí. Para matrices densas almacenadas, los métodos son algoritmos finitos. Si una matriz A es simétrica, sus valores característicos están muy bien determinados por los datos, en efecto, sea la matriz simétrica $B = A + E$, que tiene valores característicos β_i sea A (también simétrica) con valores característicos α_i

Entonces los valores característicos pueden ser numerados tal que.

$$\left| \begin{array}{c} \alpha_i - \beta_i \\ i \end{array} \right| \leq \left\| E \right\|, \text{ para toda } i$$

El análisis de error inverso se refiere a los eigen valores calculados (valores característicos) de una matriz A , al revez para una matriz $B = A + E$. Si puede probarse que E es pequeña entonces,

$$\left| \begin{array}{c} \alpha_i - \beta_i \\ i \end{array} \right| \leq \left\| E \right\|, \text{ para toda } i$$

muestra como son los errores de los eigen valores pequeños. En efecto los métodos actuales pueden llegar a los eigen valores.

que están en error por solo unos pocos dígitos, en los mismos dígitos significativos de los valores característicos grandes.

El método de Jacobi es una iteración infinita para matrices densas almacenadas, produce una serie de matrices ortogonalmente congruente con A.

$$A = U_k^t A U_k$$

Además, A_k converge a una matriz diagonal D, cuyas entradas diagonales son, por supuesto, los eigen-valores de A.

En efecto, cada A_{k+1} es calculada de una A_k previa por una rotación de coordenadas en el espacio 2 de algunos 2 índices i y j, una rotación elegida tal que a ij

Para algún k tal que A_k sea casi diagonal, las columnas de la matriz ortogonal correspondiente U_k son aproximadamente los vectores característicos columna de A. Además las columnas son ortogonales entre sí. Así el método de Jacobi da los eigen-vectores aproximados de calidad fina como un producto de la iteración básica. El solo programa es fácil de escribir, y es difícil que se haga mal. Hay algunos problemas teóricos sobre que tan buenos son los eigen-vectores y si el U_k converge actualmente.

El algoritmo original de Jacobi elige i y j para maximizar el valor absoluto de las a_{ij} de A_k . Los algoritmos de hoy en día modifican este criterio en una de dos formas:

- i).- En los métodos cíclicos de Jacobi, los elemntos fuera de la diagonal a_{ij} , se hacen cero en algún orden cíclico. ij

ii).- En los métodos de entrada de Jacobi, un elemento a_{ij} es seleccionado por eliminación solo cuando su valor absoluto está arriba de una cierta medida de entrada, la cual consigue la más pequeña, como la iteración progresa

Ejemplo:

Sea

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}, \quad S = \begin{bmatrix} 2 & \\ & 2 \end{bmatrix}, \quad T = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$S^{-1} T = \begin{bmatrix} 0 & 1/2 \\ 1/2 & 0 \end{bmatrix}$$

Si las componentes de x son v y w , el paso de Jacobi

$$Sx_{k+1} = Tx_k + b \quad \text{es}$$

$$\begin{aligned} 2v_{k+1} &= w_k + b_1 \\ 2w_{k+1} &= v_k + b_2 \end{aligned} \quad \text{ó} \quad \begin{bmatrix} v \\ w \end{bmatrix}_{k+1} = \begin{bmatrix} 0 & 1/2 \\ 1/2 & 0 \end{bmatrix} \begin{bmatrix} v \\ w \end{bmatrix}_k + \begin{bmatrix} b_1/2 \\ b_2/2 \end{bmatrix}$$

Donde tenemos que la matriz A se descompone como $A = S - T$ entonces la ecuación $Ax = b$ es la misma que $Sx = Tx + b$, por lo tanto obtenemos el paso de Jacobi anteriormente descrito.

Por supuesto que nada garantiza que el método sea bueno y una buena descomposición debe satisfacer dos condiciones:

i).- Debe ser fácil calcular el nuevo vector x_{k+1} . Por lo tanto S debe ser una matriz sencilla (! e invertible !); puede ser diagonal ó triangular.

ii).- La sucesión x_k debe converger a la solución x .

Si sustraemos la iteración $Sx_{k+1} = Tx_k + b$, de la ecuación original $Sx = Tx + b$ el resultado es una fórmula que involucra solamente los errores $e_k = x_k - x$:

$$Se_{k+1} = Te_k$$

Esta es una ecuación en diferencias. Comienza con el error inicial e_0 y después de k pasos produce un nuevo

$$e_k = (S^{-1} T)^k e_0$$

La cuestión de la convergencia es exactamente la misma que la cuestión de la estabilidad : $x_k \rightarrow x$ precisamente cuando $e_k \rightarrow 0$

Este método iterativo es convergente si y solo si cada valor propio de $S^{-1} T$ satisface $|\gamma| < 1$. Su razón de convergencia depende del tamaño máximo de $|\gamma|$, que se conoce como el radio espectral de $S^{-1} T$:

$$\rho(S^{-1} T) = \max_i |\gamma_i|$$

Recordando que una solución típica de $e_{k+1} = S^{-1} T e_k$ es

$$e_k = C_1 \lambda_1^k + \dots + C_n \lambda_n^k$$

Obviamente el mayor de los $|\lambda_i|$ será finalmente el dominante y gobernará la razón con la cual e_k converge a cero

Volviendo al ejemplo anterior tenemos la matriz decisiva
 -1

$S^{-1}T$ tiene valores propios $\pm 1/2$, lo que significa que en cada paso el error está partido por la mitad (un dígito binario más es el correcto). En este ejemplo, demasiado pequeño para ser típico, la convergencia es más rápida.

Con una matriz más grande A , existe una dificultad inmediata y muy práctica con la interacción de Jacobi expresada como

$$a_{11}(x_{1,k+1}) = (-a_{12}x_{2,k} - a_{13}x_{3,k} - \dots - a_{1n}x_{n,k}) + b_1$$

$$\vdots$$

$$a_{nn}(x_{n,k+1}) = (-a_{n1}x_{1,k} - a_{n2}x_{2,k} - \dots - a_{n,n-1}x_{n-1,k}) + b_n$$

Necesitamos conservar todas las componentes de x_k hasta terminar el cálculo de x_{k+1} . Una idea más natural, que requiere la mitad del almacenamiento, es comenzar a usar cada componente del nuevo vector x_{k+1} apenas se calcule; x_{k+1} ocupa el lugar de x_k una componente cada vez y, por tanto, podemos destruir x_k tan rápidamente como vamos creando x_{k+1} . Esto significa que la primera ecuación permanece como antes

$$a_{11}(x_{1,k+1}) = (-a_{12}x_{2,k} - a_{13}x_{3,k} - \dots - a_{1n}x_{n,k}) + b_1$$

La siguiente ecuación funciona inmediatamente con este nuevo valor de x_1 ,

$$a_{22}(x_{2,k+1}) = -a_{21}(x_{1,k+1}) + (-a_{23}x_{3,k} - \dots - a_{2n}x_{n,k}) + b_2$$

Y la última ecuación usará exclusivamente valores nuevos.

$$a_{nn} x_{n,k+1} = (-a_{n1} x_{1,k} - a_{n2} x_{2,k} - \dots - a_{n,n-1} x_{n-1,k}) + b_n$$

Este es el método de Gauss-Seidel, aun cuando aparentemente Gauss lo desconocía y Seidel no lo recomendaba.

Sorprendentemente detalle histórico, ya que no es un mal método. Nótese que cuando se mueven todos los términos x_{k+1} hacia el lado izquierdo, la matriz S es la parte triangular inferior de A . En el lado derecho, la otra matriz de la descomposición, T , es estrictamente triangular superior.

Ejemplo

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}, S = \begin{bmatrix} 2 & 0 \\ -1 & 2 \end{bmatrix}, T = \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix}, S^{-1} T = \begin{bmatrix} 0 & 1/2 \\ 0 & 1/4 \end{bmatrix}$$

Un solo paso de Gauss Seidel transforma las componentes v y w en v_{k+1} y w_{k+1}

$$\begin{aligned} 2v_{k+1} &= w_k + b_1 \\ 2w_{k+1} &= v_{k+1} + b_2 \end{aligned} \quad \text{ó} \quad \begin{bmatrix} 2 & 0 \\ -1 & 2 \end{bmatrix} x_{k+1} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} x_k + b$$

Nuevamente son decisivos los valores propios de $S^{-1} T$ y desde luego son fáciles de encontrar : son $1/4$ y 0 . El error está dividido cada vez en 4, así que un solo paso de Gauss-Seidel equivale a dos pasos de Jacobi. * Como ambos métodos requieren del mismo número de operaciones (usamos el nuevo valor en lugar del anterior y ahorramos en almacenamiento), es mejor el método de Gauss-Seidel.

* Esta regla es válida en multitud de aplicaciones, aun cuando es posible construir otros ejemplos donde Jacobi converja y Gauss-Seidel falle (ó viceversa). El caso simétrico es el más director : si todas las a_{ii} son mayores

que cero, entonces Gauss-Seidel converge si y sólo si A es positivamente definida.

Por otro lado Givens observó que, aunque tome una serie infinita de rotaciones lleva a A a la forma diagonal y, unas simples rotaciones de $1/2 (n-1)(n-2)$ pueden

llevar a A a la forma tridiagonal. Esto reduce el problema para encontrar los eigen-valores de matrices tridiagonales, y en algún caso la idea de Givens acorta el tiempo práctico de encontrar los valores característicos, por un factor de alrededor de 9 en la práctica. Pocos años después Householder introduce un nuevo método para tridiagonalizar una matriz simétrica, usando $n-2$ rotaciones en lugar de $1/2 (n-1)(n-2)$ esto reduce el tiempo por otro factor de 2 y efectivamente anula el método de Givens para negocios.

Muchos programas contemporáneos utilizan el método de Householder, pero difieren ampliamente en como son encontrados los eigen-valores de matrices tridiagonales. Consiguiendo los eigen-vectores es demasiado complicado, y la falta de conocimiento de como hacerlo es una razón para el uso ocasional continuo de los métodos de Jacobi.

Veamos brevemente la transformación de Householder, es una matriz de la forma.

$$H = I - 2 \frac{v v^T}{\|v\|^2}$$

A menudo se normaliza v para obtener un vector unitario $u = v / \|v\|$ y entonces H es $I - 2 u u^T$. En ambos casos H es tanto simétrica como ortogonal:

$$H^T H = (I - 2 u u^T)^T (I - 2 u u^T) =$$

$$= I - 4 uu^T + 4 uu^T uu^T = I$$

Así $H = H^T = H^{-1}$. En el caso complejo la matriz correspondiente $I - 2 uu^T$ es hermitiana y unitaria.

El plan de Householder es producir ceros con estas matrices y el buen éxito depende de la siguiente identidad.

Supongamos que z es el vector columna $(1, 0, \dots, 0)^T$ y $\sigma = \|x\|$ y $v = x + \sigma z$. Entonces $Hx = -\sigma z = (-\sigma, 0, \dots, 0)^T$

Comenzando con la primera columna de A y recordando que $U^{-1}AU$ final debe ser una forma tridiagonal ó de Hessenberg. Tendremos entonces que estarán involucradas solamente las $n - 1$ entradas debajo de la diagonal :

$$x = \begin{bmatrix} a_{21} \\ a_{31} \\ \vdots \\ a_{n1} \end{bmatrix} \quad z = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad Hx = \begin{bmatrix} -\sigma \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (a)$$

En este paso la matriz de Householder es solamente de orden $n - 1$, así que está sumegida en la esquina inferior derecha de una matriz U del tamaño original :

1

$$U_1 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & & & & \\ 0 & & & & \\ 0 & & & & \end{bmatrix} = U_1^{-1}, y,$$

$$U_1^{-1} AU_1 = \begin{bmatrix} a & * & * & * & * \\ -\sigma & & & & \\ 0 & & & & \\ 0 & & & & \end{bmatrix}$$

La matriz U_1 no altera la entrada a debido al 1 de la esquina superior izquierda; más aún, no toca los ceros que aparecen en la expresión (a). Se ha completado la primera parte y $U_1^{-1} AU_1$ tiene la primera columna deseada.

La segunda parte es análoga: x consta de las $n - 2$ entradas en la segunda columna, z es el vector coordenado unitario de longitud correspondiente y H es de orden $n - 2$. Cuando se inserta en U_2 , produce.

$$U_2 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & & & \\ 0 & 0 & & & \\ 0 & 0 & & & H_2 \end{bmatrix} = U_2^{-1}$$

$$U^{-1} (U^{-1} AU) U = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ 0 & * & * & * & * \\ 0 & 0 & * & * & * \\ 0 & 0 & * & * & * \end{bmatrix}$$

Finalmente U^{-1} se encargará de la tercera columna y se alcanzará la forma de Hessenberg para una matriz de 5 por 5. En general U es el producto de todas las matrices U_1, U_2, \dots, U_{n-2} y el número de operaciones requeridas es del orden de n^3 .

Ejemplo :

$$A = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 0 \end{bmatrix} \quad x = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad v = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad H = \begin{bmatrix} 0 & -1 \\ -1 & 0 \end{bmatrix}$$

el resultado es tridiagonal cuando se inserta H en U .

$$U = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 0 & -1 \\ 0 & -1 & 0 \end{bmatrix} \quad U^{-1} AU = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}$$

$U^{-1} AU$ es una matriz lista para obtener los valores propios.

C A P I T U L O

MII.

VALORES CARACTERISTICOS
DE
MATRICES NO SIMETRICAS

Un área de gran actividad en las décadas pasadas en la investigación del álgebra lineal computacional, ha sido el problema del valor característico para matrices no simétricas y solo prevalece un método de antes de la era de la computadora - - el método potencia - - y tiene solo aplicaciones limitadas hoy en día. Describamos brevemente el método de las potencias, cuyas propiedades de convergencia son fáciles de entender, comienza con una base en el principio de una ecuación de diferencias. Comienza con una elección inicial u_0 y forma sucesivamente

0

$$u_1 = Au_0, \quad u_2 = Au_1 \quad \text{y en general} \quad u_{k+1} = Au_k.$$

Cada paso es una multiplicación matriz vector y después de k

k - pasos se produce $u_k = A^k u_0$, aunque la matriz A

nunca aparecerá. El punto fundamental es que la multiplicación por la matriz A debería ser fácil de realizarse (si la matriz es grande, es mejor que sea rara ó hueca), ya que a menudo la convergencia al vector propio es muy lenta. Suponiendo que A tiene un conjunto completo de vectores propios x_1, \dots, x_n , el vector u_k estará dado

por la fórmula usual de una ecuación de diferencias:

$$u_k = c_1 \gamma_1^k x_1 + \dots + c_n \gamma_n^k x_n.$$

Imaginemos los valores propios numerados en orden creciente y supongamos que hay solo un valor propio mayor; no existe otro valor propio con la misma magnitud y γ_n no se repite.

n

Así, $\left| \gamma_1 \right| \leq \dots \leq \left| \gamma_{n-1} \right| < \left| \gamma_n \right|$. Entonces, en la

medida en que la elección inicial u_0 contenga alguna componente del vector propio x_n , de modo que $c_n \neq 0$, esta componente será dominante gradualmente:

$$\frac{u_k}{\gamma_n} = c_1 \begin{bmatrix} \gamma_1 \\ \vdots \\ \gamma_n \end{bmatrix} x_1 + \dots + c_{n-1} \begin{bmatrix} \gamma_{n-1} \\ \vdots \\ \gamma_n \end{bmatrix} x_{n-1} + c_n x_n$$

Los vectores u_k apuntan con más y más precisión hacia la dirección de x_n y el factor de convergencia es la razón

$$r = \left| \frac{\gamma_{n-1}}{\gamma_n} \right|$$

Es precisamente como la convergencia en un estado estacionario que estudiamos para los procesos de Markov, excepto que ahora el valor propio mayor de γ_n puede ser

igual a 1. De hecho no conocemos el factor de escala; de otra manera debe introducirse algún factor de escala, pues u_k puede crecer mucho ó poco, dependiendo de si $\left| \frac{\gamma_{n-1}}{\gamma_n} \right| > 1$

ó $\left| \frac{\gamma_{n-1}}{\gamma_n} \right| < 1$. Normalmente podemos dividir cada u_k por su primer componente y antes de realizar el paso siguiente;

con esta escala tan simple el método de las potencias se transforma en $u_{k+1} = Au_k / \alpha_k$ y converge a un múltiplo de

$$x_n^*$$

Ejemplo :

U_n teniendo al vector propio

$$\begin{bmatrix} 0.667 \\ 0.333 \end{bmatrix}$$

* los factores de escala α_k también convergen; tienden a γ_n

$$A = \begin{bmatrix} 0.9 & 0.2 \\ 0.1 & 0.8 \end{bmatrix} \quad \text{matriz de cambios de población}$$

$$u_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, u_1 = \begin{bmatrix} 0.9 \\ 0.1 \end{bmatrix}, u_2 = \begin{bmatrix} 0.83 \\ 0.17 \end{bmatrix}, u_3 = \begin{bmatrix} 0.781 \\ 0.219 \end{bmatrix}$$

$$u_4 = \begin{bmatrix} 0.747 \\ 0.253 \end{bmatrix}$$

La limitación más seria se ve claramente en este ejemplo : si r está cerca de 1, entonces la convergencia es muy lenta. En muchas aplicaciones $r > 0.9$, lo que significa que se tiene que hacer más de 20 interacciones para reducir

$\left(\frac{\gamma_k}{\gamma_1} \right)$ solo en un factor de 10. (El ejemplo tenía

$r = 0.7$ y fué aún muy lento). Claro que si r es igual a uno significa que $\left| \frac{\gamma_{n-1}}{\gamma_n} \right| = 1$, entonces no puede haber convergencia. Hay varias maneras de evitar esta limitación y mencionamos tres de ellas.

El método de las potencias por bloques, trabaja simultaneamente con varios vectores en lugar de con uno solo u_k .

El método de las potencias inversas.- trabaja con A^{-1} en lugar de A .

El método de las potencias inversa trasladado.- es el mejor

Supongamos que reemplazamos A por $A - \alpha I$. Entonces todos los valores propios γ_i se trasladan en la misma cantidad α

y el factor de convergencia para el método inverso cambiara

$\epsilon = 10^{-1}$, todos los valores característicos tienen modulo 0.1.

Afortunadamente, los valores característicos no son usualmente tan sensibles. De hecho diferentes valores característicos de una matriz A pueden diferir enormemente en su sensibilidad para perturbaciones de A, existen muchos resultados utiles. Son por lo general resultados a-posteriori, dando uno límites para los cambios en valores característicos como funciones de perturbaciones en una matriz e información acerca de otros valores característicos y vectores característicos.

La gran potencia de método de Jacobi para matrices simétricas, y las características extremadamente favorables del redondeo en matrices unitarias, llevaron a un deseo de usarlas para el problema de valores característicos en matrices no simétricas. El siguiente teorema básico es debido a Schur.

TEOREMA

Para una matriz arbitraria \bar{A} existe una matriz unitaria U tal que,

$$T = U^H A U \quad \text{triangular}$$

Donde U^H denota la transpuesta conjugada de U

Puesto que los valores característicos de A son los elementos diagonales de T, la espera ha sido para encontrar matrices unitarias las cuales llevan a A a una forma cercanamente triangular, y entonces deja que los elementos diagonales sirvan como valores característicos aproximados de A.

En la mayoría de métodos de atacar el problema de los valores característicos, el primer paso es condensar los datos para ahorrar tiempo y almacenaje adicionales. La forma condensada más aceptada por ahora es la matriz de Heseenberg, en la cual $a_{ij} = 0$ para $i > j$ (ó su transpuesta) es decir una matriz de la forma.

$$\begin{bmatrix}
 a & & & & & \\
 11 & a & & & & \\
 & 12 & a & & & \\
 0 & & 22 & a & & a \\
 & & & 23 & & 2n \\
 0 & 0 & & a & & a \\
 & & & 33 & & 3n \\
 \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
 \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
 0 & 0 & 0 & \cdot & \cdot & a \\
 & & & & & nn
 \end{bmatrix}$$

Es posible transformar A por congruencias ortogonales del tipo de Hausholder en una forma de Hessenberg, con solo errores de redondeo muy pequeños. Alguna condensación adicional (se dice en una forma tridiagonal) es sujeto de serias perdidas de dígitos. Una transformación a la forma de la matriz acompañante es particularmente desastrosa en la práctica, y normalmente requiere incrementos muy sustanciales en precisión para producir sucesivamente los valores característicos de A.

Como una alternativa, se puede transformar A a su forma de Hessenberg por eliminación gaussiana, con pivoteo parcial o una transformación similar.

La siguiente etapa en el problema de los eigen-valores no simétricos, es conseguir los valores característicos de una matriz de Hessenberg H. Una variedad de métodos han sido usados. Entre los cuales están.

- i).- Se puede examinar los ceros del $\det(H - zI)$ por el método de raíces encontradas, para z compleja. El método más satisfactorio parece ser el propuesto por Hyman en un número finito de programas, él hizo uso del método de Laguerre para encontrar los ceros de $f(z)$, seguido por una forma de supresión del cero.

Muchas recurrencias son usadas satisfactoriamente, para evaluar $f(z)$, $f'(z)$, y $f''(z)$, tan usadas en el proceso de Laguerre. Después que los valores característicos $\gamma_1, \dots, \gamma_r$ han sido encontrados, son suprimidos al aplicar el proceso de Laguerre a.

$$f(z) / \prod_{i=1}^r (\gamma_i - z)$$

ii).- El algoritmo de Ruthishauser, fué un importante desarrollo. Puesto que ahora ha sido suplantado por el algoritmo QR, el cual omitiremos mencionarlo.

iii).- Francis en Inglaterra y Kublanovskaja en la Union Soviética idearon el algoritmo muy interesante QR, y es considerado ampliamente como el algoritmo más satisfactorio para eigen-valores de matrices no-simétricas densas almacenadas.

El teorema básico es que una matriz real cuadrada arbitraria A puede ser factorizada en la forma $A = QR$, donde Q es una matriz ortogonal y R es una matriz triangular superior con todos los elementos diagonales r_{ij}

no negativos. Si A es no-singular, entonces ambas Q y R , son Únicas, de hecho, el cálculo esta hecho para construir una matriz Q ortogonal tal que $Q^T A = R$, donde R tiene las propiedades mencionadas.

Como una aparte, para A no-singular, el lector estará más familiarizado con la determinación del paso de una matriz triangular superior R^{-1} , tal que AR^{-1} , es una matriz ortogonal Q . Esta es la expresión matricial del proceso familiar del análisis de Gram-Schmidt.

Acaso sorprenderá al lector, que la matriz Q , que resulta del algoritmo de Gram-Schmidt, normalmente está lejos de ser ortogonal, porque por otro lado, si la misma Q está determinada tal que $Q^T A = R$, los errores de redondeo son muy pequeños.

El algoritmo básico QR, procede como sigue.

Sea $H = H_k$ una matriz de Hessenberg, para $k = 0, 1, 2, \dots$ el

factor H_k en la forma

$$H_k = Q_k R_k$$

y entonces la forma

$$H_{k+1} = R_k Q_k$$

Se prueba fácilmente que H_{k+1} , es también una matriz de Hessenberg.

Veamos el siguiente teorema básico.

TEOREMA.

Sea H , cuyos eigen-valores $\gamma_1, \dots, \gamma_n$, con.

$$|\gamma_1| < |\gamma_2| < \dots < |\gamma_n|$$

Entonces las matrices H_k convergen a una matriz triangular superior en la que los elementos de la diagonal son los valores característicos de H .

El caso más común donde el teorema anterior no satisface, es donde encontramos que H_k converge a una matriz triangular (esto significa que fuera de una matriz triangular todos los elementos de H_k tienden a cero, cuando k tiende a infinito, pero algunos elementos de la forma triangular pueden no converger). Además los grupos diagonales de 2 por 2 y 1 por 1 de la matriz H_k tienen valores característicos, los cuales en su totalidad convergen a los valores característicos de H .

Por simplicidad, consideremos una matriz H con valores característicos,

$$0 < \gamma_1 < \gamma_2 < \dots < \gamma_n ;$$

El método QR, para tales matrices converge con un error el cual es

$$O \left(\left(\frac{\gamma_1}{\gamma_2} \right)^k \right)$$

La convergencia podría ser más rápida si, en lugar de H tratáramos con la matriz $H - pI$, donde $0 < p < \gamma_1$, si p fuera prácticamente igual a γ_1 , la convergencia podría ser

extremadamente rápida. Han sido ideadas modificaciones del algoritmo, para el tan llamado "cambio original", el cual introduce a p cerca de γ_1 . Después que un valor característico γ_1 ha sido aislado, el método QR, puede

entonces ser aplicado a una matriz de $n - 1$, con valores característicos $\gamma_2, \dots, \gamma_n$. Nuevos cambios originales están entonces introducidos para sacar γ_2 tan rápidamente como sea posible etc. Con cambios originales bien ideados, el proceso ha sido observado para converger con un promedio de no menos de 2 pasos iterativos para valores característicos.

Algunas investigaciones profundizan en la investigación de cambios originales cuando alguno de los valores característicos son complejos y de módulos iguales, de lo cual no intentaremos dar ideas.

Normalmente los valores característicos son obtenidos para incrementar módulos.

Si H es una matriz simétrica banda, entonces el algoritmo QR preserva la banda de amplitud durante la iteración y es muy satisfactorio. En particular QR es un algoritmo posible para calcular valores característicos de una matriz

simétrica tridiagonal. Si H es una matriz banda no-simétrica, el algoritmo pierde las bandas cero sobre la diagonal.

Aún no hemos mencionado como conseguir los vectores característicos de una matriz de Hessenberg H . Este es el problema más difícil que mencionaremos.

El método prevaeciente es el de la iteración inversa. Los valores característicos están tomados, siendo conocidos, para algún valor característico γ , se selecciona un vector x de alguna manera. Entonces, se lleva a cabo

una iteración de la siguiente forma ;
para cada $k = 0, 1, 2 \dots$ encontramos x_{k+1}
para resolver el sistema :

$$(H - \gamma I) x_{k+1} = x_k$$

Se continua hasta que x_k es completamente grande. En algunos casos x_k es facilmente cercano al vector característico correspondiente a γ .

Si H es una matriz de Hessenberg real, γ pero es un valor característico complejo, se tiene que elegir entre uno u otro haciendo aritmética compleja ó algún proceso de selección juicioso con aritmética real.

Finalmente se transforma x_k al revéz del sistema coordenado original de A por deshacer las transformaciones de A a H .

Si alguno de los valores característicos de H son muy cercanos, los problemas reales comienzan. Un par de valores característicos cercanos pueden en casos afortunados tener distintos vectores característicos columna que están muy lejos de ser paralelos ésto representa una aproximación a un valor característico doble con divisores elementalmente lineales. Esta muy lejos como que un par de vectores característicos cercanos tendrán vectores característicos columna, estos también son paralelos. Esto representa una aproximación a los casos infinitamente más probables de un

valor característico doble con divisores no lineales.

En el caso precedente, no es difícil computar dos valores característicos que están lejos de ser paralelos. Es sólo llevar necesariamente la iteración

$$(H - \gamma I) x_{k+1} = x_k$$

con un x_0 diferente ó con valores escasamente diferentes de A .

En el último caso parece difícil obtener mucho de la iteración, pero un vector característico simple, perteneciente a A . Qué sería dado proximamente? En parte, no se sabe que el problema propuesto gustaría, pero no se sabe que es posible.

C O N C L U S I O N .

Los métodos de cálculo del álgebra lineal, se están moviendo en una etapa, donde existen métodos razonablemente satisfactorios para matrices densas almacenadas A. La principal excepción es el problema de encontrar vectores característicos con límites de error para matrices no simétricas. Los algoritmos han sido refinados varias veces y ya están siendo publicados.

Los usuarios casuales del álgebra lineal solo harán uso de tales algoritmos para sus problemas en particular. Los mejores algoritmos están escritos en algol 60. Aún cuando se pueda usar otro lenguaje, vale la pena poder leer este, con el fin justamente de poder leer estos algoritmos y adaptarlos a sus propios problemas.

Ningún método de resolver un problema computacional es realmente aprovechable para el usuario hasta que esté completamente descrito en un lenguaje computacional algebraico y hecho completamente formal. Hay aspectos indeterminados en todo algoritmo. Frecuentemente el provecho entero de un cierto proceso de computación consiste en el tratamiento de ciertas sutilezas, las cuales pueden difícilmente ser suspendidas hasta que están programados. Esta es la razón de porque un novato consultaría a uno u otro experto ó pasa grandes penalidades para retomar un algoritmo totalmente probado. Y esta es la razón de porque los profesionales difícilmente se concretarían en completar pruebas obvias de los algoritmos.

Hoy en día y desde 1979 está disponible un conjunto de programas en Fortran, llamado LINPACK, para resolver varios de los problemas básicos del álgebra lineal numérica. Y tiende a reemplazar gradualmente las miles de subrutinas existentes algunas buenas y otras no, que se ha escrito y utilizado en multitud de medios de computación. El LINPACK se diseñó independientemente de la máquina, eficiente y sencillo. El éxito que obtuvo el EISPACK para resolver problemas de valores propios estimuló el desarrollo de estos programas; gracias al EISPACK apreciamos el valor de un software matemático bien escrito.

B I B L I O G R A F I A.

- 1.-BAUER,F.L.- Optimally Scaled Matrices.-Numerische Mathematik.Vol.8,1963.- Pp.79 - 87.
- 2.-BURDEN,RICHARD L.& DOUGLAS FAIRES.-Análisis Numérico.- Grupo Editorial Iberoamericano.- 1985 .
- 3.-FORSYTHE.G.W.L.- Today's Computational Methods of Linear Algebra.- Siam Rev.9.- 1967.
- 4.-GOLUB,G.H.& J.H.WILKINSON.- Note on the Iterative Refinement of Least Squares Solution.- Numerische Mathematik ,1960 p.139 - 146 .
- 5.-NOBLE,BEN & DANIEL JAMES F.-Applied Linear Algebra.- Prentice Hall,Inc.-1977.
- 6.-STAIR,RALPH,JR.-Principles of Data Processing,Concepts Applications and Cases.-University of Florida,1981.
- 7.-STRANG,GILBERT.-Linear Algebra.-Academic Press,New York 1980.
- 8.-WILKINSON,J.H.-Rounding Errors in Algebraic Processes London,E.M.S.O.-1963.