

18
2 ej.



UNIVERSIDAD NACIONAL AUTONOMA
DE MEXICO

FACULTAD DE CIENCIAS

METODOS NUMERICOS PARA LA RESOLUCION
DE MINIMOS CUADRADOS POR EL METODO DE
GRAM-SCHMITO

TESIS PROFESIONAL
QUE PARA OBTENER EL TITULO DE:
A C T U A R I A
P R E S E N T A:
ANA MARIA HERNANDEZ MARTINEZ

Director de Tesis: José Luis Navarro Urrutia

TESIS CON
FALLA DE ORIGEN



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.

TESIS CON FALLA DE ORIGEN

INTRODUCCION

Un problema común que nos encontramos en el laboratorio de computación es el de encontrar las soluciones a problemas de Mínimos Cuadrados

Estos problemas surgen en una gran variedad de áreas y en diferentes contextos. Los problemas de mínimos cuadrados son en particular difíciles de resolver, dado que frecuentemente involucran gran cantidad de datos, y están mal condicionados por propia naturaleza. En este trabajo consideraremos métodos numéricos estables para manejar este problema. Nuestro instrumento básico es una descomposición de la matriz, basada en las transformaciones ortogonales de GRAM-SCHMIDT

Dado que un análisis general de la condición de los problemas de mínimos cuadrados está dada.

La influencia de los errores de redondeo se estudiará en detalle por una versión modificada de la ortogonalización de GRAM-SCHMIDT.

Para obtener una factorización $A = QR$ de una matriz A de $m \times n$ dada, donde R es una matriz triangular superior

y $Q^T Q = I$. Sea x el vector el cual minimiza $\|b - Ax\|_2$

y $r = b - Ax$. Esto muestra que si los productos internos son acumulados en doble precisión entonces, los errores de las x y r calculadas son menores que los errores resultantes de alguna perturbación inicial simultánea $\delta A, \delta b$ tal que :

$$\| \delta A \|_E \leq \| A \|_E \epsilon \quad \delta b \leq \| b \|_2 \epsilon \quad \epsilon \approx 2^{-24}$$

unidades de maquina.

La no reortogonalización es necesaria y el resultado es independiente de la estrategia de pivoteo usada.

CAPITULO I

DEFINICIONES BASICAS

Sea A una matriz real de $m \times n$, de rango n donde $m \geq n$ y b un vector real de $m \times 1$. Entonces existe una única x la cual resuelve el problema de mínimos cuadrados a minimizar.

$$\min_x \| b - Ax \|_2^2$$

Y es bien sabido que la solución x satisface la condición.

$$A^T (b - Ax) = 0$$

es decir el vector residual $r = b - Ax$, es ortogonal a las columnas de A , entonces se sigue que podemos calcular x de la ecuación normales.

$$A^T A x = A^T b$$

Ahora definimos siguiendo a Bauer, la condición de la matriz rectangular A .

Una definición natural de la condición de una matriz es el radio de su máximo y mínimo, con respecto a un par de normas.

Algunas veces, los renglones ó columnas de una matriz son grandes, solo para algunos factores, pero en cualquier caso escalar es computacionalmente trivial, por lo que el problema de minimizar la condición por escalonamiento óptimo tiene varias aplicaciones, lo que se desea es determinar el máximo y el mínimo para poder derivar los límites superior e inferior y para ello es necesario conocer la dirección de escalonamiento al mínimo ó último que se le acerque.

Ahora bien, si queremos mostrar condiciones para que una matriz esté optimamente escalonada.

El problema está completamente resuelto por la condición subordinada a un par de normas máximas y así cada una de

-1

las matrices A y A^{-1} tienen una distribución de signo en bloques cuadriláteras, (es decir son cuadradas), sin esta condición los últimos límites superior e inferior junto con factores apropiados de escalonamiento son obtenidos de la siguiente manera.

Sea B_1 y B_2 , dos espacios vectoriales normados, con normas

$$\|x\|_i \quad \text{para } x \in B_i \quad \text{con } i = 1, 2.$$

Para una función (mapa) $A : B_1 \rightarrow B_2$, el máximo, el cual es una norma, y el mínimo están definidos por

$$\max(A) = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} \quad \text{donde } \|x\| = 1$$

$$\min(A) = \inf_{x \neq 0} \frac{\|Ax\|}{\|x\|}$$

Entonces una condición subordinada a las normas se puede definir por

$$\text{Cond}(A) = \frac{\max(A)}{\min(A)}$$

si el $\min(A) > 0$, y si es caso entonces el $\text{cond}(A) = \infty$

el condicional puede ser totalmente finito para una A rectangular. Sin embargo $\text{cond}(A) = \max(A) \max(A)^{-1}$, si A^{-1} existe (asumiendo que A es una transformación rectangular).

Existen algunas propiedades básicas de la función condicional las cuales son válidas independientemente de las normas fundamentales.

1).- $\text{Cond}(A) > 1$

2).- si el $\text{cond}(A) = 1$, entonces $\|Ax\| = \delta \|x\|$ para una cierta $\delta > 0 \in \mathbb{R}$ donde $\delta \|A\|$ y $\|A\| = \sup\{\|Ax\|\}$ tal que $x \in A$.

3).- $\text{cond}(A^{-1}) = \text{cond}(A)$

4).- $\text{cond}(A \cdot B) = \text{cond}(A) \cdot \text{cond}(B)$

La condición 1) se cumple, ya que $\text{cond}(A) = \frac{\max(A)}{\min(A)}$, y

puesto que el infimo no excede al supremo esto significa que

$$\max(A) \cdot \min(A) = 1$$

$$\delta \quad \|Ax\| = \delta \|x\| \quad \text{y} \quad \|Ax\| = \delta$$

por lo que tenemos que la condición de una matriz rectangular A será:

$$\text{cond}(A) = \frac{\max_{\|x\|=1} \|Ax\|}{\min_{\|x\|=1} \|Ax\|}$$

En cuanto a la propiedad número 4 se cumple siempre y cuando

$$\|AB\| \leq \|A\| \|B\|$$

El número de condición para la norma L es denotado por $\kappa(A)$

En el capítulo VII se mostrará que bajo algunas restricciones $\kappa(A)$ puede ser considerada como un número condición aproximado, para el problema de calcular x . De otro modo.

$$\kappa(A^{-1}) = \kappa(A)$$

Esto muestra que en general usando aritmetica binaria de t-digitos no podremos obtener totalmente una solución aproximada a

$$A^{-1} Ax = A^{-1} b$$

t/2

Dado que $\kappa(A) \leq 2$

Ahora sea $B = AS$, donde S es cuadrada y no singular, (S es decir tiene inversa), entonces de:

$$A^T (b - Ax) = 0$$

se sigue que $B^T (b - Ax) = 0$ y así las ecuaciones

$$B^T Ax = B^T b$$

pueden ser empleadas en lugar de las ecuaciones normales. Una cuestión natural es preguntarse si podemos elegir B de tal manera que

$$\kappa(B^T A) = \kappa(A)$$

Esto es realmente posible, dado que las columnas de A son linealmente independientes y tenemos una factorización de

A , tal que $A = QR$. Donde $Q^T Q = I$ además de que Q es ortogonal R es triangular superior, por otro lado cada renglón de R y cada columna de Q están determinados de manera única, salvo un factor escalar de módulo 1

Eligiendo $B = AR^{-1} = Q$, la matriz $B^T A = Q^T QR = R$ es triangular.

De esta manera las ecuaciones

$$B^T Ax = B^T b$$

se resuelven fácilmente por sustitución hacia atrás. Además la condición

$$\kappa(B^T A) = \kappa(A)$$

se satisface dado que,

$$X(QR) = X(Q)X(R) = X(A)$$

Esta factorización se puede obtener de dos maneras diferentes y han sido propuestas para resolver los problemas de mínimos cuadrados. Un camino es eliminar los elementos de la subdiagonal de A por una serie de matrices hermitianas elementales.

$$(k) \quad H_k = \begin{pmatrix} 1 & & & \\ & 1 - 2w_k & & \\ & & \ddots & \\ & & & 1 \end{pmatrix}^T, \text{ donde } k = 1, 2, \dots, n.$$

$$(n) \quad \dots (2) \quad (1) \quad H_n \dots H_2 H_1 A = Q^T A = R$$

Recordando una Matriz Hermitiana, es aquella que es igual

a su transpuesta conjugada, se denota con A^H un ejemplo típico es :

$$A = \begin{bmatrix} 2 & 3 - 3i \\ 3 + 3i & 5 \end{bmatrix} = A^H$$

Notese que las entradas de la diagonal son números reales; no deben alterarse durante el proceso de conjugación. Cada entrada fuera de la diagonal corresponde a su reflejo al otro lado de la diagonal principal, y las dos son complejos conjugados entre sí.

En cada caso $a_{ij} = a_{ji}^*$, y este ejemplo ilustrará

claramente las cuatro propiedades de las matrices hermitianas.

se satisface dado que,

$$x(QR) = x(QR) = x(A)$$

Esta factorización se puede obtener de dos maneras diferentes y han sido propuestas para resolver los problemas de mínimos cuadrados. Un camino es eliminar los elementos de la subdiagonal de A por una serie de matrices hermitianas elementales.

$$H^{(k)} = \begin{pmatrix} I & & & \\ & 1 - 2w & & \\ & & w & \\ & & & \ddots \end{pmatrix}, \text{ donde } k = 1, 2, \dots, m.$$

$$H^{(n)} \dots H^{(2)} H^{(1)} A = Q^H A = R$$

Recordando una Matriz Hermitiana, es aquella que es igual

a su transpuesta conjugada, se denota con A^H un ejemplo típico es :

$$A = \begin{bmatrix} 2 & 3 - 3i \\ 3 + 3i & 5 \end{bmatrix} = A^H$$

Notese que las entradas de la diagonal son números reales; no deben alterarse durante el proceso de conjugación. Cada entrada fuera de la diagonal corresponde a su reflejo al otro lado de la diagonal principal, y las dos son complejos conjugados entre si.

En cada caso $a_{ij} = a_{ji}^*$, y este ejemplo ilustrará

claramente las cuatro propiedades de las matrices hermitianas.

Ahora al exponer esas cuatro propiedades, es necesario insistir otra vez en que igualmente se aplican a las matrices simétricas reales. Estas últimas representan un caso particular de matrices hermitianas, y el caso más importante.

Propiedad I.- Si $A^H = A$, entonces para todos los vectores

complejos x, x^H , Ax es real.

Cada entrada de A contribuye a la expresión $x^H Ax$.

$$x^H Ax = \begin{bmatrix} \bar{u} & \bar{v} \end{bmatrix} \begin{bmatrix} 2 & 3 - 3i \\ 3 + 3i & 5 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}$$

$= 2 \bar{u} u + 5 \bar{v} v + (3 - 3i) \bar{u} v + (3 + 3i) u \bar{v}$,
 cada uno de los términos diagonales es real ya que

$2 \bar{u} u = 2 |u|^2$ y $5 \bar{v} v = 5 |v|^2$. Los términos fuera de la diagonal son complejos conjugados entre sí, así que se combinan para dar el doble de la parte real de

$(3 - 3i) \bar{u} v$, por lo tanto, la expresión $x^H Ax$ es real.

Como una demostración general, podemos calcular $(x^H Ax)^H$. Obtendríamos el conjugado de la matriz $x^H Ax$ de 1 por 1, pero de hecho obtendríamos el mismo número:

$$(x^H Ax)^H = x^H A^H (x^H)^H = x^H Ax. \text{ Así el número debe ser real.}$$

Propiedad II.- Cada valor propio de una matriz hermitiana es real.

Demostración, suponemos que r es un valor propio y que x es su vector propio correspondiente distinto de cero: $Ax = rx$.

El truco es multiplicar por x^H , $x^H Ax = \tau x^H x$. Por la propiedad uno, el lado izquierdo es real, y como $x \neq 0$, el lado derecho $x^H x = \|x\|^2$ es real y positivo por lo tanto, τ debe ser real. En nuestro ejemplo,

$$A - \tau I = \begin{bmatrix} 2 - \tau & 3 - 3i \\ 3 + 3i & 5 - \tau \end{bmatrix} = \begin{bmatrix} 2 & \\ & 2 \end{bmatrix} = \tau - 7 \tau + 10 - (3 - 3i)^2$$

$$= \tau^2 - 7\tau - 8 = (\tau - 8)(\tau + 1)$$

Propiedad III.- Los valores propios de una matriz hermitiana, en caso de corresponder a valores propios diferentes, son ortogonales entre sí

Nuevamente la demostración comienza con la información dada, $Ax = \tau x$ y $Ay = \mu y$ y con $\tau \neq \mu$, y requiere una pequeña treta. Normalmente la transpuesta hermitiana de

$$Ax = \tau x \text{ es } x^H A = \tau x^H, \text{ pero como } A = A^H \text{ y } \tau \text{ es real,}$$

$$\text{por la propiedad II, esto es en realidad } x^H A = \tau x^H.$$

Multiplicando por y el lado derecho de esta ecuación y por

x el lado izquierdo de la otra, tenemos que :

$$x^H Ay = \tau x^H y \quad y^H Ax = \mu y^H x$$

Por lo tanto, $\tau x^H y = \mu x^H y$, y como $\mu \neq \tau$, concluimos que $x^H y = 0$: x es ortogonal a y . En el ejemplo, con $Ax=8x$ y $Ay = -y$, los vectores propios se calculan de la manera usual, de.

$$(A - 3I)x = \begin{bmatrix} -6 & 3 - 3i \\ 3 + 3i & -3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad x = \begin{bmatrix} 1 \\ 1+i \end{bmatrix}$$

$$(A + I)y = \begin{bmatrix} 3 & 3 - 3i \\ 3 + 3i & 0 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad y = \begin{bmatrix} 1 - i \\ -1 \end{bmatrix}$$

Estos dos vectores propios son ortogonales :

$$x^H y = \begin{bmatrix} 1 & 1 - i \end{bmatrix} \begin{bmatrix} 1 - i \\ -1 \end{bmatrix} = 0$$

Claramente cualesquiera múltiplos αx e βy servirán igualmente como vectores propios. Si la A original es real y simétrica, entonces sus valores propios y vectores propios son reales; además de que una matriz simétrica real puede diagonalizarse mediante una matriz ortogonal Q . En analogía con las matrices ortogonales reales que

satisfacen $Q^T Q = I$ y, por lo tanto, $Q^{-1} = Q^T$, la propiedad de tener columnas ortonormales se traduce en:

$$U^H U = \begin{bmatrix} x_1 \\ 1 \\ x_2 \\ 2 \\ \vdots \\ x_n \\ n \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \dots & x_n \\ 1 & 2 & \dots & n \end{bmatrix} = I \quad \text{o} \quad U^{-1} = U^H$$

Como en el caso complejo necesitamos de un nuevo nombre y símbolo para una matriz de columnas ortonormales : le llamaremos matriz unitaria y se denotará por U.

Esto nos conduce a la última propiedad especial de las matrices hermitianas.

H

Propiedad IV.- Si $A = A^H$, entonces existe una matriz diagonalizante que además es unitaria $S=U$. Sus columnas son ortonormales y

$$U^{-1} A U = U^H A U = \Lambda$$

Por lo tanto cualquier matriz hermitiana puede descomponerse en :

$$A = U \Lambda U^H = \tau_1 x_1 x_1^H + \tau_2 x_2 x_2^H + \dots + \tau_n x_n x_n^H$$

Esta descomposición se conoce como teorema espectral. Expresa que A es una combinación de las proyecciones unidimensionales $x_i x_i^H$ que son precisamente como las

T

proyecciones a a

Descomponen cualquier vector b en sus componentes

$p_i = x_i (x_i^H b)$ en las direcciones de los vectores propios

unitarios que constituyen un conjunto de ejes mutuamente perpendiculares. Estas proyecciones individuales p están ponderadas por las x_i y reagrupadas para formar :

$$a b = \tau_1 x_1 (x_1^H b) + \dots + \tau_n x_n (x_n^H b)$$

Si cada $r_i = 1$, hemos reagrupado b mismo, $A = U^H U$ es la

identidad. En cualquier caso podemos verificar los resultados anteriores directamente de la multiplicación de matrices :

$$\begin{aligned}
 ab &= U^H U^H b = \begin{bmatrix} x_1 & \dots & x_n \\ 1 & & \end{bmatrix} \begin{bmatrix} r_1 \\ \vdots \\ r_n \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \begin{bmatrix} b \\ \vdots \\ b \end{bmatrix} = \\
 &= \begin{bmatrix} x_1 & \dots & x_n \\ 1 & & \end{bmatrix} \begin{bmatrix} r_1 x_1 & \dots & r_1 x_n \\ \vdots & \ddots & \vdots \\ r_n x_1 & \dots & r_n x_n \end{bmatrix} \begin{bmatrix} b \\ \vdots \\ b \end{bmatrix} = r_1 x_1 x_1 b + \dots + r_n x_n x_n b
 \end{aligned}$$

Regresemos a nuestro ejemplo, ambos vectores propios tienen longitud $\sqrt{3}$ y una normalización produce la matriz que diagonaliza A :

$$U = \frac{1}{\sqrt{3}} = \begin{bmatrix} 1 & 1-i \\ 1+i & -1 \end{bmatrix} \quad U^{-1} A U = \begin{bmatrix} 2 & 0 \\ 0 & -1 \end{bmatrix}$$

entonces la descomposición espectral.

$$U^H U = r_1 x_1 x_1 + r_2 x_2 x_2 \quad \text{se transforma en :}$$

$$\frac{8}{3} \begin{bmatrix} 1 & \\ & \\ 1+i & \end{bmatrix} \begin{bmatrix} 1, 1-i \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 1-i & \\ & \\ -1 & \end{bmatrix} \begin{bmatrix} -1+i, -1 \end{bmatrix}$$

$$= \begin{bmatrix} 2 & 3-3i \\ 3+3i & 5 \end{bmatrix} = A$$

Una estimación de error detallada para el algoritmo $\hat{A} = QR$ está dada por Wilkinson y la aplicación a problemas de mínimos cuadrados está dada por Golub.

Un segundo camino es obtener la factorización por medio del proceso de ortogonalización de Gram - Schmidt a las columnas de A y esto siempre se ha sabido que tiene poca estabilidad numérica. En estas notas haremos uso de un versión ligeramente modificada del proceso de Gram - Schmidt, la cual es equivalente a la eliminación con combinación de renglones de peso, esto será aplicado a la solución de problemas de mínimos cuadrados lineales.

La propagación de errores de redondeo será estudiada en detalle y en particular una estimación de la desviación de ortogonalidad de la Q calculada y se obtendrá una cota de error para la solución calculada.

CAPITULO II

DESCRIPCIÓN DEL ALGORITMO

Se tratará de explicar la importancia de la ortogonalidad al resolver problemas de mínimos cuadrados, ya que está tiene una importancia y relación con la intuición que va más allá de los mínimos cuadrados.

Cada vez que se piensa en el plano $x-y$ y ó en el espacio tridimensional, la imaginación nos da la figura de un conjunto de ejes coordenados, nos proporciona un punto de referencia que llamamos origen, pero además los ejes coordenados son siempre ortogonales, al elegir un conjunto de ejes ortogonales lo que estamos haciendo es elegir una base para el plano $x-y$.

Si la idea es convertir cada construcción geométrica en un cálculo algebraico y es necesario una base ortogonal para que esos cálculos sean sencillos existe una especialización que hace esa base casi óptima.

Comencemos con un conjunto de vectores mutuamente ortogonales y normalizándolos los volveremos unitarios, esto es que cada v del conjunto se divide por su propia longitud y es remplazado por $v/\|v\|$. Con esto se logra que una base ortogonal se convierte en ortonormal.

Por lo que decimos que una base v_1, \dots, v_n es ortonormal si :

$$v_i \cdot v_j = \begin{cases} 0 & \text{si } i \neq j \text{ lo cual garantiza la} \\ & \text{ortogonalidad} \\ 1 & \text{si } i = j \text{ lo cual garantiza la} \\ & \text{normalización} \end{cases}$$

El ejemplo más importante es la base canónica donde los n vectores en el plano $x-y$ y en \mathbb{R}^3 además de ser perpendiculares los ejes son unitarios los v_i .

$$e_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad e_2 = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \quad \dots, \quad e_n = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

Puesto que está no es la única base ortonormal ya que se pueden rotar el conjunto de ejes sin cambiar los ángulos rectos que formen será necesario introducir matrices ortogonales ó de rotación. En el caso de los vectores canónicos e_i no se encuentran dentro del subespacio

U

n

considerado R , no será tan claro que podamos encontrar una base ortonormal, por lo que se verá que siempre es posible encontrar una de tales bases ya que siempre existe y que es posible construirla a partir de cualquier base de una manera sencilla. Este procedimiento que logra convertir cualquier conjunto de ejes oblicuos en uno perpendicular se le conoce con el nombre de GRAM-SCHMIDT.

Los temas básicos de este capítulo son :

- 1) - La solución a $Ax = b$ por mínimos cuadrados cuando las columnas de A son ortonormales.
- 2) - Definición y propiedades de las matrices ortogonales.
- 3) - El proceso de Gram - Schmidt y su interpretación como una nueva factorización de matrices.

Proyecciones y Mínimos Cuadrados : Caso (ortonormal)

Dada una matriz A de $m \times n$ y suponiendo que sus columnas son ortonormales se ve claramente que estas columnas son independientes por lo que podemos conocer la matriz de proyección sobre el espacio columna y la solución \bar{x} en

$$\text{mínimos cuadrados : } P = A (A^T A)^{-1} A^T \text{ y } \bar{x} = (A^T A)^{-1} A^T b.$$

Estas fórmulas son fácilmente expresadas ya que la matriz $A^t A$ es la identidad y son válidas en el caso de columnas ortonormales.

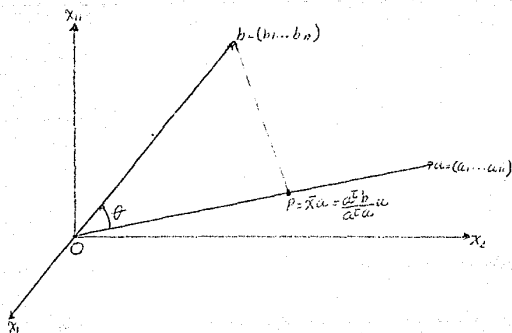
En el caso de columnas ortonormales para A se tiene:

$$\begin{bmatrix} a_1^t \\ \vdots \\ a_i^t \\ \vdots \\ a_n^t \end{bmatrix} \begin{bmatrix} a_1 & a_2 & \dots & a_n \\ \vdots & \vdots & \vdots & \vdots \\ a_1 & a_2 & \dots & a_n \\ \vdots & \vdots & \vdots & \vdots \\ a_1 & a_2 & \dots & a_n \end{bmatrix} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix} = I$$

Con esto logramos que el álgebra se simplifica ya que P y \bar{x} se cambian por:

$$P = A A^t \quad \text{y} \quad \bar{x} = A^t b.$$

Pero además es necesario simplificar la geometría y esto es posible gracias a que los ejes son perpendiculares, pues la proyección al espacio queda como proyección en cada eje.



La matriz de proyección es :

$$P = a_1 a_1^t + \dots + a_n a_n^t$$

$$P = AA^t = \begin{bmatrix} a_1 & a_2 & \dots & a_n \\ 1 & 2 & & n \end{bmatrix} \begin{bmatrix} a_1^t \\ a_2^t \\ \vdots \\ a_n^t \end{bmatrix} = \begin{bmatrix} a_1 a_1^t & a_1 a_2^t & \dots & a_1 a_n^t \\ a_2 a_1^t & a_2 a_2^t & \dots & a_2 a_n^t \\ \vdots & \vdots & \ddots & \vdots \\ a_n a_1^t & a_n a_2^t & \dots & a_n a_n^t \end{bmatrix}$$

Finalmente observamos que al hacer lo anterior el término $(A^t A)^{-1}$ desaparece en el cálculo de \bar{x} y P es la suma de las n -proyecciones por separado.

Hasta ahora se tienen a saber cinco ecuaciones que son :

- 1).- $Ax = b$ La ecuación dada probablemente inconsistente
- 2).- $A^t Ax = A^t b$ Las ecuaciones normales para \bar{x} .
- 3).- $p = A^t Ax$ La proyección de b sobre el espacio columna de A
- 4).- $P = A(A^t A)^{-1} A^t$ La matriz de proyección que nos da $p = P b$
- 5).- $\bar{x} = A^t b$ y $p = AA^t b = a_1 a_1^t b + \dots + a_n a_n^t b$ El caso en el que A tiene columnas ortonormales.

Matrices Ortogonales.

Se le llama a una matriz ortogonal a aquella que es

Lo anterior se sigue directamente de realizar $QQ^T = I$ puesto que lo que se está haciendo al realizar este producto es el producto interno de cada fila y como el resultado es la matriz identidad las filas son ortonormales no obstante las direcciones de las filas apunten en direcciones completamente diferentes de las columnas por lo que se puede afirmar que es suficiente con que las columnas sean perpendiculares para que automáticamente las filas también lo sean.

La Ortogonalización de Gram - Schmidt.

Hasta este momento se ha dicho que para resolver el problema de mínimos cuadrados $Ax = b$ se obtiene con mayores ventajas si las columnas son ortonormales ya que en cada caso se ha dicho que es fácilmente resuelto si las columnas son ortonormales por lo que a continuación daremos un método para hacerlas ortonormales. Este método es posible realizarlo antes de cualquier aplicación y la manera más fácil de visualizarlo es cuando están involucrados solo dos vectores.

Lo que se busca es producir de dos vectores independientes dados a y b dos perpendiculares v_1 y v_2 .

1 2

Lo primero que se debe hacer para lograr lo anterior es que el primer vector tenga la dirección de a esto es $v_1 = a$.

1

Nuestro problema es encontrar un segundo vector que sea perpendicular y para ello supongase que se tiene b un punto en el espacio n -dimensional y queremos encontrar la recta en la dirección del vector a busquemos entonces el punto p

1

de la recta más cercana a b por lo que la recta que une b con p es perpendicular al vector original a . Lo que se necesita es encontrar la proyección del punto p , este debe

ser un múltiplo es decir $p = \bar{x}a$ así vector a , por lo que el

problema es el cálculo de \bar{x} pero para realizar este cálculo

nos valdremos del hecho de que la recta desde b hasta el punto más cercano $p = \bar{x}a$ es perpendicular al vector a es

decir $x = a \cdot t$ por lo que finalmente se obtiene que

$p = \frac{a \cdot b}{a \cdot a}$ que es la proyección del punto b sobre la recta generada por el vector a . (ver gráfica al final del capítulo)

El vector $v = b - p$ es perpendicular al vector a por lo que el segundo eje irá en la dirección de:

$$v = b - p = b - \frac{a \cdot b}{a \cdot a} a = b - \frac{v \cdot b}{v \cdot v} v$$

y además se verifica que:

$$v \cdot v = v \cdot b - v \cdot b = 0$$

Para obtener directamente el proceso de Gram - Schmidt supongase que existe un tercer vector c independiente.

Substraemos las componentes de c en las dos direcciones v_1

y v_2 que ya hemos encontrado.

$$v_3 = c - \frac{v_1 \cdot c}{v_1 \cdot v_1} v_1 - \frac{v_2 \cdot c}{v_2 \cdot v_2} v_2$$

Nuevamente se tiene que v_3 es perpendicular a v_1 y v_2

pues lo que se está haciendo es sustraer de c sus componentes en el plano generado por a y b . La proyección de c sobre el plano es la suma de sus proyecciones en los ejes v_1 y v_2 además con esta construcción se garantiza que v_3 no puede ser el vector

cero puesto que está fuera que c estuviera dentro del plano formado por a y b contradiciendo la independencia lineal de a , b y c . Hasta ahora los vectores v_i son solo ortogonales

nos falta hacerlos ortonormales esto es unitarios, para lograrlo dividiremos cada uno por su longitud.

$$q_1 = \frac{v_1}{\|v_1\|}, \quad q_2 = \frac{v_2}{\|v_2\|}, \quad q_3 = \frac{v_3}{\|v_3\|}$$

Resumiendo el proceso anterior podemos decir que cualquier conjunto a_1, \dots, a_n de vectores independientes es posible hacerlo un conjunto v_1, \dots, v_n ortogonales mediante el método de Gram - Schmidt ; primero $v_1 = a_1$, después cada v_i es ortogonal a las v_1, \dots, v_{i-1} precedente por lo que

finalmente se obtiene que :

$$v_i = a_i - \frac{\langle a_i, v_1 \rangle}{\langle v_1, v_1 \rangle} v_1 - \dots - \frac{\langle a_i, v_{i-1} \rangle}{\langle v_{i-1}, v_{i-1} \rangle} v_{i-1} \quad (1)$$

Si observamos lo anterior se podrá ver que para cada v_i el subespacio que es generado por las a_1, \dots, a_i originales es generado también por v_1, \dots, v_i por lo que los vectores finales :

$$q_i = \frac{v_i}{\|v_i\|} \text{ son ortonormales.}$$

EJEMPLO : Supongase que los vectores originales son :

$$a_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \quad a_2 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \quad a_3 = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

entonces $a_1 = v_1$ y v_2 se calcula así :

$$v_2 = b - p = b - \frac{a_2 \cdot b}{a_2 \cdot a_2} a_2 = b - \frac{v_1 \cdot b}{v_1 \cdot v_1} v_1$$

$$\frac{a_2 \cdot v_1}{v_1 \cdot v_1} = \frac{1}{2}, \quad v_2 = a_2 - \frac{1}{2} v_1 = \begin{bmatrix} \frac{1}{2} \\ -\frac{1}{2} \\ 1 \end{bmatrix}$$

El tercer eje perpendicular viene de $1 =$:

$$v_3 = c - \frac{v_1 \cdot c}{v_1 \cdot v_1} v_1 - \frac{v_2 \cdot c}{v_2 \cdot v_2} v_2$$

Entonces :

$$\frac{a_3 \cdot v_1}{\|v_1\|} = \frac{1}{2}, \quad \frac{a_3 \cdot v_2}{\|v_2\|} = \frac{1}{3} = \frac{1}{3}$$

$$v_3 = a_3 - \frac{1}{2} v_1 - \frac{1}{3} v_2 = \begin{bmatrix} -2/3 \\ -2/3 \\ 2/3 \end{bmatrix}$$

Los vectores ortonormales finales son:

$$q_1 = \frac{v_1}{\|v_1\|} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \quad q_2 = \frac{v_2}{\|v_2\|} = \frac{1}{\sqrt{2/3}} \begin{bmatrix} 1/2 \\ -1/2 \\ 1 \end{bmatrix},$$

$$q_3 = \frac{v_3}{\|v_3\|} = \frac{1}{\sqrt{3/4}} \begin{bmatrix} -2/3 \\ 2/3 \\ 2/3 \end{bmatrix}$$

Lo cual implica que :

$$Q = \begin{bmatrix} \frac{\sqrt{2}}{2} & \frac{\sqrt{2/3}}{2} & \frac{2\sqrt{3/4}}{3} \\ \sqrt{2} & \frac{\sqrt{2/3}}{2} & \frac{2\sqrt{3/4}}{3} \\ 0 & \frac{\sqrt{2/3}}{2} & \frac{2\sqrt{3/4}}{3} \end{bmatrix}$$

y es posible recuperar las columnas originales de A a partir de los vectores q esto implica que deshaciendo las ecuaciones para v_i se tiene :

$$a_1 = v_1 \qquad a = \sqrt{2} \ q_1$$

$$a_2 = \frac{1}{2} v_1 + v_2 \qquad a = \frac{\sqrt{1/2}}{2} q_1 + \frac{\sqrt{3/2}}{2} q_2$$

$$a_3 = \frac{1}{3} v_1 + \frac{1}{3} v_2 + v_3 \qquad a = \frac{\sqrt{1/2}}{3} q_1 + \frac{\sqrt{1/6}}{2} q_2 + \frac{\sqrt{4/3}}{3} q_3$$

Y queda finalmente :

$$A = Q R$$

Donde R es una matriz triangular superior y Q una matriz ortogonal.

Por lo que finalmente se tiene :

$$\begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \frac{\sqrt{2}}{2} & \frac{\sqrt{2/3}}{2} & \frac{2\sqrt{3/4}}{3} \\ \sqrt{2} & \frac{\sqrt{2/3}}{2} & \frac{2\sqrt{3/4}}{3} \\ 0 & \frac{\sqrt{2/3}}{2} & \frac{2\sqrt{3/4}}{3} \end{bmatrix} \\
 * \begin{bmatrix} \sqrt{2} & \sqrt{1/2} & \sqrt{1/2} \\ 0 & \sqrt{3/2} & \sqrt{1/6} \\ 0 & 0 & \sqrt{4/3} \end{bmatrix}$$

El algoritmo de Gram - Schmidt es sencillo y directo, por lo que debe existir una forma sencilla de escribir el resultado final, a continuación se explicará la forma de hacerlo.

La situación es comparable con la eliminación gaussiana pues las operaciones se eligen de una manera natural conforme se van realizando. En este caso encontramos la manera correcta mediante la factorización $A = LU$ por lo que el algoritmo de Gram - Schmidt se registra mediante una factorización diferente de la matriz A .

El truco consiste en saber como es posible recobrar las columnas originales a_i a partir de los vectores finales

1

q. Si observamos el siguiente ejemplo y destacamos las
 1
 ecuaciones para v encontramos que :

$$a_1 = v_1$$

$$a_1 = \sqrt{2} q_1$$

$$a_2 = 1/2 v_1 + v_2$$

$$a_2 = \sqrt{1/2} q_1 + \sqrt{3/2} q_2$$

$$a_3 = 1/2 v_1 + v_2 + v_3$$

$$a_3 = \sqrt{1/2} q_1 + \sqrt{1/6} q_2 + \sqrt{4/3} q_3$$

Ahora si este conjunto de ecuaciones lo escribimos en forma matricial se tiene

$$\begin{bmatrix} a_1 & a_2 & a_3 \\ 1 & 2 & 3 \end{bmatrix} = \begin{bmatrix} q_1 & q_2 & q_3 \\ 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} \sqrt{2} & \sqrt{1/2} & \sqrt{1/2} \\ 0 & \sqrt{3/2} & \sqrt{1/6} \\ 0 & 0 & \sqrt{4/3} \end{bmatrix}$$

La matriz original A está factorizada en una matriz ortogonal Q multiplicada por una matriz triangular superior R, esto es A = QR.

Las columnas de Q son los vectores ortonormales que queremos : la ecuación.

$$v_i = a_i - \frac{v_1^T a_i}{v_1^T v_1} v_1 - \dots - \frac{v_{i-1}^T a_i}{v_{i-1}^T v_{i-1}} v_{i-1}$$

expresa a a como una combinación lineal de v , ..., v
 i i i

y solo se tiene que reemplazar v_1 por $\begin{pmatrix} v_1 \\ 1 \end{pmatrix} q_1$, v_2 por

$\begin{pmatrix} v_2 \\ 1 \end{pmatrix} q_2$ y así sucesivamente. Así v_i es una combinación

de q_1, \dots, q_{i-1} no involucrando las restantes q_j ; está

es la razón por la que R es triangular superior, además de que los coeficientes de esta combinación están en la i -ésima columna de R en particular las entradas de la diagonal son distintas de cero.

En la ecuación

$$v_i = a_{i1} \frac{v_1}{1} + \dots + a_{i,i-1} \frac{v_{i-1}}{1} + v_i$$

v_i se reemplaza por $\begin{pmatrix} v_i \\ 1 \end{pmatrix} q_i$ de modo que la entrada de

la diagonal es el coeficiente v_i diferente de cero.

R tiene entradas positivas en su diagonal y, por lo tanto está hecho que sea invertible, esto nos lleva al resultado más importante de este capítulo.

Cualquier matriz A con columnas linealmente independientes puede ser factorizada en un producto $A = QR$. Donde las columnas de Q son ortonormales y R es triangular superior e invertible. Si A es cuadrada sus factores Q y R también lo son y entonces Q será ortogonal.

Comenzando con $A = QR$ el problema de mínimos cuadrados $Ax = b$ es más fácil de resolver pues sabemos que :

$$\bar{x} = (A^T A)^{-1} A^T b = (R^T Q^T QR)^{-1} R^T Q^T b \text{ pero como}$$

$Q^t Q = I$ pues las columnas de Q son ortonormales tenemos que :

$$\bar{x} = (R \quad R) \begin{matrix} t & -1 & t & -1 & t \\ R & Q & b = R & Q & b \end{matrix}$$

aquí el cálculo de \bar{x} sólo requiere de la multiplicación de una matriz \times vector $Q^t b$ seguida de una sustitución regresiva en el sistema triangular $R \bar{x} = Q^t b$. El haber

ortogonalizado nos ayuda a no construir $A^t A$ y a no

resolver las ecuaciones normales $A^t A \bar{x} = A^t b$ dándonos además está mayor estabilidad numérica. (Como fué mostrado en el ejemplo de la pagina No 20 de este mismo capítulo).

Gram - Schmidt Modificado:

Nuestro instrumento básico es una descomposición de la matriz A basada en la transformación de Householder; lo que nos proporciona un " paso de preparación " para el algoritmo QR modificado.

Sea A una matriz real de $m \times n$, de rango r , y b un vector dado. Deseamos determinar un vector x tal que,

$$\| b - Ax \| = \min$$

donde $\| \cdot \|$ indica la norma euclidiana.

Si $m \geq n$ y $r < n$, entonces no hay solución única. Bajo estas condiciones tenemos que,

$$\| x \| = \min$$

La condición $\| b - Ax \| = \min$, es una condición muy natural para problemas estadísticos y numéricos. Si $m \geq n$ y $r = n$, entonces sabemos que x satisface

$$A^t Ax = A^t b$$

Desafortunadamente, la matriz $A^t A$, frecuentemente está mal condicionada e influenciada en la mayoría de los casos por errores de redondeo.

El siguiente ejemplo ilustra bien este caso. Supongamos :

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

Entonces

$${}^t A^{-1} A = \begin{bmatrix} 1 + \epsilon^2 & & & & \\ & 1 + \epsilon^2 & & & \\ & & 1 + \epsilon^2 & & \\ & & & 1 + \epsilon^2 & \\ & & & & 1 + \epsilon^2 \end{bmatrix}$$

donde claramente vemos que $\epsilon \neq 0$, el rango de ${}^t A^{-1} A$ es cinco, dado que, los eigen-valores de ${}^t A^{-1} A$ són :

$$5 + \epsilon^2, \epsilon^2, \epsilon^2, \epsilon^2, \epsilon^2$$

Ahora si los elementos de ${}^t A^{-1} A$ están calculados usando aritmética de doble precisión, entonces redondeamos a precisión simple exacta. Ahora sea n el número más grande en la computadora, tal que $fl(1.0 + n) = 1.0$ donde $fl(\cdot)$ indica la computación de punto flotante. Entonces si

$\epsilon < \sqrt{n} \div 2$, el rango de la representación computada será 1.

Consecuentemente, no importa cuán exacta sea la solución de la ecuación lineal. Es imposible resolver los ecuaciones normales :

$${}^t A A x = {}^t A b$$

Ahora vemos la descomposición de una matriz A .

Digamos que $m \geq n = r$, dado que la norma de un vector es unitariamente invariante tenemos que,

$$\| \| b - Ax \| \| = \| \| c - QAx \| \|$$

donde $c = Qb$, matriz ortogonal, elegimos A tal que :

$$QA = R = \begin{bmatrix} 1 & & & \\ & R & & \\ & & \dots & \\ & & & 0 \\ & & & & \dots & \end{bmatrix} \quad (m-n) \times n$$

Donde R es una matriz triangular superior, claramente

$$\bar{x} = R^{-1} \bar{c}$$

donde \bar{c} es el primero de los n componentes de C y por lo tanto

$$\| \| b - Ax \| \| = \left(\sum_{j=m+1}^n c_j^2 \right)^{\frac{1}{2}}$$

dado que R es una matriz triangular superior y $\bar{R}^t \bar{R} = A^t A$

donde $\bar{R}^t \bar{R}$ es la descomposición de Choleski de $A^t A$.

Ahora bien como el camino que hemos elegido es vía las transformaciones de Householder.

Sea $A^{(1)} = A$ y sean $A^{(2)}, A^{(3)}, \dots, A^{(n+1)}$ definidas como sigue:

$$A^{(k+1)} = P^{(k)} A^{(k)} \quad (k = 1, n)$$

(k) donde P es una matriz simétrica ortogonal de la forma :

$$P = I - \frac{2vv^T}{v^T v}$$

como v es ortogonal a v $\Rightarrow v^T v = 1$ a menudo se normaliza para tener una matriz ortonormal, entonces sea

$$w = \frac{v}{\|v\|} \quad \text{Por lo tanto resulta que :}$$

$$P = P^T = P^{-1} \quad \text{por ser simétrica y ortogonal así :}$$

$$PP = (I - 2ww^T)(I - 2ww^T) = (I - 4ww^T + 4w^T w w^T) = I$$

y resulta que P es una matriz hermitiana elemental y como una matriz hermitiana tiene componentes reales en la diagonal

\Rightarrow dado que cada entrada de P contribuye a la expresión

$$x^H A x \quad \text{Así cada uno de los términos diagonales es real así}$$

obteniendo el conjugado $\overline{x^H A x}$ se tiene que :

$$(x^H P x)^H = x^H P^H x = x^H P x = x^H P x$$

(\Leftarrow si hacemos $P x = y$ donde y es número real entonces, resulta que :

$$y = P x = x - 2w(w^H x)$$

$$y^H = x^H P^H = x^H P = y^H$$

$$y, \quad x^H y = x^H (P x)$$

esto es que dadas x , y si :

$$y = P x \Rightarrow x^H x = y^H y \quad \text{y por lo tanto } x^H x = y^H y, \quad \text{y es real}$$

entonces por la definición de $A^{(k+1)} = P^{(k)} A^{(k)}$ tenemos que podemos descomponer a:

ξ en $P^{(n)} P^{(n-1)} \dots P^{(1)}$, esto es en un producto de matrices hermitianas elementales y $\xi A = R$ entonces por lo que tenemos que:

$R = A^{(n-1)} \dots A^{(1)} P^{(n)} P^{(n-1)} \dots P^{(1)}$ donde $P^{(k)}$ es simétrica y ortogonal, por lo tanto ξ es una matriz ortogonal.

En el procedimiento de Gram - Schmidt modificado los elementos de R se calculan de un renglón en uno:

y definimos:

$$a_j^{(1)} = a_j \quad j = 1, 2, \dots, n$$

$$q_i = \frac{a_i^{(1)}}{d_i^{(1)}}$$

$$d_i^{(1)} = \sqrt{q_i^{(1)T} q_i^{(1)}}$$

$$r_{ij}^{(1)} = \frac{a_j^{(1)} q_i^{(1)T}}{d_i^{(1)}} \quad \text{donde } j = i+1, i+2, \dots, n$$

$$a_j^{(i+1)} = a_j^{(i)} - r_{ij}^{(i)} q_i^{(i)}$$

$\bar{A} = [\bar{A} : \bar{b}]$; y lo que deseamos obtener es la factorización:

$$\bar{A} = \bar{Q} \bar{R}$$

donde, \bar{A} es una matriz de $m \times (n+1)$:

$$r_{ij} = \frac{a_{ij}^{(k)} - q_{ij}^{(k)} r_{ij}^{(k-1)}}{d_{ij}^{(k)}} \text{ son los elementos de la matriz triangular superior } \bar{R}; \text{ y } q_{ij}^{(k)} = a_{ij}^{(k)}$$

son los elementos columna de la matriz \bar{Q} $m \times (n+1)$ y demostraremos como:

\bar{S} es una matriz diagonal cuyos elementos $(n+1) \times (n+1)$ son los elementos de la diagonal:

Según $d_{ij}^{(k)} = q_{ij}^{(k)}$ entonces la matriz $\bar{A}^{(k)}$ es de la forma:

$$\bar{A}^{(k)} = \begin{pmatrix} q_{11}^{(k)} & \dots & q_{1k}^{(k)} & \dots & a_{1n}^{(k)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ q_{i1}^{(k)} & \dots & q_{ik}^{(k)} & \dots & a_{in}^{(k)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ q_{m1}^{(k)} & \dots & q_{mk}^{(k)} & \dots & a_{mn}^{(k)} \end{pmatrix}$$

dado que $q_{ij}^{(k)} = a_{ij}^{(k)}$

y asumimos que :

$$\left. \begin{aligned} \begin{pmatrix} q \\ p \end{pmatrix}^t \begin{matrix} q & = & \delta & r \\ r & & pr & p \end{matrix} & \quad \begin{matrix} i \leq p, r \leq k-1 \\ \dots \\ \dots \\ \dots \end{matrix} \\ \begin{pmatrix} q \\ p \end{pmatrix}^t \begin{matrix} a \\ j \end{matrix}^{(k)} = 0 & \quad \begin{matrix} k \leq j \leq n \end{matrix} \end{aligned} \right\}$$

en el k -ésimo paso tomamos $q_k = a_k^{(k)}$ y calculamos :

$$d_k = \left\| \begin{matrix} q \\ k \end{matrix} \right\|_2^2$$

$$r_{kj} = \begin{pmatrix} q \\ k \end{pmatrix}^t a_j^{(k)} \div d_k$$

$$a_j^{(k+1)} = a_j^{(k)} - r_{kj} q_k \quad \text{donde } k+1 \leq j \leq n$$

y transformando el vector b del mismo modo nos da :

$$b = b^{(1)} \dots b^{(n+1)} = r$$

y en este paso k -ésimo calculamos para facilitar la notación :

$$y_k = \begin{pmatrix} q \\ k \end{pmatrix}^t b^{(k)} \div d_k$$

$$b^{(k+1)} = b^{(k)} - y_k q_k$$

tal que

$$R = \begin{bmatrix} R & | & c \\ \hline 0 & | & 1 \end{bmatrix} \begin{matrix} n \\ i \end{matrix}$$

$\underbrace{\hspace{10em}}_v$
 $\begin{matrix} n & 1 \end{matrix}$

$$D = \begin{bmatrix} D & | & 0 \\ \hline 0 & | & d \end{bmatrix} \begin{matrix} n \\ 1 \end{matrix}$$

$\underbrace{\hspace{10em}}_{n+1}$

entonces $\|Ax - b\|^2 = \|Q(Ax - b)\|^2 = \|D(Rx - c)\|^2 + d_{n+1}^2$

por lo tanto el mínimo valor de $\|Ax - b\|$ es $\sqrt{d_{n+1}^2}$

el cual se obtiene por el vector x que satisface $Rx = c$.

notece que $x(A)$ y que A (redondeada) tiene rango 1 tal que las ecuaciones normales son singulares. El procedimiento modificado tiene también la ventaja de seguir fácilmente una de las siguientes estrategias de pivoteo para ser usadas, podemos elegir como q la columna

de la matriz $A^{(k)}$ para la cual $a_{ij}^{(k)}$ donde $k \leq j \leq n$ es

máximizada, entonces los elementos de \bar{R}^k deben satisfacer

et $r < 1$ si deseamos expresar a b en las menores columnas jk de A como sea posible podemos elegir como q_k la columna para la cual P_{kj} sea maximizado donde:

$$P_{jk} = \left(\sum_j a_{jk}^{(k)} \right) b + a_{jk}^{(k)} \quad k \leq j \leq n$$

ambas $a_{jk}^{(k)}$ y p_{kj} son fácilmente calculadas por iteración.

está es probablemente la razón por la cual el procedimiento modificado se prefiere en la practica.

Demostraremos que en la práctica la superioridad del procedimiento modificado debido a otros factores independientes de la estrategia de pivoteo, esto es un argumento con resultados experimentales.

Los procedimientos sin embargo tienen propiedades numéricas totalmente diferentes cuando $n > 2$ si a está mal condicionada entonces con el clasico las columnas Q calculadas perderán su ortogonalidad, consecuentemente este procedimiento nunca deberá ser usado sin antes reortogonalizar lo cual incrementa grandemente la cantidad de evaluaciones.

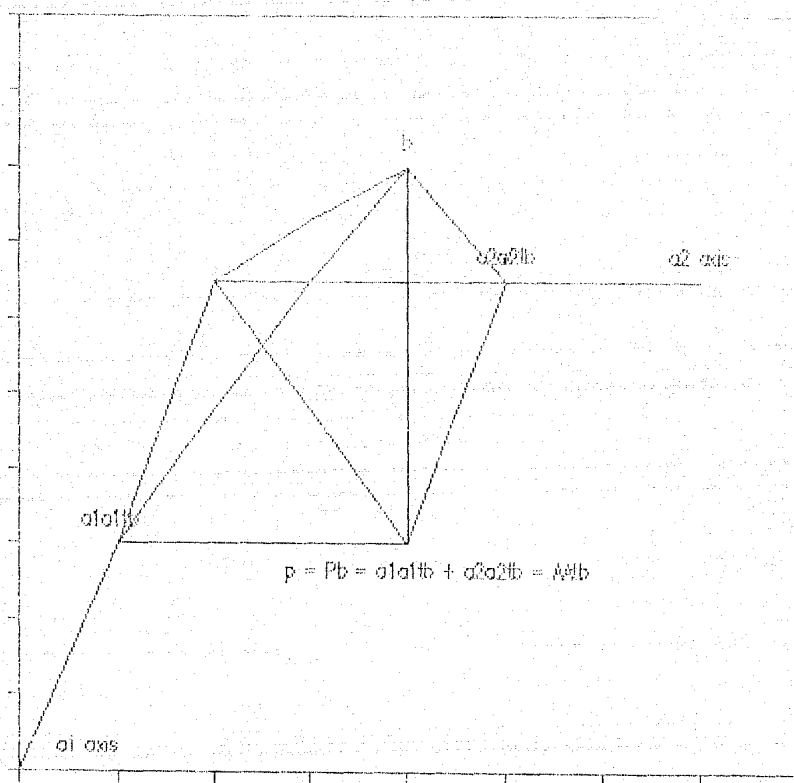
Como mostraremos la reortogonalización no es necesaria nunca cuando se usa el procedimiento modificado para resolver problemas de minimos cuadrados, el siguiente ejemplo dado nos ilustrará los diferentes aspectos de los dos procedimientos:

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

Sea ϵ una cantidad pequeña tal que debido al redondeo de la operación $1 + \epsilon$ de como resultado la unidad y puede fácilmente verificarse que si no existen otros errores de redondeo entonces la desviación máxima de ortogonalidad de las columnas de Q' calculadas está dado por :

$$\text{clasico} \quad \frac{\begin{matrix} t \\ (q'_3) & q'_2 \\ 3 & 2 \end{matrix}}{\begin{matrix} q'_3 & q'_2 \\ 3 & 2 & 2 & 2 \end{matrix}} = 4 \quad \text{modificado} \quad \frac{\begin{matrix} t \\ (q'_3) & q'_1 \\ 3 & 1 \end{matrix}}{\begin{matrix} q'_3 & q'_1 \\ 3 & 2 & 1 & 2 \end{matrix}} = 2/3$$

PROYECCION SOBRE UN PLANO



CAPITULO III

DEFINICIONES BASICAS DEL
ANALISIS DE ERROR

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 máquinas computadoras electrónicas, encarando la posibilidad de poder llevar a cabo computaciones que involucraban números muy grandes de operaciones aritméticas 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 más importante, y además, 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 calculada abandonando las más pesimistas conclusiones. Esto fue un estímulo para tratar de investigar métodos alternativos en la espera de que ellos no surgieran 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 romanas 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 fue 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 referiremos 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 , $/$, para denotar 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 pueda ser necesaria. Aunque se lo de 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

$$fl(a * b) = (a * b)(1 + \epsilon) \quad , \quad |\epsilon| < k \beta^{-t} \quad (1)$$

Donde el orden de k es la unidad y es independiente de a y b .

En algunas computadoras la relación

$$fi(a \pm b) = (a \pm b)(1 + \epsilon)$$

No siempre se toma como adición y sustracción y debemos conformarnos con resultado de la forma

$$fi(a + b) = a(1 + \epsilon_1) + b(1 + \epsilon_2), \quad |\epsilon_i| \leq \beta^{-i} \quad (2)$$

Donde en general no es posible tomar $\epsilon_1 = \epsilon_2$. Por consiguiente este RESULTADO MAS DÉBIL no tendría serias repercusiones para los resultados dados.

La relación

$$fi(a + b) = a(1 + \epsilon_1) + b(1 + \epsilon_2)$$

implica que el resultado calculado de cada operación individual tiene un error relativo bajo. Sin embargo hay un medio alternativo de interpretar la relación

$$fi(a * b) = (a * b)(1 + \epsilon)$$

la cual por consiguiente parece trivial, tiene sorprendentes implicaciones, las cuales nos pueden ilustrar en el caso de la multiplicación. La relación

$$fi(a * b) = (a * b)(1 + \epsilon)$$

establere que el producto calculado es el producto exacto de $a(1 + \epsilon)$ y b ó de a y, $b(1 + \epsilon)$, ó más simétricamente,

de $a(1 + \epsilon)$ y $b(1 + \epsilon)$. dado que ϵ es pequeña podemos decir que el resultado calculado 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 $a\sqrt{2}$ 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,

$$f_i(a \times b) = (a \times b) (1 + \epsilon) \quad \text{donde } |\epsilon| < k\delta$$

Usualmente es referido como un error de análisis inverso, y en general puede ser descrito en los siguientes términos.

Consideremos cualquier algoritmo calculado 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 calculados y no serán

exactos, pero en general, será un número infinito del conjunto de elementos $x_i (1 + \epsilon)$ para los cuales los y_j

calculados 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 son pequeños.

Encontramos que para las operaciones aritméticas básicas los valores naturales de ϵ son sugeridos por las

relaciones fundamentales satisfechas por los valores calculados; con muchos algoritmos de matrices las aplicaciones de la relación,

$$f_i(a \times b) = (a \times b) (1 + \epsilon) \quad \text{donde } |\epsilon| < k\delta$$

para todas las operaciones involucradas también lleva de una manera equitativamente natural a un conjunto de datos apropiados.

$$x_i (1 + \epsilon)$$

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

arriba para $k\epsilon$, 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 , es decir, un promedio de m operaciones por elemento de datos, un algoritmo para el cual podríamos establecer límites a-priori de $m\epsilon$, 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 calculada es decir, la distancia entre la y calculada y la y exacta correspondiente a la x dada. La estabilidad inversa no garantiza que la solución calculada 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

calculado 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, es 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 $k\epsilon_i$.

- b).- Para estimar los errores en la solución calculada 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 desastres aritméticos, tanto como los meros errores de redondeo.

Ahora volviendo a la ortogonalización.

Sean X y Y vectores de dimensión n donde $n \geq 2$ $\epsilon \leq 0.1$ implica que los siguientes límites de error para los productos internos de X y Y computados son válidos.

$$\left| f_1(x, y) - x^t y^t \right| \leq n \cdot 2^{-t} \left| x \right|_1 \left| y \right|_1 \quad (3.1)$$

donde,

$$\left| f_1(x, y) - x^t y^t \right| \leq 2 \left| x \right|_1 \left| y \right|_1 + 3/2 \cdot n 2^{-2t} \left| x \right|_1 \left| y \right|_1 \quad (3.2)$$

$$t_1 = t - \log_2(1.06), \quad 2t_2 = 2t - \log_2(1.06) \quad (3.3)$$

Aquí $\left| x \right|_1$ denota un vector con componentes $\left| x_i \right|$ y $f_1(\cdot)$ aplicado al caso cuando el producto interno es realizado en doble precisión y entonces redondeado, las cantidades calculadas serán en lo sucesivo denotadas por barras así, escribimos:

$$\bar{A}^{(k)} = \left(\bar{q}_1^{(k)}, \dots, \bar{q}_{k-1}^{(k)}, \bar{a}_k^{(k)}, \dots, \bar{a}_n^{(k)} \right)$$

Y las formulas

$$\bar{d}_k = \left\| \bar{q}_k^{(k)} \right\|_2, \quad \bar{r}_{kj}^{(k)} = \left(\bar{q}_k^{(k)}, \bar{a}_j^{(k)} \right) / \bar{d}_k$$

$$\bar{a}_j^{(k+1)} = \bar{a}_j^{(k)} - \bar{r}_{kj}^{(k)} \bar{q}_k^{(k)}, \quad k+1 \leq j \leq n$$

$$\bar{y}'_k = \left(\bar{q}_k^{(k)} \right) b_k / \bar{d}_k, \quad \bar{b}_k^{(k+1)} = \bar{b}_k^{(k)} - \bar{y}'_k \bar{q}_k^{(k)}$$

Se convierten en:

$$\bar{d}_k = r \left(\left\| \frac{\bar{q}'_k}{k} \right\| \right)^2, \quad \bar{r}'_{kj} = r \left(\left(\frac{\bar{q}'_k}{k} \right)^2 + \frac{\bar{d}_k}{k} \right)^{1/2} \quad (3.4)$$

$$\bar{a}_j^{(k+1)} = r \left(\bar{a}_j^{(k)} - \bar{r}'_{kj} \frac{\bar{q}'_k}{k} \right) \quad k+1 \leq j \leq n$$

$$y \quad (3.5)$$

$$\bar{y}'_k = r \left(\left(\frac{\bar{q}'_k}{k} \right)^2 + \frac{\bar{d}_k}{k} \right)^{1/2}, \quad \bar{b}_k^{(k+1)} = r \left(\bar{b}_k^{(k)} - \bar{y}'_k \frac{\bar{q}'_k}{k} \right)$$

este algoritmo no se rompe a menos que $\bar{d}_k = 0$ para alguna k esto puede pasar debido a la totalidad de errores de redondeo cuando $r(A) = n$ de acuerdo a

$$\left| r \left(\left\| \frac{x}{y} \right\| \right)^2 - \frac{\|x\|^2}{\|y\|^2} \right| \leq n - 2 \quad (3.1)$$

$$y \quad (3.2)$$

$$\left| r \left(\left\| \frac{x}{y} \right\| \right)^2 - \frac{\|x\|^2}{\|y\|^2} \right| \leq 2 \left(\frac{\|x\|}{\|y\|} \right)^2 + 3/2 - n^2 \quad (3.3)$$

$\bar{d}_k = 0$ implica $\bar{q}'_k = 0$ para este caso añadimos la regla trivial que cuando $\bar{q}'_k = 0$ tomamos :

$$\bar{r}'_{kj} = 0, \quad k+1 \leq j \leq n, \quad \bar{y}'_k = 0 \quad (3.6)$$

entonces \bar{R}' y \bar{Y}' están también únicamente definidas.

Para convenir en notación también introducimos las cantidades normalizadas :

$$\frac{\bar{q}}{k} = d \frac{-\bar{q}}{k}, \quad \frac{\bar{r}}{k_j} = d \frac{\bar{r}}{k_j}, \quad \frac{\bar{y}}{k} = d \frac{\bar{y}}{k} \quad (3.7)$$

donde :

$$\frac{\bar{q}}{d} = \begin{cases} \left| \left| \frac{\bar{q}}{k} \right| \right|_2, & \bar{q} \neq 0 \\ \left| \left| \frac{\bar{q}}{k} \right| \right|_2, & \bar{q} = 0 \end{cases}$$

Notese que esas cantidades nunca son calculadas y así

$$\frac{\bar{r}}{d} = \begin{cases} \left| \left| \frac{\bar{r}}{k} \right| \right|_2, & \bar{r} \neq 0 \\ \left| \left| \frac{\bar{r}}{k} \right| \right|_2, & \bar{r} = 0 \end{cases}$$

Son relaciones exactas.

CAPITULO IV

ERRDRES EN UNA PROYECCION
ELEMENTAL

Si $\bar{q}_k \neq 0$ entonces, en el k -ésimo paso calculamos los

vectores $\bar{a}_j^{(k+1)}$ con $j = k+1, \dots, n$.

Como la proyección de $\bar{a}_j^{(k)}$ en el subespacio complementario a \bar{q}_k . Si esto se ejecuta sin errores de redondeo el resultado es:

$$\bar{a}_j^{(k+1)} = \left(I - \frac{\bar{q}_k \bar{q}_k^T}{\bar{q}_k^T \bar{q}_k} \right) \bar{a}_j^{(k)} = \bar{a}_j^{(k)} - r_{kj} \frac{\bar{q}_k}{\bar{q}_k^T \bar{q}_k}$$

donde r_{kj} es el multiplicador exacto correspondiente a la \bar{q}_k calculada y $\bar{a}_j^{(k)}$. Si se usa el multiplicador calculado

\bar{r}_{kj} la sustracción se lleva a cabo exactamente, entonces el resultado es:

$$\bar{a}_j^{(k+1)} = \bar{r}_{kj} \frac{\bar{q}_k}{\bar{q}_k^T \bar{q}_k} = \bar{a}_j^{(k)} - \bar{r}_{kj} \frac{\bar{q}_k}{\bar{q}_k^T \bar{q}_k}$$

y definimos los errores $\alpha_j^{(k)}$ y $\beta_j^{(k)}$ por:

$$\bar{a}_j^{(k+1)} = \bar{a}_j^{(k)} - \bar{r}_{kj} \frac{\bar{q}_k}{\bar{q}_k^T \bar{q}_k} + \alpha_j^{(k)}$$

$$\bar{a}_j^{(k+1)} = \left(I - \frac{\bar{q}_k \bar{q}_k^T}{\bar{q}_k^T \bar{q}_k} \right) \bar{a}_j^{(k)} + \beta_j^{(k)}$$

en el caso singular cuando $\bar{q} = 0$ estas relaciones se satisfacen con :

$$\frac{\bar{a}^{(k+1)}}{j} = \frac{\bar{a}^{(k)}}{j}$$

$$\bar{a}^{(k)} = \bar{a}^{(k)} = 0$$

En el caso no singular probaremos las siguientes cantidades las cuales son básicas para nuestro análisis.

$$\left\| \frac{\bar{a}^{(k)}}{j} \right\|_2 \leq 1.45 \cdot 2^{-t} \left\| \frac{\bar{a}^{(k)}}{j} \right\|_2$$

$$\left\| \frac{\bar{a}^{(k)}}{j} \right\|_2 \leq \begin{cases} 3.23 \cdot 2^{-t} \left\| \frac{\bar{a}^{(k)}}{j} \right\|_2 \\ (2m + 3) \cdot 2^{-t} \left\| \frac{\bar{a}^{(k)}}{j} \right\|_2 \end{cases}$$

donde $\left\| \frac{\bar{a}^{(k)}}{j} \right\|_2 \leq 1.45 \cdot 2^{-t} \left\| \frac{\bar{a}^{(k)}}{j} \right\|_2$

es válida también cuando se usa precisión simple. Para simplificar la notación omitiremos los índices j y k y escribimos.

$$\frac{\bar{a}^{(k+1)}}{j} = x$$

$$\frac{\bar{a}^{(k)}}{j} = y$$

Ahora de la definición de operación en punto flotante y de las cantidades normalizadas

$$\bar{q}_k = d_k^{-1} \bar{q}'_k,$$

$$\bar{r}_{kj} = d_k^{-1} \bar{r}'_{kj}$$

$$\bar{y}_k = d_k^{-1} \bar{y}'_k$$

$$\text{donde } d_k = \begin{cases} \left\| \bar{q}'_k \right\|_2, & \bar{q}'_k \neq 0 \\ 1, & \bar{q}'_k = 0 \end{cases}$$

se sigue que

$$Z = \begin{pmatrix} y_1 & -\bar{r}_1 x_1 \\ & (1 + E_1) \end{pmatrix} \begin{pmatrix} (1 + E_2) \\ & \end{pmatrix}$$

$$\left| E_1 \right| \leq 2^{-t}, \quad \left| E_2 \right| \leq 2^{-t}$$

usando esto para eliminar y_1 de la definición de $\hat{\Omega}$ tenemos

$$\hat{\Omega} = \frac{E_2}{1 + E_1} \left\| Z \right\|_2^{-1} \begin{pmatrix} -E_1 & \bar{r}_1 x_1 \\ & \end{pmatrix},$$

y de aquí dado que $\left\| x \right\|_2 = 1$

$$\left\| \frac{\hat{z}}{2} \right\|_2 \leq \frac{1}{1-\alpha} \left\| \frac{\hat{z}}{2} + 2^{-t} \bar{r} \right\|_2$$

inmediatamente tenemos de las relaciones

$$r = x^T y, \quad z = y - \bar{r} x + \hat{z}; \quad z = y - r x + \beta$$

$$\left\| \frac{\beta}{2} \right\|_2 \leq \left\| \frac{\hat{z}}{2} \right\|_2 + |r - \bar{r}|$$

Para estimar el error en el multiplicador \bar{r} usamos la definición de la operación de punto flotante y le restamos

$$\left| - \left| \left| (x^T y) - x^T y \right| \right| \leq 2^{-t} \left| x^T y \right| + 3/2 \cdot n2^{-t} \left| x \right| \left| y \right|$$

y la identidad

$$\frac{\left| \left| (x^T y) - x^T y \right| \right|}{\left| \left| (x^T x) - x^T x \right| \right|} = r = \frac{\left| \left| (x^T y) - x^T y \right| \right|}{\left| \left| (x^T x) - x^T x \right| \right|} = r \frac{\left| \left| (x^T x) - x^T x \right| \right|}{\left| \left| (x^T x) - x^T x \right| \right|}$$

si calculamos \bar{r}' como

$$\bar{d} = \left\| \frac{1}{2} \left((x')^T x' \right) \right\|_2$$

$$\bar{r}' = \left\| \frac{1}{2} \left((x')^T y \right) \right\|_2$$

acumulando productos internos en doble precisión y dividiendo \bar{d} entre la mantisa de doble precisión, después de redondear, entonces:

$$|\bar{r} - r| \leq \frac{2^{-t} (2 + 2/3 \cdot m \cdot 2^{-t}) |r| + 2/3 \cdot m \cdot 2^{-t} \left\| \left\| y \right\| \right\|_2}{1 - 2^{-t} (2 + 2/3 \cdot m \cdot 2^{-t})}$$

Usando precisión simple tenemos

$$|\bar{r} - r| \leq \frac{(m+1) 2^{-t} |r| + m (1 + 2^{-t}) 2^{-t} \left\| \left\| y \right\| \right\|_2}{1 - m \cdot 2^{-t}}$$

Para simplificar estas últimas desigualdades hacemos enseguida las siguientes observaciones razonables.

$$m \geq 2$$

$$2 (m+1) 2^{-t} < 0.01$$

entonces ciertamente tenemos.

$$|\bar{r} - r| < \begin{cases} (2.01 \cdot |r| + 0.01 \left\| \left\| y \right\| \right\|_2) 2^{-t} \\ ((m+1) \cdot |r| + m \left\| \left\| y \right\| \right\|_2) 2^{-t} \end{cases}$$

donde los hemos tomado en cuenta que de las desigualdades:

$$e^{\frac{-t}{2}} \leq 2(1 + e^{-t}) \leq 2e^{-t/2} < 0.01$$

el factor $1.01 e^{-t/2}$ es ahora suficientemente generoso para ϵ también para los

factores $1 - \epsilon e^{-t/2}$ y $1 + \epsilon e^{-t/2}$.

Dado que $(z - \beta)$ es ortogonal a x se sigue que:

$$\|z\|_2 \leq \|y\|_2 + \|r\|_2 + \|\beta\|_2$$

sustituyendo esto y

$$\|r\|_2 \leq \frac{e^{-t/2}}{1 - e^{-t/2}} \|z\|_2 + 2\|r\|_2$$

en

$$\|\beta\|_2 \leq \|r\|_2 + \|r\|_2$$

sustituyendo el valor $\|r\|_2$ nos queda

$$\|\beta\|_2 \leq \frac{e^{-t/2}}{1 - e^{-t/2}} \|z\|_2 + 2\|r\|_2 + \|r\|_2$$

ahora sustituimos, el valor de $\left\| Z \right\|_2$

$$\left\| \beta \right\|_2 \leq \frac{2^{-t}}{1-2^{-t}} \left(\left\| y \right\|_2^{2-2^{-t}} + \left\| \beta \right\|_2^{2^{-t}} \left(|r| + |r-\bar{r}| \right) \right)$$

resolviendo para $\left\| \beta \right\|_2$ y despues de algunos arreglos nos queda.

$$(1-2 \cdot 2^{-t}) \left\| \beta \right\|_2 \leq 2^{-t} \left(\left\| y \right\|_2^{2-2^{-t}} + |r| + |r-\bar{r}| \right)$$

de aqui usando.

$$| \bar{r} - r | < \begin{cases} (2.01 \cdot |r| + 0.01 \left\| y \right\|_2^{2^{-t}}) 2^{-t} \\ ((m+1) \cdot |r| + m \left\| y \right\|_2^{2^{-t}}) 2^{-t} \end{cases}$$

nos queda

$$(1-2 \cdot 2^{-t}) \left\| \beta \right\|_2 \leq \begin{cases} \left(\left\| y \right\|_2^{2-2^{-t}} + 3|r| + \dots \right. \\ \left. + 0.02 \left\| y \right\|_2^{2^{-t}} \right) 2^{-t} \\ \left(\left\| y \right\|_2^{2-2^{-t}} + (m+2)|r| + \dots \right. \\ \left. + m \left\| y \right\|_2^{2^{-t}} \right) 2^{-t} \end{cases}$$

Maximizamos sobre β , donde:

$$0 \leq \beta \leq 1 \quad \text{teniendo}$$

$$\left(\left\| y \right\|_2^2 - r \right) + k \left\| r \right\| \leq (1+k)^{\frac{1}{2}} \left\| y \right\|_2$$

y finalmente:

$$\left\| \beta \right\|_2 \leq \begin{cases} 3 \cdot 23 \cdot 2^{-t} \left\| \frac{1}{j} \right\|_2 & \text{si } (k) \\ (2m+3) \cdot 2^{-t1} \left\| \frac{1}{j} \right\|_2 & \text{si } (k) \end{cases}$$

se sigue de:

$$\left((1-2 \cdot 2^{-t}) \left\| \beta \right\|_2 \right) \leq \begin{cases} \left(\left(\left\| y \right\|_2^2 - r \right) + 3 \left\| r \right\| + \dots \right. \\ \left. + 0.02 \left\| y \right\|_2^2 \right) 2^{-t} \\ \left(\left(\left\| y \right\|_2^2 - r \right) + (m+2) \left\| r \right\| + \dots \right. \\ \left. + m \left\| y \right\|_2^2 \right) 2^{-t1} \end{cases}$$

al retornar

$$m \geq 2, \quad 2(m+1)2^{-t} < 0.01$$

Para la estimación de $\|\hat{\alpha}\|_2$ las relaciones:

$$\|\hat{\beta}\|_2 \leq \frac{2^{-t}}{1-2^{-t}} \|\beta\|_2 \leq 2^{-t} \|\beta\|_2$$

y

$$\|\hat{z}\|_2 \leq \left(\|\hat{y}\|_2 - r^2 \right)^{1/2} + \|\beta\|_2$$

nos llevan a

$$\|\hat{\alpha}\|_2 \leq \frac{2^{-t}}{1-2^{-t}} \left(\|\hat{y}\|_2 - r^2 \right)^{1/2} + \|\beta\|_2 + |r - \hat{r}|$$

Usando la relación :

$$\left(\|\hat{y}\|_2 - r^2 \right)^{1/2} + k \|r\| \leq (1+k^2)^{1/2} \|\hat{y}\|_2$$

para los primeros dos términos y los límites en precisión simple para los últimos dos; tenemos:

$$\|\hat{\alpha}\|_2 \leq \frac{2^{-t}}{1-2^{-t}} \left(2 + 4(m+1) \cdot 2^{-t1} \right) \|\hat{y}\|_2$$

y retomando las relaciones :

$$m \geq 2; \quad 2(m+1)2^{-t1} < 0.01$$

y nuevamente esto prueba que :

$$\left\| \begin{array}{c} \frac{a^{(k)}}{j} \\ \Omega \end{array} \right\|_2 \leq 1.45 \cdot 2^{-t} \left\| \begin{array}{c} \frac{a^{(k)}}{j} \\ \Omega \end{array} \right\|_2$$

se cumple.

En un cálculo exacto, ciertamente se tiene :

$$\left\| \begin{array}{c} \frac{a^{(k+1)}}{j} \\ \Omega \end{array} \right\|_2 \leq \left\| \begin{array}{c} \frac{a^{(k)}}{j} \\ \Omega \end{array} \right\|_2$$

y esto no necesariamente tiene de ser cierto cuando son tomados en cuenta los errores de redondeo.

Estimando $\left\| \beta \right\|_2$ de la relación.

$$(1-2 \cdot 2^{-t}) \left\| \beta \right\|_2 \leq \left(\left\| y \right\|_2 \left(2^{-2} - r^2 \right)^{\frac{1}{2}} + 3 \left\| r \right\| + \dots \right. \\ \left. + 0.02 \left\| y \right\|_2 \right) 2^{-t} \\ \left(\left\| y \right\|_2 \left(2^{-2} - r^2 \right)^{\frac{1}{2}} + (m+2) \left\| r \right\| + \dots \right. \\ \left. + m \left\| y \right\|_2 \right) 2^{-t}$$

en

$$\left\| z \right\|_2 \leq \left(\left\| y \right\|_2 \left(2^{-2} - r^2 \right)^{\frac{1}{2}} + \left\| \beta \right\|_2 \right)$$

Y otra vez asociando el resultado por la derecha como una función de $\{r\}$, resulta después de algunos cálculos, que

$$\left\| \frac{\bar{a}^{(k+1)}}{j} \right\|_2 \leq \left\| \frac{\bar{a}^{(k)}}{j} \right\|_2 \begin{cases} (1 + 1.05 \cdot 2^{-k}) \\ (1 + 1.01(m+2) \cdot 2^{-k}) \end{cases}$$

Notese que estos límites solo pueden ser aprovechados

cuando $\frac{\bar{a}^{(k)}}{j}$ es también ortogonal a \bar{q} .

CAPITULO V

ERRORES EN LA FACTORIZACION

Usando el análisis de error de la proyección elemental podemos ahora derivar límites para los errores de la factorización de A y b.

Mientras podemos mostrar que el error $(QR - A)$ es pequeño independientemente de $x(A)$, esto quizá no se deba al error $(R - Q^{-1}A)$ esto nos lleva a introducir las matrices:

$$\tilde{Q} = \begin{pmatrix} \tilde{q}_1 & \tilde{q}_2 & \dots & \tilde{q}_n \end{pmatrix}, \quad \tilde{Q} = \begin{pmatrix} \tilde{q}_1 & \tilde{q}_2 & \dots & \tilde{q}_n \end{pmatrix}$$

$$\tilde{q}_k = (I - \tilde{q}_1 \tilde{q}_1^t) (I - \tilde{q}_2 \tilde{q}_2^t) \dots (I - \tilde{q}_{k-1} \tilde{q}_{k-1}^t) \tilde{q}_k,$$

$$\tilde{q}_k = (I - \tilde{q}_n \tilde{q}_n^t) (I - \tilde{q}_{n-1} \tilde{q}_{n-1}^t) \dots (I - \tilde{q}_{k-1} \tilde{q}_{k-1}^t) \tilde{q}_k, \quad (5.1)$$

cuando $\tilde{q}_j \neq 0$ la matriz $p_j^{(j)} = (I - \tilde{q}_j \tilde{q}_j^t)$ donde p es una

matriz unitaria, es la proyección del subespacio complementario a \tilde{q}_j , de otra manera es igual a la matriz unitaria y se sigue que :

$$\left\| p_j^{(j)} \right\|_2 = 1, \quad \left\| \tilde{q}_k \right\|_2 \leq 1, \quad \left\| \tilde{q}_k \right\|_2 \leq 1 \quad (5.2)$$

Primero probemos el lema simple pero importante
concerniente a \tilde{Q} y \tilde{Q}^{-1} .

LEMA (5.1)

Sea U la matriz triangular superior con elementos

$$u_{ij} = \begin{matrix} -t \\ q \end{matrix} \begin{matrix} - \\ q \end{matrix} \text{ donde } j > i, u_{ij} = 0 \text{ cuando } j \leq i \text{ entonces}$$

$$\tilde{Q} = Q(I + U), \quad \tilde{Q}^{-1} = Q^{-1}(I + U^{-1})$$

Primero probemos por inducción sobre k para $k = 1$

$$\left(\begin{matrix} -t \\ -q \end{matrix} \begin{matrix} - \\ q \end{matrix} \right)_1 \dots \left(\begin{matrix} -t \\ -q \end{matrix} \begin{matrix} - \\ q \end{matrix} \right)_k = \begin{matrix} -t \\ -q \end{matrix} \begin{matrix} - \\ q \end{matrix}_k - \dots - \begin{matrix} -t \\ -q \end{matrix} \begin{matrix} - \\ q \end{matrix}_1$$

esto es verdad para $k = 1$ y tenemos :

$$\begin{aligned} & \left(\begin{matrix} -t \\ -q \end{matrix} \begin{matrix} - \\ q \end{matrix} \right)_1 \dots \left(\begin{matrix} -t \\ -q \end{matrix} \begin{matrix} - \\ q \end{matrix} \right)_k \left(\begin{matrix} -t \\ -q \end{matrix} \begin{matrix} - \\ q \end{matrix} \right)_{k+1} = \\ & = \begin{matrix} -t \\ -q \end{matrix} \begin{matrix} - \\ q \end{matrix}_{k+1} + \left(\begin{matrix} -t \\ -q \end{matrix} \begin{matrix} - \\ q \end{matrix} \right)_1 \dots \left(\begin{matrix} -t \\ -q \end{matrix} \begin{matrix} - \\ q \end{matrix} \right)_k \end{aligned}$$

lo cual prueba la hipótesis.

Usando este resultado tenemos de la definición de q_k

$$q_k = \left(\begin{matrix} -t \\ -q \end{matrix} \begin{matrix} - \\ q \end{matrix} \right)_k - \dots - \begin{matrix} -t \\ -q \end{matrix} \begin{matrix} - \\ q \end{matrix}_1 q_k$$

y se sigue que para $k = 1, \dots, n$

$$\bar{q}_k = q_k + (q_k - q_{k-1})q_k + \dots + (q_k - q_{k-1})q_{k-1}$$

y está es la k -ésima columna de la igualdad $\bar{Q} = Q(I+U)$ y la primera parte del lema está probado, la segunda parte

es también obvia por la simetría de la definición \bar{Q} y \bar{Q} y se prueba igual.

Al tomar la transpuesta de la hipótesis de inducción usada en la prueba del lema obtenemos el siguiente

COLORARIO: si \bar{Q} está definida por (5.1)

$$\bar{Q} = (q_1, q_2, \dots, q_n), \quad \bar{Q} = (q_1, q_2, \dots, q_n)$$

$$\bar{q}_k = (I - q_1 - q_1^t) (I - q_2 - q_2^t) \dots (I - q_{k-1} - q_{k-1}^t) q_k$$

$$\bar{q}_k = (I - q_n - q_n^t) (I - q_{n-1} - q_{n-1}^t) \dots (I - q_{k-1} - q_{k-1}^t) q_k$$

tenemos la identidad:

$$I - \bar{Q} \bar{Q}^t = (I - q_n - q_n^t) \dots (I - q_2 - q_2^t) (I - q_1 - q_1^t)$$

y ahora definimos:

$$E_1 = \bar{Q} \bar{R} - A, \quad e_1 = \bar{Q} \bar{Y} + \bar{b}^{(n+1)} - b \quad (5.3)$$

$$\xi_2 = (I - \bar{Q} \bar{Q}^{-1}) A, \quad \eta_2 = (I - \bar{Q} \bar{Q}^{-1}) B - \bar{b}^{(n+1)} \quad (5.4)$$

$$\xi_3 = \bar{R} - \bar{Q} \bar{Q}^{-1} A, \quad \eta_3 = \bar{Y} - \bar{Q} \bar{Q}^{-1} b. \quad (5.5)$$

Ahora probaremos los siguientes estimadores los cuales son validos si los productos internos son acumulados en doble precision:

$$\left\| \begin{array}{c} \xi_1 \\ 1 \end{array} \right\|_2 \leq 1.5(n-1) \cdot 2^{-t} \left\| A \right\|_2 \left\| E \right\|_2 \quad (5.6)$$

$$\left\| \begin{array}{c} \eta_1 \\ 1 \end{array} \right\|_2 \leq 1.5n \cdot 2^{-t} \left\| b \right\|_2$$

$$\left\| \begin{array}{c} \xi_2 \\ 2 \end{array} \right\|_2 \leq 3.25(n-1) \cdot 2^{-t} \left\| A \right\|_2 \left\| E \right\|_2 \quad (5.7)$$

$$\left\| \begin{array}{c} \eta_2 \\ 2 \end{array} \right\|_2 \leq 3.25n \cdot 2^{-t} \left\| b \right\|_2$$

$$\left\| \begin{array}{c} \xi_3 \\ 3 \end{array} \right\|_2 \leq 1.9(n-1)^{\frac{1}{2}} n \cdot 2^{-t} \left\| A \right\|_2 \left\| E \right\|_2 \quad (5.8)$$

$$\left\| \begin{array}{c} \eta_3 \\ 3 \end{array} \right\|_2 \leq 1.9 n^{\frac{1}{2}} (n+1) \cdot 2^{-t} \left\| b \right\|_2$$

Sumando :

$$\frac{a_{(k+1)}}{j} = \frac{a_{(k)}}{j} - r \frac{q}{kj} + \delta \quad \text{para } k=1, \dots, j-1$$

y usando :

$$\frac{a_{(j)}}{j} = r \frac{q}{jj}, \quad \frac{a_{(1)}}{j} = a \quad \text{tenemos :}$$

$$\sum_{k=1}^j r \frac{q}{kj} - \frac{q}{k} - a = \sum_{k=1}^{j-1} \delta = \delta$$

de :

$$\left\| \frac{\delta}{j} \right\|_2 \leq 1.45 \cdot 2^{-t} \left\| \frac{a_{(k)}}{j} \right\|_2$$

$$\text{se sigue : } \left\| \frac{\delta}{j} \right\|_2 \leq 1.45 \cdot 2^{-t} \sum_{k=1}^{j-1} \left\| \frac{a_{(k)}}{j} \right\|_2$$

usando el límite de punto flotante fl () en doble precisión en :

$$\left\| \frac{a_{(k+1)}}{j} \right\|_2 < \left. \begin{array}{l} (1+1.05 \cdot 2^{-t}) \left\| \frac{a_{(k)}}{j} \right\|_2 \\ (1+1.01(m+2) \cdot 2^{-t+1}) \left\| \frac{a_{(k)}}{j} \right\|_2 \end{array} \right\}$$

tenemos : $\| \delta_j^{(k)} \|_2 < (1 + 1.05 \cdot 2^{-t})^{k-1} \| \delta_j \|_2$

y así recordando : $n \geq 2, \quad 2^{-(n+1)} < 0.01$

se tiene que :

$$\| \delta_j^{(k)} \|_2 < (1 + 1.05 \cdot 2^{-t})^{n-1} \| \delta_j \|_2 < 1.005 \| \delta_j \|_2 \quad (5.9)$$

esta desigualdad se toma complementariamente cuando se usa precisión simple si

$$n \geq 2, \quad 2^{-(n+1)} < 0.01$$

es sustituida por la condición más fuerte:

$$2n(n+2) \cdot 2^{-t} < 0.01 \quad \text{y entonces se tiene}$$

$$\| \delta_j \|_2 < 1.5 \cdot (j-1) \cdot 2^{-t} \| \delta_j \|_2 \quad (5.10)$$

y se tiene :

$$\| E_1 \|_E = \left[\sum_{j=1}^n \| \delta_j \|_2^2 \right]^{1/2} < 1.5 \cdot (n-1) \cdot 2^{-t} \left[\sum_{j=1}^n \| \delta_j \|_2^2 \right]^{1/2}$$

Lo cual prueba la primera parte de :

$$\| E_1 \|_E \leq 1.5(n-1) \cdot 2^{-t} \| A \|_E$$

dado que el lado derecho de b es tratado igual que las columnas de 0 , la segunda parte se sigue análogamente. Conseguirá resolveremos la ecuación diferencial.

$$\bar{a}_j^{(k+1)} = (I - q_j) \bar{a}_k^{(k)} + r_j^{(k)}$$

para $\bar{a}_j^{(k)}$ y tomamos para $k \leq j$

(5.11)

$$\bar{a}_j^{(k)} = \bar{a}_j^{(k)} + p_j^{(k+1)} \dots p_j^{(3)} p_j^{(2)} r_j^{(1)} + \dots + p_j^{(k+1)} r_j^{(k-2)} + r_j^{(k-1)}$$

donde:

$$p_j^{(k)} = p_j^{(k-1)} \dots p_j^{(2)} p_j^{(1)}$$

de aquí que $\bar{a}_j^{(k)}$ es el vector que obtenemos si ejecutamos la computación exactamente con la \bar{q}_k

Dado que la $\bar{q}_j^{(j)} = \bar{a}_j^{(j)}$ tenemos $(I - q_j) \bar{a}_j^{(j)} = 0$ consecuentemente

$$\bar{a}_j^{(k)} = \bar{a}_j^{(k)} + p_j^{(k+1)} \dots p_j^{(3)} p_j^{(2)} r_j^{(1)} + \dots + p_j^{(k+1)} r_j^{(k-2)} + r_j^{(k-1)}$$

toma también para $j+1 \leq k \leq n+1$ con:

$$\bar{a}_j^{(k)} = 0, \quad r_j^{(j)} = \dots = r_j^{(k-1)} = 0.$$

Tomando normas en :

$$\bar{a}_j^{(k)} = a_j^{(k)} + p_j^{(k+1)} + \dots + p_j^{(s)} + r_j^{(k+1)} + \dots + r_j^{(k-1)}$$

y usando :

$$\left\| \begin{pmatrix} p \\ \vdots \\ p^{(j)} \end{pmatrix} \right\|_2 = 1, \quad \left\| \begin{pmatrix} q \\ \vdots \\ q^{(k)} \end{pmatrix} \right\|_2 \leq 1, \quad \left\| \begin{pmatrix} r \\ \vdots \\ r^{(k)} \end{pmatrix} \right\|_2 \leq 1$$

llegamos a:

$$\left\| \begin{pmatrix} -a_j^{(k)} & -a_j^{(1)} \\ a_j & a_j \end{pmatrix} \right\|_2 \leq \sum_{i=1}^{s-1} \left\| \begin{pmatrix} r \\ \vdots \\ r^{(i)} \end{pmatrix} \right\|_2 = s = \min(j, k)$$

De esto se sigue usando el fl () de :

$$\left\| \begin{pmatrix} r \\ \vdots \\ r^{(k)} \end{pmatrix} \right\|_2 \leq \left. \begin{array}{l} 3.23 \cdot 2^{-t} \left\| \begin{pmatrix} a \\ \vdots \\ a^{(k)} \end{pmatrix} \right\|_2 \\ (2m+s) \cdot 2^{-t1} \left\| \begin{pmatrix} a \\ \vdots \\ a^{(k)} \end{pmatrix} \right\|_2 \end{array} \right\}$$

$$\left\| \begin{pmatrix} a_j^{(k)} - a_j^{(k)} \\ a_j - a_j \end{pmatrix} \right\|_2 \leq 3.23 \cdot 2^{-t} \sum_{i=1}^{s-1} \left\| \begin{pmatrix} a \\ \vdots \\ a^{(i)} \end{pmatrix} \right\|_2$$

$$\leq 3.25 (s-1) 2^{-t} \left\| \begin{pmatrix} a \\ \vdots \\ a \end{pmatrix} \right\|_2 \quad (5.12)$$

Usando el Corolario del Lema anterior se puede escribir para $k = n + 1$ como :

$$\left\| \begin{bmatrix} \tilde{a}_j^t \\ I - Q Q^t \end{bmatrix} a_j \right\|_2 \leq 3.25 (j-1) \cdot 2^{-t} \left\| a_j \right\|_2$$

Esto y una relación similar para b prueba:

$$\left\| \begin{matrix} E \\ 2 \end{matrix} \right\|_E \leq 3.25(n-1) \cdot 2^{-t} \left\| A \right\|_E$$

Usando la identidad :

$$\tilde{a}_{k,j}^t = \tilde{a}_{k,j}^t (I - q_{k-1}^t \tilde{a}_{k-1}^t) \dots (I - q_1^t \tilde{a}_1^t) a_{k,j}^t \quad (k)$$

podemos escribir el componente (k, j) de $\tilde{R} - Q^t A$

$$\tilde{r}_{k,j}^t - q_{k,j}^t a_{k,j}^t = (\tilde{r}_{k,j}^t - r_{k,j}^t) + q_{k,j}^t (a_{k,j}^t - a_{k,j}^t), \quad k < j$$

se sigue que :

$$\left\| \begin{matrix} \tilde{a}_{k,j}^t - a_{k,j}^t \\ a_j^t \end{matrix} \right\|_2 \leq 3.25 \cdot 2^{-t} \sum_{i=1}^{s-1} \left\| a_j^t \right\|_2$$

$$\leq 3.25 (s-1) 2^{-t} \left\| a_j \right\|_2$$

Y

$$\left\| \begin{array}{l} \bar{r} - r \\ \bar{y} - y \end{array} \right\|_2 < \left\{ \begin{array}{l} (2.01 \cdot |r| + 0.01) \left\| y \right\|_2^{-t} \\ (m+1) \cdot |r| + m \left\| y \right\|_2^{-t1} \end{array} \right.$$

Dada la norma $\| \cdot \|_3$ de la j -ésima columna en E , que está

así mismo limitada por lo cual finaliza la prueba de :

$$\left\| \begin{array}{l} E \\ 3 \end{array} \right\|_3 \leq 1.9(n-1) \cdot n \cdot 2^{-t} \left\| A \right\|_3 E$$

Notese que cuando A está mal condicionada la cancelación

ocurrirá tal que norma de $\left\| \begin{array}{l} a_j^{(k)} \\ j \end{array} \right\|_2 \leq \left\| a_j \right\|_2$

generalmente para un valor pequeño de k . Entonces al

$$\left\| \begin{array}{l} a_j^{(k)} \\ j \end{array} \right\|_2 < (1+1.05 \cdot n \cdot 2^{-t1}) \left\| a_j \right\|_2 < 1.005 \left\| a_j \right\|_2$$

estimado el cual hemos usado en :

$$\left\| \begin{array}{l} \delta_j \\ j \end{array} \right\|_2 < 1.5 \cdot (j-1) \cdot 2^{-t} \left\| a_j \right\|_2 \quad y$$

$$\left\| \frac{a_j^{(k)} - a_j^{(k-1)}}{j} \right\|_2 \leq 3.23 \cdot 2^{-t} \sum_{i=1}^{s-1} \left\| \frac{a_i^{(1)}}{j} \right\|_2$$

$$\leq 3.25 (s-1) 2^{-t} \left\| \frac{a_j}{j} \right\|_2$$

es muy bajo, y si restamos (5.6) de (5.3), esto es:

$$\left\| \frac{E}{1} \right\|_E \leq 1.5(n-1) \cdot 2^{-t} \left\| A \right\|_E \text{ menos}$$

$$\left\| \frac{E}{3} \right\|_E \leq 1.9(n-1) \cdot 2^{-t} \left\| A \right\|_E$$

el error será considerablemente sobreestimado.

Si usamos precisión simple de:

$$\left\| \frac{E}{1} \right\|_E \leq 1.5(n-1) \cdot 2^{-t} \left\| A \right\|_E$$

tomados sin cambio pero en

$$\left\| \frac{E}{2} \right\|_E \leq 3.25(n-1) \cdot 2^{-t} \left\| A \right\|_E \quad y$$

$$\left\| \frac{E}{3} \right\|_E \leq 1.9(n-1) \cdot 2^{-t} \left\| A \right\|_E$$

los límites deben ser incrementados por un factor de $(3/2m + 1)$.

A menos que una a_k ó b_k varíen, el vector q_i determinará completamente el total de los vectores p_i y q_i .

Consecuentemente q_i podría ser elegido tal que el algoritmo sea inestable en presencia del error de redondeo a no ser que bajo la reortogonalización, la línea sugerida por WILKINSON sean usadas. Esto es que el gran total de la ortogonalidad de los vectores generados por usar el método de GRAM-SCHMIDT para reortogonalizar cada uno de los vectores nuevamente generados p_i ó q_i a los vectores generados previamente

p_i ó q_i , respectivamente

Con el trabajo extra involucrado en la reortogonalización, el algoritmo es notablemente más lento, exceptuando la posibilidad de que A sea una matriz hueca.

CAPITULO VI

ORTOGONALIZACION
DE LOS
VECTORES CALCULADOS

Todos los límites de error dados en el capítulo 5 se aplican también en el caso singular, cuando $\bar{q}^k = 0$ para alguna k , es decir cuando el rango de $\bar{A} = \bar{Q}\bar{R}$ es ahora derivamos una condición suficiente para \bar{A} , para tener $r = n$.

$$E = \bar{Q}\bar{R} - A, \quad e = \bar{Q}\bar{y} + \bar{b} - b \quad (n+1)$$

podemos escribir:

$$\bar{A} - \bar{A} = (A + E) - (A + E) = R(I + F)R \quad (6.1)$$

donde:

$$F = (QER)^{-1} + QER^{-1} + (ER)^{-1}ER^{-1} \quad (6.2)$$

Tomando normas en lo anterior tenemos:

$$(6.3)$$

$$\|F\|_2 \leq 2 \left\| \begin{matrix} E \\ E \end{matrix} \right\|_2 \|R^{-1}\|_2 + \left\| \begin{matrix} E \\ E \end{matrix} \right\|_2^2 \|R^{-1}\|_2^2$$

y así la norma $\|F\|_2 < 1$ si:

$$\left\| \begin{matrix} E \\ E \end{matrix} \right\|_2 \|R\|_2^{-1} < \sqrt{2} - 1$$

usando :

$$\left\| \begin{matrix} E \\ 1 \end{matrix} \right\|_E \leq 1.5(n-1) \cdot 2^{-t} \left\| A \right\|_E$$

$$\left\| \begin{matrix} e \\ 1 \end{matrix} \right\|_2 \leq 1.5 n \cdot 2^{-t} \left\| b \right\|_2$$

se sigue que $\bar{A}^{-t} \bar{A}$ es no singular y \bar{A} tiene rango n si:

$$1.5(n-1) \cdot 2^{-t} \left\| A \right\|_E \left\| R^{-1} \right\|_2 < \sqrt{2} - 1 \quad (6.4)$$

Asumimos en lo siguiente que la fórmula anterior satisface la ortogonalidad de los vectores calculados

$\bar{q}_1, \bar{q}_2, \dots, \bar{q}_n$ pueden ser medidos por la $\| \cdot \|$ de

matriz $(I - \bar{Q}^{-t} \bar{Q})$ ahora :

$$I - \bar{Q}^{-t} \bar{Q} = - (U + U^t) \quad (6.5)$$

donde U es la matriz introducida en el LEMA 5.1.

Si resumimos :

$$\bar{a}_j^{(k+1)} = \bar{a}_j^{(k)} - \bar{r}_{kj} \bar{q}_k^{(k)} + \delta_j^{(k)}$$

para $k = i+1, i+2, \dots, j-1$ tenemos :

$$\bar{a}_j^{(i+1)} = \sum_{k=i+1}^j \bar{r}_{kj} \bar{q}_k^{(k)} - \sum_{k=i+1}^{j-1} \delta_j^{(k)} \quad (6.6)$$

$$\text{de : } \bar{a}_{j, (k+1)} = (1 - q) \bar{a}_{k, (k)} + r \bar{a}_{j, (k)}$$

$$\text{implica : } q \bar{a}_{i, j, (i+1)} = q \bar{r}_{i, j}$$

y de aquí multiplicando :

$$\bar{a}_{j, (i+1)} = \sum_{k=i+1}^j \bar{r}_{k, j} \bar{a}_{k, (i)} + \sum_{k=i+1}^{j-1} \delta_{k, j} \bar{a}_{i, (k)} \quad \text{por } \bar{a}_{i, (i)}$$

tenemos :

$$\bar{S}_{ij} = \sum_{k=i+1}^j \bar{r}_{k, j} \left(\frac{\bar{r}_{i, k}}{q} \right) = \bar{r}_{i, j} \left(r + \sum_{k=i+1}^{j-1} \delta_{k, j} \right) \quad (6.7)$$

donde \bar{S}_{ij} es la componente (i, j) de la matriz $\bar{S} = \bar{U}\bar{R}$

usando el fl. el límite de : $m \geq 2 \quad 2^{-(m+1)} 2^{-t_1} < 0.01$

$$\left| \bar{r} - r \right| < \left\{ \begin{array}{l} (2.01 - \left| r \right| + 0.01 \left\| y \right\|_2)^2 2^{-t} \\ (m+1) - \left| r \right| + m \left\| y \right\|_2^2 2^{-t_1} \end{array} \right.$$

y

$$\left\| \bar{a}_{j, (k)} \right\|_2 < (1 + 1.05 n 2^{-t_1}) \left\| a_{j, (k)} \right\|_2 < 1.006 \left\| a_{j, (k)} \right\|_2$$

tenemos :

$$\left\| \begin{matrix} S \\ j \end{matrix} \right\|_2 \leq \left(\sum_{i=1}^{j-1} 25 + \sum_{i=1}^{j-1} (j-i)^2 \right)^{-1/2} \left\| a \right\|_2$$

y dado que :

$$\sum_{i=1}^{j-1} (j-i + 7/6) \leq \sum_{i=1}^{n-1} (i+7/6) = 1/3(n-1)((n+3/2) + 11/6)$$

$$< 1.0001 \frac{2}{3} n(n+1)$$

Podemos limitar $\left\| \left\| L \right\|_2 \right\|_2$ de cada columna de S por :

$$0.87 \cdot n^{1/2} (n+1)2^{-t} \left\| a \right\|_2$$

despues de una estimación similar usando los límites para precisión simple obtenemos.

$$\left\| \left\| U R \right\|_E \right\|_E < \left\| \begin{matrix} 0.87 \cdot n^{1/2} (n+1)2^{-t} \left\| A \right\|_E \\ 0.87 \cdot n^{1/2} (n+1+2.5 \cdot m)2^{-t} \left\| A \right\|_E \end{matrix} \right\|_E$$

de :

$$I - \bar{Q}^t \bar{Q} = - (U + U^t)$$

se sigue que cuando se acumulan productos internos en doble precisión tenemos la estimación aposteriori : (6.8)

$$\left\| \left\| \begin{matrix} I - \bar{Q}^t \bar{Q} \\ U \end{matrix} \right\|_2 \right\|_2 \leq 2 \left\| U \right\|_2 \leq 1.74 n^{1/2} (n+1)2^{-t} \left\| A \right\|_E \left\| \bar{R}^{-1} \right\|_2$$

para estimar $\left\| \bar{R}^{-1} \right\|_2$ en :

$$\left\| \frac{1 - \bar{Q}^t \bar{Q}}{2} \right\|_2 \leq 2 \left\| \frac{U}{2} \right\|_2 \leq 1.74 n^{-1/2} (n+1) 2^{-t} \left\| A \right\|_2 \left\| E \right\|_2 \left\| \bar{R}^{-1} \right\|_2$$

usamos la identidad :

$$\bar{R}^t \bar{R} = (A + E) (A + E)^t - \bar{R} (U + U^t) \bar{R}$$

la cual puede ser derivada de :

$$E = \bar{Q} \bar{R} - A, \quad e = \bar{Q} \bar{Y} + B^{(n+1)} - B$$

y :

$$I - \bar{Q}^t \bar{Q} = - (U + U^t)$$

si escribimos :

$$\bar{R}^t \bar{R} = R (I + F) R \quad (6.9)$$

donde :

$$F = F - R (U \bar{R})^t \bar{R}^{-1} - R^{-1} \bar{R}^t (U \bar{R}) R^{-1}$$

y por :

$$\bar{R}^t \bar{R} = R (I + F) R$$

se sigue :

$$\left\| \bar{R}^{-1} \right\|_2 \leq 1 + \left\| \frac{F}{2} \right\|_2 \quad (5.10)$$

y de aquí usando :

$$\left\| \begin{matrix} F \\ 1 \end{matrix} \right\|_2 \leq 2 \left\| \begin{matrix} E \\ 1 \end{matrix} \right\|_E \left\| \begin{matrix} R^{-1} \\ 2 \end{matrix} \right\| + \left\| \begin{matrix} E \\ 1 \end{matrix} \right\|_E \left\| \begin{matrix} 2 \\ E \end{matrix} \right\| \left\| \begin{matrix} R^{-1} \\ 2 \end{matrix} \right\|$$

tenemos la desigualdad :

$$\left\| \begin{matrix} F \\ 2 \end{matrix} \right\|_2 \leq 2 \left\| \begin{matrix} E \\ 1 \end{matrix} \right\|_E \left\| \begin{matrix} R^{-1} \\ 2 \end{matrix} \right\| + \left\| \begin{matrix} E \\ 1 \end{matrix} \right\|_E \left\| \begin{matrix} 2 \\ E \end{matrix} \right\| \left\| \begin{matrix} R^{-1} \\ 2 \end{matrix} \right\| \\ + 2 \left\| \begin{matrix} U \bar{R} \\ E \end{matrix} \right\| \left\| \begin{matrix} R^{-1} \\ 2 \end{matrix} \right\| \left(1 + \left\| \begin{matrix} F \\ 2 \end{matrix} \right\|_2 \right)^{\frac{1}{2}}$$

usando los estimadores derivados para :

$$\left\| \begin{matrix} E \\ 1 \end{matrix} \right\|_E \text{ y } \left\| \begin{matrix} U \bar{R} \\ E \end{matrix} \right\| \text{ y la desigualdad :}$$

$(n-1) \leq 0.3 n^{\frac{1}{2}} \cdot n+1$ tenemos asumiendo que :

$$\left\| \begin{matrix} F \\ 2 \end{matrix} \right\|_2 \leq 2 \cdot 0.3 \cdot 1.5 c + (0.3 \cdot 1.5 c)^2 + \\ + 2 \cdot 0.87 c \left(1 + \left\| \begin{matrix} F \\ 2 \end{matrix} \right\|_2 \right)^{\frac{1}{2}}$$

donde :

$$c = N^{\frac{1}{2}} (n+1) \cdot 2^{-t} \left\| \begin{matrix} A \\ E \end{matrix} \right\| \left\| \begin{matrix} R^{-1} \\ 2 \end{matrix} \right\|$$

se sigue que : $\left\| \begin{matrix} F \\ 2 \end{matrix} \right\|_2 \leq b$ donde b es la raíz

cuadrada positiva de la ecuación

$$b = 0.9c + (0.45c)^2 + 1.74c(1+b)^2$$

dado que $b = 1$ correspondiente a $c = 12$ (3.419...) es realmente verificado que aquí si requerimos:

$$\beta = 3.42c = 3.42n \sqrt{(n+1)2} \left\| \begin{matrix} A \\ E \\ R^{-1} \end{matrix} \right\|_2 < 1 \quad (6.12)$$

la cual es más fuerte que la condición:

$$1.5(n-1)2^{-t} \left\| \begin{matrix} A \\ E \\ R^{-1} \end{matrix} \right\|_2 < \sqrt{2} - 1$$

Entonces $\left\| \begin{matrix} F \\ 2 \end{matrix} \right\|_2 < \beta$ y de:

$$\bar{R}^{-t} R = R(I+F)^t$$

$$\left\| \bar{R}^{-1} \right\|_2 \leq \frac{1}{1-\beta} \left\| R^{-1} \right\|_2 \quad (6.13)$$

está desigualdad se toma para precisión simple y definimos a β por:

$$\beta = 3.42n \sqrt{(n+1+2.5m)2^{-t}} \left\| \begin{matrix} A \\ E \\ R^{-1} \end{matrix} \right\|_2$$

finalmente de:

$$\left\| I - \bar{Q}^{-t} \bar{Q} \right\| \leq 2 \left\| U \right\| \leq 1.74n \sqrt{(n+1)2^{-t}} \left\| A \right\| \left\| \bar{R}^{-1} \right\|$$

y

$$\left\| \bar{R}^{-1} \right\|_2 \leq \frac{1}{1-\beta} \left\| R^{-1} \right\|_2$$

Tenemos un estimado a priori : (6.14)

$$\left\| \begin{pmatrix} 1 - \frac{1}{2} \beta \\ \beta \end{pmatrix} \right\| \approx \frac{1.74}{2} \frac{\beta}{1-\beta} n^{(n+1)} \cdot 2^{-n} \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\| + \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\| \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\|^{-1}$$

Extensivos experimentos computacionales reportan que cuando se usa el método modificado de GRAM-SMIT los resultados de pivotes se persiven pero son muy pequeños en la ortogonalidad de los vectores calculados.

Cuando $m = n$ un análisis previo indica que la solución calculada satisface la ecuación :

$$(1) \quad (1) \\ A + H \cdot x = b + e \quad (1)$$

Obviamente si A está mal condicionada $A + H$ podría ser singular mostramos que sí:

$$27 \frac{3/2}{n} \frac{-1}{2} \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\| \left\| A \right\| \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\| \left\| A \right\| \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\|^{-1} \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\|^{-p} \quad (p = 0)$$

$$\text{entonces : } \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\| \left\| x - x^{(1)} \right\| \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\| \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\| \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\|^{-p} \div \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\| \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\| \left\| \begin{pmatrix} 1 \\ \beta \end{pmatrix} \right\|^{-p-1}$$

donde : $x = A^{-1} b$ es la solución real.

de la ecuación :

$$(1) \quad (1) \\ A + H \cdot x = b + e$$

$$(1) \quad (1) \quad -1 \\ x = (A + H)^{-1} (b + e)$$

ESTA TESIS NO DEBE
SALIR DE LA BIBLIOTECA

(1)
Dado que $A + H$ es no singular, dando:

$$\left\| \begin{matrix} x \\ -x \end{matrix} \right\|_2^{(1)} = \left(\left\| A^{-1} \right\|_2 \left\| H \right\|_2^{(1)} \left\| x \right\|_2 + \left\| A^{-1} \right\|_2 \right) \div (1 - \left\| A^{-1} \right\|_2 \left\| H \right\|_2^{(1)})$$

Escribiendo a $12.5n 2^{-t}$ y recordando que $\bar{B} = \bar{U}$ tenemos

$$\left\| H \right\|_F^{(1)} = a \left\| A \right\|_F + (2^{-t} + 3/2 n 2^{-2t}) (1 + a) \left\| A \right\|_F$$

y dado que $\left\| A \right\|_2 \left\| A^{-1} \right\|_2 = 1$, se sigue de:

$$27 n^{3/2} 2^{-t} \left\| A \right\|_2 \left\| A^{-1} \right\|_2 2^{-p} \quad (p = 0)$$

(1)
que H ciertamente satisface el límite:

$$\left\| H \right\|_F^{(1)} = 13.5 n 2^{-t} \left\| A \right\|_F$$

usando las relaciones $\left\| A \right\|_F = n^{1/2} \left\| A \right\|_2$

(Una desigualdad muy débil para la mayoría de A).

$$\begin{aligned} \left\| \begin{matrix} (1) \\ H \end{matrix} \right\|_2 &= \left\| \begin{matrix} (1) \\ H \end{matrix} \right\|_F \quad \text{y} \quad \left\| e \right\|_2 = \\ &= 12.5 n^2 \left\| A \right\|_2 \times \left\| \right\|_2 \end{aligned}$$

tenemos finalmente :

$$\begin{aligned} \frac{\left\| \begin{matrix} (1) \\ x \end{matrix} \right\|_2 - x}{\left\| x \right\|_2} &= \frac{13.5 n^{3/2} \left\| A^{-1} \right\|_2 \left\| A \right\|_2 +}{1 - 13.5 n^{3/2} \left\| A^{-1} \right\|_2 \left\| A \right\|_2} \\ &+ \frac{12.5 n^{-t} \left\| A^{-1} \right\|_2 \left\| A \right\|_2}{1 - 13.5 n^{3/2} \left\| A^{-1} \right\|_2 \left\| A \right\|_2} \\ &= \frac{26 n^{3/2} \left\| A^{-1} \right\|_2 \left\| A \right\|_2}{1 - 13.5 n^{3/2} \left\| A^{-1} \right\|_2 \left\| A \right\|_2} \times (A) \end{aligned}$$

Donde $\times(A)$ es el número de condición espectral usual.
La condición.

$$27 n^{3/2} \left\| A \right\|_2 \left\| A^{-1} \right\|_2^{-p} \quad (p = 0)$$

ahora da el resultado requerido.

(1)

Dado que H está limitado uniformemente por todos los lados derechos está muestra que la A calculada no está demasiado mal condicionada, se puede tener la certeza de que la solución calculado tendrá algunas figuras correctas. El límite :

$$27 \frac{3/2}{n} \frac{-t}{2} \left\| A \right\|_2 \left\| A^{-1} \right\|_2 = 2^{-p} \quad (p = 0)$$

es extremadamente conservador. La experiencia sugiere que :

$$\left\| x - x^{(1)} \right\|_2 \leq \left\| x \right\|_2 \frac{-p}{2} \div (1 - 2^{-p-1})$$

es verdad cuando p está definido por algunas relaciones tales que :

$$\frac{1}{n} \frac{-t}{2} \left\| A \right\|_2 \left\| A^{-1} \right\|_2 = 2^{-p}$$

dado que A está demasiado mal condicionada el procedimiento del refinamiento iterativo es ciertamente un trabajo con cualquier lado derecho, dado que los errores cometidos en calcular el residual no tiene importancia.

Realmente si el residual estaba calculado exactamente, el análisis garantiza que la corrección s - ésima

(s)

calculada e - satisface la relación :

$$\left\| e - (x - x^{(s)}) \right\|_2 \leq \left\| x - x^{(s)} \right\|_2 = 2^{-p} \div (1 - 2^{-p-1})$$

si, además, ningún error fue cometido al sumar

(s) (s)
 ϵ a x tenemos:

$$\left\| \begin{matrix} x^{(s+1)} \\ x \end{matrix} \right\|_2 - \left\| x^{(s)} \right\|_2 = \left\| \begin{matrix} x^{(s+1)} \\ x \end{matrix} \right\|_2 - \left\| x^{(s)} \right\|_2 = \left\| \begin{matrix} x^{(s)} \\ x \end{matrix} \right\|_2 \left(2^{-p} + (1 - 2^{-p-1}) \right)$$

En el promedio $x^{(s)}$ se gana aproximadamente el mismo número de dígitos en cada iteración bajo este correcto para el trabajo preciso. Realmente por una modificación sutil del procedimiento es generalmente posible obtener

una $x^{(s)}$ de cualquier precisión usando la factorización original de A. En la práctica es adecuado calcular el residual $r^{(s)}$ - ésimo, usando acumulación de productos

internos, dado que si $r^{(s)}$ es el residual calculado y r el verdadero residual entonces $r^{(s)} = r + f^{(s)}$ donde:

$$\left\| r^{(s)} \right\|_2 = \left\| \begin{matrix} -t^{(s)} \\ 2 \cdot r \end{matrix} \right\|_2 + \frac{3/2}{2} n \frac{-2t}{2} \left\| A \right\|_2 \left\| x^{(s)} \right\|_2$$

Este resultado es consecuencia inmediata del análisis de error dado por: WILKINSON.

La acumulación de productos internos no es vital en cualquier otro paso del procedimiento. El fracaso de usarlo donde quiera meramente deja un límite de condición:

$$27 n \frac{3/2}{2} \frac{-t}{2} \left\| A \right\|_2 \left\| A^{-1} \right\|_2 \frac{-p}{2} \quad (p = 0)$$

reduciendo el rango de números condición para el cual el refinamiento iterativo proceda, y también baja la proporción de mejoramiento en algún grado. Dado que $x_{(s)}$

satisface el límite requerido y $x_{(s)}$ finalmente alcanzados por lo cual.

$$x - x_{(s)} \leq \epsilon \quad \text{y} \quad x = 2$$

a dicho paso x es "correcto para el trabajo preciso".

CAPITULO VII

CONDICION DEL PROBLEMA
DE
MINIMOS CUADRADOS

Antes de estimar los errores en la solución calculada, estudiaremos la sensibilidad para perturbaciones en A y b del problema de mínimos cuadrados, el cual siguiendo a GOLUB, escribimos.

$$\begin{bmatrix} I & A \\ A^t & 0 \end{bmatrix} \begin{bmatrix} r \\ x \end{bmatrix} = \begin{bmatrix} b \\ 0 \end{bmatrix} \quad (7.1)$$

Sea el sistema perturbado.

$$\begin{bmatrix} I & A + \delta A \\ (A + \delta A)^t & 0 \end{bmatrix} \begin{bmatrix} r + \delta r \\ x + \delta x \end{bmatrix} = \begin{bmatrix} b + \delta b \\ 0 \end{bmatrix} \quad (7.2)$$

y sea $A = V R V^t$ y $A + \delta A = \tilde{V} \tilde{R} \tilde{V}^t$

sea la factorización de la matriz perturbada, entonces :

$$\begin{bmatrix} \delta r \\ \delta x \end{bmatrix} = \begin{bmatrix} I & A \\ A^t & 0 \end{bmatrix}^{-1} \begin{bmatrix} 0 & \delta A \\ \delta A^t & 0 \end{bmatrix} \begin{bmatrix} r \\ x \end{bmatrix} + \begin{bmatrix} \delta b \\ 0 \end{bmatrix}$$

Notese que deseamos la solución de mínimos cuadrados del problema original sin turbaciones y frecuentemente se desearía comparar este vector con la solución final x. Este vector z, por supuesto sería calculado por los procedimientos de ortogonalización discutidos anteriormente.

Usando .

$$\begin{bmatrix} I & A + \delta A \\ A + \delta A & 0 \end{bmatrix} \begin{bmatrix} r + \delta r \\ x + \delta x \end{bmatrix} = \begin{bmatrix} b + \delta b \\ 0 \end{bmatrix}$$

la inversa puede ser calculada rápidamente y tenemos:

$$\begin{bmatrix} \delta r \\ R \delta x \end{bmatrix} = \begin{bmatrix} I - V V^t & V \\ V^t & -I \end{bmatrix} \begin{bmatrix} -\delta A x + \delta b \\ -R \delta A x + r \end{bmatrix} \quad (7.3)$$

Los eigenvalores de la matriz simétrica $(I - V V^t)$ son

todos iguales a 0 ó 1 de aquí $\left\| I - V V^t \right\| \leq 1$ y

$$\begin{bmatrix} \delta r \\ R \delta x \end{bmatrix} = \begin{bmatrix} I - V V^t & V \\ V^t & -I \end{bmatrix} \begin{bmatrix} -\delta A x + \delta b \\ -R \delta A x + r \end{bmatrix}$$

nos lleva a:

$$\left\| \begin{bmatrix} \delta r \\ R \delta x \end{bmatrix} \right\|_2 = \left\| \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \right\| \left\| \begin{bmatrix} R^{-1} \delta A x + r \\ \delta A x + \delta b \end{bmatrix} \right\|_2 \quad (7.4)$$

Para estimar $\|R^{-1}\|_2$ sea la factorización de A , $A = QR$

entonces podemos escribir :

$$R^{-1}R = A^{-1}A = (A + \delta A)^{-1}(A + \delta A) = R^{-1}(I + k)R$$

donde :

$$k = (Q^{-1} \delta A R^{-1}) + Q \delta A R^{-1} + (\delta A R^{-1})^{-1} \delta A R^{-1}$$

Ponemos :

$$\alpha = (\sqrt{2}+1) \left\| R^{-1} \right\|_2 \left\| \delta A \right\|_2 \quad (7.5)$$

y asumamos que $\alpha < 1$, entonces, es facilmente mostrar que :

$$\begin{aligned} \left\| k \right\|_2 &\leq 2 \left\| R^{-1} \right\|_2 \left\| \delta A \right\|_2 + \\ &+ \left\| R^{-1} \right\|_2^2 \left\| \delta A \right\|_2^2 < \alpha < 1, \end{aligned}$$

y de aqui :

$$\left\| R^{-1} \right\|_2 \leq \left\| R^{-1} \right\|_2 \left\| (I+k)^{-1} \right\|_2 \left\| R^{-1} \right\|_2 < \frac{1}{1-\alpha} \left\| R^{-1} \right\|_2$$

De aqui por : $x(R) = x(QR) = x(A)$

$$\left\| A \right\|_2 \left\| R^{-1} \right\|_2 \leq (1-\delta)^{-1/2} \left\| R \right\|_2 \left\| R^{-1} \right\|_2 = (1-\alpha)^{-1/2} x(A)$$

Dado que $\begin{pmatrix} \delta x \\ 2 \\ 2 \end{pmatrix} \leq \begin{pmatrix} R^{-1} \\ 2 \\ 2 \end{pmatrix} \begin{pmatrix} R \delta x \\ 2 \\ 2 \end{pmatrix}$ se sigue

ahora de

$$\begin{pmatrix} \delta r \\ 2 \\ 2 \\ R \delta x \\ 2 \\ 2 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} R^{-1} \\ 2 \\ 2 \\ \delta A \\ 2 \\ 2 \\ \times \\ 2 \\ 2 \\ \delta B \\ 2 \\ 2 \end{pmatrix}$$

que :

(7.6)

$$\begin{pmatrix} \delta r \\ 2 \\ 2 \\ A \\ 2 \\ 2 \\ \delta x \\ 2 \\ 2 \end{pmatrix} \leq \begin{bmatrix} \frac{\times(A)}{\sqrt{1-\alpha}} & 1 \\ \frac{\times(A)}{1-\alpha} & \frac{\times(A)}{\sqrt{1-\alpha}} \end{bmatrix}$$

$$* \left[\begin{array}{c} \begin{pmatrix} r \\ 2 \\ 2 \end{pmatrix} \frac{\begin{pmatrix} \delta A \\ 2 \\ 2 \end{pmatrix}}{\begin{pmatrix} A \\ 2 \\ 2 \end{pmatrix}} \\ \begin{pmatrix} A \\ 2 \\ 2 \end{pmatrix} \times \begin{pmatrix} \delta A \\ 2 \\ 2 \end{pmatrix} \frac{\begin{pmatrix} \delta A \\ 2 \\ 2 \end{pmatrix}}{\begin{pmatrix} A \\ 2 \\ 2 \end{pmatrix}} + \begin{pmatrix} \delta b \\ 2 \\ 2 \end{pmatrix} \end{array} \right]$$

Es posible construir ejemplos donde esos límites sean cercanamente obtenidos:

El resultado da :

$$\left\| \begin{array}{c} \left\| \delta r \right\|_2 \\ \left\| A \right\|_2 \left\| \delta x \right\|_2 \end{array} \right\|_2 \leq \left\| \begin{array}{cc} \frac{x(A)}{\sqrt{1-\alpha}} & 1 \\ \frac{x(A)}{1-\alpha} & \frac{x(A)}{\sqrt{1-\alpha}} \end{array} \right\|_2 *$$

$$* \left\| \begin{array}{c} \left\| r \right\|_2 \frac{\left\| \delta A \right\|_2}{\left\| A \right\|_2} \\ \left\| b \right\|_2 \frac{\left\| \delta A \right\|_2}{\left\| A \right\|_2} + \left\| \delta b \right\|_2 \end{array} \right\|_2$$

tiene las siguientes interpretaciones interesantes.

notese que : $x = R^{-1} b$ $Rx = R^{-1} Q Ax$

y así :

$$\begin{aligned} \left\| \begin{matrix} A \\ 2 \end{matrix} \right\| \times \left\| \begin{matrix} x \\ 2 \end{matrix} \right\| + \alpha(A) \left\| \begin{matrix} r \\ 2 \end{matrix} \right\| &\leq \alpha(A) \left(\left\| \begin{matrix} Ax \\ 2 \end{matrix} \right\| + \left\| \begin{matrix} r \\ 2 \end{matrix} \right\| \right) \\ &\leq \sqrt{2} \alpha(A) \left\| \begin{matrix} b \\ 2 \end{matrix} \right\| \end{aligned}$$

Entonces del primer renglon de :

$$\left\| \begin{matrix} \delta r \\ 2 \end{matrix} \right\| \leq \begin{bmatrix} \frac{\alpha(A)}{\sqrt{1-\alpha}} & 1 \\ \frac{2}{1-\alpha} & \frac{\alpha(A)}{\sqrt{1-\alpha}} \end{bmatrix} \left\| \begin{matrix} A \\ 2 \end{matrix} \right\| \left\| \begin{matrix} \delta x \\ 2 \end{matrix} \right\|$$

$$\left[\begin{array}{l} \left\| \begin{matrix} r \\ 2 \end{matrix} \right\| \frac{\left\| \delta A \right\|}{\left\| \begin{matrix} A \\ 2 \end{matrix} \right\|} \\ \left\| \begin{matrix} A \\ 2 \end{matrix} \right\| \left\| \begin{matrix} x \\ 2 \end{matrix} \right\| \frac{\left\| \delta A \right\|}{\left\| \begin{matrix} A \\ 2 \end{matrix} \right\|} + \left\| \begin{matrix} \delta b \\ 2 \end{matrix} \right\| \end{array} \right]$$

Se sigue que :

$$\frac{\left\| \begin{matrix} \delta r \\ \delta b \end{matrix} \right\|_2}{\left\| b \right\|_2} \leq \frac{\kappa(A)}{\sqrt{1-\alpha}} \frac{\left\| \delta A \right\|_2}{\left\| A \right\|_2} + \frac{\left\| \delta b \right\|_2}{\left\| b \right\|_2} \quad (7.7)$$

De aquí como número condición para la descomposición de b en componentes ortogonales Ax y r podemos tomar

$$\frac{1}{\sqrt{2(1-\alpha)}} \kappa(A)$$

Si asumimos que $\left\| x \right\|_2 \neq 0$, entonces podemos también derivar un número condición para la determinación de x. Del segundo renglon de :

$$\left\| \begin{matrix} \left\| \delta r \right\|_2 \\ \left\| A \right\|_2 \left\| \delta x \right\|_2 \end{matrix} \right\|_2 \leq \left[\begin{matrix} \frac{\kappa(A)}{\sqrt{1-\alpha}} & 1 \\ \frac{\kappa(A)}{1-\alpha} & \frac{\kappa(A)}{\sqrt{1-\alpha}} \end{matrix} \right] \quad *$$

$$* \left[\begin{matrix} \left\| r \right\|_2 & \frac{\left\| \delta A \right\|_2}{\left\| A \right\|_2} \\ \left\| A \right\|_2 \left\| x \right\|_2 & \frac{\left\| \delta A \right\|_2}{\left\| A \right\|_2} + \left\| \delta b \right\|_2 \end{matrix} \right]$$

se sigue :

$$\frac{\left\| \delta x \right\|_2}{\left\| x \right\|_2} \leq \frac{\kappa(A)}{\sqrt{1-\alpha}} \left(\frac{1 + \kappa(A)}{\sqrt{1-\alpha}} \frac{\left\| r \right\|_2}{\left\| A \right\|_2 \left\| x \right\|_2} \right) \frac{\left\| \delta A \right\|_2}{\left\| A \right\|_2} + \frac{\kappa(A)}{\sqrt{1-\alpha}} \frac{\left\| \delta b \right\|_2}{\left\| A \right\|_2 \left\| x \right\|_2}$$

Notese en particular la importancia de la proporción

$$\frac{\left\| r \right\|_2}{\left\| A \right\|_2 \left\| x \right\|_2}$$

si las ecuaciones son cercanamente compatibles en el sentido de que : $\left\| r \right\|_2 < \left\| A \right\|_2 \left\| x \right\|_2$, entonces la

situación es mucho muy parecida cuando $m > n$ como cuando $m = n$ y $\kappa(A)$ es un número condición aproximado para la determinación de κ .

Sin embargo, si $\left\| r \right\|_2$ es de la misma magnitud de

$\left\| A \right\|_2 \left\| x \right\|_2$ entonces el número condición (en un

caso mal condicionado) es más que $\kappa^2(A)$.

CAPITULO VIII

ERRORES EN LA SOLUCION
CALCULADA

Frecuentemente se desea determinar una x tal que $b - Ax$ sea minimizado, sujeto a la condición de que $Ax = b$.

Como la solución calculada tomamos :

$$\bar{r} = \bar{b}^{(n+1)}, \quad \bar{x} = f((\bar{R}^T)^{-1} \bar{y}^T) \quad (8.1)$$

Si definimos R tal que x es la solución exacta de :

$$\tilde{R}x = (\bar{R} + \delta \bar{R}) x = \bar{y}$$

entonces usando las definiciones :

$$E = (I - \bar{Q}\bar{Q}^T) A, \quad e = (I - \bar{Q}\bar{Q}^T) b - \bar{b}^{(n+1)}$$

$$E = \bar{R} - Q^T A, \quad \bar{e} = y - Q^T b$$

podemos introducir :

$$\bar{r} = \bar{b}^{(n+1)}, \quad \bar{x} = f((\bar{R}^T)^{-1} \bar{y}^T)$$

como :

$$\bar{r} = (I - \bar{Q}\bar{Q}^T) b - e, \quad \bar{x} = R^{-1} (Q^T b + e) \quad (8.2)$$

sustituyendo $b = r + Ax$ en :

$$\bar{r} = (I - \bar{Q}\bar{Q}^T) b - e, \quad \bar{x} = R^{-1} (Q^T b + e)$$

tenemos :

$$\begin{bmatrix} \bar{r} \\ \bar{x} \end{bmatrix} = \begin{bmatrix} I - \bar{Q}\bar{Q}^{-1} & (I - \bar{Q}\bar{Q}^{-1})A \\ \bar{R}^{-1} & \bar{R}^{-1}\bar{Q}^{-1}A \end{bmatrix} \begin{bmatrix} r \\ x \end{bmatrix} + \begin{bmatrix} -e_2 \\ \bar{R}^{-1}e_3 \end{bmatrix} \quad (8.3)$$

Ahora dado que $A^{-1}r = 0$ tenemos de :

$$E = \bar{Q}\bar{R}^{-1}A, \quad e_3 = \bar{Q}\bar{y} + b_3 - b_3$$

que :

$$\bar{Q}^{-1}r = \bar{R}^{-1}(A + E)^{-1}r = \bar{R}^{-1}E^{-1}r$$

y del LEMA (5.1) ver capítulo 5.
Se sigue que :

$$\bar{Q}^{-1} = (I + U)^{-1}\bar{Q}^{-1}, \quad \bar{Q} = \bar{Q}(I + U)^{-1}$$

Así los errores en la solución calculada pueden ser escritos como :

$$\begin{bmatrix} \bar{r} - r \\ \bar{x} - x \end{bmatrix} = \begin{bmatrix} -\bar{Q}^{-1}\bar{R}^{-1}E \\ \bar{R}^{-1}(I+U)^{-1}\bar{R}^{-1}E \end{bmatrix} (I - \bar{Q}^{-1}\bar{Q})A \begin{bmatrix} r \\ x \end{bmatrix} + \begin{bmatrix} -e_2 \\ \bar{R}^{-1}e_3 \end{bmatrix} \quad *$$

$$* \begin{bmatrix} r \\ x \end{bmatrix} + \begin{bmatrix} -e_2 \\ \bar{R}^{-1}e_3 \end{bmatrix} \quad (8.4)$$

Ahora estimaremos la norma 1 de las submatrices aparecidas en :

$$\begin{pmatrix} \tilde{r} \\ \tilde{x} \end{pmatrix} = \begin{pmatrix} \tilde{z} & -\tilde{t} & \tilde{t} \\ -\tilde{Q} & \tilde{R} & \tilde{E} \\ & & 1 \end{pmatrix} \begin{pmatrix} (I - \tilde{Q} \tilde{Q}^{-1}) A \\ \tilde{R}^{-1} \tilde{t} \\ \tilde{R}^{-1} (\tilde{Q} A - \tilde{R}) \end{pmatrix} + \begin{pmatrix} \tilde{e} \\ \tilde{z} \\ \tilde{e} \end{pmatrix}$$

$$\tilde{x} = \begin{pmatrix} \tilde{r} \\ \tilde{x} \end{pmatrix} + \begin{pmatrix} -\tilde{e} \\ \tilde{z} \\ \tilde{e} \end{pmatrix}$$

Asumimos para abreviar que el modo fl () de computación es usada en ambos, en la descomposición y en la sustitución hacia atrás. Los estimadores relevantes cuando usamos precisión simple no son muy diferentes y pueden ser derivados de manera similar. Dado que de acuerdo a :

$$\left\| \begin{pmatrix} p \\ 2 \end{pmatrix} \right\|_2 = 1 \quad \left\| \begin{pmatrix} \tilde{q} \\ 2 \end{pmatrix} \right\|_2 \leq 1 \quad \left\| \begin{pmatrix} \tilde{q} \\ 2 \end{pmatrix} \right\|_2 \leq 1$$

$$\left\| \begin{pmatrix} \tilde{z} \\ \tilde{Q} \\ \tilde{E} \end{pmatrix} \right\|_2 < n \quad \text{tenemos inmediatamente de :}$$

$$\left\| \begin{pmatrix} \tilde{E} \\ 1 \\ \tilde{E} \end{pmatrix} \right\|_2 \leq 1.5 (n-1) \cdot 2^{-t} \left\| \begin{pmatrix} A \\ \tilde{E} \end{pmatrix} \right\|_2$$

$$\left\| \begin{pmatrix} \tilde{e} \\ 1 \\ \tilde{e} \end{pmatrix} \right\|_2 \leq 1.5 n \cdot 2^{-t} \left\| \begin{pmatrix} b \\ 2 \end{pmatrix} \right\|_2$$

$$\left\| \begin{array}{c} \bar{r} \\ 2 \\ E \end{array} \right\|_2 \leq 3.25 (n-1) \cdot 2^{-t} \left\| \begin{array}{c} A \\ E \end{array} \right\|_F$$

$$\left\| \begin{array}{c} e \\ 2 \\ 2 \end{array} \right\|_2 \leq 3.25 n \cdot 2^{-t} \left\| \begin{array}{c} b \\ 2 \end{array} \right\|_2$$

(8.5)

$$\left\| \begin{array}{c} \bar{r} \\ -Q \bar{R}^{-1} E \\ 1 \\ 2 \end{array} \right\|_2 \leq 1.5 n^{1/2} (n-1) \cdot 2^{-t} \left\| \begin{array}{c} \bar{R}^{-1} \\ 2 \end{array} \right\|_2 \left\| \begin{array}{c} A \\ E \end{array} \right\|_F$$

(8.6)

$$\left\| (I - \bar{Q} \bar{Q}^t) A \right\|_2 \leq 3.25 \cdot (n-1) \cdot 2^{-t} \left\| \begin{array}{c} A \\ E \end{array} \right\|_F$$

Escribamos el residuo de dos submatrices en :

$$\left[\begin{array}{c} \bar{r} - r \\ x - x \end{array} \right] = \left[\begin{array}{c} -Q \bar{R}^{-1} E \\ 1 \\ (\bar{R} (I+U) R) E \end{array} \quad \begin{array}{c} (I - \bar{Q} \bar{Q}^t) A \\ \bar{R}^{-1} \bar{Q}^t (Q A - R) \end{array} \right]^*$$

$$* \left[\begin{array}{c} r \\ x \end{array} \right] + \left[\begin{array}{c} -e \\ 2 \\ \bar{R}^{-1} e \\ 3 \end{array} \right]$$

como

$$\left(\bar{R} (I+U) R \right) E = \bar{R} \bar{R}^{-1} \bar{R} (I+U) R E \quad (8.7)$$

$$\bar{R}^{-1} \bar{Q}^t (Q A - R) = \bar{R} \bar{R}^{-1} \bar{R}^{-1} \bar{Q}^t (Q A - R) - \bar{R}^{-1} \delta \bar{R} \quad (8.8)$$

Además de :

$$\bar{A}^{-t} \bar{A} = (A + E)^{-1} (A + E) = R^{-1} (I + F) R$$

$$I - \bar{Q}^{-t} \bar{Q} = -(U + U)$$

se sigue :

$$\bar{R} (I + U) \bar{R} = \bar{A}^{-t} \bar{A} - \bar{R}^{-t} U \bar{R} = R^{-1} (I + F) R$$

donde :

$$F = F - R^{-1} U \bar{R}^{-1} R^{-1}$$

Comparando está con :

$$\bar{R}^{-t} \bar{R} = R^{-1} (I + F) R$$

$$y \quad \left\| \bar{R}^{-t} \right\|_2^2 \leq \frac{1}{1 - \beta} \left\| R^{-1} \right\|_2^2$$

observamos que ciertamente $\left\| F \right\|_3 \leq \beta$

$$y \quad \left\| (\bar{R}^{-t} (I + U) \bar{R})^{-1} \right\|_2 \leq \frac{1}{1 - \beta} \left\| R^{-1} \right\|_2$$

si ponemos :

$$\left\| R^{-1} \bar{R} \right\|_2 = \left\| (I + \bar{R}^{-1} \delta \bar{R})^{-1} \right\|_2 = 1 + r \quad (8.9)$$

Entonces usando :

$$\left\| \begin{pmatrix} \bar{c} \\ 1 \end{pmatrix} \right\|_2 \leq 1.5 (n-1) \cdot 2^{-t} \left\| \begin{pmatrix} A \\ E \end{pmatrix} \right\|_2$$

$$\left\| \begin{pmatrix} \bar{e} \\ 1 \end{pmatrix} \right\|_2 \leq 1.5 n \cdot 2^{-t} \left\| \begin{pmatrix} b \\ 2 \end{pmatrix} \right\|_2$$

tenemos de :

$$\begin{pmatrix} \bar{r} \\ (I + V)R \end{pmatrix} \begin{matrix} \sim -1 \\ \sim -1 \end{matrix} \begin{matrix} \bar{t} \\ \bar{t} \end{matrix} = R \begin{matrix} \sim -1 \\ \sim -1 \end{matrix} \begin{matrix} \bar{r} \\ R \end{matrix} \begin{pmatrix} I + V \\ R \end{pmatrix} \begin{matrix} \bar{t} \\ \bar{t} \end{matrix} E$$

$$\left\| \begin{pmatrix} \bar{r} \\ (I + U)R \end{pmatrix} \begin{matrix} \sim -1 \\ \sim -1 \end{matrix} \begin{matrix} \bar{t} \\ \bar{t} \end{matrix} \right\|_2 \leq \frac{1}{1 - \beta} 1.5 (n-1) \cdot 2^{-t}$$

$$\pm 2^{-t} \left\| \begin{pmatrix} R^{-1} \\ 2 \end{pmatrix} \right\|_2 \left\| \begin{pmatrix} A \\ E \end{pmatrix} \right\|_2 \quad (8.10)$$

y de :

$$E = \bar{R} - Q A, \quad \bar{e} = y - Q b$$

y de :

$$\bar{R}^{-1} \begin{pmatrix} \bar{t} \\ (Q A - R) \end{pmatrix} \begin{matrix} \sim -1 \\ \sim -1 \end{matrix} = R \begin{matrix} \sim -1 \\ \sim -1 \end{matrix} \begin{pmatrix} \bar{r} \\ R \end{matrix} \begin{pmatrix} Q A - \bar{R} \\ -\bar{R}^{-1} \bar{R} \end{pmatrix} \begin{matrix} \bar{t} \\ \bar{R} \end{matrix}$$

(8.11)

$$\left\| \begin{pmatrix} \bar{R}^{-1} \begin{pmatrix} \bar{t} \\ (Q A - R) \end{pmatrix} \begin{matrix} \sim -1 \\ \sim -1 \end{matrix} \\ 2 \end{pmatrix} \right\|_2 \leq (n + r) \left(\left\| \begin{pmatrix} E \\ 3 \\ E \end{pmatrix} \right\|_2 \left\| \begin{pmatrix} R^{-1} \\ 2 \end{pmatrix} \right\|_2 + \left\| \begin{pmatrix} \bar{R}^{-1} \\ \bar{R} \end{pmatrix} \right\|_2 \right)$$

WILKINSON muestra que si la acumulación en doble precisión es usada entonces,

$$\begin{aligned} \left\| \delta R' \right\|_E &\leq (1.001 \cdot 2^{-t} + 3/2 (n+1) 2^{-2t}) \left\| R' \right\|_E \\ &\leq 2^{-t} \left\| \bar{R}' \right\|_E \end{aligned} \quad (8.12)$$

dado que $\left\| \bar{R} \right\|_E \leq n^{1/2} \left\| \bar{R}' \right\|_2$

tenemos :

$$\begin{aligned} \left\| \begin{matrix} -1 \\ \bar{R} \end{matrix} \delta \bar{R} \right\|_2 &= \left\| \begin{matrix} -1 \\ (\bar{R}')^{-1} \end{matrix} \delta \bar{R}' \right\|_2 \leq n^{1/2} \cdot 2^{-t} \times (\bar{R}') * \\ * \left\| \bar{R} \right\|_2 &\leq (1 + \beta)^{1/2} \left\| R \right\|_2 \end{aligned} \quad (8.13)$$

cuando $\beta < 1$, se sigue :

$$(8.14)$$

$$\left\| \begin{matrix} -1 \\ \bar{R} \end{matrix} \delta \bar{R} \right\|_2 \leq n^{1/2} \cdot 2^{-t} \times (\bar{R}) \leq (1 + \beta)^{1/2} n^{1/2} \cdot 2^{-t} \left\| \bar{R}^{-1} \right\|_2 \left\| A \right\|_E$$

sustituyendo en :

$$\left\| \begin{matrix} \sim -1 \\ \bar{R} \end{matrix} (Q A - R) \right\|_2 \leq (n + r) \left(\left\| \begin{matrix} E \\ 3 \\ E \end{matrix} \right\|_E \left\| \bar{R}^{-1} \right\|_2 + \left\| \begin{matrix} \sim -1 \\ \bar{R} \end{matrix} \delta \bar{R} \right\|_2 \right)$$

de :

$$\left\| \begin{matrix} E \\ \delta \\ E \end{matrix} \right\|_2 \leq 1.9 (n-1)^{\frac{1}{2}} n^{-t} \left\| \begin{matrix} A \\ E \end{matrix} \right\|_E$$

$$\left\| \begin{matrix} e \\ \delta \\ 2 \end{matrix} \right\|_2 \leq 1.9 n^{\frac{1}{2}} (n+1)^{-t} \left\| \begin{matrix} b \\ 2 \end{matrix} \right\|_2$$

y :

$$\left\| \begin{matrix} \bar{R}^{-1} \\ \delta \bar{R} \end{matrix} \right\|_2 \leq n^{\frac{1}{2}} 2^{-t1} x(\bar{R}) \leq (1+\beta)^{\frac{1}{2}} n^{\frac{1}{2}} 2^{-t1} \left\| \begin{matrix} R^{-1} \\ 2 \end{matrix} \right\|_2 \left\| \begin{matrix} A \\ E \end{matrix} \right\|_E$$

podemos demostrar que :

$$\left\| \begin{matrix} \bar{R}^{-1} \\ \delta \bar{R} \end{matrix} \right\|_2 \leq (1+r) - 1.9 n^{\frac{1}{2}} (n+1)^{\frac{1}{2}} \\ + 2^{-t} \left\| \begin{matrix} R^{-1} \\ 2 \end{matrix} \right\|_2 \left\| \begin{matrix} A \\ E \end{matrix} \right\|_E \quad (8.15)$$

Para la estimación de :

$$\left\| \begin{matrix} \bar{R}^{-1} \\ \delta \bar{R} \end{matrix} \right\|_2 = \left\| (I + \bar{R}^{-1} \delta \bar{R})^{-1} \right\|_2 = 1 + r$$

nos lleva a:

$$r \leq \left\| \begin{matrix} \bar{R}^{-1} \\ \delta \bar{R} \end{matrix} \right\|_2 \div (1 - \left\| \begin{matrix} \bar{R}^{-1} \\ \delta \bar{R} \end{matrix} \right\|_2) \quad (8.16)$$

de :

$$\left\| \begin{matrix} \bar{R}^{-1} \\ \delta \bar{R} \end{matrix} \right\|_2 \leq n^{\frac{1}{2}} 2^{-t1} x(\bar{R}) \leq (1+\beta)^{\frac{1}{2}} n^{\frac{1}{2}} 2^{-t1} \left\| \begin{matrix} R^{-1} \\ 2 \end{matrix} \right\|_2 \left\| \begin{matrix} A \\ E \end{matrix} \right\|_E$$

$$\left\| \bar{R}^{-1} \right\|_2 \leq \frac{1}{1-\beta} \left\| R^{-1} \right\|_2$$

y

$$\beta = 3.42 \epsilon = 3.42 n^{-1/2} \left\| A \right\|_2 \left\| R^{-1} \right\|_2 < 1$$

tenemos :

$$\left\| \bar{R}^{-1} \delta R \right\|_2 \leq 1.05 \left(\frac{1+\beta}{1-\beta} \right)^{1/2} n^{-1/2} \left\| R^{-1} \right\|_2 \left\| E \right\|_2 \leq \frac{1.05 \beta}{3.42(n+1)} \left(\frac{1+\beta}{1-\beta} \right)^{1/2}$$

y si hacemos la suposición adicional $\beta \leq 0.9$ entonces,

$$\left\| \bar{R}^{-1} \delta R \right\|_2 \leq 1.22 / (n+1) \leq 0.41$$

usando esto para estimar el denominador en :

$$r \leq \left\| \bar{R}^{-1} \delta R \right\|_2 \div \left(1 - \left\| \bar{R}^{-1} \delta R \right\|_2 \right)$$

podemos mostrar que :

$$r \leq 1.8 n^{-1/2} \epsilon x(\bar{R}') \quad (8.17)$$

$$r \leq \frac{0.53}{n+1} \left(\frac{1+\beta}{1-\beta} \right)^{1/2} \beta \quad (8.18)$$

Aunque no generalmente es verdad que $x(\bar{R}') \leq x(A)$

ocasionalmente $x(\bar{R}') < x(R') < x(\bar{R}) \approx x(A)$.

En efecto, para muchos sistemas muy mal condicionados

$x(\bar{R}')$ serán de orden unitario.

Entonces :

$$\left\| \begin{matrix} -1 \\ \bar{R} \delta R \end{matrix} \right\|_2 \leq n \cdot 2^{-t} \cdot x(R) \leq (1 + \beta) \cdot n \cdot 2^{-t} \left\| \begin{matrix} -1 \\ \bar{R} \end{matrix} \right\|_2 \left\| A \right\|_E$$

y consecuentemente de :

$$r \leq \frac{0.53}{n+1} \left(\frac{1 + \beta}{1 - \beta} \right) \cdot \beta$$

se sobreestimar\u00e1 considerablemente el error en la sustituci\u00f3n hacia atr\u00e1s, el cual siempre es de poca importancia.

Ahora tomaremos normas en :

$$\left\| \begin{matrix} \bar{r} - r \\ x - x \end{matrix} \right\|_E = \left\| \begin{matrix} \approx -1 \cdot t & & \\ -Q \bar{R} E & & (1 - \bar{Q} \bar{Q}) A \\ 1 & & \\ \bar{r} & \bar{r}^{-1} \cdot t & \\ (\bar{R} (1+U) R) E & E & \bar{R} (Q A - R) \end{matrix} \right\|_E^*$$

$$* \left\| \begin{matrix} r \\ x \end{matrix} \right\|_E + \left\| \begin{matrix} -e \\ 2 \\ \bar{r}^{-1} \\ R e \end{matrix} \right\|_E \quad (8.4)$$

y usando :

(8.5)

$$\left\| \begin{matrix} \approx -1 \cdot t \\ -Q \bar{R} E \\ 1 \end{matrix} \right\|_2 \leq 1.5 \cdot n \cdot (n-1) \cdot 2^{-t} \left\| \begin{matrix} -1 \\ \bar{R} \end{matrix} \right\|_2 \left\| A \right\|_E$$

(8.6)

$$\left\| (1 - \bar{Q} \bar{Q}) A \right\|_2 \leq 3.25 \cdot (n-1) \cdot 2^{-t} \left\| A \right\|_E$$

$$\left\| \left(\bar{R}^{-1} (1 + U) R \right)^{-1} E \right\|_2 \leq \frac{1}{1 - \beta} \cdot 1.5 (n-1) \cdot \epsilon$$

$$\cdot 2^{-t} \left\| \left\| R^{-1} \right\|_2 \left\| A \right\|_2 \right\|_E \quad (8.10)$$

y

$$\left\| \left(\bar{R}^{-1} (Q A - R) \right)^{-1} E \right\|_2 \leq (1 + \rho) \cdot 1.9 n^{1/2} (n+1) \cdot \epsilon$$

$$\cdot 2^{-t} \left\| \left\| \bar{R}^{-1} \right\|_2 \left\| A \right\|_2 \right\|_E \quad (8.15)$$

junto con los estimadores para $\left\| \left\| e \right\|_2 \right\|_2$ y $\left\| \left\| e \right\|_3 \right\|_2$

en :

$$\left\| \left\| E \right\|_2 \right\|_E \leq 3.25 (n-1) \cdot 2^{-t} \left\| \left\| A \right\|_2 \right\|_E$$

$$\left\| \left\| e \right\|_2 \right\|_2 \leq 3.25 n \cdot 2^{-t} \left\| \left\| L \right\|_2 \right\|_2$$

y

$$\left\| \left\| \bar{R}^{-1} \right\|_2 \left\| A \right\|_2 \right\|_2 \leq (1 - \beta)^{-1/2} \cdot \epsilon(A)$$

eliminando : $\left\| \left\| \bar{R}^{-1} \right\|_2$ por medio de la desigualdad :

$$\left\| \bar{R}^{-1} \right\|_2 \left\| A \right\|_2 \leq (1 - \beta)^{-1/2} \kappa(A)$$

y retomando las desigualdades :

$$n - 1 \leq 0.2n \sqrt{n+1}, \quad n \leq 2/3 n \sqrt{n+1}$$

y llegamos al resultado :

$$\left\| \begin{array}{c} \left\| \bar{r} - r \right\|_2 \\ \left\| A \right\|_2 \left\| \bar{x} - x \right\|_2 \end{array} \right\|_2 \leq \left[\begin{array}{c} \frac{0.79 - \kappa(A)}{\sqrt{1 - \beta}} \quad 0.81 \\ 0.24 \kappa(A) \quad \kappa(A) \\ \frac{1}{1 - \beta} \quad \frac{1}{\sqrt{1 - \beta}} \end{array} \right] \quad * \quad \dagger$$

$$* \quad \left[\begin{array}{c} \left\| r \right\|_2 \frac{\left\| A \right\|_2 E}{\left\| A \right\|_2} \\ \left\| A \right\|_2 \left\| x \right\|_2 \frac{\left\| A \right\|_2 E}{\left\| A \right\|_2} + \left\| b \right\|_2 \end{array} \right] \quad -t \quad (8.19)$$

donde :

$$f(n) = (1 + r) 1.9 n \sqrt{n+1}$$

Observamos que :

$$\left\| \begin{array}{l} \| \bar{r} - r \|_2 \\ \| A \|_2 \| \bar{x} - x \|_2 \end{array} \right\|_2 \leq \begin{bmatrix} \frac{0.79 + x(A)}{\sqrt{1-\beta}} & 0.81 \\ \frac{0.24 x(A)}{1-\beta} & \frac{x(A)}{\sqrt{1-\beta}} \end{bmatrix} \cdot$$

$$\cdot \left[\begin{array}{l} \| r \|_2 \frac{\| A \|_2 E}{\| A \|_2} \\ \| A \|_2 \| x \|_2 \frac{\| A \|_2 E}{\| A \|_2} + \| b \|_2 \end{array} \right] \cdot f(n)^{-t}$$

da un límite más pequeño para el error que aquel

obtenido en :

$$\left\| \begin{array}{c} \left\| \delta r \right\|_2 \\ \left\| A \right\|_2 \left\| \delta x \right\|_2 \end{array} \right\|_2 \leq \left[\begin{array}{cc} \frac{\kappa(A)}{\sqrt{1-\alpha}} & \alpha \\ \frac{\kappa(A)}{1-\alpha} & \frac{\kappa(A)}{\sqrt{1-\alpha}} \end{array} \right] *$$

$$* \left[\begin{array}{c} \left\| r \right\|_2 \frac{\left\| \delta A \right\|_2}{\left\| A \right\|_2} \\ \left\| A \right\|_2 \left\| x \right\|_2 \frac{\left\| \delta A \right\|_2}{\left\| A \right\|_2} + \left\| \delta b \right\|_2 \end{array} \right]$$

cuando :

$$\left\| \delta A \right\|_2 = f(n) \cdot 2^{-t} \left\| A \right\|_E, \quad \left\| \delta b \right\|_2 = f(n) \cdot 2^{-t} \left\| b \right\|_2$$

y esto de acuerdo a :

$$\beta = 3.42 \quad c = 3.42 \cdot n \cdot (n+1) \cdot 2^{-t} \left\| A \right\|_2 \left\| R^{-1} \right\|_2 < 1$$

Y :

$$\alpha = (\sqrt{2} + 1) \left\| \left\| R^{-1} \right\| \right\|_2 \left\| \left\| \delta A \right\| \right\|_2$$

implica : (nota 161)

$$\alpha = (1 + \delta) (\sqrt{2} + 1) \cdot 1.9 n^{1/2} (n + 1) \kappa$$

$$\approx 2^{-t} \left\| \left\| R^{-1} \right\| \right\|_2 \left\| \left\| A \right\| \right\|_E > \beta$$

Como el factor $(1 + \delta)$ esencialmente es de importancia, podemos decir que la desviación de la solución, de la verdadera solución es aproximadamente menor que la desviación resultante, de alguna perturbación $\delta A, \delta b$ tal que :

$$\left\| \left\| \delta A \right\| \right\|_E \div \left\| \left\| A \right\| \right\|_E \approx \left\| \left\| \delta b \right\| \right\|_2 \div \left\| \left\| b \right\| \right\|_2 \approx 2^{-t} \cdot n^{1/2}$$

Además cuando A está mal condicionada podemos, como remarcado en el final de la solución 5, esperar un resultado mucho mayor que este.

Finalmente observamos que cuando el procedimiento modificado es usado para computar R, entonces el procedimiento clásico, no será contado para ser usado para calcular y si esto fue hecho tendremos :

$$Y \approx \bar{Q}^{-t} b = (\bar{Q} - \tilde{Q})^{-t} b + \tilde{Q}^{-t} b$$

y consecuentemente :

$$\begin{aligned} \left\| \left\| y - \bar{Q}^{-t} b \right\| \right\|_2 &\approx \left\| \left\| \bar{Q} - \tilde{Q} \right\| \right\|_2 \left\| \left\| b \right\| \right\|_2 \\ &= \left\| \left\| U Q \right\| \right\|_2 \left\| \left\| b \right\| \right\|_2 = \text{Const.} \times (A)^{-t} \left\| \left\| b \right\| \right\|_2 \end{aligned}$$

puede ser verificado que si $\|r\|_2 \leq \|b\|_2$, esto significa que los errores en la solución calculada generalmente serán multiplicados por $\kappa(A)$.

CONCLUSIONES

El procedimiento modificado de GRAM-SCHMIDT ha sido usado sucesivamente, desde 1960 como parte del sistema de datos procesados para espectros ópticos.

Tiene más ó menos la misma gran estabilidad numérica como otros algoritmos comparables basados en transformaciones ortogonales y generalmente será preferible el procedimiento clásico de GRAM-SCHMIDT. La reortogonalización es entonces no siempre necesaria, y la solución es obtenida en aproximadamente $m n^2$ multiplicaciones. El pivotes puede ser hecho fácilmente, pero no es requisito para la estabilidad, el requerimiento de almacenaje es aproximadamente :

$$m n + n^2 \div 2$$

Esto sería comparado a $2 m n + 4 n^2 \div 3$ multiplicaciones y n^2 lugares de almacenaje necesitadas para formar y resolver las ecuaciones normales en doble precisión.

Soluciones más precisas para problemas de Mínimos Cuadrados pueden ser obtenidas por usar el refinamiento iterativo, este fue el primero propuesto por GOLUB y usado también en SAHVER. Esquemas diferentes para el refinamiento son discutidos en las notas de refinamiento iterativo de mínimos cuadrados. Los resultados obtenidos son sin embargo, solo satisfactorios cuando las ecuaciones son también compatibles es decir cuando :

$$\|x^2(A)\|_2 \|r\|_2 < \|A\|_2 \|x\|_2$$

En notas posteriores de SAHVER donde también se resuelven ejemplos numéricos se aplica el refinamiento iterativo a problemas de mínimos cuadrados para calcular los residuales del sistema :

$$\begin{pmatrix} r + Ax \\ A r \end{pmatrix} = \begin{pmatrix} I & A \\ A & 0 \end{pmatrix} \begin{pmatrix} r \\ x \end{pmatrix} = \begin{pmatrix} b \\ 0 \end{pmatrix}$$

en precisión doble, y usando la inversa:

$$\begin{pmatrix} 1 & A \\ t & 0 \end{pmatrix}^{-1} \approx \begin{pmatrix} I - \bar{Q} \bar{Q}^t & \bar{Q} \bar{R}^{-t} \\ \bar{R}^{-1} \bar{Q}^t & \bar{R}^{-1} (I + U) \bar{R}^{-t} \end{pmatrix}$$

Para los incrementos en r y x . Esta inversa resulta como una pequeña generalización del procedimiento modificado de GRAM-SCHMIDT y como se mostró, que la iteración generalmente convergerá a la solución exacta, correctamente redondeada cuando:

$$\|x\|_2^2 (A) \left\| \begin{matrix} r \\ 2 \end{matrix} \right\|_2 \text{ es del mismo orden de magnitud que:}$$

$$\frac{t}{2} \left\| \begin{matrix} A \\ 2 \end{matrix} \right\|_2 \left\| \begin{matrix} x \\ 2 \end{matrix} \right\|_2$$

BIBLIOGRAFIA

- 1.- Bauer, F. L. Optimally Scaled Matrices.
Numerische Mathematik. Vol .5 1963
p. 73 - 87
- 2.- Golub. G. H. & J. H. Wilkinson.- Note on the
Iterative Refinement of Least Squares Solution
Numerische Mathematik 1965 p. 139 - 148
- 3.- Strang, Gilbert.- Linear Algebra.- Academic Press,
New - York. 1980
- 4.- Valdez Bautista Beatriz.- Tesis Profesional
" Aspectos de Cálculo del Algebra Lineal "
Cd. Universitaria de México 1988, p.100 - 135

Introducción	I
--------------------	---

C A P I T U L O I

Definiciones Básicas	1
----------------------------	---

C A P I T U L O II

Descripción del Algoritmo	12
---------------------------------	----

C A P I T U L O III

Definiciones Básicas del Análisis de Error	28
--	----

C A P I T U L O IV

Errores en una Proyección Elemental	47
---	----

C A P I T U L O V

Errores en la Factorización	59
-----------------------------------	----

C A P I T U L O VI

Ortogonalización de los Vectores Calculados	71
---	----

C A P I T U L O VII

Condición del Problema de Mínimos Cuadrados	84
---	----

C A P I T U L O VIII

Errores en la Solución Calculada	92
Conclusiones	103