

## PRACTICA 3a

Resolución de sistemas lineales.  
Métodos directos.

Se abre el telón y se ven dos sistemas lineales incompatibles. ¿Cómo se llama la obra?

*Kramer vs. Kramer.*

**Eliminación gaussiana. Factorización LU.** Sea  $A\mathbf{x} = \mathbf{b}$ , un sistema de  $n$  ecuaciones lineales con  $n$  incógnitas, donde  $A$  es una matriz cuadrada de orden  $n$  de elementos reales y no singular (por lo cual el sistema admite una solución y ésta es única). Los métodos numéricos *directos* para resolver dicho sistema son aquellos que, en ausencia de errores de redondeo, obtienen la solución exacta en un número *finito* de pasos. El método fundamental es el procedimiento de *eliminación gaussiana* con la estrategia de *pivoteo parcial*, la cual permite el intercambio de filas. Dicho método muestra que la matriz  $A$  admite una factorización de la forma

$$A = PLU.$$

Aquí,  $L$  es una *matriz triangular inferior* con elementos diagonales iguales a la unidad y cuyos elementos bajo la diagonal son los *multiplicadores* utilizados durante la eliminación.  $U$  es una *matriz triangular superior* cuyos elementos son los coeficientes del sistema final equivalente, siendo los elementos de la diagonal no nulos. Si los intercambios de fila durante el proceso de eliminación son registrados en un vector  $\mathbf{p}$  (cuyo valor inicial es  $(1, 2, \dots, n)$ ), entonces la *matriz de permutación*  $P$  se obtiene a partir de la matriz identidad  $I$  de orden  $n$  permutando sucesivamente la columna  $i$  por la columna  $j = \mathbf{p}(i)$  para  $i = 1, 2, \dots, n$ . Por ejemplo, si para  $n = 4$ ,  $\mathbf{p} = (3, 4, 3, 4)$ , entonces

$$P = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}.$$

Conocida la factorización, el sistema de ecuaciones lineales es equivalente a  $PLU\mathbf{x} = \mathbf{b}$ , el cual conduce a dos sistemas *triangulares*:

$$L\mathbf{y} = P^t\mathbf{b}, \quad U\mathbf{x} = \mathbf{y}.$$

Así, para obtener el vector solución  $\mathbf{x}$ , se resuelve en primer lugar el primer sistema por *sustitución hacia adelante* para obtener  $\mathbf{y}$  y luego se resuelve el segundo sistema por *sustitución hacia atrás*.

☞ La factorización  $LU$  de una matriz de orden  $n$  es un procedimiento finito de  $n - 1$  pasos cuyo *costo computacional*, medido por el número de operaciones aritméticas necesarias llevarla a cabo, es  $\mathcal{O}(n^3)$ . En tanto que el costo computacional de la resolución de los dos sistemas triangulares es  $\mathcal{O}(n^2)$ , el cual es un orden de magnitud menor

que el costo de la factorización. (Notar también que en la factorización el costo computacional de la determinación de los pivotes requiere de  $\mathcal{O}(n^2)$  comparaciones, el cual es entonces despreciable frente al costo de las operaciones aritméticas).

☞ Si tenemos varios sistemas de ecuaciones con la misma matriz de coeficientes  $A$ ,

$$A\mathbf{x}_1 = \mathbf{b}_1, \quad A\mathbf{x}_2 = \mathbf{b}_2, \quad \dots, \quad A\mathbf{x}_m = \mathbf{b}_m,$$

los mismos pueden ser tratados simultáneamente como el sistema matricial

$$AX = B,$$

donde  $X = [\mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_m]$  y  $B = [\mathbf{b}_1 | \mathbf{b}_2 | \dots | \mathbf{b}_m]$  son matrices rectangulares  $n \times m$ .

**Ejercicio 1.** Resolver *a mano* por eliminación gaussiana con pivoteo parcial el sistema:

$$\begin{pmatrix} 0 & 4 & 1 \\ 1 & 1 & 3 \\ 2 & -2 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 9 \\ 6 \\ -1 \end{pmatrix}.$$

**Ejercicio 2.** El siguiente ejercicio muestra que la estrategia de pivoteo parcial es crucial para mejorar la *estabilidad numérica* del procedimiento de eliminación gaussiana. El sistema

$$\begin{pmatrix} 0.0003 & 1.566 \\ 0.3454 & -2.436 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1.569 \\ 1.018 \end{pmatrix},$$

tiene la solución exacta  $\mathbf{x} = (10, 1)^t$ . Utilice aritmética decimal de cuatro dígitos para resolverlo sin y con pivoteo parcial. Indique el porque de la diferencia en las soluciones numéricas obtenidas.

☞ La estabilidad numérica del método de eliminación gaussiana depende del tamaño del llamado *factor de crecimiento*, definido como la razón de las magnitudes del mayor coeficiente de  $U$  al mayor coeficiente de  $A$ . Sin pivoteo, este factor puede ser arbitrariamente grande y por lo tanto la eliminación gaussiana sin pivoteo es numéricamente inestable. Sin embargo, aún con pivoteo parcial este factor puede ser tan grande como  $2^{n-1}$  (el cual es enorme aún para modestos valores de  $n$ ) y por lo tanto el método, estrictamente hablando, es numéricamente inestable. Sin embargo, la experiencia ha mostrado que éste comportamiento del factor de crecimiento es extremadamente raro. En la práctica este factor se mantiene moderado. En consecuencia, *el método de eliminación gaussiana con pivoteo parcial puede considerarse numéricamente estable en la práctica.*

Para la implementación en una computadora del procedimiento de eliminación gaussiana, en vez de programar nuestras propias subrutinas, utilizaremos la biblioteca de rutinas LAPACK (Linear Algebra PACKage), la cual es una colección de subrutinas para resolver numéricamente problemas matemáticos que se enmarcan en el campo del álgebra lineal, junto con su interfaz LAPACK95, la cual explota las características propias de Fortran 95 (asignación dinámica de memoria, arreglos de tamaño ajustable en la lista

de argumentos y argumentos opcionales). Para la resolución de sistemas lineales, LAPACK95 proporciona un amplio conjunto de subrutinas dependiendo de la naturaleza particular de la matriz de coeficientes. En particular, la subrutina LA\_GESV permite resolver un sistema de ecuaciones lineales para una matriz  $A$  general almacenada en un arreglo bidimensional y uno o varios términos independientes almacenados en un arreglo bidimensional  $B$  según el esquema discutido. La cuestión de utilizar las rutinas de LAPACK95 consta de dos partes: cómo llamar a las subrutinas en nuestro programa, y como compilar el mismo. A continuación discutimos éstos dos puntos.

El uso de las rutinas de LAPACK95 requiere de la invocación de dos módulos a través de la sentencia USE:

- LA\_PRECISION, quien define la precisión de los tipos de datos reales, ya sea de simple o doble precisión, via el parámetro WP cuya asignación será SP ó DP, respectivamente. Por ejemplo, si (como haremos) trabajamos en doble precisión, la sentencia apropiada es

```
USE la_precision, ONLY: WP => DP
```

- F95\_LAPACK, quien define las interfaces explícitas de las subrutinas a utilizar. Por ejemplo, si vamos a utilizar la subrutina LA\_GESV, la sentencia apropiada es

```
USE F95_LAPACK, ONLY: la_gesv
```

La invocación de la subrutina LA\_GESV procede entonces como sigue, donde indicamos con corchetes los argumentos opcionales,

```
CALL la_gesv(A,B,[IPIV=p],[INFO=info])
```

Como dato de entrada  $A$  es un arreglo bidimensional que almacena la matriz de coeficientes del sistema y  $B$  un arreglo unidimensional o bidimensional que contiene los términos independientes. La subrutina devuelve en  $A$  los factores  $L$  y  $U$  de la factorización y en  $B$  las soluciones de los sistemas. Opcionalmente devuelve también el vector de permutaciones  $\mathbf{p}$ , con el cual puede construirse la matriz de permutación  $P$  que completa la factorización y una clave entera de error, `info`, cuyo valor igual a cero implica que la resolución del sistema se efectuó con éxito y otros valores algún tipo de error.

☞ Una discusión completa de los argumentos necesarios para invocar cada subrutina de LAPACK95 está dada en el manual de usuario, el cual puede ser consultado *on-line* en <http://www.netlib.org/lapack95/lug95/>.

Finalmente para compilar el programa fuente debemos informar explícitamente al compilador que utilizaremos la biblioteca de rutinas LAPACK y su interfaz LAPACK95 como sigue<sup>1</sup>:

<sup>1</sup>Para otros sistemas el comando puede variar.

```
$ gfortran -Wall -o programa programa.f90 \
-I/usr/local/include/lapack95 \
-llapack -llapack95
```

**Ejercicio 3.** Utilice las rutina LA\_GESV para resolver los sistemas  $AX = B$  para las matrices de coeficientes  $A$  dadas por:

$$(a) \begin{pmatrix} 0 & 1 & 2 \\ 1 & 2 & 3 \\ 2 & 3 & 2 \end{pmatrix}, \quad (b) \begin{pmatrix} 2 & 1 & 0 \\ 0 & 3 & 2 \\ 1 & 2 & 4 \end{pmatrix},$$

$$(c) \begin{pmatrix} 1 & 2 & 1 \\ 2 & 0 & -2 \\ -1 & 2 & 3 \end{pmatrix},$$

y términos independientes dados por

$$B = \begin{pmatrix} 2 & 4 & 1 \\ 4 & 12 & 0 \\ 5 & 17 & 1 \end{pmatrix},$$

Explicitar la factorización  $PLU$  de las matrices  $A$ .

**Matrices especiales.** Cuando la matriz  $A$  del sistema tiene ciertas propiedades especiales es posible reducir ya sea el costo computacional del procedimiento de eliminación gaussiana, el costo de almacenamiento de la matriz o ambos. En primer lugar, para algunas clases de matrices la eliminación gaussiana es numéricamente estable *sin* ninguna estrategia de pivoteo, a saber cuando:

- $A$  es *diagonalmente dominante* por columnas o filas, esto es,

$$|a_{jj}| \geq \sum_{i \neq j} |a_{ij}|, \text{ para cada } j = 1, 2, \dots, n, \text{ ó}$$

$$|a_{ii}| \geq \sum_{j \neq i} |a_{ij}|, \text{ para cada } i = 1, 2, \dots, n,$$

- $A$  es *simétrica*  $A = A^t$  y *definida positiva*, esto es,  $\mathbf{x}^t A \mathbf{x} > 0$  para todo  $\mathbf{x} \neq 0$ .

Por otra parte, para otras matrices conocidas como *matrices de banda* la presencia de un gran número de elementos nulos en un patrón regular permite simplificar el algoritmo de factorización. Entre ellas se encuentran, en particular, las *matrices tridiagonales*. Estas variantes serán discutidas a continuación.

**Matrices simétricas definidas positivas. Método de Cholesky.** Si la matriz de coeficientes  $A$  es simétrica y definida positiva, existe una factorización de la forma  $A = LL^t$  donde  $L$  es una matriz triangular inferior con elementos positivos en la diagonal. Dicha factorización, determinada por el *método de Cholesky*, no requiere intercambio de filas y es estable frente a los errores de redondeo. Esta factorización es utilizada para resolver el sistema e implementada en LAPACK95 por la subrutina LA\_POSV, donde  $A$  es almacenada en un arreglo bidimensional como es usual, aunque, por la simetría, solo se accede al triángulo superior o inferior.

**Ejercicio 4.** Resolver los sistemas de ecuaciones  $AX = B$  donde

$$A = \begin{pmatrix} 38 & 3 & 4 & 6 & 5 \\ 3 & 48 & 6 & 7 & 7 \\ 4 & 6 & 52 & 1 & 6 \\ 6 & 7 & 1 & 56 & 10 \\ 5 & 7 & 6 & 10 & 74 \end{pmatrix}$$

$$B = \begin{pmatrix} 56 & 112 & 168 \\ 71 & 142 & 213 \\ 69 & 138 & 207 \\ 80 & 160 & 240 \\ 102 & 204 & 306 \end{pmatrix}$$

Explicitar la factorización de  $A$ .

**Matrices tridiagonales.** Las *matrices tridiagonales* tienen todos sus elementos nulos por fuera de la diagonal principal y la subdiagonal y superdiagonal adyacentes. En tal caso, los elementos relevantes de la matriz pueden ser entonces almacenados en tres arreglos unidimensionales independientes, uno por cada diagonal. Esto permite que sólo se requiera un espacio de almacenamiento  $\mathcal{O}(3n)$  en vez del arreglo  $n^2$  completo y, además, el costo computacional de la factorización resulta ser  $\mathcal{O}(9n)$ , lo cual implica que el procedimiento es muy eficiente respecto al caso general. LAPACK95 proporciona la subrutina LA\_GTSV para resolver sistemas de ecuaciones lineales con matrices de coeficientes tridiagonales almacenados de este modo.

**Ejercicio 5.** Resolver los sistemas de ecuaciones  $AX = B$  donde

$$A = \begin{pmatrix} 2 & 2 & 0 & 0 & 0 & 0 \\ 1 & 5 & 4 & 0 & 0 & 0 \\ 0 & 4 & 8 & 4 & 0 & 0 \\ 0 & 0 & 6 & 8 & 6 & 0 \\ 0 & 0 & 0 & 6 & 4 & 3 \\ 0 & 0 & 0 & 0 & 4 & 9 \end{pmatrix}$$

y

$$B = \begin{pmatrix} 4 & 8 & 12 \\ 10 & 20 & 30 \\ 16 & 32 & 48 \\ 20 & 40 & 60 \\ 13 & 26 & 39 \\ 13 & 26 & 39 \end{pmatrix}$$

**Normas vectoriales y matriciales.** Toda solución numérica de un sistema lineal  $A\mathbf{x} = \mathbf{b}$  debe considerarse como una solución aproximada  $\hat{\mathbf{x}}$  debido a los errores de redondeo introducidos en el cómputo de la misma. Para medir cuán buena es la aproximación  $\hat{\mathbf{x}}$  a la solución exacta  $\mathbf{x}$  se recurre al concepto de *norma vectorial*. La *distancia entre los vectores*  $\hat{\mathbf{x}}$  y  $\mathbf{x}$  es entonces definida como la norma de la diferencia de los mismos:  $\|\mathbf{x} - \hat{\mathbf{x}}\|$ . En la práctica, las normas vectoriales en  $\mathbb{R}^n$  de uso frecuente son, siendo  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  un vector de  $\mathbb{R}^n$ :

- norma  $l_\infty$ :  $\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} \{|x_i|\}$ ,
- norma  $l_1$ :  $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$ ,
- norma  $l_2$ :  $\|\mathbf{x}\|_2 = [\sum_{i=1}^n x_i^2]^{1/2}$ ,

Asimismo se define una *norma matricial natural* subordinada a la respectiva norma vectorial, siendo entonces para  $A = (a_{ij})$  una matriz de  $\mathbb{R}^{n \times n}$ :

- norma  $l_\infty$ :  $\|A\|_\infty = \max_{1 \leq i \leq n} \{\sum_{j=1}^n |a_{ij}|\}$ ,
- norma  $l_1$ :  $\|A\|_1 = \max_{1 \leq j \leq n} \{\sum_{i=1}^n |a_{ij}|\}$
- norma  $l_2$ :  $\|A\|_2 = [\rho(A^t A)]^{1/2}$ ,

donde, el *radio espectral* de una matriz  $A$ , de  $n \times n$ , se define como  $\rho(A) = \max \{|\lambda_i|, \lambda_i \text{ autovalor de } A\}$ .

**Ejercicio 6.** Calcular las normas  $l_\infty$ ,  $l_1$  y  $l_2$  del vector  $\mathbf{x} = (-1, 1, -2)$  de  $\mathbb{R}^3$  y las normas  $l_\infty$  y  $l_1$  de la matriz

$$A = \begin{pmatrix} 1 & 2 & -1 \\ 0 & 3 & -1 \\ 5 & -1 & 1 \end{pmatrix}$$

**Ejercicio 7.** Escribir sentencias Fortran para calcular las normas vectoriales  $l_\infty$ ,  $l_1$  y  $l_2$  y las normas matriciales  $l_\infty$  y  $l_1$ . *Ayuda:* considerar el uso de las funciones intrínsecas SUM, ABS, MAXVAL y DOT\_PRODUCT.

**Estimación a posteriori del error. Número de condición.** En la práctica, una solución *numérica*  $\hat{\mathbf{x}}$  de un sistema de ecuaciones lineales  $A\mathbf{x} = \mathbf{b}$  tendrá cierto error  $\|\mathbf{x} - \hat{\mathbf{x}}\|$ , el cual, por supuesto, no puede ser calculado. Ahora bien, siempre podemos calcular el *vector residuo*  $\mathbf{r} = \mathbf{b} - A\hat{\mathbf{x}}$ . Puesto que  $\mathbf{r} = \mathbf{0}$  implica que  $\hat{\mathbf{x}}$  es la solución exacta del sistema, resulta natural sugerir que si  $\|\mathbf{r}\|$  es pequeña, entonces  $\hat{\mathbf{x}}$  es una solución aproximada precisa. Sin embargo, esto no es siempre el caso. La exactitud de la solución calculada depende también del *número de condición*  $\kappa(A) = \|A\| \|A^{-1}\|$  de la matriz  $A$  del sistema, de acuerdo a la siguiente desigualdad:

$$\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \kappa(A) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}.$$

Esta desigualdad nos dice que el error relativo en una solución numérica puede ser tan grande como  $\kappa(A)$  veces su residuo relativo. De este modo, solo si  $\kappa(A) \simeq 1$  el error relativo y el residuo relativo son del mismo tamaño, y el residuo podrá ser utilizado con seguridad como una estimación del error. En esta situación se dice que la matriz  $A$  del sistema es una *matriz bien condicionada*. Si, por el contrario,  $\kappa(A) \gg 1$ , el intervalo en que se puede situar el error relativo es muy amplio y por lo tanto no se puede asegurar que la solución numérica es precisa, aún cuando el residuo sea pequeño. En esta situación, se dice que la matriz  $A$  es una *matriz mal condicionada*.

**Ejercicio 8.** Mostrar que para toda matriz no singular,  $\kappa(A) \geq 1$ .

**Ejercicio 9.** El ejemplo canónico de matrices mal condicionadas son las *matrices de Hilbert*. La matriz de Hilbert  $H^{(n)}$  de orden  $n$  está definida por

$$H_{ij}^{(n)} = \frac{1}{i+j-1}, \quad 1 \leq i, j \leq n$$

y es una matriz simétrica definida positiva. El siguiente fragmento de código Fortran permite calcular su inversa  $T^{(n)} = [H^{(n)}]^{-1}$ .

```
t(1,1) = REAL(n*n,WP)
i=1
DO j=2,n
  t(i,j) = - t(i,j-1) * &
    & REAL((n+j-1)*(i+j-2)*(n+1-j),WP) / &
    & REAL((i+j-1)*(j-1)*(j-1),WP)
END DO
DO i=2,n
  DO j=1,n
    t(i,j) = - t(i-1,j) * &
      & REAL((n+i-1)*(i+j-2)*(n+1-i),WP) / &
      & REAL((i+j-1)*(i-1)*(i-1),WP)
  END DO
END DO
```

Utilice dicho código para calcular el número de condición  $\kappa_\infty$  de las diez primeras matrices de Hilbert.

**Ejercicio 10.** Sea  $A\mathbf{x} = \mathbf{b}$  el sistema lineal dado por

$$A = \begin{pmatrix} 0.89 & 0.53 \\ 0.47 & 0.28 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 0.36 \\ 0.19 \end{pmatrix}.$$

Dada la solución aproximada  $\hat{\mathbf{x}} = (-11.5, 20)^t$ , calcule el vector residual y estime el error relativo. Sabiendo que la solución exacta es  $\mathbf{x} = (1, -1)^t$ , calcule el error exacto. Compare los resultados.

**Estimación a priori del error. Análisis perturbativo.** En la sección anterior hemos supuesto que  $A$  y  $\mathbf{b}$  están libres de error. Sin embargo, si el sistema lineal proviene de la modelización de un problema físico debemos esperar que los coeficientes del sistema estén sujetos a error bien porque provienen de otros cálculos o de mediciones. Además, aún por la sola causa de la representación de los números en punto flotante en la computadora,  $A$  y  $\mathbf{b}$  se encuentran afectados de errores. El número de condición juega también aquí un papel importante en el análisis de tales errores.

**Ejercicio 11.** Dado el sistema  $A\mathbf{x} = \mathbf{b}$ , suponga que el término  $\mathbf{b}$  se perturba por una cantidad  $\delta\mathbf{b}$ . Mostrar que el efecto de esta perturbación sobre la solución será  $\mathbf{x} + \delta\mathbf{x}$ , con  $\delta\mathbf{x}$  acotada por

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \kappa(A) \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|}.$$

Si ahora es la matriz  $A$  la que es perturbada por una cantidad  $\delta A$ , mostrar que el efecto de tal perturbación sobre la solución será  $\mathbf{x} + \delta\mathbf{x}$  con  $\delta\mathbf{x}$  acotada por

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x} + \delta\mathbf{x}\|} \leq \kappa(A) \frac{\|\delta A\|}{\|A\|}.$$

**Ejercicio 12.** Sea

$$A = \begin{pmatrix} 6 & 13 & -17 \\ 13 & 29 & -38 \\ -17 & -38 & 50 \end{pmatrix}$$

Si se desea resolver el sistema  $A\mathbf{x} = \mathbf{b}$ , donde  $\mathbf{b} = (1.9358, -2.3412, 3.5718)^t$  tiene sus elementos redondeados a 5 dígitos, ¿qué puede decirse del error relativo en la solución?

En general, si  $A$  y  $\mathbf{b}$  están perturbados de modo que  $(A + \delta A)(\mathbf{x} + \delta\mathbf{x}) = \mathbf{b} + \delta\mathbf{b}$ , entonces

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\kappa(A)}{1 - \kappa(A)(\|\delta A\|/\|A\|)} \left( \frac{\|\delta A\|}{\|A\|} + \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} \right),$$

siempre que  $\|\delta A\| \|A^{-1}\| < 1$ .

**Ejercicio 13.** Mostrar que si  $A$  y  $\mathbf{b}$  son conocidos exactamente, al ser representados como número de puntos flotante en la computadora la perturbación  $\delta\mathbf{x}$  correspondiente está acotada aproximadamente en la norma  $l_\infty$  por

$$\frac{\|\delta\mathbf{x}\|_\infty}{\|\mathbf{x}\|_\infty} \lesssim 2\mathbf{u} \kappa(A),$$

donde  $\mathbf{u}$  es la unidad de redondeo de la máquina y suponemos que  $A$  está bien condicionada.