

Sistemas de Ecuaciones Lineales II

- **Métodos Iterativos:** Matrices dispersas. Esquema general. Métodos de Jacobi y de Gauss-Seidel.

Matrices dispersas

En muchas aplicaciones los sistemas de ecuaciones lineales que deben resolverse involucran matrices de gran tamaño, pero la mayor parte de sus entradas son nulas. Estas matrices se denominan **dispersas** o **ralas** (en inglés y en OCTAVE, **sparse**) y existen técnicas para almacenarlas que sólo requieren una cantidad de posiciones de memoria aproximadamente igual al número de entradas no nulas de la matriz. Por ejemplo:

```
>> n = 7;  
>> A = diag(ones(n-1,1),-1)+2*diag(ones(n,1))+diag(ones(n-1,1),1);  
>> A  
A =  
  
    2    1    0    0    0    0  
    1    2    1    0    0    0  
    0    1    2    1    0    0  
    0    0    1    2    1    0  
    0    0    0    1    2    1  
    0    0    0    0    1    2
```

Para visualizar la gráficamente la estructura de la matriz podemos usar el comando `spy(A)`.

Ejemplo llenado de matrices dispersas

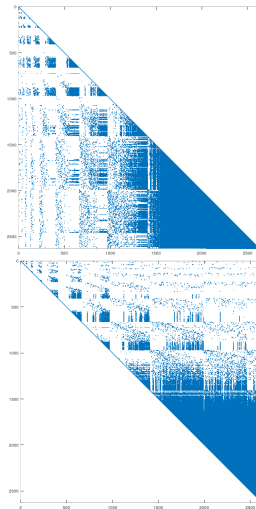
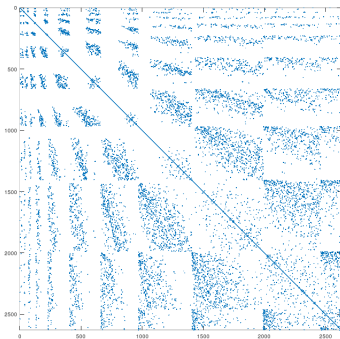
```
>> load DATA_025.mat
```

```
>> [L,U,P]=lu(A);
```

```
>> whos
```

Attr	Name	Size	Bytes	Class
====	====	====	=====	=====
	A	2629x2629	314928	double
	L	2629x2629	9048128	double
	P	2629x2629	21032	double
	U	2629x2629	8804992	double

Ejemplo llenado de matrices dispersas (cont.)



- ▶ Cuando la matriz \mathbf{A} del sistema a resolver es dispersa, pero no banda, los métodos (directos) estudiados hasta ahora (eliminación de Gauss, factorización LU o Cholesky) presentan el defecto denominado **llenado (fill-in)**.
- ▶ El llenado consiste en que, a medida que el proceso de eliminación avanza, se van creando elementos no nulos en posiciones de \mathbf{L} y \mathbf{U} en donde la matriz \mathbf{A} tiene ceros.
- ▶ Como consecuencia del llenado se tiene, por una parte, el aumento del número de flop y con ello el aumento del error de redondeo. Por otra parte se tiene el aumento en las necesidades de memoria para almacenar las matrices \mathbf{L} y \mathbf{U} , lo que puede llegar a hacer imposible aplicar estos métodos cuando \mathbf{A} es de gran tamaño.
- ▶ Los métodos que estudiaremos en el siguiente capítulo, llamados **iterativos**, evitan el llenado y sus consecuencias, al trabajar resolviendo reiteradamente sistemas con matriz diagonal o triangular-dispersa.

Esquema general

- Considere el sistema de ecuaciones

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

con $\mathbf{A} \in \mathbb{R}^{n \times n}$ no singular y $\mathbf{b} \in \mathbb{R}^n$.

Un **método iterativo** para resolver el sistema construye, a partir de un vector inicial $\mathbf{x}^{(0)}$, una sucesión de vectores $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(k)}, \dots$ la que, bajo condiciones apropiadas, resultará convergente a \mathbf{x} .

- Si suponemos $\mathbf{A} = \mathbf{N} - \mathbf{P}$, donde \mathbf{N} debe ser invertible, entonces

$$\begin{aligned}\mathbf{A}\mathbf{x} = \mathbf{b} &\iff (\mathbf{N} - \mathbf{P})\mathbf{x} = \mathbf{b} \\ &\iff \mathbf{N}\mathbf{x} = \mathbf{P}\mathbf{x} + \mathbf{b} \\ &\iff \mathbf{x} = \mathbf{N}^{-1}\mathbf{P}\mathbf{x} + \mathbf{N}^{-1}\mathbf{b}.\end{aligned}$$

- Se usa la igualdad $\mathbf{N}\mathbf{x} = \mathbf{P}\mathbf{x} + \mathbf{b}$ para definir un esquema general para construir la sucesión $\{\mathbf{x}^{(k)}\}$. Este esquema se da en el siguiente algoritmo.

Esquema general (cont.)

► Algoritmo del esquema general:

Dado el vector inicial $\mathbf{x}^{(0)}$,
para $k = 1, 2, \dots$ resolver:
 $\mathbf{N}\mathbf{x}^{(k)} = \mathbf{P}\mathbf{x}^{(k-1)} + \mathbf{b}$,
hasta que se satisfaga un criterio de detención.

- Definiendo $\mathbf{M} := \mathbf{N}^{-1}\mathbf{P}$ (**matriz de iteración**) y $\mathbf{e}^{(k)} := \mathbf{x} - \mathbf{x}^{(k)}$ (**error de $\mathbf{x}^{(k)}$**), para cada $k = 1, 2, \dots$ se tiene

$$\begin{aligned}\mathbf{e}^{(k)} &= \mathbf{x} - \mathbf{x}^{(k)} \\ &= (\mathbf{N}^{-1}\mathbf{P}\mathbf{x} + \mathbf{N}^{-1}\mathbf{b}) - (\mathbf{N}^{-1}\mathbf{P}\mathbf{x}^{(k-1)} + \mathbf{N}^{-1}\mathbf{b}) \\ &= \mathbf{N}^{-1}\mathbf{P}(\mathbf{x} - \mathbf{x}^{(k-1)}) \\ &= \mathbf{M}\mathbf{e}^{(k-1)}\end{aligned}$$

y, recursivamente,

$$\mathbf{e}^{(k)} = \mathbf{M}^k \mathbf{e}^{(0)}, \quad k = 1, 2, \dots$$

Convergencia de métodos iterativos.

- **Teorema.** (Convergencia) La sucesión $\{\mathbf{x}^{(k)}\}$ converge a la solución \mathbf{x} de $\mathbf{Ax} = \mathbf{b}$, si y sólo si, $\rho(\mathbf{M}) < 1$.
- **Observación.** Si la sucesión $\{\mathbf{x}^{(k)}\}$ converge, necesariamente lo hace a la solución \mathbf{x} de $\mathbf{Ax} = \mathbf{b}$.
- **Lema.** (Cota para el radio espectral) Sea \mathbf{A} una matriz cuadrada. Para cualquier norma matricial se tiene que

$$\rho(\mathbf{A}) \leq \|\mathbf{A}\|.$$

- **Corolario.** (Condición suficiente de convergencia) Una condición suficiente para que la sucesión $\{\mathbf{x}^{(k)}\}$ sea convergente a la solución \mathbf{x} de $\mathbf{Ax} = \mathbf{b}$ es que

$$\|\mathbf{M}\| < 1,$$

donde \mathbf{M} es la matriz de iteración del método que genera a $\{\mathbf{x}^{(k)}\}$.

Criterio de detención

- **Detención del proceso.** Cuando el proceso iterativo es convergente, éste se debe detener para un $\mathbf{x}^{(k)}$ tal que

$$\frac{\|\mathbf{M}\|}{1 - \|\mathbf{M}\|} \left\| \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \right\| \leq \text{tol},$$

donde **tol** indica un nivel de tolerancia prefijado para el error.

- **Lema.** Para $\|\mathbf{M}\| < 1$, se tiene que $\rho(\mathbf{M}) < 1$ y

$$\left\| \mathbf{x} - \mathbf{x}^{(k)} \right\| \leq \frac{\|\mathbf{M}\|}{1 - \|\mathbf{M}\|} \left\| \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \right\|.$$

- **Consecuencia:** Si detenemos el algoritmo cuando

$$\frac{\|\mathbf{M}\|}{1 - \|\mathbf{M}\|} \left\| \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \right\| \leq \text{tol},$$

nos aseguramos que el error $\left\| \mathbf{x} - \mathbf{x}^{(k)} \right\|$ es menor que **tol**.

Criterio de detención (cont.)

- El criterio de detención implica calcular $\|\mathbf{M}\|$, lo que en general es difícil. En lugar de ello, calculamos

$$m_k := \frac{\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|}{\|\mathbf{x}^{(k-1)} - \mathbf{x}^{(k-2)}\|}.$$

que constituye una estimación de $\|\mathbf{M}\|$. En efecto,

- **Lema.** Para $k = 2, 3, \dots$ se tiene que $m_k \leq \|\mathbf{M}\|$.

Demostración.

$$\begin{aligned} m_k &= \frac{\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|}{\|\mathbf{x}^{(k-1)} - \mathbf{x}^{(k-2)}\|} = \frac{\|[\mathbf{x}^{(k)} - \mathbf{x}] - [\mathbf{x}^{(k-1)} - \mathbf{x}]\|}{\|[\mathbf{x}^{(k-1)} - \mathbf{x}] - [\mathbf{x}^{(k-2)} - \mathbf{x}]\|} \\ &= \frac{\|\mathbf{e}^{(k)} - \mathbf{e}^{(k-1)}\|}{\|\mathbf{e}^{(k-1)} - \mathbf{e}^{(k-2)}\|} = \frac{\|\mathbf{M}[\mathbf{e}^{(k-1)} - \mathbf{e}^{(k-2)}]\|}{\|\mathbf{e}^{(k-1)} - \mathbf{e}^{(k-2)}\|} \leq \max_{\mathbf{y} \in \mathbb{R}^n: \mathbf{y} \neq \mathbf{0}} \frac{\|\mathbf{M}\mathbf{y}\|}{\|\mathbf{y}\|} = \|\mathbf{M}\|. \end{aligned}$$

- Así, proceso iterativo detenemos cuando:

$$\frac{m_k}{1 - m_k} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \text{tol.}$$

Métodos de Jacobi y Gauss-Seidel

- Se considera resolver un sistema $\mathbf{A}\mathbf{x} = \mathbf{b}$ con $a_{ii} \neq 0$, para $i = 1, \dots, n$.
Sea $\mathbf{x}^{(0)} = (x_1^{(0)}, \dots, x_n^{(0)})^t$ arbitrario y escribamos la matriz \mathbf{A} en la forma

$$\mathbf{A} = \mathbf{D} - \mathbf{E} - \mathbf{F},$$

donde $\mathbf{D} = \text{diag}(\mathbf{A})$, $-\mathbf{E}$ es la matriz triangular inferior de \mathbf{A} y $-\mathbf{F}$ es la matriz triangular superior de \mathbf{A} :

$$\mathbf{A} = \begin{pmatrix} \ddots & & -\mathbf{F} \\ & \mathbf{D} & \\ -\mathbf{E} & & \ddots \end{pmatrix}$$

- Notemos que tanto \mathbf{D} como $\mathbf{D} - \mathbf{E}$ son matrices inversibles, ya que $a_{ii} \neq 0$ para $i = 1, \dots, n$.

- El **método de Jacobi** corresponde al esquema iterativo general con

$$N := D \quad \text{y} \quad P := E + F.$$

- El **método de Gauss-Seidel** corresponde al esquema iterativo general con

$$N := D - E \quad \text{y} \quad P := F.$$

Convergencia del Método de Jacobi

- La matriz de iteración del método de Jacobi verifica:

$$M = D^{-1}(E + F) = \begin{pmatrix} 0 & -\frac{a_{12}}{a_{11}} & \cdots & \cdots & -\frac{a_{1n}}{a_{11}} \\ -\frac{a_{21}}{a_{22}} & 0 & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & 0 & -\frac{a_{n-1,n}}{a_{n-1,n-1}} \\ -\frac{a_{n1}}{a_{nn}} & \cdots & \cdots & -\frac{a_{nn-1}}{a_{nn}} & 0 \end{pmatrix}$$

- En este caso $\|M\|_{\infty} = \max_{1 \leq i \leq n} \left\{ \frac{1}{|a_{ii}|} \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \right\}.$

- Cuando A es de **diagonal dominante estricta**, es decir,

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n,$$

entonces $\|M\|_{\infty} < 1$ y el método de Jacobi resulta convergente.

Convergencia del Método de Gauss-Seidel

- El método de Gauss-Seidel es convergente cuando \mathbf{A} es de **diagonal dominante estricta**, es decir, cuando

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n.$$

- El método de Gauss-Seidel también es convergente cuando \mathbf{A} es **definida positiva**.

Resumen

- ▶ **Teorema.** Si A es de diagonal dominante estricta, entonces los métodos de Jacobi y de Gauss-Seidel convergen.
- ▶ **Observación.** Para una matriz arbitraria A , la convergencia de uno de estos métodos no implica la convergencia del otro.
- ▶ **Teorema.** Si A es simétrica y definida positiva, el método de Gauss-Seidel es convergente.
- ▶ **Observación.** Aunque A sea simétrica y definida positiva, el método de Jacobi puede ser divergente.

Convergencia de los métodos (cont.)

- **Ejemplo.** Para $s \in \mathbb{R}$, considere la matriz simétrica

$$\mathbf{A} = \begin{pmatrix} 1 & s & s \\ s & 1 & s \\ s & s & 1 \end{pmatrix},$$

cuyos valores propios son: $1 - s$ (con multiplicidad 2) y $1 + 2s$.

La matriz \mathbf{A} es definida positiva cuando $s \in (-1/2, 1)$ y es de diagonal dominante estricta para $s \in (-0.5, 0.5)$.

- Se resolvió el sistema $\mathbf{Ax} = \mathbf{b}$ para un par de valores de s , usando los métodos de Jacobi y de Gauss-Seidel, con $\mathbf{b} = (1, 1, 1)^t$ y $\mathbf{x}^{(0)} = (0.5, 0.5, 0.5)^t$.
- Para $s = 0.3$, Jacobi itera 37 veces y Gauss-Seidel 12 veces (en ambos casos se implementó el criterio de detención visto en clase, con una tolerancia de 10^{-8}).
Ambos métodos entregan como solución $\mathbf{x} = (0.6250, 0.6250, 0.6250)^t$, que es la solución exacta.

Convergencia de los métodos (cont.)

- Para $s = 0.8$, en las mismas condiciones anteriores, Gauss-Seidel converge en la iteración 53 a

$$(0.384\,615\,391\,735, 0.384\,615\,381\,035, 0.384\,615\,381\,784)^t$$

que difiere de la solución exacta

$$\mathbf{x} = (0.384\,615\,384\,615, 0.384\,615\,384\,615, 0.384\,615\,384\,615)^t$$

en menos de $\text{tol} = 10^{-8}$ en cada componente.

- En cambio, al cabo de 100 iteraciones Jacobi entrega

$$10^{19} \times (-1.862\,199\,431\,313, -1.862\,199\,431\,313, -1.862\,199\,431\,313)^t,$$

vector que no tiene ninguna relación con la solución del sistema.

Se nota claramente que, en este caso, el método de Jacobi diverge.