

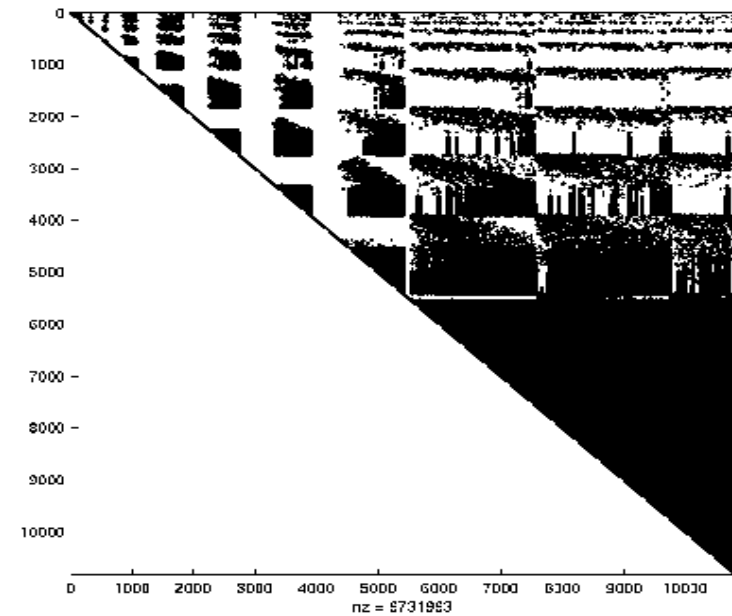
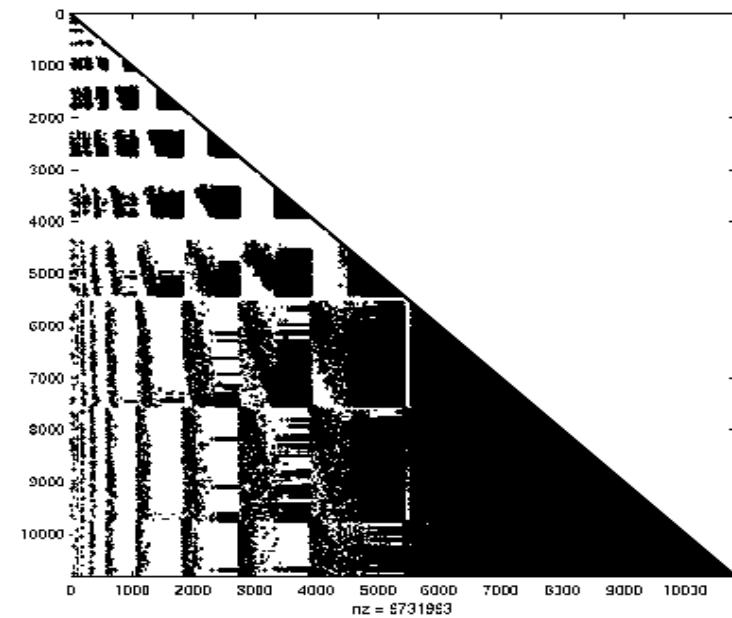
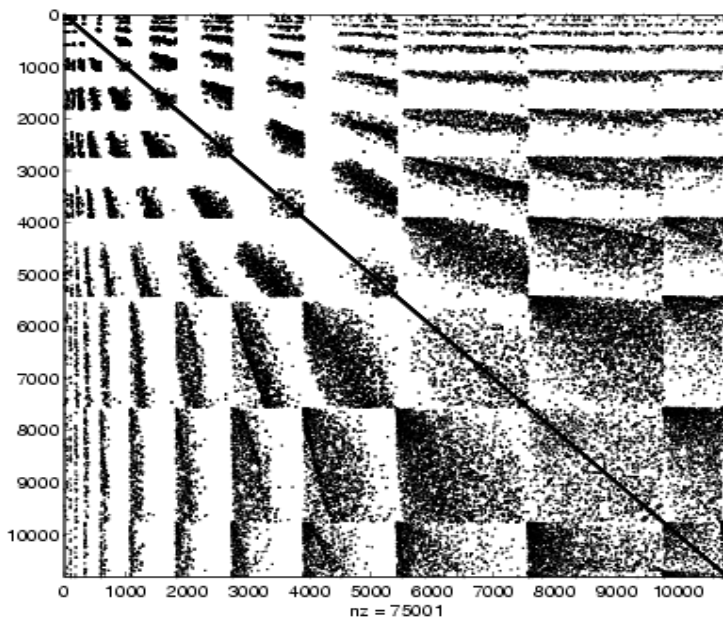
# Sistemas de Ecuaciones Lineales V

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

# Matrices dispersas

- Cuando la matriz  $A$  del sistema a resolver es dispersa, pero no banda, los métodos (directos) estudiados hasta ahora (eliminación de Gauss 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  $L$  y  $U$  en donde la matriz  $A$  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  $L$  y  $U$ , lo que puede llegar a hacer imposible aplicar estos métodos cuando  $A$  es de gran tamaño.
- Los métodos que estudiaremos en seguida, llamados **iterativos**, evitan el llenado y sus consecuencias, al trabajar resolviendo reiteradamente sistemas con matriz diagonal o triangular-dispersa.

# Llenado de matrices dispersas



## Llenado de matrices dispersas (cont.)

```
>> load data.0125.mat
```

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

```
>> whos
```

Name	Size	Bytes	Class
A	10821x10821	951580	sparse array
L	10821x10821	151325524	sparse array
U	10821x10821	151325524	sparse array
b	10821x1	129860	sparse array

```
Grand total is 25300219 elements using 303732496 bytes
```

## Esquema general

- Considere el sistema de ecuaciones

$$Ax = b,$$

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

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

- Si suponemos  $A = N - P$ , donde  $N$  debe ser invertible, entonces

$$Ax = b \iff Nx = Px + b \iff x = N^{-1}Px + N^{-1}b.$$

## Esquema general (cont.)

- Se usa la igualdad  $Nx = Px + b$  para definir un esquema general para construir la sucesión  $\{x^{(k)}\}$ .

- **Algoritmo del esquema general:**

Dado el vector inicial  $x^{(0)}$ ,

para  $k = 1, 2, \dots$  resolver:

$$| \quad Nx^{(k)} = Px^{(k-1)} + b,$$

hasta que se satisfaga un criterio de detención.

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

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

# Convergencia de métodos iterativos.

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

$$\rho(A) \leq \|A\|.$$

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

$$\|M\| < 1,$$

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

## Criterio de detención

- **Detención del proceso.** Cuando el proceso iterativo es convergente, éste se debe detener para un  $\mathbf{x}^{(k+1)}$  tal que  $\|e^{(k+1)}\| = \|\mathbf{x} - \mathbf{x}^{(k+1)}\| \leq \text{tol}$ , donde  $\text{tol}$  indica un nivel de tolerancia prefijado para el error.

- **Lema.** Para  $\|M\| < 1$ , se tiene que

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

- Un criterio de detención usual consiste en detener el proceso cuando

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

Sin embargo, **este criterio resulta muchas veces inadecuado!**

En efecto, si  $\frac{\|M\|}{1 - \|M\|} \gg 1$ , usualmente,

$$\frac{\|M\|}{1 - \|M\|} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| > \text{tol}.$$



## Criterio de detención (cont.)

- **Observación.** Al graficar

$$F(\mathbf{M}) := \frac{\|\mathbf{M}\|}{1 - \|\mathbf{M}\|},$$

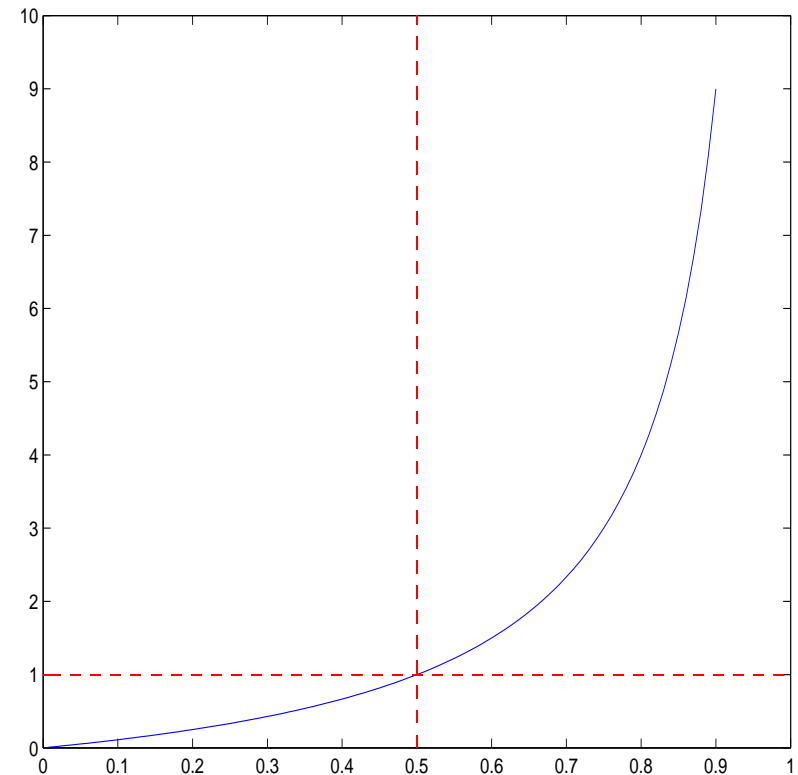
se puede ver que:

$$\|\mathbf{M}\| \leq \frac{1}{2} \implies F(\mathbf{M}) \leq 1,$$

$$\|\mathbf{M}\| > \frac{1}{2} \implies F(\mathbf{M}) > 1.$$

Luego, si  $\|\mathbf{M}\| > \frac{1}{2}$ , puede ser incorrecto detener el proceso cuando sólo se tiene

$$\left\| \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} \right\| \leq \text{tol}.$$



## Criterio de detención (cont.)

- El criterio de detención implica calcular  $\|\mathbf{M}\|$ , lo que en general es difícil. El siguiente lema indica una manera de estimar  $\|\mathbf{M}\|$ .
- **Lema.** Para  $k = 1, 2, \dots$  se tiene:

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

### Demostración.

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

- En el criterio de detención puede utilizarse  $m_k$  como una estimación de  $\|\mathbf{M}\|$ .

En tal caso, el proceso iterativo se detendrá cuando:

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

## Descomposición de una matriz

- Se considera resolver un sistema  $Ax = b$  con  $a_{ii} \neq 0$ , para  $i = 1, \dots, n$ .

Sea  $x^{(0)} = (x_1^{(0)}, \dots, x_n^{(0)})^t$  arbitrario y escribamos la matriz  $A$  en la forma

$$A = D - E - F,$$

donde  $D = \text{diag}(A)$ ,  $-E$  y  $-F$  son:

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

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

# Método de Jacobi

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

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

- Algoritmo de Jacobi:**

Dado el vector inicial  $\mathbf{x}^{(0)}$ ,

para  $k = 1, 2, \dots$  resolver:

$$D\mathbf{x}^{(k)} = (E + F)\mathbf{x}^{(k-1)} + \mathbf{b},$$

hasta que se satisfaga un criterio de detención.

- En la iteración  $k$ , el vector  $\mathbf{x}^{(k)}$  puede obtenerse **por componentes** como sigue:

Para  $i = 1, \dots, n$  :

$$x_i^{(k)} = \frac{1}{a_{ii}} \left[ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k-1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k-1)} \right].$$

## Método de Jacobi (cont.)

- La matriz de iteración del método de Jacobi verifica:

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

- Para la norma infinito de  $M$  se tiene que  $\|M\|_{\infty} = \max_{1 \leq i \leq n} \left\{ \frac{1}{|a_{ii}|} \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \right\}.$

- Cuando  $A$  es de **diagonal dominante estricta**, es decir, cuando se tiene que

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

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

# Método de Gauss-Seidel

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

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

La matriz de iteración es entonces  $M = (D - E)^{-1}F$ .

- Algoritmo de Gauss-Seidel:**

Dado el vector inicial  $\mathbf{x}^{(0)}$ ,

para  $k = 1, 2, \dots$ , resolver:

$$(D - E)\mathbf{x}^{(k)} = F\mathbf{x}^{(k-1)} + \mathbf{b},$$

hasta que se satisfaga un criterio de detención.

- En la iteración  $k$ , el vector  $\mathbf{x}^{(k)}$  puede obtenerse **por componentes** como sigue:

Para  $i = 1, \dots, n$ :

$$x_i^{(k)} = \frac{1}{a_{ii}} \left[ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k)} - \sum_{j=i+1}^n a_{ij} x_j^{(k-1)} \right].$$

Notemos que esto corresponde a aprovechar en el paso  $k$ , los valores  $x_j^{(k)}$  ya calculados.

# Convergencia de los métodos de Jacobi y de Gauss-Seidel

- **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  $s \in (-0.5, 1)$  y es de diagonal dominante estricta para  $s \in (-0.5, 0.5)$ .

- Se resolvió el sistema  $\mathbf{A}\mathbf{x} = \mathbf{b}$  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

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

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.