

## ▷ Lezione IV

# Metodi iterativi per sistemi lineari

**Libro** Capitolo 5, 5.9, 5.10, 5.11.

## Metodi iterativi per sistemi lineari

Dati una matrice  $A \in \mathbb{R}^{n \times n}$  e un vettore  $\mathbf{b} \in \mathbb{R}^n$ , vogliamo determinare  $\mathbf{x} \in \mathbb{R}^n$  tale che risolva il sistema lineare.

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

Un metodo iterativo costruisce una successione  $\{\mathbf{x}^k\}$ , dove ogni  $\mathbf{x}^k \in \mathbb{R}^n$ , che converge alla soluzione  $\mathbf{x}$  per  $k \rightarrow \infty$ . Abbiamo quindi

$$\mathbf{x}^0 \xrightarrow{\text{un passo}} \mathbf{x}^1 \xrightarrow{\text{un passo}} \mathbf{x}^2 \xrightarrow{\text{un passo}} \dots \xrightarrow{\text{un passo}} \mathbf{x}^k \xrightarrow{k \rightarrow \infty} \mathbf{x}.$$

Più formalmente possiamo introdurre la seguente definizione.

### Definizione 4.1 - Metodo iterativo convergente

Dato un vettore iniziale  $\mathbf{x}^0$ , un metodo iterativo si dice convergente se, detta  $\mathbf{x}$  la soluzione esatta del sistema lineare,

$$\lim_{k \rightarrow \infty} \mathbf{x}^k = \mathbf{x} \iff \lim_{k \rightarrow \infty} \mathbf{e}^k = \mathbf{0},$$

dove  $\mathbf{e}^k = \mathbf{x} - \mathbf{x}^k$  è l'errore commesso al passo  $k$ .

Una forma generale per i metodi iterativi è data dall'espressione seguente

$$\mathbf{x}^{k+1} = B\mathbf{x}^k + \mathbf{g},$$

dove  $B \in \mathbb{R}^{n \times n}$  è detta matrice di iterazione, che generalmente può dipendere da  $A$ , e  $\mathbf{g}$  un vettore che viene costruito partendo da  $A$  e  $\mathbf{b}$ . La matrice  $B$  e il vettore  $\mathbf{g}$  non possono essere presi in maniera arbitraria, infatti devono essere scelti in modo che il metodo sia consistente.

### Definizione 4.2 - Metodo iterativo consistente

Un metodo della forma

$$\mathbf{x}^{k+1} = B\mathbf{x}^k + \mathbf{g},$$

è detto consistente se  $B$  e  $\mathbf{g}$  sono tali che

$$\mathbf{x} = B\mathbf{x} + \mathbf{g},$$

ossia se la soluzione esatta soddisfa esattamente il metodo numerico.

Notiamo che se il metodo è consistente, se per un certo  $k$  abbiamo che  $\mathbf{x}^k = \mathbf{x}$  allora anche  $\mathbf{x}^{k+1} = B\mathbf{x} + \mathbf{g} = \mathbf{x}$ , ossia il metodo si "ferma" sulla soluzione esatta.

#### Teorema 4.1

Denotiamo con  $\mathbf{e}^k = \mathbf{x} - \mathbf{x}^k$  è l'errore commesso al passo  $k$ , ogni metodo consistente soddisfa

$$\mathbf{e}^{k+1} = B\mathbf{e}^k$$

*Proof.* A partire dalla definizione di  $\mathbf{e}^{k+1}$  sfruttiamo la consistenza del metodo per ottenere

$$\mathbf{e}^{k+1} = \mathbf{x} - \mathbf{x}^{k+1} = B\mathbf{x} + \mathbf{g} - B\mathbf{x}^k - \mathbf{g} = B(\mathbf{x} - \mathbf{x}^k) = B\mathbf{e}^k.$$

□

La matrice di iterazione  $B$  applicata all'errore al passo  $k$ ,  $\mathbf{e}^k$ , restituisce l'errore al passo successivo  $\mathbf{e}^{k+1}$ . Prima di introdurre il seguente teorema introduciamo cosa intendiamo per raggio spettrale della matrice, è un numero positivo che indichiamo con  $\rho(B)$  ed è dato da

$$\rho(B) = \max_{i=1,\dots,n} |\lambda_i(B)|,$$

con  $\lambda_i(B)$  autovalore  $i$ -esimo di  $B$ . Il raggio spettrale soddisfa la seguente proprietà

$$\lim_{k \rightarrow \infty} B^k = 0 \quad \Longleftrightarrow \quad \rho(B) < 1.$$

ed inoltre per ogni norma di matrice  $\|\cdot\|$  otteniamo che

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

#### Teorema 4.2 - Condizione necessaria e sufficiente di convergenza

Un metodo iterativo consistente della forma

$$\mathbf{x}^{k+1} = B\mathbf{x}^k + \mathbf{g},$$

è convergente se e solo se

$$\rho(B) < 1,$$

*Proof.* Mostriamo l'implicazione che se il raggio spettrale soddisfa  $\rho(B) < 1$  allora il metodo è convergente. Essendo per ipotesi un metodo consistente, allora vale

$$\mathbf{e}^{k+1} = B\mathbf{e}^k,$$

inoltre, ripetendo il ragionamento per  $k-1, k-2 \dots$  si ottiene che

$$\mathbf{e}^{k+1} = B^{k+1}\mathbf{e}^0.$$

Dato che  $\rho(B) < 1$  sappiamo che  $\lim_{k \rightarrow \infty} B^k = 0$ , quindi otteniamo che

$$\lim_{k \rightarrow \infty} \mathbf{e}^{k+1} = \lim_{k \rightarrow \infty} B^{k+1}\mathbf{e}^0 = \mathbf{0}.$$

Mostriamo ora l'implicazione inversa, ovvero vogliamo mostrare che se il raggio spettrale è maggiore di 1 allora il metodo non è convergente. Supponiamo che  $\rho(B) \geq 1$ , cioè esiste almeno un autovalore  $\lambda$  di  $B$  tale che  $|\lambda| \geq 1$ . Scegliamo  $\mathbf{x}^0$  tale che l'errore iniziale,  $\mathbf{e}^0 = \mathbf{x} - \mathbf{x}^0$ , sia uguale all'autovettore associato a  $\lambda$ . Abbiamo quindi

$$B\mathbf{e}^0 = \lambda\mathbf{e}^0$$

e otteniamo che l'errore al passo  $k$  si esprime come

$$\mathbf{e}^k = B\mathbf{e}^{k-1} = B^2\mathbf{e}^{k-2} = \dots = B^k\mathbf{e}^0 = \lambda^k\mathbf{e}^0$$

Quest'ultima espressione non può convergere a zero per  $k \rightarrow \infty$  dato che  $|\lambda| > 1$ . Quindi il raggio spettrale della matrice di iterazione per un metodo convergente deve necessariamente essere minore di 1.  $\square$

#### Esempio 4.1

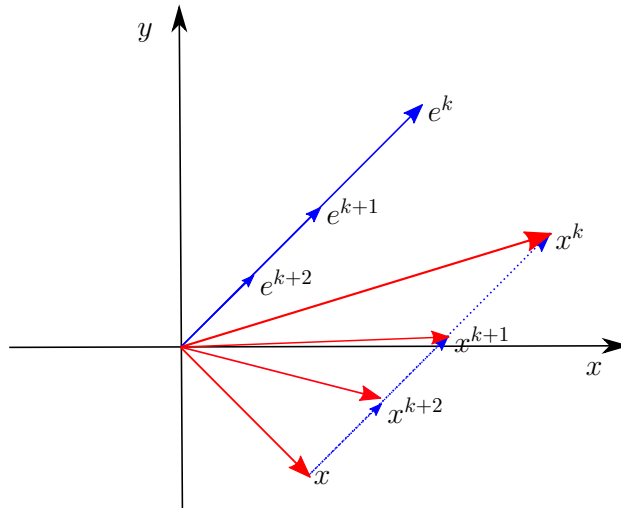
Sia data la matrice di iterazione  $B$  seguente

$$B = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{bmatrix},$$

chiaramente gli autovalori associati a  $B$  sono  $\lambda_1 = \lambda_2 = \frac{1}{2}$ , che implica che il raggio spettrale è dato da  $\rho(B) = \frac{1}{2}$ . Supponendo che l'errore al passo  $k$  sia dato dal vettore  $\mathbf{e}^k = [5, 7]^\top$  possiamo calcolare l'errore al passo  $k+1$  come

$$\mathbf{e}^{k+1} = B\mathbf{e}^k = \frac{1}{2} \begin{bmatrix} 5 \\ 7 \end{bmatrix}.$$

Graficamente, proseguendo con le iterazioni, otteniamo il seguente andamento



in cui si osserva che, nel caso in cui  $\rho(B) < 1$ , applicare la matrice  $B$  ad un vettore lo fa accorciare. Di conseguenza, se l'errore diminuisce, la soluzione numerica si avvicina a quella esatta.

Consideriamo ora la matrice  $A$  scomposta tramite il seguente splitting additivo, ovvero data  $P \in$

$\mathbb{R}^{n \times n}$  abbiamo che

$$A = P - (P - A), \quad (1.1)$$

dove la matrice  $P$  viene detta matrice preconditionatore. Deriviamo ora la matrice di iterazione  $B$  e il vettore  $\mathbf{g}$  data la scomposizione proposta, risulta

$$A\mathbf{x} = \mathbf{b} \rightarrow [P - (P - A)]\mathbf{x} = \mathbf{b} \rightarrow P\mathbf{x} = (P - A)\mathbf{x} + \mathbf{b} \rightarrow \mathbf{x} = P^{-1}(P - A)\mathbf{x} + P^{-1}\mathbf{b}$$

Se definiamo  $B = P^{-1}(P - A)$  e  $\mathbf{g} = P^{-1}\mathbf{b}$ , allora possiamo scrivere il metodo iterativo, che risulta consistente per costruzione, come

$$\mathbf{x}^{k+1} = B\mathbf{x}^k + \mathbf{g}.$$

Definiamo il residuo del problema al passo  $k$  come

$$\mathbf{r}^k = \mathbf{b} - A\mathbf{x}^k,$$

chiaramente il residuo è nullo nel caso in cui  $\mathbf{x}^k$  è uguale alla soluzione esatta del problema lineare  $\mathbf{x}$ . Abbiamo che

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

Ricordiamo che la formulazione di metodi diretti e iterativi per sistemi lineari ha lo scopo di evitare di dover invertire la matrice  $A$  per poi applicarla al sistema  $A^{-1}A\mathbf{x} = A^{-1}\mathbf{b}$  per ottenere  $\mathbf{x} = A^{-1}\mathbf{b}$ . Questa operazione avrebbe infatti un costo computazionale proibitivo. Tuttavia, notiamo che considerando la scomposizione (1.1), ad ogni passo  $k$  dobbiamo comunque risolvere un sistema lineare dato da

$$P\mathbf{x}^{k+1} = (P - A)\mathbf{x}^k + \mathbf{b}$$

quindi questo approccio ha senso solo se la risoluzione del sistema lineare in  $P$  è estremamente meno costosa che la risoluzione del sistema originale in  $A$ . Devo quindi scegliere  $P$  in modo che sia invertibile ma allo stesso tempo che il sistema lineare associato sia facile da risolvere. Diversi metodi iterativi sono caratterizzati da diverse scelte di  $P$ , sempre richiedendo che  $P$  sia invertibile e che il sistema lineare associato sia rapido da risolvere. Tipiche scelte di  $P$  sono matrice diagonale e matrice triangolare.

## Metodo di Jacobi

Il metodo di Jacobi definisce la matrice di preconditionamento  $P$  come la matrice diagonale  $D$ , che ha per diagonale la stessa della matrice  $A$ . Ovvero  $P = D = \text{diag}(A)$  dove

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ \dots & & & & \\ a_{n1} & & \dots & a_{n,n-1} & a_{nn} \end{bmatrix} \xrightarrow{D=\text{diag}(A)} D = \begin{bmatrix} a_{11} & 0 & \dots & 0 \\ 0 & a_{22} & 0 & \dots & 0 \\ \dots & & & & \\ 0 & & \dots & 0 & a_{nn} \end{bmatrix}$$

Abbiamo quindi che il metodo iterativo di Jacobi risulta scritto come

$$D\mathbf{x}^{k+1} = (D - A)\mathbf{x}^k + \mathbf{b} \rightarrow \mathbf{x}^{k+1} = D^{-1}(D - A)\mathbf{x}^k + D^{-1}\mathbf{b}.$$

Abbiamo sfruttato il fatto che  $D^{-1}$  è una matrice diagonale con elementi  $\frac{1}{a_{ii}}$  e inoltre, premoltiplicare per una matrice diagonale equivale a effettuare uno scaling delle righe. Possiamo quindi definire la matrice di iterazione  $B_J$  e il vettore  $\mathbf{g}$  associati al metodo di Jacobi abbiamo

$$\mathbf{g} = D^{-1}\mathbf{b}. \quad \text{e} \quad B_J = D^{-1}(D - A) = I - D^{-1}A = \begin{bmatrix} 0 & -\frac{a_{12}}{a_{11}} & \dots & -\frac{a_{1n}}{a_{11}} \\ -\frac{a_{21}}{a_{22}} & 0 & -\frac{a_{23}}{a_{22}} & \dots & -\frac{a_{2n}}{a_{22}} \\ \dots & & & & \\ -\frac{a_{n1}}{a_{nn}} & \dots & -\frac{a_{n,n-1}}{a_{nn}} & 0 \end{bmatrix}$$

Data l'espressione di  $B_J$  posso scrivere la soluzione del problema componente per componente, ovvero per il passo  $k$  come

$$x_i^{k+1} = \frac{1}{a_{ii}} \left( b_i - \sum_{j=1, j \neq i}^n a_{ij} x_j^k \right) \quad \text{per } i = 1, \dots, n$$

notiamo che l'elemento  $i$ -esimo del vettore  $\mathbf{x}^{k+1}$  è calcolabile indipendentemente dalle altri componenti, quindi l'algoritmo è facilmente parallelizzabile ottenendo la soluzione del problema ancora più velocemente.

Essendo un metodo iterativo della forma  $\mathbf{x}^{k+1} = B_J \mathbf{x}^k + \mathbf{g}$ , allora il metodo di Jacobi converge se e solo se  $\rho(B_J) < 1$ . Esistono condizioni solo sufficienti per la convergenza che sono molto più facili da verificare, in particolare abbiamo che, se  $A$  è una matrice a dominanza diagonale stretta (per righe o per colonne), allora  $|\lambda_i| < 1$  per ogni  $i$  e quindi il metodo converge.

### Metodo di Gauss-Seidel

Scomponiamo la matrice  $A$  come  $A = D - E - F$  dove

$$D = \begin{bmatrix} a_{11} & 0 & \dots & 0 \\ 0 & a_{22} & 0 & \dots & 0 \\ \dots & & & & \\ 0 & & \dots & 0 & a_{nn} \end{bmatrix} \quad E = \begin{bmatrix} 0 & 0 & \dots & 0 \\ -a_{21} & 0 & 0 & \dots & 0 \\ \dots & & & & \\ -a_{n1} & \dots & -a_{n,n-1} & 0 \end{bmatrix}$$

$$F = \begin{bmatrix} 0 & -a_{12} & \dots & -a_{1n} \\ 0 & 0 & -a_{23} & \dots & -a_{2n} \\ \dots & & & & \\ 0 & \dots & 0 & 0 \end{bmatrix}.$$

Il metodo di Gauss-Seidel pone la matrice di preconditionamento  $P = D - E$ , ossia la parte triangolare inferiore di  $A$ . Infatti, notiamo che  $D - E = T$  dove  $t_{ij} = a_{ij}$  se  $j \leq i$ ,  $t_{ij} = 0$  se  $j > i$ , e  $T - A = F$ . ottenere

$$P \mathbf{x}^{k+1} = (P - A) \mathbf{x}^k + \mathbf{b} \quad \rightarrow \quad (D - E) \mathbf{x}^{k+1} = F \mathbf{x}^k + \mathbf{b}.$$

La forma componente per componente del metodo è data da

$$x_i^{k+1} = \frac{1}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{k+1} - \sum_{j=i+1}^n a_{ij} x_j^k \right) \quad \text{per } i = 1, \dots, n.$$

Anche in questo caso, grazie alla struttura di  $P$ , ad ogni iterazione abbiamo un sistema semplice da risolvere, utilizzando una sostituzione in avanti. Dato che il calcolo di  $x_i^{k+1}$  richiede di aver già calcolato i valori di  $x_j^{k+1}$  per  $j = 1, \dots, i-1$  allora l'algoritmo di Gauss-Seidel è più difficile da parallelizzare rispetto al metodo di Jacobi. Tuttavia, il metodo di Gauss-Seidel spesso converge più velocemente rispetto al metodo di Jacobi, ossia raggiunge una certa soglia di tolleranza sull'errore con un numero minore di iterazioni.

La matrice di iterazione di Gauss-Seidel è data da

$$B_{GS} = (D - E)^{-1} F$$

(in alternativa,  $B_{GS} = T^{-1}(T - A)$ ) e sappiamo che se  $\rho(B_{GS}) < 1$  allora il metodo risulta convergente. Nel caso particolare in cui  $A$  è una matrice tridiagonale, ovvero se  $a_{ij} = 0$  per  $|i| > j + 1$ , allora abbiamo che  $\rho(B_{GS}) = [\rho(B_J)]^2$ . Quindi il metodo di Jacobi converge se e solo se il metodo di Gauss-Seidel converge, inoltre se convergono allora quest'ultimo converge più velocemente. Come negli altri metodi studiati, verificare che  $\rho(B_{GS}) < 1$  è solitamente molto costoso e quindi da evitare, ci affidiamo quindi a condizioni solo sufficienti per sapere se a priori il metodo di Gauss-Seidel converge o meno. Queste condizioni sono le seguenti:

- se la matrice  $A$  è a dominanza diagonale stretta per righe, o per colonne
- oppure se la matrice  $A$  è simmetrica e definita positiva.

## Metodo di Richardson

Consideriamo la scomposizione  $A = P - (P - A)$  usata precedentemente. Possiamo esprimere un'iterazione del generico metodo iterativo come

$$P\mathbf{x}^{k+1} = (P - A)\mathbf{x}^k + \mathbf{b} = P\mathbf{x}^k - A\mathbf{x}^k + \mathbf{b} = P\mathbf{x}^k + \mathbf{r}^k \quad \rightarrow \quad \mathbf{x}^{k+1} = \mathbf{x}^k + P^{-1}\mathbf{r}^k.$$

Il metodo di Richardson stazionario è definito generalizzando la precedente con un valore  $\alpha$  costante come

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha P^{-1}\mathbf{r}^k,$$

mentre il metodo di Richardson dinamico o non-stazionario è definito dato un valore  $\alpha^k$  variabile per ogni passo  $k$  come

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha^k P^{-1}\mathbf{r}^k.$$

Per i metodi di Richardson la matrice di iterazione risulta quindi

$$B = I - \alpha_k P^{-1}A.$$

### Teorema 4.3

Supponiamo che  $P$  sia non singolare e che  $P^{-1}A$  sia definita positiva, quindi con  $\lambda_i > 0$ . Ordiniamo gli autovalori di  $P^{-1}A$  nel seguente modo

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0.$$

Allora il metodo di Richardson stazionario converge se e solo se

$$0 < \alpha < \frac{2}{\lambda_1}$$

inoltre esiste un valore ottimale di  $\alpha$  detto  $\alpha_{opt}$  di valore

$$\alpha_{opt} = \frac{2}{\lambda_1 + \lambda_n}$$

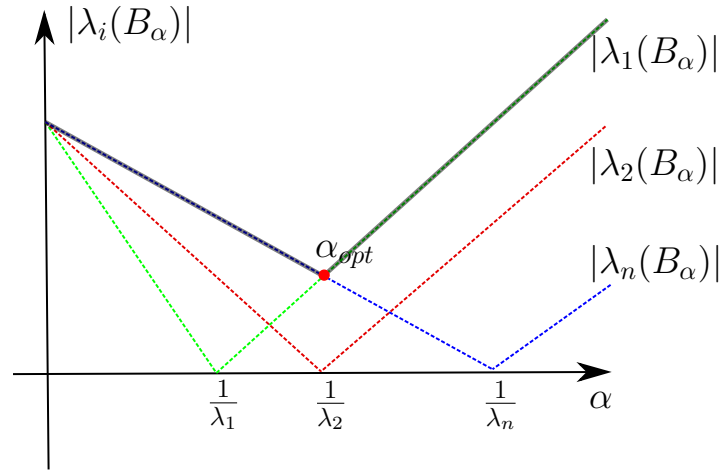
tale che il raggio spettrale della matrice di iterazione  $B$  è minimo, pari a

$$\rho_{opt} = \min_{\alpha} \rho(B_{\alpha}) = \frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n}.$$

*Proof.* Data l'espressione di  $B$  abbiamo che i suoi autovalori si possono esprimere come  $\lambda_i(B_{\alpha}) = 1 - \alpha\lambda_i$ . Sappiamo che un metodo iterativo è convergente quando  $|\lambda_i(B_{\alpha})| < 1$  per tutti gli  $i = 1, \dots, n$ . Nel nostro caso abbiamo

$$|1 - \alpha\lambda_i| < 2 \quad \Rightarrow \quad 0 < \alpha < \frac{2}{\lambda_1}$$

Verifichiamo ora il valore di  $\alpha_{opt}$ : come detto gli autovalori della matrice  $B_{\alpha}$  sono dati da  $\lambda_i(B_{\alpha}) = 1 - \alpha\lambda_i$  per  $i = 1, \dots, n$ . Ricordando che  $\rho(B_{\alpha}) = \max_{i=1, \dots, n} |\lambda_i(B_{\alpha})|$ , allora consideriamo il seguente grafico che rappresenta ogni  $|\lambda_i(B_{\alpha})|$  al variare di  $\alpha$



Essendo il raggio spettrale il massimo tra tutti gli autovalori, allora per ogni  $\alpha$  i suoi valori sono dati dall'involuppo evidenziato in grigio nel grafico. Il valore ottimale di  $\alpha$ , ovvero quello che rende minimo il raggio spettrale di  $B_\alpha$ , è identificato dall'intersezione tra  $1 - \alpha\lambda_n = \alpha\lambda_1 - 1$ . Otteniamo quindi il valore ottimale ricercato e il corrispondente raggio spettrale.  $\square$

## Metodo del gradiente

Consideriamo una matrice  $A \in \mathbb{R}^{n \times n}$  simmetrica e definita positiva, si vuole risolvere il problema di determinare  $\mathbf{x} \in \mathbb{R}^n$  tale che

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

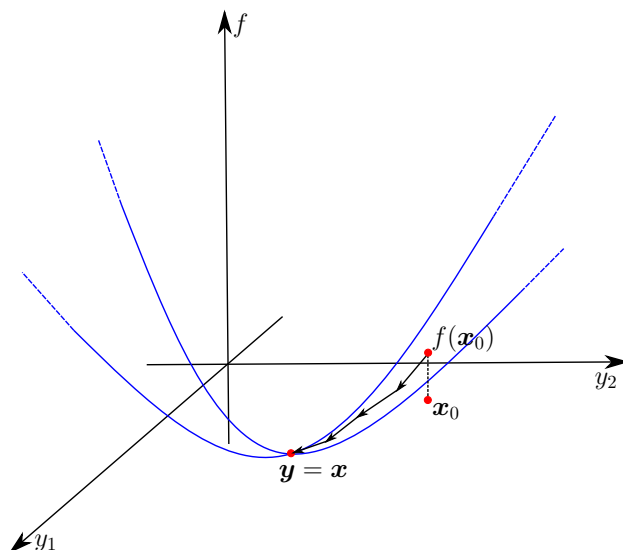
Notiamo che se  $A$  è simmetrica e definita positiva allora anche  $A^{-1}$  è simmetrica e definita positiva, infatti  $\lambda(A) = 1/\lambda(A^{-1})$  ovvero

$$A\mathbf{w} = \lambda\mathbf{w} \quad \rightarrow \quad A^{-1}A\mathbf{w} = \lambda A^{-1}\mathbf{w} \quad \rightarrow \quad \lambda^{-1}\mathbf{w} = A^{-1}\mathbf{w}.$$

Costruiamo ora la funzione  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  quadratica associata alla risoluzione del sistema lineare data da

$$f(\mathbf{y}) = \frac{1}{2}\mathbf{y}^\top A\mathbf{y} - \mathbf{y}^\top \mathbf{b},$$

dove  $\mathbf{y} \in \mathbb{R}^n$  è un vettore. Graficamente essa può essere rappresentata come un paraboloide se  $\mathbf{y} \in \mathbb{R}^2$



Vogliamo mostrare che la soluzione del sistema lineare equivale alla ricerca del punto di minimo della funzione  $f$ . Volendo calcolare il minimo di  $f$  dobbiamo risolvere il problema di determinare  $\mathbf{y}$  tale che

$$\mathbf{0} = \nabla f(\mathbf{y}) = A\mathbf{y} - \mathbf{b}.$$

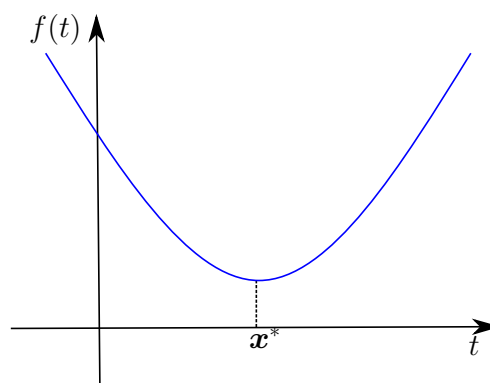
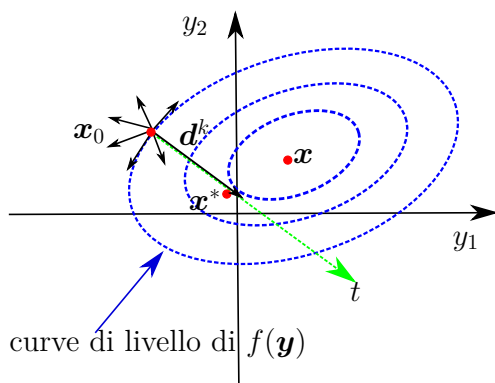
ovvero calcolare il minimo  $\mathbf{y}$  di  $f$  è equivalente a determinare  $\mathbf{x}$  che risolvere il sistema lineare  $A\mathbf{x} = \mathbf{b}$ . Abbiamo quindi l'equivalenza  $\mathbf{y} = \mathbf{x}$ . Il metodo del gradiente costruisce una sequenza di soluzioni che converge verso il minimo di  $f$ .

La funzione  $f$  può essere definita in modi diversi, ottenendo lo stesso punto di minimo ma un valore minimo differente (in altre parole, si può utilizzare un paraboloide traslato).

Il metodo del gradiente fa parte dei metodi di discesa, tali per cui partendo da un  $\mathbf{x}^0$  ad ogni  $k \geq 1$  prima si ricerca una direzione di discesa  $\mathbf{d}^k$ , successivamente viene scelto un passo  $\alpha_k$  con cui avanzare, ed infine si aggiorna il valore  $\mathbf{x}^{k+1}$

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k \quad \text{tale per cui} \quad f(\mathbf{x}^{k+1}) < f(\mathbf{x}^k).$$

Un metodo di discesa è caratterizzato principalmente dalla scelta della direzione  $\mathbf{d}^k$  e dal passo  $\alpha_k$ .

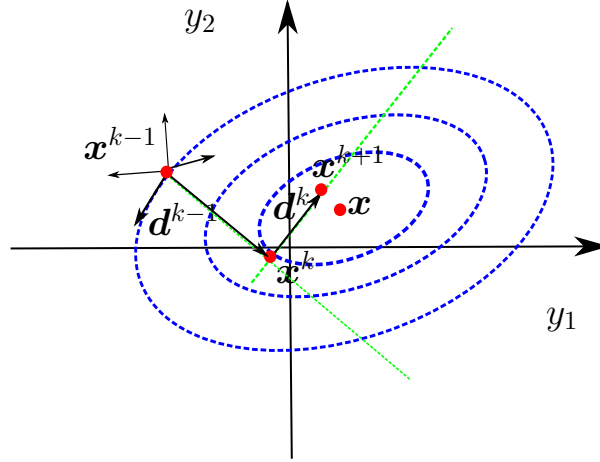


Data una direzione di discesa  $\mathbf{d}^k$  allora possiamo calcolare il valore del passo ottimale, richiedendo che la funzione  $f$  sia minima lungo la direzione data da  $\mathbf{d}^k$ . Tale funzione rimane quadratica ma dipendente solo dall'ascissa curvilinea  $t$ , essa è definita come  $f(\mathbf{x}^k + t\mathbf{d}^k)$ . Il valore minimo si



realizza per un certo  $t$  pari al passo  $\alpha_k$  che soddisfa

$$\alpha_k = \frac{(\mathbf{d}^k)^\top \mathbf{r}^k}{(\mathbf{d}^k)^\top A \mathbf{d}^k}.$$



#### Approfondimento 4.1

Dimostriamo che il valore di  $\alpha_k$  corrisponde al minimo di  $f(\mathbf{x}^k + t\mathbf{d}^k)$ :

$$\begin{aligned} f(\mathbf{x}^k + t\mathbf{d}^k) &= f(t) = \frac{1}{2}(\mathbf{x}^k + t\mathbf{d}^k)^\top A(\mathbf{x}^k + t\mathbf{d}^k) - (\mathbf{x}^k + t\mathbf{d}^k)^\top \mathbf{b} \\ \frac{df(t)}{dt} &= (\mathbf{d}^k)^\top A(\mathbf{x}^k + t\mathbf{d}^k) - (\mathbf{d}^k)^\top \mathbf{b} = 0 \\ t &= \frac{(\mathbf{d}^k)^\top (\mathbf{b} - A\mathbf{x}^k)}{(\mathbf{d}^k)^\top A \mathbf{d}^k} = \frac{(\mathbf{d}^k)^\top \mathbf{r}^k}{(\mathbf{d}^k)^\top A \mathbf{d}^k}. \end{aligned}$$

Per quanto riguarda  $\mathbf{d}^k$  una scelta che sembrerebbe ottimale è quella di massima discesa partendo dal punto  $\mathbf{x}^k$ , ovvero data da

$$\mathbf{d}^k = -\nabla f(\mathbf{x}^k) = -(A\mathbf{x}^k - \mathbf{b}) = \mathbf{b} - A\mathbf{x}^k = \mathbf{r}^k,$$

quindi se scelgo  $\mathbf{d}^k = \mathbf{r}^k$ , il metodo di discesa è detto metodo del gradiente. Inoltre possiamo calcolare il nuovo valore del residuo che è dato da

$$\mathbf{r}^{k+1} = \mathbf{b} - A\mathbf{x}^{k+1} = \mathbf{b} - A(\mathbf{x}^k + \alpha_k \mathbf{r}^k) = \mathbf{r}^k - \alpha_k A \mathbf{r}^k,$$

dove il termine  $A\mathbf{r}^k$ , che computazionalmente costa  $\mathcal{O}(n^2)$ , è già stato calcolato per determinare il valore di  $\alpha_k$ . L'algoritmo è il seguente.

**Algoritmo 4.1 - Metodo del gradiente**


---

**Data:** Dato  $\mathbf{x}^0$ , calcolo  $\mathbf{r}^0 = \mathbf{b} - A\mathbf{x}^0$ ;  
**while** (*criterio di convergenza*) **do**  
     $\alpha_k \leftarrow \frac{(\mathbf{r}^k)^\top \mathbf{r}^k}{(\mathbf{r}^k)^\top A\mathbf{r}^k};$   
     $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \alpha_k \mathbf{r}^k;$   
     $\mathbf{r}^{k+1} \leftarrow \mathbf{r}^k - \alpha_k A\mathbf{r}^k;$   
     $k \leftarrow k + 1;$   
**end**

---

È possibile dimostrare che per il metodo del gradiente otteniamo che le direzioni generate sono, a coppie, ortogonali ovvero

$$\mathbf{d}^k \cdot \mathbf{d}^{k+1} = 0 \quad \Longleftrightarrow \quad \mathbf{r}^k \cdot \mathbf{r}^{k+1} = 0 \quad k = 0, \dots,$$

tuttavia non è garantito che la direzione  $k+1$ -esima risulti ortogonale a tutte le direzioni da  $k-1$  a 0. Infatti sostituendo l'espressione prima di  $\mathbf{r}^{k+1}$  e successivamente di  $\alpha_k$  otteniamo

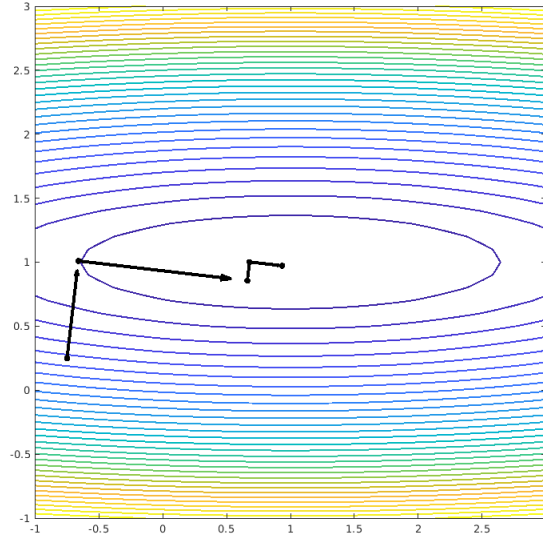
$$\mathbf{r}^k \cdot \mathbf{r}^{k+1} = \mathbf{r}^k \cdot (\mathbf{r}^k - \alpha_k A\mathbf{r}^k) = \mathbf{r}^k \cdot \mathbf{r}^k - \frac{(\mathbf{r}^k)^\top \mathbf{r}^k}{(\mathbf{r}^k)^\top A\mathbf{r}^k} (\mathbf{r}^k)^\top A\mathbf{r}^k = 0.$$

**Esempio 4.2**

Consideriamo il sistema  $A\mathbf{x} = \mathbf{b}$  con

$$A = \begin{bmatrix} \frac{1}{20} & 0 \\ 0 & 1 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} \frac{1}{20} \\ 1 \end{bmatrix}.$$

Ovviamente la soluzione esatta del sistema è  $\mathbf{x} = [1, 1]^\top$ . Supponiamo di scegliere come *guess* iniziale  $\mathbf{x}^0 = [-\frac{3}{4}, \frac{1}{4}]^\top$  e osserviamo in figura le iterazioni del metodo del gradiente sovrapposte alle linee di livello della funzione  $f$ . Osserviamo che le direzioni di discesa sono perpendicolari alle linee di livello e, a due a due fra loro.



#### Teorema 4.4 - Convergenza del metodo del gradiente

Se  $A \in \mathbb{R}^{n \times n}$  è simmetrica e definita positiva, allora il metodo del gradiente converge alla soluzione del sistema  $A\mathbf{x} = \mathbf{b}$  per ogni dato iniziale  $\mathbf{x}^0$ . Inoltre, vale la seguente stima dell'errore  $\mathbf{e}^k = \mathbf{x} - \mathbf{x}^k$  data da

$$\|\mathbf{e}^k\|_A \leq C^k \|\mathbf{e}^0\|_A \quad \text{con} \quad C = \frac{\text{cond}(A) - 1}{\text{cond}(A) + 1}$$

dove  $\|\mathbf{x}\|_A = \sqrt{\mathbf{x}^\top A \mathbf{x}}$  (detta norma dell'energia o  $A$ -norma) e  $C$  è detto fattore di abbattimento dell'errore.

Ricordiamo che, sotto le ipotesi del teorema, il numero di condizionamento della matrice  $A$  è dato da  $\text{cond}(A) = \lambda_{\max}/\lambda_{\min}$ . Inoltre, dato che il fattore di abbattimento dell'errore è minore strettamente di 1 allora abbiamo che

$$\lim_{k \rightarrow \infty} \|\mathbf{e}^k\|_A = 0.$$

Notiamo inoltre che  $C$  è costante e dipende solamente dalla matrice  $A$ , risulta tanto più piccolo quanto il numero di condizionamento di  $A$  è vicino a 1, mentre  $C$  è tanto più grande e vicino a 1 quanto più  $\text{cond}(A) \gg 1$ . Possiamo quindi concludere che se una matrice ha condizionamento  $\text{cond}(A) \approx 1$ , la matrice  $A$  è detta ben condizionata e il metodo del gradiente converge velocemente. Al contrario se  $\text{cond}(A) \gg 1$  allora  $A$  è detta matrice mal condizionata e il metodo del gradiente converge molto lentamente dato che  $C \approx 1$ . Infatti se consideriamo in cui  $A \in \mathbb{R}^{2 \times 2}$  allora l'eccentricità degli ellissi, curve di livello di  $f$ , è legata al condizionamento di  $A$ : più  $\text{cond}(A)$  è alto e più l'ellisse risulta schiacciata e quindi le direzioni generate dal gradiente, essendo in questo specifico caso tutte ortogonali tra di loro, convergono molto più lentamente al minimo di  $f$ .

## Metodo del gradiente coniugato

Il metodo del gradiente genera le direzioni  $\mathbf{d}^k$  partendo dal gradiente della funzione, ci chiediamo se tale direzioni siano ottimali rispetto al problema considerato. Osserviamo che

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k = \mathbf{x}^{k-1} + \alpha_{k-1} \mathbf{d}^{k-1} + \alpha_k \mathbf{d}^k = \dots = \mathbf{x}^0 + \sum_{j=0}^k \alpha_j \mathbf{d}^j$$

e quindi il termine  $\mathbf{x}^{k+1} - \mathbf{x}^0$  è una combinazione lineare delle  $k$  direzioni di discesa considerate  $\mathbf{d}^j = \mathbf{r}^j$ . Per  $k = n$  tali vettori non sono, in generale, linearmente indipendenti ma anche se lo fossero non è detto che i coefficienti  $\alpha_j$  trovati corrispondano allo sviluppo di  $\mathbf{x} - \mathbf{x}^0$  su tale base. Quindi potrebbero servire molte più iterazioni di  $n$  per poter ottenere tale richiesta.

Il metodo del gradiente coniugato garantisce la scelta ottimale delle direzioni di discesa  $\mathbf{d}^k$ , tali da essere  $A$ -ortogonali o  $A$ -coniugate tra loro. Tale condizione permette di avere le direzioni di discesa tutte linearmente indipendenti tra loro essendo  $A$  non singolare.

### Definizione 4.3 - $A$ -ortogonalità

Data una matrice  $A \in \mathbb{R}^{n \times n}$  simmetrica, due vettori  $\mathbf{x}$  e  $\mathbf{y} \in \mathbb{R}^n$  sono detti  $A$ -ortogonali se

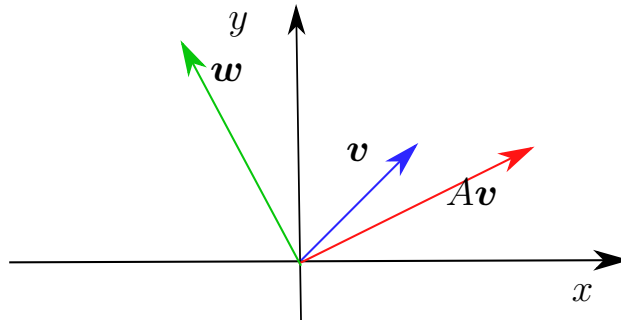
$$(\mathbf{A}\mathbf{y})^\top \mathbf{x} = 0 \quad \Longleftrightarrow \quad \mathbf{y}^\top \mathbf{A}\mathbf{x} = 0.$$

### Esempio 4.3

Consideriamo la matrice

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

simmetrica definita positiva, con autovalori  $\lambda_1 = 2$ ,  $\lambda_2 = 1$  e autovettori paralleli a  $\mathbf{i}$  e  $\mathbf{j}$ . Consideriamo il vettore  $\mathbf{v} = [1, 1]^\top$ . Abbiamo che  $\mathbf{A}\mathbf{v} = [2, 1]^\top$ . Un vettore  $\mathbf{y}$   $A$ -ortogonale a  $\mathbf{v}$  è tale che  $2y_1 + y_2 = 0$ , per esempio  $\mathbf{w} = [-1, 2]^\top$ . Notiamo che  $\mathbf{v} \cdot \mathbf{w} \neq 0$ . Osserviamo inoltre che se consideriamo un vettore  $\mathbf{v} = [2, 0]^\top$ , parallelo agli autovettori di  $A$ , allora  $\mathbf{A}\mathbf{v}$  risulta parallelo a  $\mathbf{v}$  e in questo caso il vettore  $A$ -ortogonale è anche semplicemente ortogonale.



Abbiamo che il metodo del gradiente coniugato, in aritmetica esatta, converge alla soluzione esatta in al più  $n$  iterazioni. Quindi, seppure il metodo del gradiente coniugato è stato derivato come metodo iterativo esso convergerà in un numero finito di passi. Nel caso in cui  $n$  sia grande, lo schema viene comunque bloccato prima di fare  $n$  passi ottenendo quindi una soluzione approssimata.

**Algoritmo 4.2 - Metodo del gradiente coniugato**


---

**Data:** Dato  $x^0$ , calcolo  $r^0 = b - Ax^0$  e impongo  $d^0 = r^0$ ;  
**while** (*criterio di convergenza*) **do**  
 $\alpha_k \leftarrow \frac{(d^k)^\top r^k}{(d^k)^\top A d^k};$   
 $x^{k+1} \leftarrow x^k + \alpha_k d^k;$   
 $r^{k+1} \leftarrow r^k - \alpha_k A d^k;$   
 $\beta_k \leftarrow \frac{(A d^k)^\top r^{k+1}}{(A d^k)^\top d^k};$   
 $d^{k+1} \leftarrow r^{k+1} - \beta_k d^k;$   
 $k \leftarrow k + 1;$   
**end**

---

Notiamo che la direzione di discesa coincide con il residuo solo alla prima iterazione, successivamente viene "corretta" sottraendo le direzioni precedentemente utilizzate. Le direzioni  $\{d^{k+1}\}$  calcolate sono quindi  $A$ -ortogonali tra loro.

**Teorema 4.5 - Convergenza del metodo del gradiente coniugato**

Sia  $A \in \mathbb{R}^{n \times n}$  una matrice simmetrica e definita positiva, allora il metodo del gradiente coniugato converge alla soluzione del sistema  $Ax = b$ , per ogni  $x^0 \in \mathbb{R}^n$  in al più  $n$  iterazioni. Inoltre, vale la stima

$$\|e^k\|_A \leq \frac{2c^k}{1+c^{2k}} \|e^0\|_A \quad \text{dove} \quad c = \frac{\sqrt{\text{cond}(A)} - 1}{\sqrt{\text{cond}(A)} + 1}.$$

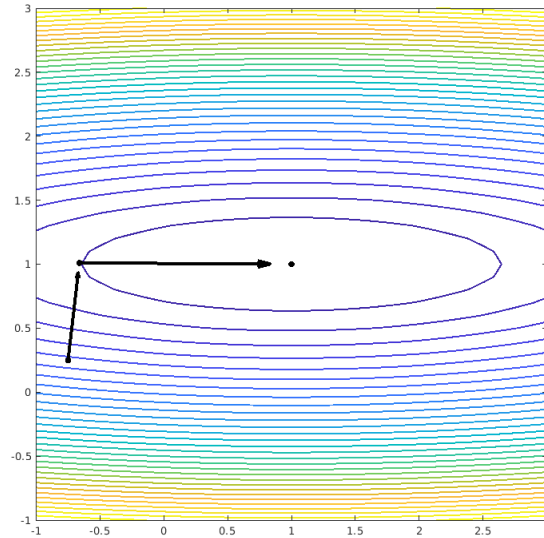
Abbiamo anche in questo caso che il fattore di abbattimento dell'errore converge a zero e lo fa più velocemente che nel caso del metodo del gradiente.

**Esempio 4.4**

Consideriamo il sistema dell'esempio precedente,  $Ax = b$  con

$$A = \begin{bmatrix} \frac{1}{20} & 0 \\ 0 & 1 \end{bmatrix} \quad b = \begin{bmatrix} \frac{1}{20} \\ 1 \end{bmatrix}.$$

Ovviamente la soluzione esatta del sistema è  $x = [1, 1]^\top$ . Supponiamo di scegliere come *guess* iniziale  $x^0 = [-\frac{3}{4}, \frac{1}{4}]^\top$  e osserviamo in figura le iterazioni del metodo del gradiente coniugato sovrapposte alle linee di livello della funzione  $f$ . Osserviamo che le direzioni di discesa sono tali per cui la soluzione viene ottenuta in sole 2 iterazioni.



Sia il metodo del gradiente che il metodo del gradiente coniugato possono essere migliorati usando un preconditionatore, cioè trovando una matrice  $P$  invertibile, detto preconditionatore, tale che

$$\text{cond}(P^{-1}A) \ll \text{cond}(A).$$

#### Approfondimento 4.4

Esistono diverse scelte possibili per la matrice  $P$ . Lo scopo è trovare  $P$  tale che il condizionamento di  $P^{-1}A$  sia il più piccolo possibile, e in questo senso la scelta più efficace sarebbe  $P = A$ : infatti  $A^{-1}A = I$  ha condizionamento unitario. Tuttavia, la matrice  $A$  non è in generale facile da invertire, e, se lo fosse, non avremmo bisogno di ricorrere a un metodo iterativo per risolvere il sistema. Un buon preconditionatore  $P$  sarà quindi i) facile da invertire ii) "simile" alla matrice  $A$ . Consideriamo un esempio:

$$A = \begin{bmatrix} 2 & 1 \\ 1 & 100 \end{bmatrix}$$

Il condizionamento di  $A$  è  $\text{cond}(A) = 50.26$  dato da  $\frac{\lambda_{\max}}{\lambda_{\min}} = \frac{100.0102}{1.9898}$ . Proviamo a utilizzare come  $P$  la diagonale di  $A$ : una matrice diagonale è infatti facile da invertire. Otteniamo

$$P^{-1}A = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{100} \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 1 & 100 \end{bmatrix} = \begin{bmatrix} 1 & \frac{1}{2} \\ \frac{1}{100} & 1 \end{bmatrix}$$

e il condizionamento di  $P^{-1}A$  è ridotto a  $\text{cond}(P^{-1}A) = 2.63$ .

Vogliamo quindi risolvere il seguente problema preconditionato utilizzando il metodo del gradiente o del gradiente coniugato seguente

$$A\mathbf{x} = \mathbf{b} \quad \xrightarrow{\text{moltiplico per } P^{-1}} \quad P^{-1}A\mathbf{x} = P^{-1}\mathbf{b}.$$

Avendo quindi il fattore di abbattimento dell'errore molto minore del sistema di partenza, possiamo ottenere la soluzione con meno passi e quindi in meno tempo.

I metodi di discesa generano una successione  $\{\mathbf{x}^k\}$  tale che converge alla soluzione esatta  $\mathbf{x}$  per  $k \rightarrow \infty$ . Chiaramente dobbiamo terminare l'algoritmo utilizzando un opportuno criterio d'arresto. Una possibilità è un criterio basato sul residuo  $\mathbf{r}^k$ , data una tolleranza  $\epsilon$  abbiamo

$$\frac{\|\mathbf{r}^k\|}{\|\mathbf{b}\|} < \epsilon.$$

Tale criterio è efficace se il numero di condizionamento di  $A$  non è troppo elevato, infatti abbiamo che

$$\frac{\|\mathbf{e}^k\|}{\|\mathbf{x}\|} \leq \text{cond}(A) \frac{\|\mathbf{r}^k\|}{\|\mathbf{b}\|}.$$

Un'altra possibilità è basare il criterio d'arresto sull'incremento della soluzione calcolata, ovvero

$$\frac{\|\mathbf{x}^k - \mathbf{x}^{k+1}\|}{\|\mathbf{b}\|} < \epsilon.$$