



RISOLUZIONE DI SISTEMI LINEARI: METODI ITERATIVI (Relazione tra residuo ed errore e principi dei metodi iterativi)

In questa lezione introdurremo il concetto di residuo per un sistema lineare e analizzeremo la relazione fondamentale tra residuo ed errore. Successivamente, esamineremo i principi che guidano la costruzione dei metodi iterativi e, infine, presenteremo alcuni esempi concreti di tali metodi.

Iniziamo osservando che esiste una relazione algebrica semplice tra l'errore commesso nella risoluzione di un sistema lineare e il residuo ad esso associato. Consideriamo un generico sistema lineare

$$Ax = b.$$

Sia x^* la soluzione approssimata ottenuta mediante un generico metodo iterativo e sia r^* il corrispondente residuo, definito come il vettore

$$r^* = b - Ax^*.$$

Dunque, prendiamo il termine noto b e sottraiamo il vettore di Ax^* (ovvero il prodotto matrice A per vettore x^* , e quindi ancora un vettore). Definiamo la differenza $b - Ax^*$ come il vettore residuo. Se x^* coincidesse con la soluzione esatta $x^* = x$, allora il residuo risulterebbe nullo

$$r^* = b - Ax^* = 0,$$

perché la soluzione esatta soddisfa il problema $Ax = b$. Pertanto, sulla soluzione esatta il residuo è pari a zero. Tuttavia, quando la soluzione approssimata differisce da quella esatta, il residuo sarà generalmente diverso da zero. Il nostro obiettivo è quindi analizzare la relazione tra il residuo e l'errore. Ricordando che l'errore è definito come la differenza tra la soluzione esatta x e la soluzione calcolata (o approssimata) x^*

$$e^* = x - x^*,$$

possiamo ricavare che

$$r^* = b - Ax^* = Ax - Ax^* = A(x - x^*) = Ae^*,$$

dove abbiamo sostituito $b = Ax$, da cui

$$r^* = Ae^*,$$

Dunque, possiamo trovare il residuo, applicando la matrice A all'errore rispetto alla soluzione approssimata $e^* = x - x^*$.

Inoltre, invertendo la matrice A posso trovare l'errore applicando A^{-1} al residuo

$$e^* = A^{-1}r^*.$$

Abbiamo dunque trovato la relazione che intercorre fra residuo ed errore.

Vogliamo sfruttare questa relazione per ottenere stime dell'errore in funzione del residuo, che saranno poi utilizzate nell'analisi dei metodi iterativi. Considerando la norma della relazione tra errore e residuo $e^* = A^{-1}r^*$, otteniamo

$$\|e^*\| = \|A^{-1}r^*\|.$$

Per la proprietà delle norme abbiamo che (la norma del prodotto è minore uguale al prodotto delle norme)

$$\|e^*\| = \|A^{-1}r^*\| \leq \|A^{-1}\| \|r^*\|.$$

Moltiplicando e dividendo per la norma di A otteniamo



$$\|e^*\| = \|A^{-1}r^*\| \leq \|A^{-1}\| \|A\| \frac{\|r^*\|}{\|A\|},$$

e ricordando che la norma di A per la norma di A^{-1} è il numero di condizionamento di A

$$K(A) = \|A^{-1}\| \|A\|,$$

troviamo che

$$\|e^*\| = \|A^{-1}r^*\| \leq \|A^{-1}\| \|A\| \frac{\|r^*\|}{\|A\|} = K(A) \frac{\|r^*\|}{\|A\|}.$$

Abbiamo quindi trovato che l'errore in norma è minore o uguale del numero di condizionamento di A per la norma del residuo, diviso per la norma di A

$$\|e^*\| \leq K(A) \frac{\|r^*\|}{\|A\|}.$$

Se volessimo tradurre questa relazione per l'errore relativo

$$\frac{\|x - x^*\|}{\|x\|} = \frac{\|e^*\|}{\|x\|},$$

basterebbe dividere la relazione precedentemente trovata $\|e^*\| \leq K(A) \frac{\|r^*\|}{\|A\|}$ per $\|x\|$ (ovvero, la norma della soluzione esatta) a sinistra e a destra, ottenendo

$$\frac{\|e^*\|}{\|x\|} \leq K(A) \frac{\|r^*\|}{\|A\| \|x\|}.$$

Ricordando che $Ax = b$ e quindi $\|Ax\| = \|b\|$ da cui $\|b\| = \|Ax\| \leq \|A\| \|x\|$ e dato che questa quantità è a denominatore possiamo usarla come maggiorazione dell'ultima relazione trovata, ovvero

$$\frac{\|e^*\|}{\|x\|} \leq K(A) \frac{\|r^*\|}{\|A\| \|x\|} \leq K(A) \frac{\|r^*\|}{\|b\|}.$$

Quindi l'errore relativo è minore o uguale del numero di condizionamento per il residuo relativo $\frac{\|r^*\|}{\|b\|}$

$$\frac{\|e^*\|}{\|x\|} \leq K(A) \frac{\|r^*\|}{\|b\|}.$$

Infatti, possiamo associare alla norma di b una misura della dimensione del residuo e quindi questo possiamo interpretare $\frac{\|r^*\|}{\|b\|}$ come residuo relativo.

Il numero di condizionamento della matrice A svolge, anche in questo contesto, il ruolo di possibile amplificatore del residuo nella determinazione dell'errore. In conclusione, se $K(A)$ è relativamente piccolo e, quindi, se la matrice A è ben condizionata, a un residuo di piccola entità corrisponderà un errore ridotto. Questo risultato è particolarmente incoraggiante, poiché in molti casi avremo la possibilità di controllare quantitativamente i residui, mentre, in generale, l'errore sulla soluzione non è direttamente accessibile, dato che la soluzione esatta non è nota a priori.

Di conseguenza, se il residuo è noto e la matrice è ben condizionata, si può immediatamente concludere che anche l'errore sarà contenuto, a condizione che il residuo sia piccolo. Al contrario, se la matrice A è mal condizionata, ossia se $K(A)$ assume valori elevati, può accadere che, pur in presenza di un residuo di piccola entità, l'errore risulti comunque grande.

Vediamo a questo proposito un esempio su un semplice sistema 2×2

$$Ax = b,$$



dove

$$A = \begin{pmatrix} 1.2969 & 0.8648 \\ 0.2161 & 0.1441 \end{pmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0.8642 \\ 0.1440 \end{bmatrix}.$$

La soluzione di questo sistema è facilmente ricavabile e risulta

$$\mathbf{x} = \begin{bmatrix} 2 \\ -2 \end{bmatrix}.$$

Supponendo di aver calcolato una soluzione \mathbf{x}^* uguale a

$$\mathbf{x}^* = \begin{bmatrix} 0.991 \\ -0.487 \end{bmatrix},$$

l'errore assoluto che si commette è un errore significativo

$$\mathbf{e}^* = \begin{bmatrix} 1.01 \\ -1.52 \end{bmatrix},$$

e quindi l'errore relativo

$$\frac{\|\mathbf{e}^*\|_2}{\|\mathbf{x}\|_2} \approx 39\%,$$

quindi un errore grande. Invece, il residuo associato a \mathbf{x}^* è estremamente piccolo

$$\mathbf{r}^* = \mathbf{b} - A\mathbf{x}^* = \begin{bmatrix} -10^{-8} \\ 10^{-8} \end{bmatrix}.$$

Dunque, in corrispondenza della soluzione \mathbf{x}^* , a un residuo estremamente piccolo può corrispondere un errore significativo. Questo rappresenta un caso in cui, anche in un sistema 2×2 , si osserva come residui di piccola entità possano essere associati a errori decisamente elevati. Sulla base della teoria generale sviluppata in precedenza, ci si aspetta che la causa di questo comportamento sia il mal condizionamento della matrice A . In effetti, in questo caso A^{-1} si può calcolare facilmente a mano ed ha questa espressione

$$A^{-1} = \begin{pmatrix} 0.1441 & -0.8648 \\ -0.2161 & 1.2969 \end{pmatrix} 10^8.$$

Quindi, il numero di condizionamento di A , calcolato nella norma infinito, è

$$K_{\infty}(A) = \|A^{-1}\|_{\infty} \|A\|_{\infty} \approx 3.3 \cdot 10^8.$$

Abbiamo dunque individuato un esempio di matrice fortemente mal condizionata, per la quale è possibile determinare un opportuno termine noto \mathbf{b} tale che la soluzione del sistema associato presenti un residuo estremamente piccolo, mentre l'errore commesso risulti invece significativamente elevato.

Esaminiamo ora i principi fondamentali che guidano la costruzione dei metodi iterativi. Ricordiamo che l'obiettivo dei metodi iterativi è quello di generare una successione di vettori $\mathbf{x}^{(k)}$ tali che il limite per k che tende all'infinito sia uguale a \mathbf{x} , essendo \mathbf{x} la soluzione esatta del sistema $A\mathbf{x} = \mathbf{b}$

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}.$$

Dal punto di vista pratico, non è desiderabile proseguire le iterazioni k indefinitamente, ma si preferisce arrestare il processo al minimo indice $k = n$, per cui l'errore risulti inferiore o uguale a una tolleranza ϵ fissata a priori

$$\|\mathbf{x}^{(n)} - \mathbf{x}\| \leq \epsilon.$$

Pertanto, possiamo decidere di fermarci non appena l'errore raggiunga un valore minore o uguale a una tolleranza predefinita, che potrebbe essere ad esempio 10^{-2} , 10^{-3} , 10^{-6} .



Vediamo ora come costruire metodi iterativi, ossia quale formula ci permette di generare vari tipi di metodi. La strategia chiave consiste nello scomporre la matrice A del sistema come differenza di due matrici P ed N

$$A = P - N.$$

È importante precisare che tale approccio non corrisponde a una fattorizzazione. Nella fattorizzazione, infatti, A viene rappresentata come prodotto di matrici, ad esempio nella fattorizzazione LU, mentre qui A viene espressa come differenza di due matrici.

L'unica ipotesi che facciamo è che P sia non singolare. Allora, da $Ax = b$, e considerando la scomposizione $A = P - N$, segue che

$$Ax = (P - N)x = Px - Nx = b,$$

e dunque

$$Px = Nx + b,$$

dove nell'ultimo passaggio abbiamo semplicemente portato Nx a secondo membro. In questa relazione $Px = Nx + b$, troviamo la soluzione del sistema sia a sinistra che a destra dell'uguale. Pertanto, possiamo pensare di impostare un processo iterativo, in cui, dato un valore iniziale $x^{(0)}$ scelto arbitrariamente, si possa generare una successione di vettori $x^{(k)}$ così costruita

Assegnato $x^{(0)}$, si generi una successione $\{x^{(k)}\}$ risolvendo

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

Quindi, supponendo noto $x^{(0)}$, per $k = 1$, calcoliamo

$$Px^{(1)} = Nx^{(0)} + b, \quad k = 1.$$

Noto $x^{(1)}$, possiamo calcolare $x^{(2)}$ (per $k = 2$), e poi $x^{(3)}$, $x^{(4)}$ e così via. In generale, per ogni $k \geq 1$, possiamo ottenere il vettore $x^{(k)}$ in funzione del vettore precedente $x^{(k-1)}$. In questo modo, abbiamo costruito una formula generale che ci consente di generare iterativamente una successione di vettori, seguendo la struttura del problema iniziale.

A questo punto, è opportuno esaminare se tale metodo iterativo così costruito sia in grado di convergere alla soluzione esatta del sistema. Poiché P e N sono scelte arbitrarie, con l'unica condizione che P sia non singolare, è possibile costruire una famiglia infinita di metodi iterativi. Tuttavia, non sempre è possibile garantire che questi metodi convergano, ossia che la successione $\{x^{(k)}\}$ si avvicini alla soluzione esatta x del sistema lineare. Pertanto, è necessario esaminare i principi teorici alla base della convergenza dei metodi iterativi.

Ripartendo dalla formula

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

introduciamo l'errore al passo k tra la soluzione calcolata $x^{(k)}$ e la soluzione esatta x

$$e^{(k)} = x^{(k)} - x.$$

Sottraendo membro a membro le seguenti due equazioni

$$Px = Nx + b$$

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

otteniamo

$$Px^{(k)} - Px = Nx^{(k-1)} - Nx,$$

da cui

$$P(x^{(k)} - x) = N(x^{(k-1)} - x).$$



Ricordando che $\mathbf{e}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}$, segue che

$$P\mathbf{e}^{(k)} = N\mathbf{e}^{(k-1)}.$$

Infine, moltiplicando a sinistra e destra per P^{-1} , otteniamo

$$\mathbf{e}^{(k)} = P^{-1}N\mathbf{e}^{(k-1)}.$$

Abbiamo ottenuto una relazione ricorsiva che lega l'errore al passo $k - 1$ con l'errore al passo k . In particolare, tale relazione stabilisce che l'errore al passo k può essere espresso come il prodotto (matrice per vettore) di una matrice $B = P^{-1}N$ per il vettore errore al passo precedente $k - 1$

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

La matrice B è quindi definita come

$$B = P^{-1}N.$$

Da questa definizione si evince il motivo per cui abbiamo richiesto che P fosse non singolare. Infatti, è evidente che, se vogliamo risolvere

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

per calcolare $\mathbf{x}^{(k)}$, è necessario essere certi che il sistema sia risolubile, ovvero che P sia non singolare. La matrice $B = P^{-1}N$ è chiamata matrice di interazione associata alla scomposizione (o splitting) $A = P - N$.

È possibile verificare ricorsivamente (partendo dalla relazione $\mathbf{e}^{(k)} = B\mathbf{e}^{(k-1)}$) che esiste una relazione che lega l'errore al passo k con l'errore al passo iniziale, ovvero

$$\mathbf{e}^{(k)} = B^k\mathbf{e}^{(0)}.$$

Abbiamo concluso questa prima parte di analisi affermando che l'errore al passo k del metodo iterativo può essere espresso come l'errore al passo iniziale ($k = 0$) moltiplicato per la matrice di iterazione B , e più precisamente come il prodotto della matrice B elevata alla potenza k -esima. Ricordiamo che la potenza k -esima di una matrice non è altro che il prodotto della matrice per sé stessa ripetuto k volte.

A questo punto, desideriamo comprendere quali conclusioni possiamo trarre riguardo alla convergenza del metodo iterativo. È infatti fondamentale che il metodo iterativo sia convergente per ogni scelta possibile del vettore iniziale $\mathbf{x}^{(0)}$, e quindi per ogni possibile errore iniziale $\mathbf{e}^{(0)}$. In altre parole, vogliamo garantire che la convergenza si verifichi indipendentemente dal vettore iniziale scelto, il che equivale a richiedere che il limite di $\mathbf{e}^{(k)}$ per k che tende all'infinito sia uguale zero

$$\text{Convergenza } \forall \mathbf{x}^{(0)} \Leftrightarrow \lim_{k \rightarrow \infty} \mathbf{e}^{(k)} = 0.$$

Ciò implica che l'errore, al progressivo aumentare dei passi iterativi, deve tendere a zero, garantendo così la convergenza della successione $\{\mathbf{x}^{(k)}\}$ verso la soluzione esatta del sistema.

Poiché abbiamo visto che $\mathbf{e}^{(k)} = B^k\mathbf{e}^{(0)}$, allora, per avere convergenza del metodo iterativo, dovremmo avere che la matrice B^k per k che tende all'infinito tende a 0

$$\lim_{k \rightarrow \infty} B^k = 0.$$

Questo implica che le potenze k -esime della matrice B all'infinito tendono alla matrice nulla

$$B^k \rightarrow 0 \quad \text{per } k \rightarrow \infty.$$

Tale condizione è vera se e solo se la norma di B (per una qualsiasi norma) è minore di uno

$$\|B\| < 1.$$



Inoltre, ciò è equivalente a richiedere che il raggio spettrale della matrice B è minore di 1
$$\|B\| < 1 \quad \Leftrightarrow \quad \rho(B) < 1.$$

È infatti possibile verificare che il raggio spettrale $\rho(B)$ è sempre minore o uguale alla norma di B per ogni norma possibile

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

Pertanto, possiamo concludere una condizione sufficiente per la convergenza di un metodo iterativo per la risoluzione di sistema lineari è che il raggio spettrale della matrice B sia minore di 1

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

Inoltre, la convergenza sarà tanto più rapida quanto più piccola è la norma della matrice di iterazione B (o, equivalentemente, quanto più piccolo è il raggio spettrale di B).

In sintesi, abbiamo introdotto una famiglia generale di metodi iterativi per la risoluzione di sistemi lineari attraverso lo splitting della matrice A , ossia $A = P - N$. Questo genera una successione di vettori $\mathbf{x}^{(k)}$ che soddisfano la relazione $P\mathbf{x}^{(k)} = N\mathbf{x}^{(k-1)} + \mathbf{b}$ ($k \geq 1$). Per garantire la convergenza del metodo iterativo, è necessario che esista una norma della matrice di iterazione B tale che $\|B\| < 1$, o equivalentemente che il raggio spettrale di B sia minore di 1 ($\rho(B) < 1$).