

Che cos'è la precisione di macchina in un sistema floating-point $F(b, t, L, U)$ e come si calcola? (si ricavi il valore)

Definiamo arrotondamento a t cifre di un numero reale scritto in notazione floating-point

$$x = \text{sign}(x)(0, d_1 d_2 \dots d_t \dots) \cdot b^p$$

il numero

$$fl^t(x) = \text{sign}(x)(0, d_1 d_2 \dots \tilde{d}_t) \cdot b^p$$

dove la mantissa è stata arrotondata alla t -esima cifra

$$\tilde{d}_t = \begin{cases} d_t & \text{se } d_{(t+1)} < \frac{b}{2} \\ d_t + 1 & \text{se } d_{(t+1)} \geq \frac{b}{2} \end{cases}$$

Definiamo:

$$\text{Errore Relativo} \leftarrow \frac{\overbrace{|x - fl^t(x)|}^{\text{Errore Assoluto}}}{|x|} \quad \text{per } x \neq 0$$

Stimiamo il numeratore

$$\begin{aligned} |x - fl^t(x)| &= b^p \cdot \overbrace{|(0, d_1 d_2 \dots d_t \dots) - (0, d_1 d_2 \dots \tilde{d}_t)|}^{\text{Errore di arrotondamento a } t \text{ cifre dopo la virgola} \leq \frac{b^{-t}}{2}} \\ &\leq b^p \cdot \frac{b^{-t}}{2} = \frac{b^{p-t}}{2} \end{aligned}$$

Notiamo subito un aspetto: l'errore dipende da p , cioè dall'ordine di grandezza del numero (in base b).

Stimiamo da sopra $\frac{1}{|x|}$, ovvero da sotto $|x|$:

$$|x| = (0, d_1 d_2 \dots d_t \dots) \cdot b^p$$

Poiché $d_1 \neq 0$, p fissato, il minimo valore della mantissa è $0,100\dots = b^{-1}$. Quindi:

$$|x| \geq b^{-1} \cdot b^p = b^{p-1} \iff \frac{1}{|x|} \leq \frac{1}{b^{p-1}}$$

Otteniamo

$$\frac{|x - fl^t(x)|}{|x|} \leq \frac{\frac{b^{p-t}}{2}}{b^{p-1}} = \frac{b^{p-t+1-p}}{2} = \frac{b^{1-t}}{2} = \varepsilon_M$$

Perché moltiplicazione e addizione sono operazioni stabili? Si ricavano le stime dell'errore

Dati $x, y \in \mathbb{R}$ e $*$ operazione binaria fondamentale, si ha che in $F(b, t, L, U)$ l'operazione risulta: $x * y = fl^t(fl^t(x) * fl^t(y))$

Dati $x, y \neq 0$, gli errori relativi su tali quantità saranno rispettivamente

$\varepsilon_x = \frac{|x - \tilde{x}|}{|x|}$ e $\varepsilon_y = \frac{|y - \tilde{y}|}{|y|}$ mentre quello sul risultato dell'operazione

sarà: $\varepsilon_{x*y} = \frac{|x*y - \tilde{x}*\tilde{y}|}{|x*y|}$

Un'operazione è stabile se l'ordine di grandezza di ε_{x*y} è vicino a quello di ε_x e ε_y

► **MOLTIPLICAZIONE**

$$\begin{aligned} \varepsilon_{x*y} &= \frac{|x*y - \tilde{x}\tilde{y}|}{|x*y|} = \frac{|x*y - \tilde{x}y + \tilde{x}y - \tilde{x}\tilde{y}|}{|x*y|} = \frac{|(x-\tilde{x})y + (\tilde{x}-\tilde{y})\tilde{y}|}{|x*y|} \leq \frac{|(x-\tilde{x})y| + |(\tilde{x}-\tilde{y})\tilde{y}|}{|x*y|} = \\ &= \frac{|x-\tilde{x}| |y|}{|x| |y|} + \frac{|\tilde{x}-\tilde{y}| |\tilde{y}|}{|x| |y|} = \varepsilon_x \frac{|y|}{|x|} + \varepsilon_y \frac{|\tilde{y}|}{|x|} \Rightarrow \varepsilon_{x*y} \approx \varepsilon_x + \varepsilon_y \quad \text{STABILE} \end{aligned}$$

► ADDIZIONE

$$\varepsilon_{x+y} = \frac{|(x+y) - (\tilde{x} + \tilde{y})|}{|x+y|} = \frac{|x - \tilde{x} + y - \tilde{y}|}{|x+y|} \leq \frac{|x - \tilde{x}|}{|x+y|} + \frac{|y - \tilde{y}|}{|x+y|} = \frac{|x|}{|x+y|} \varepsilon_x + \frac{|y|}{|x+y|} \varepsilon_y = w_1 \varepsilon_x + w_2 \varepsilon_y$$

$w_1 \varepsilon_x + w_2 \varepsilon_y \geq \varepsilon_{x+y}$ ma $w_1, w_2 \leq 1$ quindi STABILE

2.2 Divisione

La divisione è la moltiplicazione per il reciproco $\frac{x}{y} = x \cdot \frac{1}{y}$.

Analizzando quindi l'operazione di reciproco

$$\varepsilon_{\frac{1}{y}} = \frac{\left| \frac{1}{y} - \frac{1}{\tilde{y}} \right|}{\left| \frac{1}{y} \right|} = \frac{\left| \frac{\tilde{y} - y}{y\tilde{y}} \right|}{\left| \frac{1}{y} \right|} = \frac{|\tilde{y} - y|}{|y|} \cdot \frac{|y|}{|\tilde{y}|} \approx \varepsilon_y \quad \left(\text{questo perché } \frac{|\tilde{y} - y|}{|y|} = \varepsilon_y \right)$$

Perché la sottrazione tra numeri approssimati può essere instabile? (si ricavi la stima dell'errore e si faccia un esempio)

► SOTTRAZIONE

$$\varepsilon_{x-y} = \frac{|(x-y) - (\tilde{x} - \tilde{y})|}{|x-y|} = \frac{|x - \tilde{x} + \tilde{y} - y|}{|x-y|} \leq \frac{|x - \tilde{x}|}{|x-y|} + \frac{|\tilde{y} - y|}{|x-y|} = \frac{|x|}{|x-y|} \varepsilon_x + \frac{|y|}{|x-y|} \varepsilon_y = w_1 \varepsilon_x + w_2 \varepsilon_y$$

$w_1 \varepsilon_x + w_2 \varepsilon_y \geq \varepsilon_{x-y}$ ma $w_1, w_2 > 1$ quindi INSTABILE

Es.:

Con $x = 0,100017$, $y = -0,100014$, $\bar{x} = Fl^5(x) = 0,10002$, $\bar{y} = Fl^5(y) = -0,10001$, si ha:

$$\varepsilon_{x+y} = \frac{|(x+y) - (\bar{x} + \bar{y})|}{|x+y|} = \frac{|0,000003 - 0,00001|}{|0,000003|} = \frac{0,000009}{0,000003} = 2,3$$

Quindi l'errore relativo sarebbe del $233,3\%$.

Perché il residuo non pesato non può essere non essere una buona stima dell'errore nel metodo di bisezione? (si ricavi la stima del residuo pesato in modo rigoroso).

Vogliamo stimare l'errore di bisezione, applicato nelle seguenti ipotesi:

$$\left. \begin{array}{l} f \in C^1[a, b] \\ \{x_n\} \in [c, d] \subseteq [a, b] \\ f'(x) \neq 0, \forall x \in [c, d] \end{array} \right\} \Rightarrow e_n = |x_n - \xi| = \frac{|f(x_n)|}{|f'(z_n)|}, \quad n \geq n_0, \quad z_n \in \begin{cases} (x_n, \xi) \\ (\xi, x_n) \end{cases}$$

C^1 indica derivabile 1 volta con derivata continua.

Dimostriamolo utilizzando il teorema del valor medio

$$\text{Sia } f \in C[a, b] \text{ derivabile in } [a, b] \Rightarrow \exists z \in [a, b] : \frac{f(b) - f(a)}{b - a} = f'(z)$$

Consideriamo il caso $\xi < x_n$ (se $x_n < \xi$ la dimostrazione è analoga)

$$f(x_n) - f(\xi) = f'(z_n)(x_n - \xi), \quad z_n \in (\xi, x_n)$$

con $f(\xi) = 0$, cioè

$$|f(x_n)| = |f'(z_n)||x_n - \xi|$$

che si può riscrivere come

$$e_n = |x_n - \xi| = \frac{|f(x_n)|}{|f'(z_n)|}$$

Osserviamo che:

- e_n è un "residuo pesato"
- $f'(x) \neq 0 \Rightarrow$ zero è semplice
- e_n è una stima a posteriori (serve aver calcolato x_n)

Siccome non conosciamo z_n , diamo delle stime pratiche dell'errore:

- Se è noto che $|f'(x)| \geq k > 0 \Rightarrow e_n = \frac{|f(x_n)|}{|f'(z_n)|} \leq \frac{|f(x_n)|}{k}$
- Se f' è nota, per n abbastanza grande si ha

$$\underbrace{f'(x_n) \approx f'(z_n)}_{\approx f'(\xi)} \Rightarrow e_n \approx \frac{|f(x_n)|}{|f'(z_n)|}$$

- Se f' non è nota, si può approssimare con

$$f'(z_n) \approx \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}, \quad \text{per } n \text{ abbastanza grande}$$

Convergenza metodo di bisezione

Il metodo di bisezione si basa sull'applicazione iterativa del Teorema degli zeri di funzioni continue:

Se $f(x) \in C[a, b]$ e $f(a)f(b) < 0$ (cioè f cambia segno) allora

$$\exists \xi : f(\xi) = 0, \quad \xi \in (a, b)$$

Il procedimento consiste nel passare da $[a_n, b_n] \rightarrow [a_{n+1}, b_{n+1}]$ in cui uno degli estremi è diventato il punto medio

$$x_n = \frac{a_n + b_n}{2}$$

A meno che per qualche n non risulti $f(x_n) = 0$, si tratta di un processo infinito che ci permette di costruire tre successioni $\{a_n\}, \{b_n\}, \{x_n\}$ tali che:

- $|\xi - a_n|, |\xi - b_n| \leq b_n - a_n = \frac{b-a}{2^n}$
- $|\xi - x_n| < \frac{b_n - a_n}{2} = \frac{b-a}{2^{n+1}}$

È semplice dimostrare che tutte e tre le successioni convergono ad uno zero $\xi \in (a, b)$

- $0 \leq |\xi - a_n|, |\xi - b_n| < \frac{b-a}{2^n} \xrightarrow{n \rightarrow \infty} 0 \xRightarrow{\text{Teor. Carabinieri}} |\xi - a_n|, |\xi - b_n| \rightarrow 0, \quad n \rightarrow \infty$
- $0 \leq |\xi - x_n| < \frac{b-a}{2^{n+1}} \Rightarrow |\xi - x_n| \rightarrow 0, \quad n \rightarrow \infty$

Si dimostri la convergenza del metodo di Newton nel caso strettamente convesso o concavo

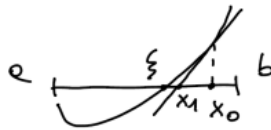
Metodo di Newton: Linearizzare iterativamente la funzione con la tangente nel punto

$$\begin{cases} y = 0 \\ y = f(x_n) + f'(x_n)(x - x_n) \end{cases} \Rightarrow x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

Convergenza metodo di Newton:

$$\begin{cases} f \in C^2[a, b] \\ f(a)f(b) < 0 \\ f''(x) > 0 \quad \forall x \in [a, b] \\ x_0 : f(x_0)f''(x_0) > 0 \end{cases} \Rightarrow \begin{array}{l} \text{Il metodo di Newton è ben definito (cioè } f'(x_n) \neq 0) \\ \text{e converge all'unico zero } \xi \text{ di } f \text{ in } [a, b] \end{array}$$

Si considera il caso:



- $f(a) < 0, f(b) > 0$
- $f''(x) > 0 \quad \forall x \in [a, b]$
- $x_0 \in [a, b]$

Dimostriamo come prima cosa: $x_n \in (\xi, b] \Rightarrow x_{n+1} \in (\xi, b]$

f è esattamente convessa \Rightarrow La tangente sta "sotto al grafico" $\forall x \in [a, b]$
 \Rightarrow La tangente in un punto $\in (\xi, b]$ interseca l'asse x "a destra" di ξ

Dimostriamo quindi: $x_{n+1} < x_n$ (cioè $\{x_n\}$ è decrescente)

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} > 0$$

Poiché per $x_n \in (\xi, b]$ si ha $f(x_n) > 0$. Inoltre $f'(x_n) > 0$ in $(\xi, b]$ altrimenti per avere uno zero f'' in $(\xi, b]$ dovrebbe cambiare segno.

Abbiamo quindi che $\{x_n\}$ è una successione decrescente, con $x_n > \xi \quad \forall n$.

Allora

$$\exists \lim_{n \rightarrow \infty} x_n = \inf \{x_n\} = \eta \quad \text{con} \quad \eta \geq \xi$$

Infine

$$\begin{aligned} \eta &= \lim x_{n+1} = \lim \left(x_n - \frac{f(x_n)}{f'(x_n)} \right) \\ &= \lim x_n - \lim \frac{f(x_n)}{f'(x_n)} \\ &= \lim x_n - \frac{\lim f(x_n)}{\lim f'(x_n)} \\ &= \lim x_n - \frac{f(\lim x_n)}{f'(\lim x_n)} \leftarrow \lim x_n = \eta \\ &= \eta - \frac{f(\eta)}{f'(\eta)} \end{aligned}$$

Quindi

$$\eta = \eta - \frac{f(\eta)}{f'(\eta)} \quad \text{con} \quad f'(\eta) \neq 0 \Rightarrow \frac{f(\eta)}{f'(\eta)} = 0 \Rightarrow f(\eta) = 0 \Rightarrow \eta = \xi$$

Velocità (ordine) di convergenza del metodo di Newton

Applicando Taylor centrato in x_n e calcolata in ξ con il resto di Lagrange del secondo ordine:

$$f(\xi) = f(x_n) + f'(x_n)(\xi - x_n) + \frac{f''(z_n)}{2}(\xi - x_n)^2 \quad \text{con } z_n \in \text{int}(x_n, \xi) \subset [c, d]$$

$$-\frac{f(x_n)}{f'(x_n)} = \xi - x_n + \frac{f''(z_n)}{2f'(x_n)}(\xi - x_n)^2 \quad \text{ma } -\frac{f(x_n)}{f'(x_n)} = x_{n+1} - x_n$$

$$x_{n+1} - x_n = \xi - x_n + \frac{f''(z_n)}{2f'(x_n)}(\xi - x_n)^2$$

$$e_{n+1} = |x_{n+1} - \xi| = c_n e_n^2 \quad \text{con } c_n = \frac{1}{2} \frac{f''(z_n)}{f'(x_n)}$$

ma c_n è succ. limitata:

$$c_n = \frac{\max |f''(x)|}{\min |f'(x)|} \cdot \frac{1}{2} = c \quad \forall x \in [c, d] \Rightarrow e_{n+1} \leq c e_n^2 \quad \text{e } \lim \frac{e_{n+1}}{e_n^2} = c \quad \text{quindi}$$

il metodo converge con velocità quadratica

Perché il metodo di Newton per zeri semplici ha ordine di convergenza almeno 2? Quando ha ordine esattamente 2? (si dimostri la relazione fondamentale che lega e_{n+1} ed e_n)

Usando Taylor:

$$f(\xi) = f(x_n) + f'(x_n)(\xi - x_n) + \frac{f''(z_n)}{2}(\xi - x_n)^2$$

$$-\frac{f(x_n)}{f'(x_n)} = \xi - x_n + \frac{1}{2f'(x_n)} \cdot f''(z_n)(\xi - x_n)^2 \quad \text{con } -\frac{f(x_n)}{f'(x_n)} = x_{n+1} - x_n \quad \text{per def.}$$

quindi:

$$x_{n+1} - \xi = \frac{f''(z_n)}{2f'(x_n)}(\xi - x_n)^2 \rightarrow |x_{n+1} - \xi| = \frac{|f''(z_n)|}{|2f'(x_n)|} |\xi - x_n|^2 \quad \text{che diventa:}$$

$$e_{n+1} = \frac{1}{2} \cdot \left| \frac{f''(z_n)}{f'(x_n)} \right| \cdot e_n^2 \quad \text{si ha quindi: } \lim \frac{e_{n+1}}{e_n^2} = \frac{1}{2} \left| \frac{f''(\xi)}{f'(\xi)} \right| \quad \text{da cui si deriva}$$

che se $f''(\xi) \neq 0$ e $f'(\xi) \neq 0$ l'ordine di convergenza è esattamente 2. Se

$f''(\xi) = 0$, $f'(\xi) \neq 0$ e $\exists f'''(\xi)$ allora l'ordine è almeno 3.

Ordine di convergenza delle iterazioni di punto fisso

Sia ξ punto fisso di $\phi \in C^p(I)$, $p \geq 1$ con intervallo I di \mathbb{R} , supponiamo di essere in ipotesi che garantiscano la convergenza a ξ di $x_{n+1} = \phi(x_n)$, $n \geq 0$ con $x_n \in I$, allora $\{x_n\}$ ha:

① ordine esattamente $p=1 \iff 0 < |\phi'(\xi)| < 1$

② ordine esattamente $p > 1 \iff \phi^{(j)}(\xi) = 0, 1 \leq j \leq p-1$ e $\phi^{(p)}(\xi) \neq 0$

Dim ①: $e_{n+1} = |\phi'(z_n)| e_n$ con $z_n \in \text{int}(\xi, x_n)$ per il teorema del valor medio

$$\text{risulta: } \lim_{n \rightarrow \infty} \frac{e_{n+1}}{e_n} = |\phi'(\lim z_n)| = |\phi'(\xi)|$$

Dim ②: si utilizza la formula di Taylor di grado $p-1$ centrata in ξ con resto p -esimo in forma di Lagrange:

$$x_{n+1} = \phi(x_n) = \phi(\xi) + \phi'(\xi)(x_n - \xi) + \frac{\phi''(\xi)}{2!}(x_n - \xi)^2 + \dots + \frac{\phi^{(p-1)}(\xi)}{(p-1)!}(x_n - \xi)^{p-1} + \frac{\phi^{(p)}(u_n)}{p!}(x_n - \xi)^p$$

con $u_n \in \text{int}(\xi, x_n)$ e $u_n \xrightarrow{n \rightarrow \infty} \xi$

Per dimostrare la condizione sufficiente " \Leftarrow " si applicano le condizioni

$\phi^{(j)}(\xi) = 0$ per $1 \leq j \leq p-1$ e $\phi^{(p)}(\xi) \neq 0$. Inoltre $\phi(\xi) = \xi$. Risulta, quindi:

$$x_{n+1} - \xi = \frac{\phi^{(p)}(u_n)}{p!}(x_n - \xi)^p \text{ che passando ai moduli diventa:}$$

$$\frac{e_{n+1}}{e_n^p} = \frac{|\phi^{(p)}(u_n)|}{p!} \xrightarrow{n \rightarrow \infty} \frac{|\phi^{(p)}(\xi)|}{p!} \neq 0 \text{ quindi } \{x_n\} \text{ ha ordine } p$$

Per dimostrare la condizione necessaria " \Rightarrow " si suppone per assurdo che

$\exists j < p$ t.c. $\phi^{(j)}(\xi) \neq 0$ e $\{x_n\}$ ha ordine p . Si avrebbe quindi, tramite Taylor:

$$\frac{e_{n+1}}{e_n^k} \rightarrow \frac{|\phi^{(k)}(\xi)|}{k!} = L' \neq 0 \text{ con } k = \min \{j < p : \phi^{(j)}(\xi) \neq 0\}$$

$$\text{ma } \frac{e_{n+1}}{e_n^p} = \frac{e_{n+1}}{e_n^k} \cdot e_n^{k-p} \text{ e so che } \frac{e_{n+1}}{e_n^p} \rightarrow L \neq 0$$

Quindi:

$$\frac{e_{n+1}}{e_n^k} \cdot e_n^{k-p} \rightarrow \infty \text{ perche' } \frac{e_{n+1}}{e_n^k} \rightarrow L' \text{ e } e_n^{k-p} \rightarrow \infty \text{ perche' } e_n \rightarrow 0 \text{ e } k-p < 0$$

$$\text{cioe' } \frac{e_{n+1}}{e_n^p} \rightarrow \infty, \quad n \rightarrow \infty \text{ che contraddice l'ipotesi dell'ordine finito.}$$

Esistenza del polinomio interpolatore di grado $\leq n$ su $n+1$ nodi distinti

ESISTENZA

Dati $n+1$ nodi distinti $\{x_i\}_{0 \leq i \leq n}$ considero per ogni nodo fissato il polinomio elementare di Lagrange:

$$l_i(x) = \frac{N_i(x)}{N_i(x_i)} \quad \text{con} \quad N_i(x) = \prod_{\substack{j=0, j \neq i \\ j \leq n}} (x - x_j) \neq 0 \quad \text{e} \quad N_i(x) \text{ polinomio grado } n$$

Quindi $l_i(x)$ ha grado n e $l_i(x_k) = \delta_{ik} = \begin{cases} 0 & i \neq k \\ 1 & i = k \end{cases}$ (delta di Kronecker)

Possiamo quindi definire il polinomio interpolatore di Lagrange:

$$P_n(x) = \prod_n(x) = \sum_{i=0}^n y_i l_i(x) \quad \text{con} \quad \prod_n(x) \in P_n$$

Verifica dell'interpolazione:

$$\prod_n(x_k) = \sum_{i=0}^n y_i l_i(x_k) = \sum_{i=0}^n y_i \delta_{ik} = y_k \delta_{kk} = y_k, \quad 0 \leq k \leq n \quad (\delta_{ik} = 0)$$

Unicità del grado del polinomio interpolatore di grado $\leq n$ su $n+1$ nodi distinti (esempio con grado $< n$)

UNICITÀ

Supponiamo $\exists p, q \in P_n$ t.c. $p(x_i) = y_i = q(x_i)$ con $0 \leq i \leq n$.

Allora il polinomio $p - q \in P_n$ e si ha:

$$(p - q)(x_i) = p(x_i) - q(x_i) = 0, \quad 0 \leq i \leq n \quad \text{cioè } p - q \text{ avrebbe } n+1 \text{ zeri dist.}$$

Ma $p - q$ può avere al massimo n zeri distinti, a meno che non sia un p. nullo. Quindi $(p - q)(x) = 0 \quad \forall x \Rightarrow p(x) = q(x) \quad \forall x$

Es.:

Un esempio di polinomio di grado $< n$ che interpola $n+1$ punti si ha con $f(x) = x^2$ dove, prendendo $n \geq 3$ punti, si avrà sempre il polinomio interpolatore uguale a $f(x)$ e, quindi, di grado 2 che è quindi $< n$.

Perché l'interpolazione a tratti a passo costante converge uniformemente con errore $O(h^2)$ se $f \in C^2[a, b]$? (si ricavi una stima dell'errore)

Siano $f \in C^{s+1}[a, b]$, $s \geq 0$ e $\{x_i\} \subset [a, b]$ $n+1$ nodi distinti, con n multiplo di s . Allora $\exists K_s > 0 : \text{dist}(f, \Pi_s^c) \leq K_s h^{s+1}$, $h = \max \Delta x_i$.

Dimostrazione per $s=1$ ($f \in C^2[a, b]$)

$$\text{dist}(f, \Pi_1^c) = \max_{x \in [a, b]} |f(x) - \Pi_1^c(x)| = \max_{0 \leq i \leq n-1} \max_{x \in [x_i, x_{i+1}]} |f(x) - \Pi_{1,i}^c(x)| =$$

$$= \max_{1 \leq i \leq n} \max_{x \in [x_{i-1}, x_i]} |f(x) - \Pi_{1,i}^c(x)|$$

Si ricava quindi la stima dell'errore: $\max_{x \in [x_{i-1}, x_i]} \left| \frac{f''(x)}{2} \right| \cdot \frac{h^2}{8} = M_{2,i} \frac{h^2}{8}$ da cui

$$\text{dist}(f, \Pi_1^c) = \max_{1 \leq i \leq n} \max_{x \in [x_{i-1}, x_i]} |f(x) - \Pi_{1,i}^c(x)| \leq \frac{h^2}{8} \cdot \max_{1 \leq i \leq n} M_{2,i} = \frac{M_2}{8} h^2$$

$$\text{con } M = \max_{x \in [a, b]} |f''(x)|$$

Sistema delle equazioni normali per approssimazione polinomiale/ Sistema delle equazioni normali per approssimazione ai minimi quadrati

Dati N punti $\{(x_i, y_i)\}$, $y_i = f(x_i)$, $1 \leq i \leq N$ e $m < N$, il vettore $a \in \mathbb{R}^{m+1}$ minimizza $\phi(a) = \sum_{i=1}^N (y_i - \sum_{j=0}^m a_j \cdot x_i^j)^2 \iff$ risolve il sistema $V^T V a = V^T y$

Dire che $a \in \mathbb{R}^{m+1}$ è un minimo assoluto per $\phi(a)$ equivale a dire:
 $\phi(a+b) \geq \phi(a) \quad \forall b \in \mathbb{R}^{m+1}$ ma $\phi(a+b) = \phi(a) + 2(V^T(Va - y), b) + (Vb, Vb)$

Dim " \Leftarrow ": assumendo che $V^T V a = V^T y$ abbiamo:

$$V^T V a - V^T y = V^T (Va - y) = 0 \quad \text{e} \quad (V^T (Va - y), b) = (0, b) \quad \text{vettore nullo in } \mathbb{R}^{m+1} = 0$$

Dim " \Rightarrow ": assumendo che $\phi(a+b) \geq \phi(a) \quad \forall b \in \mathbb{R}^{m+1}$ allora:

$$\phi(a+b) = \phi(a) + 2(V^T (Va - y), b) + (Vb, Vb) \geq \phi(a) \quad \forall b \quad \text{cioè:}$$

$$2(V^T (Va - y), b) + (Vb, Vb) \geq 0 \quad \forall b.$$

Ponendo $b = \varepsilon v$, con v vettore versore $((v, v) = 1)$. Dividendo per $\varepsilon > 0$ e considerando $\varepsilon \rightarrow 0$ otteniamo:

$$(V^T (Va - y), v) \geq 0 \quad \forall v$$

Ma essendo la disuguaglianza valida per ogni v , allora possiamo sostituire $-v$ a v e otteniamo:

$$(V^T (Va - y), -v) = -(V^T (Va - y), v) \geq 0 \quad \forall v \quad \text{cioè}$$

$$0 \leq (V^T (Va - y), v) \leq 0 \Rightarrow (V^T (Va - y), v) = 0 \quad \forall v$$

Ma essendo il vettore nullo l'unico vettore ortogonale a tutti i vettori:

$$V^T (Va - y) = 0 \iff V^T V a = V^T y \quad \text{ovvero "a" soluzione del sistema}$$

Derivazione numerica rapporto incrementale simmetrico

Sistema delle equazioni normali

DERIVAZIONE SIMMETRICA CON RAPPORTO INCREMENTALE SIMMETRICO

Assumiamo che $f \in C^3(I)$, con I intervallo di derivazione, scriviamo la formula di Taylor "da destra" e "da sinistra" centrate, in x

$$f(x+h) = f(x) + hf'(x) + \frac{h^2}{2} f''(x) - \frac{h^3}{3} f'''(\xi)$$

$$f(x-h) = f(x) - hf'(x) + \frac{h^2}{2} f''(x) - \frac{h^3}{3} f'''(\eta) \text{ dove } \xi \in (x, x+h), \eta \in (x-h, x)$$

da cui si ottiene, sottraendo membro a membro

$$f(x+h) - f(x-h) = 2hf'(x) + O(h^3) \text{ e anche } S(h) = \frac{f(x+h) - f(x-h)}{2h} = f'(x) + O(h^2)$$

$$\text{con } |f'(x) - S(h)| = \frac{1}{2} |f'''(\xi) + f'''(\eta)| h^2 \leq \frac{1}{2} (|f'''(\xi)| + |f'''(\eta)|) h^2 \leq dh^2$$

$$\text{dove } d = \frac{1}{6} \max_{t \in I} |f'''(t)|$$

Questo dimostra che l'errore commesso dall'approssimazione della derivata in x con il rapporto incrementale simmetrico

$$S(h) = \frac{f(x+h) - f(x-h)}{2h} \text{ è } O(h^2) \text{ per } f \in C^3(I)$$

Rapporto incrementale standard

derivazione numerica standard col rapporto incrementale

Supponiamo $f \in C^2(I)$, con I intervallo di derivazione, f derivabile, dato il rapporto incrementale:

$$S_+(h) = \frac{f(x+h) - f(x)}{h}, \quad h > 0$$

si avrà che $S_+(h) = f'(x) + O(h)$ ed esprimendolo la formula di Taylor si ottiene:

$$f(x+h) = f(x) + f'(x)h + \frac{f''(\xi)}{2}h^2, \quad \xi \in \text{int}(x, x+h)$$

$$\text{quindi } S_+(h) = f'(x) + \frac{f''(\xi)}{2}h$$

Se \tilde{S}_+ approssima S_+ , si ha la seguente stima dell'errore:

$$|f'(x) - \tilde{S}_+(h)| \leq |f'(x) - S_+(h)| + |S_+(h) - \tilde{S}_+(h)|$$

da cui

$$\begin{aligned} |S_+(h) - \tilde{S}_+(h)| &= \frac{f(x+h) - f(x) - \tilde{f}(x+h) + \tilde{f}(x)}{h} \\ &\leq \frac{1}{h} (|f(x+h) - \tilde{f}(x+h)| + |\tilde{f}(x) - f(x)|) \\ &\leq \frac{1}{h} (|f(x+h) - \tilde{f}(x+h)| + |\tilde{f}(x) - f(x)|) \\ &\leq \frac{2}{h} \|f - \tilde{f}\|_\infty \end{aligned}$$

ponendo $\varepsilon \geq \|f - \tilde{f}\|_\infty$ si ottiene che

$$|f'(x) - \tilde{S}_+(h)| \leq ch + \frac{2\varepsilon}{h} := E(h), \quad c = \frac{f''(\xi)}{2}$$

se $h \rightarrow 0$, $E(h) \rightarrow \infty$ e poiché $E'(h) = c - \frac{2\varepsilon}{h^2}$ il suo punto minimo è $\sqrt{\frac{2\varepsilon}{c}} = O(\sqrt{\varepsilon})$

Dunque l'errore minimo commesso dall'approssimazione della derivata utilizzando il rapporto incrementale $\frac{f(x+h) - f(x)}{h}$ è $O(\sqrt{\varepsilon})$ per $f \in C^2$

Formule di quadratura composte

Per le formule di quadratura composte ci riconduciamo a due casi, facendo riferimento nel caso di calcolo di nodi equispaziati:

- Per $s = 1$ alla formula dei trapezi, in cui l'integrale viene approssimato alla somma delle aree dei trapezi lineari corrispondenti all'interpolazione lineare a tratti.

$$\begin{aligned} I_n^{trap}(f) &= I(\Pi_1^c) \\ &= \int_a^b \Pi_1^c(x) dx \\ &= \sum (\text{aree trapezi}) \\ &= \underbrace{\frac{h}{2} \cdot (f(x_0) + f(x_1))}_{\text{area trapezio 1}} + \underbrace{\frac{h}{2} \cdot (f(x_1) + f(x_2))}_{\text{area trapezio 2}} + \dots + \\ &\quad + \underbrace{\frac{h}{2} \cdot (f(x_{n-2}) + f(x_{n-1}))}_{\text{area trapezio (n-1)-esimo}} + \underbrace{\frac{h}{2} \cdot (f(x_{n-1}) + f(x_n))}_{\text{area trapezio n-esimo}} \\ &= \frac{h}{2} (f(x_0) + f(x_n)) + \sum_{i=1}^{n-1} h \cdot f(x_i) \end{aligned} \quad I_n(f) = \sum_{i=0}^n w_i \cdot f(x_i), \text{ con } w_i = \begin{cases} \frac{h}{2}, & i = 0, n \\ h, & 1 \leq i \leq n-1 \end{cases}$$

- Per $s = 2$, integrando la funzione quadratica a tratti, si ottiene la formula delle parabole:

$$I_n^{parab}(f) = I(\Pi_2^c) = \sum (\text{aree trapezi parabolici}) = \sum_{i=0}^n w_i \cdot f(x_i), \text{ con } w_i = \begin{cases} h/3, & i = 0, n \text{ pari} \\ 4h/3, & i \text{ dispari} \\ 2h/3, & i \text{ pari } 2 \leq i \leq n-2 \end{cases}$$

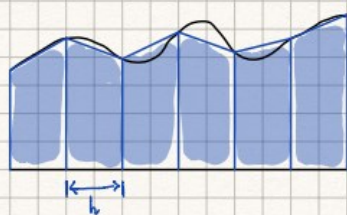
$$\int_{\alpha}^{\beta} \Pi_2(x) dx = \frac{h}{3} \cdot f(\alpha) + \frac{4}{3} \cdot h \cdot f\left(\frac{\alpha + \beta}{2}\right) + \frac{h}{3} \cdot f(\beta)$$

Errore formula trapezi

La formula dei trapezi utilizza l'interpolazione lineare a tratti, imponendo $s=1$ l'integrale viene approssimato con la somma delle aree dei trapezi lineari. L' i -esimo trapezio ha altezza $h = \frac{b-a}{n}$ e basi $f(x_{i-1})$ e $f(x_i)$ con $1 \leq i \leq n$, si avrà quindi l'area $A = \frac{h}{2} (f(x_{i-1}) + f(x_i))$ quindi:

$\frac{h}{2} (f(x_{i-1}) + f(x_i)) + \sum_{i=1}^{n-1} h \cdot f(x_i)$, ottenendo così la formula dei trapezi:

$$I_n(f) = \sum_{i=0}^n w_i f(x_i) \quad \text{con} \quad w_i = \begin{cases} \frac{h}{2}, & i=0, n \\ h, & 1 \leq i \leq n-1 \end{cases}$$



$$I_n^{\text{trap}}(f) = I(\Pi_1^c) = \sum (\text{area trapezi lineari})$$

Per ricavare una stima dell'errore possiamo usare la stima $|I(f) - I_n(f)| = |I(f) - I(\Pi_n^c)| \leq |I(f - \Pi_n^c)| \leq (b-a) \text{dist}(f, \Pi_n^c)$. Se $\text{dist} \rightarrow 0$ allora ci sarà convergenza, altrimenti potrebbero presentarsi problemi di divergenza. Per quanto riguarda le formule di quadratura composte ottenute come $I_n(f) = I(\Pi_n^c)$, con n multiplo di s : $|I(f) - I_n(f)| \leq (b-a) \text{dist}(f, \Pi_n^c) \leq \leq (b-a) K_s \cdot h^{s+1}$ se $f \in C^{s+1}[a, b]$ con $h = \max \Delta x_i$. Quindi per qualsiasi distribuzione dei nodi per cui $h \rightarrow 0$ se $f \in C^{s+1}[a, b]$ le formule sono sempre convergenti con un errore proporzionale a h^{s+1} , ma $s=1$ per i trapezi quindi per $f \in C^2$ sarà convergente con un errore $O(h^2)$.

Condizionamento di sistemi lineari

La prima risposta delle due risponde anche ad un'altra domanda:

Si ricavi una stima dell'errore relativo sulla soluzione di un sistema lineare non singolare in caso di vettore termine noto effetto da errore:

Perturbazione del termine noto

Sia $A \in \mathbb{R}^{n \times n}$ una matrice non singolare, $x \in \mathbb{R}^n$ soluzione del sistema $Ax=b$, $b \neq 0$ e $\tilde{x} = x + \delta x$ soluzione del sistema perturbato $A\tilde{x} = \tilde{b} = b + \delta b$

Fissata una norma vettoriale $\|\cdot\|$ in \mathbb{R}^n vale la seguente stima dell'errore relativo su x : $\frac{\|\delta x\|}{\|x\|} \leq k(A) \frac{\|\delta b\|}{\|b\|}$ dove $k(A) = \|A\| \cdot \|A^{-1}\|$.

Dim.

Dato che $x = A^{-1}b \neq 0$ calcoliamo l'errore relativo in norma.

$$\tilde{x} = x + \delta x = A^{-1}b + A^{-1}\delta b \Rightarrow \|\delta x\| = \|A^{-1}\delta b\| \leq \|A^{-1}\| \cdot \|\delta b\|.$$

Stimiamo quindi da sopra $\frac{1}{\|x\|}$, cioè da sotto $\|x\|$.

Essendo x la soluzione: $\|b\| = \|Ax\| \leq \|A\| \cdot \|x\|$ da cui $\|x\| \geq \frac{\|b\|}{\|A\|}$ e $\frac{1}{\|x\|} \leq \frac{\|A\|}{\|b\|}$

$$\text{perciò: } \frac{\|\delta x\|}{\|x\|} \leq \frac{\|A^{-1}\| \cdot \|\delta b\|}{\|x\|} \leq \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\delta b\|}{\|b\|} = k(A) \cdot \frac{\|\delta b\|}{\|b\|}$$

Perturbazione sulla matrice

Sia $A \in \mathbb{R}^{n \times n}$ una matrice non singolare, $x \in \mathbb{R}^n$ soluzione del sistema $Ax=b$, $b \neq 0$ e $\tilde{x} = x + \delta x$ soluzione del sistema perturbato $\tilde{A}\tilde{x} = b$, $\tilde{A} = A + \delta A$

Fissata una norma vettoriale $\|\cdot\|$ in \mathbb{R}^n vale la seguente stima dell'errore relativo su x : $\frac{\|\delta x\|}{\|\tilde{x}\|} \leq k(A) \cdot \frac{\|\delta A\|}{\|A\|}$ dove $k(A) = \|A\| \cdot \|A^{-1}\|$.

Dim.

Da $\tilde{A}\tilde{x} = (A + \delta A)(x + \delta x) = b$ otteniamo:

$$Ax + A\delta x + \delta A\tilde{x} = b \quad \text{cioè:}$$

$$\delta x = A^{-1}(-\delta A\tilde{x}) = -A^{-1}\delta A\tilde{x} \quad \text{quindi:}$$

$$\|\delta x\| \leq \|A^{-1}\| \cdot \|\delta A\tilde{x}\| \leq \|A^{-1}\| \cdot \|\delta A\| \cdot \|\tilde{x}\| \quad \text{perciò:}$$

$$\frac{\|\delta x\|}{\|\tilde{x}\|} \leq \|A^{-1}\| \cdot \|\delta A\| = \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\delta A\|}{\|A\|} = k(A) \cdot \frac{\|\delta A\|}{\|A\|}$$

Sistemi sovradeterminati

Essi sono sistemi $Ax=b$ su A in $\mathbb{R}^{m \times n}$ e b in \mathbb{R}^m con $m > n$ (più equazioni che incognite):

$$\begin{pmatrix} a_{11} & \dots & a_{1n} \\ a_{21} & \dots & a_{2n} \\ \vdots & & \vdots \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}$$

Essi hanno soluzione se e solo se b sta nel sottospazio di \mathbb{R}^m generato dalle colonne di A :

$$b \in \langle C_1(A), \dots, C_n(A) \rangle \subset \mathbb{R}^m$$

Vogliamo quindi fare in modo che esista una soluzione x tale che $\text{dist}_2(b, Ax) = \|b - Ax\|_2 = 0$

Nuovamente abbiamo un problema di minimo $\phi(x)$ se e solo se $\phi(x+z) \geq \phi(x)$.

Ora

$$\begin{aligned} \phi(x+z) &= \|b - A(x+z)\|_2^2 \\ &= (b - A(x+z), b - A(x+z)) \quad \leftarrow \text{prod. scalare in } \mathbb{R}^m \\ &= (b - Ax - Az, b - Ax - Az) \\ &= (b - Ax, b - Ax) - 2(Az, b - Ax) + (Az, Az) \\ &= \phi(x) + 2(z, A^t(Ax - b)) + \|Az\|_2^2 \end{aligned}$$

- “ \Uparrow ” Se $A^t Ax = A^t b$ allora

$$\phi(x+z) = \phi(x) + \overbrace{\|Az\|_2^2}^{\geq 0} \geq \phi(x) \quad \forall z$$

cioè $\phi(x)$ è minimo

- “ \Downarrow ” Se $\phi(x)$ è un minimo allora $\forall \varepsilon > 0$ e $\forall v \in \mathbb{R}^n, \|v\|_2 = 1$

$$\phi(x + \varepsilon v) = \phi(x) + 2(\varepsilon v, A^t(Ax - b)) + \|\varepsilon v\|_2^2 \geq \phi(x)$$

cioè

$$2\varepsilon(v, A^t(Ax - b)) + \varepsilon^2 \geq 0$$

e dividendo per ε

$$2(v, A^t(Ax - b)) + \varepsilon \geq 0 \quad \forall \varepsilon, v$$

da cui per $\varepsilon \rightarrow 0$

$$(v, A^t(Ax - b)) \geq 0 \quad \forall v$$

Ma allora prendendo $-v$

$$(-v, A^t(Ax - b)) \geq 0 \quad \forall v$$

cioè

$$(v, A^t(Ax - b)) \leq 0 \quad \forall v$$

e quindi

$$(v, A^t(Ax - b)) = 0 \quad \forall v$$

da cui $A^t(Ax - b) = 0$ perché l'unico vettore ortogonale a tutti i versori è il vettore nullo,

cioè $A^t Ax = A^t b$

In conclusione se A ha rango max = n la soluzione ai minimi quadrati del sistema sovradeterminato è unica.

Sistema delle equazioni normali per approssimazione lineare

Sapendo che dati N punti $\{(x_i, y_i)\}$, $y_i = f(x_i)$, $1 \leq i \leq N$ e $m < N$, se il vettore $a \in \mathbb{R}^{m+1}$ minimizza $\phi(a) = \sum_{i=1}^N (y_i - \sum_{j=0}^m a_j x_i^j)^2$ allora risolve il sistema

$V^T V a = V^T y$, si possono usare le proprietà di $V^T V$ per trovare il sistema relativo alla retta dei minimi quadrati. $V^T V$ è una matrice simmetrica e semidefinita positiva. Inoltre $(Vv, Vv) = 0 \iff Vv = 0$ e $(Vv, Vv) = (V^T V v, v)$

quindi $v = 0$ se V ha rango max cioè se ha almeno $m+1$ punti distinti tra i nodi di campionamento. Si ricava quindi una matrice V t.c.

$V = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m+1} & x_{m+1}^2 & \dots & x_{m+1}^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & x_N^2 & \dots & x_N^m \end{pmatrix}$ La sottomatrice $V \in \mathbb{R}^{(m+1) \times (m+1)}$ è matrice di Vandermonde per l'interpolazione di grado $\leq m$ su $m+1$ nodi distinti, quindi è non singolare.

Questo evidenzia che, quindi, il rango della sottomatrice è $m+1$ e che le intere colonne $m+1$ di V sono linearmente indipendenti come vettori di \mathbb{R}^N . Quindi si possono calcolare gli elementi della matrice $V^T V$

e del vettore noto $V^T y$, con $m=1$

$$V^T V = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_N \end{pmatrix} \cdot \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{pmatrix} = \begin{pmatrix} N & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} \in \mathbb{R}^{2 \times 2}$$

$$V^T y = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_N \end{pmatrix} \cdot \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix} \in \mathbb{R}^2$$

quindi il sistema è: $\begin{pmatrix} N & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$

Costo computazionale MEG/Metodo di Gauss (indicato come possibile domanda file 20/21)

Il costo computazione del meg è dato dall'analisi tra ciclo interno, composto da n moltiplicazioni ed n somme, scritto come:

$$c_n^{meg} \sim \sum_{i=1}^{n-1} \sum_{k=i+1}^n 2n$$

$$= 2n \sum_{i=1}^{n-1} (n-i)$$

$$= 2n \sum_{j=1}^{n-1} j$$

$$= 2n \cdot \frac{n(n-1)}{2}$$

$$= n^3 - n^2 \sim n^3, \quad n \rightarrow \infty$$

Vedendo però che le operazioni vettoriali non ha senso farle sui vettori riga, le facciamo solo sul segmento di vettori con indici da $i+1$ ad n , verificando che otteniamo:

$$c_n^{meg} \sim \sum_{i=1}^{n-1} \sum_{k=i+1}^n 2(n-i)$$

$$= 2 \sum_{i=1}^{n-1} (n-i)^2$$

$$= 2 \sum_{j=1}^{n-1} j^2 \sim \frac{2}{3} n^3$$

Ottenendo infine:

$$\frac{(n-1)^3}{3} < \sum_{j=1}^{n-1} j^2 < \frac{n^3}{3} - 1$$

Fattorizzazione QR per soluzioni di sistemi sovradeterminati (indicato come possibile domanda file 20/21)

Sia $A \in \mathbb{R}^{m \times n}$, $m \geq n$, $\text{rango}(A) = n$. Fattorizzando $A = QR$, si ha che

$$A^t A = (QR)^t QR = R^t Q^t QR = R^t I R = R^t R$$

e

$$A^t b = R^t Q^t b$$

quindi il sistema $A^t A x = A^t b$ diventa

$$R^t R x = R^t Q^t b$$

ma essendo R (e quindi R^t) invertibile

$$(R^t)^{-1} R^t R x = R x = (R^t)^{-1} R^t Q^t b = Q^t b$$

cioè il sistema $A^t A x = A^t b$ equivale al sistema triang. sup.

$$R x = d = Q^t b$$

che si può facilmente risolvere con la sostituzione all'indietro.

Nella pratica andiamo a risolvere un sistema perturbato del tipo $\tilde{R}\tilde{x}=\tilde{d}$, con R molto meglio condizionata di $A^t A$.

Applicazione del MEG nel calcolo del determinante (indicato come possibile domanda file 20/21)

Innanzitutto diciamo che il MEG è basato su due proprietà fondamentali in merito a trasformazioni della matrice. Esse sono:

- la sostituzione alla riga k della somma della riga k con la riga i moltiplicata per uno scalare
- lo scambio di due righe porta il determinante a cambiare segno

Esso si applica ripetutamente per mettere zeri in ogni colonna sotto la diagonale principale, in maniera tale da ottenere una matrice triangolare superiore.

Partendo ad esempio da una matrice A 3×3 (a destra) si arriva facilmente alla matrice U (a sinistra):

$$A = \begin{pmatrix} 1 & 2 & 1 \\ 2 & 2 & 3 \\ -1 & -3 & 0 \end{pmatrix} \quad U = \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{pmatrix}$$

Questo afferma che il MEG consiste in una sequenza di $n-1$ trasformazioni:

$$A^{(1)} = A \rightarrow A^{(2)} \rightarrow A^{(3)} \rightarrow \dots \rightarrow A^{(n)} = U$$

ottenendo una matrice U finale che è triangolare superiore con la seguente struttura schematica:

$A^{(i)} = \begin{pmatrix} \text{trapezio di zeri} & a_{ki}^{(i)} \\ 0 & a_{ii}^{(i)} \end{pmatrix}$

lo scopo è mettere zeri al posto di a_{ki} , $i+1 \leq k \leq n$

$a_{ki}^{(i)} = 0$
 $k > j = 1, \dots, i-1$

Se $a_{ii}^{(i)} \neq 0$ (cioè se l'elemento diagonale di $A^{(i)}$ è non nullo) si può propagare a destra la struttura del trapezio di zeri con le trasformazioni

$$\mathcal{R}_k^{(i+1)} := \mathcal{R}_k^{(i)} + \left(-\frac{a_{ki}^{(i)}}{a_{ii}^{(i)}} \right) \cdot \mathcal{R}_i^{(i)}, \quad i+1 \leq k \leq n$$

che mettono lo zero al posto k, i visto che

$$a_{ki}^{(i+1)} = a_{ki}^{(i)} + \left(-\frac{a_{ki}^{(i)}}{a_{ii}^{(i)}} \right) \cdot a_{ii}^{(i)} = 0$$

Meg e relazione con la fattorizzazione LU (indicato come possibile domanda file 20/21)

Per alcune classi di matrici, come quelle a diagonale strettamente dominante e le simmetriche definite positive, ciò già vale. Partendo dall'esempio di matrici 3 x 3, definiamo U (ottenuta alla fine del MEG):

$$U = A^{(3)} = T_{3,2}(-m_{3,2}) T_{3,1}(-m_{3,1}) T_{2,1}(-m_{2,1}) A$$

Posto $\mathcal{L} = T_{3,2}(-m_{3,2}) T_{3,1}(-m_{3,1}) T_{2,1}(-m_{2,1})$ abbiamo

$$U = \mathcal{L}A \Rightarrow A = LU \text{ con } L = \mathcal{L}^{-1}$$

Avendo L che ha ordine invertito:

$$L = \mathcal{L}^{-1} = \overbrace{(T_{2,1}(-m_{2,1}))^{-1} (T_{3,1}(-m_{3,1}))^{-1} (T_{3,2}(-m_{3,2}))^{-1}}^{\text{ordine invertito}} \quad \text{ricordando che } (T_{k,i}(\alpha))^{-1} = T_{k,i}(-\alpha)$$

$$= T_{2,1}(m_{2,1}) T_{3,1}(m_{3,1}) T_{3,2}(m_{3,2})$$

ed L è matrice triangolare inferiore (Lower Triangular):

$$L = \begin{pmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{pmatrix}$$

Otteniamo quindi una fattorizzazione $PA = LU$ con:

- U matrice triangolare superiore
- P matrice di permutazione ottenuta come prodotto di matrici di scambio
- L matrice triangolare inferiore e con i moltiplicare rimescolati nel triangolo inferiore sotto la diagonale.

In sintesi grafica, si ha una fattorizzazione di costo dato dagli scambi e dai moltiplicatori usati, che sarà $2/3 \cdot n^3$:

$$PA = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{pmatrix} \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{pmatrix}$$

Calcolo della matrice inversa (indicato come possibile domanda file 20/21)

Un modo semplice per calcolare l'inversa si basa su una proprietà del prodotto matrice-vettore che abbiamo usato spesso, cioè che $z = Bc$ si può interpretare come combinazione lineare delle colonne di B tramite i coeff. di c, cioè:

$$z = Bc = c_1 \underset{\text{colonna 1 di } B}{C_1(B)} + \dots + c_n \underset{\text{colonna n di } B}{C_n(B)}$$

Applicando l'osservazione con:

cioè $c_i = 1$ e $c_j = 0, j \neq i$, otteniamo

$$Be^{(i)} = C_i(B)$$

(cioè con $c_i = 1, c_j = 0$ e $j \neq i$)

Allora per $B = A^{-1}$

$$C_i(A^{-1}) = A^{-1}e^{(i)} \Leftrightarrow AC_i(A^{-1}) = e^{(i)}$$

cioè $C_i(A^{-1})$ è la soluzione di

$$Ax^{(i)} = e^{(i)}, \quad 1 \leq i \leq n$$

$$c = e^{(i)} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Possiamo quindi calcolare l'inversa "per colonne" risolvendo $M = n$ sistemi lineari, tutti con matrice A, in cui il termine noto varia tra gli n vettori coordinati della base canonica.

Se usassimo n volte il metodo standard avremmo un costo $2/3n^4$.

Invece calcolando una volta col meg la fattorizzazione LU e risolvendo le $M = n$ coppie di

sistemi triangolari come segue, si ha il calcolo dell'inversa a costo cubico:

$$\begin{cases} Ly^{(i)} = Pe^{(i)} \\ Ux^{(i)} = y^{(i)} \end{cases}$$

si ha un costo

$$c_n^{(2)} \sim \frac{2}{3}n^3 + 2n^2n = \frac{8}{3}n^3 \text{ flops}$$