



Corso Luigi Einaudi, 55 - Torino

Appunti universitari

Tesi di laurea

Cartoleria e cancelleria

Stampa file e fotocopie

Print on demand

Rilegature

NUMERO: 1190

DATA: 22/10/2014

A P P U N T I

STUDENTE: Parisi

MATERIA: Calcolo Numerico

Prof. Scuderi

Il presente lavoro nasce dall'impegno dell'autore ed è distribuito in accordo con il Centro Appunti.

Tutti i diritti sono riservati. È vietata qualsiasi riproduzione, copia totale o parziale, dei contenuti inseriti nel presente volume, ivi inclusa la memorizzazione, rielaborazione, diffusione o distribuzione dei contenuti stessi mediante qualunque supporto magnetico o cartaceo, piattaforma tecnologica o rete telematica, senza previa autorizzazione scritta dell'autore.

**ATTENZIONE: QUESTI APPUNTI SONO FATTI DA STUDENTIE NON SONO STATI VISIONATI DAL DOCENTE.
IL NOME DEL PROFESSORE, SERVE SOLO PER IDENTIFICARE IL CORSO.**

I N D I C E

- Lezione 1: Aritmetica dei calcolatori ed errori di calcolo
 - Rappresentazione floating-point
 - Underflow e Overflow
 - Errore assoluto, errore relativo, errore di arrotondamento
 - Precisione di macchina & errore di Round-off
 - Operazioni di macchina
- Lezione 2: Aritmetica dei calcolatori ed errori di calcolo
 - Cancellazione numerica
 - Problema numerico
 - Numero di condizionamento
 - Stabilità di un algoritmo
 - Richiami su vettori e matrici
 - Risoluzione numerica di sistemi lineari
- Lezione 3: Metodi numerici per la risoluzione di sistemi lineari
 - Sostituzione all'indietro per sistemi triangolare superiore
 - Sostituzione in avanti per sistemi triangolari inferiori
 - Metodo dell'eliminazione di Gauss.
- Lezione 4: Pivoting parziale e fattorizzazione delle matrici.
 - Pivoting parziale
 - Fattorizzazione matrici $PA = LU$
- Lezione 5: Applicazioni fattorizzazione $PA = LU$ e metodi iterativi
 - Applicazioni fattorizzazione $PA = LU$
 - Fattorizzazione Choleski
 - Metodo di Jacobi
 - Metodo di Gauss-Seidel
- Lezione 6: Metodi di Jacobi e Gauss-Seidel
 - Criteri di arresto iterazioni
 - Comandi matlab
 - Convergenza dei metodi iterativi
 - Convergenza dei metodi di Jacobi e Gauss-Seidel
 - Proprietà metodi iterativi.

Lezione n. 1

Aritmetica calcolatori ed errori calcolo.

Consideriamo il sistema di numerazione decimale ω in base 10

Un numero reale " ω " si può rappresentare con le cifre che vanno da 0 a 9 ed ogni numero reale si può esprimere come combinazione di potenze in base 10

$$\bullet \omega = 109,34 = 1 \times 10^2 + 0 \times 10^1 + 9 \times 10^0 + 3 \times 10^{-1} + 4 \times 10^{-2}$$

RAPPRESENTAZIONE FLOATING-POINT:

$$\omega = pN^q \quad p = \text{numero reale} \quad q = \text{intero} \quad N = \text{base (per noi sono' sempre 10)}$$

$$\bullet \omega = 0,015 \cdot 10^{-2} = 0,15 \cdot 10^{-2} = 0,0015 \cdot 10^0$$

Spostando la virgola a destra e a sinistra bisogna modificare anche opportunamente la potenza in base 10 = N

Quindi per un numero reale " ω " esistono più rappresentazioni floating-point, per cui la rappresentazione non è unica.

" ω " è univocamente determinato se: $N^{-1} \leq |p| < 1 \Rightarrow$ in $N=10 \quad 0,1 \leq |p| < 1$

In questo caso la prima cifra di p dopo il punto decimale deve essere diversa da zero.

Se questa condizione risulta soddisfatta, la rappresentazione floating point si dice normalizzata.

p e q sono quindi unici e li possiamo nominare come: mantissa e esponente o caratteristica rispettivamente.

$$\bullet \omega = 0,15 \cdot 10^{-2}$$

$$p = 0,15 \quad q = -2$$

$$\omega = 12,4 \cdot 10^0 = 0,124 \cdot 10^2$$

$$p = 0,124 \quad q = 2$$

$$\omega = 0,0013 \cdot 10^4 = 0,13 \cdot 10^2$$

$$p = 0,13 \quad q = 2$$

Ricapitolando: se conosciamo N , la coppia (p, q) individua univocamente il numero reale $\omega = pN^q$ nella rappresentazione normalizzata.

Per memorizzare " ω " è sufficiente memorizzare $p^{\pm q}$ oppure il segno di ω , $|p|$ e q .

La memoria di un calcolatore è però limitata e quindi:

- $|p|$ può avere al massimo E cifre

- $m \leq q \leq M$ con $m < 0 =$ minimo esponente rappresentabile,

$M > 0 =$ massimo esponente rappresentabile.

UNDERFLOW e OVERFLOW:

- Sia $a = p \cdot N^q$ con $N^{-1} \leq |p| < 1$; ← numeri normalizzati
- Se $q < m$, il numero non è rappresentabile e si ha un problema di underflow.
Il numero reale viene approssimato allo zero e il processo di calcolo prosegue.
 - Se $q > M$, il numero non è rappresentabile e si ha un problema di overflow.
In questo caso il calcolatore emette il processo di calcolo e segnala all'utente un errore.

Ci sono dei casi in cui possiamo risolvere il problema dell'overflow.

Supponiamo di avere: $N = 10$ t $-127 \leq q \leq 128$

$$y = \sqrt{x_1 \cdot x_2} \quad \text{con } x_1 = 10^{100} \quad x_2 = 10^{50}$$

$$x_1 \cdot x_2 = 10^{150} \quad q > 128 \quad \rightarrow \text{overflow}$$

Per aggirare il problema di overflow si può:

$$y = \sqrt{x_1 \cdot x_2} = \sqrt{x_1} \cdot \sqrt{x_2}$$

Se invece: $y = \sqrt{x_1^2 + x_2^2}$ con $x_1 = 10^{100}$ $x_2 = 10^{50}$

mantenendo l'espressione così, si verificherebbe overflow. Se invece:

$$y = \sqrt{x_1^2 \left(1 + \left(\frac{x_2}{x_1}\right)^2\right)} = x_1 \sqrt{1 + \left(\frac{x_2}{x_1}\right)^2}$$

non si ottiene overflow. Si potrebbe però ottenere un fenomeno di underflow se x_2 è un numero piccolo. Se si ottiene underflow il numero verrà approssimato a zero.

SE |p| HA PIU' DI t CIFRE: normalizzate

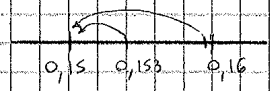
Sia $a = p \cdot N^q$ con $N^{-1} \leq |p| < 1$ e $m \leq q \leq M$

Per rappresentare questo numero in un sistema floating-point con t cifre, per la mantissa si ricorre ad una delle due tecniche di arrotondamento:

1- tecnica di troncamento: si esclude la parte a destra della t-esima cifra.

• $N = 10$ t = 2 $a = 0,153 \cdot 10^0$

$\bar{a} = 0,15 \cdot 10^0$ → troncamento.



2- tecnica di arrotondamento: si aggiunge $\frac{1}{2} N^{-t}$ a p e poi si tronca il risultato alla t-esima cifra.

• $N = 10$ t = 2 $a = 0,153$ $\frac{1}{2} N^{-t} = \frac{1}{2} 10^{-2} = 0,005$

$$\begin{array}{r} 0,153 \\ + 0,005 \\ \hline 0,158 \end{array}$$
 $\bar{a} = 0,15$ → troncamento.

Se utilizzo la tecnica di troncamento: $\bar{p} = \bar{p}_1 \Rightarrow |p - \bar{p}| < N^{-t}$

Se utilizzo la tecnica di arrotondamento: $\bar{p} = \bar{p}_1 \Rightarrow |p - \bar{p}| \leq \frac{1}{2} N^{-t}$

$\bar{p} = \bar{p}_2 \Rightarrow |p - \bar{p}| \leq \frac{1}{2} N^{-t}$

L'errore assoluto non supera N^{-t} e $\frac{1}{2} N^{-t}$ rispettivamente.

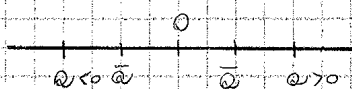
Le tecniche di arrotondamento soddisfano le seguenti disuguaglianze:

$$|a - \bar{a}| = \begin{cases} < N^{q-t} & \text{per troncamento} \\ \leq \frac{1}{2} N^{q-t} & \text{per arrotondamento} \end{cases}$$

$$\frac{|a - \bar{a}|}{|a|} \leq \frac{|a - \bar{a}|}{N^{q-t}} \begin{cases} < N^{-t} & \text{per troncamento} \\ \leq \frac{1}{2} N^{-t} & \text{per arrotondamento} \end{cases}$$

L'errore della tecnica di arrotondamento è la metà di quella di troncamento, quindi la prima è preferibile alla seconda in quanto la seconda può generare un errore maggiore (che può essere il doppio).

Inoltre il segno di $\frac{a - \bar{a}}{a}$ è costante (positivo) ossia:



PRECISIONE DI MACCHINA & ERRORE DI ROUND-OFF:

È la quantità:

$$\text{eps} = \begin{cases} N^{1-t} & \text{per troncamento (1)} \\ \frac{1}{2} N^{1-t} & \text{per arrotondamento (2)} \end{cases}$$

È una costante caratteristica di ogni aritmetica floating-point, dipende dalle cifre rappresentabili e dalla tecnica di arrotondamento (1) & (2).

Questa costante rappresenta la massima precisione di calcolo su un calcolatore.

Se " \bar{a} " denota il numero di macchina corrispondente al numero reale " a ", posto

$\epsilon = \frac{\bar{a} - a}{a}$ possiamo scrivere:

$$\bar{a} = a(1 + \epsilon) \quad \text{con } |\epsilon| \leq \text{eps}$$

$$\epsilon = \frac{\bar{a} - a}{a} \quad \rightarrow \quad a\epsilon = \bar{a} - a \quad \rightarrow \quad \bar{a} = a(1 + \epsilon)$$

$$|\epsilon| = \frac{|\bar{a} - a|}{|a|} \leq \text{eps}$$

Lezione 2

Aritmetica calcolatori ed errori calcolo

CANCELLAZIONE NUMERICA:

Consiste in una perdita di cifre della mantissa e si verifica quando si esegue un'operazione di sottrazione tra due numeri "quasi uguali" ossia con un certo numero di cifre in comune ed arrotondati a numeri di macchina.

- Per $N=10$ $t=5$ e tecnica di arrotondamento.

$$a_1 = 0,157826831$$

$$a_2 = 0,157368212$$

$$\text{In aritmetica esatte: } a_1 - a_2 = 0,000476619 = 0,476619 \cdot 10^{-3}$$

$$\text{In aritmetica finite: } \bar{a}_1 = 0,15782 \quad \bar{a}_2 = 0,15735$$

$$\bar{a}_1 \ominus \bar{a}_2 = 0,00047 = 0,47000 \cdot 10^{-3}$$

Se i due operandi iniziali non fossero stati affetti da errore, il risultato sarebbe stato accurato. Gli errori iniziali amplificano l'errore dell'operazione finale.

$$\text{Consideriamo: } a_1 = p_1 N^q \quad \text{e} \quad a_2 = p_2 N^q$$

$$\text{Supponiamo che: } \bar{a}_1 = \bar{p}_1 N^q \quad \text{e} \quad \bar{a}_2 = \bar{p}_2 N^q \quad \text{numeri di macchina.}$$

I numeri reali sono arrotondati e quindi $p_1 \neq \bar{p}_1$ e/o $p_2 \neq \bar{p}_2$.

Supponiamo che \bar{p}_1 e \bar{p}_2 abbiano t cifre e che le prime s cifre coincidano tra loro, allora:

$$\bar{a}_1 \ominus \bar{a}_2 = (\bar{p}_1 \ominus \bar{p}_2) N^q = (\bar{p}_1 - \bar{p}_2) N^q = \bar{p} N^{q+s}$$

ove solamente le prime $t-s$ cifre in \bar{p} provengono dalle mantisse p_1 e p_2 ; le restanti s (poste uguali a zero) non hanno alcun significato.

E' quindi poco raccomandabile eseguire sottrazioni su numeri quasi uguali poiché il risultato sarà poco preciso, se i numeri sono stati arrotondati.

Talvolta manipolando le espressioni matematiche che definiscono un problema è possibile evitare il fenomeno della cancellazione numerica.

Quando questo non è possibile si dice che la cancellazione è insita nel problema.

PROBLEMA NUMERICO:

Si intende una descrizione chiara e non ambigua di una commessione funzionale f :

$$f: \begin{array}{ccc} & x & \\ & \text{input} & \\ & \xrightarrow{f} & \\ & y & \\ & \text{output} & \end{array}$$

La commessione fra x e y può essere esplicita ($y = f(x)$) oppure implicita ($f(x, y) = 0$).
Quando immetto i dati di input (x) nel calcolatore essi vengono seguiti da un errore e quindi sono rappresentati con \bar{x}

$$x \xrightarrow{\infty} f(x)$$

$$\bar{x} \xrightarrow{\infty} f(\bar{x})$$

∞ = precisione infinita.

Quindi solo i dati di ingresso vengono perturbati poiché i calcoli vengono eseguiti in entrambi i casi con precisione infinita senza commettere quindi errori e senza approssimare o arrotondare i calcoli.

Se accade che: $x \approx \bar{x}$ e $f(x) \approx f(\bar{x})$ allora il problema è ben condizionato. Ormai a piccole perturbazioni relative o assolute nei dati, corrispondono piccole perturbazioni nel risultato.

Se accade che: $x \approx \bar{x}$ e $f(x) \neq f(\bar{x})$, a piccole variazioni nei dati corrispondono grandi variazioni nel risultato, allora il problema si dice mal condizionato.

Quindi se: $\frac{|f(x) - f(\bar{x})|}{|f(x)|} \approx \frac{|x - \bar{x}|}{|x|}$, $x, f(x) \neq 0$, ossia gli errori associati a \bar{x} e $f(\bar{x})$ hanno lo stesso ordine di grandezza, allora il problema si dice ben condizionato. Altrimenti mal condizionato.

NUMERO DI CONDIZIONAMENTO:

Per studiare il condizionamento di un problema occorre determinare relazioni del tipo:

$$- \frac{|f(x) - f(\bar{x})|}{|f(x)|} \approx K(f, x) \frac{|x - \bar{x}|}{|x|}$$

$$- \frac{|f(x) - f(\bar{x})|}{|f(\bar{x})|} \leq K(f, x) \frac{|x - \bar{x}|}{|x|}$$

K = numero di condizionamento.

Se $K \approx 1$ allora il problema è ben condizionato e viceversa.

RICHIAMI SU VETTORI E MATRICI:

Sia $x = (x_1, \dots, x_n)^T$ (vettore colonna trasposto per questioni di spazio). Le norme di vettore più usate sono:

- $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2} = \sqrt{x^T \cdot x}$ → NORMA EUCLIDEA $\| \cdot \| = \text{norma}$

- $\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$ ^{valore assoluto} → NORMA INFINITA

- $\|x\|_1 = \sum_{i=1}^n |x_i|$ ^{valore assoluto} → NORMA UNO

Le norme di matrice più usate sono:

- $\|A\|_2 = \sqrt{\rho(A^T \cdot A)}$ → NORMA SPETTRALE

con $\rho(B) = \max_i |\lambda_i|$, con λ_i autovalore di B, e detto raggio spettrale di B

- $\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$ ^{valore assoluto}

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix}$$

Si sommano gli elementi in valore assoluto di ciascuna riga. Dopodiché il massimo tra i risultati ottenuti rappresenta la norma infinito ($\|A\|_\infty$).

- $\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|$ ^{valore assoluto}

Si sommano gli elementi di ciascuna colonna in valore assoluto. Dopodiché il massimo tra i risultati ottenuti rappresenta la norma 1 ($\|A\|_1$).

Per le norme 1, 2 e ∞ valgono le seguenti proprietà:

- $\|AB\| \leq \|A\| \|B\|$

- $\|I\| = 1$ con I = matrice identità

- $\|Ax\| \leq \|A\| \|x\|$

Una matrice si dice triangolare superiore quando gli elementi sotto la diagonale principale sono uguali a zero.

Si dice triangolare inferiore quando gli elementi sopra la diagonale principale sono uguali a zero.

Una matrice si dice diagonale quando gli elementi sopra e sotto la diagonale principale sono uguali a zero.

RISOLUZIONE NUMERICA DI SISTEMI LINEARI:

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1m}x_m = b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2m}x_m = b_2 \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mm}x_m = b_m \end{cases}$$

in forma matriciale:

$$\begin{pmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mm} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix}$$

in forma compatta: $Ax = b$

Assumeremo che questo sistema lineare abbia una e una sola soluzione.

Cosa succede se perturbiamo ^{di poco} i dati? andiamo ad esaminare il condizionamento del sistema:

Consideriamo i dati A e b e i dati perturbati $\bar{A} = A + \delta A$ e $\bar{b} = b + \delta b$.
 $\bar{x} = x + \delta x$ è la soluzione in aritmetica esatta del sistema perturbato $\bar{A}\bar{x} = \bar{b}$.

Per esaminare il condizionamento confrontiamo:

- l'errore relativo associato alla soluzione $\rightarrow \frac{\|\delta x\|}{\|x\|}$ con
- l'errore relativo associato ai dati $\rightarrow \frac{\|\delta A\|}{\|A\|} \quad \frac{\|\delta b\|}{\|b\|}$

Se $A + \delta A$ è non singolare, quindi il sistema assume una sola soluzione, e $\|\delta A\|$ (perturbazione in A) risulta $< 1/\|A^{-1}\|$ (quindi piccola), si dimostra che:

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{K(A)}{1 + K(A)} \frac{\|\delta A\|}{\|A\|} \left(\frac{\|A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right) \quad \text{dove} \quad K(A) = \|A\| \|A^{-1}\| \geq \|A \cdot A^{-1}\| = \|\mathbb{I}\| = 1$$

Suppongo inoltre: $\|\delta A\| < \frac{1}{2\|A^{-1}\|}$ si ha che:

$$\frac{\|\delta x\|}{\|x\|} \leq 2K(A) \left(\frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right)$$

In questo caso se $2K(A)$ è piccolo, e piccole perturbazioni sui dati, corrispondono a una piccola perturbazione sul risultato, ed il problema è ben condizionato.

$K(A)$ viene pertanto definito come numero di condizionamento.

Per le norme 1, 2 e ∞ otteniamo:

$$K(A) = \|A\| \|A^{-1}\| \geq \|AA^{-1}\| = \|\mathbb{I}\| = 1$$

Lezione 3

Metodi numerici per risoluzione sistemi lineari.

I metodi numerici per risolvere sistemi lineari sono essenzialmente di due tipi:

- diretti
- iterativi

METODI DIRETTI:

① Tecniche di sostituzione all'indietro per sistemi triangolari superiori.

Dato il sistema triangolare superiore:

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1m}x_m = b_1 \\ \vdots \\ a_{m-1,m-1}x_{m-1} + a_{m-1,m}x_m = b_{m-1} \\ a_{mm}x_m = b_m \end{cases}$$

com'è noto è invertibile
 con $a_{ii} \neq 0$ $\forall i$ questo
 fa sì che $\det \neq 0 \rightarrow$ 1 soluzione

Risoluzione:

$$\begin{cases} x_m = \frac{b_m}{a_{mm}} \\ x_{m-1} = \frac{b_{m-1} - a_{m-1,m}x_m}{a_{m-1,m-1}} \\ \vdots \\ x_1 = \frac{b_1 - \sum_{j=2}^m a_{1j}x_j}{a_{11}} \end{cases}$$

sostituisco x_m nella penultima eq. e trovo x_{m-1} e così via.

$$\begin{cases} x_m = \frac{b_m}{a_{mm}} \\ x_i = \frac{b_i - \sum_{j=i+1}^m a_{ij}x_j}{a_{ii}} \end{cases}$$

con $i = m-1, \dots, 1$

Costo computazionale:

$$\frac{m^2}{2}$$

(in termini di operazioni aritmetiche)

Comandi matlab:

$$A = [\dots, \dots, \dots];$$

$$b = [\dots, \dots];$$

$$m = \text{length}(b);$$

$$x = \text{zeros}(m, 1);$$

$$x(m) = b(m) / A(m,m);$$

si sostituisce perché
 in questo modo si fanno
 delle funzioni di matlab
 e quindi è più veloce.

for $i = m-1 : -1 : 1$

è da $m-1$ a 1 con passo -1

Si può sostituire con:

```

s = 0;
for j = i+1:m
    s = s + A(i,j) * x(j);
end

```

per $i+1$ a m con passo +1

$$s = A(i, i+1:m) * x(i+1:m);$$

$$x(i) = (b(i) - s) / A(i,i);$$

end

$$\sum_{j=i+1}^m a_{ij}x_j = a_{i,i+1}x_{i+1} + \dots + a_{i,m}x_m = \underbrace{(a_{i,i+1}, \dots, a_{i,m})}_{A(i, i+1:m)} \underbrace{(x_{i+1}, \dots, x_m)}_{x(i+1:m)}^T$$

③ Metodo dell'eliminazione di Gauss: *Guarante (da risolvere) il sistema Gauss*

ci consente di risolvere sistemi la cui matrice dei coefficienti è una matrice generica ossia non diagonale o triangolare ecc.

Dato il sistema lineare: $Ax = b$ di ordine n con 1 soluzione

Risoluzione: $Ax = b \xrightarrow{n-1} Ux = \bar{b}$
 triangolare superiore

Quindi questo metodo trasforma inizialmente il sistema assegnato in uno equivalente, ossia con la stessa soluzione, ma il sistema equivalente è associato ad una matrice triangolare superiore U . Una volta ricondotti al sistema equivalente possiamo utilizzare la tecnica ④ e ricavare l'incognita x .
 ↳ sostituzione all'indietro

- Trasformazioni:
- la soluzione rimane invariata se si scambiano tra loro due eq. del sistema.
 - la soluzione rimane invariata se si sostituisce ad un eq. del sistema una combinazione lineare dell'eq. stessa con un'altra.

Esempio: $n=4$

$$Ax = b \iff \begin{cases} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + a_{14}x_4 = b_1 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 + a_{24}x_4 = b_2 \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 + a_{34}x_4 = b_3 \\ a_{41}x_1 + a_{42}x_2 + a_{43}x_3 + a_{44}x_4 = b_4 \end{cases}$$

- si elimina l'incognita x_1 dalla 2°, 3°, 4° eq.
- si elimina l'incognita x_2 dalla 3°, 4° eq.
- si elimina l'incognita x_3 dalla 4° eq.

3 passi = $n-1 = 4-1 = 3$

Passo $K=1$: si pone $a_{ij}^{(1)} := a_{ij}$ e $b_i^{(1)} := b_i$. Supponiamo $a_{11}^{(1)} \neq 0$ altrimenti si scambiano la prima eq. con l'eq. K -esima tale che $a_{K1}^{(1)} \neq 0$.

$$a_{2j}^{(2)} = a_{2j}^{(1)} + \left(-\frac{a_{21}^{(1)}}{a_{11}^{(1)}}\right) a_{1j}^{(1)} \quad b_2^{(2)} = b_2^{(1)} + \left(-\frac{a_{21}^{(1)}}{a_{11}^{(1)}}\right) b_1^{(1)}$$

Allo stesso modo elimino x_1 dalla 3° e 4° equazione

Proseguo allo stesso modo con il passo $K=2$ e $K=3$ ottenendo alla fine di $K=3$ la matrice semplificata U .

Si risolve il sistema $Ux = \bar{b}$ con la tecnica di sostituzione all'indietro
 ↳ Upper (matrice triangolare superiore)

④

Comandi Matlab:

$$A = [\dots; \dots; \dots];$$

$$b = [\dots; \dots; \dots];$$

$$m = \text{length}(b);$$

for $k = 1 : m-1$

Si può sostituire con

$$A(k+1:m) = -A(k+1:m, k) / A(k, k)$$

for $i = k+1 : m$

$$A(i, k) = -A(i, k) / A(k, k);$$

Si può sostituire con

$$A(i, k+1:m) = A(i, k+1:m) + A(i, k) \cdot$$

for $j = k+1 : m$

$$A(i, j) = A(i, j) + A(i, k) * A(k, j);$$

$$\cdot A(k, i+1:m)$$

$$b(i) = b(i) + A(i, k) * b(k);$$

end

- Consideriamo il sistema di ordine $n=18$:

$$Ax = b \quad a_{ij} = \cos((j-1)Q_i) \quad Q_i = \frac{2i-1}{2n} \frac{\pi}{n}$$

$$b_i = \sum_{j=1}^m a_{ij}$$

la cui soluzione esatta è: $x = (1, \dots, 1)^T$

Il condizionamento in norma infinito è pari a: $K_\infty(A) \approx 17$
(sistema lineare ben condizionato).

$$\frac{\|x - \bar{x}\|_\infty}{\|x\|_\infty} \approx 10^{-3} \quad \leftarrow \text{errore relativo}$$

con \bar{x} calcolata mediante il metodo di eliminazione di Gauss senza Pivoting.

$$\frac{\|x - \bar{x}\|_\infty}{\|x\|_\infty} \approx 10^{-15} \quad \leftarrow \text{errore relativo}$$

con \bar{x} calcolata mediante il metodo di eliminazione di Gauss con Pivoting parziale.

Questo esempio mostra come l'errore relativo sia molto più grande se non si utilizza il pivoting parziale.

FATTORIZZAZIONE MATRICIA:

- $P_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$ scambio $2 \leftrightarrow 3$ righe

→ $P_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$ moltiplico per una matrice A

→ $P_2 A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \cdot \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{31} & a_{32} & a_{33} \\ a_{21} & a_{22} & a_{23} \end{pmatrix}$

Quindi se voglio realizzare uno scambio tra le righe o le equazioni k e r del sistema lineare basta che io moltiplichi a sinistra per un'opportuna matrice ottenuta scambiando le righe k e r della matrice identità.

$$P_k A^{(k)} x = P_k b^{(k)}$$

dove

$$k \leftrightarrow r$$

$$P_k = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \quad \begin{matrix} k \\ \updownarrow \\ r \end{matrix}$$

Al passo k la sostituzione delle eq. $i = k+1, \dots, m$ con le eq. che si ottengono rispettivamente sommando alle eq. stesse la k -esima eq. moltiplicata per m_{ik} , si può realizzare nel seguente modo:

$$M_k P_k A^{(k)} x = M_k P_k b^{(k)}$$

Lezione 5

Applicazioni fattorizzazione $PA=LU$, Metodi iterativi

- Calcolo del determinante di A: calcolo prima $PA=LU$ e poi

$$\det(A) = (-1)^N \prod_{i=1}^m u_{ii}$$

$$U = \begin{pmatrix} u_{11} & & & \\ 0 & \ddots & & \\ 0 & & u_{ii} & \\ 0 & & & \ddots & \\ & & & & & u_{mm} \end{pmatrix}$$

$N =$ numero totale degli scambi di eq effettuati

- Calcolo dell'inversa di A: $A^{-1} = ?$

$$PA=LU \rightarrow (PA)^{-1} = (LU)^{-1} \rightarrow A^{-1}P^{-1} = U^{-1}L^{-1} \rightarrow$$

$$A^{-1}P^{-1}P = U^{-1}L^{-1}P \rightarrow A^{-1} = U^{-1}L^{-1}P$$

Costo computazionale: m^3 (minimo costo per l'inversione)

- Risoluzione di p sistemi lineari aventi tutti la stessa matrice dei coefficienti:

$$Ax_1 = b_1 \quad Ax_2 = b_2 \quad Ax_3 = b_3 \quad p=3$$

risolvendo tramite eliminazione di Gauss: (con comandi Matlab)

$$x_1 = A \setminus b_1 \quad \text{con costo computazionale } m^3/3$$

$$x_2 = A \setminus b_2 \quad \text{con costo computazionale } m^3/3$$

$$x_3 = A \setminus b_3 \quad \text{con costo computazionale } m^3/3$$

Costo computazionale totale: m^3

Possiamo però cercare di risolvere, con minore costo computazionale, questo problema:

determiniamo $PA=LU$

$$Ax_i = b_i \rightarrow \underbrace{PA}_{LU} x_i = P b_i \rightarrow \underbrace{LU}_{y_i} x_i = P b_i \rightarrow$$

$$\begin{cases} L y_i = P b_i \\ U x_i = y_i \end{cases} \text{ per } i=1,2,3$$

$$\begin{aligned} \text{Costo computazionale: } & \frac{m^3}{3} \text{ della fattorizzazione } PA=LU & + \\ & m^2 \text{ per risoluzione di ciascun sistema} & = \\ & \frac{m^3}{3} + 3m^2 \end{aligned}$$

$3m^2$ è trascurabile per m sufficientemente grande e quindi il costo computazionale è: $\frac{m^3}{3}$

A SIMMETRICA DEFINITA POSITIVA: FATTORIZZAZIONE CHOLESKI

In questo caso il pivoting è superfluo poiché il metodo ③ è già stabile. Il pivoting non migliora ulteriormente la accuratezza della soluzione.

Quindi in questo caso sussiste questa decomposizione:

$$A = LU$$

È possibile riscrivere: $LU = L_1 L_1^T$

con L_1 matrice triangolare inferiore con elementi positivi sulla diagonale principale. Questa fattorizzazione è detta fattorizzazione di Cholesky.

$A = L_1 L_1^T$ si ottiene con costo computazionale $\frac{n^3}{6}$ ovvero la metà del costo computazionale del metodo di Gauss.

Comandi Matlab:

$R = \text{chol}(A)$ calcola il fattore $R = L_1^T$ triangolare superiore della fattorizzazione di Cholesky $A = R^T R = L_1 L_1^T$, della matrice simmetrica definita positiva A .

- Risoluzione del sistema lineare $Ax = b$:

A definita simmetrica positiva

$$R = \text{chol}(A)$$

R = triangolare sup

R^T = triang. inf.

$$A = R^T \cdot R$$

$$R^T \cdot \underbrace{R \cdot x}_y = b$$

$$\rightarrow \begin{cases} R^T y = b & \Rightarrow y \\ R x = y & \Rightarrow x \end{cases}$$

Costo computazionale:

- se R è nota: n^2

- altrimenti: $O\left(\frac{n^3}{6} + n^2\right) = O\left(\frac{n^3}{6}\right)$

- Calcolo inversa di A :

A simmetrica definita positiva.

$$A = R^T \cdot R$$

$$\rightarrow (A)^{-1} = (R^T R)^{-1} \rightarrow A^{-1} = R^{-1} \cdot (R^T)^{-1} = R^{-1} (R^{-1})^T$$

Nel primo caso avrei dovuto calcolare l'inversa di due matrici, mentre con solo di una, a guadagno di un minor costo computazionale.

Costo computazionale: $\frac{2n^3}{3}$

Quindi: $x^{(0)} = \begin{pmatrix} x_1^{(0)} \\ \vdots \\ x_m^{(0)} \end{pmatrix} \xrightarrow{\text{G}} x^{(1)} = \begin{pmatrix} x_1^{(1)} \\ \vdots \\ x_m^{(1)} \end{pmatrix} \xrightarrow{\text{G}} x^{(2)} = \begin{pmatrix} x_1^{(2)} \\ \vdots \\ x_m^{(2)} \end{pmatrix} \text{ ecc...}$

② Metodo di Gauss-Seidel:

$$x_i^{(k+1)} = \frac{b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^m a_{ij} x_j^{(k)}}{a_{ii}} \quad \text{con } k=0, 1, \dots$$

In questo modo si ottiene una soluzione del problema più accurata.

• $\begin{cases} 3x_1 - x_2 = 2 \\ x_1 + 2x_2 = 3 \end{cases} \quad \begin{cases} x_1 = \frac{2+x_2}{3} \\ x_2 = \frac{3-x_1}{2} \end{cases} \quad x = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

- Metodo di Jacobi:

$$\begin{cases} x_1^{(k+1)} = \frac{2+x_2^{(k)}}{3} \\ x_2^{(k+1)} = \frac{3-x_1^{(k)}}{2} \end{cases}$$

	k=0	k=1	k=2	...
$x_1^{(k)}$	$x_1^{(0)}$	$x_1^{(1)}$	$x_1^{(2)}$	$x_1^{(3)}$
0	$2/3$	$7/6$	$19/8$...
0	$3/2$	$7/6$	$1/2$...

$$\frac{\|x - x^{(3)}\|_{\infty}}{\|x\|_{\infty}} \approx 0,8 \cdot 10^{-1}$$

Per avere un errore con una tolleranza specifica allora vado avanti con l'iterazione finché non la raggiungo

- Metodo di Gauss-Seidel:

$$\begin{cases} x_1^{(k+1)} = \frac{2+x_2^{(k)}}{3} \\ x_2^{(k+1)} = \frac{3-x_1^{(k+1)}}{2} \end{cases}$$

	k=0	k=1	k=2	...
$x_1^{(k)}$	$x_1^{(0)}$	$x_1^{(1)}$	$x_1^{(2)}$	$x_1^{(3)}$
0	$2/3$	$19/8$	$107/108$...
0	$7/6$	$35/16$	$27/216$...

$$\frac{\|x - x^{(3)}\|_{\infty}}{\|x\|_{\infty}} \approx 0,8 \cdot 10^{-2}$$

l'errore con questo secondo metodo è minore.

CRITERI DI ARRESTO ITERAZIONI:

Quando si implementa un procedimento iterativo occorre definire dei criteri di arresto. Fissare un numero massimo (K_{max}) di iterazioni ne rappresenta uno.

Un altro criterio d'arresto, che riguarda la lontananza di approssimazione $x^{(k)}$, consiste nel fissare una tolleranza relativa ($toll_r$ ($\geq \epsilon_r$)) o assoluta ($toll_a$) e nell'arrestare il processo all'iterazione \bar{k} se $x^{(\bar{k})}$ soddisfa la seguente disuguaglianza:

$$\|x^{(k+1)} - x^{(k)}\| \leq toll_r \|x^{(k+1)}\|$$

oppure

$$\|x^{(k+1)} - x^{(k)}\| \leq toll_a$$

COMANDI MATLAB: METODO JACOBI

$$A = [\dots ; \dots ; \dots];$$

$$b = [\dots ; \dots ; \dots];$$

$$K_{max} = \dots;$$

$$toll = \dots;$$

$$m = \text{length}(b);$$

$$x_0 = \text{zeros}(m, 1);$$

$$D = \text{diag}(\text{diag}(A));$$

$$C = A - D;$$

for $K = 1:K_{max}$

$$x_1 = D \setminus (b - C * x_0)$$

$$\text{if } \text{abs}(x_1 - x_0) \leq toll;$$

break

end

$$x_0 = x_1;$$

end

METODO GAUSS-SEIDEL:

$$A = [\dots ; \dots ; \dots];$$

$$b = [\dots ; \dots ; \dots];$$

$$K_{max} = \dots;$$

$$toll = \dots;$$

$$m = \text{length}(b);$$

$$x_0 = \text{zeros}(m, 1);$$

$$D = \text{tril}(A);$$

$$C = A - D;$$

for $K = 1:K_{max}$

$$x_1 = D \setminus (b - C * x_0)$$

$$\text{if } \text{abs}(x_1 - x_0) \leq toll;$$

break

end

$$x_0 = x_1;$$

end

CONVERGENZA DEI METODI DI JACOBI E GAUSS-SEIDEL:

- Si dimostra che se A è a diagonale dominante per righe allora: $\|I - D^{-1}A\|_{\infty} < 1$
 Se A è a diagonale dominante per colonne: $\|I - D^{-1}A\|_1 < 1$.
- Pertanto se A è a diagonale dominante per righe e colonne, allora i metodi di Jacobi e di Gauss-Seidel applicati al sistema lineare $Ax = b$, convergono.
- Si dimostra che se A è simmetrica definita positiva, allora il metodo di Gauss-Seidel converge.

Osserviamo inoltre che la convergenza del metodo di Gauss-Seidel non implica quella del metodo di Jacobi e viceversa.

Se entrambi convergono, in generale Gauss-Seidel converge più rapidamente.

La velocità di convergenza dipende dal raggio spettrale: tanto più esso è piccolo tanto più rapidamente converge.

PROPRIETÀ METODI ITERATIVI: $Ax = b$

- x viene determinate come limite di una successione di vettori convergente;
- x in aritmetica con precisione infinita di calcolo viene determinate in maniera approssimate; in aritmetica con precisione finita di calcolo può essere determinate con una precisione soddisfacente le richieste dell'utente;
- non modificano la matrice dei coefficienti A ;
- sono efficienti per matrici sparse (con la maggior parte degli elementi uguali a zero) e di grandi dimensioni.

- somme esponenziali di ordine m :

$$\tilde{f}(x) = e_m(x) = \sum_{k=0}^m a_k e^{kx}$$

$$\tilde{f}(x) = e_m(x) = \sum_{k=0}^m a_k e^{-kx}$$

È utilizzato per l'approssimazione di funzioni con comportamento di tipo esponenziale ma è approssimabile con dei polinomi.

- funzioni spline di ordine m :

Sono funzioni polinomiali a tratti dotate di derivate continue fino all'ordine $m-1$.

Sono utilizzate per approssimare funzioni che generalmente non sono ben approssimabili da polinomi poiché per raggiungere un certo grado di approssimazione richiedono un grado elevato del polinomio e i polinomi di grado elevato oscillano fortemente.

Osserviamo che ciascun elemento dei suddetti sottospazi si esprime come combinazione lineare di funzioni, che definiscono una base per il sottospazio in questione, mediante i parametri a_k oppure a_k, b_k .

- le funzioni che definiscono una base per il sottospazio lineare dei polinomi algebrici di grado m , sono:

$$\varphi_0(x) = 1, \quad \varphi_1(x) = x, \quad \dots \quad \varphi_m(x) = x^m$$

Ogni polinomio si esprime quindi come combinazione di questi monomi.

Se il suddetto sottospazio ha dunque dimensione finita $m+1$.

INDIVIDUARE UNA PARTICOLARE FUNZIONE APPROSSIMANTE:

Si adatterà uno dei seguenti criteri:

- criterio dell'interpolazione:

È utilizzato per l'approssimazione di una funzione.

- criterio dei minimi quadrati:

È utilizzato per l'approssimazione di un insieme di dati sperimentali.

Ora come funzione f_m fissiamo un polinomio e utilizziamo il criterio dell'interpolazione per determinare il polinomio f_m di grado minimo che passa per i n punti (x_i, y_i) . Per determinare il polinomio di grado minimo a grado $n-1$:

$$p_{n-1}(x) = c_1 x^{n-1} + c_2 x^{n-2} + \dots + c_{n-1} x + c_n$$

definito dagli n coefficienti c_1, \dots, c_n generici che vogliamo determinare in modo tale che:

$$p_{n-1}(x_i) = y_i \quad i = 1, \dots, n$$

Si dimostra che se $x_i \neq x_j$ per $i \neq j$ allora esiste uno ed un solo polinomio soddisfacente le condizioni di interpolazione.

Imponiamo il passaggio per i punti x_i :

$$\begin{cases} c_1 x_1^{n-1} + \dots + c_{n-1} x_1 + c_n = y_1 \\ c_1 x_2^{n-1} + \dots + c_{n-1} x_2 + c_n = y_2 \\ \vdots \\ c_1 x_n^{n-1} + \dots + c_{n-1} x_n + c_n = y_n \end{cases}$$

Si ottiene così un sistema lineare di n equazioni e c_1, \dots, c_n incognite.

In forma matriciale:

$$\begin{pmatrix} x_1^{n-1} & \dots & x_1 & 1 \\ x_2^{n-1} & \dots & x_2 & 1 \\ \vdots & & \vdots & \vdots \\ x_n^{n-1} & \dots & x_n & 1 \end{pmatrix} \begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \iff VC = y$$

La matrice dei coefficienti è la matrice di Van der Monde, che è notoriamente mal condizionata quindi la soluzione non è accurata in quanto a piccole variazioni nei dati corrispondono non piccole variazioni nel risultato.

In caso $x_i \neq x_j \Rightarrow \det(V) = \prod_{i > j} (x_i - x_j) \neq 0 \Rightarrow$ la matrice V è non singolare \Rightarrow il sistema lineare assume una sola soluzione:

$$\exists! c = (c_1 \dots c_n)^T : Vc = y.$$

Quindi esiste un solo polinomio interpolante i dati assegnati: $\exists! p_{n-1}$ passante per i punti (x_i, y_i) .

- Calcoliamo i coefficienti del polinomio interpolante la funzione $f(x) = 1/(1+x^2)$ in 13 punti equispaziati nell'intervallo $[0, 1]$

$$m = 13;$$

$$x = \text{linspace}(0, 1, m);$$

$$f = \text{inline}('1./(1+x.^2)');$$

$$y = f(x);$$

$$c = \text{polyfit}(x, y, m-1);$$

Warning: Polynomial is badly conditioned...

serve per definire l'espressione di una funzione.

RAPPRESENTAZIONE DI LAGRANGE DEL POLINOMIO DI INTERPOLAZIONE:

Prima di definire questa rappresentazione dobbiamo definire dei particolari polinomi interpolanti:

si considerino i dati: $(x_1, 0), \dots, (x_{j-1}, 0), (x_j, 1), (x_{j+1}, 0), \dots, (x_m, 0)$

supponiamo: $x_i \neq x_j \quad i \neq j$

Notiamo come le ordinate sono sempre 0 tranne quando ha ascissa x_j :

$$l_j(x_i) = \delta_{ij} = \begin{cases} 1 & \text{se } i=j \\ 0 & \text{se } i \neq j \end{cases}$$

$$l_1 \rightarrow (x_1, 1) (x_2, 0) \dots (x_m, 0)$$

$$l_2 \rightarrow (x_1, 0) (x_2, 1) \dots (x_m, 0) \quad \text{ecc.}$$

Si può dimostrare che:

$$l_j(x_j) = \frac{(x_j - x_1) \dots (x_j - x_{j-1})(x_j - x_{j+1}) \dots (x_j - x_m)}{(x_j - x_1) \dots (x_j - x_{j-1})(x_j - x_{j+1}) \dots (x_j - x_m)} = 1 \quad \text{per } j = 1, \dots, m$$

Questi polinomi $l_j(x)$ sono detti polinomi fondamentali di Lagrange.

A questo punto il polinomio di interpolazione lo si può scrivere:

$$P_{m-1}(x) = \sum_{j=1}^m y_j l_j(x)$$

$$P_{m-1}(x_i) = \sum_{j=1}^m y_j l_j(x_i) = \sum_{j=1}^m y_j \delta_{ij} = y_i$$

Lezione 8

Rappresentazione di Newton del polinomio di interpolazione

RAPPRESENTAZIONE DI NEWTON DEL POLINOMIO DI INTERPOLAZIONE:

Dati n punti distinti, definiamo le seguenti quantità:

$$f[x_1, x_2] = \frac{f(x_2) - f(x_1)}{x_2 - x_1} \quad \text{= differenza divisa}$$

differenza divisa di ordine 1 di $f(x)$,

$$f[x_1, x_2, x_3] = \frac{f[x_2, x_3] - f[x_1, x_2]}{x_3 - x_1}$$

differenza divisa di ordine 2 di $f(x)$, e più in generale:

$$f[x_1, \dots, x_m] = \frac{f[x_2, \dots, x_m] - f[x_1, \dots, x_{m-1}]}{x_m - x_1}$$

differenza divisa di ordine $m-1$.

Le differenze divise sono invarianti rispetto a permutazioni degli argomenti cioè:

$$- f[x_1, x_2] = f[x_2, x_1]$$

$$- f[x_1, x_2, x_3] = f[x_2, x_1, x_3] = f[x_3, x_2, x_1] \dots \text{ecc}$$

Si definisce rappresentazione di Newton del polinomio interpolante ai punti (x_i, y_i) con $y_i = f(x_i)$, la seguente espressione:

$$P_{m-1}(x) = f(x_1) + f[x_1, x_2](x-x_1) + f[x_1, x_2, x_3](x-x_1)(x-x_2) + \dots + f[x_1, x_2, \dots, x_m](x-x_1) \dots (x-x_{m-1})$$

I polinomi $1, x-x_1, \dots, (x-x_1) \dots (x-x_{m-1})$ sono detti polinomi fondamentali di Newton.

Verifichiamo che per $P_{m-1}(x)$ verifichi le proprietà di interpolazione

$$P_{m-1}(x_i) = y_i \quad \text{per ogni } i=1, \dots, m.$$

Verifichiamo che: $P_{m-1}(x_1) = y_1 = f(x_1)$

$$P_{m-1}(x_1) = f(x_1) + f[x_1, x_2](x_1 - x_1) + f[x_1, x_2, x_3](x_1 - x_1)(x_1 - x_2) + \dots$$

$$P_{m-1}(x_1) = f(x_1) \quad \checkmark$$

Verifichiamo che: $P_{m-1}(x_2) = f(x_2)$

$$P_{m-1}(x_2) = f(x_1) + f[x_1, x_2](x_2 - x_1) + f[x_1, x_2, x_3](x_2 - x_1)(x_2 - x_2) + \dots$$

- Scriviamo la rappresentazione di Newton del polinomio $p_3(x)$ interpolante i punti $(0, 1)$, $(1, -1)$, $(2, 1)$, $(\frac{1}{2}, 2)$

x	y		$i=1$	$i=2$	$i=3$
0	1				
1	-1	\swarrow	$\frac{-1-1}{1-0} = -2$		
2	1	\swarrow	$\frac{1-(-1)}{2-1} = 2$	\swarrow	$\frac{2-(-2)}{2-0} = \frac{2}{2} = 2$
$\frac{1}{2}$	2	\swarrow	$\frac{2-1}{\frac{1}{2}-2} = -\frac{2}{3}$	\swarrow	$\frac{-\frac{2}{3}-2}{\frac{1}{2}-1} = \frac{16}{3}$
					\swarrow
					$\frac{\frac{16}{3}-2}{\frac{1}{2}-0} = \frac{20}{3}$

$$p_3(x) = 1 + (-2)(x-0) + 2(x-0)(x-1) + \frac{20}{3}(x-0)(x-1)(x-2)$$

Comandi Matlab: mo

function a = diffdiv(x,y)

m = length(x);

for i=1:m-1

for j=m:-1:i+1

y(j) = (y(j)-y(j-1))/(x(j)-x(j-1));

end

end

a=y;

Descriviamo ora l'algoritmo che consente di valutare in un punto oppure in un vettore di punti il polinomio interpolante. A tal scopo consideriamo la rappresentazione di Newton del polinomio di grado 3 interpolante i punti $(x_1, f(x_1))$, $(x_2, f(x_2))$, $(x_3, f(x_3))$, $(x_4, f(x_4))$ e, indicando con a_i le differenze divise in esso presenti, scriviamo:

$$p_3(x) = a_1 + a_2(x-x_1) + a_3(x-x_1)(x-x_2) + a_4(x-x_1)(x-x_2)(x-x_3)$$

$a_i =$ differenze divise.

$$p_3(z) = a_1$$

$$p_3(z) = a_1 + a_2(z-x_1)$$

$$p_3(z) = p_3(z) + a_3(z-x_1)(z-x_2)$$

$$p_3(z) = p_3(z) + a_4(z-x_1)(z-x_2)(z-x_3)$$

	•	+
	1	2
	2	2
	2	2

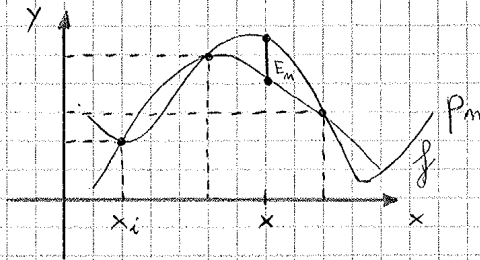
↳ In questo modo otteniamo ad ogni passaggio quasi sempre 2 moltiplicazioni e 2 addizioni.

CONVERGENZA DEL POLINOMIO DI INTERPOLAZIONE:

Denotato con p_m il polinomio interpolante la funzione f nei punti x_1, \dots, x_{m+1} , si definisce errore di interpolazione la funzione:

$$E_m(x) = f(x) - p_m(x)$$

Ci chiediamo se p_m è una buona approssimazione della funzione f e quanto piccolo è l'errore di interpolazione E_m



osserviamo che $E_m(x) = 0$ per $x = x_i$ ovvero $E_m(x_i) = f(x_i) - p_m(x_i) = 0$ e $E_m(x) = 0$ per ogni $x \in [a, b]$ e per $f = p$ con p polinomio di grado n ovvero, per unicità del polinomio interpolante, $p_m = p$.

Per quale x_i e f : $\lim_{m \rightarrow \infty} E_m(x) = 0$?

La risposta a questa domanda non è sempre affermativa; infatti si dimostra che, fissate una matrice di interpolazione ovvero una matrice di nodi:

$$\begin{matrix} x_1 \\ x_1 & x_2 \\ x_1 & x_2 & x_3 \\ x_1 & x_2 & x_3 & x_4 \\ \dots \end{matrix}$$

con $x_i \in [a, b]$ e $x_i \neq x_j$, $i \neq j$, è sempre possibile determinare una funzione $f \in [a, b]$ per la quale si abbia:

$$\lim_{m \rightarrow \infty} \|E_m\|_{\infty} \neq 0$$

Non è quindi detto che: $p_m(x) \rightarrow f(x)$ per $m \rightarrow +\infty \quad \forall f$

Se i nodi di interpolazione sono scelti coincidenti con gli zeri dei polinomi di Chebichev, $z_i = -\cos((2i-1)/(2m))\pi) \in (-1, 1)$, $i=1, \dots, m$ la condizione $f \in C^1[-1, 1]$ garantisce che:

$$\lim_{m \rightarrow \infty} \|E_m\|_{\infty} = 0 \quad \leftarrow \text{convergenza.}$$

- Se i nodi x_i sono scelti coincidenti con gli zeri dei polinomi di Chebichev traslati in $[s, S]$, cioè $x_i = \frac{b-a}{2} z_i + \frac{b+a}{2}$, $a = -s$ e $b = s$ e $f(x) = \frac{1}{s^2+x^2}$ (funzione di Runge), allora

$$\lim_{m \rightarrow \infty} \|E_m(x)\|_{\infty} = 0$$

$$P_m \rightarrow f \quad \text{per } m \rightarrow \infty$$

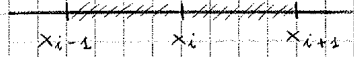
La convergenza dipende quindi da:

- scelta dei nodi di interpolazione peggior caso: equispaziati migliore: Chebichev
- regolarità funzione f

- $S_3(x_i) = y_i \quad i = 0, \dots, m$

Sono le condizioni di interpolazione.

$4m \quad \{a_i, b_i, c_i, d_i\} \quad i = 1, \dots, m$



$\lim_{x \rightarrow x_i^-} [a_i + b_i x + c_i x^2 + d_i x^3] = \lim_{x \rightarrow x_i^+} [a_{i+1} + b_{i+1} x + c_{i+1} x^2 + d_{i+1} x^3]$

$[a_i + b_i x_i + c_i x_i^2 + d_i x_i^3] = [a_{i+1} + b_{i+1} x_i + c_{i+1} x_i^2 + d_{i+1} x_i^3] \quad m-1 \text{ eq.}$

Calcolo derivate prima: $b_i + 2c_i x_i + 3d_i x_i^2 = b_{i+1} + 2c_{i+1} x_i + 3d_{i+1} x_i^2 \quad m-1 \text{ eq.}$

Calcolo derivate seconda: $2c_i + 6d_i x_i = 2c_{i+1} + 6d_{i+1} x_i \quad m-1 \text{ eq.}$

Imponendo la terza condizione abbiamo: $m+1$ equazioni

In definitiva abbiamo: $4m-2$ eq. in $4m$ incognite

Ma quindi bisogno di oltre due eq. per risolvere le $4m$ incognite $\{a_i, b_i, c_i, d_i\}$

Pertanto per definire una spline cubica occorre imporre oltre due ulteriori condizioni.

Fra tutte le possibili condizioni aggiuntive, le tipiche sono:

- $S_3^{(1)}(x_0) = 0 \quad S_3^{(2)}(x_m) = 0 \quad \rightarrow$ spline cubiche naturali.

- $S_3^{(3)}(x_i^-) = S_3^{(3)}(x_i^+) \quad S_3^{(3)}(x_{m-1}^-) = S_3^{(3)}(x_{m-1}^+) \quad \rightarrow$ condizione not-a-knot.

- $S_3^{(1)}(x_0) = f'(x_0) \quad S_3^{(1)}(x_m) = f'(x_m)$

Imponendo una delle suddette condizioni, la spline corrispondente esiste ed è unica.

Per determinare la spline cubica S_3 soddisfacente le condizioni suddette e una delle condizioni aggiuntive, non si procede andando a risolvere il sistema di $4m$ eq. in $4m$ incognite, ma si procede assumendo come incognite non i parametri a_i, b_i, c_i, d_i ma le quantità $\Pi_i = S_3^{(1)}(x_i)$, $i = 0, \dots, m$. In questo modo risolveremo un sistema lineare di ordine $m+1$ anziché $4m$ che gode di proprietà particolari e si impone una delle condizioni aggiuntive. Il sistema è simmetrico, a diagonale dominante e tridiagonale e si può risolvere con il metodo dell'eliminazione di Gauss.

o Lezione 10

Criterio dei minimi quadrati

Per individuare un'unica funzione approssimante vi sono due criteri:

- l'interpolazione
- dei minimi quadrati

Abbiamo già notato che l'interpolazione polinomiale non garantisce al crescere del grado del polinomio una maggiore accuratezza nell'approssimazione di una funzione. Quindi non è detto che all'aumentare del numero dei punti di interpolazione, la convergenza sia garantita.

Questo può essere superato con l'interpolazione mediante funzioni polinomiali a tratti. Essa tuttavia mal si presta ad essere utilizzata per estrapolare informazioni dai dati noti, ovvero per generare volutezioni in punti che giacciono al di fuori dell'intervallo di interpolazione.

o In quest'ultimo caso conviene utilizzare il criterio dei minimi quadrati.

Tale criterio viene adottato anche quando si vuole approssimare un insieme di dati (x_i, y_i) , $i=1, \dots, m$, con una funzione approssimante f_m che dipende da un numero $n \ll m$ di parametri.

Consideriamo, come funzione approssimante f_m , una combinazione lineare (a coefficienti costanti) delle funzioni di base $\varphi_k(x)$, $k=1, \dots, n$:

$$f_m(x) = C_1 \varphi_1(x) + C_2 \varphi_2(x) + \dots + C_n \varphi_n(x)$$

Il criterio dei minimi quadrati consiste allora nel determinare i coefficienti C_k del modello, imponendo che:

$$E^2(C_1, \dots, C_n) = \sum_{i=1}^m [y_i - f_m(x_i)]^2 = \sum_{i=1}^m \left[y_i - \sum_{k=1}^n C_k \varphi_k(x_i) \right]^2 = \text{volere minimo}$$

E^2 è denominata come residuo e noi vogliamo, imponendo che questo residuo risulti minimo, che la funzione passi il più vicino possibile ai punti.

Questa funzione f_m che minimizza il residuo prende il nome di migliore approssimazione dei dati secondo il criterio dei minimi quadrati.

osserviamo che per $n = m-1$, la migliore approssimazione f_m dei dati (x_i, y_i) , $i=1, \dots, m$, coincide con la funzione interpolante i punti (x_i, y_i) .

In fatti per essa si ha $E^2 = 0$.

Questa matrice $B = A^T A$, si dimostra una matrice simmetrica e semi-definita positiva, cioè $x^T B x \geq 0 \quad \forall x \neq 0$.
vettore nullo

B è definita positiva se e solo se le colonne di A sono linearmente indipendenti.

Pertanto, se le colonne di A sono linearmente indipendenti, il sistema delle eq. normali ammette una ed una sola soluzione e questa rende minima ϵ^2 .

↳ In tal caso, per la risoluzione del sistema $Bc = d$ si potrebbe utilizzare la decomposizione di Choleski $B = L_1 L_1^T$, ma risulta più efficiente procedere all'altro verso la decomposizione $A = QR$.

↳ Se le colonne ^{di A} sono invece linearmente dipendenti, il sistema ammette infinite soluzioni; si dimostra che solo una di queste ha norma euclidea minima. Pertanto in questo caso assumiamo come soluzione del nostro problema ai minimi quadrati, il vettore c^* con definito:

$$\sum_{i=1}^m [y_i - (c_1^* \varphi_1(x_i) + c_2^* \varphi_2(x_i) + \dots + c_m^* \varphi_m(x_i))]^2 = \text{minimo}$$

$$\|c^*\|_2^2 = \text{minimo}$$

- Assoguiti i dati (x_i, y_i) , $i=1, \dots, m$, determiniamo il polinomio di grado 1, $p_1(x) = c_1 x + c_2$, che meglio approssima i suddetti dati secondo il criterio dei minimi quadrati.

↳ coefficienti c_1 e c_2 si ottengono come soluzione del sistema delle eq. normali:

$$\begin{pmatrix} \sum_{i=1}^m x_i^2 & \sum_{i=1}^m x_i \\ \sum_{i=1}^m x_i & m \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^m y_i x_i \\ \sum_{i=1}^m y_i \end{pmatrix}$$

e definiscono la cosiddetta retta di regressione lineare.

$$f_m = p_1 = c_1 x + c_2$$

$$\begin{pmatrix} \sum_{i=1}^m \varphi_1(x_i) \varphi_1(x_i) & \sum_{i=1}^m \varphi_1(x_i) \varphi_2(x_i) \\ \sum_{i=1}^m \varphi_1(x_i) \varphi_2(x_i) & \sum_{i=1}^m \varphi_2(x_i) \varphi_2(x_i) \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^m y_i \varphi_1(x_i) \\ \sum_{i=1}^m y_i \varphi_2(x_i) \end{pmatrix}$$

- Data la funzione:
$$S(x) = \begin{cases} 3+x-9x^3 & 0 \leq x \leq 1 \\ a+b(x-1)+c(x-1)^2+d(x-1)^3 & 1 \leq x \leq 2 \end{cases}$$
- Determinazione a, b, c in modo che $S(x)$ sia una spline cubica in $[0, 2]$.
- Determinazione d , se esiste, che rende $S(x)$ una spline cubica naturale.

$$-S_3(1^-) = S_3(1^+) \quad 3+1-9 = a \quad \Rightarrow \quad a = -5$$

$$S_3^{(1)}(1^-) = S_3^{(1)}(1^+) \quad 1-27x^2 \Big|_{x=1} = b + 2c(x-1) + 3d(x-1)^2 \Big|_{x=1}$$

$$1-27 = b \quad \Rightarrow \quad b = -26$$

$$S_3^{(2)}(1^-) = S_3^{(2)}(1^+) \quad -54x \Big|_{x=1} = 2c + 6d(x-1) \Big|_{x=1}$$

$$-54 = 2c \quad \Rightarrow \quad c = -27$$

$$S(x) \in C^2$$

$-S''(0) = 0 \quad S''(2) = 0$. Se sono soddisfatte queste due condizioni la funzione è una spline cubica naturale.

$$S''_{[0,1]} = 54x \quad \Rightarrow \quad S''(0) = 0$$

$$S''_{[1,2]} = 2c + 6d(x-1) \quad \Rightarrow \quad 2c + 6d = 0 \quad \Rightarrow \quad -54 + 6d = 0$$

$$\Rightarrow \quad d = 9$$

$$S''(2) = 0 \quad \text{per} \quad d = 9$$

Per stimare la rapidità di convergenza di un metodo, si introduce la seguente definizione:

si definisce ordine di convergenza di una successione $\{x_n\}$ convergente ad α , il numero reale $p \geq 1$ tale che:

$$\lim_{n \rightarrow \infty} \frac{|\alpha - x_{n+1}|}{|\alpha - x_n|^p} = C \neq 0 \text{ e } +\infty$$

α rappresenta il limite della successione x_n .

Se questa condizione di limite è verificata, allora la successione $\{x_n\}$ ha ordine di successione p .

- Se $p=1 \Rightarrow C < 1 \Rightarrow$ convergenza lineare
- Se $1 < p < 2 \Rightarrow$ convergenza superlineare
- Se $p=2 \Rightarrow$ convergenza quadratica
- Se $p=3 \Rightarrow$ convergenza cubica.

Se la successione $\{x_n\}$ è definita da un metodo iterativo e ha ordine p , allora si dice che il metodo iterativo ha ordine p di convergenza.

Supponiamo che un metodo abbia ordine di convergenza p allora vale la definizione di limite vista precedentemente, ossia per n sufficientemente grande

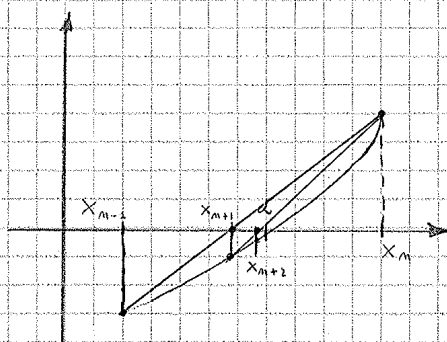
$$\lim_n \frac{|\alpha - x_{n+1}|}{|\alpha - x_n|^p} = C \Rightarrow |\alpha - x_{n+1}| \approx C |\alpha - x_n|^p \Rightarrow$$

$$\Rightarrow \text{se } |\alpha - x_n| \leq 10^{-k} \text{ allora } |\alpha - x_{n+1}| \approx C \cdot 10^{-kp}$$

Dunque per n sufficientemente grande se x_n ha k decimali corretti, il numero di decimali corretti presenti in x_{n+1} è dell'ordine di $k \cdot p$.

METODO DELLE SECANTI:

graficamente può essere rappresentato così:



Nel metodo delle secanti si parte da due valori a e b che inalterano con x_0 e x_1 e si definisce x_{m+1} con $m=1, 2, \dots$ come l'intersezione con l'asse delle ascisse della retta passante per i punti $(x_{m-1}, f(x_{m-1}))$ e $(x_m, f(x_m))$:

$$\begin{cases} y = f(x_m) + \frac{f(x_m) - f(x_{m-1})}{x_m - x_{m-1}} (x - x_m) \\ y = 0 \end{cases}$$

$$x_{m+1} \equiv x = x_m - f(x_m) \frac{x_m - x_{m-1}}{f(x_m) - f(x_{m-1})} \quad m = 1, 2, \dots$$

Se la funzione $f(x)$ è tale per cui tutte le successive approssimazioni appartengono all'intervallo di partenza $[a, b]$ e se d è una radice semplice di $f(x) = 0$, si dimostra che se $f \in C^2([a, b])$ e $|d - x_0|, |d - x_1|$ sono sufficientemente piccoli, allora:

$$\lim_{m \rightarrow \infty} x_m = d$$

e la convergenza è garantita.

L'ordine di convergenza del metodo è $p = \frac{1 + \sqrt{5}}{2} \approx 1,618$ ed è una convergenza di ordine superlineare $p > 2$.

Comandi Matlab:

```
f = inline('...');
x0 = ...;          x1 = ...;
mmax = ...;       toll = ...;
for m = 1:mmax
    x2 = x1 - f(x1) * (x1 - x0) / (f(x1) - f(x0));
    if abs(x2 - x1) <= toll
        break
    end
    x0 = x1;        x1 = x2;    end    x2
```


Comandi Matlab:

```
f = inline('---');  
f0l = inline('---');  
x0 = -;  
mmax = -;  
toll = -;  
for m = 1:mmax  
    x1 = x0 - f(x0) / f0l(x0);  
    if abs(x1 - x0) <= toll;  
        break  
    end  
    x0 = x1;  
end  
x1
```

Teorema:

Sia $I := [a, b]$. Se:

① $g(x) \in I$ per ogni $x \in I$ ($g(I) \subseteq I$)

② $g \in C^1(I)$

③ $|g'(x)| \leq m < 1 \quad \forall x \in I$

Se ①, ② e ③ sono soddisfatte allora esiste uno e uno solo punto fisso d di g in I e per $x_{n+1} = g(x_n)$ si ha:

$$\lim_{n \rightarrow \infty} x_n = d \quad \forall x_0 \in I$$

quindi la convergenza di x_n in d .

Inoltre, se $\exists I_\alpha : |g'(x)| > 1 \quad \forall x \in I_\alpha$, allora:

$$\lim_{n \rightarrow \infty} x_n \neq d$$

quindi x_n non converge ad d .

• Determiniamo la radice $\sqrt{2} \approx 1,4142135623$ dell'equazione

$f(x) = x^4 - 4 = 0$. Consideriamo:

i) $g(x) = x + \frac{f(x)}{f'(x)} \Rightarrow x_{n+1} = x_n + \frac{(x_n^4 - 4)}{4x_n^3}$

ii) $g(x) = x - \frac{f(x)}{f'(x)} \Rightarrow x_{n+1} = x_n - \frac{(x_n^4 - 4)}{4x_n^3}$ con $K = -11$

A partire da $x_0 = 1$, si ha:

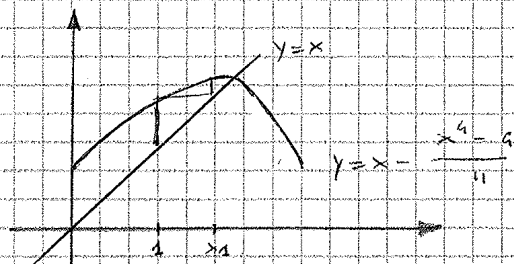
m	i	ii
1	-2	1,272727
2	10	1,337830
3	10 ²	1,414330
4	10 ⁴	1,414208
5	10 ⁶⁴	1,414213 $\approx \sqrt{2}$

Il diverso comportamento dei metodi i e ii dipende dal fatto che risulta:

i) $g(x) = x + \frac{(x^4 - 4)}{4x^3} \Rightarrow |g'(x)| = |1 + \frac{4x^3}{4x^3}| > 1$ per ogni $x > 0$

ii) $g(x) = x - \frac{(x^4 - 4)}{4x^3} \Rightarrow |g'(x)| = |1 - \frac{4x^3}{4x^3}| < 1$ per ogni $x \in (0, \sqrt[3]{\frac{11}{2}})$.

Pertanto, posto $I = (0, \sqrt[3]{\frac{11}{2}})$ il teorema precedente garantisce la convergenza a partire da qualsiasi $x_0 \in I$. Graficamente si ha:



$$f(x) = e^x - 2x^2$$

$$x_{m+1} = x_m - \frac{e^{x_m} - 2x_m^2}{e^{x_m} - 4x_m}$$

da applicare e partire da $x_0 = 0$

$$x_0 = 0, \quad x_1 = -1, \quad x_2 = -0,6263\dots, \quad x_3 = -0,5441\dots \approx \alpha_1$$

- Proporre un metodo iterativo ^{convergente} del tipo $x_{m+1} = g(x_m)$ per approssimare la radice negativa dell'eq. $e^x - 2x^2 = 0$.

Calcolare le prime tre iterate a partire da $x_0 = 0$.

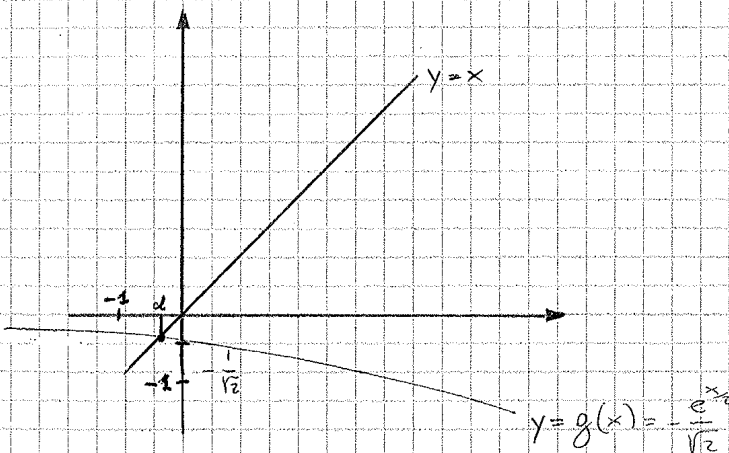
$$f(x) = e^x - 2x^2 \quad f(x) \equiv 0$$

Usando il metodo del punto fisso:

$$2x^2 = e^x \Rightarrow x = \pm \sqrt{\frac{e^x}{2}}$$

$$g(x) = -\sqrt{\frac{e^x}{2}} = -\frac{e^{x/2}}{\sqrt{2}}$$

$$x_{m+1} = g(x_m) \rightarrow x_{m+1} = -\frac{e^{x_m/2}}{\sqrt{2}}$$



$$x = g(x) \\ g(x) = -\frac{e^{x/2}}{\sqrt{2}}$$

$$I = (-\infty, 0]$$

$$i) g(I) \subseteq I$$

$$g(I) = \left[-\frac{1}{\sqrt{2}}, 0\right] \subseteq (-\infty, 0] \quad \checkmark$$

$$ii) g(I) \in C^1(I)$$

$$g = -\frac{e^{x/2}}{\sqrt{2}} \quad \checkmark$$

$$iii) |g'(x)| \leq m < 1 \quad \forall x \in I$$

$$|g'(x)| = \left| -\frac{e^{x/2}}{2\sqrt{2}} \right| = \frac{e^{x/2}}{2\sqrt{2}} \leq \frac{1}{2\sqrt{2}} < 1$$

$$m = \frac{1}{2\sqrt{2}}$$

$$\frac{x}{2} < 0 \quad \exp\left(\frac{x}{2}\right) \leq \exp(0) = 1 \quad \checkmark$$

Allora il teorema mi dice che:

$$\exists! \alpha : g(\alpha) = \alpha$$

$$x_m : x_{m+1} = g(x_m)$$

$$x_m \rightarrow \alpha \quad \forall x_0 \in I$$

Lezione 13

Calcolo degli integrali

FORMULE DI QUADRATURA INTERPOLATORIE:

Consideriamo il seguente problema: il calcolo dell'integrale definito su un intervallo chiuso e limitato:

$$\int_a^b f(x) dx$$

Ricordiamo che non è sempre possibile calcolare analiticamente un integrale, per esempio $\int_0^1 e^{-x^2} dx$ non è risolvibile mediante il calcolo della funzione primitiva con funzioni elementari.

Quando anche ciò fosse possibile potrebbe risultare complicato calcolarlo, per esempio l'integrale coincide con la somma di una serie oppure è definito mediante funzioni elementari e non, approssimate.

In queste situazioni e quando f è nota solo per punti oppure è valutabile per ogni valore di x mediante un algoritmo, si rende allora necessario l'utilizzo di un metodo numerico in grado di restituire un valore approssimato dell'integrale, indipendentemente dalla complessità della funzione integranda.

Un metodo numerico per il calcolo di un integrale definito assume la seguente espressione:

$$\int_a^b f(x) dx \approx \sum_{j=1}^m w_j f(x_j)$$

che è definita come formula di quadratura.

I punti x_j sono detti odi di quadratura e i numeri w_j coefficienti o pesi della formula di quadratura.

Fissati m modi distinti x_1, \dots, x_m nell'intervallo $[a, b]$, si definiscono interpolatorie le formule di quadratura che si ottengono approssimando l'integrale assegnato con l'integrale del polinomio $p_{m-1}(f; x)$ che interpola la funzione integranda f nei suddetti modi, cioè:

$$\int_a^b f(x) dx \approx \int_a^b p_{m-1}(f; x) dx \quad f \approx p_{m-1}(f; x)$$

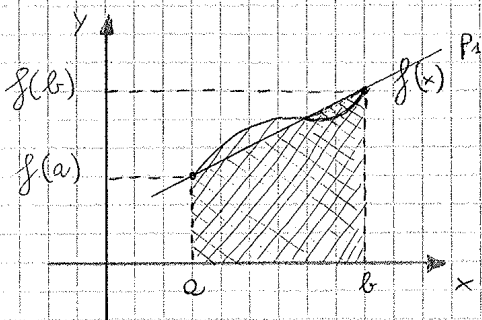
Il vantaggio sta nel fatto che l'integrale di un polinomio lo si può calcolare analiticamente. Quindi:

$$\int_a^b p_{m-1}(f; x) dx = \int_a^b \sum_{j=1}^m l_j(x) f(x_j) dx =$$

rappresentazione di Lagrange del polinomio interpolante

- per $n=2$, $x_1=a$ e $x_2=b$

si ha la formula del trapezio:



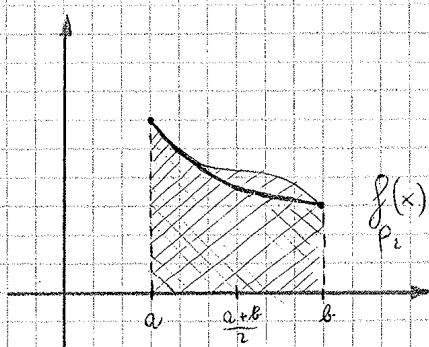
$$\int_a^b f(x) dx \approx \int_a^b P_1(f; x) dx =$$

$$= \int_a^b \left[\underbrace{\frac{x-b}{a-b}}_{L_1(x)} f(a) + \underbrace{\frac{x-a}{b-a}}_{L_2(x)} f(b) \right] dx =$$

$$= \frac{b-a}{2} [f(a) + f(b)]$$

Si approssima l'integrale $\int_a^b f(x) dx$ con l'area del trapezio di base maggiore $f(b)$ e base minore $f(a)$ per l'altezza $b-a$ diviso 2.

- per $n=3$, $x_1=a$, $x_2=\frac{a+b}{2}$ e $x_3=b$ si ha la formula di Simpson:



$$\int_a^b f(x) dx \approx \int_a^b P_2(f; x) dx =$$

$$= \int_a^b \left[\frac{(x-\frac{a+b}{2})(x-b)}{(a-\frac{a+b}{2})(a-b)} f(a) + \frac{(x-a)(x-b)}{(\frac{a+b}{2}-a)(\frac{a+b}{2}-b)} f(\frac{a+b}{2}) + \frac{(x-a)(x-\frac{a+b}{2})}{(b-a)(b-\frac{a+b}{2})} f(b) \right] dx =$$

$$= \frac{b-a}{2} \left[f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right]$$

Si approssima $\int_a^b f(x) dx$ con l'arco di parabola passante per $a, (a+b)/2$ e b .

Supponiamo che la funzione integranda $f(x)$ sia una funzione integrabile, ma non abbastanza regolare (per esempio $f(x)$ o una delle sue prime derivate presenta delle irregolarità o punti di discontinuità in $[a, b]$).

È evidente che approssimare una funzione irregolare con un polinomio che è una funzione regolare non è vantaggioso.

Supponiamo inoltre che $f(x)$ sia fattorizzabile nella forma:

$$f(x) = w(x)g(x)$$

dove $w(x)$ è una funzione di forma semplice contenente le irregolarità di $f(x)$, mentre $g(x)$ rappresenta la parte regolare della funzione integranda $f(x)$.

Ricordiamo inoltre che, per ogni n , P_n possiede esattamente n zeri, tutti reali, distinti e appartenenti ad (a, b) .

Questo vuol dire che questi zeri del polinomio P_n ortogonale in (a, b) rispetto alla funzione peso $w(x)$, possono essere scelti come nodi di quadratura, poiché reali, distinti e appartenenti ad (a, b) .

Se io scelgo questi come nodi x_i su cui interpolare e vado a costruire le formule di quadratura pesate interpolatorie associate a questi nodi, vado a costruire delle formule dal nome di formule Gaussiane.

FORMULE GAUSSIANE:

Sono delle formule interpolatorie pesate che utilizzano come nodi x_j gli zeri dei polinomi ortogonali in $[a, b]$ rispetto alla funzione peso $w(x)$.

- per $w(x) = 1$, $[a, b] = [-1, 1]$ e x_j^{GL} coincidenti con gli zeri del polinomio (di Legendre) di grado n ortogonale in $[-1, 1]$ rispetto a $w(x)$, si ha la formula di Gauss-Legendre:

$$\int_{-1}^1 g(x) dx \approx \sum_{j=1}^n \bar{w}_j^{GL} g(x_j^{GL})$$

$$\text{con } \bar{w}_j^{GL} = \int_{-1}^1 l_j(x) dx \quad j = 1, \dots, n$$

- per $w(x) = \frac{1}{\sqrt{1-x^2}}$, $[a, b] = [-1, 1]$ e x_j^{GC} coincidenti con gli zeri del polinomio (di Chebichev) di grado n ortogonale in $[-1, 1]$ rispetto a $w(x)$, si ha la formula di Gauss-Chebichev:

$$\int_{-1}^1 \frac{1}{\sqrt{1-x^2}} g(x) dx \approx \sum_{j=1}^n \bar{w}_j^{GC} g(x_j^{GC})$$

$$\text{con } \bar{w}_j^{GC} = \frac{\pi}{n} \quad x_j^{GC} = -\cos\left(\frac{2j-1}{2n}\pi\right) \quad j = 1, \dots, n$$

- per $w(x) = (1-x)^\alpha (1+x)^\beta$, $\alpha, \beta > -1$, $[a, b] = [-1, 1]$ e x_j^{GS} coincidenti con gli zeri del polinomio (di Jacobi) di grado n ortogonale in $[-1, 1]$ rispetto a $w(x)$, si ha la formula di Gauss-Jacobi:

$$\int_{-1}^1 (1-x)^\alpha (1+x)^\beta g(x) dx \approx \sum_{j=1}^n \bar{w}_j^{GS} g(x_j^{GS})$$

$$\bar{w}_j^{GS} = \int_{-1}^1 (1-x)^\alpha (1+x)^\beta l_j(x) dx \quad j = 1, \dots, n$$

Lezione 14

Calcolo degli integrali

CONVERGENZA:Si definisce errore di quadratura della formula:

$$\int_a^b w(x) g(x) dx \approx \sum_{j=1}^m \bar{w}_j g(x_j)$$

la quantità:

$$R_m(g) = \int_a^b w(x) g(x) dx - \sum_{j=1}^m \bar{w}_j g(x_j)$$

Per le formule interpolatorie pesate:

$$R_m(g) = \int_a^b w(x) E_m(g; x) dx$$

dove $E_m(g; x) = g(x) - p_{m-1}(g; x)$ definisce l'errore di interpolazione

Osserviamo che fissato m , numero dei nodi di interpolazione, le formule interpolatorie sono esatte, cioè l'errore $R_m \equiv 0$, ogni qualvolta $g(x)$ è un polinomio di grado $m-1$. Infatti, in tal caso, per l'unicità del polinomio interpolante, g è il polinomio $p_{m-1}(g)$ (interpolante g in m nodi distinti) coincidono e, di conseguenza $E_m(g) \equiv 0$.

Pertanto, mediante le formule interpolatorie pesate definite a partire da m nodi distinti, possiamo calcolare esattamente gli integrali del tipo $\int_a^b w(x) g(x) dx$, in cui g è un qualsiasi polinomio di grado $m-1$.

Quindi:

$$\int_a^b w(x) g(x) dx = \sum_{j=1}^m \bar{w}_j g(x_j) \quad (\text{e non più } \approx)$$

Si dimostra che le formule di Newton-Cotes a m nodi e chiuse (ovvero quelle che includono tra i loro nodi anche gli estremi di integrazione) sono esatte (cioè $R_m(f) \equiv 0$) ogni qual volta f è un polinomio di grado $m-1$ se m è pari, di grado m se m è dispari.

- la formula del trapezio ($m=2$) è esatta per tutti i polinomi di grado 1.
- la formula di Simpson ($m=3$) è esatta per tutti i polinomi di grado 3.

Si dimostra che le formule gaussiane a m nodi sono esatte ogni qualvolta g è un polinomio di grado $2m-1$.

Una formula di quadratura si dice convergente se:

$$\lim_{m \rightarrow \infty} R_m(g) = 0$$

Teorema:

Se: - $g \in C[a, b]$ con $\int_a^b w(x)g(x)dx \approx \sum_{j=1}^m \bar{w}_j g(x_j)$
 - $\sum_{j=1}^m |\bar{w}_j| \leq M$ con M dipendente da m

allora si può dimostrare: $\lim_{m \rightarrow \infty} R_m(g) = 0$

Inoltre se: - $g \in C^k[a, b]$

allora:

$$|R_m(g)| = O\left(\frac{1}{m^k}\right), m \rightarrow \infty$$

Questo ci dice che, più è regolare g e più è rapida la convergenza.
 (più k è grande e più è rapida la convergenza).

Questo teorema garantisce la convergenza delle formule gaussiane per ogni $g \in C[a, b]$. Infatti si ha:

$$0 < \int_a^b w(x) l_i^2(x) dx = \sum_{j=1}^m \bar{w}_j l_i^2(x_j) = \bar{w}_i$$

Pertanto, i pesi di una Gaussiana sono tutti positivi e inoltre:

$$\sum_{j=1}^m |\bar{w}_j| = \sum_{j=1}^m \bar{w}_j = \int_a^b w(x) dx = M$$

• Calcoliamo l'integrale: $I = \int_0^1 x^3 e^x dx = -2e^1 + 6 \approx 0,563436...$
 mediante formula di Gauss Legendre l_m^{GL} con m nodi di quadratura.

Cambiamento di variabile: $x = \frac{b+a}{2} t + \frac{b-a}{2} = \frac{1}{2} t + \frac{1}{2}$

m	l_m^{GL}	$ I - l_m^{GL} $
2	0,5455...	$1,730 \cdot 10^{-2}$
3	0,56323634...	$1,42 \cdot 10^{-4}$
4	0,56343534...	$4,00 \cdot 10^{-7}$
5	0,56343634...	$5,66 \cdot 10^{-10}$
6	0,563436343081...	$4,76 \cdot 10^{-13}$
7	0,56343634308130...	$6,66 \cdot 10^{-16}$

- Scriviamo la versione composta della formula del trapezio, detta formula dei trapezi, utilizzando una partizione uniforme (ovvero costante) dell'intervallo di integrazione. A tal scopo poniamo $h = \frac{b-a}{m}$, consideriamo i punti:

$$a \equiv x_1 < x_2 < \dots < x_m < x_{m+1} \equiv b$$

$$x_i = a + (i-1)h \quad i = 1, \dots, m+1$$

ed approssimiamo, mediante formula del trapezio, l'integrale esteso all'intervallo di integrazione $[x_i, x_{i+1}] \quad i = 1, \dots, m$:

$$\begin{aligned} \int_a^b f(x) dx &= \sum_{i=1}^m \int_{x_i}^{x_{i+1}} f(x) dx \approx \sum_{i=1}^m \frac{x_{i+1} - x_i}{2} [f(x_i) + f(x_{i+1})] = \\ &= \frac{h}{2} [f(x_1) + f(x_2) + f(x_2) + f(x_3) + \dots + f(x_m) + f(x_{m+1})] = \\ &= \frac{h}{2} [f(x_1) + 2 \sum_{i=2}^m f(x_i) + f(x_{m+1})] \end{aligned}$$

Comandi Matlab:

function t = trapezi (f, a, b, m)

%

x = linspace (a, b, m+1);

y = f(x);

t = (b-a) / (2*m) * (y(1) + 2*sum(y(2:m)) + y(m+1));

- Calcoliamo l'integrale:

$$I = \int_0^1 x^3 e^x dx = -2e^1 + 6 \approx 0,5634363430819089$$

mediante la formula dei trapezi I_m^T con $m = 2^m + 1$ e $m = 2, \dots, 15$

m	I_m^T	$ I - I_m^T $
5,0000 e + 00	0,6136008380528142	5,62 e - 02
8,0000 e + 00	0,5775647357738433	1,61 e - 02
9,7000 e + 01	0,5663739616442370	3,54 e - 03
3,3000 e + 01	0,5643210859873768	8,85 e - 04
6,5000 e + 01	0,5636575502632232	2,21 e - 04
1,2300 e + 02	0,5634316462201426	5,53 e - 05
2,5700 e + 02	0,5634501689503072	1,38 e - 05
5,1300 e + 02	0,5634397385542483	3,46 e - 06
⋮	⋮	⋮
3,2763 e + 01	0,5634363433257720	8,44 e - 10

Si dimostra che:

- le formule composte convergono per $m \rightarrow \infty$, ovvero per $h \rightarrow 0$, per $f \in C[a, b]$
- per la formula dei trapezi si ha:

$$R_m(f) = - \frac{(b-a)}{12} h^2 f''(\xi) \quad \text{con } \xi \in (a, b)$$

e per $f \in C^2[a, b]$; se f è periodica e il periodo coincide con $b-a$ (a è un sottomultiplo di $b-a$), allora l'ordine di convergenza è superiore a $O(h^2)$. Infatti:

$$|R_m(f)| = O(h^{2k+1})$$

per $f \in C^{2k+1}[a, b]$ e $f^{(2j-1)}(a) = f^{(2j-1)}(b)$ $j=1 \dots k$

- per la formula composta di Simpson si ha:

$$|R_m(f)| = O(h^4) \quad \text{per } f \in C^4[a, b]$$

Si nota come le formule composte convergano certamente purché f sia continua, però l'ordine di convergenza non dipende dalla regolarità della funzione f , è quindi fisso: $O(h^2)$ per la formula dei trapezi e $O(h^4)$ per la formula di Simpson.

Per funzioni poco regolari conviene considerare non più le partizioni uniformi ma quelle di tipo adattativo, per le quali si ha un maggiore addensamento di nodi in prossimità delle irregolarità.

Dunque là dove la funzione è poco regolare consideriamo una partizione più densa di nodi e viceversa là dove la funzione è regolare.

Comandi Matlab:

$$[q, N] = \text{quad}(f, a, b, \text{toll})$$

Calcola, mediante la formula composta di Simpson ed una partizione di tipo adattativo, il valore approssimato q dell'integrale $\int_a^b f(x) dx$ con una tolleranza assoluta toll ed utilizzando N valutazioni di funzione ($\text{toll} = 10^{-6}$ default)