

# MATEMATICA COMPUTAZIONALE

libro: "Fondamenti di calcolo numerico" G. Honegatto, CLUT 5/10/2005

Argomento di oggi interamente preso dal libro.

Sistema di equazioni lineari

$$Ax = b \quad A \in \mathbb{R}^{n \times m} \xrightarrow{\text{MATRICI}}, \quad b \in \mathbb{R}^n \xrightarrow{\text{VETTORE}} x?$$

metodi diretti (soluz. esatta in un numero finito di passi, alterando A)

iterativi (non alterano A, costruiscono una successione di vettori, convergono sotto opportune ipotesi alla soluzione esatta, arresta il processo appena ottenuta la precisione richiesta)

$$A = \begin{pmatrix} \emptyset & \text{||} & \emptyset \\ \text{||} & \emptyset & \text{||} \\ \emptyset & \text{||} & \emptyset \end{pmatrix}$$

matrice sparsa = quando il numero degli elementi nulli è molto grande. Per questa è meglio usare il metodo iterativo.

Quando la matrice è piena è meglio usare il metodo diretto.

Metodo di Gauss, eliminazione

A triangolare superiore

$$A = \begin{pmatrix} \times & & \\ \emptyset & \times & \\ & \emptyset & \times \end{pmatrix}$$

$$Ax = b$$

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n = b_1 \\ \dots \\ a_{22}x_2 + a_{23}x_3 + \dots + a_{2n}x_n = b_2 \\ \vdots \\ a_{nn}x_n = b_n \end{cases}$$

Partendo dall'ultima

$$x_n = \frac{b_n}{a_{nn}} \quad \text{si va a } x_{n-1} \text{ fino a } x_1.$$

$$x_k = \frac{b_k - \sum_{s=k+1}^n a_{ks}x_s}{a_{kk}}, \quad k = n-1, \dots, 1$$

Per la triangolare inferiore sarebbe la stessa cosa  $A = \begin{pmatrix} \times & & \\ & \times & \\ & & \times \end{pmatrix}$ .

In sintesi vogliamo trasformare un generico sistema in uno equivalente di forma triangolare.

1) Quando ad una equaz. sostituiamo una comb. lineare dell'equazione con un'altra dello stesso sistema, il nuovo sistema risulta equivalente al precedente.

2) Il metodo di Gauss dimostra che è sempre possibile mediante un numero finito di combinazioni lineari di quel tipo ed eventuali permutazioni di equazioni, trasformare un generico sistema in un sistema equivalente di forma triangolare.

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + a_{14}x_4 + \dots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 + a_{24}x_4 + \dots + a_{2n}x_n = b_2 \end{cases} \quad \text{assumiamo } a_{11} \neq 0$$

elimino  $x_1$  dalle ultime  $m-1$  equazioni, sommando alla  $i$ -esima eq. la prima moltiplicata per  $m_{i1} = -\frac{a_{i1}}{a_{11}}$ ,  $i = 2 \dots m$

$$a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2$$

$$-\frac{a_{21}}{a_{11}} (a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1)$$

] somma di queste due

$$\emptyset x_1 + a_{22}^{(2)}x_2 + a_{23}^{(2)}x_3 + \dots = b_2^{(2)}$$

RESULTATO DELLA SOMMA

Quindi ora abbiamo:

$$a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1$$

$$\emptyset a_{22}^{(2)}x_2 + a_{23}^{(2)}x_3 + \dots + a_{2n}^{(2)}x_n = b_2^{(2)}$$

Applicando ripetutamente questa procedura ottengo un nuovo sistema dove  $x_1$  sarà sempre  $\emptyset$ .

$$a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1$$

$$\emptyset a_{22}^{(2)}x_2 + \dots + a_{2n}^{(2)}x_n = b_2^{(2)}$$

$$\emptyset a_{32}^{(2)}x_2 + \dots + a_{3n}^{(2)}x_n = b_3^{(2)}$$

Possiamo applicare la procedura alla matrice più interna, partendo dalla seconda equazione.

Lo facciamo per tutte le matrici interne finché non otteniamo una matrice triangolare inferiore, ovvero  $m-1$  volte.

Algoritmo

1) elimino le variabili in  $m-1$  passi

al passo  $k$ -esimo,  $k = 1 \dots m-1$

$$a_{ij}^{(k+1)} = a_{ij}^{(k)} + m_{ik} a_{kj}^{(k)}, \quad j = k+1 \dots n$$

$$m_{ik} = -\frac{a_{ik}^{(k)}}{a_{kk}^{(k)}}, \quad b_i^{(k+1)} = b_i^{(k)} + m_{ik} b_k^{(k)}$$

2) risolvere sistema triangolare superiore.

Matematica computazionale NO!

6/10/2005

Non dà risultati specifici, ma modelli di calcolo.

La soluzione trovata sarà sempre approssimata, vincoli: accuratezza, velocità di calcolo, margine di errore.

ESEMPIO

$y(f) = \sin \log f$  equazione esplicita dove l'errore può essere nullo.

$y = e^{f^4} \sin y f$  equazione implicita ( $f$  e  $y$  non sono separate), da risolvere ma in modo approssimato. ( $f$  nota non esplicitamente).

$y(f) = \sum_{k=0}^{\infty} a_k e^{-b_k f^2}$  anche se la  $f$  è nota esplicitamente non pensiamo cal

colata (a causa della somma degli oo termini).

$y(\varphi) = \sin(\varphi)$  f. nota esplicitamente e calcolabile non direttamente (ma anche in questo caso c'è un errore)

$y(\varphi) = 3\varphi$  f. nota esplicitamente, calcolabile direttamente (anche in questo caso abbiamo un errore, perché non sappiamo come è  $\varphi$  e come si può rappresentare con un numero finito di bit).

Quindi avremo sempre e comunque un margine di errore.

Il computer calcola il log con le serie.

Se dovessimo calcolare:  $\sin \log \sqrt{\varphi^{2000}}$  bisognerebbe procedere dall'interno verso l'esterno, invece il computer ne usa una espressione approssimata tipo  $\sum a_k x^k$  ma estremamente più veloce da calcolare. Bisogna quindi sempre conoscere il margine di errore.

ESEMPIO

Quanto occupa una foto da  $1024 \times 768$  a colori?

$3 \cdot 1024 \cdot 768 (\cdot 1) = 2,4 \text{ Mb}$  da 0 a 255  $256 = 2^8$   
↓  
RGB 3 valori da 255 perché abbiamo supposto che un pixel occupa un byte di memoria

Se devo spedire la foto con un modem a  $40 \text{ KB/s}$ ?

$(2,4 \cdot 8) \cdot 1024 = 500 \text{ s} \approx 8 \text{ min}$   
40  
ha convertito i byte in bit

Possiamo convertirla in Jpeg che si basa sulla tecnica lossy, basata sulla trasformata di Fourier.

Gauss 19/10/2005

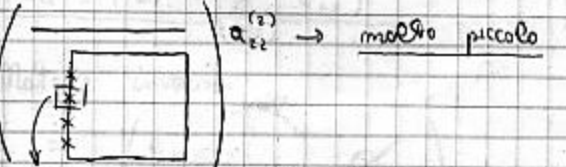
Il procedimento delle eliminazioni di Gauss può essere portato a termine senza permutare l'ordine iniziale delle eq. se  $A$  diag. dominante (per righe o colonne) o simmetrica def. positiva.

$-\frac{a_{ik}^{(k)}}{|a_{kk}^{(k)}} = m_{ik}$   $\left| \begin{matrix} a_{kk} & \text{pivot} \\ \# & \text{MAI IN QUESTO CASO} \\ \emptyset & \end{matrix} \right| \quad |a_{ii}| > \sum_{j \neq i} |a_{ij}| \quad \text{DIAG. DOMINANTE}$

$a_{ij} = a_{ji}$  simmetrica

$A$  def. pos.  $x^T A x > 0 \quad \forall x \neq 0$

tecnica del pivot parziale



scego valore abs. più grande e scambio le due righe  
 $a_{kk}^{(k)} = \max_i |a_{ik}^{(k)}|$

Fattorizzazione LU (lower-upper)

Interpretare il metodo di Gauss come successione finita di trasformazioni di A e b, cioè moltip. di A e b per opportune matrici.

Riformulare l'algoritmo di Gauss in 2 parti:

1)  $G \cdot \rightarrow GA = U$  triangolare sup.

2)  $GAX = Gb$ ,  $\bar{b} = Gb$ ,  $Ux = \bar{b}$

Vediamo 1) scambio eq. i-esima eq. j-esima

$$P_{i,j} = \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & \circ & \\ & & & 1 \end{pmatrix}$$

i-esima  
j-esima

$$I = \begin{pmatrix} 1 & \circ \\ \circ & 1 \end{pmatrix}$$

$$P_{i,j} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \rightarrow \text{ho scambiato le righe di } I$$

$$A = \begin{pmatrix} 4 & 3 \\ 2 & 5 \end{pmatrix} \xrightarrow{P_{2,1}} P_{2,1}A = \begin{pmatrix} 2 & 5 \\ 4 & 3 \end{pmatrix}$$

la sostituzione dell' eq. i-esima con la stessa più la j-esima moltiplicata per  $m_{i,j}$

$$M_{i,j} = \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & \circ & \\ & & & 1 \end{pmatrix}$$

$m_{i,j}$   
↓  
obiettivo

$$P_1, P_2, \dots, P_{m-1} \quad (I \text{ se non scambio})$$

$$M_1, M_2, \dots, M_{m-1}, \quad M_j = M_{m,j} \dots M_{2,j} = \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & \circ & \\ & & & 1 \end{pmatrix}$$

$m_{21,j}$   
 $m_{32,j}$   
 $\vdots$   
 $m_{m-1,j}$

Di fatto Gauss dice che:

$$\left. \begin{array}{l} Ax = b \\ M_1 P_1 Ax = M_1 P_1 b \end{array} \right] \text{1° passo}$$

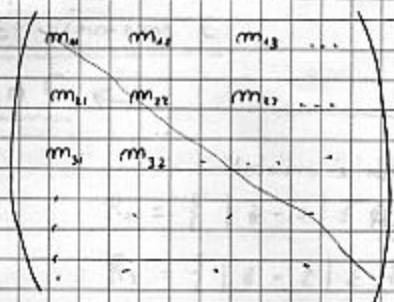
$$\left. M_2 P_2 M_1 P_1 Ax = M_2 P_2 M_1 P_1 b \right] \text{2° passo}$$

$$\vdots$$

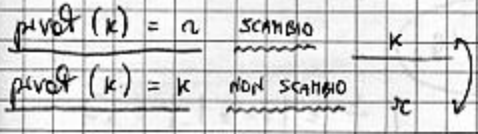
$$\left. M_{m-1} P_{m-1} \dots M_1 P_1 Ax = M_{m-1} P_{m-1} \dots M_1 P_1 b \right] \text{ultimo passo}$$

Quindi  $G$  sarà  $= M_{m-1} P_{m-1} \dots M_1 P_1$  costo computazionale  $\frac{n^3}{3}$  relativamente basso

g vena utilizzata per trasformare vettori e matrici, o per memorizzare  $m_{ij}$  e eventuali permutazioni



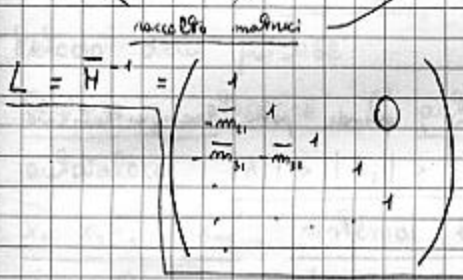
Non serve memorizzare tutte le matrici P, basta memorizzare vettore pivot:



$$H_{n-1} P_{n-1} \dots H_1 A = H_{n-1} P_{n-1} \dots H_1 P_1 b$$

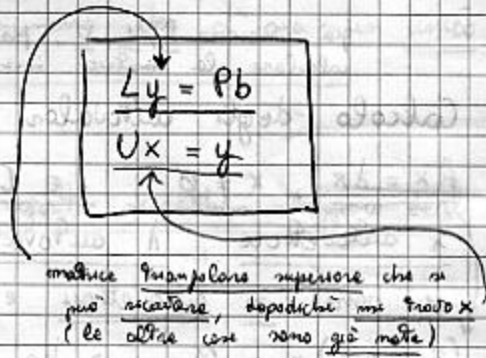
$$\underbrace{H_{n-1} \dots H_1}_{\bar{H}} \underbrace{P_{n-1} \dots P_1 A}_P = U$$

$$PA = \bar{H}^{-1} U$$



$$PA = LU$$

$$PAx = Pb \rightarrow LUx = Pb$$



Ho scomposto le cose in due matrici triangolari, una superiore e l'altra inferiore.

27/10/2005

Fattorizzazione LU

Se non ci sono scambi  $P = I \Rightarrow A = LU = LUD$ , (o)

con  $L$  e  $D$  triang. unitarie

Se  $A$  è simmetrica  $U_i = L_i^T$  se poi  $A$  è def. pos.  $\Rightarrow d_{ii} > 0 \Rightarrow A$  simmetrica

def. pos.  $A = L_i L_i^T$ ,  $L_i = LD^{1/2}$  ( $D^{1/2} = \sqrt{d_{ii}}$ )

$$l_{ii} = \sqrt{d_{ii}}$$

$$l_{ij} = (a_{ij} - \sum_{k=1}^{j-1} l_{ik} l_{jk}) / l_{jj}, \quad j = 1 \dots i-1$$

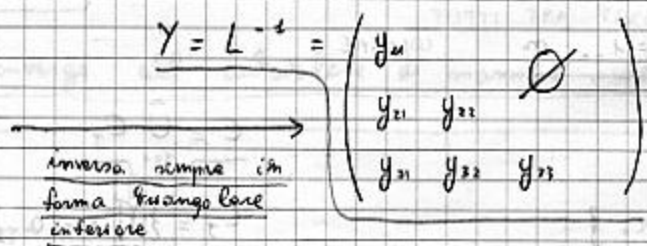
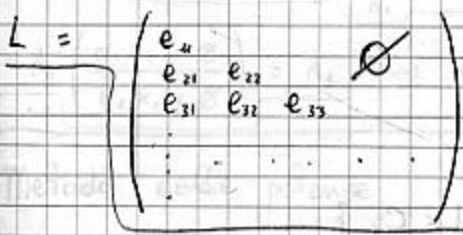
ALGORITHM CHOLESKI

lo chiede all'esame

$$l_{ii} = (a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2)^{1/2}$$

CHOLESKI  $\frac{m^3}{6}$  COMPLESSITA' COMPUTAZIONALE

Matrice inversa



$$y_{jj} = \frac{1}{l_{jj}} \quad j = 1 \dots n$$

$$y_{ij} = - \left( \sum_{k=j}^{i-1} l_{ik} y_{kj} \right) / l_{ii}$$

$i = j+1 \dots n$

Per calcolare l'inversa di una matrice generica è molto difficile, ma con le cose viste precedentemente diventa banale.

Come calcolare l'inversa

$$GA = U$$

$$U^{-1}GA = U^{-1}U = I \quad \rightarrow \text{moltiplicato per } U^{-1} \text{ entrambi i membri}$$

$$U^{-1}GA = I \Rightarrow A^{-1} = U^{-1}G \quad \rightarrow \text{matrice triangolare superiore}$$

$$\downarrow$$

se moltiplico per A entrambi i membri

$$PA = LU \quad \text{ci può servire per calcolare } A^{-1}!$$

$$P^{-1}PA = P^{-1}LU$$

$$A = P^{-1}LU$$

$$AU^{-1} = P^{-1}L \quad (\text{moltiplicato per } U^{-1})$$

$$AU^{-1}L^{-1} = P^{-1} \quad (\text{moltiplicato per } L^{-1})$$

$$AU^{-1}L^{-1}P = I$$

$$A^{-1} = U^{-1}L^{-1}P \quad \rightarrow \text{porta } A \text{ dall'altra parte che diventa } A^{-1}$$

già conoscevo P, L, U, per le inverse abbiamo già visto gli algoritmi, quindi posso tranquillamente calcolare la matrice inversa di A

Calcolo degli autovalori di A

$$Ax = \lambda x, \quad x \neq \emptyset, \quad \lambda \in \mathbb{C}$$

x autovettore,  $\lambda$  autovalore

ci sono m autovettori e m autovalori

$$(A - \lambda I)x = \emptyset \quad \rightarrow \text{ha portato } \lambda \text{ dall'altra parte e } \text{RACCOLTO } x$$

Voglio che abbia anche una soluzione non nulla.  $\rightarrow \det(A - \lambda I) = 0 \rightarrow$  vuol dire che ha più di una soluzione.

$$= \lambda^m + \alpha_1 \lambda^{m-1} + \alpha_2 \lambda^{m-2} + \dots$$

POLINOMIO CARATTERISTICO

Sto cercando le m radici di POL. CAR.

$\lambda$  incognita nei reali o nei complessi.

ci sono vari algoritmi, a seconda di quali autovalori vogliamo trovare.

teorema di localizzazione (Gerschgorin)

localizza l'area dei complessi dove posso cercare gli autovalori

$$r_i = \sum_{j \neq i} |a_{ij}|, \quad i = 1 \dots m \quad \text{RIGHE}$$

$$c_j = \sum_{i \neq j} |a_{ij}|, \quad j = 1 \dots m \quad \text{COLONNE}$$

$$R = \bigcup_{i=1}^m R_i$$

$$R_i = \{ z \in \mathbb{C} \mid |z - a_{ii}| \leq r_i \}$$

$$C = \bigcup_{j=1}^m C_j$$

$$C_j = \{ z \in \mathbb{C} \mid |z - a_{jj}| \leq c_j \}$$

ogni autovalore  $\in$  insieme  $\mathbb{R}$   
 $\in$  insieme  $\mathbb{C}$

$\mathbb{R} \cap \mathbb{C}$   $\leftarrow$  SOMMA DEI MODULI DEI RESTANTI

Esempio

$$A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 7 \\ \frac{1}{2} & 0 & 4 \end{pmatrix}$$

$$R_1 = \{ |z-1| \leq \emptyset \} \quad (*)$$

$$R_2 = \{ |z-2| \leq 7 \}$$

$$R_3 = \{ |z-4| \leq \frac{1}{2} \}$$

$$C_1 = \{ |z-1| \leq \frac{1}{2} \}$$

$$C_2 = \{ |z-2| \leq \emptyset \} \quad (*)$$

$$C_3 = \{ |z-4| \leq 7 \} \quad \text{UNIONE DEI TRE CERCHI}$$

l' ~~insieme~~ <sup>UNIONE</sup> dei tre cerchi data l'insieme dove cercare gli autovalori

(\*) quando il cerchio collana ad un unico punto  $\lambda=1$ ,  $\lambda=2$   
 d'altro autovalore va cercato nell'area risultata.

Metodo delle potenze

serve a trovare l'autovalore di massimo modulo, ammesso che ci sia un unico  
 autovalore  $|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq |\lambda_4| \dots$

$x_1, x_2, \dots, x_n$  sistema linearmente indipendente

$\lambda_i \in \mathbb{R} \rightarrow$  perché se fosse complesso anche il suo conjugato sarebbe maggiore degli altri

$$v_0 = \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n \rightarrow \text{generico vettore}$$

$$A v_0 = \alpha_1 A x_1 + \alpha_2 A x_2 + \dots + \alpha_n A x_n \rightarrow \text{moltiplico per matrice } A$$

$$A x_i = \lambda_i x_i \dots$$

$$\equiv \alpha_1 \lambda_1 x_1 + \alpha_2 \lambda_2 x_2 + \dots + \alpha_n \lambda_n x_n =$$

$$A^2 v_0 = A A v_0 = A (\alpha_1 \lambda_1 x_1 + \alpha_2 \lambda_2 x_2 + \dots + \alpha_n \lambda_n x_n)$$

$$\equiv \alpha_1 \lambda_1 A x_1 + \alpha_2 \lambda_2 A x_2 + \dots + \alpha_n \lambda_n A x_n$$

$$\equiv \alpha_1 \lambda_1 \lambda_1 x_1 + \alpha_2 \lambda_2 \lambda_2 x_2 + \dots + \alpha_n \lambda_n \lambda_n x_n$$

$$\equiv \alpha_1 \lambda_1^2 x_1 + \alpha_2 \lambda_2^2 x_2 + \dots + \alpha_n \lambda_n^2 x_n =$$

$$A^k v_0 = \alpha_1 \lambda_1^k x_1 + \alpha_2 \lambda_2^k x_2 + \dots + \alpha_n \lambda_n^k x_n$$

$$\equiv \lambda_1^k (\alpha_1 x_1 + \alpha_2 \frac{\lambda_2^k}{\lambda_1^k} x_2 + \dots + \alpha_n \frac{\lambda_n^k}{\lambda_1^k} x_n)$$

$$\frac{\lambda_2^k}{\lambda_1^k} < 1$$

$$\lim_{k \rightarrow \infty} \frac{v_{k+1}}{v_k} = \frac{\lambda_1^{k+1} (\alpha_1 x_1 + \alpha_2 \frac{\lambda_2^{k+1}}{\lambda_1^{k+1}} x_2 + \dots)}{\lambda_1^k (\alpha_1 x_1 + \alpha_2 \frac{\lambda_2^k}{\lambda_1^k} x_2 + \dots)} =$$

$$\frac{v_{k+1}}{v_k} = \lambda_1 \left( \frac{\alpha_1 x_1 + \emptyset}{\alpha_1 x_1 + \emptyset} \right) = \lambda_1 \rightarrow \text{converge all'autovalore di massimo modulo}$$

I TERMINI SUCCESSIVI AL PRIMO SI SEMPLIFICANO TUTTI TRA LORO

Metodo delle potenze

2/11/2005

Converge tanto più rapidamente quanto più  $|λ_1| > |λ_2|$  (o discosta dai successivi)

\* il autovalore per  $A \Rightarrow λ^{-1}$  autovalore per  $A^{-1}$  perché

$$Ax = λx \Rightarrow A^{-1}Ax = λA^{-1}x \Rightarrow x = λA^{-1}x \Rightarrow λ^{-1}x = A^{-1}x$$

Il metodo delle potenze può essere utilizzato per migliorare l'approssimazione di un dato autovalore  $p$  approssimato di  $λ$  autovalore di  $A$ .

$$(A - pI)x = Ax - px = λx - px = (λ - p)x \rightarrow \text{MATRICE} \times \text{VETTORE} = \text{VETTORE} \times \text{SCALARE}$$

$$λ - p \text{ autovalore di } A - pI \Rightarrow (A - pI)^{-1} \text{ autovalore di } (A - pI)^{-1}$$

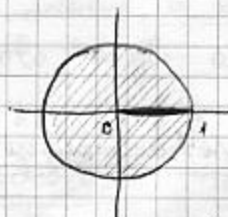
Chiamiamo  $μ = (λ - p)^{-1}$  immaginiamo che  $p$  sia approssimazione buona  $ACCURATA$   
 $\approx λ$   
 $\downarrow$   
 $λ - p$  è quasi  $0$  e il suo reciproco è molto grande

quindi  $μ$  autovalore di max modulo

$$(A - pI)^{-1}$$

quindi uso il metodo delle potenze e trovo una quantità più accurata per  $μ$ .

$$\frac{1}{μ} = λ - p \quad λ = p + \frac{1}{μ}$$



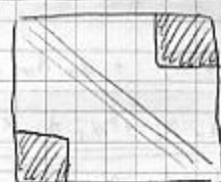
l'autovalore  $λ$  trova nel cerchio di raggio 1 vicino a  $p=1$  e con questo metodo si fa una migliore approssimazione (raffinamento)  $\rightarrow$  METODO DELLE POTENZE INVERSE

Torniamo a Gauss e ai sistemi iterativi su esso

$$Ax = b \quad x? \quad \text{con i 3 metodi iterativi?}$$

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2 \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 + \dots = b_3 \\ \vdots \end{cases}$$

applicati a matrici con struttura regolare con pochi elementi non nulli quasi tutti  $0$



$$\begin{aligned} x_1 &= \frac{b_1 - \sum_{j=2}^n a_{1j}x_j}{a_{11}} \\ x_2 &= \frac{b_2 - \sum_{j=3}^n a_{2j}x_j}{a_{22}} \\ x_3 &= \frac{b_3 - \sum_{j=4}^n a_{3j}x_j}{a_{33}} \end{aligned}$$

$\rightarrow$  input, vettore di partenza  
 $x^{(0)} = (x_1^{(0)}, \dots, x_n^{(0)})$

Dato un vettore di partenza  $\rightarrow$  devo trovare un vettore approssimato decente.

$$\begin{aligned} x^{(1)}, x^{(2)}, x^{(3)} &\rightarrow x \\ x_1^{(1)} &= \frac{b_1 - \sum_{j=2}^n a_{1j}x_j^{(0)}}{a_{11}} \end{aligned}$$



$$x_2^{(1)} = \frac{b_2 - \sum_{j \neq 2} a_{2j} x_j^{(0)}}{a_{22}}$$

$$x_3^{(1)} = \frac{b_3 - \sum_{j \neq 3} a_{3j} x_j^{(0)}}{a_{33}}$$

$$x_4^{(1)} = \dots$$

$$x_5^{(1)} = \dots$$
 ecc...

informazioni date in input → senza correzione si chiama JACOBI

si può scrivere anche:  $x_i^{(k)} = \frac{b_i - a_{ii} x_i^{(k-1)} - \sum_{j \neq i} a_{ij} x_j^{(k-1)}}{a_{ii}}$

con questa correzione si chiama metodo di GAUSS-SEIDEL (più facile)

al primo giro devo ottenere un vettore  $x^{(1)} = (x_1^{(1)}, x_2^{(1)}, \dots, x_n^{(1)})$

di fatto alla fine posso scrivere: e con via a tutti gli altri passi

$$x_1^{(k+1)} = \frac{b_1 - \sum_{j \neq 1} a_{1j} x_j^{(k)}}{a_{11}}$$

$$x_2^{(k+1)} = \dots$$

$$x_3^{(k+1)} = \dots$$
 ecc...

ora faccio un ulteriore raffinamento al secondo passo:

$$x_1^{(2)} = \frac{b_1 - \sum_{j \neq 1} a_{1j} x_j^{(1)}}{a_{11}}$$

per tutte le altre incognite...

dopo un po' di giri devo avere una buona approssimazione

terzo metodo (del sovrarilassamento)

$$x_i^{(k+1)} = x_i^{(k)} + \omega r_i^{(k)}$$

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

ancora è un mistero per gli studenti quale sia il  $\omega$  perfetto

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

Esercizio (per casa) FARLO BENE!

Verificare che scegliendo  $\omega = 1$  il metodo del sovrarilassamento si riduce a metodo di Gauss-Seidel. →  $x_i^{(k+1)} = \frac{b_i - \sum_{j \neq i} a_{ij} x_j^{(k+1)}}{a_{ii}}$

Approfondimento sulla lezione scorsa

PER RIGHE:  $|a_{ii}| \geq \sum_{j \neq i} |a_{ij}|$  (PER COLONNE:  $|a_{ii}| \geq \sum_{k \neq i} |a_{ki}|$ )

9/11/2005

i metodi convergono quando la matrice A è a diagonale dominante per righe o per colonne. → LO CHIEDE SEMPRE ALL'ESAME!!!

$$x_2^{(1)} = \frac{b_2 - a_{21} x_1^{(1)} - \sum_{j \neq 1,2} a_{2j} x_j^{(0)}}{a_{22}} \rightarrow \text{GAUSS-SEIDEL}$$

EQUAZIONI NON LINEARI

$f(x) = \phi$ 

$$x^{2.1} - 7x + 3 = \phi \rightarrow \text{come la risolvo?}$$

$$x^{10} - 7x + 7x^3 - 15x^2 = \phi$$

$$x^2 - \log_2 \log_2 x - 7 = 0$$

SICURAMENTE PIU' DIFFICILI DELLE EQ. LINEARI

Capita spesso di trovarle in fisica ad es...

Ci vogliono algoritmi che approssimano per poterle risolvere.

$x_1, x_2, x_3, \dots, x_n \rightarrow x$   $f(x) = 0$   
si risolvono con i metodi iterativi

$x_{n+1} = g(x_n)$   $g$  funzione di iterazione

PROCED. ITERATIVO

$x_1 = g(x_0)$

$x_2 = g(x_1)$

$x_3 = g(x_2)$

Noi presenteremo l'algoritmo di NEWTON e NEWTON-RAPHSON

Newton-Raphson

Si parte da  $x_0$ , e basato sulla formula precedente.

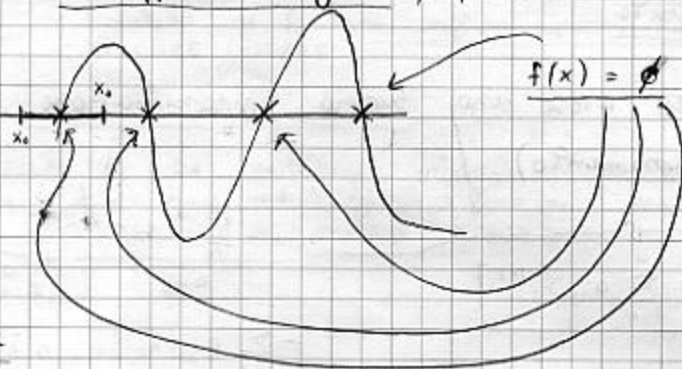
$g(x) = x - \frac{f(x)}{f'(x)}$   $f'(x)$  derivata

$x_{n+1} = g(x_n) = x_n - \frac{f(x_n)}{f'(x_n)}$   $\rightarrow$  FORMULA

1)  $f' \neq 0$

2) si basa sullo sviluppo di Taylor, quindi la scelta di  $x_0$  non deve essere casuale:

bisogna scegliere un intervallo molto piccolo in cui sono sicura che ricada una sola radice.

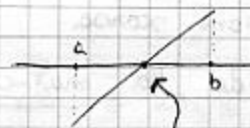


3)  $f'' > 0$  nell'intervallo dove sto lavorando (esclude casi in cui ci sono flessi - quindi non ci sono più radici, ma una -).

4)  $f(a)f(b) < 0$   $[a, b]$

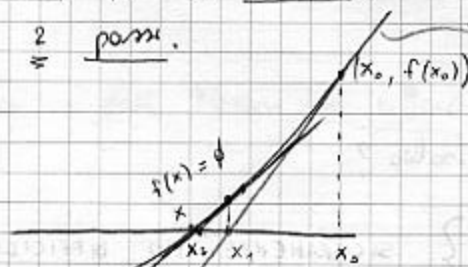
ipotesi a priori

così mi assicuro di avere almeno una radice.



almeno un punto taglia l'asse x

Dopo aver fissato questi punti applico Newton (che è molto rapido degli algoritmi lineari), si fa in 1 o 2 passi.



$x_{n+1}$  intersezione con asse x nella tangente alla curva  $f(x)$   $(x_n, f(x_n))$

seconda tangente, secondo giro (molto vicino alla

relazione x)

Metodo molto rapido, ma devo scegliere attentamente l'intervallo di partenza.  
Ma quando mi devo fermare con le iterazioni?

Test di convergenza

Sono due.

$$|x_m - x_{m+1}| \leq \epsilon \cdot \rho_{oll}$$

↓  
Tolleranza

$$|x_m - x_{m+1}| \leq \epsilon \cdot \rho_{oll} |x_{m+1}|$$

DEVO APPLICARE (IMPLEMENTARE) TUTTI E DUE 'STI TEST

23/10/05

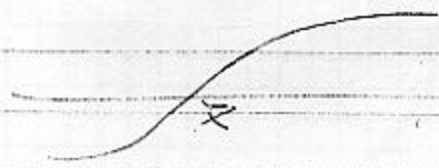
Eq non lineari

$$F(x) = 0$$

$$x'' + 7x + \cos x$$

Def.  $f: A \subseteq \mathbb{R} \rightarrow \mathbb{R}$

$x \in A$  è  
punto di accumulazione



$$(1) f(a) f(b) < 0$$

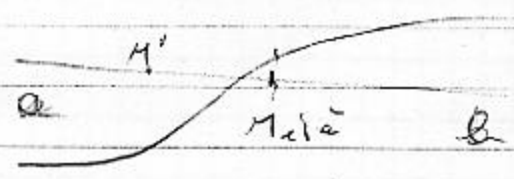
↑ IPOTESI  
↓ IPOTESI

$$f(\bar{x}) = 0$$

(2) f continua CONTINUA

C'è Metodo Bisezione per risolvere eq non lineari  
↓  
METODO BISEZIONE PER RISOLVERE EQ NON LINEARI

Divido ogni volta una parte es.



[a, M]

$$f(a) f(M) < 0$$

ok!

# TEST DI CONSERVAZIONE

$$(1) \quad |f(x_n)| < \epsilon$$

Successo  $x_1, x_2, x_3$   
SUCCESSIVO

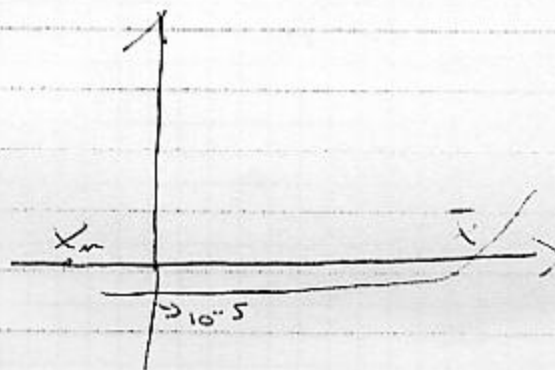
tale da

$$f(\bar{x}) = 0$$

↓  
f tende  
TENDE A 0

Vale il test (1)?

es)

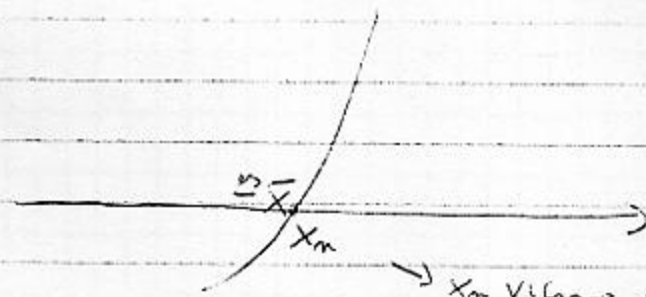


$$|f(x_n)| < \epsilon$$

→ NO!

$x_n$  ANCHE SE È PASSATO  
ma se è passato  
il valore di  $\bar{x}$   
È LONTANO DA  $\bar{x}$

Esempio



→  $x_n$  vicino a  $\bar{x}$

ma  $f(x_n)$

NO!

NON È PASSATO  
Non è passato  
è  
più infinito  
DURSI INFINITO

• Come viene una curva usata?

$$\underline{|x_m - x_{m+1}| \leq \epsilon_{toll}}$$

$$\underline{|x_m - x_{m+1}| \leq \epsilon_{toll} |x_{m+1}|}$$

Argomenti assistite

GIOVEDÌ  
CONTINUATO TRASFORMATA E MINIMI QUADRATI  
CONTINUATO Trasformata e Minimi Quadrati

• Macchine 2000/2005

• Doppio teorema

LEZIONE  
PROFESSORE

BR ALGEBRICH

$$\underline{P_N(x) = x^N + a_1 x^{N-1} + a_2 x^{N-2} \dots + a_N}$$

$$\underline{\pi = \max |a_n|}$$

$$\underline{C = \{ |z| \leq 1 + \pi \}}$$

N add

2N-2 moltiplic

$$\underline{x^2 = x \cdot x}$$

es

GRADO 3

GRADO 3

$$\underline{P_3(x) = x^3 + a_1 x^2 + a_2 x + a_3}$$

$$\underline{= a_3 + x(a_2 + x(a_1 + x))}$$

<sup>DERIVATO</sup> <sup>OGNI</sup> <sup>PARENTESI</sup> <sup>TORNA</sup> <sup>ABBIAHO</sup>  
 Derivato ogni funzione / Derivato allora  
<sup>VA</sup> <sup>TERMINALE</sup> <sup>NOTO</sup> <sup>E</sup> <sup>X</sup>  
 in termini noto e x

<sup>PER</sup> <sup>PUNTO</sup> <sup>X</sup> <sup>FISSATO</sup>  
Metodo Horner → per punto x fissato

$$\begin{cases} P_0 = 1 \\ P_n = x P_{n-1} + e_k, \quad k = 1 \dots n \end{cases}$$

<sup>ADDEDO</sup> <sup>CASO</sup>  
 • Nuovo caso

$$P_1 = x + e_1$$

$$P_2 = x(x + e_1) + e_2 = x^2 + e_1 x + e_2$$

$$P_3 = x(x^2 + e_1 x + e_2) + e_3$$

<sup>PERCHÉ</sup> <sup>HO</sup> <sup>SOMMA</sup>  
 • Perché no serve?

<sup>PER</sup> <sup>CALCOLARE</sup> <sup>LO</sup> <sup>DEIVATO</sup>  
 per calcolare lo derivato

<sup>NEWTON</sup> <sup>RICHIEDUTA</sup>  
<sup>CA</sup> <sup>DERIVATA</sup> <sup>PRIMA</sup>  
 Co derivata prima  
 Newton richiede

<sup>OGNI</sup> <sup>PUNTO</sup> <sup>X</sup> <sup>FISSATO</sup> <sup>DI</sup> <sup>CALCOLARE</sup>  
 ogni passaggio di calcolo

<sup>DERIVATA</sup>  
 Se voglio calcolare  $P_n'(x_n)$  per Newton

<sup>CONVENIENDO</sup> <sup>POLINOMIO</sup> <sup>AGGIUNTIVO</sup>  
 Costruisco polinomio aggiunto  $q_{n-1}(x) = P_n(x) - P_n(x_1)$   
 $x - x_1$

Con  $x \rightarrow x_1$   $q_{n-1}(x) \rightarrow P_n'(x)$

<sup>È</sup> <sup>INTAGGIOSO</sup> <sup>RISPETTO</sup> <sup>NEWTON</sup>  
 • È vantaggioso rispetto a Newton

Si dimostra che  $p_{n-1}(x) = x^{n-1} + b_1 x^{n-2} + \dots + b_{n-1}$

dalle  $b_0 = 1$

$$b_n = b_{n-1}x + a_n, \quad n = 1, \dots, n-1$$

• Supponiamo di aver trovato approssimazione grossolana per  $x_1$  come lo trovo per le altre?

Vediammo prossima volta.

## LAGRANGE

30/12/2005

$x_i^{(0)}$  approssimazione delle radici  $x_i$ .

Calcolo  $x_i$  con il metodo di Newton. Uso  $p_n(x_i)$  e  $p_n'(x_i)$ .

Costruisco il polinomio ridotto.

$q_{n-1}(x) = \frac{p_n(x) - p_n(x_1)}{x - x_1}$  ha per radici le stesse radici di  $p_n$  tranne  $x_1$ . Calcolo con  $x_2$ .

Per:

$q_{n-2}(x) = \frac{q_{n-1}(x) - q_{n-1}(x_2)}{x - x_2}$  e così via; fino ad ottenere un polinomio di grado 2 di cui

si sa approssimazione le radici. Le radici vanno calcolate in ordine crescente  $|x_1| \leq |x_2| \leq \dots$  (gli errori si accumulano nel corso della procedura).

## Interpolazione lagrangiana

Supponiamo di avere dei punti  $p, q$  nel piano. Vogliamo trovare la retta  $L_1$  (interpolante di Lagrange di grado 1) che passa per  $p$  e  $q$ . ( $L_1 = ax + b$ ).

Per due punti otteniamo un sistema di due equazioni in due incognite la cui soluzione sono i coefficienti  $a$  e  $b$ . Consideriamo un ulteriore punto  $T$ , la curva  $L_2$  che passa per  $p, q$  e  $T$  ha equazione:  $L_2(x) = a_0 + a_1x + a_2x^2$ .

Quindi, dati  $n$  punti del piano,  $L_n(x) = a_0 + a_1x + \dots + a_nx^n$  tale che  $L_n(x_i) = y_i$  (polinomio lagrangiano).

$(x_0, y_0) \dots (x_n, y_n)$

Si dimostra che il polinomio interpolante di Lagrange esiste ed è unico.

In generale, per  $n$  punti, questa procedura è la peggiore.

$$L_n(x_0) = y_0 = a_0 + a_1x_0 + \dots + a_nx_0^n$$

$$L_n(x_1) = y_1 = a_0 + a_1x_1 + \dots + a_nx_1^n$$

$$L_n(x_n) = y_n = a_0 + a_1x_n + \dots + a_nx_n^n$$

sistema di  $n+1$  equazioni in  $n+1$  incognite ( $a_0, a_n$ )

$$\begin{bmatrix} 1 & x_0 & x_0^2 & \dots & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^n \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^n \end{bmatrix} \neq \emptyset \iff \underline{\underline{\text{soluzione esiste ed \u00e9 unica}}}$$

Questa \u00e9 una matrice del tipo VAN DER MONDE, ovvero \u00e9 malcondizionata.

Esistono perci\u00f2 altre tecniche, pi\u00f9 stabili.

Interpolazione Lagrangiana con metodo delle differenze divise (o metodo di Newton) (FINITE) \u2192 per le parentesi quadre

$$L_n(x) = y_0 + (x-x_0) [x_0, x_1] + (x-x_0)(x-x_1) [x_0, x_1, x_2] + \dots + (x-x_0) \dots (x-x_{n-1}) \cdot [x_0, x_1, \dots, x_n]$$

se  $a \rightarrow b$   $[a, b]$  \u00e9 la derivata prima di  $f(x)$

$$[a, b] = \frac{f(b) - f(a)}{b - a} \quad f(x_i) = y_i \quad f'(x) \quad a \rightarrow b$$

Conosciamo  $f(x)$  solo nei punti  $(x_i, y_i)$  e vogliamo sapere come si comporta  $f$  nei punti non assegnati.

$$[x_0, x_1, x_2] = \frac{[x_1, x_2] - [x_0, x_1]}{x_2 - x_0} \quad f''(x)$$

Questa rappresentazione ha senso anche se i punti sono vicini fra loro.

$f(x) - L_n(x) = R_n(x)$  errore o resto dell'interpolazione Lagrangiana.

$$R_n(x) = (x-x_0) \dots (x-x_n) \cdot \underbrace{[x, x_0, x_1, \dots, x_n]}_{\text{DIFFERENZA DIVISA}}$$

Esiste un teorema che afferma  $[x, x_0, \dots, x_n] = \frac{f^{(n+1)}(\xi)}{n+1}$

$\xi$  giace all'interno dell'intervallo, ma non \u00e9 possibile conoscerne la posizione esatta.

Supponiamo  $f$  \u00e9  $P_n$  polinomio di grado  $n$ , es.  $f(x) = x^n \Rightarrow f^{(n+1)}(x) = \emptyset$

$\Rightarrow R_n(x) = \emptyset \Rightarrow f(x) = L_n(x) \rightarrow$  se resto \u00e9  $\emptyset$ , la funzione \u00e9 uguale al polinomio Lagrangiano

In generale, l'interpolante di Lagrange non \u00e9 un buono strumento per l'approssimazione di una funzione, perch\u00e9 aumentano i dati di input la funzione trovata non tende a "chiacciarsi" su quella reale.

$$L_{n,k} = \frac{(x-x_0) \dots (x-x_{k-1})(x-x_{k+1}) \dots (x-x_n)}{(x_k-x_0) \dots (x_k-x_{k-1})(x_k-x_{k+1}) \dots (x_k-x_n)}$$

$$d_n = \max_{[x_0, \dots, x_n]} \sum_{k=0}^n |L_{n,k}(x)| \quad \text{costanti di Lebesgue}$$

$$|R_n(x)| \leq (1 + d_n) E_n(f) \quad d_n E_n \xrightarrow{?} \emptyset$$

errore di migliore approssimazione uniforme

Dalle teoria sappiamo che  $d_n \geq c \log n \iff d_n$  non \u00e9 limitata (cresce al crescere di  $n$ )

Domanda all'esame



Nel caso migliore,  $d_n \approx c \log n$   $x_i$  oeu di Cheb

altrimenti  $d_n \approx e^n$

$x_i$  equidistanti

PROGETTO  $\rightarrow I_{i,j} \rightarrow I_{i,j}^e = I_{i,j} + \epsilon_{i,j}$

generazione di numeri casuali.

## Interpolazione di Lagrange

7/12/2005

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

$$[a, b, c] = \frac{[b, c] - [a, b]}{c - a}$$

$$[a, b, c, d] = \frac{[b, c, d] - [a, b, c]}{d - a}$$

$$[a, a] = f'(a)$$

$$[a, a, a] = \frac{f''(a)}{2}$$

$$l_{m,k}(x) = \frac{(x-x_0) \dots (x-x_{k-1})(x-x_{k+1}) \dots (x-x_m)}{(x_k-x_0) \dots (x_k-x_{k-1})(x_k-x_{k+1}) \dots (x_k-x_m)}$$
 polinomio fondamentale di Lagrange

$$L_m(f, x) = \sum_{k=0}^m l_{m,k}(x) y_k$$
 utile se molti  $y_k$  sono  $= \phi$   
vantaggio: se  $x_i \approx x_{i+1}$

$$d_n = \max_{[x_0, x_n]} \sum_{k=0}^n |l_{n,k}(x)| \approx c \log n$$

$$|f(x) - L_m(f, x)| \leq c(1+d_n) E_n(f)$$
 errore di migliore approssimazione uniforme

il nostro scopo è che questa quantità vada a zero

$$\text{se } d_n E_n(f) = \phi \Rightarrow f(x) - L_m(f, x) \rightarrow \phi$$
$$d_n E_n(f) \rightarrow \phi$$

$d_n$  su nodi equidistanti  $d_n \approx e^n$

$E_n(f) \rightarrow \phi$  regolarità di  $f$

$E_n(f) = \min_{P \in \mathcal{P}_n[x_0, x_n]} \max |f(x) - P(x)| \rightarrow$  errore di migliore approssimazione uniforme

Una domanda all' esame può essere: quali sono vantaggi e svantaggi dei due approcci di Lagrange?

Esempio diretto all' esame

$$f(x) = x^{3/2} [0, 1]$$

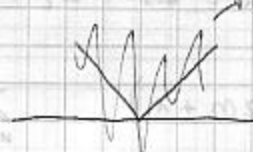
$$f \in C^1$$

$$E_n(f) \sim \frac{1}{m}$$

$$d_n E_n(f) \sim$$

$$\log \frac{1}{m} \rightarrow \phi$$

interpolante di Lagrange  
già con 5 punti c'è un oscillazione di overshoot e undershoot



Commento sulla convergenza (fenomeno di Runge)

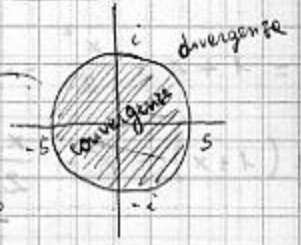
$$f(x) = \frac{1}{1+x^2}$$

RUNGE

$$L_n(f) \not\rightarrow f$$

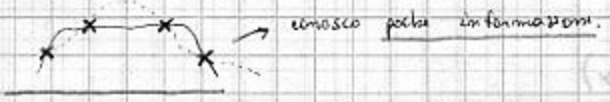
non avere tutto le derivate che voglio

Se estendiamo lo studio nel campo dei complessi  $\rightarrow i^2 = -1$   
avere divisione per zero (ha delle singolarità)



dentro questo cerchio non ci sono punti di singolarità

Supponiamo che  $f(x)$  continua in  $[a, b]$   
appross.  $f$  con  $g$

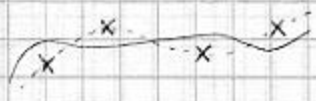


Come faccio a sapere di quanto mi discosto?

Approssimazione uniforme (o errore, norma)

$$\|f - g\|_\infty = \max_{[a, b]} |f(x) - g(x)|$$

Ho  $f$  cont. :  $\exists P_n$  q.c.  $\|f - P_n\| \rightarrow 0$ ?



WEIERSTRASS dice che si possono trovare polinomi che al crescere di  $n$  si avvicinano al vero grafico della funzione.

Un esempio esplicito:  $\exists! E_n(f) = \min_{P \in \mathcal{P}_n} \|f - P\| \rightarrow$  polinomio di migliore approssimazione uniforme, ma non sappiamo costruirlo, lo sappiamo solo in teoria.

Errore quadratico medio

$$\int_a^b (f(x) - g(x))^2 dx = \text{errore quadratico medio}$$

$w(x) \rightarrow$  si può aggiungere questa funzione peso

$$f_n(x) = \sum_{k=0}^n c_k \phi_k(x) \rightarrow \text{migliore } g \text{ che possiamo scegliere}$$

$\phi_n$  cont. quadratico integrabile, ortogonale

$$\int_a^b w(x) (f(x) - f_n(x))^2 dx \quad \text{minimo}$$

$$c_k = a_k = \int_a^b w(x) f(x) \phi_k(x) dx$$

eff. di Fourier

Risultato più che soddisfacente: metodo dei medi quadratici con qualche riferimento al metodo uniforme.

Definizione di ortogonalità:  $\int_a^b w(x) \phi_r(x) \phi_s(x) dx = \begin{cases} 0, & r \neq s \\ 1, & r = s \end{cases}$

Approssimazione in media quadratica (non c'è su Homegate)

$$\int_a^b w(x) [f(x) - f_n(x)]^2 dx \rightarrow \text{trovare } f_n \text{ tale che questa quantità sia minima}$$

14/12/2005

$$w(x) = \frac{1}{\sqrt{1-x^2}}, [a, b] = [-1, 1]$$

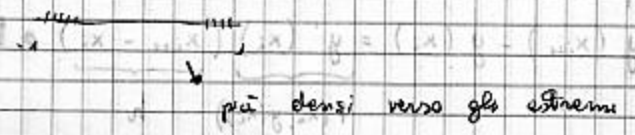
$T_n$  pol. di Chebicev di prima specie

$T_0 = 1, T_1(x) = x, T_{m+1}(x) = 2x \cdot T_m(x) - T_{m-1}(x) \rightarrow$  relazione di ricorrenza a 3 termini  
 zeri di  $T_m$  ( $T_m = 0$ )

$x_k = \cos \frac{\pi}{2m} (2k+1) \quad k=0, \dots, m-1$  reali distinti e in  $(-1, 1)$

$$\int_{-1}^1 \frac{T_m(x) T_n(x)}{\sqrt{1-x^2}} dx = \begin{cases} 0, & m \neq n \\ \frac{\pi}{2}, & m = n > 0 \\ \pi, & m = n = 0 \end{cases}$$

$T_m$  ortogonale



$\left\{ \frac{T_0}{\sqrt{\pi}}, T_n \frac{\sqrt{2}}{\sqrt{\pi}} \right\}_m$  sistema ortogonale in  $(-1, 1)$   $\rightarrow$   $\int_{-1}^1 \frac{1}{\sqrt{1-x^2}} \tilde{T}_m(x) \tilde{T}_n(x) dx = \begin{cases} 0, & m \neq n \\ \pm, & m = n \end{cases}$

$f_m(x) = c_0 \frac{T_0}{\sqrt{\pi}} + \sum_{k=1}^m c_k T_k(x) \frac{\sqrt{2}}{\sqrt{\pi}} \rightarrow$  realizza la minima approssimazione in media quadratica se:

$$c_0 = a_0 = \frac{1}{\sqrt{\pi}} \int_{-1}^1 f(x) T_0 dx$$

$$c_k = a_k = \frac{\sqrt{2}}{\sqrt{\pi}} \int_{-1}^1 f(x) T_k(x) dx \rightarrow$$

polinomio di quasi migliore approssimazione

### Equazioni differenziali:

usate in fisica, ingegneria, economia, chimica. Studiano fenomeni evolutivi, per esempio come  $x$  propaga un suono, come  $x$  propaga il calore in una stanza. Sempre difficili da risolvere con l'analisi.

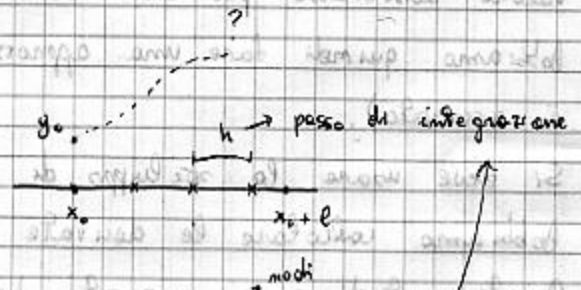
esempio

$$\begin{cases} y'(x) = y(x) \\ y(0) = 1 \end{cases} \quad y(x)? \text{ in } [0, l] ? \rightarrow y(x) = e^x \text{ (caso banale)}$$

Cos'è una eq. diff. del 1° ordine?

$$\begin{cases} y'(x) = f(x, y(x)) \\ y(x_0) = y_0 \end{cases} \quad [x_0, x_0 + l]$$

condizione iniziale

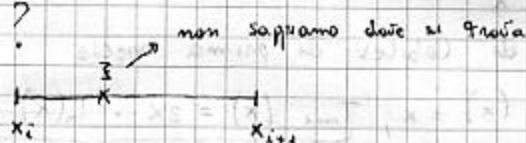


Vogliamo trovare la soluzione in generici punti equidistanti.  $\rightarrow x_i = x_0 + i h$

dove  $h = \frac{l}{m}$

## Metodo di Eulero

$$y(x_{i+1}) - y(x_i) = y'(\xi) (x_{i+1} - x_i), \quad \xi ?$$



faccio un' approssimazione grossolana  $\rightarrow y'(\xi) \approx y'(x_i)$

quindi:

$$y(x_{i+1}) - y(x_i) = \underbrace{y'(x_i)}_{\substack{f(x_i, y(x_i)) \\ \text{per la relazione} \\ \text{non invariabile}}} (x_{i+1} - x_i) = h$$

$y(x_{i+1}) = y(x_i) + h f(x_i, y(x_i)) \rightarrow$  scrivendo  $y(x_i) = y_i$ , abbiamo  $\rightarrow$

$$\rightarrow y_{i+1} = y_i + h f(x_i, y_i)$$

Se poniamo  $i=0 \rightarrow y_1 = y_0 + h f(x_0, y_0) \rightarrow$  di questa conosciamo tutto:  $y_0$  (cond. iniziale),  $h$  (fissata da noi),  $f$  (data).

Ma questa formula è affetta da un grosso errore dato da:  $y'(\xi) \approx y'(x_i)$

Calcoliamola ora per  $i=1$ :

$y_2 = y_1 + h f(x_1, y_1) \rightarrow$  questa è due volte affetta da errore, perché uso qualcosa che era già affetto da errore ( $y_1$ ).

Sempre peggio per  $i=2, i=3, i=4, \dots$

Quindi questo metodo è molto semplice, ma ha la propagazione dell' errore (accumulazione dell' errore).

Ci servono metodi più accurati.

## Metodo di Eulero - Cauchy

Basato su questa approssimazione grossolana  $\rightarrow y'(\xi) \approx \frac{y_{i+1} - y_{i-1}}{2h} = y'(x_i)$

↓  
più accurata della precedente

$$y_{i+1} = y_{i-1} + 2h y'(x_i) = y_{i-1} + 2h f(x_i, y_i)$$

Calcoliamola per  $i=1$ :

$$y_2 = y_0 + 2h f(x_1, y_1) \rightarrow y_0 \text{ lo conosciamo (cond. iniziale), ma } y_1 ?$$

Quindi questo metodo di Eulero - Cauchy non è auto partente (ha bisogno di un valore ausiliario per innescare il processo).

Abbiamo quindi fare una approssimazione per  $y_1$  (non col metodo di Gauss, che non è accurato).

Si deve usare lo sviluppo di Taylor ricordando:  $y'(x) = f(x, y(x))$ .

Abbiamo calcolare le derivate parziali  $\rightarrow f_x, f_y$ .

Questi metodi sono superati, hanno + di 100 anni.

Metodi passo passo (step by step)

Classe di metodi più utile per risolvere le eq. di ff:

$$y_{i+k} = \sum_{j=0}^{k-1} \{ a_j y_{i+j} + h b_j f(x_{i+j}, y_{i+j}) \} + h b_k f(x_{i+k}, y_{i+k})$$

$b_k = \emptyset$  esplicita (predictor)  $\rightarrow$  dai valori precedenti riesco a prevedere il valore corrente

$b_k \neq \emptyset$  implicita (corrector)

Due famiglie di metodi:

Metodo di Milne (corrector)

$$y_{i+1} = y_{i-3} + \frac{4}{3} h (2y'_{i-2} - y'_{i-1} + 2y'_i)$$

$$y_{i+1} = y_{i-1} + \frac{h}{3} (y'_{i-1} + 4y'_i + y'_{i+1})$$

$$y(x)_{i+1} = f(x_{i+1}, y(x)_{i+1})$$

Non è autoportante.

Metodo di Runge-Kutta (esplicito)

$$y_{i+1} = y_i + \frac{1}{2} (k_1 + k_2) \leftarrow \text{formula}$$

$$k_1 = h f(x_i, y_i)$$

$$k_2 = h f(x_i + h, y_i + k_1)$$

$\swarrow$  prima volta che dentro la  $f$  non comparano solo  $x_i, y_i$

È autoportante.

Matematica computazionale in generale

Non dai risultati specifici ma modelli di calcolo.

la soluzione trovata sarà sempre approssimata, vincoli: accuratezza, velocità di calcolo, marginale di errore.

Esempio

$$y(f) = \sin \log f$$

equazione esplicita, dove l'errore può essere nullo.

$$y = e^{f^2} \sin y f$$

equazione implicita ( $f$  e  $y$  non sono separate). da risolvere ma in modo approssimato ( $f$  nota non esplicitamente)

$$y(f) = \sum_{k=0}^{\infty} a_k e^{-b_k f^k}$$

anche se la  $f$  è nota esplicitamente, non possiamo calcolarla (a causa della somma degli  $\infty$  termini)

$$y(f) = \sin f$$

$f$  nota esplicitamente e calcolabile non direttamente (ma anche in questo caso c'è un errore)

$$y(f) = 3f$$

$f$  nota esplicitamente calcolabile direttamente (anche in questo caso abbiamo un errore, perché non sappiamo come è  $f$  e come  $\pi$  può rappresentare con un numero finito di bit)

Quindi avremo sempre e comunque un marginale di errore.

Il computer calcola il log con le serie.

Se dovessimo calcolare  $\sin \log \sqrt{f^{2+x}}$  bisognerebbe procedere dall'interno verso l'esterno, invece il computer ne crea un'espressione approssimata tipo  $\sum a_n x^n$  ma estremamente più veloce da calcolare. Bisogna quindi sempre conoscere il marginale di errore.

Esempio

Quanto occupa una foto da  $1024 \times 768$  a colori?

$$3 \cdot 1024 \cdot 768 \cdot 1 = 2,4 \text{ Mb}$$

da 0 a 255

$$256 = 2^8$$

RGB  $\downarrow$  valori da 255  $\downarrow$  perché abbiamo supposto che un pixel occupa un byte di memoria

Se devo spedire la foto con un modem a  $40 \text{ kb/s}$ ?

$$\frac{(2,4 \cdot 8) \cdot 1024}{40} = 500 \text{ sec} \approx 8 \text{ min}$$

ha convertito i byte in bit

possiamo convertirla in jpeg che si basa sulla tecnica lossy, basata sulla trasformata di Fourier.

Matematica computazionale in grafica

Risoluzione in megapixel: quanti milioni di pixel.

Metodi per comprimere l'immagine  $\rightarrow$  trattiamo l'immagine

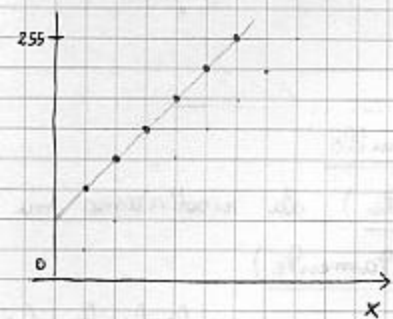
Non si perdono nel numero di pixel, ma nella qualità dell'immagine.

Per esempio si vedono dei blocchi nell'immagine.

# Trasformata di Fourier

Algoritmo per comprimere basato sulla Trasformata di Fourier.

Ammettiamo di aver immagine e di vedere solo una riga di shell. Esso corrisponde a vedere una riga di pixel (formato da 3 matrici R, G, B).



Un'immagine la cui componente rossa da sinistra a destra tende ad aumentare.

PIXEL	R
1	23
2	24
3	25
4	26
⋮	⋮

È sprecato usare memoria per degli tutti i pixel, quando in realtà è una retta.

$$R = a + bx$$

$$R = 22 + x$$

$$\begin{cases} a = 22 \\ b = 1 \end{cases}$$

lo possiamo fare con l'equazione, fornendo termini no  $\theta_0$  e coeff. angolare.

Numero di parametri da dare molto più basso.

Se noi zomiamo continuamente vediamo i pixel che formano il pezzo di immagine, e non si può più andare oltre. Sappiamo le informazioni dei pixel, ma non tra pixel e pixel. Se usiamo le equazioni possiamo stimarlo, l'equazione è continua e può stimare anche tra pixel e pixel (rappresentazione retinale invece che discretizzata → stessa cosa che accade ai fonti true type).

Il problema è che le equazioni non sono così facili come quelle di una retta.

Per esempio:  $R = a + bx + cx^2$  (parabola) → siamo andati a complicare, ma non è sufficiente.

Possiamo quindi incrementare l'ordine:  $R = a + bx + cx^2 + dx^3 + ex^4 + \dots$

da compressione funziona facendo la rappresentazione tramite una funzione invece che pixel per pixel, cercando di usare il numero più basso di coefficienti.

Si possono usare per esempio degli esponenziali:  $a + be^x + ce^{2x} + de^{3x} + ee^{4x} + \dots$

Non è questo il problema, ma quanti coefficienti uso, devo usarne il meno poss. bile.

Devo trovare una rappresentazione sparsa (con il più basso numero di coefficienti).

$$f(x) = a_0 + a_1 \varphi_1(x) + a_2 \varphi_2(x) + a_3 \varphi_3(x) + \dots \rightarrow \text{combinazione lineare}$$

Come si scelgono  $\varphi(i)$ : primo criterio → sparsità

Ma noi dobbiamo sapere subito quanto valgono  $a_0, a_1, \dots$  i coeff.. Facciamo quindi il processo di ANALISI: noti  $f(x_1), f(x_2), \dots \rightarrow a_0, a_1, a_2, \dots$  dobbiamo trovare subito i coeff. (in tempo reale) → neanche 2 sec per avere il jpeg, in una videocamera digitale devo fare 24 foto al secondo.

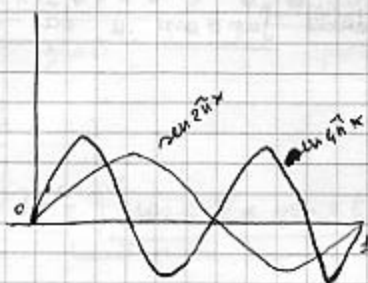
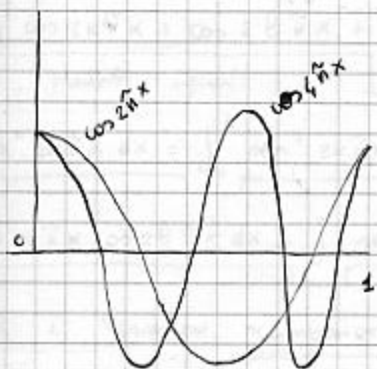
Per poi visualizzare l'immagine dobbiamo fare la  sintesi:  $a_0, a_1, a_2, \dots \rightarrow f(x_1), f(x_2), \dots$

ma non è necessario che accada velocissima come l'analisi, ma deve essere comunque veloce.

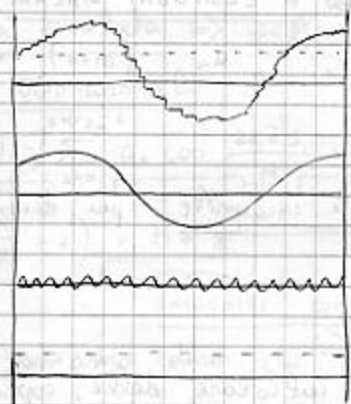
la compressione deve dare un'immagine equivalente con una funzione  $\rightarrow$  spartiti e veloci  
 tra i due algoritmi (ANALISI e SINTESI).

Sistema  $\varphi$

$$\varphi \equiv \{ \cos 2\pi x, \cos 4\pi x, \cos 6\pi x, \cos 8\pi x, \sin 2\pi x, \sin 4\pi x, \sin 6\pi x \dots \} \quad x \in [0, 1]$$



Q che serie rappresentarle così?



$\rightarrow$  se dovessi rappresentare queste dovrei usare pixel per ogni minima perturbazione ma mi accorgo che essa è l'unione di due funzioni + semplici

0,5

$\rightarrow \cos 2\pi x$  (1)

$\rightarrow$  tanti  $\sin 4\pi x$  (0,1)

$\rightarrow$  alta al livello del cos (0,5)

$$f(x) = a_0 + a_1 \cos 2\pi x + b_{10} \sin 2\pi x$$

$$a_0 = 0,5$$

$$a_1 = 1$$

$$b_{10} = 0,1$$

d'equalizzatore dello stereo rubisce lo stesso principio, onde alta e bassa frequenza, carap. prossima dei segnali ad alta e bassa frequenza (ma non usa la trasformata di Fourier).

Sistema di Fourier

$$f(x) = a_0 + \sum_{k=1}^{\infty} a_k \cos 2k\pi x + \sum_{k=1}^{\infty} b_k \sin 2k\pi x, \quad x \in [0, 1]$$

$$f(0) = a_0 + \sum_{k=1}^{\infty} a_k \cos 2k\pi \cdot 0 + \sum_{k=1}^{\infty} b_k \sin 2k\pi \cdot 0$$

$$= a_0 + \sum_{k=1}^{\infty} a_k$$

$$f(1) = a_0 + \sum_{k=1}^{\infty} a_k$$

$$f(x+1) = a_0 + \sum_{k=1}^{\infty} a_k \cos 2k\pi (x+1) + \sum_{k=1}^{\infty} b_k \sin 2k\pi (x+1)$$

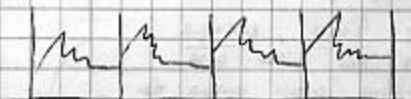
$$= a_0 + \sum_{k=1}^{\infty} a_k \cos (2k\pi x + 2k\pi) + \sum_{k=1}^{\infty} b_k \sin (2k\pi x + 2k\pi)$$

$$= a_0 + \sum_{k=1}^{\infty} a_k \cos 2k\pi x + \sum_{k=1}^{\infty} b_k \sin 2k\pi x$$

UGUALE ALL' INIZIO

$\rightarrow$  AGGIUNGO UN GIRO CHE NON SERVE

$$f(x+1) = f(x) \quad \text{PERIODICA}$$





Per comprimere sostituiamo un sistema discreto (pixel) o un sistema funzionale (con pochi coeff.), ANALISI e SINTESI e viceversa molto rapidi. Fourier è in grado di soddisfare questi requisiti.

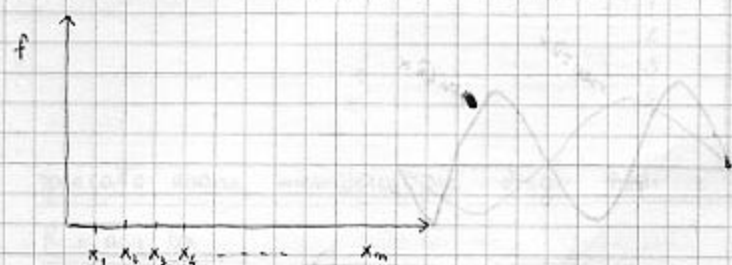
26/10/2005

Funzione di compressione

$$f(x) = a_0 \varphi_0(x) + a_1 \varphi_1(x) + \dots + a_m \varphi_m(x) + \dots$$

ANALISI → trovare i coefficienti (non immediato)  $a_0, a_1, \dots, a_m$

SINTESI → trovare  $f(x)$  dalle funzioni e dai coefficienti



$$f(x_0) = a_0 \varphi_0(x_0) + a_1 \varphi_1(x_0) + \dots + a_m \varphi_m(x_0)$$

$$f(x_1) = a_0 \varphi_0(x_1) + a_1 \varphi_1(x_1) + \dots + a_m \varphi_m(x_1)$$

$$f(x_m) = a_0 \varphi_0(x_m) + a_1 \varphi_1(x_m) + \dots + a_m \varphi_m(x_m)$$

SISTEMA di EQUAZIONI LINEARI

con  $a_0, a_1, \dots, a_m$   $(m+1)$  incognite

Si risolve con Gauss in tempo  $O(m^3)$ . Più aumentano le incognite, più diventa complicato ed è difficile avere risultato in tempo reale.

Abbiamo scegliere  $\varphi_m$  che permette di abbattere questo ostacolo.

Fourier permette di passare da  $O(m^3) \rightarrow O(m^2)$ , arriviamo a calcolare delle espressioni, però ancora non è sufficiente, per esempio per le foto.

Allora esiste la trasformata veloce di Fourier ( $O(m \log m)$  di complessità).

Ma ancora non è sufficiente per le operazioni in tempo reale.

La trasformata di Fourier viene implementata a livello hardware (in chip) alla velocità della luce. Prima era fatto via software ed era lento. Dal pentium 8486 fu fatto nell'hardware. Una tecnologia nuova (10 anni) permette di farlo in  $O(N)$ , (ware letto), → algoritmo jpeg2000.

Fourier

$$f(x) = a_0 + a_1 \cos 2\pi x + a_2 \cos 4\pi x + \dots + a_k \cos 2k\pi x + \dots + b_1 \sin 2\pi x + b_2 \sin 4\pi x + \dots + b_k \sin 2k\pi x + \dots$$

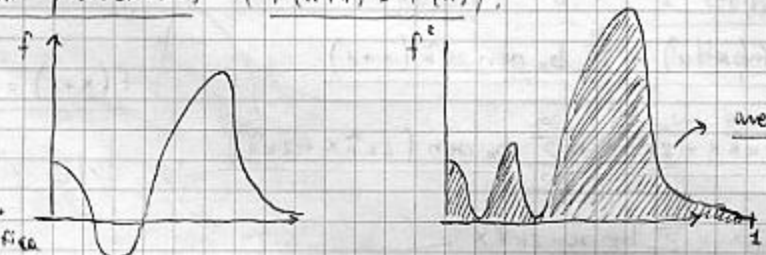
$x \in [0, 1]$  → polinomio trigonometrico.

Rappresentiamo le frequenze oscillatrici delle funzioni, più sono basse meglio è.

1) Si applica a funzioni periodiche, ( $f(x+1) = f(x)$ ).

$$2) \int_0^1 f^2(x) dx$$

in fisica, energia che può avere una funzione, oppure energia dei pixel nella macchina fotografica



area ha un valore finito

$$\int_0^1 f^2(x) dx = \sum_{k \geq 0} a_k^2 + b_k^2$$

limita

limita, quindi  $a_k$  e  $b_k$  devono essere finite, la luce della macchina fotografica non ci obbliga

$a_k$  e  $b_k$  devono essere  $\neq 0$  o tendere a 0, tali che li posso trascurare nella somma (ecco la complessione).

$$\int_0^1 f^2(x) dx = \int_0^1 \left[ \sum_k a_k \cos 2k\tilde{\omega}x + b_k \sin 2k\tilde{\omega}x \right]^2 dx \approx \int_0^1 \left[ \sum_k \left( a_k^2 \cos^2 2k\tilde{\omega}x + b_k^2 \sin^2 2k\tilde{\omega}x \right) + \sum_{k, \ell} \dots \right] dx$$

$$a_k a_\ell \cos 2k\tilde{\omega}x + \cos 2\ell\tilde{\omega}x + \sum_{\substack{k, \ell \\ k \neq \ell}} b_k b_\ell \sin 2k\tilde{\omega}x \sin 2\ell\tilde{\omega}x + \sum_{k, \ell} a_k b_\ell \cos 2k\tilde{\omega}x \sin 2\ell\tilde{\omega}x$$

prodotti uguali      prodotti uguali      prodotti misti

$$\int_0^1 \cos^2 2k\tilde{\omega}x dx = \int_0^1 \sin^2 2k\tilde{\omega}x dx = 1 \rightarrow \text{REGOLA}$$

$$\int_0^1 \cos 2k\tilde{\omega}x \cos 2\ell\tilde{\omega}x dx = \int_0^1 \sin 2k\tilde{\omega}x \sin 2\ell\tilde{\omega}x dx = \int_0^1 \cos 2k\tilde{\omega}x \sin 2\ell\tilde{\omega}x dx = 0 \rightarrow \text{REGOLA}$$

quindi i termini scompaiono tutti (vale per ogni base ortogonale - ortonormale) e abbiamo:

$$\int_0^1 f^2(x) dx = \sum_{k \geq 0} a_k^2 + b_k^2$$

Energia limitata  $\Rightarrow$  weff.  $\rightarrow 0$

3) f. ortonormali

$$f\left(\frac{j}{N}\right) = \sum_{k=0}^{N/2-1} a_k \cos \frac{2k\tilde{\omega}j}{N} + \sum_{k=0}^{N/2-1} b_k \sin \frac{2k\tilde{\omega}j}{N}$$

$$j = 0, \dots, N-1, \quad f(x_j) = \frac{j}{N}$$

il valore con  $\frac{j}{N}$  è uguale a quello con  $0$  però non si mette

$$2 \left( \frac{N}{2} + 1 \right) = N + 2$$

INCOGNITE CON N EQUAZIONI (NON SI PUÒ FARE) PIÙ INCOGNITE RISPETTO ALLE EQ.

Però possiamo dire che  $b_0 = 0$ , ed anche che  $b_{N/2} = 0$  quindi abbiamo tolto due incognite e possiamo risolvere il sistema.

$$\begin{cases} a_k = \frac{1}{N} \sum_j f_j \cos \frac{2k\tilde{\omega}j}{N} \\ b_k = \frac{1}{N} \sum_j f_j \sin \frac{2k\tilde{\omega}j}{N} \end{cases}$$

FOURIER IN FORMA ANALITICA

Quante operazioni sono necessarie per trovare  $a_k, b_k$ ?

$\frac{N-1}{2}$   
SOMMATORIE

$\frac{N}{2}$   
MOLT. DENTRO  
SOMMATORIA

$\frac{1}{2}$   
PRODOTTO  
FUORI

per un coefficiente  $\rightarrow$  in totale  $2N$ .

$$\frac{4N \cdot \frac{N}{2}}{2} = \frac{4N^2}{2} = O(2N^2) \rightarrow O(N^2)$$

WEFF.  
SOMMATORIA PRINCIPALE

$$\begin{cases} a_k = \frac{1}{N} \sum_{s=0}^{N/2} f_s \cos \frac{2k\tilde{\omega}s}{N} \\ b_k = \frac{1}{N} \sum_{s=0}^{N/2} f_s \sin \frac{2k\tilde{\omega}s}{N} \end{cases}$$

TRASFORMATA VELOCE DI FOURIER

Bisogna usare i complessi:

$$C_k = \frac{1}{N} \sum_{s=0}^{N-1} f_s e^{-i \frac{2\tilde{\omega}js}{N}}$$

dove  $e^{ix} = \cos x + i \sin x$

$$(a+ib) \pm (c+id) = (a \pm c) + (b \pm d)i \rightarrow \text{REGOLA}$$

$$(a+ib) \cdot (c+id) = ac + iad + ibc - bd \rightarrow \text{REGOLA}$$

$$= (ac - bd) + (ad + bc)i$$

$$i^2 = -1$$

$$e^{-ix} = \cos(-x) + i \sin(-x) = \cos x - i \sin x \rightarrow \text{REGOLA} \leftarrow$$

$$\omega = e^{-i \frac{2\pi}{N}} \rightarrow \text{SE VALE QUESTA} \rightarrow C_k = \frac{1}{N} \sum_{j=0}^{N-1} f_j \omega^{jk}$$

$$C_k = \frac{1}{N} \sum_{j=0}^{N-1} f_j e^{-i \frac{2\pi jk}{N}} \leftarrow \begin{matrix} \text{PERCHÉ È} \\ \text{COSÌ} \end{matrix}$$

$$C_k^{(N)} = \frac{1}{N} \sum_{j=0}^{N-1} f_j \omega^{jk} \rightarrow \text{dividiamolo in due parti (numero di punti è una potenza di 2)}$$

$$= \frac{1}{N} \sum_{j=0}^{N/2-1} [f_j \omega^{jk} + f_{j+\frac{N}{2}} \omega^{(j+\frac{N}{2})k}] =$$

$$= \frac{1}{N} \sum_{j=0}^{N/2-1} [f_j \omega^{jk} + f_{j+\frac{N}{2}} \omega^{jk} \cdot \omega^{\frac{N}{2}k}] =$$

$$\omega^{\frac{N}{2}k} = e^{-i \frac{2\pi}{N} \cdot \frac{N}{2} \cdot k} = e^{-i \pi k} = \cos(\pi k) - i \sin(\pi k) \Rightarrow \begin{cases} k \text{ pari} \rightarrow 1 \\ k \text{ dispari} \rightarrow -1 \end{cases} \rightarrow (-1)^k$$

$$= \frac{1}{N} \sum_{j=0}^{N/2-1} [f_j + f_{j+\frac{N}{2}} (-1)^k] \omega^{jk}$$

quindi cambia a seconda che  $k$  sia pari o di pari

$$C_{2m}^{(N)} = \frac{1}{N} \sum_{j=0}^{N/2-1} [f_j + f_{j+\frac{N}{2}}] \omega^{2m_j}$$

pari

$$C_{2m+1}^{(N)} = \frac{1}{N} \sum_{j=0}^{N/2-1} [f_j - f_{j+\frac{N}{2}}] \omega^j \omega^{2m_j}$$

dispari

$$C_{2m}^{(N)} = \frac{1}{2} \left[ \frac{1}{N/2} \sum_{j=0}^{N/2-1} (f_j + f_{j+\frac{N}{2}}) \omega_1^{m_j} \right]$$

$$C_{2m+1}^{(N)} = \frac{1}{2} \left[ \frac{1}{N/2} \sum_{j=0}^{N/2-1} (f_j - f_{j+\frac{N}{2}}) \omega^j \omega_1^{m_j} \right]$$

3/11/2005

Fourier

$$C_k = \frac{1}{N} \sum_{j=0}^{N-1} f_j \omega^{jk}$$

$k = 0, \dots, N-1$

$$\omega := e^{-i \frac{2\pi}{N}}$$

$$C_{2m} = \frac{1}{2} \left[ \frac{1}{N/2} \sum_{j=0}^{N/2-1} f_j^{(1)} \omega_1^{m_j} \right]$$

$$C_{2m+1} = \frac{1}{2} \left[ \frac{1}{N/2} \sum_{j=0}^{N/2-1} f_{j+\frac{N}{2}}^{(1)} \omega_1^{m_j} \right]$$

$$f_j^{(1)} = f_j^{(0)} + f_{j+\frac{N}{2}}^{(0)}$$

$$f_{j+\frac{N}{2}}^{(1)} = (f_j^{(0)} - f_{j+\frac{N}{2}}^{(0)}) \omega^j$$

$$j = 0 \dots \frac{N}{2} - 1$$

$$m_j = 0 \dots \frac{N}{2} - 1$$

$$\omega_1 = \omega^2 = e^{-i \frac{2\pi}{N}} = e^{-i \frac{2\pi}{N/2}}$$

Tutti quei calcoli danno  $2N^2$  complessità

Così non abbiamo dimostrato nulla.

Alli riduco da  $n$  calcoli a  $\frac{N}{2}$  calcoli pari e  $\frac{N}{2}$  calcoli dispari, ho guadagnato pochissimo  $\frac{1}{2} N^2$  contro  $2N^2$  di quella intera.

Le due somme si possono ancora dividere  $\left\langle \begin{matrix} N/4 - 1 \\ N/4 - 1 \end{matrix} \right\rangle \left\langle \begin{matrix} N/4 - 1 \\ N/4 - 1 \end{matrix} \right\rangle$

Queste si possono ancora dividere in due somme di  $N/8 - 1$  termini, e così via...

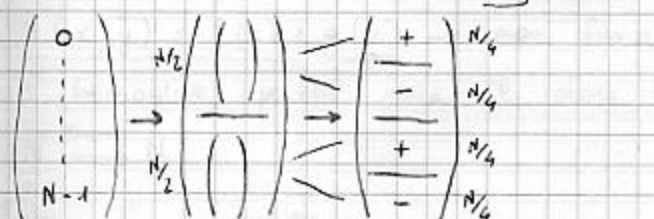
In ogni passaggio scivo i coefficienti in un modo diverso.

Quante operazioni ci vogliono per risolvere la trasformata?

Ma dobbiamo ridurre alla somma di un solo termine, che significa il termine

seno.

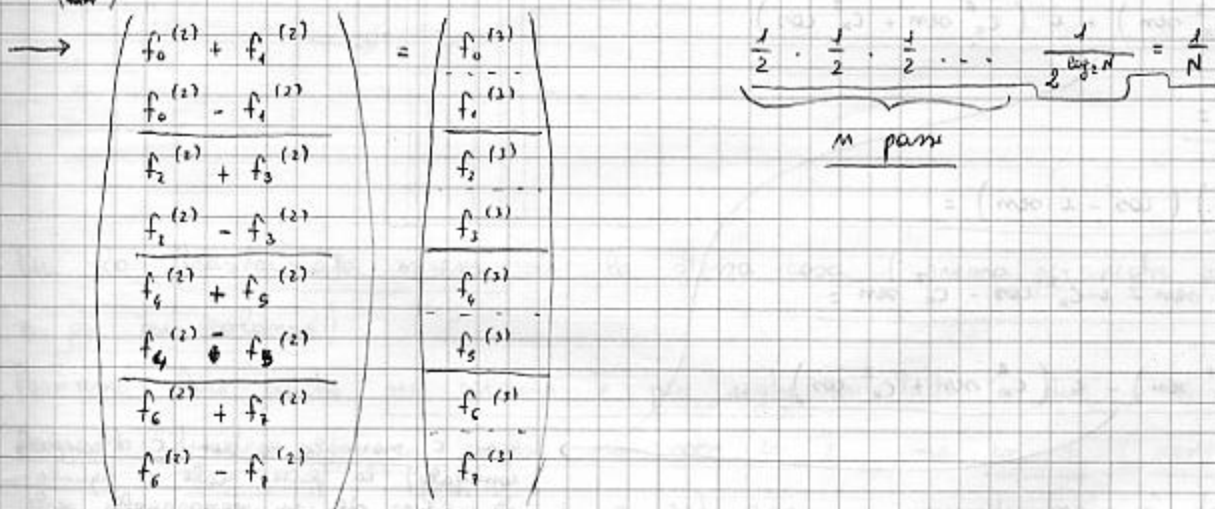
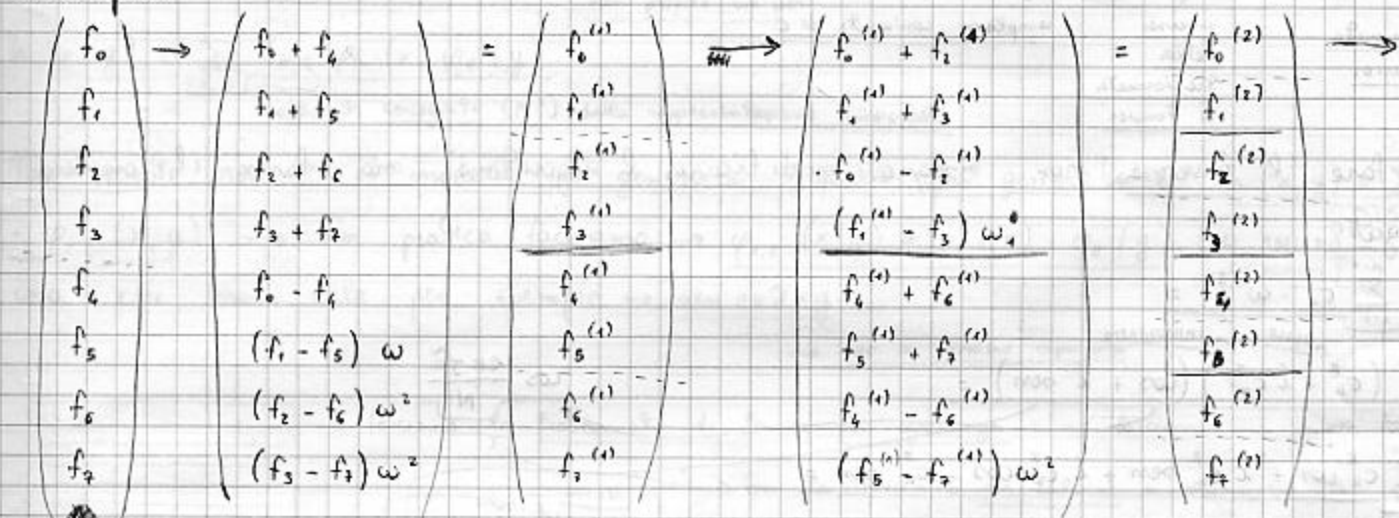
Al primo passaggio faccio  $\frac{3}{2} N$  operazioni.



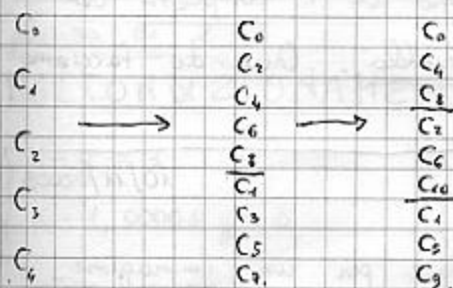
Anche nel passaggio dopo dovremo fare  $\frac{3}{2} N$  operazioni, consideriamo però l'altezza dell'albero che andiamo a creare,  $\log_2 N \rightarrow \frac{3}{2} N \cdot \log_2 N \rightarrow O(N \log N)$ .

Sembra che l'algoritmo sia nato nel 1965 da Cooley e Tukey, la trasformata era già usata ma con grossi dati c'erano problemi, già ai tempi di Fourier questo algoritmo fu sviluppato.

Esempio



Però anche calcolare prima i pari, poi i dispari:



Quindi i coefficienti saranno tutti scambiati, per ricoverarli dobbiamo scrivere il codice binario dei coefficienti:

$$\begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix} \rightarrow \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} \rightarrow \begin{pmatrix} 0 \\ 4 \\ 2 \\ 6 \\ 1 \\ 5 \\ 3 \\ 7 \end{pmatrix} \begin{matrix} k=0 \\ k=4 \\ k=2 \\ k=6 \\ k=1 \\ k=5 \\ k=3 \\ k=7 \end{matrix}$$

Questa è solamente l'ANALISI di Fourier.

Abbiamo fatto la analisi e trovare la funzione, esiste la formula inversa:

$$f_j = \sum_{k=0}^{N-1} c_k \omega^{-jk} \quad \cos \frac{2\pi k n}{N} + i \sin \frac{2\pi k n}{N}$$

$\uparrow$  complesso coniugato       $\uparrow$  differenza nel segno

$$\text{FT}(f) \xrightarrow{\text{trasformata di Fourier}} \text{IFT} = \text{FT}(\bar{c})$$

$\downarrow$  inversa della trasformata di Fourier       $\downarrow$  complesso coniugato di c

Per fare l'inversa non bisogna modificare l'algoritmo, ma farne il complesso coniugato.

$$f_j = \sum_{k=0}^{N-1} c_k \cdot \omega^{-jk} = \sum_{k=0}^{N-1} (c_k^R + i c_k^I) (\cos + i \sin) = \sum_{k=0}^{N-1} c_k^R \cos + i c_k^R \sin + i c_k^I \cos - c_k^I \sin = \sum_{k=0}^{N-1} (c_k^R \cos - c_k^I \sin) + i (c_k^R \sin + c_k^I \cos)$$

$\cos \frac{2\pi j n}{N}$   
stessa cosa al seno

$$f_j = \sum_{k=0}^{N-1} \bar{c}_k \cdot \omega^{jk} = \sum_{k=0}^{N-1} (c_k^R - i c_k^I) (\cos - i \sin) = \sum_{k=0}^{N-1} c_k^R \cos - i c_k^R \sin - i c_k^I \cos - c_k^I \sin = \sum_{k=0}^{N-1} (c_k^R \cos - c_k^I \sin) - i (c_k^R \sin + c_k^I \cos)$$

con c normale e un  $\bar{c}$  (complesso coniugato) la parte reale è uguale e cambia solo un segno nella parte immaginaria.

Anche se le immagini hanno valori reali, ma lavoriamo con i complessi che hanno parte reale (quella data) e parte immaginaria nulla. Quando facciamo la cosa inversa è un guoto.

Fourier

10/11/2005

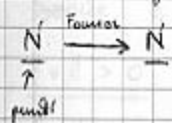
Ma abbiamo fatto Fourier per una siga, la dobbiamo fare per un'immagine

(matrice); ad es:  $f(x) = a_0 + a_1 \cos 2k\tilde{u}x + a_2 \cos 4k\tilde{u}x + \dots + b_1 \sin 2k\tilde{u}x + b_2 \sin 4k\tilde{u}x + \dots =$

$$= \sum_i \alpha_i \varphi_i(x)$$

$$g(y) = c_0 + c_1 \cos 2k\tilde{u}y + c_2 \cos 4k\tilde{u}y + \dots + d_1 \sin 2k\tilde{u}y + d_2 \sin 4k\tilde{u}y + \dots = \sum_e \beta_e \varphi_e(y)$$

$h(x, y) = f(x) + g(y) \rightarrow$  non funziona se le uniamo così! Troppo semplice per l'immagine, perché non il colore si modifica uniformemente!



$N^2 \rightarrow 2N$  (componente x + componente y)

↑  
valori per ogni immagine (punti  $\cdot N \times N =$  parametri)

Ma vogliamo riassumere  $N^2$  parametri in  $2N$  (compressione); ma il modello così è troppo semplice.

Proviamo a fare:  $h(x, y) = f(x) \odot g(y) \rightarrow$  così avremo  $N^2$  prodotti, ovvero: non proprio un per

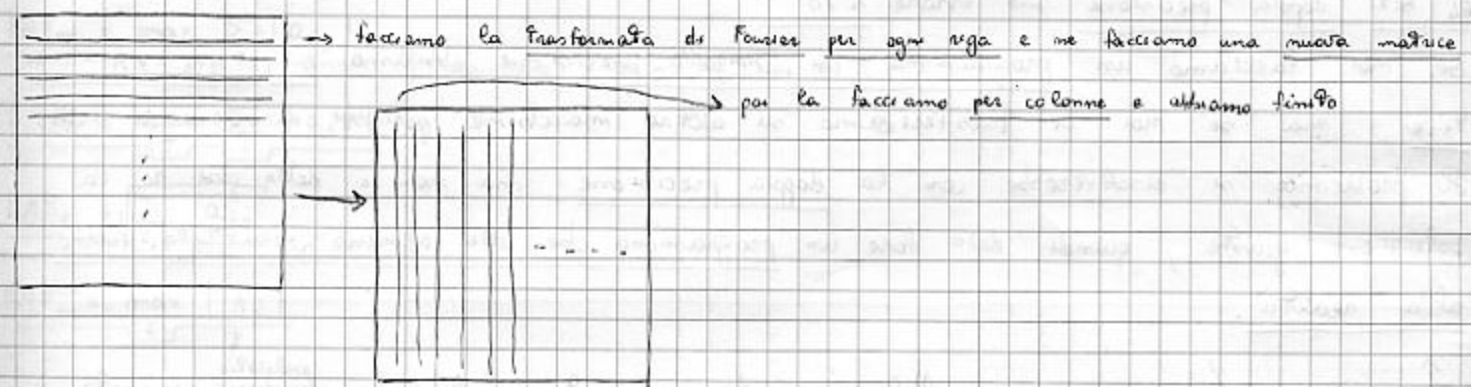
$$h(x, y) = \sum_{k, e} \gamma_{k, e} \varphi_k(x) \varphi_e(y)$$

coefficienti incogniti ( $N^2$ ) della rappresentazione funzionale

Ma questo non è il modello più generale, quello più generale è:  $h(x, y) = \sum_{k, e} \gamma_{k, e} \cdot$

$\varphi_{k, e}(x, y) \rightarrow$  in pratica facciamo  $\varphi_{k, e}(x, y) = \varphi_k(x) \cdot \varphi_e(y)$  (LI ASSIEMO SEPARATI)

cosa che non vale per esempio  $\rightarrow \sin 2k\tilde{u}xy$  qui non si possono separare



Per la trasformata inversa è la stessa cosa (prima per righe applichiamo l'inversa e poi per colonne).

Possiamo fare anche per colonne e per righe in tutti e due i casi.

Possiamo dividere Fourier anche non solo in 2, ma in 3 (resto 0, resto 1, resto 2), oppure per 4 (resto 0, 1, 2, 3) ecc... quindi non è detto che i punti devono essere per forza una potenza di 2. L'efficienza però è migliore con una potenza di 2.

### MALCONDIZIONAMENTO

$$\begin{cases} x - y = 1 \\ x - 1,00001y = 0 \end{cases} \quad \begin{cases} x = 1,00001y \\ 0,00001y = 1 \end{cases} \quad \begin{cases} y = 10^5 \\ x = 100,001 \end{cases} \quad \begin{cases} x = 100,000 \\ y = 100,001 \end{cases}$$

$$Ax = b$$

$$A = \begin{pmatrix} 1 & -1 \\ 1 & -1,00001 \end{pmatrix}$$

$$\begin{cases} x - y = 1 \\ x - 0,99999y = 0 \end{cases} \rightarrow \begin{cases} x = 0,99999y \\ y = x - 1 \end{cases} \rightarrow \begin{cases} x = 0,99999y \\ y = -1 + 0,99999y \end{cases}$$

$$\begin{cases} x = -0,99999 \\ y = -\frac{1}{0,00001} = -100'000 \end{cases}$$

$$\begin{cases} x = 99'999 \\ y = -100'000 \end{cases}$$

$$Ax = b$$

$$A = \begin{pmatrix} 1 & -1 \\ 1 & -0,99999 \end{pmatrix}$$

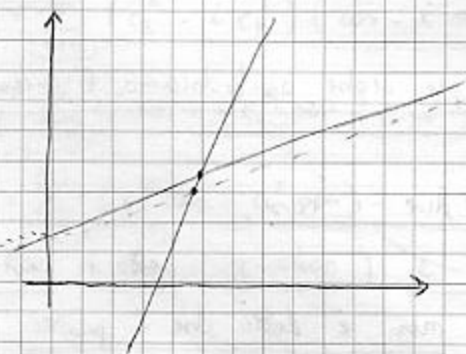
LE DUE MATRICI SONO VERAMENTE MOLTO SIMILI, MA IL PROBLEMA È CHE LE SOLUZIONI SONO COMPLETAMENTE DIVERSE.

Questo vuol dire che commettendo un piccolissimo errore abbiamo dei risultati diversi, mi. Quindi non sappiamo se i nostri risultati sono veri e quanto lontani dalla verità.

32 bit singola precisione → errore  $\sim 10^{-6}$

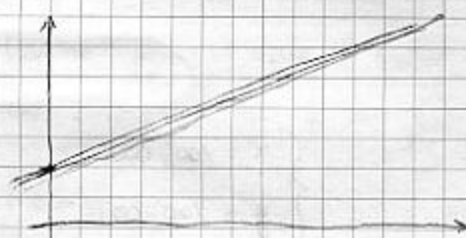
64 bit doppia precisione → errore  $\sim 10^{-16}$

Se noi facciamo un programma in singola precisione abbiamo 6 cifre significative, ma se noi ci trasferiamo su altre macchine potrebbe accadere di fatto il problema si risolverebbe con la doppia precisione, ma non è detto che ho la soluzione giusta, quindi devo fare un programma che dia almeno risultato vicino alla realtà.



Nel caso di un sistema lineare di due incognite, dove troviamo l'intersezione delle due rette, se le spostiamo di poco il risultato non cambia molto.

Ma non era il nostro caso.



Le nostre rette erano quasi parallele:

$$\begin{pmatrix} 1 & -1 \\ 1 & -1,00001 \end{pmatrix} \quad \begin{pmatrix} 1 & -1 \\ 1 & -0,99999 \end{pmatrix}$$

$$\begin{cases} x - y = 1 \\ x - 2y = 0 \end{cases}$$

Dobbiamo avere strumenti che ci dicono se una matrice è malcondizionata o bencondizionata.

## Norma

Grandezza di uno scalare, valore assoluto.  $\rightarrow \|v\|$

Serve a misurare la grandezza del vettore.

$$\|v\| > 0 \quad \text{se} \quad v \neq \emptyset$$

$$\|v\| = 0 \quad \Leftrightarrow \quad v = \emptyset$$

$$\|\alpha \cdot v\| = |\alpha| \cdot \|v\|$$

$$\|u + v\| \leq \|u\| + \|v\|$$

Come si calcola la norma di un vettore?

$$\|u\|_{\infty} = \max_i |u_i|$$

$$\|u\|_1 = \sum_i |u_i|$$

$$\|u\|_2 = \sqrt{\sum_i u_i^2}$$

Le norme si definiscono anche per le matrici:

$$\|\alpha \cdot A\| = |\alpha| \cdot \|A\|$$

$$\|A + B\| \leq \|A\| + \|B\|$$

Come si calcolano le norme di matrici?

$$\|A\|_{\infty} = \max_i \sum_j |a_{ij}|$$

$$\|A\|_1 = \max_j \sum_i |a_{ij}|$$

$$\|A\|_F = \sqrt{\sum_{i,j} a_{ij}^2}$$

$$\|A\|_2 = \max_i |\lambda_i|$$

↑  
autovalore

La norma mi dice quanto è grande la matrice, quanto questa matrice ~~amplifica~~ amplifica il vettore con cui dobbiamo fare l'operazione.

$Ax = b \rightarrow$  diciamo che commettiamo un certo errore su essa

$A(x + \Delta x) = b + \Delta b \rightarrow$  se la collego a quella prima ho:

$$A\Delta x = \Delta b$$

allora è vero anche  $\Delta x = A^{-1} \Delta b$  (su matrici compatibili)  $\rightarrow$  questo mi dice quant'è l'errore che ho commesso

$$\|\Delta x\| = \|A^{-1} \Delta b\| \leq \|A^{-1}\| \cdot \|\Delta b\|$$

Ricordandomi, dico che  $\|b\| = \|Ax\| \leq \|A\| \cdot \|x\|$

$$\|\Delta x\| \leq \|A^{-1}\| \cdot \|\Delta b\|$$

$$\|A\| \cdot \|x\| \geq \|b\|$$



$$\frac{\|\Delta x\|}{\|A\| \cdot \|x\|} \leq \frac{\|A^{-1}\| \cdot \|\Delta b\|}{\|b\|}$$

$$\frac{\|\Delta x\|}{\|x\|} \leq \|A\| \cdot \|A^{-1}\| \cdot \frac{\|\Delta b\|}{\|b\|}$$

errore relativo su  $x$        $\gamma(A)$       questo da errore relativo su  $b$   
 a noi interessa questa parte (INDICE DI CONDIZIONAMENTO)

$$\gamma(A) = \|A\| \cdot \|A^{-1}\|$$

$$1 = \|I\| = \|A \cdot A^{-1}\| \leq \|A\| \cdot \|A^{-1}\| \rightarrow \text{norma della matrice identica}$$

$\|A\| \cdot \|A^{-1}\| \geq 1 \rightarrow$  l'errore non si può ridurre, ma solo amplificare

$\frac{\|\Delta b\|}{\|b\|} \rightarrow$  errore relativo rispetto a  $b$ , non da altro che le sue cifre significative

es:  $\Delta b = 0,001$        $b = 10,00000000 \dots$

Supponiamo di scrivere  $\gamma(A) = 10^9$ , la formula sarà:

$$\frac{\|\Delta x\|}{\|x\|} \leq 10^9 \cdot \frac{\|\Delta b\|}{\|b\|} \rightarrow P \text{ cifre significative}$$

delle  $P$  cifre significative noi ce ne mangiamo  $9$ .

Con la singola precisione si possono tollerare matrici malcondizionate fino a  $10^6$  cifre significative. Altrimenti dobbiamo passare alla doppia precisione.

Per vedere ciò basta una addizione e una  sottrazione.

Con l'addizione, anche se tronchiamo la parte decimale, il risultato è simile, nella sottrazione il risultato ha meno cifre significative degli operati.

Quindi un metodo è stabile se non propaga gli errori del calcolatore.

### Fourier

16/11/2005

Si da una rappresentazione funzionale, di una discretizzata.

$$f(x) = a_0 + a_1 \cos 2\pi \tilde{f} x + a_2 \cos 4\pi \tilde{f} x + \dots$$

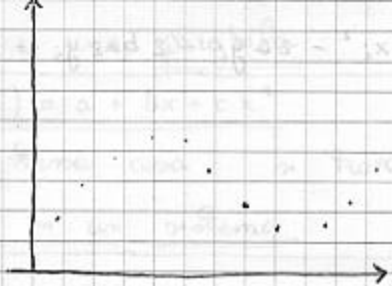
$$f(x) = \sum_i a_i \varphi_i(x)$$

Il obiettivo dell'analisi di Fourier è trovare i coefficienti in un numero finito di passi.

$$f(x_1) = a_0 \varphi_0(x_1) + a_1 \varphi_1(x_1) + a_2 \varphi_2(x_1) + \dots$$

$$f(x_2) = a_0 \varphi_0(x_2) + a_1 \varphi_1(x_2) + a_2 \varphi_2(x_2) + \dots$$

Questo sistema è risolvibile in modo analitico (più veloce).



→ conosco la funzione in un set di punti (dove  
possiamo passare infinite curve, non sappiamo cosa  
accade in mezzo - problema di interpolazione -).

Impongo che la funzione dove io conosco i punti valga:  $f(x_1), f(x_2) \dots$

Cosa accade in mezzo dipende da cosa abbiamo scelto per  $\varphi$ .

Potrebbe essere una retta (quasi corretto se i punti sono molto fitti). Molto più probabilmente è una curva. Seno e coseno sono delle curve con infinite derivate, quindi sono molto usate, le funzioni scelte non devono essere troppo lisce. (lisce: che non hanno bruschi cambiamenti).

In natura ci sono comportamenti lineari, quadratici (gravità), cubici...

Se una funzione può essere approssimata alla serie di Taylor, allora può essere ricondotta ad un polinomio, ma ci sono dei problemi quando il grado della funzione è grande: 1° è che nei punti che non sappiamo ha un comportamento oscillatorio.

Fare una funzione di grado molto alto è complicata e spesso non dà il risultato corretto, potrebbe bastarci  $f(x) = a_0 + a_1 x$  e mi servono due punti per trovare la; come lo scelgo? Usando il metodo dei minimi quadrati!

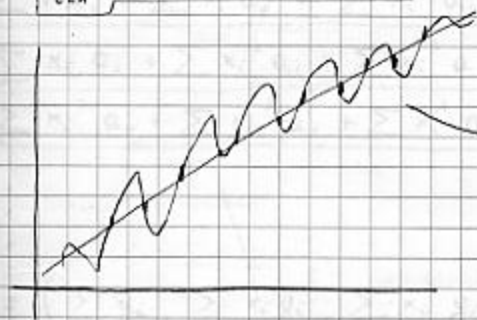
Minimi quadrati

$$S := \sum_{i=1}^N (y_i - a - b x_i)^2$$

definito come

$$y_i \equiv f(x_i)$$

$$f(x) = a + b x$$



si fa la somma dei quadrati delle parti ammesse

il parametro  $a$  e  $b$  migliori sono quelli per cui  $S$  è minimo.

Come facciamo a trovare il minimo di una funzione? Facciamo la derivata e vediamo dove sono i punti di massimo o minimo.

E a due variabili? Imponiamo che la derivata sia nulla con quei parametri:

$$\frac{\partial f(x, y)}{\partial x} = 0$$

$$\frac{\partial f(x, y)}{\partial y} = 0$$

DERIVATE PARZIALI → E POI DOBBIAMO CALCOLARE IL DET. HESSIANO:

$$\left| \frac{\partial^2 f(x, y)}{\partial x^2} \right| > 0$$

imporre le derivate rispetto ad  $a$  e  $b$  uguali a 0

$$S(a, b) \rightarrow \frac{\partial S}{\partial a} = 0 \quad \frac{\partial S}{\partial b} = 0$$

$$S(a, b) := \sum_{i=1}^N (y_i - a - b x_i)^2$$

$$S(a, b) := \sum_{i=1}^n (y_i - a - bx_i)^2 = \sum_{i=1}^n (y_i^2 + a^2 + b^2 x_i^2 - 2ay_i - 2bx_i y_i + 2abx_i)$$

$$\begin{cases} \frac{\partial S}{\partial a} = \sum_{i=1}^n 2a - 2y_i + 2bx_i = 0 \\ \frac{\partial S}{\partial b} = \sum_{i=1}^n 2bx_i^2 - 2x_i y_i + 2ax_i = 0 \end{cases}$$

$$\begin{cases} Na + \sum x_i b = \sum y_i \\ \sum x_i a + \sum x_i^2 b = \sum x_i y_i \end{cases}$$

$$Ah = z \quad A := \begin{pmatrix} N & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix}; \quad z = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$

RSOLVIAMO IL SISTEMA; RSCRIVENDOLO COSÌ:

$$\begin{cases} S_0 a + S_1 b = T_0 \\ S_1 a + S_2 b = T_1 \end{cases} \quad \text{dove} \quad S_i = \sum_{k=1}^n x_k^i \quad T_i = \sum_{k=1}^n x_k^i y_k$$

$$a = \frac{(T_0 - S_1 b)}{S_0}$$

$$b = \left( T_1 - \frac{S_1}{S_0} (T_0 - S_1 b) \right) \frac{1}{S_2 - S_1^2/S_0} \Rightarrow$$

$$\Rightarrow b = \left( S_0 T_1 - S_1 T_0 + S_1^2 b \right) \frac{1}{S_2 S_0 - S_1^2} \Rightarrow$$

$$\Rightarrow b (S_2 S_0 - S_1^2) = S_0 T_1 - S_1 T_0 \Rightarrow$$

$$\Rightarrow b = \frac{(S_0 T_1 - S_1 T_0)}{(S_2 S_0 - S_1^2)}$$

$$a = \left[ T_0 - S_1 \frac{(S_0 T_1 - S_1 T_0)}{S_2 S_0 - S_1^2} \right] \frac{1}{S_0} \Rightarrow$$

$$\Rightarrow a = \frac{(S_2 S_0 - S_1^2) T_0 - S_1 (S_0 T_1 - S_1 T_0)}{S_0 (S_2 S_0 - S_1^2)} \Rightarrow$$

$$\Rightarrow a = \frac{S_2 T_0 - S_1 T_1 + S_1^2 T_0}{S_2 S_0 - S_1^2}$$

$$\text{quindi: } \begin{cases} a = \frac{S_2 T_0 - S_1 T_1}{S_2 S_0 - S_1^2} \\ b = \frac{S_0 T_1 - S_1 T_0}{S_2 S_0 - S_1^2} \end{cases}$$

Chi ci dice che è un max o un min? Dobbiamo andare a vedere l'hessiano:

$$\begin{vmatrix} \frac{\partial^2 S}{\partial a^2} & \frac{\partial^2 S}{\partial a \partial b} \\ \frac{\partial^2 S}{\partial b \partial a} & \frac{\partial^2 S}{\partial b^2} \end{vmatrix} \rightarrow \text{o calcoliamo l' Hessiano, o facciamo un metodo che ci dice chi è il max.}$$

$$S = \sum_{i=1}^n (y_i - a - bx_i)^2 \rightarrow \text{scegliamo sempre } \underline{a} \text{ e } \underline{b} \text{ tale che } S \rightarrow \infty$$

Quindi la quant.  $S(a, b)$  non ha massimo, quindi è per forza un minimo.

Se noi abbiamo una parabola:

$$f(a, b, c; x) = a + bx + cx^2$$

o fa la stessa cosa, o troviamo le tre incognite facendo le derivate parziali e mettendole in un sistema.

17/11/2005

Minimi quadrati

$$S = \sum_{i=1}^N (y_i - a_0 - a_1 x_i - a_2 x_i^2 - a_3 x_i^3 - \dots)^2$$

$$\frac{\partial S}{\partial a_0} = -2 \sum_{i=1}^N (y_i - a_0 - a_1 x_i - a_2 x_i^2 - a_3 x_i^3 - \dots) =$$

REGOLA

$$\frac{d g^2(x)}{dx} = 2g(x) \cdot \frac{dg(x)}{dx}$$

$$\frac{\partial S}{\partial a_1} = -2 \sum_{i=1}^N (y_i - a_0 - a_1 x_i - a_2 x_i^2 - a_3 x_i^3 - \dots) x_i =$$

$$\frac{\partial S}{\partial a_2} = -2 \sum_{i=1}^N (y_i - a_0 - a_1 x_i - a_2 x_i^2 - a_3 x_i^3 - \dots) x_i^2 =$$

$$= -2 \sum_{i=1}^N y_i + 2a_0 \sum 1 + 2a_1 \sum x_i + 2a_2 \sum x_i^2 + 2a_3 \sum x_i^3 + \dots = 0$$

$$= -2 \sum_{i=1}^N x_i y_i + 2a_0 \sum x_i + 2a_1 \sum x_i^2 + 2a_2 \sum x_i^3 + 2a_3 \sum x_i^4 + \dots = 0$$

$$= -2 \sum_{i=1}^N x_i^2 y_i + 2a_0 \sum x_i^2 + 2a_1 \sum x_i^3 + 2a_2 \sum x_i^4 + 2a_3 \sum x_i^5 + \dots = 0$$

Alla fine avremo:

$$\begin{cases} \sum 1 a_0 + \sum x_i a_1 + \sum x_i^2 a_2 + \sum x_i^3 a_3 + \dots = \sum y_i \\ \sum x_i a_0 + \sum x_i^2 a_1 + \sum x_i^3 a_2 + \sum x_i^4 a_3 + \dots = \sum x_i y_i \\ \sum x_i^2 a_0 + \sum x_i^3 a_1 + \sum x_i^4 a_2 + \sum x_i^5 a_3 + \dots = \sum x_i^2 y_i \\ \vdots \\ \vdots \end{cases}$$

ogni volta si comincia con un grado in più nella x

$$\underline{h} = (\sum y_i, \sum x_i y_i, \sum x_i^2 y_i, \dots)^T$$

$$\underline{z} = (a_0, a_1, a_2, a_3, \dots)^T$$

Avremo una matrice:

$$A \cdot \underline{z} = \underline{h}$$

$$S_i = \sum_k x_k^i$$

$$T_i = \sum_k y_k x_k^i$$

$$\begin{pmatrix} S_0 & S_1 & S_2 & S_3 & \dots \\ S_1 & S_2 & S_3 & S_4 & \dots \\ S_2 & S_3 & S_4 & S_5 & \dots \\ S_3 & S_4 & S_5 & S_6 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$

ma questa matrice è totalmente mal condizionata, e più m va avanti col grado peggio è

Noi potremmo scegliere un polinomio di grado m-1, rispetto agli m punti che abbiamo, dimodoché i parametri saranno esattamente n: a<sub>0</sub>, a<sub>1</sub>, a<sub>2</sub>, ..., a<sub>m</sub>.

polinomio  $N-1$

$$f(x) = a_0 + a_1 x + \dots + a_{N-1} x^{N-1}$$

$$\begin{cases} f(x_1) = y_1 = a_0 + a_1 x_1 + \dots + a_{N-1} x_1^{N-1} \\ \vdots \\ f(x_m) = y_m = a_0 + a_1 x_m + \dots + a_{N-1} x_m^{N-1} \end{cases}$$

Riprendiamo:

$$S = \sum (y_i - \sum_{j=0}^{N-1} a_j x_i^j)^2 \rightarrow \text{se va a } \phi \text{ vuol dire che la soluzione coincide perfettamente con i punti dati}$$

↑  
tende al polinomio di interpolazione

Quando la funzione oscilla molto, vuol dire che è malcondizionata.

Nel malcondizionamento anche un errore piccolissimo  $\rightarrow$  propaga e porta a risultati assolutamente senza senso.

Malcondizionamento nelle macchine

$$a = m \cdot 10^i \quad 0 \leq m < 1 \quad 0, \dots$$

↓  
mantissa

$$327,4 \rightarrow 0,3274 \cdot 10^3$$

$$0,0022 \rightarrow 0,22 \cdot 10^{-2}$$

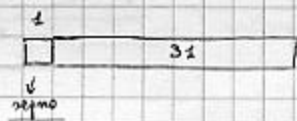
Ma non esiste solo la base dieci.

$$a = m \cdot B^i \text{ nel sistema di base } B$$

Nella macchina si usano i transistor (ON/OFF) per rappresentare gli interi.

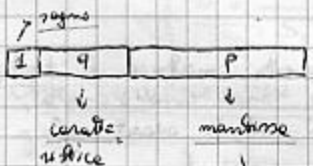
1	2	3	$\rightarrow 2^n$
0	00	000	
1	01	001	
	10	010	
	11	011	
		100	
		101	
		110	
		111	

Si riservano 32 bit:



$$\text{si arriva a } \frac{2^{32}-1}{2} = 2 \cdot 147 \cdot 483 \cdot 647$$

Ma molti numeri non rientrano nel range.



$\rightarrow$  altro metodo per rappresentarli

MANTISSA COMPRESA TRA 0,1 e 1

$$a = 2^s (x_1 \cdot 2^{-1} + x_2 \cdot 2^{-2} + x_3 \cdot 2^{-3} + \dots + x_p \cdot 2^{-p})$$

dove  $|s| \leq 2^{q-1}$

Quindi il numero più alto sarà:  $2^{2^{q-1}}$  (quando tutti i bit della caratteristica sono pari a 1).

$$\begin{matrix} 1 & 8 & & 23 \end{matrix} \rightarrow 2^{128} \approx 10^{38}$$

Qual è il numero più piccolo?  $\rightarrow s = -2^{p-1}$  e tutte le  $x$  sono a  $\phi$ .

Il numero più piccolo dopo  $\phi$  è quello dove solo  $x_p$  è  $\pm \rightarrow 2^s \cdot 2^{-p}$   
 con  $s = \phi$

$1 + \epsilon = 1$   
 perché  $\epsilon$  è molto piccolo (tipo 0,000000100...) il calcolatore  $\approx$  perde un sacco di cifre.  
 Io devo trovare l' $\epsilon$  più piccolo per cui  $1 + \epsilon > 1$  (zero macchina).  
 Esso è  $2^{-p}$  (ovvero tutte le  $x$  a  $\phi$  e  $x_p$  a  $\pm$ ).  
 Un fermarsi reale questo è  $2^{-23} \approx 10^{-7}$

$1 + \epsilon = 1$   
 $\epsilon < \epsilon_0 \approx 10^{-7}$

Ma  $10^{-7}$  non è ancora sufficiente.

Altro modo:  $\rightarrow$  aggiunto 5 bit alla caratteristica

1	11	20	32 BIT	$2^{1023} \approx 10^{308}$	MAX
				$2^{-20} \approx 10^{-6}$	ZERO MACCHINA

Ma oggi  $\approx$  usa doppia precisione:

1	11	52	64 BIT	$2^{1023} \approx 10^{308}$	MAX
				$2^{-52} \approx 10^{-16}$	ZERO MACCHINA

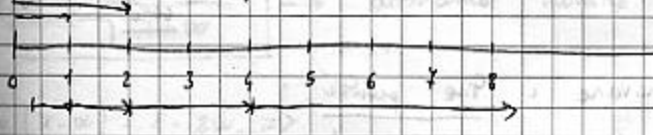
Underflow: quando il numero è più piccolo del valore assoluto di  $10^{-308}$ , la soluzione è che  $\approx$  mette il numero a  $\phi$ .

Per l'overflow la soluzione non è così semplice.

Consideriamo tre casi:

5	2 <sup>s</sup>	$x_1 = 1, x_2 = \dots = x_{p-1} = \phi$	$x_1 = \dots = x_{p-1} = \phi, x_p = 1$	$x_1 = \dots = x_p = 1$	$2^s \cdot 2^{-p} (1 + 2 + \dots + 2^{p-1})$ $= 2^s \cdot 2^{-p} (2^p - 1)$ $= 2^s \cdot \frac{2^p - 1}{2^p}$ $= 2^s (1 - 2^{-p})$
4		$2^s \cdot 2^{-1}$	$2^s \cdot 2^{-p}$	$2^s (2^{-1} + 2^{-2} + \dots + 2^{-p})$	
0	1	$2^{-1}$	$2^{-p}$	$1 - 2^{-p}$	
1	2	1	$2^{-(p-1)}$	$2 - 2^{-(p-1)}$	
2	4	2	$2^{-(p-2)}$	$4 - 2^{-(p-2)}$	
3	8	4	$2^{-(p-3)}$	$8 - 2^{-(p-3)}$	

$\rightarrow [0, 1]$   
 $\rightarrow [0, 2]$   
 $\rightarrow [0, 4]$   
 $\rightarrow [0, 8]$



$\textcircled{2} \div \textcircled{3}$   
 $\textcircled{1} \div \textcircled{3}$

Perché il primo bit è sempre 1, non  $\approx$  rappresenta per niente nelle macchine, ma esso rimane virtualmente a disposizione della mantissa. (NOTAZIONE IMPLICITA 1)

Quando i bit della mantissa passano da 52 a 53 e ho  $2^{-53}$ .

Vediamo come gli errori si propagano:

$$\sqrt{1+\delta} - 1 \quad \delta < \epsilon \quad \sqrt{1+\delta} - 1 = \phi \rightarrow \text{ma noi sappiamo che non è vero}$$

$$\frac{(\sqrt{1+\delta} - 1)(\sqrt{1+\delta} + 1)}{(\sqrt{1+\delta} + 1)} = \frac{1+\delta - 1}{\sqrt{1+\delta} + 1} = \frac{\delta}{\sqrt{1+\delta} + 1} \xrightarrow{\delta \rightarrow 0} \frac{\delta}{1+1} = \frac{\delta}{2}$$

Soluzione dei minimi quadrati

24/11/2005

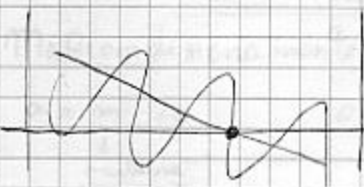
Polinomio con grado basso.

REGG Y 2002 @ CIBANO . IT

$$S := \sum_{i=1}^n (y_i - f(\underline{a}; x_i))^2 \quad f(\underline{a}) := a_0 + a_1 x + \dots + a_m x^m$$

la matrice dovrà avere determinante maggiore di zero.

A volte capita che minimizzando una funzione non si riesce a trovare la soluzione, ma dobbiamo approssimare.



Newton e bisezione sono dei metodi iterativi dove lo spazio in cui si poteva trovare lo zero veniva sempre ridotto.

Col Newton si fa prima, bisogna trovare le derivate parziali.

$$S = f(x; \theta)$$

$$S = \sum_{i=1}^n (y_i - \log(a + x_i) + x_i)^2$$

$$y = \log(a + x) + x$$

$$\frac{\partial S}{\partial a} = -2 \sum_{i=1}^n (y_i - \log(a + x_i) + x_i) \cdot \frac{1}{a + x_i} = \phi \rightarrow \text{eq. non lineare}$$

Ma non sappiamo se è un MAX o un MIN, ~~potrebbe essere un MIN locale.~~

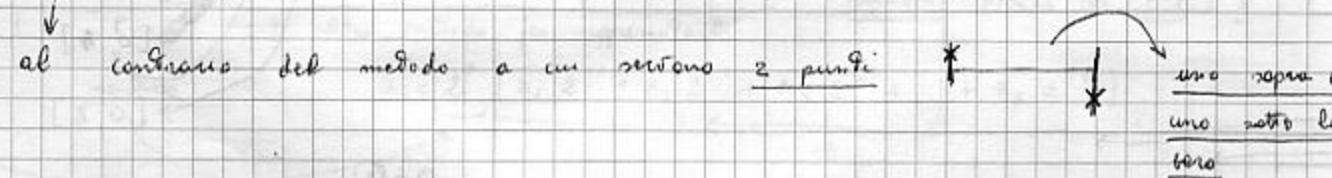
potrebbe però essere un MIN locale.

Allora è meglio usare un metodo globale.

Un metodo ad hoc è:



abbiamo bisogno di tre punti in questo caso.

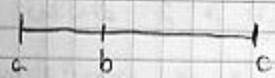


Piano piano si restringe lo spazio dove trovare i tre punti.

Noi parliamo solo di minimizzazione di funzioni perché la massimizzazione è la minimizzazione col segno meno.

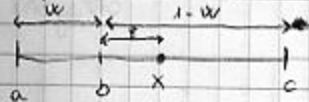
Metodo della sezione aurea

Metodo lento ma non richiede l'uso delle derivate. È un metodo robusto (funziona abbastanza spesso). Più lento di quello della bisezione.



$$f(b) < \{f(a), f(c)\} \rightarrow \text{condizione di partenza}$$

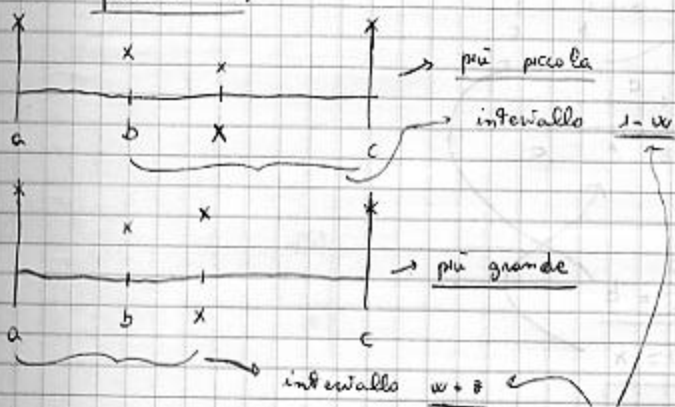
Dopo mettiamo un punto  $x$ :



$$w := \frac{b-a}{c-a} \rightarrow \text{dice quanto } b \text{ sia vicino ad } a$$

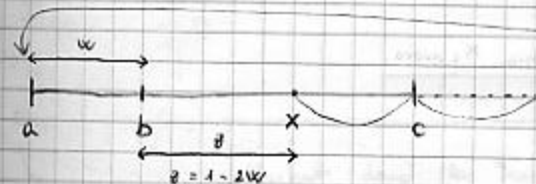
$$z := \frac{x-b}{c-a} \quad 1-w = \frac{c-b}{c-a}$$

Dobbiamo vedere se la funzione nel punto  $x$  è più piccola o più grande che nel punto  $b$ .



Dobbiamo cercare di evitare intervalli troppo piccoli, altrimenti nessuno ci garantisce che dall'altra parte non ci siano dei minimi. Bisogna ridurre l'intervallo, ma con prudenza; quindi  $1-w = w+z \rightarrow z = 1-2w$

usa nel dare?



Ma questo  $w$  chi ce lo ha detto?

Ammettiamo che sia stato calcolato come noi calcoliamo  $z$ .

$\left. \begin{array}{l} \text{distanza di } b \text{ da } a \rightarrow w \\ \text{distanza di } x \text{ da } b \rightarrow \frac{z}{1-w} \end{array} \right\} \text{relative!} \rightarrow \text{dobbiamo imposte uguali } w = \frac{z}{1-w} \Rightarrow$

$$\Rightarrow w = \frac{1-2w}{1-w} \rightarrow \text{dobbiamo trovare quanto vale } w$$

$$w(1-w) = 1-2w \Rightarrow$$

$$w - w^2 = 1 - 2w \Rightarrow$$

$$w^2 - 2w - w + 1 = 0$$

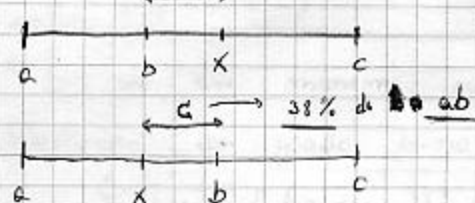
$$w^2 - 3w + 1 = 0 \Rightarrow w = \frac{3 \pm \sqrt{9-4}}{2} = \frac{3 \pm 5}{2} \rightarrow \text{scartiamo la radice positiva perché supererebbe l'intervallo.}$$



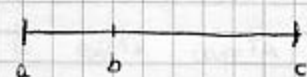
$$x = \frac{3 - \sqrt{5}}{2} = 0,38197... = G \rightarrow \text{numero uscente a Pitagora}$$

Mez punto bc scegliamo x al 40% vicino a b

al 60% vicino a c

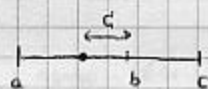


Algoritmo

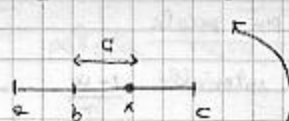


$$f(b) < \{f(a), f(c)\}$$

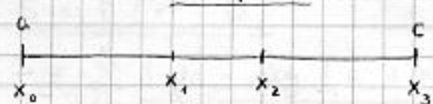
$$(b-a) > c-b \quad x = b - G(b-a)$$



$$(c-b) > (b-a) \quad x = b + G(c-b)$$

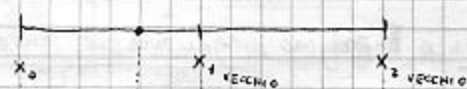


Avremo 4 punti:



$$\begin{aligned} x_1 &= x & x_2 &= b \\ x_1 &= b & x_2 &= x \end{aligned}$$

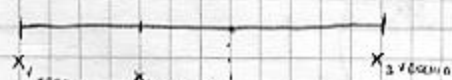
CASO  $f(x_1) < f(x_2)$ :



$$x = x_1_{\text{vecchio}} - G(x_2_{\text{vecchio}} - x_0_{\text{vecchio}}) = (1-G)x_2_{\text{vecchio}} + Gx_0_{\text{vecchio}}$$

Ora si fa la stessa cosa con l'intervallo  $x_0_{\text{nuovo}} x_2_{\text{nuovo}}$ .

CASO  $f(x_2) < f(x_1)$ :



$$x = x_2_{\text{vecchio}} + G(x_3_{\text{vecchio}} - x_2_{\text{vecchio}}) = (1-G)x_3_{\text{vecchio}} + Gx_2_{\text{vecchio}}$$

Andiamo avanti con l'intervallo giusto.

- fissati a, b, c

$$x_0 = a; \quad x_3 = c$$

$$b-a > c-b \Rightarrow x_1 = b - G(b-a); \quad x_2 = b$$

$$c-b > b-a \Rightarrow x_2 = b + G(c-b); \quad x_1 = b$$

calcola f0, f1, f2, f3

loop fino a che  $x_3 - x_0 > \epsilon$

$$f_2 < f_2 \Rightarrow x_0^{\text{nuovo}} = x_0; x_1^{\text{nuovo}} = x_1; x_2^{\text{nuovo}} = x_2; x_3^{\text{nuovo}} = G x_2 + (1-G) x_1 \rightarrow \text{procedimento ricorsivo}$$

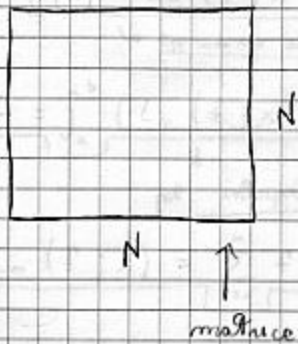
$$f_2 < f_1 \Rightarrow x_0^{\text{nuovo}} = x_1; x_1^{\text{nuovo}} = x_2; x_2^{\text{nuovo}} = x_3; x_3^{\text{nuovo}} = (1-G) x_2 + G x_3$$

calcola  $f_0, f_1, f_2, f_3$

vide finisci  $x_0, x_1, x_2, x_3$

## Progetto

1/12/2005

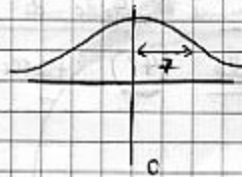


immagini  
 $i, j = 1 \dots 256$

$$I_{ij} \rightarrow I_{ij}^e = I_{ij} + E_{ij}$$

livelli di grigio

generiamo numeri casuali (rand) e lo aggiungiamo al pixel



numero a caso estratto dalla distribuzione gaussiana



$$0 \leq \eta_{i,j} \leq 1$$

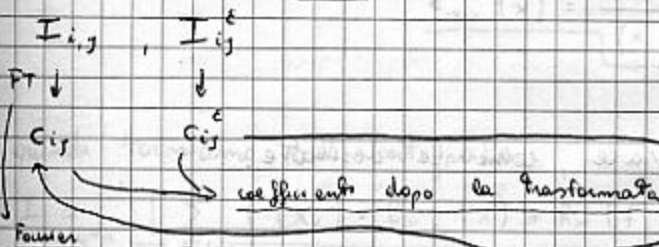
$$-10 \leq E_{i,j} \leq 10$$

$$E_{i,j} = 20 (\eta_{i,j} - 0,5)$$

0 - 255  $\rightarrow$  toni di grigio

i num. casuali generati non potrebbero essere consentiti allora min di 0  $\rightarrow$  li mettiamo a 0, mag. di 255  $\rightarrow$  li mettiamo a 255.

arrotondiamo  $I_{i,j}$  ottenuto all'intero piú vicino  $\rightarrow 7,4 \dots \rightarrow 7$ .



Fourier toglie il disturbo dal segnale: si buttano tutti i coefficienti piccoli ( $= 0$ )

$$P_{i,j}^e(\theta) := \begin{cases} C_{i,j}^e, & |C_{i,j}^e| > \theta \\ \emptyset, & |C_{i,j}^e| \leq \theta \end{cases}$$

$$|a| := \sqrt{R_i^2(a) + I_n^2(a)}$$

$\downarrow$  reale                       $\downarrow$  immaginaria

$$E(\theta) := \sum_{i,j=1}^n (P_{i,j}^e(\theta) - C_{i,j}^e)^2$$

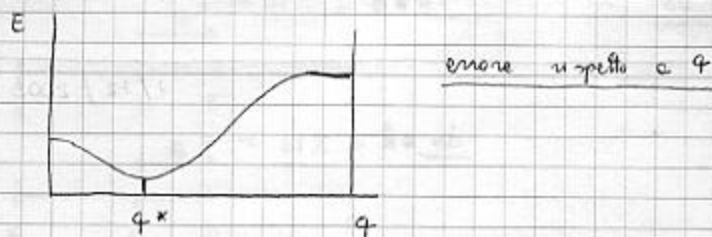
$\uparrow$  errore dipende da  $\theta$                        $\uparrow$  con errore                       $\uparrow$  imm. vera

Vediamo che succede per  $E(0)$ :

mi tengo tutti i coefficienti.

Se  $E(\infty)$ :

azzerò tutto perché saranno tutti più piccoli. (unico colore)



voglio  $q$  che risulta più piccolo

$$E(q) = \sum_{i,j} (p_{ij}^e(q) - c_{ij})^2$$

devo trovare il minimo di questo

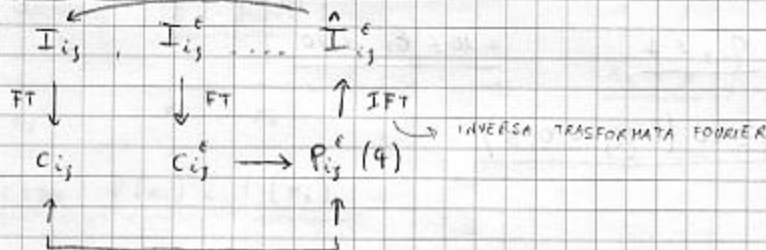
si fa con l'algoritmo "regolabile" che nel Matlab non c'è

si può usare un qualunque metodo di minimizzazione

Dopo che abbiamo ottenuto  $q$  ottimale ( $q^*$ ) facciamo la trasformata inversa:

$$p_{ij}^e(q) \rightarrow \hat{I}_{ij}^e$$

RIASSUNTO:



Pacchetto: Image J → legge le imm già in formato binario (non lo usa nessuno!)

Java linguaggio più usato. Max 2 persone.

Formule di quadratura

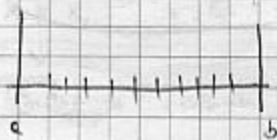
Per approssimare gli integrali.

Non è sempre immediato risolvere l'integrale, come le derivate.

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

$$I = I_n + E_n$$

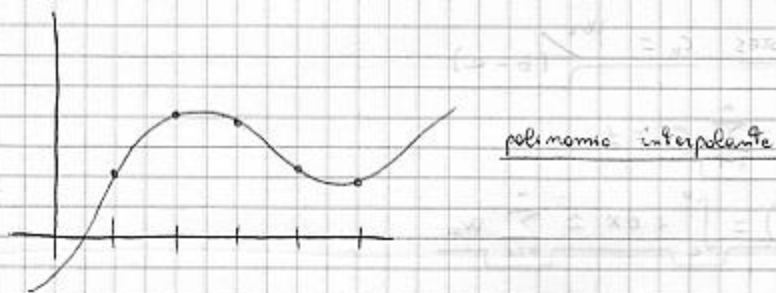
$\downarrow$  integrale approssimato       $\downarrow$  errore



$$I_n = \sum_{k=1}^n w_k \cdot f(x_k)$$

le formule di quadratura servono a scegliere i modi o i punti

$$I = \sum_{k=1}^n w_k f(x_k) + E_n$$



$$f(x) \approx L_n(f; x)$$

↳ approssimo al pol. interpolante

$$I = \int_a^b (L_n(f; x) + R_n(f)) dx$$

pol. interpolante di Lagrange

$$L_n(f; x) = \sum_{k=0}^n f(x_k) \cdot l_{n,k}(x)$$

$$L_n(f; x)$$

$$L_n(f; x) = \sum l_{n,k}(x) \cdot f(x_k)$$

$$l_{n,k}(x) = \begin{cases} 1 & k=l \\ 0 & k \neq l \end{cases}$$

$$l_{n,0}(x) \rightarrow l_{n,0}(x_0), l_{n,0}(x_1) \dots l_{n,0}(x_n)$$

$$l_{n,1}(x) \rightarrow l_{n,1}(x_0), l_{n,1}(x_1) \dots l_{n,1}(x_n)$$

$$l_{n,n}(x) \rightarrow l_{n,n}(x_0), l_{n,n}(x_1) \dots l_{n,n}(x_n)$$

$$L_n(f; x_k) = \sum l_{n,k}(x_k) \cdot f(x_k)$$

$$= f(x_k)$$

verifica la proprietà di interpolazione

$$l_{n,k}(x) = \frac{(x-x_0) \dots (x-x_{k-1}) \cdot (x-x_{k+1}) \cdot \dots \cdot (x-x_n)}{(x_k-x_0) \cdot \dots \cdot (x_k-x_{k-1}) \cdot (x_k-x_{k+1}) \cdot \dots \cdot (x_k-x_n)}$$

può valere 0

↳ tutti i termini  $x_k$

quindi formiamo all' integrale:

$$I = \int_a^b \left[ \sum_{k=0}^n f(x_k) \cdot l_{n,k}(x) + R_n(f, x) \right] dx =$$

$$= \sum_{k=0}^n f(x_k) \int_a^b l_{n,k}(x) dx + \int_a^b l_{n,k}(x) dx + \int_a^b R_n(f; x) dx =$$

$$= \sum_{k=0}^n w_k f(x_k) + E_n$$

$w_k$  (pesi)  
 $E_n$  (errore)

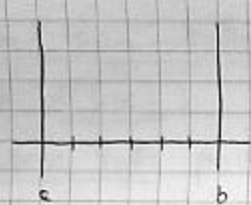
Una formula di quadratura (scelta di un insieme di nodi + pesi) ha un grado di precisione d dove:

$$R_n = 0 \quad n \leq d$$

Metiamoci in una situazione in cui i nodi sono equidistanti:

$$I = \sum_{k=0}^n w_k f(x_k) + E_n \quad w_k = \int_a^b l_{n,k}(x) dx$$

FORMULE DI NEWTON - COTES



NUMERI DI COTES  $c_k = \frac{w_k}{(b-a)}$

$$\sum_{k=1}^n c_k = 1$$

$$(b-a) = \int_a^b 1 dx \approx \sum_{k=1}^n w_k$$

$$\lim_{n \rightarrow \infty} \sum_{k=1}^n |c_k| = \infty$$

↓  
dipende dall'interpolazione (oscillanti)

Le formule di Cotes hanno senso per un numero di punti piccolo.

Esempio (2 punti)



$$I = \frac{1}{2} (f(a) + f(b)) \cdot (b-a)$$

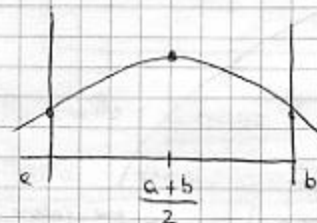
$$I = \frac{h}{2} (f(a) + f(b)) \text{ dove } h = (b-a)$$

quindi  $I = w_1 f(x_1) + w_2 f(x_2)$   $w_1, w_2 = \frac{h}{2}$   
funzione malissimo perché abbiamo approssimato a una retta

TRAPEZI

$$w_k = \left\{ \frac{h}{2}, \frac{h}{2} \right\}$$

Esempio (3 punti)

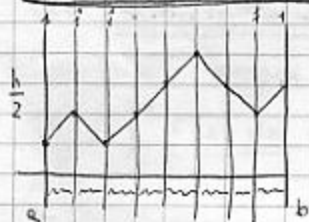


$$I = \frac{h}{3} \left\{ f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right\} \text{ dove } h = \frac{b-a}{2}$$

CAVALIERI-SIMPSON

$$w_k = \left\{ \frac{h}{3}, \frac{4}{3}h, \frac{h}{3} \right\} \text{ dove } h = (b-a)/2$$

Formule quadratiche composte → utilizzando questi esempi si risolvono



$$I = \sum_{i=1}^m I_i = \sum_{i=1}^m \frac{h}{2} (f(a + (i-1)h) + f(a + ih))$$

avremo le medie fra tutti i punti

$$I = \frac{h}{2} [f(a) + 2f(a+h) + 2f(a+2h) + \dots + 2f(a+(m-1)h)]$$

$$= \frac{h}{2} [f(a) + 2 \sum_{i=1}^{m-1} f(a+ih) + f(b)]$$