

Laboratorio di Algoritmi
Numerici di interesse
per l'Ingegneria

Garau Andrea

Matematica Computazionale

Laboratorio di Algoritmi
Numerici di interesse
per l'Ingegneria

Garau Andrea
47017



Italia • Sardegna • Cagliari • via Ospedale



Prof. *Giuseppe Rodriguez*

Studente: *Garau Andrea*

Stampato in Sardegna, Assemini

Copyright © *tutti i diritti sono riservati*

Questo documento è stato composto dall'autore mediante L^AT_EX, su un sistema GNU OS: Arch Linux *x86_64* con Kernel Release 4.10.2 - 1 - *ARCH*. Per la stesione del documento, così come per il codice prodotto, si è usato il programma **free software Emacs** . L'unico software commerciale utilizzato, per i soli esperimenti numerici, è Matlab[®] . I grafici sono stati prodotti con il pacchetto TikZ.

Per maggiori informazioni in merito garau.and@gmail.com

Introduzione



LOBIETTIVO di questo elaborato è quello di analizzare alcuni algoritmi numerici, e simularli praticamente. Questo lavoro si basa principalmente sulle lezioni del corso di *Matematica Computazionale* tenuto dal Prof. Rodriguez a.a. 2016/2017¹. Non c'è nulla che la matematica non possa descrivere. Conoscere il **modello matematico** del fenomeno, un complesso di formule che descrivono il suo comportamento, vuol dire conoscere **tutto**² del fenomeno. Talvolta, diciamo pure il 99.99...9% della volte, le equazioni così ottenute sono troppo complesse per essere risolte, si associa allora un **problema numerico** allo scopo di renderlo risolvibile numericamente su un calcolatore, attraverso una sequenza: **l'algoritmo**.

Definizione 0.0.1: Algoritmo

Una sequenza univoca di un numero finito di operazioni elementari che stabilisce come calcolare la soluzione di un problema, assegnati certi dati iniziali.

Un concetto **fondamentale** è quello di **problema ben posto**. J.Hadamard ne diede una definizione.

Definizione 0.0.2: Problema ben posto

Un problema è **ben posto** se esso possiede, in un prefissato **campo di definizione**, **una e una sola** soluzione e questa dipende con **continuità** dai dati: In caso contrario, viene detto mal posto.

Nota storica 0.0.1: J.Hadamard



Jacques Solomon Hadamard (Versailles, 8 dicembre 1865 – Parigi, 17 ottobre 1963) è stato un matematico francese, conosciuto principalmente per la sua dimostrazione del teorema dei numeri primi. [Wiki16c].

¹Ovviamente solo quelle che ha seguito l'autore!

²Occorrerebbe dare la rigorosa definizione di tutto.

Organizzazione del lavoro

Capitolo 1. Nel primo capitolo si esamineranno gli errori inevitabilmente presenti quando si fanno misure e calcoli. Come vengono memorizzati i numeri reali su un calcolatore e quali sono le principali cause di origine e propagazione degli errori.

Capitolo 2. In questo secondo capitolo si richiamano alcuni concetti di Algebra lineare, delle nozioni che verranno utilizzate nel seguito del lavoro.

Capitolo 3. La fattorizzazione QR.

Capitolo 4 – 5 – 6. Sono rispettivamente: i metodi diretti, autovalori e autovettori e i SVD. In ogni capitolo si svolge la trattazione teorica³ Nel 4 in particolar modo è presente l'Algoritmo di Cholesky.

Capitolo 7. Metodi iterativi, spazi di Krylov, Arnoldi Lanczos e GMRES. Non viene trattata⁴ la parte dei metodi del gradiente.

Capitolo 8. La presentazione sottoforma di grafici di alcuni Algoritmi simulati e confrontati con quelli di Matlab[®] .

Nel testo sono *sparse* delle *Note storiche*⁵ prima di iniziare si è ritenuto doveroso riportare una Nota sul *Princeps mathematicorum*⁶ .

Nota storica 0.0.2: Johann Friedrich Carl Gauss



Johann Friedrich Carl Gauss (Braunschweig, 30 aprile 1777 – Gottinga, 23 febbraio 1855) è stato un matematico, astronomo e fisico tedesco, che ha dato contributi determinanti in analisi matematica, teoria dei numeri, statistica, calcolo numerico, geometria differenziale, geodesia, geofisica, magnetismo, elettrostatica, astronomia e ottica.

Talvolta definito *il Principe dei matematici* (Princeps mathematicorum) come Eulero o *il più grande matematico della modernità* (in opposizione ad Archimede, considerato dallo stesso Gauss come il maggiore fra i matematici dell'*antichità*), è annoverato fra i più importanti matematici della storia avendo contribuito in modo decisivo all'evoluzione delle scienze matematiche, fisiche e naturali. Definì la matematica come *la regina delle scienze*. [Wiki17d].

³Sono gli appunti presi a lezione.

⁴Anche se stata svolta.

⁵Come quella di sopra di J.Hadamard.

⁶Dove il suo nome, anche indirettamente, appare in molti testi di matematica.

Indice

- 1** | Errori
 - 1.1 Numeri macchina 1
 - 1.1.1 Insieme dei numeri macchina 2
 - 1.1.2 Operazioni 3

- 2** | Richiami di Algebra Lineare
 - 2.1 Spazi 5
 - 2.2 Matrici 8
 - 2.2.1 Autovalori e autovettori 9
 - 2.3 Matrici di forma particolare 10

- 3** | Fattorizzazione QR (anni '60)
 - 3.1 Ortogonalizzazione di Gram-Schmidt 14

- 4** | Sistemi sovradeterminati e sottoderminati
 - 4.1 Classificazione dei sistemi 17
 - 4.2 Caso *facile* $m > n = k$ sovradeterminato rango pieno 19
 - 4.2.1 Sistema delle equazioni normali 21
 - 4.3 $k = m < n$ 24

- 5** | Autovalori e autovettori
 - 5.1 Trasformazione di similitudine 28
 - 5.2 Metodo iterativi 30
 - 5.2.1 Metodo delle potenze 30
 - 5.2.2 Eigenvector centrality 33

6 | Singular Value Decomposition

- 6.1 Un problema facile 37
- 6.2 Un altro problemino ... meno facile 37
- 6.3 Pseudo-inversa 38
- 6.4 Usando la fattorizzazione QR 40
- 6.5 Un problema facile 40
- 6.6 Minimi quadrati 41
- 6.7 Sistema normale 41
 - 6.7.1 Caratteristiche degli algoritmi 41

7 | Metodi Iterativi per la risoluzione di sistemi lineari

- 7.1 Metodi iterativi del prim'ordine 43
- 7.2 Costruzione di metodi iterativi lineari 45
- 7.3 Criterio di arresto 46
- 7.4 Precondizionamento 47
- 7.5 Iterazioni in sottospazi di Krylov 49
 - 7.5.1 Il gradiente coniugato come metodo di Krylov 49
 - 7.5.2 L'iterazione di Arnoldi 50
 - 7.5.3 L'iterazione di Lanczos 53
 - 7.5.4 Il metodo GMRES 54

8 | Risultati ottenuti

- 8.1 Le simulazioni fattorizzazione Cholesky 58
- 8.2 Le simulazioni fattorizzazione QR 60
 - 8.2.1 Householder 60
 - 8.2.2 Givens 62
 - 8.2.3 Gram-Schmidt 64
- 8.3 Problemi minimi quadrati, fattorizzazioni QR 66
- 8.4 Metodi iterativi metodo GMRES 69
 - Riferimenti bibliografici 72

Elenco delle figure

4.1	Dominio e codominio.	17
5.1	Rete	33
8.1	Fattorizzazione Cholesky per colonne. Vengono confrontati i tempi di calcolo di due Algoritmi, l'Algoritmo implementato e l'Algoritmo <code>chol()</code> di Matlab [®] . Si risolve il sistema lineare, considerando una matrice piena casuale random, si risolvono così due sistemi uno triangolare inferiore e uno triangolare superiore, Nel grafico in ascisse è riportata la dimensione della matrice da 3 a 200, in ordinate in scala logaritmica i tempi di calcolo in [s].	58
8.2	Fattorizzazione Cholesky. In ascisse sono riportate le dimensioni della matrice A da 3 a 200, in ordinate, in scala logaritmica è riportata la norma dell'errore calcolato come differenza tra soluzione calcolata e soluzione esatta.	59
8.3	Fattorizzazione QR di Householder ($m = n$). Vengono confrontati i tempi di calcolo di due Algoritmi, l'Algoritmo è quello di pagina 102 libro di testo [Rod08] e l'Algoritmo <code>qr()</code> di Matlab [®]	60
8.4	Fattorizzazione QR di Householder ($m = n$). Errori commessi.	61
8.5	Fattorizzazione QR di Givens ($m = n$). Tempi di calcolo.	62
8.6	Fattorizzazione QR di Givens ($m = n$). Errori di calcolo.	63
8.7	Fattorizzazione QR di Gram-Schmidt ($m = n$). Tempi di calcolo.	64
8.8	Fattorizzazione QR di Gram-Schmidt ($m = n$). Errori di calcolo.	65
8.9	Tempo di calcolo della Soluzione Minimi Quadrati $A(m = n)$	66
8.10	Calcolo dell'errore sulla soluzione di un problema ai Minimi Quadrati $A(m = n)$	67
8.11	Tempo di calcolo della Soluzione Minimi Quadrati $A(m = n)$	69
8.12	Calcolo dell'errore sulla soluzione di un problema ai Minimi Quadrati $A(m = n)$	70

Elenco delle tabelle

8.1	Tempo medio di calcolo [s]	68
8.2	Tempo medio di calcolo [s]	71

Capitolo 1

Errori

La presenza di errori nei calcoli può essere dovuta a varie cause:

- il problema in studio potrebbe essere affetto da errori sperimentali
- un errata modellizzazione del fenomeno
- il problema deve essere posto in una forma numerica più semplice per potere essere realizzabile
- la memorizzazione dei dati su un calcolatore digitale introduce errori di arrotondamento.

Occorre poter quantificare l'errore, sia a la quantità da stimare e a^* la sua approssimazione, definiamo le quantità:

Definizione 1.0.1: Errore assoluto

$$\epsilon = |a - a^*|$$

di difficile interpretazione se non si hanno informazioni sull'ordine di grandezza delle quantità da stimare.

Definizione 1.0.2: Errore relativo

$$\rho = \frac{|a - a^*|}{|a|}$$

1.1 Numeri macchina

I numeri di macchina sono quei numeri che possono essere rappresentati esattamente su un calcolatore, essi includono:

- 1 i numeri interi, non troppo grandi

2 i numeri con parte decimale, che possono essere rappresentati mediante due differenti tecniche:

- numeri a virgola fissa riserva t cifre per la parte intera, s per la parte decimale:

$$(\alpha_{t-1}\alpha_{t-2}\dots\alpha_1\alpha_0.\alpha_{-1}\alpha_{-2}\dots\alpha_{-s})_\beta$$

non possono essere espressi numeri maggiori o uguali a β^t o minori di β^{-s}

- numeri a virgola mobile, si memorizzano separatamente le cifre significative e l'ordine di grandezza del numero:

$$(0.\alpha_1\alpha_2\dots\alpha_t)_\beta\beta^p$$

3 per i numeri complessi si memorizzano separatamente parte reale e immaginaria separatamente, in due variabili reali.

1.1.1 Insieme dei numeri macchina

L'insieme dei numeri macchina in base β , con t cifre significative ed esponente nell'intervallo $[L, U]$ è l'insieme:

$$\mathbb{F}(\beta, t, L, U) = \{0\} \cup \{x \in \mathbb{R} : x = \text{sign}(x) \cdot m \cdot \beta^p\}$$

dove

- t e β son interi positivi, con $\beta \geq 2$;
- $\text{sign}(x) = \pm 1$ a seconda del segno di x ;
- la quantità:

$$m = \sum_{i=1}^t d_i \beta^{-i} = (0.d_1d_2\dots d_t)_\beta$$

è la mantissa ($0 \leq d_i \leq \beta - 1$);

- p è un intero compreso tra L e U , detto esponente o caratteristica;
- $d_1 \neq 0$ normalizzazione, non sempre viene applicata, in quel caso si parla di virgola mobile normalizzata.

L'insieme dei numeri macchina \mathbb{F} è un sottoinsieme proprio di \mathbb{R}^1 , dato che l'insieme è discreto e finito, **non** può contenere tutti i numeri reali. Ridefiniamo la funzione

$$\begin{aligned} fl : \mathbb{R} &\longrightarrow \mathbb{F} \\ x &\longmapsto fl(x) \end{aligned}$$

che ad ogni numero reale fa corrispondere un rappresentante in \mathbb{F} . Tale funzione è suriettiva², ma non iniettiva³. Dato $x \in \mathbb{R}$ possiamo avere:

¹Dati due insiemi A e B non vuoti, diremo che A è un sottoinsieme proprio di B se tutti gli elementi dell'insieme A appartengono anche all'insieme B e almeno un elemento dell'insieme B non appartiene all'insieme A : $A \subset B$ cioè è valida l'inclusione stretta.

²Una funzione suriettiva quando ogni elemento del codominio è immagine di **almeno** un elemento del dominio.

³Una funzione iniettiva è una funzione che associa ad elementi distinti del dominio, elementi distinti del codominio.

1. $x \in \mathbb{F}$ Questo vuol dire che l'esponente è nell'intervallo $[L, U]$ e la mantissa no ha più di t cifre significative se rappresentata in base β .
2. $|x| < \beta^{L-1}$ L'esponente p è minore di L . Un numero troppo piccolo che genera **underflow**.
3. $|x| \geq \beta^U$ L'esponente p è maggiore di U . Questa situazione viene detta di **overflow**.
4. $|x| \in [\beta^{L-1}, \beta^U), x \notin \mathbb{F}$ In questo caso il numero di cifre significative di x è superiore a t : occorre scegliere una rappresentazione che approssimi x con sufficiente accuratezza, sono possibili alcune tecniche:

Troncamento le cifre in eccesso vengono trascurate

Arrotondamento questo procedimento equivale ad aggiungere $\frac{1}{2}\beta^{-t}$ alla mantissa di x e successivamente troncare il risultato.

Ha interesse considerare la **precisione macchina**:

Definizione 1.1.1: Precisione macchina

la distanza tra 1 e il successivo numero di macchina.

$$\epsilon_M = \beta^{1-t}$$

1.1.2 Operazioni

Le normali operazioni matematiche $+$, $-$, $*$, $/$ nel calcolatore diventano le **operazioni macchina**: \oplus , \ominus , \otimes , \oslash . L'errore relativo commesso da un'operazione di macchina deve quindi essere dello stesso ordine di grandezza di quello di arrotondamento commesso nel memorizzare il risultato esatto. Quando un sistema in virgola mobile non verifica questa ipotesi si parla di: **aritmetica aberrante**.

Capitolo 2

Richiami di Algebra Lineare



N questo capitolo si concentrano dei risultati importanti utili successivamente nel testo. Si darà prima la definizione di spazio lineare e le sue proprietà, si vedranno gli spazi normati e quelli di Hilbert. Si analizzeranno le matrici e alcune proprietà importanti s esse. Verranno anche studiati gli autovalori e autovettori, e si vedranno alcune matrici che possiedono una particolare struttura. Infine l enorme matriciali e il numero di condizionamento, di relativa importanza per la soluzione di sistemi lineari.

2.1 Spazi

Definizione 2.1.1: Spazio Lineare

Uno spazio lineare o vettoriale reale è un insieme V su cui sono definite due operazioni di somma e prodotto per uno scalare

$$\begin{aligned} + : V \times V &\longrightarrow V \\ (x, y) &\longmapsto x + y \end{aligned}$$

$$\begin{aligned} \cdot : \mathbb{R} \times V &\longrightarrow V \\ (\alpha, x) &\longmapsto \alpha x \end{aligned}$$

che, per ogni $\alpha, \beta \in \mathbb{R}$ e $x, y, z \in V$ godono delle seguenti proprietà:

- 1 $x + y \in V$ (chiusura rispetto alla somma)
- 2 $\alpha x \in V$ (chiusura rispetto al prodotto)
- 3 $y + x = x + y$ (proprietà commutativa)
- 4 $(x + y) + z = x + (y + z)$ (proprietà associativa)
- 5 $\exists 0 \in V : x + 0 = x$ (elemento neutro)
- 6 $\exists -x \in V : x + (-x) = 0$ (inverso additivo)

- 7 $\alpha(\beta\mathbf{x}) = (\alpha\beta)\mathbf{x}$ (proprietà associativa)
- 8 $\alpha(\mathbf{x} + \mathbf{y}) = \alpha\mathbf{x} + \alpha\mathbf{y}$ (proprietà distributiva in V)
- 9 $(\alpha + \beta)\mathbf{x} = \alpha\mathbf{x} + \beta\mathbf{x}$ (proprietà distributiva in \mathbb{R})
- 10 $1\mathbf{x} = \mathbf{x}$ (elemento neutro)

Si chiama **combinazione lineare** dei vettori $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ con coefficienti $\alpha_i \in \mathbb{R}$, $i = 1, \dots, k$ il vettore $\mathbf{x} \in V$ dato da

$$\mathbf{x} = \sum_{i=1}^k \alpha_i \mathbf{x}_i$$

sottospazio generato dai vettori $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ è formato da tutte le loro combinazioni lineari

$$\text{span}(\mathbf{x}_1, \dots, \mathbf{x}_k) := \left\{ \mathbf{x} \in V : \mathbf{x} = \sum_{i=1}^k \alpha_i \mathbf{x}_i \right\}.$$

I vettori $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ sono **linearmente indipendenti** se

$$\sum_{i=1}^k \alpha_i \mathbf{x}_i = \mathbf{0} \Rightarrow \alpha_i = 0, i = 1, \dots, k.$$

Questa condizione equivale al fatto che nessuno dei vettori sia combinazione lineare degli altri. Una **base** è un insieme di vettori linearmente indipendenti tali che ogni vettore dello spazio possa essere espresso come combinazione lineare. La **dimensione** di uno spazio lineare è la cardinalità di una base.

Definizione 2.1.2: Spazio normato

Uno **spazio normato** è uno spazio lineare su cui è definita una funzione

$$\begin{aligned} \|\cdot\| : V &\longrightarrow \mathbb{R} \\ \mathbf{x} &\longmapsto \|\mathbf{x}\| \end{aligned}$$

detta **norma**, tale che per ogni $\mathbf{x}, \mathbf{y} \in V$ e $\alpha \in \mathbb{R}$ si abbia

- 1 $\|\mathbf{x}\| \geq 0$ e $\|\mathbf{x}\| = 0 \iff \mathbf{x} = \mathbf{0}$ (positività)
- 2 $\|\alpha\mathbf{x}\| = |\alpha|\|\mathbf{x}\|$ (omogeneità)
- 3 $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ (disuguaglianza triangolare)

uno spazio normato è anche uno **spazio metrico**, in cui distanza tra due vettori è misurata mediante la funzione:

$$d(\mathbf{x}, \mathbf{y}) := \|\mathbf{x} - \mathbf{y}\|.$$

Le più importanti norme in uno spazio \mathbb{R}^n sono:

$$\text{norma due: } \|\mathbf{x}\|_2 = \left(\sum_{i=1}^n x_i^2 \right)^{\frac{1}{2}}$$

$$\text{norma uno: } \|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$$

$$\text{norma } \infty : \|\mathbf{x}\|_\infty = \max_{i=1, \dots, n} |x_i|$$

Gli spazi funzionali $C[a, b]$ e $L^e[a, b]$ vengono dotati delle norme

$$\|f\|_\infty = \max_{x \in [a, b]} |f(x)|, \quad f \in C[a, b]$$

$$\|g\|_2 = \left(\int_a^b |g(x)|^2 dx \right)^{\frac{1}{2}}, \quad g \in L^2[a, b].$$

Un risultato importante che ci servirà di seguito, si veda 6.5 a pagina 40 è il seguente teorema:

Teorema 2.1.1: Equivalenza delle norme

Tutte le norme in \mathbb{R}^n sono equivalenti, nel senso che per ogni coppia di norme $\|\cdot\|_\alpha, \|\cdot\|_\beta$ esistono due costanti positive m e M tali che, per ogni $\mathbf{x} \in \mathbb{R}^n$, si ha:

$$m\|\mathbf{x}\|_\beta \leq \|\mathbf{x}\|_\alpha \leq M\|\mathbf{x}\|_\beta.$$

Questo teorema consente di studiare indifferentemente qualsiasi norma per la convergenza di una successione di vettori.

Definizione 2.1.3: Spazio di Hilbert

Uno spazio di Hilbert è uno spazio lineare su cui è definito un prodotto scalare (o prodotto interno)

$$\begin{aligned} \langle \cdot, \cdot \rangle : V \times V &\longrightarrow \mathbb{R} \\ (\mathbf{x}, \mathbf{y}) &\longmapsto \langle \cdot, \cdot \rangle \end{aligned}$$

tale che per ogni $\mathbf{x}, \mathbf{y}, \mathbf{z} \in V$ e $\alpha \in \mathbb{R}$

- 1 $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ e $\langle \mathbf{x}, \mathbf{x} \rangle = 0 \iff \mathbf{x} = 0$ (positività)
- 2 $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$ (proprietà commutativa)
- 3 $\langle \alpha \mathbf{x}, \mathbf{x} \rangle = \alpha \langle \mathbf{x}, \mathbf{x} \rangle$ (omogeneità)
- 4 $\langle \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z} \rangle$ (linearità)

Tale spazio deve essere completo rispetto alla norma indotta dal prodotto interno

$$\|\mathbf{x}\| := \langle \mathbf{x}, \mathbf{x} \rangle^{\frac{1}{2}}$$

Due vettori sono ortogonali quando il loro prodotto scalare è nullo. Il prodotto scalare comunemente utilizzato in \mathbb{R}^n è: $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n x_i y_i$.

Nota storica 2.1.1: Hilbert



David Hilbert (Königsberg, 23 gennaio 1862 – Gottinga, 14 febbraio 1943) è stato un matematico tedesco. È stato uno dei più eminenti ed influenti matematici del periodo a cavallo tra il XIX secolo e il XX secolo. [Wiki17f]

2.2 Matrici

Una matrice $m \times n$ è un quadro di mn numeri reali o complessi di m righe e n colonne. Sia $A \in \mathbb{C}^{m \times n}$ la matrice **aggiunta** A^* si ottiene scambiando righe si A con le sue colonne e coniugandone gli elementi, se $A \in \mathbb{R}$ allora si dice **trasposta**. Date due matrici $A \in \mathbb{R}^{m \times n}$ e $B \in \mathbb{R}^{n \times p}$ il prodotto tra matrici è la matrice $C = AB \in \mathbb{R}^{m \times p}$ definita da:

$$c_{ij} = \sum_{k=1}^n a_{ik}b_{kj}, \quad i = 1, \dots, m; j = \dots, p.$$

Una matrice reale $m \times n$ rappresenta la generica **trasformazione lineare** tra gli spazi lineari \mathbb{R}^n e \mathbb{R}^m . Il Prodotto matriciale corrisponde alla composizione di due trasformazioni lineari.

Una matrice quadrata A si dice **invertibile** o **non singolare** se \exists una matrice A^{-1} detta matrice inversa tale che $A^{-1}A = AA^{-1} = I$. Definiamo alcune proprietà delle matrici:

- $(AB)^T = B^T A^T$;
- $(AB)^{-1} = B^{-1} A^{-1}$;
- $(A^T)^{-1} = (A^{-1})^T$;
- A è invertibile se e solo se **ha righe (e colonne) linearmente indipendenti**.

Il **determinante** è una funzione che associa a ciascuna matrice quadrata un numero reale, esso può essere calcolato mediante la formula di Laplace;

$$\det(A) = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det(A_{ij})$$

essendo A_{ij} la sotto-matrice che si ottiene sa A eliminando la i -esima riga e j -sima colonna. alcune proprietà dei determinanti:

- $\det(A^T) = \det(A), \det(A^*) = \overline{\det(A)}, \det(A^{-1}) = \det(A)^{-1}$;
- $\det(AB) = \det(A)\det(B), \det(\alpha A) = \alpha^n \det(A)$;
- uno scambio di due righe produce un cambio di segno del determinante;
- A è non singolare se e solo se $\det(A) \neq 0$

Nota storica 2.2.1:

Laplace



Pierre-Simon Laplace, marchese di Laplace (Beaumont-en-Auge, 23 marzo 1749 – Parigi, 5 marzo 1827), è stato un matematico, fisico, astronomo e nobile francese. Fu uno dei principali scienziati nel periodo napoleonico. [\[Wiki17\]](#)

Definizione 2.2.1: Rango

Il **rango** di una matrice $\text{rank}(A)$ può essere definito come il massimo numero di righe (o colonne) linearmente indipendenti o come l'ordine della più grande sotto-matrice quadrata con determinante non nullo (non singolare).

2.2.1 Autovalori e autovettori

Definizione 2.2.2: Autovalore e autovettore

Si dicono **autovalore** ed **autovettore** di una matrice A uno scalare λ ed un vettore $\mathbf{x} \neq 0$ che verificano la relazione

$$A\mathbf{x} = \lambda\mathbf{x} \quad (2.1)$$

La 2.1 può essere riscritta:

$$(A - \lambda I)\mathbf{x} = 0.$$

Questo è un sistema lineare omogeneo dove il vettore $\mathbf{x} \neq 0$, pertanto per ammettere una soluzione non nulla deve annullarsi il determinante della matrice $(A - \lambda I)$, questo risulta essere un polinomio di grado n in λ , dove n è il numero di righe della matrice A quadrata. Tale polinomio è noto come **polinomio caratteristico**:

$$p_A(\lambda) = \det(A - \lambda I)$$

quindi calcolare gli zeri di $p_A(\lambda)$ vuol dire calcolare gli n autovalori di A .

La matrice $(A - \lambda I)$ ha **sempre** $\det(A - \lambda I) \neq 0$ anche se A è singolare.

Per ciascun autovalore $\lambda_k, k = 1, \dots, n$, una soluzione non nulla del **sistema singolare** omogeneo¹

$$(A - \lambda_k I)\mathbf{x} = 0$$

fornisce il corrispondente autovettore, il rango della matrice $(A - \lambda_k I) \leq n - 1$, dunque l'autovettore rimane determinato a meno di una costante.

Definizione 2.2.3: Spettro

Lo spettro di una matrice è l'insieme degli autovalori

$$\sigma(A) = \{\lambda_1, \dots, \lambda_n\} \quad (2.2)$$

Definizione 2.2.4: Raggio spettrale

Il raggio spettrale è il massimo dei moduli degli autovalori

$$\rho(A) = \max_{k=1, \dots, n} |\lambda_k|. \quad (2.3)$$

La **molteplicità algebrica** di un autovalore è la sua molteplicità come zero del polinomio caratteristico. Si definisce, **molteplicità geometrica** di un autovalore il massimo numero di autovettori linearmente indipendenti ad esso corrispondenti, in generale per ogni autovalore si ha sempre:

$$\text{molteplicità geometrica} \leq \text{molteplicità algebrica}$$

Se per un autovalore vale il minore stretto, la matrice viene detta **difettiva**.

Proprietà:

¹Il sistema è singolare $\det(A - \lambda_k I) = 0$, stiamo cercando soluzioni non banali del vettore \mathbf{x} .

- $\det(\mathbf{A}) = \prod_{k=1}^n \lambda_k$;
- $\sigma(\mathbf{A}^T) = \sigma(\mathbf{A}), \sigma(\mathbf{A}^{-1}) = \{\lambda_1^{-1}, \dots, \lambda_n^{-1}\}, \sigma(\mathbf{A}^p) = \{\lambda_1^p, \dots, \lambda_n^p\}$;
- **ad autovalori distinti corrispondono autovettori indipendenti**;
- se un autovettore \mathbf{x} è noto il **quoziente di Rayleigh**: $\frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$ fornisce il corrispondente autovalore.

2.3 Matrici di forma particolare

Alcune matrici posseggono una **struttura** quando hanno delle proprietà che rendono più agevole la risoluzione di un problema che le coinvolge, ecco alcune di esse.

Hermitiane

Una matrice è Hermitiana se coincide con la sua aggiunta, se la matrice è reale si dice simmetrica. Una matrice Hermitiana \mathbf{A} è definita positiva **se**:

$$\mathbf{x}^* \mathbf{A} \mathbf{x} > 0, \quad \forall \mathbf{x} \in \mathbb{C}^n$$

valida per matrici simmetriche reali

Nota storica 2.3.1: Hermite



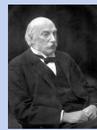
Charles Hermite (Dieuze, 24 dicembre 1822 – Parigi, 14 gennaio 1901) è stato un matematico francese che diede rilevanti contributi a campi quali teoria dei numeri, forme quadratiche, teoria degli invarianti, polinomi ortogonali, funzioni ellittiche e algebra. Egli fu il primo a dimostrare che la costante e , è un numero trascendente. [Wiki16d].

Teorema 2.3.1: Matrice Hermitiana

Se \mathbf{A} è Hermitiana, i suoi autovalori sono reali ed esiste per \mathbb{C}^n una base di autovettori ortogonali. Se \mathbf{A} è anche definita positiva, gli autovalori sono positivi.

Unitarie

Nota storica 2.3.2: Rayleigh



John William Strutt 3° barone di Rayleigh (Langford Grove, 12 novembre 1842 – Witham, 30 giugno 1919) è stato un fisico britannico. [Wiki16e].

Una matrice complessa è unitaria se la sua inversa coincide con l'aggiunta:

$$\mathbf{Q} \mathbf{Q}^* = \mathbf{Q}^* \mathbf{Q} = \mathbf{I}$$

Una matrice reale è **ortogonale** se l'inverso coincide con la trasposta:

$$\mathbf{Q}^T \mathbf{Q} = \mathbf{Q} \mathbf{Q}^T = \mathbf{I}$$

Teorema 2.3.2

Se \mathbf{Q} è unitaria, allora:

- 1 $|\det(\mathbf{Q})| = 1$, se $\mathbf{Q} \in \mathbb{R} \Rightarrow \det(\mathbf{Q}) = \pm 1$;
- 2 $\|\mathbf{Q} \mathbf{x}\|_2 = \|\mathbf{x}\|_2$, per ogni $\mathbf{x} \in \mathbb{C}^n$.

Triangolari

Una matrice U si dice triangolare superiore se i suoi elementi verificano $u_{ij} = 0$, per $i > j$. Una matrice L si dice triangolare inferiore se i suoi elementi verificano $l_{ij} = 0$, per $i < j$. Una matrice D si dice diagonale se è simultaneamente triangolare superiore e inferiore.

Teorema 2.3.3

Se T è triangolare il suo determinante è dato dal prodotto degli elementi diagonali. Gli autovalori coincidono con gli elementi diagonali.

Le matrici triangolari formano un'algebra: l'inversa di una matrice triangolare, quando esiste, e il prodotto di due matrici triangolari sono ancora matrici triangolari dello stesso tipo. La stessa proprietà vale per matrici unitarie, ma questo caso non si può parlare di algebra dato che l'unitarietà non si trasmette attraverso la somma di matrici ed il prodotto per uno scalare.

Banda

Gli elementi di una matrice a banda verificano la relazione:

$$a_{ij} = 0, \text{ se } i - j \geq k \text{ oppure } i - j \leq -m$$

Sparse

Le matrici sparse sono dotate di pochi elementi diversi da zero, meno del 10%. Pertanto si possono memorizzare solo gli elementi non nulli e della loro posizione all'interno della matrice.

Norme matriciali

Lo spazio lineare $\mathcal{M}_{m \times n}$ cui appartengono le matrici, può essere dotato della struttura di spazio normato definendo una qualsiasi norma $\|\cdot\|$ che verifichi gli assiomi richiesti.

Definizione 2.3.1: Norma indotta

Una norma matriciale si dice **indotta** da una norma vettoriale $\|\cdot\|$, se:

$$\|\mathbf{A}\| = \sup_{\mathbf{x} \neq 0} \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|}$$

o se, equivalentemente,

$$\|\mathbf{A}\| = \max_{\|\mathbf{x}\|=1} \|\mathbf{Ax}\|.$$

Una norma così definita viene anche detta **naturale**, o subordinata.

Seguono alcuni teoremi sulle norme matriciali

Teorema 2.3.4: Norma matriciale indotta dalla norma vettoriale ∞

$$\|\mathbf{A}\|_{\infty} = \max_{i=1, \dots, m} \sum_{j=1}^n |a_{ij}|.$$

Teorema 2.3.5: Norma matriciale indotta dalla norma 1

$$\|\mathbf{A}\|_1 = \max_{i=1, \dots, n} \sum_{j=1}^m |a_{ij}|.$$

Teorema 2.3.6: Norma matriciale indotta dalla norma 2

$$\|\mathbf{A}\|_2 = \sqrt{\rho(\mathbf{A}^* \mathbf{A})}.$$

Né consegue che se \mathbf{A} è Hermitiana:

$$\|\mathbf{A}\|_2 = \rho(\mathbf{A}).$$

Numero di condizionamento

Un parametro di fondamentale importanza, per la soluzione di un sistema lineare è il **numero di condizionamento**.

Definizione 2.3.2: Numero di condizionamento

Il numero di condizionamento di una matrice relativamente alla risoluzione di un sistema lineare, la quantità:

$$k(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|. \quad (2.4)$$

Misura il massimo fattore di amplificazione dell'errore relativo sulla soluzione rispetto all'errore relativo sui dati. Esso è influenzato dalla norma matriciale adottata.

Verifica le seguenti proprietà:

- 1 per ogni matrice \mathbf{A} , $k(\mathbf{A}) \geq 1$;
- 2 per ogni norma naturale, $k(\mathbf{I}) = 1$;
- 3 $k(\mathbf{AB}) \leq k(\mathbf{A})k(\mathbf{B})$;
- 4 se \mathbf{Q} è ortogonale $k_2(\mathbf{Q}) = 1$;
- 5 $k_2(\mathbf{A}) = \sqrt{\frac{\lambda_{max}(\mathbf{A}^T \mathbf{A})}{\lambda_{min}(\mathbf{A}^T \mathbf{A})}}$;
- 6 se \mathbf{A} è simmetrica $k_1(\mathbf{A}) = \frac{|\lambda_{max}(\mathbf{A})|}{|\lambda_{min}(\mathbf{A})|}$.

Capitolo 3

Fattorizzazione QR (anni '60)

Prima di introdurre gli algoritmi dei *metodi iterativi* occorre studiare la fattorizzazione QR.

Teorema 3.0.1: Fattorizzazione QR

Data una matrice **qualsiasi**: $A_{m \times n}$, $\exists Q : Q^T Q = Q Q^T = I$ (Ortagonale), e $\exists R$ triangolare superiore : $A = QR$

La complessità computazionale di questa fattorizzazione è $O\left(\frac{2}{3}n^3\right)^1$.

$$\begin{bmatrix} A_{(m \times n)} \end{bmatrix} = \begin{bmatrix} Q_{(m \times m)} \end{bmatrix} \begin{bmatrix} * & * & * \\ & * & * \\ & & * \\ R_{(m \times n)} \end{bmatrix}$$

La matrice R essendo triangolare superiore può essere espressa come: $R_{(m \times n)} = \begin{bmatrix} R_1_{(n \times n)} \\ O_{(m-n)} \end{bmatrix}$.

Le matrici ortogonali sono *comode* quando si ha a che fare con la norma 2, infatti:

$$\|Qx\|^2 = (Qx)^T (Qx) = x^T Q^T Q x = x^T x = \|x\|^2. \quad (3.1)$$

Quando A è una matrice quadrata non singolare, la fattorizzazione QR può essere utilizzata per la risoluzione del sistema lineare $Ax = b$.

Sostituendo ad A la **sua** QR si ottiene:

$$\begin{cases} Qc = b \\ Rx = c. \end{cases} \quad (3.2)$$

Grazie alla 3.1 questa fattorizzazione ha un notevole vanataggio sulla crescita del condizionamento.

¹Peggio di Gauss.

Teorema 3.0.2: Condizionamento QR

Se $A = QR$, oppure $A = RQ$ allora $\|A\|_2 = \|R\|_2$

Dimostrazione. Valendo la 3.1 da $A = QR$ otteniamo:

$$\|A\|_2 \leq \|Q\|_2 \|R\|_2 = \|R\|_2.$$

In maniera analoga partendo da $R = Q^T A$ e dunque:

$$\|R\|_2 \leq \|Q\|_2 \|A\|_2 = \|A\|_2.$$

Le due relazioni implicano la tesi. □

Teorema 3.0.3: Numero di condizionamento

La matrice R ha lo stesso numero di condizionamento rispetto alla norma-2 della matrice A .

Dimostrazione. Poiché $A^{-1} = R^{-1}Q^T$, applicando il teorema 3.0.2 si ottiene:

$$\kappa_2(A) = \|A\|_2 \|A^{-1}\|_2 = \|R\|_2 \|R^{-1}\|_2 = \kappa_2(R).$$

□

Si sono analizzate diverse fattorizzazioni QR: Householder, Givens e Gram-Schmidt. Sono entrambe fondamentali e ampiamente utilizzate nelle simulazioni sperimentali svolte, ma per non dilungarmi eccessivamente sulla trattazione ci si focalizzerà, per ora, solo su Gram-Schmidt. Perché sono la **base** degli algoritmi iterativi, che sono stati usati anche in *altri lavori*. **ortogonalizzazione**. Si è **sempre** fatto riferimento al libro di testo [Rod08].

3.1 Ortogonalizzazione di Gram-Schmidt

Il processo di Gram-Schmidt consente di ortogonalizzare n vettori di \mathbb{R}^m con ($m \geq n$). Iniziamo pensando di poter scrivere la matrice A come composta da vettori colonna:

$$A = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]$$

allora questo procedimento permette di scrivere:

$$Q = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n]$$

tale che

$$\text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_n\} = \text{span}\{\mathbf{a}_1, \dots, \mathbf{a}_n\} \quad \text{e} \quad \mathbf{q}_i^T \mathbf{q}_j = \delta_{ij}.$$

Nota storica 3.1.1: Gram-Schmidt



Jørgen Pedersen Gram (Nustrup, 27 giugno 1850 – Copenhagen, 29 aprile 1916) è stato un matematico ed attuario danese. Tra i suoi studi si ricordano le espansioni in serie determinate dai metodi dei minimi quadrati, i numeri primi minori di un

dato numero e le serie per la funzione zeta di Riemann. Il processo di ortogonalizzazione che porta il suo nome assieme a quello di Erhard Schmidt venne pubblicato per la prima volta nel 1883. Morì a 66 anni dopo essere stato investito da una bicicletta. [Wiki17e]

Erhard Schmidt (Dorpat, 13 gennaio 1876 – Berlino, 6 dicembre 1959) è stato un matematico tedesco. Conseguì il dottorato presso l'Università Georg-August di Gottinga nel 1905 sotto la supervisione di David Hilbert, con cui fornì importanti contributi per l'analisi funzionale. [Wiki17k]



Al primo passo si pone:

$$\tilde{\mathbf{q}}_1 = \mathbf{a}_1, \quad r_{11} = \|\tilde{\mathbf{q}}_1\|_2 \quad \mathbf{q}_1 = \frac{\tilde{\mathbf{q}}_1}{r_{11}}$$

Al generico passo k si calcola

$$\tilde{\mathbf{q}}_k = \mathbf{a}_k - \sum_{j=1}^{k-1} r_{jk} \mathbf{q}_j \quad (3.3)$$

con

$$r_{jk} = \mathbf{q}_j^T \mathbf{a}_k \quad j = 1, 2, \dots, k-1,$$

ottenendo

$$\mathbf{q}_k = \frac{\tilde{\mathbf{q}}_k}{r_{kk}}, \quad \text{per } r_{kk} = \|\tilde{\mathbf{q}}_k\|.$$

La scelta effettuata per le costanti r_{jk} permette che ogni vettore \mathbf{q}_k risulti ortogonale ai precedenti.

Applicando l'algoritmo, pedissequamente, per esempio al passo $k = 3$ si ha.

$$k = 3 \quad \tilde{\mathbf{q}}_3 = \mathbf{a}_3 - r_{13} \mathbf{q}_1 - r_{12} \mathbf{q}_2 \quad (3.4)$$

$$r_{13} = \mathbf{q}_1^T \mathbf{a}_3$$

$$r_{12} = \mathbf{q}_1^T \mathbf{a}_2$$

$$\mathbf{q}_3 = \frac{\tilde{\mathbf{q}}_3}{r_{33}} \quad (3.5)$$

e unendo le 3.4 e 3.5 si ottiene:

$$\mathbf{q}_3 = \frac{\mathbf{a}_3 - r_{13} \mathbf{q}_1 - r_{12} \mathbf{q}_2}{r_{33}} \Rightarrow \mathbf{a}_3 = r_{13} \mathbf{q}_1 + r_{12} \mathbf{q}_2 + r_{33} \mathbf{q}_3,$$

che possiamo riscrivere come:

$$\mathbf{a}_k = \sum_{j=1}^k r_{jk} \mathbf{q}_j, \quad k = 1, \dots, n \quad (3.6)$$

e dunque la matrice \mathbf{A} può essere rappresentata come:

$$\mathbf{A} = [\mathbf{q}_1 \dots \mathbf{q}_k \dots \mathbf{q}_n] \begin{bmatrix} r_{11} & \dots & r_{1k} & \dots & r_{1n} \\ & \ddots & \vdots & & \vdots \\ & & r_{kk} & & \vdots \\ & & & \ddots & \vdots \\ & & & & r_{nn} \end{bmatrix}.$$

L'ortogonalizzazione delle colonne della matrice \mathbf{A} costruisce la fattorizzazione \mathbf{QR} . Questo procedimento è noto come **CGS** (*Classical Gram-Schmidt*). La matrice \mathbf{R} è costruita per colonne e quando $n < m$ la matrice \mathbf{Q} pur avendo colonne ortonormali, **non** è ortogonale, non essendo quadrata. Può essere completata per ottenere una matrice ortogonale aggiungendo $m - n$ vettori che completino la base ortonormale.

È stata sudita una una differente versione dell'algoritmo, perché la colonne di \mathbf{Q} tendono a perdere la reciproca ortogonalità al procedere delle iterazioni, rendendo numericamente instabile l'algoritmo.

La versione prende il nome di **MGS** (*Modified Gram-Schmidt*), genera una successione di $n + 1$ matrici tali che $\mathbf{A}^{(1)} = \mathbf{A}$ e $\mathbf{A}^{(n+1)} = \mathbf{Q} \in \mathbb{R}^{m \times n}$.

La k -esima matrice assume la forma

$$\mathbf{A}^{(k)} = [\mathbf{q}_1, \dots, \mathbf{q}_{k-1}, \mathbf{a}_k^{(k)}, \dots, \mathbf{a}_n^{(k)}]$$

con

$$\begin{cases} \mathbf{q}_i^T \mathbf{q}_j = \delta_{ij}, & i, j = 1, \dots, k-1, \\ \mathbf{q}_j^T \mathbf{a}_j^{(k)} = 0, & i = 1, \dots, k-1; j = k, \dots, n. \end{cases}$$

Al passo k , la matrice $\mathbf{A}^{(k+1)}$ viene ottenuta attraverso la seguente trasformazione:

$$[\mathbf{q}_1, \dots, \mathbf{q}_{k-1}, \mathbf{a}_k^{(k)}, \dots, \mathbf{a}_n^{(k)}] \rightarrow [\mathbf{q}_1, \dots, \mathbf{q}_k, \mathbf{a}_{k+1}^{(k+1)}, \dots, \mathbf{a}_n^{(k+1)}]$$

ponendo

$$\begin{aligned} r_{kk} &= \|\mathbf{a}_k^{(k)}\|, \\ \mathbf{q}_k &= \frac{\mathbf{a}_k^{(k)}}{r_{kk}}, \\ r_{kj} &= \mathbf{q}_k^T \mathbf{a}_j^{(k)}, \quad j = k+1, \dots, n, \\ \mathbf{a}_j^{(k+1)} &= \mathbf{a}_j^{(k)} - r_{jk} \mathbf{q}_k, \quad j = k+1, \dots, n. \end{aligned}$$

Si ottengono due famiglie di vettori:

- 1 $\{\mathbf{q}_1, \dots, \mathbf{q}_k\}$ formata da vettori tra loro ortonormali,
- 2 $\{\mathbf{a}_{k+1}^{(k+1)}, \dots, \mathbf{a}_n^{(k+1)}\}$ gli elementi sono ortogonali a tutti i vettori della prima famiglia.

All' n -esima iterazione la prima famiglia conterrà n vettori tra loro ortonormali, la seconda sarà vuota. L'algoritmo termina n passi, la matrice \mathbf{R} viene costruita per righe, anzichè per colonne. La complessità computazionale è $O(mn^2)$. Anche in **MGS** può verificarsi una perdita di ortogonalità nelle colonne di \mathbf{A} .

Capitolo 4

Sistemi sovradeterminati e sottodeterminati

4.1 Classificazione dei sistemi

PONIAMOCI il seguente problema:

$$\begin{aligned} \mathbf{Ax} = \mathbf{b}, \quad \mathbf{A} \in \mathbb{R}^{n \times n}, \quad \mathbf{b}, \mathbf{x} \in \mathbb{R}^n, \\ k = \text{rank}(\mathbf{A}) = n \end{aligned}$$

esso ammette soluzione, se e solo se, il vettore \mathbf{b} appartiene allo spazio lineare generato dalle colonne di \mathbf{A} , si veda la figura 4.1.

Partendo da uno spazio vettoriale, il **dominio**, di dimensione n , attraverso l'operatore \mathbf{A} , cerchiamo le \mathbf{x} nel **codominio**, spazio vettoriale di dimensione \mathbb{R}^n .

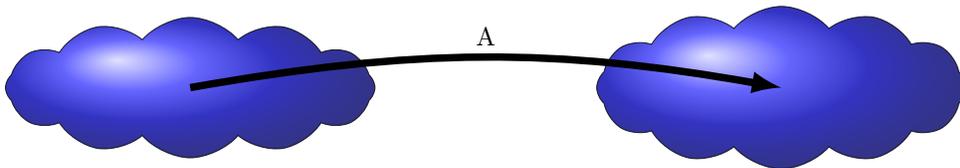


Figura 4.1: Dominio e codominio.

Nel caso in cui il dominio o il codominio non hanno più dimensione \mathbb{R}^n , la soluzione potrebbe non esistere o esistere infinite. Inoltre il problema dipende dal rango della matrice \mathbf{A} che potrebbe non essere pieno ovvero:

$$k = \text{rank}(\mathbf{A}) < \min(m, n).$$

$k < n < m$ caso sovradeterminato

Il rango della matrice \mathbf{A} indica quante colonne sono linearmente dipendenti, se riesco a conoscere quali sono, posso pensare di eliminarle, per farlo pongo a 0 le corrispondenti componenti di \mathbf{x} ,

ma questa operazione potrebbe non essere semplice. Per esempio una matrice A ha due colonne linearmente dipendenti¹, il vettore x risulta allora:

$$\begin{bmatrix} \spadesuit & \spadesuit & \spadesuit & * & * \\ \spadesuit & \spadesuit & \spadesuit & * & * \\ \spadesuit & \spadesuit & \spadesuit & * & * \\ \spadesuit & \spadesuit & \spadesuit & * & * \\ \spadesuit & \spadesuit & \spadesuit & * & * \\ \spadesuit & \spadesuit & \spadesuit & * & * \\ \spadesuit & \spadesuit & \spadesuit & * & * \end{bmatrix} \begin{bmatrix} \spadesuit \\ \spadesuit \\ \spadesuit \\ \spadesuit \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} \spadesuit \\ \spadesuit \\ \spadesuit \\ \spadesuit \\ \spadesuit \\ \spadesuit \\ \spadesuit \end{bmatrix}$$

$k < m < n$ caso sottodeterminato

Posso avere ∞ soluzioni.

Esempio: dato il sistema lineare:

$$\begin{cases} 2x + y + z = 4 \\ 4x + 2y + 2z = 8 \end{cases} \Rightarrow \begin{bmatrix} 2 & 1 & 1 \\ 4 & 2 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 4 \\ 8 \end{bmatrix}$$

le due equazioni sono linearmente dipendenti, dato che la seconda è il doppio della prima, quindi abbiamo ∞ soluzioni, ma il rango vale 1 dunque sono **ancora più** ∞ soluzioni. Ipotizziamo ora di cambiare il secondo membro della seconda equazione²:

$$\underbrace{4x + 2y + 2z}_{\text{l.d.}} = \underbrace{5}_{\text{l.i.}}$$

si divide i due parti: una linearmente dipendente e una linearmente indipendente, si è dunque passati da ∞ soluzioni a **nessuna soluzione**.

$k < m = n$

Il sistema è quadrato, ma non a rango pieno allora si possono verificare i casi:

- ① nessuna soluzione, poiché si arriva ad una palese incongruenza:

$$\begin{cases} x + 2y = 1 \\ 2x + 4y = 3 \end{cases} \Rightarrow \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix};$$

- ② ∞ soluzioni

$$\begin{cases} x + 2y = 1 \\ 2x + 4y = 2 \end{cases} \Rightarrow \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}.$$

¹Le ultime due colonne in *.

²Ora non sono più una il doppio dell'altra.

Allora possiamo pensare di procedere lavorando con una matrice a rango ridotto: k . Questo tipo di problema è noto come **rank deficient**. Questo porta ad avere k autovalori: $\lambda_1, \lambda_2, \dots, \lambda_k$, gli altri sono **nulli**³, ma facendo **calcolo numerico** essendoci errori, potrebbe non essere 0 \Rightarrow il concetto di rango dipende dalla precisione:

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k > \epsilon > \lambda_{k+1} \geq \dots \geq \lambda_n$$

quindi il problema si sposta su **quanto** posso scegliere questo errore, in alcuni casi potrebbe essere semplice, in altri impossibile.

Definizione 4.1.1: range di una matrice o immagine

Il range rappresenta il codominio della funzione:

$$R(\mathbf{A}) = \left\{ \mathbf{y} \in \mathbb{R}^m : \mathbf{y} = \sum_{j=1}^n x_j \mathbf{a}_j \right\} \subset \mathbb{R}^m$$

Definizione 4.1.2: nucleo di una matrice

$$N(\mathbf{A}) = \{ \mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{0} \} \subset \mathbb{R}^n$$

entrambi sono sottospazi.

4.2 Caso *facile* $m > n = k$ sovradeterminato rango pieno

Dato il sistema lineare

$$\mathbf{A}_{(n \times m)} \mathbf{x}_{(m \times 1)} = \mathbf{b}_{n \times 1} \quad (4.1)$$

Il caso è sovraderminato a rango pieno: $rank(\mathbf{A}) = n$. Ovvero \nexists soluzione. Se non posso rendere *uguali* i vettori \mathbf{Ax} e \mathbf{b} , allora cerco di renderli *quasi uguali*, per far ciò utilizzo la norma, la norma è *qualcosa* che mi dice quanto lungo è un vettore⁴.

Se i vettori non possono essere uguali, né cerco la minima distanza. Ovvero cerco di risolvere un problema del tipo:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Ax} - \mathbf{b}\|_2 \quad (4.2)$$

il pedice 2 indica che è la norma euclidea. La norma euclidea è così definita:

$$\begin{aligned} \|\mathbf{x}\| &= \sqrt{\sum_{i=1}^n x_i^2} = \sqrt{\mathbf{x}^T \mathbf{x}} \Rightarrow \\ \|\mathbf{Ax} - \mathbf{b}\| &= \sqrt{\sum_{i=1}^m [(\mathbf{Ax})_i - b_i]^2} \end{aligned} \quad (4.3)$$

³O comunque molto piccoli.

⁴Infatti è uno scalare.

la radice è sempre scomoda, eliminiamola facendo il quadrato, minimizzeremo appunto la norma al quadrato. Ricordando $(\mathbf{Ax})_i = \sum_{j=1}^n a_j x_j$ sostituendo ed eliminando la radice:

$$\|\mathbf{Ax} - \mathbf{b}\|^2 = \sum_{i=1}^m \left[\sum_{j=1}^n a_j x_j - b_i \right]^2 = \varphi(x_1, x_2, \dots, x_n)$$

Studiando geometria di questa funzione abbiamo: in 1 variabile è una parabola, in 2 variabili è un paraboloide, sicuramente φ ha un solo minimo. Da un problema ad ∞ soluzioni passo a un problema a una ed una sola soluzione è dunque ben posto, questo tipo di problemi prende il nome: LSP (Least Square Problem). Come lo risolviamo? Per trovare un minimo occorre annullare la derivata, pertanto deriviamo i singoli termini.

Richiami

Ricordando alcune formule che saranno utili di seguito, consideriamo la seguente funzione $f(\mathbf{x}) = \mathbf{b}^T \mathbf{x}$, con \mathbf{b} vettore di costanti, \mathbf{x} vettore di variabili, possiamo certamente scrivere

$$\mathbf{b}^T \mathbf{x} = \sum_{i=1}^n b_i x_i = \mathbf{x}^T \mathbf{b}$$

se deriviamo questa espressione otteniamo

$$\begin{aligned} \frac{\partial}{\partial x_k} \sum_{i=1}^n b_i x_i &= b_k \Rightarrow \\ \nabla \mathbf{b}^T \mathbf{x} &= \mathbf{b} = \nabla \mathbf{x}^T \mathbf{b} \end{aligned} \quad (4.4)$$

dove l'operatore ∇ indica le derivate parziali per ogni variabile x_i .

Consideriamo ora questa matrice quadrata: $\mathbf{x}^T \mathbf{Ax}$, ricordando che il prodotto matrice vettore è **non** commutativo, ma associativo possiamo *mettere* le parentesi come meglio riteniamo e sviluppare:

$$\begin{aligned} \mathbf{x}^T (\mathbf{Ax}) &= \sum_{i=1}^n x_i (\mathbf{Ax})_i = \\ &= \sum_{i=1}^n x_i \left(\sum_{j=1}^n a_{ij} x_j \right) = \text{tutti scalari} \Rightarrow \\ &= \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j. \end{aligned} \quad (4.5)$$

Deriviamo la 4.5

$$\begin{aligned} \frac{\partial}{\partial x_k} \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j &= \sum_{j=1}^n a_{kj} x_j + \sum_{i=1}^n a_{ik} x_i = \\ &= (\mathbf{Ax})_k + (\mathbf{A}^T \mathbf{x})_k \end{aligned} \quad (4.6)$$

i risultati ottenuti di nostro interesse sono la 4.2 4.6 che riportiamo. Che per una matrice simmetrica diventano:

$$\begin{aligned} \mathbf{A} &= \mathbf{A}^T \\ \nabla(\mathbf{x}^T \mathbf{A} \mathbf{x}) &= \mathbf{A} \mathbf{x} + \mathbf{A}^T \mathbf{x} = 2\mathbf{A} \mathbf{x} \end{aligned} \quad (4.7)$$

$$\nabla \mathbf{b}^T \mathbf{x} = \nabla \mathbf{x}^T \mathbf{b} = \mathbf{b}. \quad (4.8)$$

4.2.1 Sistema delle equazioni normali

Ricordando la 4.3 eseguiamo il quadrato della norma:

$$\begin{aligned} \|\mathbf{A} \mathbf{x} - \mathbf{b}\|^2 &= (\mathbf{A} \mathbf{x} - \mathbf{b})^T (\mathbf{A} \mathbf{x} - \mathbf{b}) = [(\mathbf{A} \mathbf{x})^T - \mathbf{b}^T] [\mathbf{A} \mathbf{x} - \mathbf{b}] = \\ &= [\mathbf{x}^T \mathbf{A}^T - \mathbf{b}^T] [\mathbf{A} \mathbf{x} - \mathbf{b}] = \mathbf{x}^T \underbrace{\mathbf{A}^T \mathbf{A}}_{\text{simm.}} \mathbf{x} - \underbrace{\mathbf{x}^T \mathbf{A}^T \mathbf{b} - \mathbf{b}^T \mathbf{A} \mathbf{x}}_{\text{sonouguale}} + \underbrace{\mathbf{b}^T \mathbf{b}}_{\text{costante}} = \\ &= \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} - 2\mathbf{x}^T \mathbf{A}^T \mathbf{b} + \mathbf{b}^T \mathbf{b} = \varphi(\mathbf{x}). \end{aligned}$$

$\mathbf{x}^T \mathbf{A}^T \mathbf{b}$ e $\mathbf{b}^T \mathbf{A} \mathbf{x}$ sono uguali perché: $(\mathbf{b}^T \mathbf{A} \mathbf{x})^T = \mathbf{x}^T \mathbf{A}^T \mathbf{b}$, essendo scalari, il trasposto di uno scalare è lo scalare stesso. Ora annulliamo e deriviamo $\varphi(\mathbf{x})$, ricordando le regole 4.2 e 4.6:

$$\begin{aligned} \nabla \varphi(\mathbf{x}) &= 0 \\ 2\mathbf{A}^T \mathbf{A} \mathbf{x} - 2\mathbf{A}^T \mathbf{b} &= 0 \Rightarrow \\ \mathbf{A}^T \mathbf{A} \mathbf{x} &= \mathbf{A}^T \mathbf{b} \end{aligned} \quad (4.9)$$

siamo giunti al Sistema delle equazioni normali (ortogonali) la 4.9 che gode delle seguenti:

- 1 è un sistema lineare
- 2 quadrato, infatti:
 - ▬ $\mathbf{A}, \in \mathbb{R}^{m \times n}$
 - ▬ $\mathbf{A}^T, \in \mathbb{R}^{n \times m}$
 - ▬ $\mathbf{A}^T \mathbf{A}, \in \mathbb{R}^{n \times n}$
 - ▬ $\mathbf{b}, \in \mathbb{R}^m$
 - ▬ $\mathbf{A}^T \mathbf{b}, \in \mathbb{R}^n$
- 3 e la matrice $\mathbf{A}^T \mathbf{A}$ ha $k = \text{rank}(\mathbf{A}^T \mathbf{A}) = \min(n, n) = n \Rightarrow$ è non singolare

La matrice $\mathbf{A}^T \mathbf{A}$ è:

- 1 quadrata $n \times n$
- 2 simmetrica \Rightarrow gli autovalori $\in \mathbb{R}$.

- 3 definita positiva, ovvero: $\mathbf{x}^T(\mathbf{A}^T\mathbf{A})\mathbf{x} > 0, \forall \mathbf{x} \neq 0$ questo ci permette di affermare che gli autovalori sono maggiori di 0, e che la matrice può essere invertita.

Possiamo risolvere la 4.9, come un normale sistema, soluzione delle equazioni normali: Least Square Solution (LSS):

$$\mathbf{x} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{b}. \quad (4.10)$$

La matrice $\mathbf{A}^\dagger = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T$ viene detta **pseudo-inversa** è la matrice risolvente del problema ai minimi quadrati, è inversa *solo* a sinistra:

$$\begin{aligned} \mathbf{A}^\dagger\mathbf{A} &= (\mathbf{A}^T\mathbf{A})^{-1}(\mathbf{A}^T\mathbf{A}) = \mathbf{I} \\ \mathbf{A}\mathbf{A}^\dagger &= \mathbf{A}(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T \neq \mathbf{I} \end{aligned}$$

trovare la soluzione come nella 4.10 non è in generale facile numericamente, si preferisce scomporre la matrice, o meglio fattorizzarla, per cercare di ottenere sistemi più facili.

Fattorizzazione Cholesky

Teorema 4.2.1: Cholesky

Se \mathbf{B} è simmetrica definita positiva $\Rightarrow \exists$ una unica matrice \mathbf{L} triangolare inferiore : $\mathbf{B} = \mathbf{L}\mathbf{L}^T$. Questa è nota come *fattorizzazione di Cholesky*.

Nota storica 4.2.1: Cholesky



André-Louis Cholesky (Montguyon, 15 ottobre 1875 – Bagneux, 31 agosto 1918) è stato un ingegnere francese. Divenne Ufficiale di artiglieria nella prima guerra mondiale, trovò un modo per migliorare i procedimen-

ti per il calcolo di minimi quadrati definendo il procedimento di decomposizione di matrici. Muore per le ferite riportate sul campo di battaglia. Il suo procedimento venne pubblicato solo nel 1924 da un ufficiale di nome Benoit. [Wiki16b]. Il Prof. Claude Brezinski^a ha contattato un nipote di Cholesky riuscendo a visionare gli appunti originali, questi risultano uguali a quelli sviluppati successivamente dai matematici.

^aCollega del Prof. Rodriguez?

Abbiamo visto che nella 4.9 la matrice $\mathbf{A}^T\mathbf{A}$ è definita positiva, pertanto è possibile applicare il teorema 4.2.1, questo ci permette di scrivere: $\mathbf{B} = \mathbf{A}^T\mathbf{A} = \mathbf{L}\mathbf{L}^T \Rightarrow$

$$\mathbf{L} \underbrace{(\mathbf{L}^T \mathbf{x})}_{\mathbf{y}} = \underbrace{\mathbf{A}^T \mathbf{b}}_{\mathbf{c}} \quad (4.11)$$

l'equazione 4.11 può essere scomposta in un sistema:

$$\begin{cases} \mathbf{L}\mathbf{y} = \mathbf{c} \\ \mathbf{L}^T\mathbf{x} = \mathbf{y} \end{cases} \quad (4.12)$$

Analizzando il prodotto matriciale otteniamo:

$$\mathbf{A} = \mathbf{L}\mathbf{L}^T = \begin{bmatrix} l_{11} & & & & \\ \vdots & \ddots & & & \\ l_{1i} & \dots & l_{ii} & & \\ \vdots & & & \ddots & \\ l_{1n} & \dots & \dots & \dots & l_{nn} \end{bmatrix} \begin{bmatrix} l_{1n} & \dots & l_{1j} & \dots & l_{1n} \\ & \ddots & \vdots & & \vdots \\ & & l_{ii} & & \vdots \\ & & & \ddots & \vdots \\ & & & & l_{nn} \end{bmatrix}$$

eseguendo questo prodotto si arriva:

$$a_{ij} = \sum_{k=1}^i l_{ki} l_{kj}, \quad i \leq j$$

$$a_{ij} = l_{jj}^2 + \sum_{k=1}^{i-1} l_{ki} l_{kj}, \quad i \leq j.$$

avendo esplicitato l'ultimo termine: $i = j$ gli elementi di L in funzione di A , e *separando* i due casi: $i < j$, $i = j$, pertanto:

$$l_{ij} = \frac{1}{r_{ii}} \left(a_{ij} - \sum_{k=1}^{i-1} l_{ki} l_{kj} \right), \quad i < j,$$

$$l_{jj} = \left(a_{jj} - \sum_{k=1}^{j-1} r_{kj}^2 \right)^{\frac{1}{2}}, \quad i = j.$$

Il fatto che la matrice sia simmetrica definita positiva garantisce che il radicando risulti positivo. Il costo computazionale della fattorizzazione di Cholesky è pari a $O\left(\frac{n^3}{6}\right)^5$. Lo svantaggio di questo approccio riguarda la **stabilità**, infatti la matrice $A^T A$ ha un numero di condizionamento pari al quadrato di $A \Rightarrow \kappa(A^T A) = (\kappa(A))^2$, e ricordando che:

$$\frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|} \leq \kappa(A) \frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|}$$

$$\kappa(A) = \|A\| \|A\|^{-1} \geq 1$$

dato che il numero di condizionamento di una matrice è sempre $\kappa(A) \geq 1$ il quadrato sarà ancora peggiore, questo vuol dire che *l'operatore* amplifica maggiormente gli errori, potrebbe portare a dei valori completamente inutilizzabili, cosa che non accade con la fattorizzazione QR, come visto precedentemente.

Quindi riassumendo:

Algoritmo: Cholesky

Input: $A^T A$

Output: minima soluzione

1. $A^T A = LL^T$
2. risolvi l'equazione $L(L^T \mathbf{x}) = A^T \mathbf{b}$
3. ovvero il sistema:

$$\begin{cases} L\mathbf{y} = \mathbf{c} \\ L^T \mathbf{x} = \mathbf{y} \end{cases}$$

⁵Quello di Gauss è dell'ordine $O\left(\frac{n^3}{3}\right)$, il doppio di Cholesky.

Si può migliorare . . .

Ricordando la 4.2, vogliamo *minimizzare* il quadrato della norma 2 di una matrice $m \times n$, cosa che possiamo fare applicando il teorema 3.0.1 a pagina 13, che si utilizza per questo tipo di matrici, e dunque fattorizzando QR. Si ha pertanto: $A = QR$:

$$\begin{aligned}
 \|\mathbf{Ax} - \mathbf{b}\|^2 &= \|\mathbf{QRx} - \mathbf{b}\|^2 = \\
 &= \|\mathbf{Q}(\mathbf{Rx} - \mathbf{Q}^T\mathbf{b})\|^2 = \|\mathbf{Rx} - \mathbf{Q}^T\mathbf{b}\|^2 = \\
 &= \left\| \begin{bmatrix} \mathbf{R}_1 \\ 0 \end{bmatrix} \mathbf{x} - \begin{bmatrix} \mathbf{Q}_1^T \\ \mathbf{Q}_2^T \end{bmatrix} \mathbf{b} \right\|^2 = \left\| \begin{bmatrix} \mathbf{R}_1\mathbf{x} \\ 0 \end{bmatrix} - \begin{bmatrix} \mathbf{Q}_1^T\mathbf{b} \\ \mathbf{Q}_2^T\mathbf{b} \end{bmatrix} \right\|^2 = \\
 &= \underbrace{\|\mathbf{R}_1\mathbf{x} - \mathbf{Q}_1^T\mathbf{b}\|^2}_{\mathbf{c}_1} + \underbrace{\|\mathbf{Q}_2^T\mathbf{b}\|^2}_{\mathbf{c}_2}.
 \end{aligned}$$

minimizzando:

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|^2 = \min_{\mathbf{x}} \|\mathbf{R}_1\mathbf{x} - \mathbf{c}_1\|^2 + \|\mathbf{c}_2\|^2 \quad (4.13)$$

minimizzare vuol dire trovare la \mathbf{x} tale che $\mathbf{R}_1\mathbf{x} - \mathbf{c}_1 = 0 \Rightarrow \mathbf{R}_1\mathbf{x} = \mathbf{c}_1$, allora la 4.13 diventa: $\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|^2 = \|\mathbf{c}_2\|^2$ che risulta il minimo. Riassumendo:

Algoritmo: Fattorizzazione QR

Input: A

Output: minima soluzione

1. $A = QR$
2. $R = \begin{bmatrix} R_1 \\ 0 \end{bmatrix} \quad c_1 = Q_1^T b \quad c_2 = Q_2^T b$
3. risolvo $R_1 x = c_1 \rightarrow x$ **teoricamente** coincide con la 4.10, nella pratica no, perché sono presenti gli errori
4. $\min \|\mathbf{Ax} - \mathbf{b}\|^2 = \|\mathbf{c}_2\|^2$.

Il vantaggio di questa fattorizzazione rispetto al metodo di Cholesky è che il condizionamento rimane lineare, non va al quadrato.

4.3 $k = m < n$

Teorema 4.3.1

$$S = \{ \mathbf{x} \in \mathbb{R}^n : \|\mathbf{Ax} - \mathbf{b}\|^2 = \min \} \quad \mathbf{x} \in S \iff \mathbf{A}^T \underbrace{(\mathbf{b} - \mathbf{Ax})}_{\mathbf{r}=\text{residuo}} = 0 \quad \mathbf{A}^T\mathbf{b} - \mathbf{A}^T\mathbf{Ax} = 0 \Rightarrow \mathbf{A}^T\mathbf{Ax} = \mathbf{A}^T\mathbf{b}$$

Il teorema 4.3.1 ci assicura che il minimo può essere trovato come $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$ ⁶.

Data la matrice $\mathbf{A} \in \mathbb{R}^{n \times m}$, il suo rango vale $m = \min(n, m)$, il prodotto $\mathbf{A}^T \mathbf{A}$ ha dimensione $n \times n$, dunque ha rango: $k = \text{rank}(\mathbf{A}^T \mathbf{A}) = \min(n, n) = n$ ⁷. Quindi certamente $\mathbf{A}^T \mathbf{A}$ è singolare. Il sistema è consistente dunque posso **sempre** raggiungere l'uguaglianza, ovvero ammette sempre una soluzione. La soluzione \mathbf{r} dell'equazione $\mathbf{A}^T (\mathbf{b} - \mathbf{A} \mathbf{x}) = \mathbf{A}^T \mathbf{r} = 0$ equivale al nucleo di \mathbf{A}^T , infatti:

$$\mathcal{N}(\mathbf{A}^T) = \{ \mathbf{x} \in \mathbb{R}^m : \mathbf{A}^T \mathbf{x} = 0 \} \Rightarrow \mathbf{r} \in \mathcal{N}(\mathbf{A}^T).$$

il vettore residuo è per definizione pari:

$$\mathbf{r} = \mathbf{b} - \mathbf{A} \mathbf{x} \Rightarrow \mathbf{b} = \underbrace{\mathbf{A} \mathbf{x}}_{\in \mathcal{R}(\mathbf{A})} + \underbrace{\mathbf{r}}_{\in \mathcal{N}(\mathbf{A}^T)} \quad (4.14)$$

i due vettori sono perpendicolari:

$$\mathbf{A} \mathbf{x} \perp \mathbf{r}$$

il codominio, vettore \mathbf{b} , $\in \mathbb{R}^m$, mentre $\mathbf{A} \mathbf{x} \in \mathcal{R}(\mathbf{A})$, quindi:

$$\mathbb{R}^m = \mathcal{R}(\mathbf{A}) \oplus \mathcal{N}(\mathbf{A}^T)$$

dove \oplus indica la somma diretta:

$$\oplus \equiv \begin{cases} \mathcal{R}(\mathbf{A}) \cup \mathcal{N}(\mathbf{A}^T) = \mathbb{R}^m \\ \mathcal{R}(\mathbf{A}) \perp \mathcal{N}(\mathbf{A}^T). \end{cases}$$

Il problema è che abbiamo ∞ vettori che soddisfano la 4.14, allora cerchiamo il più corto. In sostanza cerchiamo soluzioni:

$$\begin{aligned} \min \|\mathbf{x}\| \\ \text{s.a. } \mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}. \end{aligned}$$

Il problema è noto come Minimum Norm Least Squares Solution (MNLSS) o LSS duale.

Analizzando il problema da più vicino, sostituendo $\mathbf{y} = \mathbf{A} \mathbf{z}$ può essere riscritto:

$$\begin{cases} \min \|\mathbf{y}\|_2^2 \\ \text{s.a. } \mathbf{A}^T \mathbf{y} = \mathbf{c} \end{cases} \quad \begin{cases} \min \frac{1}{2} \mathbf{y}^T \mathbf{y} \\ \text{s.a. } \mathbf{A}^T \mathbf{y} = \mathbf{c} \end{cases}$$

con $\mathbf{y} \perp \mathcal{N}(\mathbf{A})$, inoltre il fattore $\frac{1}{2}$, è utile numericamente. Per rispettare i vincoli scriviamo e risolvere il problema scriviamo i moltiplicatori di Lagrange, che sono tanti quanti sono le equazioni del sistema, con $\mathbf{z} \in \mathbb{R}^n$, calcolo la distanza dall'origine:

$$\begin{aligned} \mathcal{L}(\mathbf{y}, \mathbf{z}) &= \frac{1}{2} \mathbf{y}^T \mathbf{y} - \sum z_i [(\mathbf{A}^T \mathbf{y})_i - c_i] = \\ &= \frac{1}{2} \mathbf{y}^T \mathbf{y} + \mathbf{z}^T (\mathbf{c} - \mathbf{A}^T \mathbf{y}) = \\ &= \frac{1}{2} \mathbf{y}^T \mathbf{y} + \mathbf{z}^T \mathbf{c} - \mathbf{z}^T \mathbf{A}^T \mathbf{y}. \end{aligned}$$

⁶Che coincide con la 4.9.

⁷Si faccia riferimento al punto 3 a pagina 21, il rango della matrice per quel caso è n , per questo caso è sempre n , ma i due ranghi sono diversi, infatti per quel caso 3 a pagina 21 $n < m$, in questo caso $n > m$, quindi questa n è maggiore di quel caso. Per essere chiari se chiamassi n_1 il rango della matrice nel caso 3 a pagina 21, e n_2 il rango in studio possiamo scrivere $n_2 > n_1$.

derivando, ricordando le opportune regole di derivazione per matrici 4.8 a pagina 21 e imponendole nulle le derivate:

$$\begin{aligned}\nabla_y \mathcal{L}(y, z) &= \frac{1}{2} 2y - Az = 0 \\ \nabla_z \mathcal{L}(y, z) &= c - A^T y = 0\end{aligned}$$

quindi si ottiene:

$$\begin{cases} A^T Az = c \\ y = Az. \end{cases} \quad (4.15)$$

Invertendo la seconda equazione delle 4.15 si ottiene:

$$z = \underbrace{A(A^T A)^{-1}}_{(A^T)^\dagger} c$$

dove $(A^T)^\dagger = A(A^T A)^{-1}$ è la matrice pseudo-inversa duale che è solo inversa destra:

$$\begin{aligned}A^T (A^T)^\dagger &= A^T A (A^T A)^{-1} = \mathbf{I} \\ (A^T)^\dagger A^T &= A (A^T A)^{-1} A^T \neq \mathbf{I}.\end{aligned}$$

Capitolo 5

Autovalori e autovettori



CONSIDERIAMO il seguente problema¹

$$\mathbf{Ax} = \lambda \mathbf{x} \quad \mathbf{A} \in \mathbb{R}^{n \times n}, \mathbf{x} \in \mathbb{R}^n \setminus \{0\}, \lambda \in \mathbb{R}$$

$$\mathbf{Ax} - \lambda \mathbf{x} = 0$$

$$\mathbf{Ax} - \lambda \mathbf{Ix} = 0$$

$$(\mathbf{A} - \lambda \mathbf{I})\mathbf{x} = 0$$

dove λ è detto autovalore della, matrice \mathbf{A} ; il sistema lineare omogeneo è soddisfatto se la matrice $(\mathbf{A} - \lambda \mathbf{I})$ è singolare ovvero:

$$\det(\mathbf{A} - \lambda \mathbf{I}) = 0$$

$$p_n(\lambda) = \det(\mathbf{A} - \lambda \mathbf{I})$$

e questo viene chiamato **polinomio caratteristico** della matrice \mathbf{A} . La soluzione di questo polinomio, le radici, forniscono gli autovalori di \mathbf{A} .

Teorema 5.0.1: Fondamentale dell'Algebra

$p_n(x)$ ha n radici complesse contate con la loro molteplicità.

Questo teorema ci assicura che $\exists n$ radici, il problema è che non esiste una formula risolutiva per polinomi di ordine superiore al sesto, quindi **funziona** solo per polinomi di grado basso, che si traduce in matrici di ordine n basso, il che ovviamente rende impraticabile, la determinazione degli autovalori percorrendo questa strada.

¹Riprenderemo alcuni concetti visti nella sezione 2.2.1 a pagina 9.

5.1 Trasformazione di similitudine

Definizione 5.1.1: Similitudine

La matrice B è **simile** a A se $\exists X$ non singolare : $X^{-1}AX = B$ e si indica $A \sim B$

Teorema 5.1.1

$$A \sim B \iff \sigma(A) = \sigma(B)$$

Dimostrazione. Teorema 5.1.1

$$\begin{aligned} Bx = \lambda x &\Rightarrow X^{-1}AXv = \lambda v \\ A \underbrace{Xv}_w &= \lambda \underbrace{Xv}_w \\ Aw &= \lambda w \end{aligned}$$

□

Definizione 5.1.2: Matrice diagonalizzabile

$$\begin{aligned} A \sim D &= \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) \\ \exists X : X^{-1}DX &= D \\ A &= XDX^{-1} \quad \nexists \text{ per tutte le matrici.} \end{aligned}$$

Definizione 5.1.3: Matrice unitariamente diagonalizzabile

Se: $X = U$ dove $U^*U = UU^* = I$

$$U^{-1}AU = D \Rightarrow U^*AU = D$$

Si ricordi che mentre fare X^{-1} è un calcolo che porta inevitabilmente ad errori farne la trasposta (aggiunta) non **introduce errori**.

Teorema 5.1.2

A è diagonalizzabile \iff gli autovettori sono **linearmente indipendenti**.

Dimostrazione. Teorema 5.1.2

$$\begin{aligned} X^{-1}AX &= D, \quad X = [v_1|v_2|\dots|v_n] \\ AX &= XD \\ A[v_1|v_2|\dots|v_n] &= [v_1|v_2|\dots|v_n]diag(\lambda_1, \lambda_2, \dots, \lambda_n) \\ X^{-1}AX &= D \end{aligned}$$

X potrà essere invertita solo se le colonne sono linearmente indipendenti. \square

Teorema 5.1.3: Schur

$$A \in \mathbb{C}^{n \times n} \exists U \quad U^*AU = T = \begin{bmatrix} \lambda_1 & * & \dots & * \\ 0 & \lambda_2 & \dots & * \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix}$$

Corollario: Se $A = A^*$ (*Hermitiana*) \Rightarrow è unitariamente diagonale.

$$\begin{aligned} A &= UTU^* \quad A^* = (UTU^*)^* = UT^*U^* \Rightarrow \\ A &= A^* \Rightarrow T = T^* \Rightarrow T = D \\ \lambda_i &\in \mathbb{R} \quad v_i \perp v_j \end{aligned}$$

gli autovalori sono reali e gli autovettori sono perpendicolari.

Molti problemi si concentrano sulla *localizzazione* che permette di sapere dove *non* sono gli autovalori. Altri algoritmi si concentrano sulla *separazione* degli autovalori. Un motore di ricerca utilizza questi algoritmi, si è interessati alla conoscenza dell'autovalore più grande in modulo, viceversa una rete sociale è interessata alla conoscenza dell'autovalore più piccolo (in modulo), un microblogging può essere interessato al penultimo autovalore etc.

Cerchiamo l'autovalore più grande, possiamo raffinare la ricerca effettuando una deflazione la matrice *smontandola*, per esempio prendo una matrice 10×10 , trovo l'autovalore più grande in modulo, elimino l'autovalore trovato ottenendo una matrice la matrice che diventa 9×9 , ritrovo l'autovalore più grande e la matrice diventa 8×8 , etc ...

Dall'autovalore posso trovare l'autovettore corrispondente, mentre se conosco l'autovettore posso trovare l'autovalore sfruttando il quoziente di Rayleigh:

$$Av = \lambda v \Rightarrow \frac{v^T Av}{v^T v} = \frac{\lambda v^T v}{v^T v} = \lambda \quad (5.1)$$

Nota storica 5.1.1: Schur



Issai Schur (Mahilëu, 10 gennaio 1875 - Tel Aviv 10 gennaio 1941 (66 anni giusti giusti)) è stato un matematico tedesco. È conosciuto per i suoi lavori in rappresentazione dei gruppi, per quelli in combinatoria, teoria dei numeri e fisica teorica. [Wiki15b]

Teorema 5.1.4: Bauer-Fike

$$\begin{aligned} A \in \mathbb{C}^{n \times n} \quad \exists X = X^{-1}AX = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) \\ B = A + E, \mu \in \sigma(B) \\ \min_{\lambda \in \sigma(A)} |\lambda - \mu| < k_2(X) \|E\|_2. \end{aligned}$$

dove la matrice E e la matrice degli errori, $k_2(X)$ è il numero di condizionamento della matrice: $\|X\|_2 \|X^{-1}\|_2$ e μ è l'autovalore calcolato.

Corollario: Se $A = A^*$ X è unitaria $\Rightarrow k_2(X) = 1$
 $\min_{\lambda \in \sigma(A)} |\lambda - \mu| \leq \|E\|_2$. E è stabile quando gli autovettori sono perpendicolari: $\lambda_i \in \mathbb{R} \quad v_i \perp v_j$.

Nota storica 5.1.2: Bauer

Friedrich Ludwig "Fritz" Bauer (10 giugno 1924 - 26 Marzo 2015) è stato un informatico tedesco e professore. Bauer ha lavorato anche nei comitati che hanno sviluppato il linguaggio di programmazione *ALGOL60*. Nel 1968, Bauer ha coniato il termine **Software Engineering**. [Wiki17a].

5.2 Metodo iterativi

Consideriamo il sistema $Ax = b$, vogliamo risolverlo, possiamo usare dei metodi

iterativi iniziamo a vederne uno: il Metodo delle potenze.

5.2.1 Metodo delle potenze

Consideriamo una matrice $A [n \times n]$ prendiamo un vettore a caso $x^{(0)} \in \mathbb{R}^2$ e *costruiamo* un algoritmo:

$$\begin{aligned} x^{(1)} &= Ax^{(0)} \\ x^{(2)} &= Ax^{(1)} \\ x^{(3)} &= Ax^{(2)} \\ &\vdots \\ x^{(n+1)} &= Ax^{(n)} \end{aligned}$$

si verifica sempre che il vettore $x^{(n)}$ tende a diventare l'autovettore, poiché in realtà l'autovettore indica una direzione. Prendendo spunto da questo algoritmo³, possiamo svilupparne un altro, in particolare facendo delle ipotesi su A :

- ① diagonalizzabile: $A = XDX^{-1} \iff v_i$ sono indipendenti;
- ② $|\lambda_1| > |\lambda_i| \quad i = 2, 3, \dots, n$ sono ordinati in maniera decrescente;

²Per indicare che il procedimento è iterativo si mette un numero tra parentesi tonde all'esponente, le parentesi si mettono per considerare che non è un esponente.

³Che in genere crea problemi di under e over flow

3 $\mathbf{x}^{(0)}$ ha una componente non nulla nella \mathbf{v}_1 .

Il primo vettore può essere scritto come un combinazione lineare di termini $\mathbf{x}^{(0)} = \sum_{i=1}^n \alpha_i \mathbf{v}_i$, pertanto:

$$\begin{aligned} \mathbf{x}^{(k)} &= \mathbf{A}\mathbf{x}^{(k-1)} = \mathbf{A}\mathbf{A}\mathbf{x}^{(k-2)} = \mathbf{A}^2\mathbf{x}^{(k-2)} = \dots = \mathbf{A}^k\mathbf{x}^{(0)} \\ \mathbf{x}^{(k)} &= \mathbf{A}^k\mathbf{x}^{(0)} = \mathbf{A}^k\left(\sum_{i=1}^n \alpha_i \mathbf{v}_i\right) = \sum_{i=1}^n \alpha_i \mathbf{A}^k \mathbf{v}_i = \\ &= \sum_{i=1}^n \alpha_i \lambda_i^k \mathbf{v}_i = \alpha_1 \lambda_1^k \mathbf{v}_1 + \sum_{i=2}^n \alpha_i \lambda_i^k \mathbf{v}_i \\ &= \text{avendo sfruttato l'ipotesi 3} \Rightarrow \alpha_1 \neq 0 \\ &= \alpha_1 \lambda_1^k \left(\mathbf{v}_1 + \sum_{i=2}^n \frac{\alpha_i}{\alpha_1} \left(\frac{\lambda_i}{\lambda_1}\right)^k \mathbf{v}_i\right) \end{aligned}$$

poiché vale l'ipotesi 2 $|\frac{\lambda_i}{\lambda_1}| < 1$ e per $k \rightarrow \infty \Rightarrow |\frac{\lambda_i}{\lambda_1}| \rightarrow 0$

$$\mathbf{x}^{(k)} \rightarrow \alpha_1 \lambda_1^k \mathbf{v}_1$$

il limite tende ad essere un autovettore, il vettore: $\mathbf{x}^{(k)}$ tende ad essere parallelo a \mathbf{v}_1 . Prendiamo in considerazione i seguenti casi:

$$\begin{aligned} |\lambda_1| > 1 \quad \lambda_1^k &\rightarrow \infty \quad \text{il vettore si allunga} \\ |\lambda_1| < 1 \quad \lambda_1^k &\rightarrow 0 \quad \text{il vettore sparisce} \\ |\lambda_1| = 1 \quad \cancel{\neq} & \end{aligned}$$

come usare l'algorithm? Ad ogni iterazione normalizzo \mathbf{x} e *richiamo* il vettore:

$$\begin{aligned} \mathbf{q}^{(0)} &= \frac{\mathbf{x}^{(0)}}{\|\mathbf{x}^{(0)}\|} \\ &\vdots \\ \mathbf{q}^{(k+1)} &= \frac{\mathbf{x}^{(k+1)}}{\|\mathbf{x}^{(k+1)}\|} \end{aligned}$$

risolvo il sistema: $\mathbf{x}^{(k+1)} = \mathbf{A}\mathbf{q}^{(k)}$, ovvero trovo l'autovettore successivamente calcolo l'autovalore l'equazione 5.1 di Rayleigh, ricordando che $\|\mathbf{v}^T \mathbf{v}\| = \|\mathbf{q}^T \mathbf{q}\| = 1$:

$$\begin{aligned} \lambda &= \frac{\mathbf{v}^T \mathbf{A} \mathbf{v}}{\mathbf{v}^T \mathbf{v}} \Rightarrow \\ \lambda^{(k+1)} &= \frac{\mathbf{q}^{(k+1)T} \mathbf{A} \mathbf{q}^{(k+1)}}{\|\mathbf{q}^{(k+1)T} \mathbf{q}^{(k+1)}\|} \\ &= \frac{\mathbf{q}^{(k+1)T} \mathbf{A} \mathbf{q}^{(k+1)}}{1} \\ \lambda^{(k+1)} &= \mathbf{q}^{(k+1)T} \mathbf{A} \mathbf{q}^{(k+1)} \end{aligned}$$

Nota storica 5.2.1: Cauchy

Augustin-Louis Cauchy (Parigi, 21 agosto 1789 – Sceaux, 23 maggio 1857) è stato un matematico e ingegnere francese. Ha avviato il progetto della formulazione e dimostrazione rigorosa dei teoremi dell'analisi infinitesimale basato sull'utilizzo delle nozioni di limite e di continuità. [Wiki17b].

come posso terminare l'algoritmo? Fisso un limite. Criterio di Cauchy:

$$|\lambda^{(k+1)} - \lambda^{(k)}| < \tau$$

... e dal momento dal punto di vista numerico gli errori relativi sono *migliori* di quelli assoluti allora preferiamo:

$$\frac{|\lambda^{(k+1)} - \lambda^k|}{|\lambda^{(k+1)}|} < \tau \quad (5.2)$$

... e dal momento che la 5.2 non è proprio *buona* dal punto di vista numerico, preferiamo:

$$|\lambda^{(k+1)} - \lambda^k| < \tau |\lambda^{(k+1)}|. \quad (5.3)$$

Potrebbe accadere che la 3 nella pagina precedente non sia verificata e cioè: $\alpha_1 = 0$ questo comprometterebbe l'algoritmo ovvero potremmo trovare **altri** autovettori. Per **cercare** di evitare questo, prendiamo $\mathbf{x}^{(0)}$ assolutamente a caso. Questa è una condizione necessaria, ma non sufficiente, per evitare lo sfortunato caso in cui: $\alpha_1 = 0$.

Un'altra considerazione sulle ipotesi, se la 2 a pagina 30 non fosse verificata ovvero gli autovalori, in modulo, sono uguali, per esempio se accade che $\lambda_1 = 10; \lambda_2 = -10$ l'algoritmo **non converge**.

Algoritmo Metodo delle potenze

Si riporta l'algoritmo pronto per essere implementato.

Algoritmo: Metodo delle potenze

Input: $A, \mathbf{x}^{(0)}, \tau, N_{max}$ (numero massimo di iterazioni)

Output: λ, \mathbf{q}, k

- 1 $\mathbf{q}^{(0)} = \frac{\mathbf{x}^{(0)}}{\|\mathbf{x}^{(0)}\|}$
 - 2 $k = 0$
 - 3 $\lambda = \mathbf{q}^T \mathbf{A} \mathbf{q}$, flag = true
 - 4 while flag
 - a $k = k + 1, \lambda_{old} = \lambda$
 - b $\mathbf{x} = \mathbf{A} \mathbf{q}$
 - c $\mathbf{q} = \frac{\mathbf{x}}{\|\mathbf{x}\|}$
 - d $\lambda = \mathbf{q}^T \mathbf{A} \mathbf{q}$
 - e flag ($|\lambda - \lambda_{old}| \geq \tau |\lambda|$) AND ($k \leq N_{max}$)
- end while

5.2.2 Eigenvector centrality

Questo algoritmo è la base di noto motore di ricerca, consideriamo la rete in figura 5.1, come possiamo considerare l'importanza di ciascun nodo? E come possiamo quantificarla? Scriviamo la matrice di adiacenza corrispondente alla rete, che risulta simmetrica poiché il grafo non è orientato.

$$\begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \begin{cases} x_1 = x_2 + x_3 \\ x_2 = x_1 + x_3 + x_4 \\ x_3 = x_1 + x_2 \\ x_4 = x_2 \end{cases} \quad (5.4)$$

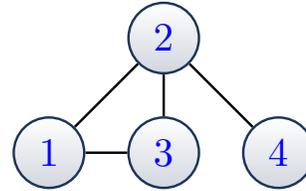


Figura 5.1: Rete

più in generale possiamo esprimere questa forma come $x_i = \frac{1}{\lambda} \sum_{j=1}^n a_{ij} x_j$ dove il fattore $\frac{1}{\lambda}$ è una costante, questa relazione può esprimersi come: $\mathbf{x} = \frac{1}{\lambda} \mathbf{A} \mathbf{x} \Rightarrow \mathbf{A} \mathbf{x} = \lambda \mathbf{x}$ dove x_i non può essere negativa.

Definizione 5.2.1: Matrice non negativa

Una matrice è non negativa se **tutti** gli elementi sono maggiori o uguali a 0.

Definizione 5.2.2: Primitiva

La matrice $\mathbf{A} \geq 0$ si dice primitiva se $\exists m > 0: \mathbf{A}^m \gg 0$.

Definizione 5.2.3: Riducibile

Una matrice quadrata \mathbf{A} di dimensione n si dice **riducibile** se esiste una matrice di permutazione \mathbf{P} tale che

$$\mathbf{P} \mathbf{A} \mathbf{P}^T = \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ 0 & \mathbf{B}_{22} \end{bmatrix}$$

essendo \mathbf{B}_{11} e \mathbf{B}_{22} matrici di dimensione $k \times k$ e $(n - k) \times (n - k)$, rispettivamente. Una matrice è **irriducibile** se non è riducibile.

Teorema 5.2.1: Perron-Frobenius

Data una matrice quadrata \mathbf{A} , se questa è non negativa, primitiva e irriducibile allora:

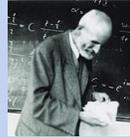
- 1 L'autovalore λ di valore massimo in modulo è reale e positivo
- 2 λ è un autovalore semplice
- 3 l'autovettore corrispondente ha tutte le componenti positive
- 4 l'autovettore corrispondente è l'unico non negativo
- 5 $\rho(\mathbf{A})$ è una funzione strettamente crescente in ognuno dei suoi elementi.

Nota storica 5.2.2: Perron-Frobenius



Ferdinand Georg Frobenius (26 ottobre 1849 - 3 Agosto 1917) è stato un matematico tedesco, noto per i suoi contributi alla teoria delle funzioni ellittiche, equazioni differenziali e alla teoria dei gruppi. [Wiki17c]
 Oskar Perron (7 maggio 1880 - 22 febbraio 1975) è stato un matematico tedesco. Ha scritto un libro enciclopedico sulle frazioni continue *Die Lehre von den*

Kettenbrüchen. Ha introdotto il paradosso di Perron: Sia N il più grande intero. Se $N > 1$, allora $N^2 > N$, contraddice la definizione di N . Quindi $N = 1$. Per illustrare il pericolo di supporre che la soluzione di un problema di ottimizzazione possa esistere. [Wiki15a].



Sfruttiamo il Metodo delle potenze e calcoliamoci, per la matrice 5.4, i pesi partendo da un vettore iniziale: $\mathbf{x}^{(0)} = [1, 1, 1, 1]^T$, allora il vettore $\mathbf{x}^{(1)} = [2, 3, 2, 1]^T$ etc. ... giungeremo all'autovettore più grande.

Capitolo 6

Singular Value Decomposition



La scomposizione **SVD** riveste un ruolo fondamentale, non né verrà riportato l'algoritmo per la sua determinazione, ma si sfrutteranno direttamente i risultati. Prendiamo una **qualsiasi** matrice **A** essa può **sempre** essere decomposta come: $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, $\mathbf{A} \in \mathbb{R}^{m \times n}$; $\mathbf{U} = [\mathbf{u}_1 | \mathbf{u}_2 | \dots | \mathbf{u}_m]$; $\mathbf{V} = [\mathbf{v}_1 | \mathbf{v}_2 | \dots | \mathbf{v}_n]$; $\mathbf{\Sigma}_{m \times n} =$

$$\begin{bmatrix} \sigma_1 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_n & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 \end{bmatrix} \text{ dove: } \mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}_m, \mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}_n,$$

$$\sigma_i \in \mathbb{R}^+ \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0, \quad k = \text{rank}(\mathbf{A}).$$

Per un sistema lineare vale:

$$\begin{aligned} \mathbf{Ax} &= [\mathbf{u}_1 | \dots | \mathbf{u}_k | \mathbf{u}_{k+1} | \dots | \mathbf{u}_m] \begin{bmatrix} \sigma_1 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_n & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^T \mathbf{x} \\ \vdots \\ \mathbf{v}_k^T \mathbf{x} \\ \mathbf{v}_{k+1}^T \mathbf{x} \\ \vdots \\ \mathbf{v}_n^T \mathbf{x} \end{bmatrix} = \\ &= [\mathbf{u}_1 | \dots | \mathbf{u}_k | \mathbf{u}_{k+1} | \dots | \mathbf{u}_m] \begin{bmatrix} \sigma_1 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_n & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^T \mathbf{x} \\ \vdots \\ \mathbf{v}_k^T \mathbf{x} \\ \mathbf{v}_{k+1}^T \mathbf{x} \\ \vdots \\ \mathbf{v}_n^T \mathbf{x} \end{bmatrix} = \\ &= [\mathbf{u}_1 \dots \mathbf{u}_k | \mathbf{u}_{k+1} \dots \mathbf{u}_m] \begin{bmatrix} \sigma_1 \mathbf{v}_1^T \mathbf{x} \\ \vdots \\ \sigma_k \mathbf{v}_k^T \mathbf{x} \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \sum_{i=1}^k \sigma_i (\mathbf{v}_i^T \mathbf{x}) \mathbf{u}_i \end{aligned}$$

da queste forme così come sono scritte si vedono i sottospazi vettoriali, in particolare:

■ il nucleo: $\mathcal{N}(\mathbf{A}) = \text{span}[\mathbf{v}_{k+1} \dots \mathbf{v}_n]$

- ▀ il range: $\mathcal{R}(\mathbf{A}) = \text{span}[\mathbf{u}_1 \dots \mathbf{u}_k]$
- ▀ il nucleo trasposto: $\mathcal{N}(\mathbf{A}^*) = \mathcal{N}(\mathbf{A}^T) = \text{span}[\mathbf{u}_{k+1} \dots \mathbf{u}_n]$
- ▀ il range trasposto: $\mathcal{R}(\mathbf{A}^*) = \mathcal{R}(\mathbf{A}^T) = \text{span}[\mathbf{v}_1 \dots \mathbf{v}_k]$
- ▀ $\mathcal{N}(\mathbf{A})^\perp = \mathcal{R}(\mathbf{A}^T)$, $\mathcal{R}(\mathbf{A})^\perp = \mathcal{N}(\mathbf{A}^T)$

Sistema omogeneo

Pensiamo di voler risolvere il problema lineare omogeneo, **immediatamente** possiamo saperne la soluzione, se conosciamo la scomposizione SVD, infatti

$$\mathbf{A}\mathbf{x} = \mathbf{0} \Rightarrow \mathcal{N}(\mathbf{A}) \neq \{0\} \Rightarrow \mathbf{x} = \sum_{i=k+1}^n \alpha_i \mathbf{v}_i$$

Norma

Vogliamo conoscere la norma 2:

$$\begin{aligned} \|\mathbf{A}\|_2 &= \sqrt{\rho(\mathbf{A}\mathbf{A}^T)} = \sqrt{\rho(\mathbf{V}\Sigma^T\mathbf{U}^T\mathbf{U}\Sigma\mathbf{V}^T)} = \\ &= \sqrt{\rho(\underbrace{\mathbf{V}\Sigma^2\mathbf{V}^T}_{\text{fatt.spett.}})} = \sigma_1 \end{aligned}$$

Numero di condizionamento

Il numero di condizionamento in norma 2, grazie alla scomposizione SVD può essere calcolato semplicemente:

$$k_2(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2 = \frac{\sigma_1}{\sigma_n}$$

Prendendo come riferimento il teorema di Bauer-Fike 5.1.4 a pagina 30 possiamo generalizzarlo ad una matrice $m \times n$:

Teorema 6.0.1

$$\begin{aligned} \mathbf{A} &\in \mathbb{C}^{m \times n} \\ |\sigma_i(\mathbf{A} + \mathbf{E}) - \sigma_i(\mathbf{A})| &\leq \|\mathbf{E}\|_2. \end{aligned} \tag{6.1}$$

Il problema è sempre stabile. La relazione 6 considera gli errori assoluti, questo potrebbe essere un problema se σ_i è piccolo, allora possiamo pensare di **scartare** i σ_i troppo piccoli e pensare di **troncare** la serie:

$$\mathbf{A} = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T \approx \sum_{i=1}^l \sigma_i \mathbf{u}_i \mathbf{v}_i^T = \mathbf{A}_l.$$

Un teorema ci assicura la migliore approssimazione è data:

Teorema 6.0.2: Migliore approssimazione

$$\arg \min_{\text{rank}(X)=l} \|A - X\|_2 = A_l$$

quindi:

$$\|A - A_l\|_2 = \sigma_{l+1}$$

$$\|A - A_l\|_F = \sqrt{\sigma_{l+1}^2 + \dots + \sigma_k^2}$$

6.1 Un problema facile

Prendiamo un sistema lineare, con matrice $A \in \mathbb{R}^{n \times n}$, potremmo usare Gauss, che essendo diretto è **corretto**, ma qua useremo il metodo iterativo SVD, pertanto:

$$Ax = b \Rightarrow U \Sigma V^T x = b \Rightarrow \Sigma \underbrace{V^T x}_y = \underbrace{U^T b}_c$$

$$c = U^T b = \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_n^T \end{bmatrix} b = \begin{bmatrix} u_1^T b \\ u_2^T b \\ \vdots \\ u_n^T b \end{bmatrix}$$

$$\Sigma y = c \Rightarrow y_i = \frac{c_i}{\sigma_i} = \frac{u_i^T b}{\sigma_i}$$

$$y = V^T x \Rightarrow Vy = x = [v_1, v_2, \dots, v_n] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \Rightarrow \text{troviamo dunque il vettore } x :$$

$$x = \sum_{i=1}^n \frac{u_i^T b}{\sigma_i} v_i \quad (6.2)$$

6.2 Un altro problemino ... meno facile

Sicuramente ci sarà una matrice, che chiameremo con poca fantasia $A \in \mathbb{R}^{m \times n}$: rettangolare cerchiamo di risolvere il problema dei minimi quadrati con minima norma: MNLSS (Minimum

Norm Least Square Solution) rappresentabile:

$$\begin{aligned} & \min \|\mathbf{x}\|_2 \\ \text{s.a. } & \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2 \end{aligned} \quad (6.3)$$

data: $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$

$$\begin{aligned} \|\mathbf{Ax} - \mathbf{b}\|_2 &= \|\mathbf{U}\Sigma\mathbf{V}^T\mathbf{x} - \mathbf{b}\|_2 = \|\Sigma \underbrace{\mathbf{V}^T\mathbf{x}}_{\mathbf{y}} - \underbrace{\mathbf{U}^T\mathbf{b}}_{\mathbf{c}}\|_2 = \|\Sigma\mathbf{y} - \mathbf{c}\|_2 = \\ &= \left\| \begin{bmatrix} \Sigma_k & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} - \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} \right\|_2 = \left\| \begin{bmatrix} \Sigma_k y_1 \\ 0 \end{bmatrix} - \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} \right\|_2 = \\ &= \left\| \begin{bmatrix} \Sigma_k y_1 - c_1 \\ -c_2 \end{bmatrix} \right\|_2 \end{aligned}$$

il vettore \mathbf{c}_2 è costante, posso minimizzare, annullandolo, il termine: $\Sigma_k y_1 - c_1 = 0 \Rightarrow \Sigma_k y_1 = c_1 \Rightarrow y_i = \frac{c_i}{\sigma_i}, i = 1 \dots k$. Scegliendo queste y annulliamo la prima riga, dunque il minimo risulta: $\|\mathbf{c}_2\|_2$, ma \mathbf{y} può essere scritta:

$$\mathbf{y} = \mathbf{V}^T \mathbf{x} \Rightarrow \mathbf{x} = \mathbf{V}\mathbf{y} = \mathbf{V} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

la seconda condizione delle 6.3 richiede $\min \|\mathbf{x}\|$ questo implica dato che ho ∞ soluzioni **devo** scegliere quella che minimizza \mathbf{x} allora prendo $y_2 = 0$, la soluzione normale. Allora:

$$\begin{aligned} \mathbf{x} &= \mathbf{V} \begin{bmatrix} y_1 \\ 0 \end{bmatrix} = [\mathbf{V}_1 \mathbf{V}_2] \begin{bmatrix} y_1 \\ 0 \end{bmatrix} = \mathbf{V}_1 y_1 = \\ &= \sum_{i=1}^k y_i \mathbf{v}_i = \sum_{i=1}^k \frac{c_i}{\sigma_i} \mathbf{v}_i = \end{aligned}$$

$$\mathbf{x} = \sum_{i=1}^k \frac{\mathbf{u}_i^T \mathbf{b}}{\sigma_i} \mathbf{v}_i \quad (6.4)$$

la stessa dell'equazione 6.2.

6.3 Pseudo-inversa

Considerando la scomposizione SVD di una matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$ vogliamo calcolare la pseudo-inversa:

$$\mathbf{x} = \mathbf{A}^\dagger \mathbf{b} = \mathbf{V}\mathbf{y} = \mathbf{V} \begin{bmatrix} \Sigma^{-1} \mathbf{c}_1 \\ 0 \end{bmatrix} = \mathbf{V} \begin{bmatrix} \frac{1}{\sigma_1} & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & \frac{1}{\sigma_k} & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \mathbf{V} \underbrace{\begin{bmatrix} \Sigma_k^{-1} & 0 \\ 0 & 0 \end{bmatrix}}_{\mathbf{A}^\dagger} \mathbf{U}^T \mathbf{b}$$

quindi possiamo riassumere come:

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T = \mathbf{U} \begin{bmatrix} \sigma_1 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_k & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 0 \end{bmatrix} \mathbf{V}^T$$

$$\mathbf{A}^\dagger = \mathbf{V}\Sigma^\dagger\mathbf{U}^T = \mathbf{V} \begin{bmatrix} \frac{1}{\sigma_1} & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\sigma_k} & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 0 \end{bmatrix} \mathbf{U}^T$$

Considerando ora una matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$ **a rango pieno** $= n$:

$$\begin{aligned} \mathbf{A}^\dagger (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T &= (\mathbf{V}\Sigma^T \underbrace{\mathbf{U}^T \mathbf{U}}_{\mathbf{I}} \Sigma \mathbf{V}^T)^{-1} \mathbf{V}\Sigma^T \mathbf{U}^T = (\mathbf{V}\Sigma^2 \mathbf{V}^T)^{-1} \mathbf{V}\Sigma^T \mathbf{U}^T = \\ &= \mathbf{V}\Sigma^{-2} \underbrace{\mathbf{V}^T \mathbf{V}}_{\mathbf{I}} \Sigma^T \mathbf{U}^T = \mathbf{V} \begin{bmatrix} \frac{1}{\sigma_1^2} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\sigma_n^2} \end{bmatrix} \begin{bmatrix} \sigma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_n \end{bmatrix} \mathbf{U}^T \end{aligned}$$

Alcune proprietà della pseudo-inversa

- 1 $(\mathbf{A}^\dagger)^\dagger = \mathbf{A}$
- 2 $(\mathbf{A}^*)^T = (\mathbf{A}^T)^*$
- 3 $(\mathbf{A}^T \mathbf{A})^\dagger = \mathbf{A}^\dagger (\mathbf{A}^T)^\dagger$
- 4 $(\mathbf{AB})^\dagger \neq \mathbf{B}^\dagger \mathbf{A}^\dagger$
- 5 $\mathbf{AA}^\dagger \neq \mathbf{A}^\dagger \mathbf{A}$

6.4 Usando la fattorizzazione QR

Teorema 6.4.1

Data una matrice $A \in \mathbb{R}^{m \times n}$, non necessariamente a rango pieno, che per semplicità possiamo considerare con $m > n \Rightarrow \exists Q : QQ^T = Q^T Q = I$

$$R = \begin{bmatrix} \sigma_1 & * & * \\ \vdots & \ddots & * \\ 0 & \cdots & \sigma_n \\ 0 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 0 \end{bmatrix} = \begin{bmatrix} R_1 \\ 0 \end{bmatrix}$$

$$A = QR$$

o in maniera compatta:

$$A = [Q_1 Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix} = Q_1 R_1$$

R e A hanno la stessa SVD, infatti, se $R = U\Sigma V^T$:

$$A = QR = \underbrace{(QU)}_{\tilde{U}} \Sigma V^T = \tilde{U} \Sigma V^T$$

ricordando che il prodotto tra due matrici ortogonali è una matrice ortogonale: $\tilde{U}\tilde{U}^T = \tilde{U}^T\tilde{U} = I$

6.5 Un problema facile

Data $A \in \mathbb{R}^{n \times n}$ non singolare, può essere decomposta: $Ax = b \Rightarrow Q(Rx) = b \Rightarrow Qy = b$

$$\begin{cases} Qy = b \\ Rx = y \end{cases}$$

Usare la fattorizzazione QR per un sistema **rettangolare** è sicuramente preferibile rispetto a Gauss poiché la QR è più **stabile**, ma la preferiamo anche nel caso quadrato, come questo, per la propagazione degli errori, infatti:

- 1 **QR**: $k_2(R) = k_2(A)$
- 2 **Gauss**: $k_\infty(U) \leq 4^{n-1} k_\infty(A)$ ¹.

Usando LU si ha amplificazione degli errori, se non si usa il pivoting si ha ulteriore propagazione.

¹Si potrebbe pensare che quanto detto sia inutile perché stiamo confrontando due numeri di condizionamento diversi, 2 per uno e ∞ per l'altro, è vero, ma il discorso rimane valido anche con questo confronto, si usa questo solo per una semplicità di calcolo.

6.6 Minimi quadrati

Consideriamo un problema, già visto, con la QR: la matrice $A \in \mathbb{R}^{m \times n}$ con $m > n$ e rango = n (pieno), è noto che abbiamo una sola soluzione:

$$\begin{aligned} \min_x \|Ax - b\|_2^2 &= \min_x \|QRx - b\|_2^2 = \min_x \|Rx - \underbrace{Q^T b}_c\|_2^2 \\ &= \left\| \begin{bmatrix} R_1 \\ 0 \end{bmatrix} x - \begin{bmatrix} C_1 \\ C_2 \end{bmatrix} \right\|_2^2 = \left\| \begin{bmatrix} R_1 x - c_1 \\ -c_2 \end{bmatrix} \right\|_2^2 = \end{aligned}$$

... e quindi scegliendo una x tale che $R_1 x = c_1$ allora il minimo vale $\|c_2\|_2^2$.

6.7 Sistema normale

Prendiamo la stessa matrice di sopra e la trasponiamo, abbiamo ∞ soluzioni, abbiamo un sistema del tipo: $A^T y = c$, dove possiamo scomporre la matrice A come: $A = QR$, e giungiamo al problema:

$$\begin{cases} \min \|y\| \\ \min \|A^T y - c\| \end{cases} \quad \begin{cases} A^T A z = c \\ y = A z \end{cases}$$

svolgendo il prodotto delle matrici $A^T A$:

$$\begin{aligned} A^T A &= R^T Q^T Q R = R^T R = \begin{bmatrix} R_1^T & 0 \end{bmatrix} \begin{bmatrix} R_1 \\ 0 \end{bmatrix} = R_1^T R_1 \Rightarrow \\ R_1^T R_1 z &\Rightarrow R_1 z = (R_1^T)^{-1} c \Rightarrow y = Q \begin{bmatrix} R_1 \\ 0 \end{bmatrix} z = Q \begin{bmatrix} R_1 z \\ 0 \end{bmatrix} = Q \begin{bmatrix} (R_1^T)^{-1} c \\ 0 \end{bmatrix} \end{aligned}$$

6.7.1 Caratteristiche degli algoritmi

Alcune Caratteristiche degli algoritmi:

- 1 **Gauss** $\mathcal{O}(\frac{1}{3}n^3)$ più veloce di QR;
- 2 **QR** $\mathcal{O}(\frac{2}{3}n^3)$ più stabile di Gauss.
- 3 **SVD** \nexists è iterativo.

Capitolo 7

Metodi Iterativi per la risoluzione di sistemi lineari

I metodi iterativi per la risoluzione di un sistema lineare

$$\mathbf{Ax} = \mathbf{b} \quad (7.1)$$

generano, a partire da un vettore iniziale $\mathbf{x}^{(0)}$, una successione di vettori $\mathbf{x}^{(k)}$, $k = 0, 1, \dots$, che, sotto opportune ipotesi, converge alla soluzione del problema. Nei metodi diretti che valutano la soluzione in un numero *finito* di passi, gli errori presenti nei risultati nascono esclusivamente dall'aritmetica finita e/o dalla presenza di errori sui dati. In quelli iterativi, invece, agli errori sperimentali e di arrotondamento si aggiungono gli errori di **troncamento**, derivati dal fatto che il limite cercato deve essere necessariamente approssimato troncando la successione per un indice sufficientemente grande. Si aggiunge, quindi, il problema della determinazione di un *criterio di arrotondamento*.

I metodi iterativi risultano convenienti per matrici di *grandi dimensioni*, specialmente se *strutturate* o *sparse*. Infatti, mentre un metodo diretto opera modificando la matrice del sistema, spesso alterandone la struttura e aumentando il numero degli elementi non nulli, un metodo iterativo, non richiede una modifica della matrice del sistema e in alcuni casi neanche la sua memorizzazione. È solo necessario poter accedere in qualche modo ai suoi elementi.

Questi metodi possono, inoltre, essere utilizzati nei casi in cui si voglia raffinare una soluzione approssimata ottenuta con altri algoritmi o mediante informazioni a priori sul problema.

Infine è possibile ridurre il tempo di elaborazione, eseguendo un minor numero di iterazioni in quei casi in cui non sia richiesta un'elevata accuratezza.

7.1 Metodi iterativi del prim'ordine

Definizione 7.1.1: Metodo iterativo del primo ordine

Il calcolo di un termine della successione coinvolge solo il termine precedente.

assume la forma

$$\mathbf{x}^{(k+1)} = \psi(\mathbf{x}^{(k)}).$$

Definizione 7.1.2: Globalmente convergente

Un metodo si dice globalmente convergente se per ogni vettore iniziale $\mathbf{x}^{(0)} \in \mathfrak{R}^n$ si ha che

$$\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}\| = 0,$$

dove \mathbf{x} è la soluzione e $\|\cdot\|$ denota una qualsiasi norma.

Definizione 7.1.3: Consistente

Un metodo iterativo si dice **consistente** se

$$\mathbf{x}^{(k)} = \mathbf{x} \Rightarrow \mathbf{x}^{(k+1)} = \mathbf{x}.$$

Teorema 7.1.1

a consistenza è una condizione necessaria per la convergenza.

Un metodo **lineare, stazionario del primo ordine** assume la forma

$$\mathbf{x}^{(k+1)} = \mathbf{B}\mathbf{x}^{(k)} + \mathbf{f}. \quad (7.2)$$

La 7.2 è:

lineare perché è tale la relazione che la esprime;

stazionario perché la *matrice di iterazione* \mathbf{B} e il vettore \mathbf{f} non variano al variare dell'indice di iterazione k ;

del prim'ordine perché il calcolo del vettore $\mathbf{x}^{(k+1)}$ dipende solo dal termine precedente $\mathbf{x}^{(k)}$.

Teorema 7.1.2

Un metodo iterativo lineare, stazionario del prim'ordine è consistente se e solo se:

$$\mathbf{f} = (\mathbf{I} - \mathbf{B})\mathbf{A}^{-1}\mathbf{b}. \quad (7.3)$$

Definizione 7.1.4: Errore

Si definisce **errore** al passo k il vettore:

$$\mathbf{e}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}. \quad (7.4)$$

Teorema 7.1.3

Condizione sufficiente per la convergenza del metodo iterativo 7.2 è che esista una norma consistente $\|\cdot\|$ tale che $\|\mathbf{B}\| < 1$

Dimostrazione. Per le proprietà delle norme matriciali si ha che

$$\|\mathbf{e}^{(k)}\| \leq \|\mathbf{B}^k\| \|\mathbf{e}^{(0)}\| \leq \|\mathbf{B}\|^k \|\mathbf{e}^{(0)}\|,$$

quindi $\|\mathbf{B}\| < 1 \Rightarrow \|\mathbf{e}^{(k)}\| \rightarrow 0$. □

Teorema 7.1.4

n metodo iterativo lineare è convergente se e solo se il raggio spettrale $\rho(\mathbf{B})$ della matrice di iterazione \mathbf{B} è minore di 1.

Dimostrazione. Sia $\rho(\mathbf{B}) < 1$, per ogni $\epsilon > 0 \exists$ una norma di matrice tale che

$$\|\mathbf{B}\| \leq \rho(\mathbf{B}) + \epsilon.$$

Essendo $\rho(\mathbf{B})$ strettamente minore a 1, è possibile scegliere ϵ in modo tale che $\rho(\mathbf{B}) + \epsilon < 1$. Il metodo iterativo risulta allora convergente. Se fosse $\rho(\mathbf{B}) \geq 1 \Rightarrow \exists$, un autovalore λ di \mathbf{B} avente modulo maggiore o uguale a 1. Sia \mathbf{v} l'autovettore corrispondente e scegliamo il vettore iniziale $\mathbf{x}^{(0)} = \mathbf{x} + \mathbf{v}$, essendo \mathbf{x} la soluzione del sistema. In questo caso si avrebbe $\mathbf{e}^{(0)} = \mathbf{v}$ e quindi:

$$\begin{aligned} \mathbf{B}\mathbf{e}^{(0)} &= \lambda\mathbf{e}^{(0)} \quad \Rightarrow \\ \mathbf{e}^{(k)} &= \mathbf{B}^k\mathbf{e}^{(0)} = \\ &= \mathbf{B}^{k-1}\mathbf{B}\mathbf{e}^{(0)} = \mathbf{B}^{k-1}\lambda\mathbf{e}^{(0)} = \\ &= \mathbf{B}^{k-2}\mathbf{B}\lambda\mathbf{e}^{(0)} = \mathbf{B}^{k-2}\lambda\mathbf{B}\mathbf{e}^{(0)} = \mathbf{B}^{k-2}\lambda^2\mathbf{e}^{(0)} = \\ &\quad \vdots \\ &= \lambda^k\mathbf{e}^{(0)}. \end{aligned}$$

Poiché risulta $\|\mathbf{e}^{(k)}\| = |\lambda|^k \cdot \|\mathbf{e}^{(0)}\|$, nell'ipotesi che sia $|\lambda| \geq 1$ l'errore $\mathbf{e}^{(0)}$ non converge a zero. □

La convergenza è globale, quindi non dipende dal vettore iniziale $\mathbf{x}^{(0)}$.

7.2 Costruzione di metodi iterativi lineari

Una strategia spesso usata per la costruzione di metodi iterativi lineari è quella dello *splitting additivo*. Consiste nello scrivere la matrice del sistema nella forma:

$$\mathbf{A} = \mathbf{P} - \mathbf{N}, \tag{7.5}$$

dove la *matrice di preconditionamento* \mathbf{P} è non singolare. Il sistema equivalente che si ottiene:

$$\mathbf{P}\mathbf{x} = \mathbf{N}\mathbf{x} + \mathbf{b},$$

che porta alla definizione del metodo iterativo

$$\mathbf{P}\mathbf{x}^{(k+1)} = \mathbf{N}\mathbf{x}^{(k)} + \mathbf{b}. \tag{7.6}$$

Il metodo costruito è consistente, e poiché $\det(\mathbf{P}) \neq 0$ si può scrivere:

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

ponendo $\mathbf{B} = \mathbf{P}^{-1}\mathbf{N}$ e $\mathbf{f} = \mathbf{P}^{-1}\mathbf{b}$. Per i Teoremi 7.1.3 e 7.1.4 una condizione di sufficienza per la convergenza è data da $\|\mathbf{P}^{-1}\mathbf{N}\| < 1$, dove $\|\cdot\|$ è una qualsiasi norma matriciale consistente, mentre una condizione necessaria e sufficiente è $\rho(\mathbf{P}^{-1}\mathbf{N}) < 1$. Perché il metodo risulti operativo è **necessario** che la matrice \mathbf{P} sia più semplice da invertire di \mathbf{A} .

7.3 Criterio di arresto

Nell'applicazione di un metodo iterativo, risulta cruciale il criterio utilizzato per arrestare le iterazioni, in quanto da esso dipende la qualità dell'approssimazione della soluzione del sistema.

1. Un metodo per rilevare la convergenza consiste nel verificare che la successione delle iterate soddisfi il criterio di Cauchy, controllando lo scarto tra due iterazioni successive. Scegliendo:
 - una tolleranza $\tau > 0$;
 - norma vettoriale.

La condizione di stop allora risulta:

$$\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \tau, \quad (7.7)$$

o meglio:

$$\frac{\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|}{\|\mathbf{x}^{(k)}\|} \leq \tau, \quad (7.8)$$

o meglio:

$$\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \tau \|\mathbf{x}^{(k)}\|. \quad (7.9)$$

Questa condizione non sarà mai verificata se il metodo non converge.

Quando la matrice \mathbf{A} è **simmetrica, definita positiva** può essere dimostrato che l'errore al passo k di un metodo iterativo del tipo 7.2 verifica la maggiorazione:

$$\|\mathbf{e}^{(k)}\|_2 \leq \frac{1}{1 - \rho(\mathbf{B})} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|_2$$

il che mostra che il test di arresto 7.7 è affidabile solo quando il raggio spettrale della matrice di iterazione \mathbf{B} è sufficientemente inferiore a 1. Questo in generale è vero anche per una generica matrice \mathbf{A} .

La scelta di τ deve essere ragionevole, dipende dalla precisione richiesta nel problema in esame, ma che sia superiore alla precisione di macchina.

2. Per considerare il caso in cui il metodo **non converga**, per evitare quell'effetto spiacevole di un *loop* ∞ è cosa buona e giusta fissare un limite massimo di iterazioni N .

3. Nei metodi iterativi per la risoluzione di sistemi lineari al passo k è possibile ottenere una maggiorazione per l'errore $\mathbf{e}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}$ in termini del *vettore residuo* $\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}$.

Può essere scritto:

$$\begin{aligned}\mathbf{A}\mathbf{e}^{(k)} &= \mathbf{A}\mathbf{x}^{(k)} - \mathbf{A}\mathbf{x} = \\ &= \mathbf{A}\mathbf{x}^{(k)} - \mathbf{b} = \\ &= -\mathbf{r}^{(k)},\end{aligned}$$

da cui si ricava:

$$\|\mathbf{e}^{(k)}\| = \|\mathbf{A}^{-1}\mathbf{r}^{(k)}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{r}^{(k)}\|.$$

Ricordando poi che $\|\mathbf{b}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$, si arriva alla seguente maggiorazione per l'errore relativo al passo k :

$$\frac{\|\mathbf{e}^{(k)}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\mathbf{r}^{(k)}\|}{\|\mathbf{b}\|} = \kappa(\mathbf{A}) \frac{\|\mathbf{r}^{(k)}\|}{\|\mathbf{b}\|}.$$

Il criterio di stop:

$$\frac{\|\mathbf{r}^{(k)}\|}{\|\mathbf{b}\|} \leq \tau \tag{7.10}$$

la stima dell'errore fornita dal residuo normalizzato è tanto attendibile quanto più la matrice \mathbf{A} è mal condizionata.

In generale le condizioni di arresto per un criterio sono una combinazione di più criteri di arresto.

7.4 Precondizionamento

I metodi analizzati sono del tipo:

$$\mathbf{x}^{(k+1)} = \mathbf{B}\mathbf{x}^{(k)} + \mathbf{f}, \tag{7.11}$$

sfruttando la 7.5 a pagina 45 si ha:

$$\begin{aligned}\mathbf{B} &= \mathbf{P}^{-1}\mathbf{N} \\ &= \mathbf{P}^{-1}(\mathbf{P} - \mathbf{A}) \\ &= (\mathbf{I} - \mathbf{P}^{-1}\mathbf{A}), \\ \mathbf{f} &= \mathbf{P}^{-1}\mathbf{b}.\end{aligned}$$

Per la consistenza, ricordando la 7.1 e il teorema 7.1.2, questa relazione equivale al sistema lineare

$$(\mathbf{I} - \mathbf{B})\mathbf{x} = \mathbf{f},$$

ossia:

$$\begin{aligned}(\mathbf{I} - \mathbf{B})\mathbf{x} &= \mathbf{f} \\ (\mathbf{I} - \mathbf{P}^{-1}\mathbf{N})\mathbf{x} &= \mathbf{P}^{-1}\mathbf{b} \\ (\mathbf{I} - \mathbf{I} + \mathbf{P}^{-1}\mathbf{A})\mathbf{x} &= \mathbf{P}^{-1}\mathbf{b} \\ \mathbf{P}^{-1}\mathbf{A}\mathbf{x} &= \mathbf{P}^{-1}\mathbf{b}.\end{aligned}$$

Questo equivale a dire che il metodo iterativo risolve il sistema lineare **precondizionato** della matrice \mathbf{P} .

La scelta di un buon precondizionamento è un compito tutt'altro che facile, d'altronde sono possibile ∞ scelte. Per matrici di dimensione molto elevate si cerca di introdurre precondizionatori che preservino la struttura o la sparsità di \mathbf{A} .

Per accelerare la convergenza del metodo, è necessario fare in modo che \mathbf{P} :

1. approssimi la matrice \mathbf{A} , cosicché $\mathbf{P}^{-1}\mathbf{A} \simeq \mathbf{I}$,
2. facile da invertire.

Ricordano il teorema 7.1.4: è condizione necessaria e sufficiente per la convergenza del metodo iterativo 7.11 è

$$\rho(\mathbf{B}) = \rho(\mathbf{I} - \mathbf{P}^{-1}\mathbf{A}) < 1.$$

Indicando con μ_i e \mathbf{v}_i gli autovalori e gli autovettori della matrice $\mathbf{P}^{-1}\mathbf{A}$ si ha che:

$$(\mathbf{I} - \mathbf{P}^{-1}\mathbf{A})\mathbf{v}_i = (1 - \mu_i)\mathbf{v}_i,$$

quindi la matrice di iterazione \mathbf{B} ha autovalori $(1 - \mu_i)$ e autovettori \mathbf{v}_i . La condizione per convergenza diventa quindi:

$$\max_{i=1, \dots, n} |1 - \mu_i| < 1$$

cioè

$$0 < \mu_i < 2, \quad i = 1, \dots, n.$$

Dal momento che la convergenza sarà tanto più veloce quanto più il raggio spettrale di \mathbf{B} è vicino a zero, questo significa che un buon precondizionatore deve far sì che gli autovalori della matrice precondizionata $\mathbf{P}^{-1}\mathbf{A}$ presentino un **cluster** a 1, cioè che si accumulino in un intorno di 1.

I metodi lineari del prim'ordine possono essere espressi in funzione del residuo $\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}$. Fissato lo *splitting* 7.5

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

essendo $\mathbf{N} = \mathbf{P} - \mathbf{A}$,

$$\begin{aligned} \mathbf{x}^{(k+1)} &= (\mathbf{I} - \mathbf{P}^{-1}\mathbf{A})\mathbf{x}^{(k)} + \mathbf{P}^{-1}\mathbf{b} = \\ &= \mathbf{x}^{(k)} + \mathbf{P}^{-1}(-\mathbf{A}\mathbf{x}^{(k)} + \mathbf{b}) = \\ &= \mathbf{x}^{(k)} + \mathbf{P}^{-1}\mathbf{r}^{(k)}. \end{aligned} \tag{7.12}$$

Dalla 7.12 si vede che **se fosse** $\mathbf{P} = \mathbf{A}$ si avrebbe:

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

la convergenza avverrebbe in un passo. Quando $\mathbf{P} \simeq \mathbf{A}$, il **residuo precondizionato** $\mathbf{P}^{-1}\mathbf{r}^{(k)}$ fornisce solo un'approssimazione dell'errore.

Per motivi numerici, il metodo viene applicato risolvendo il sistema:

$$\mathbf{P}\mathbf{z}^{(k)} = \mathbf{r}^{(k)},$$

senza invertire esplicitamente \mathbf{P} e calcolando poi:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{z}^{(k)}.$$

7.5 Iterazioni in sottospazi di Krylov

Nota storica 7.5.1: Krylov



Nikolay Mitrofanovich Krylov (29 Novembre 1879 San Pietroburgo – Maggio 11, 1955 Mosca) è stato un matematico russo noto per i suoi lavori sulla interpolazione, meccanica non-lineare e sulla fisica matematica. [Wiki17g].

Sia $A \in \mathbb{R}^n$. Per $r = 1, 2, \dots$, si definisce **sottospazio di Krylov** con indice r , generato da A e da \mathbf{b} , l'insieme:

$$\mathcal{K}_r := \text{span}(\mathbf{b}, A\mathbf{b}, \dots, A^{r-1}\mathbf{b}) \subseteq \mathbb{R}^n.$$

I **metodi di proiezione in sottospazi di Krylov** determinano, per ogni r , una soluzione approssimata $\mathbf{x}^{(r)}$ del sistema lineare $A\mathbf{x} = \mathbf{b}$ che appartenga al sottospazio \mathcal{K}_r . La scelta del vettore $\mathbf{x}^{(r)}$ può essere fatta applicando diversi criteri di ottimalità, che conducono a differenti algoritmi. Se i vettori $A^k\mathbf{b}$, $k = 0, 1, \dots, n-1$, risultassero tutti indipendenti si avrebbe

$$\mathcal{K}_n \equiv \mathbb{R}^n$$

e l'algoritmo terminerebbe in n passi. Se questa condizione non fosse verificata potrebbe avvenire che per un dato $l < n$ il vettore $A^l\mathbf{b}$ risulti linearmente indipendente dai vettori precedenti. In questo caso esisterebbero scalari $\alpha_i, i = 0, \dots, l$ non tutti nulli tali che

$$\sum_{i=0}^l \alpha_i A^i \mathbf{b} = \mathbf{0}.$$

Supponendo $\alpha_0 \neq 0$, è possibile esplicitare in tale relazione il vettore \mathbf{b} , ottenendo:

$$A\mathbf{x}^* = \mathbf{b} = -\frac{1}{\alpha_0} \sum_{i=1}^l \alpha_i A^i \mathbf{b} = A \left(-\frac{1}{\alpha_0} \sum_{i=1}^l \alpha_i A^{i-1} \mathbf{b} \right),$$

dove \mathbf{x}^* indica la soluzione esatta del sistema lineare, il che implica

$$\mathbf{x}^* = -\frac{1}{\alpha_0} \sum_{i=1}^l \alpha_i A^{i-1} \mathbf{b} \in \mathcal{K}_l$$

il che implica che il metodo terminerebbe in l passi. Questi metodi terminando a n passi vengono comunque considerati iterativi (*pseudo-iterativi*) in quanto, se opportunamente preconditionati forniscono una soluzione sufficientemente precisa in un numero inferiore ad n .

7.5.1 Il gradiente coniugato come metodo di Krylov

Il metodo del gradiente coniugato è costituito da una iterazione in sottospazi di Krylov, come riportato dal teorema:

Teorema 7.5.1

ia $\mathbf{Ax} = \mathbf{b}$ con A simmetrica definita positiva. Se inizializzato con $\mathbf{x}^{(0)} = \mathbf{0}$, allora il metodo del gradiente coniugato procede fintantoché il residuo $\mathbf{r}^{(k)}$ non si annulla e valgono le seguenti identità per $l = 1, \dots, k$

$$\begin{aligned}\mathcal{K}_l &:= \text{span}(\mathbf{b}, \mathbf{Ab}, \dots, \mathbf{A}^{l-1}\mathbf{b}) = \text{span}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(l)}) \\ &= \text{span}(\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(l-1)}) = \text{span}(\mathbf{r}^{(0)}, \dots, \mathbf{r}^{(l-1)}).\end{aligned}$$

Per ogni l si ha:

$$\mathbf{r}^{(l)T}\mathbf{r}^{(j)} = 0 \quad \text{e} \quad \mathbf{p}^{(l)T}\mathbf{Ap}^{(j)} = 0, \quad j = 0, 1, \dots, l-1.$$

il teorema 7.5.1 mostra che il metodo del gradiente coniugato opera sostituendo la base \mathcal{K}_l numericamente instabile

$$\{\mathbf{b}, \mathbf{Ab}, \dots, \mathbf{A}^{l-1}\mathbf{b}\}$$

con quella formata dalle direzioni A -coniugate $\{\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(l-1)}\}$. La generica iterata $\mathbf{x}^{(l)}$ del metodo del gradiente coniugato ha una proprietà di ottimalità. Se la matrice A è simmetrica definita positiva la relazione

$$\|\mathbf{y}\|_A = \sqrt{\mathbf{y}^T \mathbf{A} \mathbf{y}}$$

definisce una norma matriciale di \mathbb{R}^n .

Teorema 7.5.2

Il vettore $\mathbf{x}^{(l)}$ minimizza in \mathcal{K}_l la norma- A dell'errore cioè

$$\mathbf{x}^{(l)} = \arg \min_{\mathbf{x} \in \mathcal{K}_l} \|\mathbf{x}^* - \mathbf{x}\|_A$$

essendo \mathbf{x}^* la soluzione esatta del sistema lineare. Inoltre se $\mathbf{e}^{(l)} = \mathbf{x}^* - \mathbf{x}^{(l)}$ allora per ogni l si ha

$$\|\mathbf{e}^{(l)}\|_A \leq \|\mathbf{e}^{(l-1)}\|_A$$

e si ottiene $\mathbf{e}^{(k)} = \mathbf{0}$ per ogni $k \leq n$.

7.5.2 L'iterazione di Arnoldi

Nota storica 7.5.2: Arnoldi

Walter Edwin Arnoldi (New York, Dicembre 14, 1917 - Ottobre 5, 1995) è stato un ingegnere meccanico Americano, principalmente noto per l'iterazione di Arnoldi, un algoritmo usato nell'algebra lineare numerica. [Wiki16a].

L'iterazione di Arnoldi consente di costruire una base $\{\mathbf{q}_1, \dots, \mathbf{q}_l\}$ per il sottospazio \mathcal{K}_l formata da vettori tra loro ortonormali, cioè tali che $\mathbf{q}_i^T \mathbf{q}_j = \delta_{ij}$, sostituendo così la base numericamente instabile $\{\mathbf{b}\}, \mathbf{Ab}, \dots, \mathbf{A}^{l-1}\mathbf{b}$. Ciò viene fatto applicando il metodo di ortogonalizzazione di

Gram-Schmidt:

$$\mathbf{q}_1 = \frac{\mathbf{b}}{\|\mathbf{b}\|},$$

$$\tilde{\mathbf{q}}_{s+1} = \mathbf{A}^s \mathbf{b} - \sum_{i=1}^s h_{is} \mathbf{q}_i, \quad s = 1, \dots, l-1$$

$$\mathbf{q}_{s+1} = \frac{\tilde{\mathbf{q}}_{s+1}}{\|\tilde{\mathbf{q}}_{s+1}\|}$$

dove le norme vettoriali utilizzate sono euclidee e gli scalari h_{is} possono essere ricavati imponendo l'ortogonalità di $\tilde{\mathbf{q}}_{s+1}$ con i vettori precedentemente calcolati. Per motivi numerici, per la costruzione di $\tilde{\mathbf{q}}_{s+1}$ si utilizza il vettore $\mathbf{A}\mathbf{q}_s$ in luogo di $\mathbf{A}^s \mathbf{b}$. Infatti, pur trattandosi di due differenti vettori, questa scelta consente di passare da una base di \mathcal{K}_s ad una di \mathcal{K}_{s+1} mediante un algoritmo più stabile. Si giunge così alle seguenti formule, che consentono di costruire una base ortonormale per \mathcal{K}_{l+1} :

$$\mathbf{q}_1 = \frac{\mathbf{b}}{\|\mathbf{b}\|},$$

$$\tilde{\mathbf{q}}_{s+1} = \mathbf{A}\mathbf{q}_s - \sum_{i=1}^s h_{is} \mathbf{q}_i, \quad s = 1, \dots, l$$

$$\mathbf{q}_{s+1} = \frac{\tilde{\mathbf{q}}_{s+1}}{h_{s+1,s}}, \quad s = 1, \dots, l,$$

$$h_{is} = \mathbf{q}_i^T \mathbf{A}\mathbf{q}_s, \quad i = 1, \dots, s,$$

$$h_{s+1,s} = \|\tilde{\mathbf{q}}_{s+1}\|.$$

Di seguito l'Algoritmo:

Algoritmo: Iterazione di Arnoldi

Input: A, \mathbf{b}, τ

Output:

1. $\mathbf{q}_1 = \mathbf{b}/\|\mathbf{b}\|$
2. $s = 1$
3. repeat
 - $\mathbf{v} = \mathbf{A}\mathbf{q}_s$
 - for $i = 1, \dots, s$
 - ♣ $h_{is} = \mathbf{q}_i^T \mathbf{v}$
 - ♣ $\mathbf{v} = \mathbf{v} - h_{is} \mathbf{q}_i$
 - $h_{s+1,s} = \|\mathbf{v}\|$
 - $\mathbf{q}_{s+1} = \mathbf{v}/h_{s+1,s}$
 - $s = s + 1$
4. until $s = n + 1$ OR $|h_{s,s-1}| < \tau$

L'annullarsi, per un dato $l < n$, della norma $h_{l+1,l}$ interrompe l'algoritmo si parla allora di *breakdown*, e indica che il vettore $\mathbf{A}\mathbf{q}_l$ è **linearmente dipendente** dai precedenti e pertanto lo spazio \mathcal{K}_l contiene la soluzione del sistema lineare $\mathbf{A}\mathbf{x} = \mathbf{b}$.

Una delle formule trovate può essere riscritta come:

$$\mathbf{A}\mathbf{q}_s = \sum_{i=1}^{s+1} h_{i,s} \mathbf{q}_i, \quad s = 1, \dots, l,$$

esplicitando:

$$\begin{aligned} s = 1 \quad \mathbf{A}\mathbf{q}_1 &= \sum_{i=1}^2 h_{i,1} \mathbf{q}_i = h_{1,1} \mathbf{q}_1 + h_{2,1} \mathbf{q}_2 \\ s = 2 \quad \mathbf{A}\mathbf{q}_2 &= \sum_{i=1}^3 h_{i,2} \mathbf{q}_i = h_{1,2} \mathbf{q}_1 + h_{2,2} \mathbf{q}_2 + h_{3,2} \mathbf{q}_3 \\ &\vdots \\ s = l \quad \mathbf{A}\mathbf{q}_l &= \sum_{i=1}^{l+1} h_{i,l} \mathbf{q}_i = h_{1,l} \mathbf{q}_1 + h_{2,l} \mathbf{q}_2 + \dots + h_{l+1,l} \mathbf{q}_{l+1} \end{aligned}$$

si vede che può essere riscritta in maniera compatta:

$$\mathbf{A}\mathbf{Q}_l = \mathbf{Q}_{l+1} \tilde{\mathbf{H}}_l. \quad (7.13)$$

Dove $\mathbf{Q} = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_l] \in \mathbb{R}^{n \times l}$ e

$$\tilde{\mathbf{H}}_l = \begin{bmatrix} h_{1,1} & h_{1,2} & \dots & \dots & h_{1,l} \\ h_{2,1} & h_{2,2} & \dots & \dots & h_{2,l} \\ & h_{3,2} & \ddots & & h_{3,l} \\ & & \ddots & \ddots & \vdots \\ & & & h_{l,l-1} & h_{l,l} \\ & & & & h_{l+1,l} \end{bmatrix} \in \mathbb{R}^{(l+1) \times l}.$$

Dalla 7.13 si ottiene:

$$\mathbf{Q}_l^T \mathbf{A} \mathbf{Q}_l = \mathbf{H}_l := \begin{bmatrix} h_{1,1} & h_{1,2} & \dots & \dots & h_{1,l} \\ h_{2,1} & h_{2,2} & \dots & \dots & h_{2,l} \\ & h_{3,2} & \ddots & & h_{3,l} \\ & & \ddots & \ddots & \vdots \\ & & & h_{l,l-1} & h_{l,l} \end{bmatrix} \in \mathbb{R}^{(l) \times l}.$$

Che per $l = n$ diventa:

$$\mathbf{Q}_n^T \mathbf{A} \mathbf{Q}_n = \mathbf{H}_n,$$

l'iterazione di Arnoldi porta alla matrice di Hessemberg mediante n trasformazioni ortogonali. Essendo le matrici \mathbf{A} e \mathbf{H}_n unitariamente simili, questo potrebbe essere visto come un passo preliminare all'applicazione di un metodo per la ricerca degli autovalori che tragga vantaggio dalla struttura di Hessemberg, come l'algoritmo QR.

in realtà la matrice H_l contiene importanti informazioni sulla matrice A anche per l molto minore di n . In molti casi gli autovalori estremali di H_l convergono rapidamente agli autovalori estremali di A e continuando con le iterazioni un numero sempre maggiore di autovalori delle due matrici tende a coincidere. Questo consente, nel caso di matrici di grandi dimensioni di ottenere dopo *poche* iterazioni approssimazioni sufficientemente accurate degli autovalori estremali dalla matrice in esame, i quali spesso consentono di determinare importanti informazioni sulla matrice stessa (raggio spettrale, condizionamento, etc.) o sul problema fisico da cui la matrice deriva.

7.5.3 L'iterazione di Lanczos

Nota storica 7.5.3: Lanczos



Cornelius (Cornel) Lanczos era un matematico e fisico ungherese ebreo, nato il 2 febbraio 1893 (Székesfehérvár) e morto il 25 giugno 1974 (Budapest). [Wiki17h].

Supponendo di applicare l'iterazione di Arnoldi ad una matrice A simmetrica. Per questo particolare caso si ricava, che la matrice $H_l = Q_l^T A Q_l$, detta matrice di *Ritz*, dovendo essere simultaneamente simmetrica e in forma di Hessemberg, deve necessariamente essere tridiagonale. Essa sarà tale quindi che i suoi autovalori, i *valori di Ritz*, sono reali. Il metodo di Lanczos consente di ottenere una fattorizzazione $AQ_l = Q_{l+1}\tilde{H}_l$ dove $Q_l = [q_1, \dots, q_l] \in \mathbb{R}^{n \times l}$ è costituita da colonne ortonormali che costituiscono una base

per il sottospazio di Krylov di dimensione l .

Ponendo $h_{i,i} = \alpha_i$ e $h_{i+1,i} = h_{i,i+1} = \beta_i$ possiamo scrivere la matrice H_l :

$$H_l = \begin{bmatrix} \alpha_1 & \beta_1 & & & \\ \beta_1 & \alpha_2 & \ddots & & \\ & \ddots & \ddots & \beta_{l-1} & \\ & & & \beta_{l-1} & \alpha_l \end{bmatrix}$$

Lo sfruttamento della particolare struttura della matrice porta ad un abbattimento della complessità computazionale. Il metodo, che prende il nome di **iterazione di Lanczos**, assume la forma:

$$\begin{aligned} q_1 &= \frac{\mathbf{b}}{\|\mathbf{b}\|}, \\ \tilde{q}_{s+1} &= Aq_s - \alpha_s q_s - \beta_{s-1} q_{s-1}, \quad s = 1, \dots, l \\ q_{s+1} &= \frac{\tilde{q}_{s+1}}{\beta_s}, \quad s = 1, \dots, l, \end{aligned}$$

con

$$\begin{aligned} \alpha_s &= q_s^T A q_s, \\ \beta_s &= \|\tilde{q}_{s+1}\|. \end{aligned}$$

Nell'algoritmo, se si dovesse verificare, per un certo passo $s < n$ un *breakdown* in output avrà le matrici Q e H ridotte e cioè di dimensioni $n \times s$ e $s \times s$ rispettivamente, dove s è l'indice dove si verifica il breakdown.

Il presentarsi del breakdown nel caso dell'applicazione del metodo alla risoluzione di un sistema lineare $\mathbf{Ax} = \mathbf{b}$ implica che il sottospazio di Krylov \mathcal{K}_l conterrà la soluzione esatta del sistema lineare simmetrico.

Algoritmo: Iterazione di Lanczos

Input: A, \mathbf{b}, τ

Output:

1. $\beta_0 = 0, \mathbf{q}_0 = \mathbf{0}$
2. $\mathbf{q}_1 = \mathbf{b}/\|\mathbf{b}\|$
3. $s = 1$
4. repeat
 - $\mathbf{v} = \mathbf{A}\mathbf{q}_s$
 - $\alpha_s = \mathbf{q}_s^T \mathbf{v}$
 - $\mathbf{v} = \mathbf{v} - \alpha_s \mathbf{q}_s - \beta_{s-1} \mathbf{q}_{s-1}$
 - $\beta_s = \|\mathbf{v}\|$
 - $\mathbf{q}_{s+1} = \mathbf{v}/\beta_s$
 - $s = s + 1$
5. until $s = n + 1$ OR $|h_{s,s-1}| < \tau$

7.5.4 Il metodo GMRES

Dato un sistema lineare:

$$\mathbf{Ax} = \mathbf{b},$$

con \mathbf{A} non singolare. I metodi GMRES (*Generalized Minimal RESidual*) consiste nel costruire iterativamente una base ortonormale per lo spazio di Krylov \mathcal{K}_l , per $l = 1, \dots, n$, approssimando, ad ogni iterazione, la soluzione del sistema \mathbf{x}^* col vettore $\mathbf{x}^{(l)} \in \mathcal{K}_l$ che minimizza la norma-2 del residuo: $\mathbf{r}^{(l)} = \mathbf{b} - \mathbf{Ax}^{(l)}$.

Ponendo

$$\mathbf{x}^{(l)} = \arg \min_{\mathbf{x} \in \mathcal{K}_l} \|\mathbf{b} - \mathbf{Ax}\|_2 \quad (7.14)$$

proiettiamo la soluzione nello spazio di Krylov \mathcal{K}_l e la valutiamo risolvendo un problema ai minimi quadrati. Pur trattandosi di un metodo che termina la più in n passi, le proprietà dell'iterazione di Arnoldi fanno sì che in molti casi si ottenga un'approssimazione sufficientemente accurata della soluzione del sistema in un numero di passi molto inferiore a n .

Un qualunque vettore di \mathcal{K}_l può essere espresso nella base ortonormale $\{ \mathbf{q}_1, \dots, \mathbf{q}_l \}$ nella forma $\mathbf{Q}_l \mathbf{y}$, essendo \mathbf{y} un generico vettore di \mathbb{R}^l . Dunque

$$\begin{aligned} \|\mathbf{b} - \mathbf{Ax}\| &= \|\mathbf{AQ}_l \mathbf{y} - \mathbf{b}\| = \\ &= \|\mathbf{Q}_{l+1} \tilde{\mathbf{H}}_l \mathbf{y} - \mathbf{b}\| = \\ &= \|\mathbf{Q}_{l+1} \tilde{\mathbf{H}}_l \mathbf{y} - \mathbf{Q}_{l+1} \mathbf{Q}_{l+1}^T \mathbf{b}\| = \\ &= \|\mathbf{Q}_{l+1} (\tilde{\mathbf{H}}_l \mathbf{y} - \mathbf{Q}_{l+1}^T \mathbf{b})\| = \\ &= \|\tilde{\mathbf{H}}_l \mathbf{y} - \|b\| \mathbf{e}_1\|. \end{aligned}$$

Avendo sfruttato la 7.13 e ricordando che $\mathbf{Q}_{l+1}^T \mathbf{Q}_{l+1} = \mathbf{I}_{l+1}$, $\mathbf{Q}_{l+1}^T \mathbf{b} = \|b\| \mathbf{e}_1$ dove \mathbf{e}_1 il vettore 1 della base canonica: $\mathbf{e}_1 = [1, 0, \dots, 0]^T$.

Per il problema 7.14 è possibile calcolare la soluzione $\mathbf{y}^{(l)}$ di:

$$\min_{\mathbf{y} \in \mathbb{R}^l} = \|\tilde{\mathbf{H}}_l \mathbf{y} - \|b\| \mathbf{e}_1\|,$$

ponendo successivamente $\mathbf{x}^{(l)} = \mathbf{Q} \mathbf{y}^{(l)}$. Questo problema di minimo oltre a essere non vincolato, è caratterizzato da una matrice dei coefficienti in $\mathbb{R}^{(l+1) \times l}$, quindi di dimensioni contenute fintantoché l'indice di iterazione non cresce. La forma di Hessemberg rende agevole la soluzione del problema ai minimi quadrati.

Capitolo 8

Risultati ottenuti

Le simulazioni svolte sono centrate sulla risoluzione dei sistemi lineari equazione 7.1, per ciascuna si procede in tal modo:

- 1 si prende una matrice \mathbf{A} random,
- 2 si fissa il vettore soluzione $\hat{\mathbf{x}}$ come vettore di tutti *uno*,
- 3 si *costruisce* il vettore di termini noti $\mathbf{b} = \mathbf{A}\hat{\mathbf{x}}$,
- 4 si risolve il sistema lineare: $\mathbf{Ax} = \mathbf{b}$, si trova la \mathbf{x} (che è diversa da $\hat{\mathbf{x}}$),
- 5 si calcola l'errore: $\|\mathbf{x} - \hat{\mathbf{x}}\|$,
- 6 si calcola anche il tempo di calcolo con in comandi: `tic`, `toc`,
- 7 si confrontano tempi e errori tra gli algoritmi presentati e gli algoritmi di Matlab[®] .

Si riportano anche i risultati sulla fattorizzazione QR.

Ogni **simulazione** può essere divisa in 3 **fasi** distinte: `preprocessor`, `processor`, `postprocessor`.

La prima è quella che comprende i dati iniziali, gli input: la dimensione massima della matrice, i margini di errore che si vuole commettere etc.

La seconda: la vera simulazione, nel codice qui sono sicuramente presenti: gli algoritmi veri e propri, e sicuramente le `tic`, `toc` utilizzate per calcolare i tempi.

La terza scrittura su disco i risultati ottenuti esportazione dei risultati in file **leggibile** da L^AT_EX, per l'importazione.

8.1 Le simulazioni fattorizzazione Cholesky

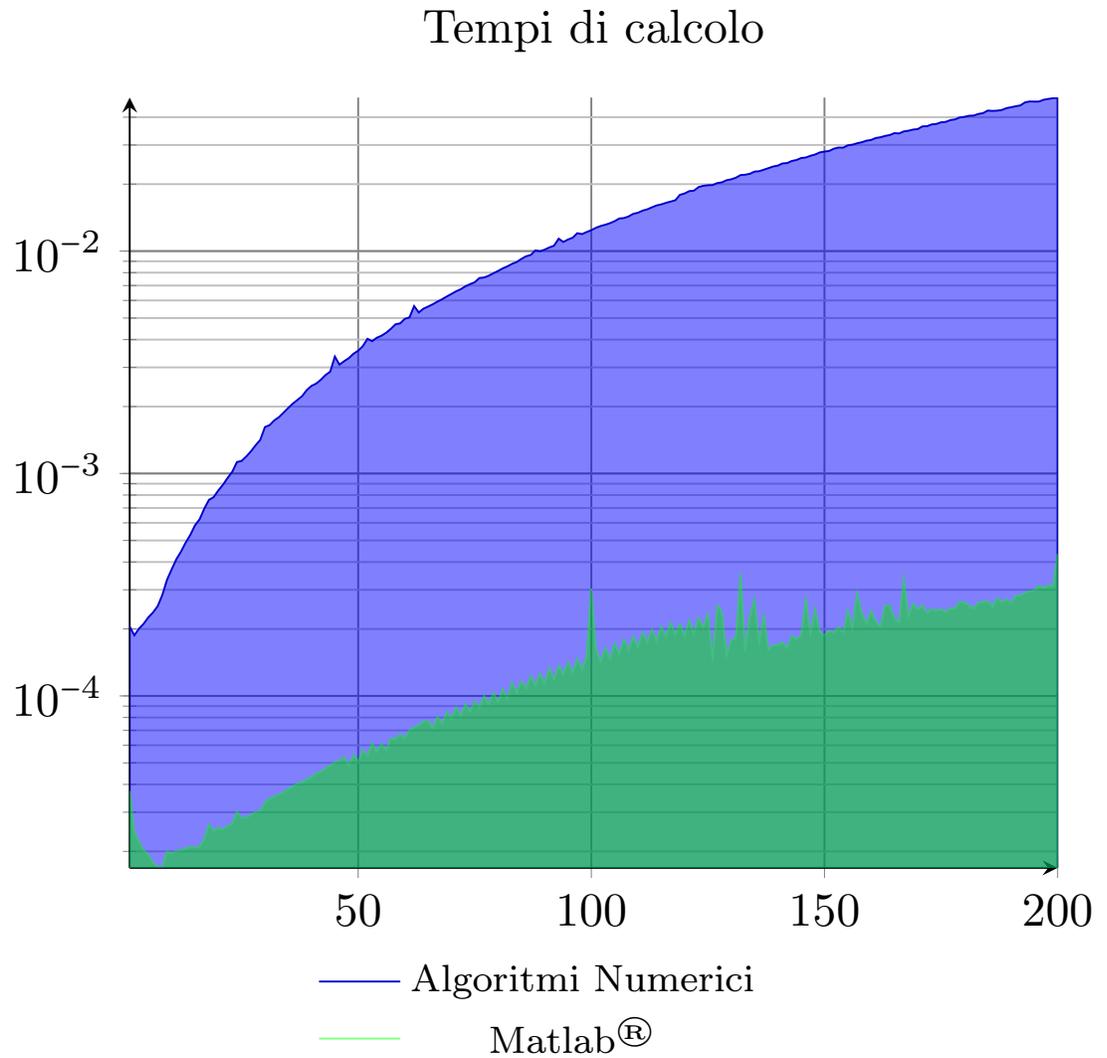


Figura 8.1: Fattorizzazione Cholesky per colonne. Vengono confrontati i tempi di calcolo di due Algoritmi, l'Algoritmo implementato e l'Algoritmo `chol()` di Matlab®. Si risolve il sistema lineare, considerando una matrice piena casuale random, si risolvono così due sistemi uno triangolare inferiore e uno triangolare superiore. Nel grafico in ascisse è riportata la dimensione della matrice da 3 a 200, in ordinate in scala logaritmica i tempi di calcolo in [s].

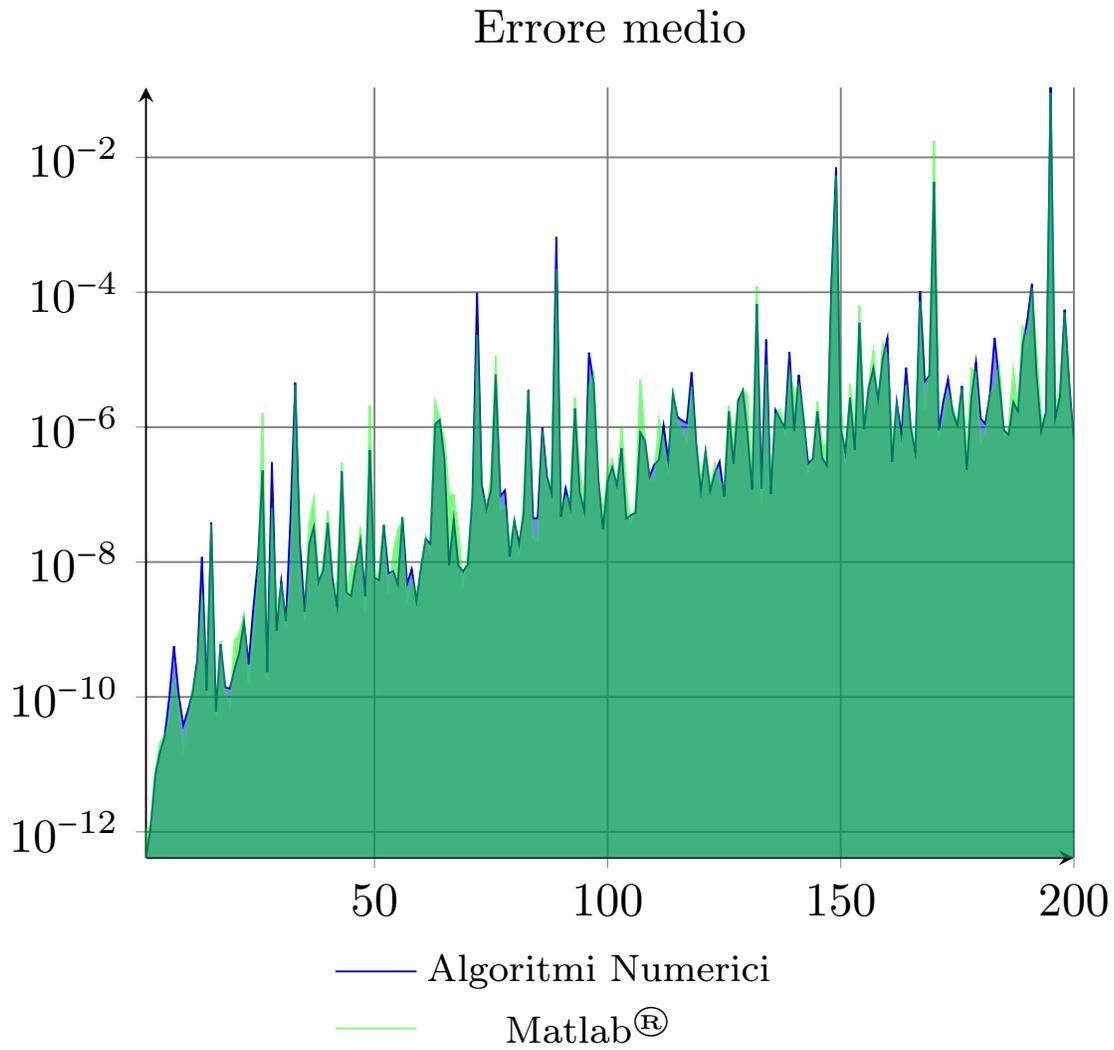


Figura 8.2: Fattorizzazione Cholesky. In ascisse sono riportate le dimensioni della matrice A da 3 a 200, in ordinate, in scala logaritmica è riportata la norma dell'errore calcolato come differenza tra soluzione calcolata e soluzione esatta.

8.2 Le simulazioni fattorizzazione QR

Prima di procedere con la sperimentazione sul calcolo dei minimi quadrati con fattorizzazione QR si sono analizzate e confrontate le fattorizzazioni QR: Householder, Givens, CGD, MGS, Matlab[®]. Si riportano le singole fattorizzazioni confrontate **tutte** con `qr()`. In tutte le simulazioni le matrici sono sempre prese quadrate.

8.2.1 Householder

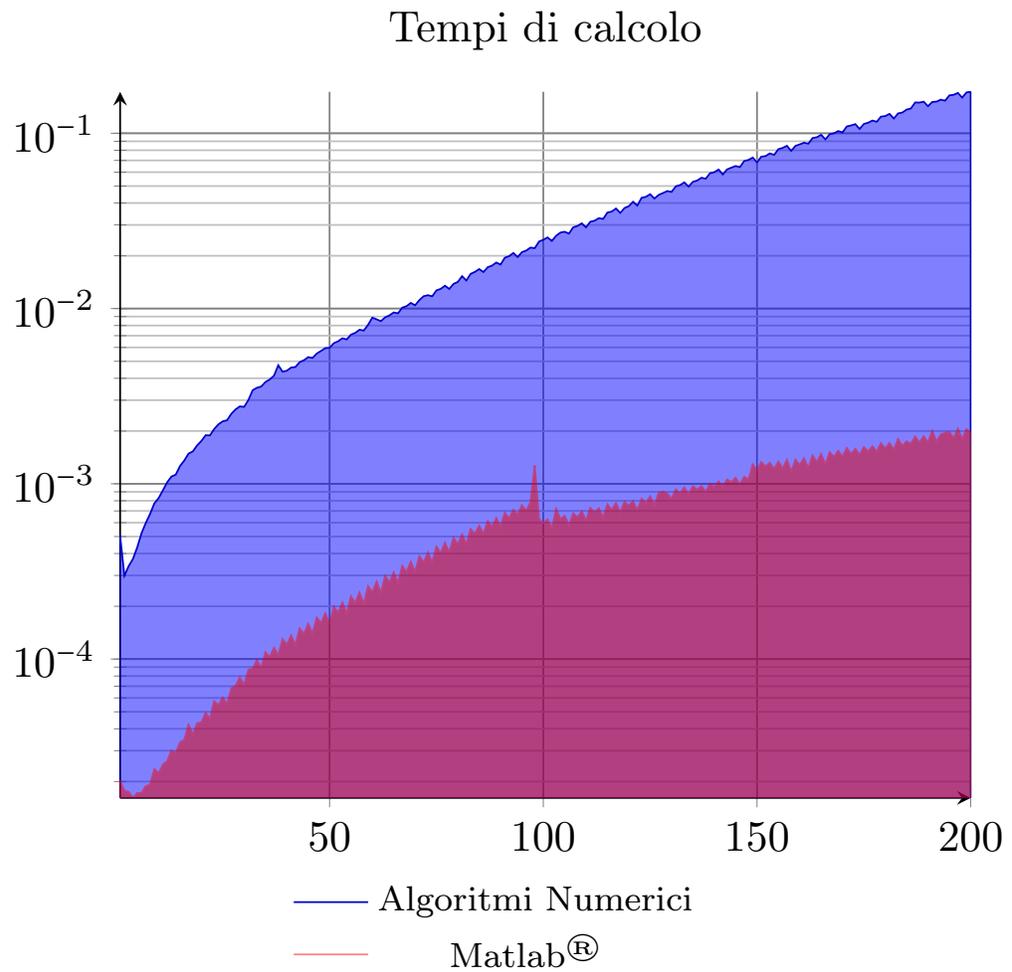


Figura 8.3: Fattorizzazione QR di Householder ($m = n$). Vengono confrontati i tempi di calcolo di due Algoritmi, l'Algoritmo è quello di pagina 102 libro di testo [Rod08] e l'Algoritmo `qr()` di Matlab[®].

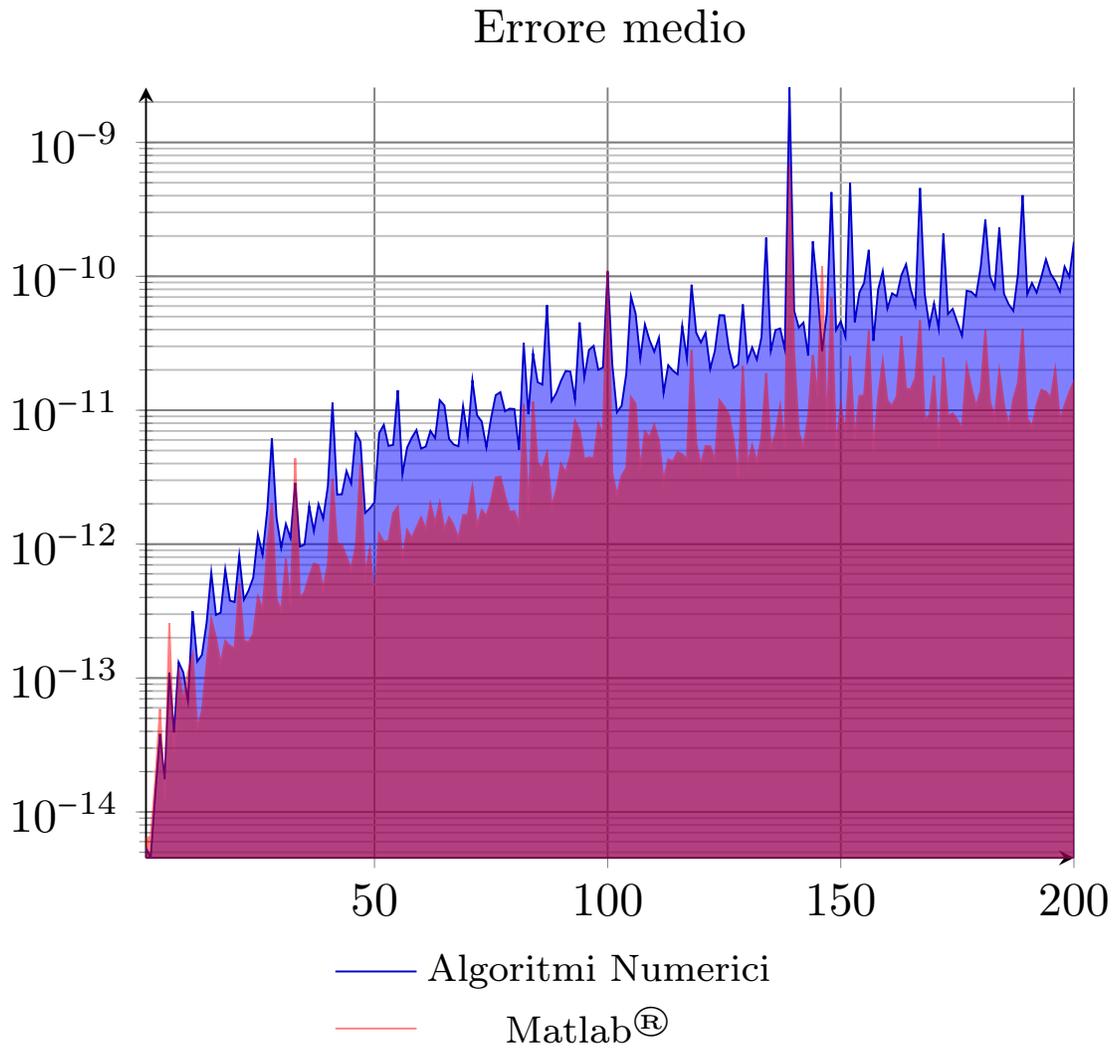


Figura 8.4: Fattorizzazione QR di Householder ($m = n$). Errori commessi.

8.2.2 Givens

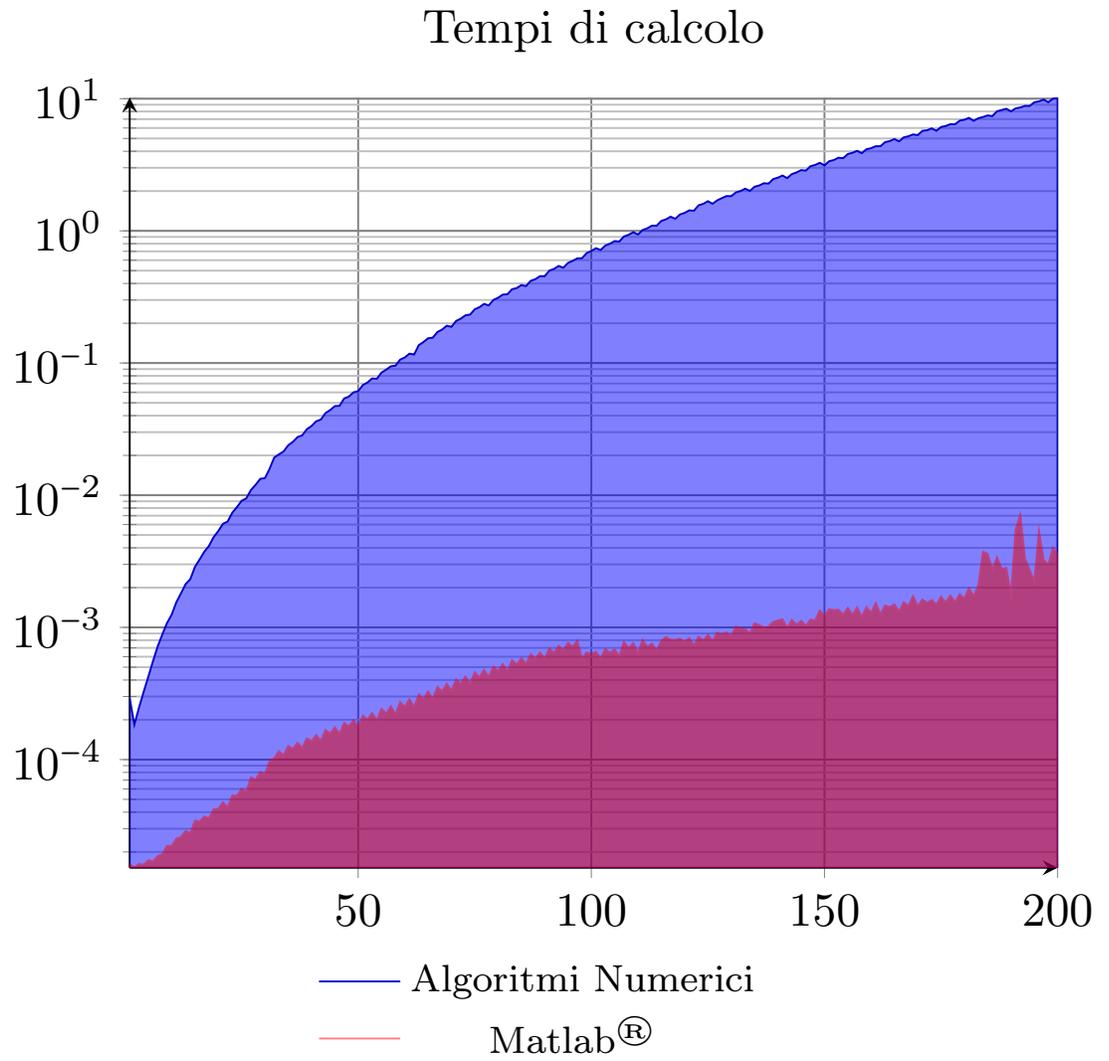


Figura 8.5: Fattorizzazione QR di Givens ($m = n$). Tempi di calcolo.

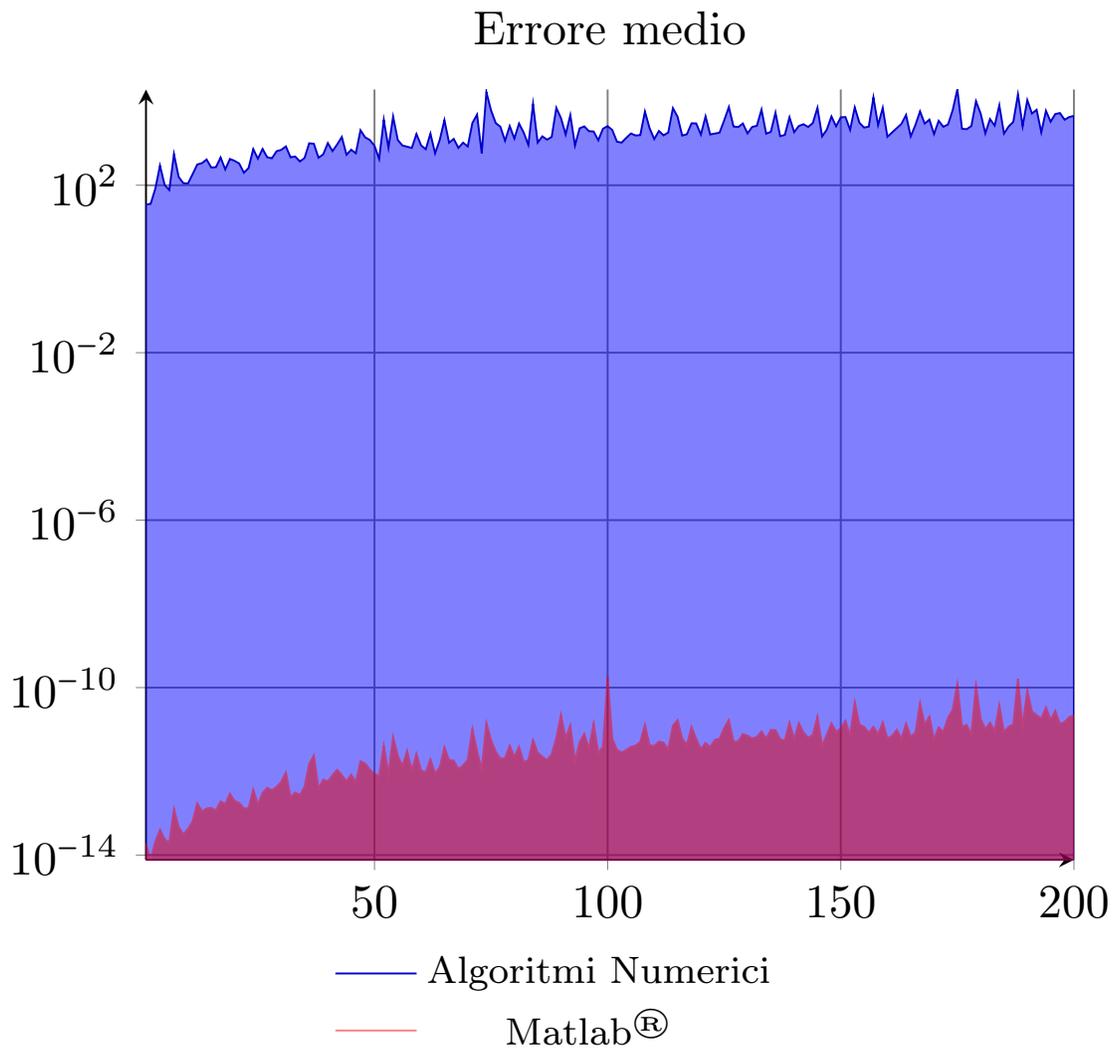


Figura 8.6: Fattorizzazione QR di Givens ($m = n$). Errori di calcolo.

8.2.3 Gram-Schmidt

Per questo caso,ma solo per questo caso, si è simulato **solo** classic Gram-Schmidt.

Tempi di calcolo

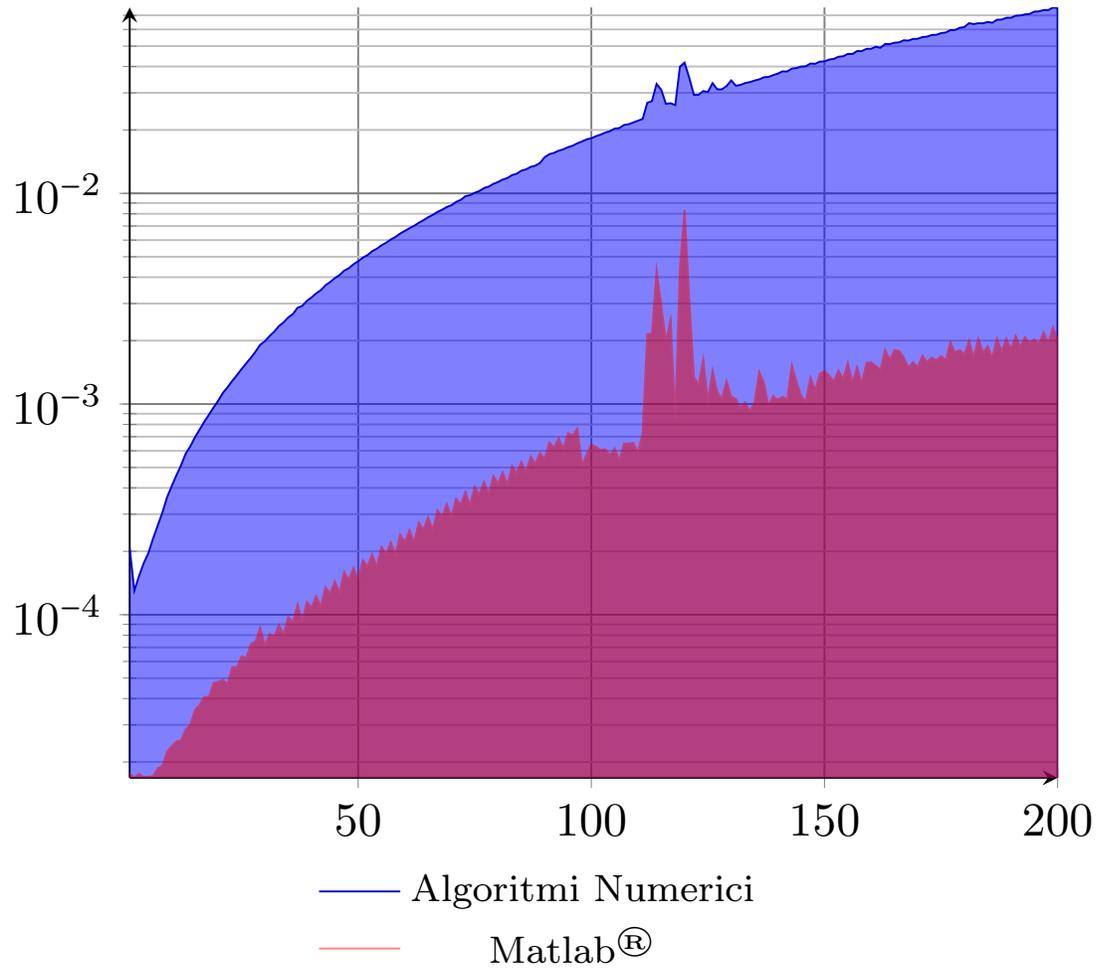


Figura 8.7: Fattorizzazione QR di Gram-Schmidt ($m = n$). Tempi di calcolo.

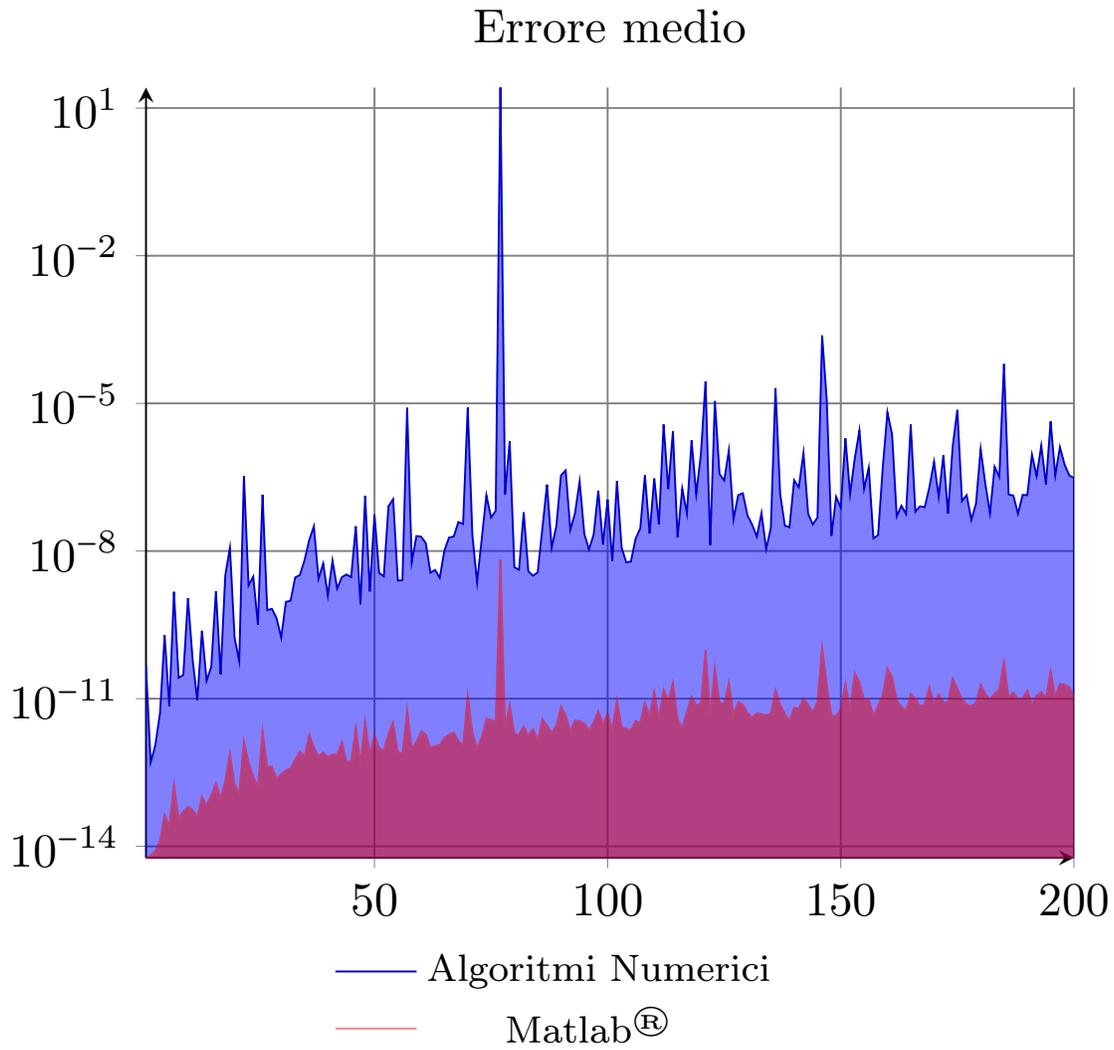


Figura 8.8: Fattorizzazione QR di Gram-Schmidt ($m = n$). Errori di calcolo.

8.3 Problemi minimi quadrati, fattorizzazioni QR

Dopo aver scritto simulato, le singole fattorizzazioni QR, ora le si usa per risolvere un sistema lineare. I grafici mettono a confronto le varie fattorizzazioni, che la legenda riporta.

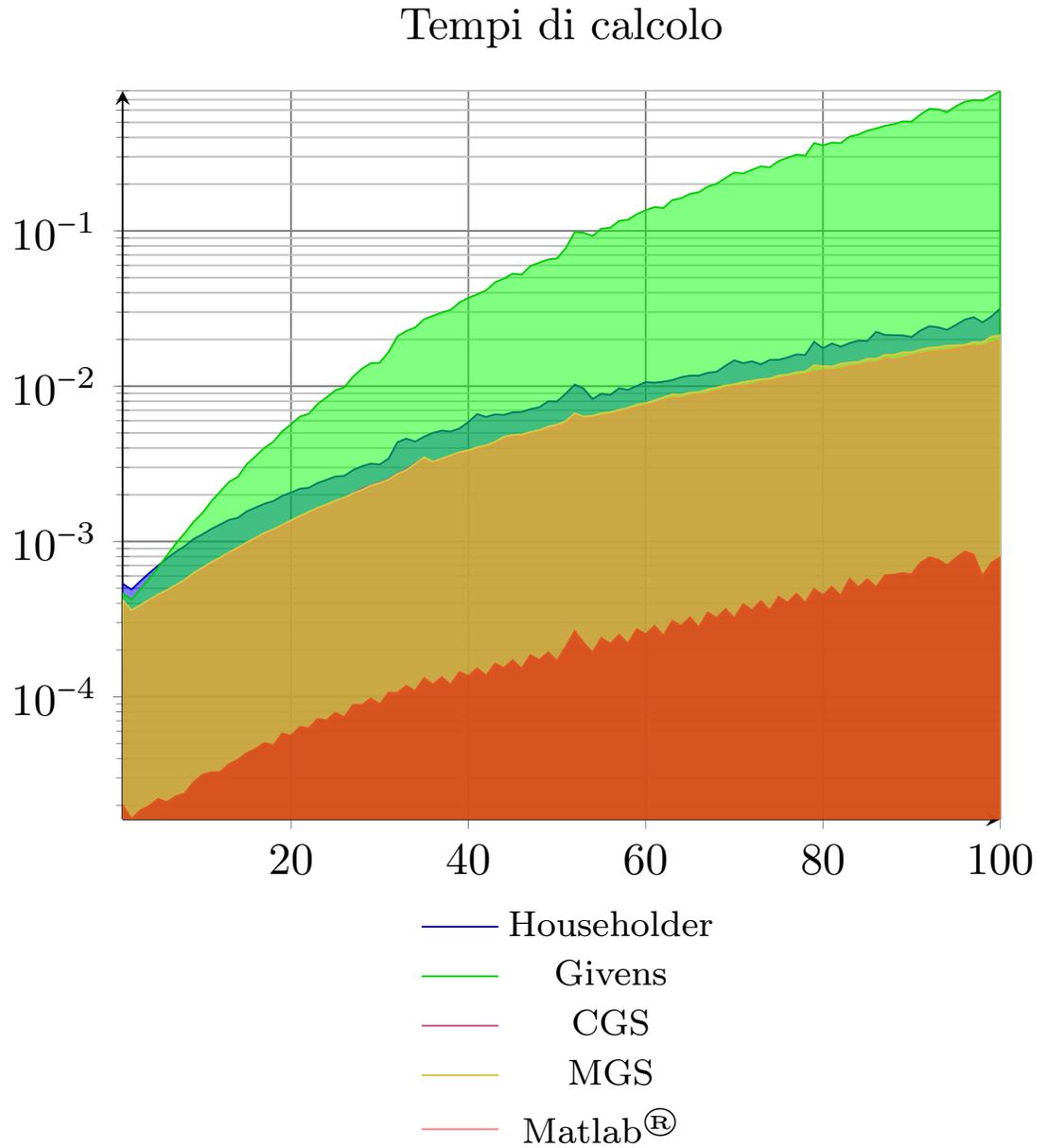


Figura 8.9: Tempo di calcolo della Soluzione Minimi Quadrati $A(m = n)$.

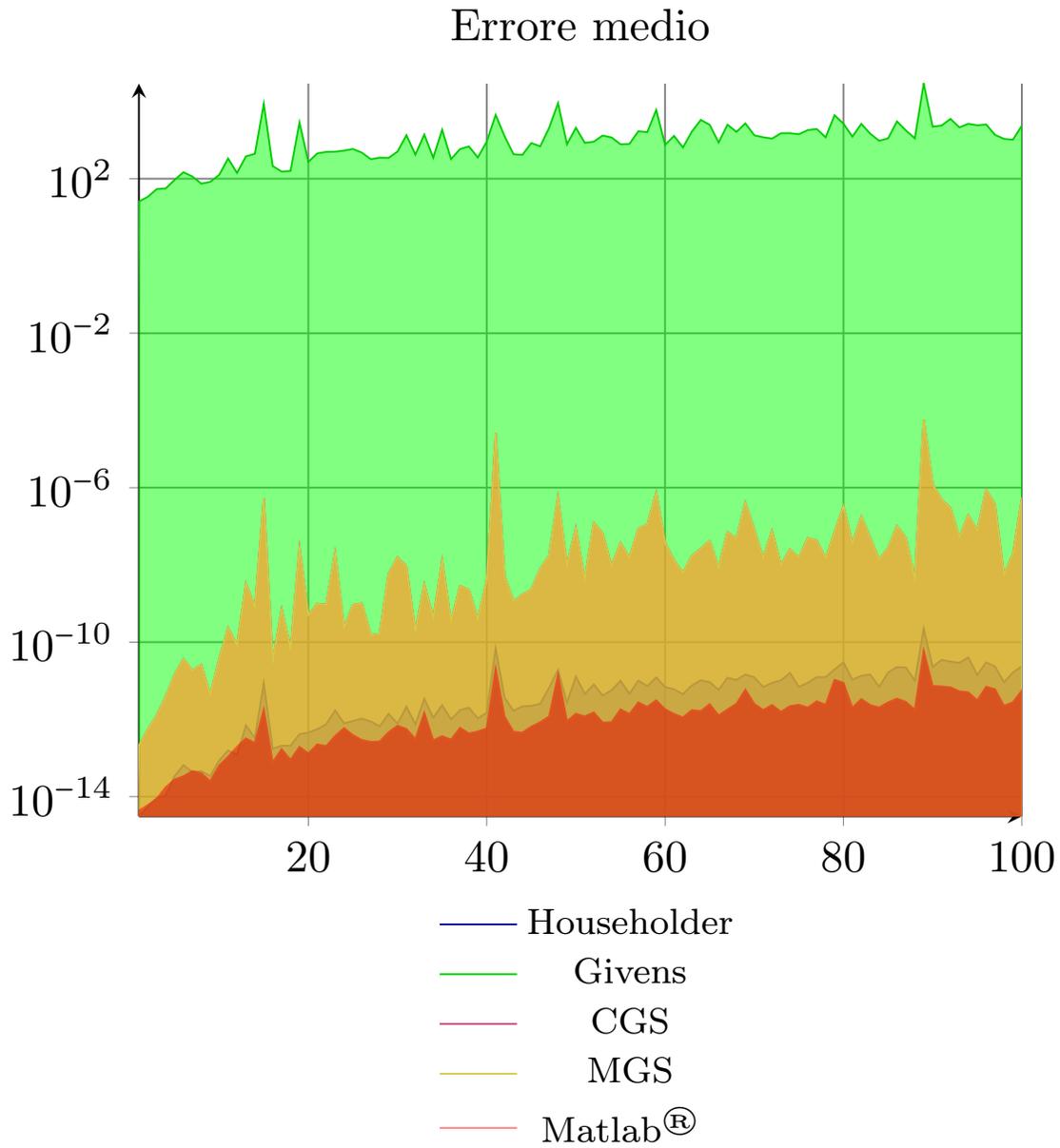


Figura 8.10: Calcolo dell'errore sulla soluzione di un problema ai Minimi Quadrati $A(m = n)$.

Tabelle riassuntive

Tabella 8.1: Tempo medio di calcolo [s]

Dimensione Matrice	Tempo medio [s]				
	Householder	Givens	CGS	MGS	Matlab
10	0.001112	0.001525	0.000672	0.000680	0.000031
20	0.002066	0.005676	0.001337	0.001369	0.000056
40	0.005890	0.037051	0.003837	0.003863	0.000136
60	0.010618	0.135898	0.007465	0.007783	0.000252
80	0.017573	0.355477	0.012593	0.013468	0.000451
100	0.031643	0.797802	0.019703	0.021267	0.000791

8.4 Metodi iterativi metodo GMRES

Si è voluto confrontare `gmres()` di Matlab[®] con gli Algoritmi implementati, la matrice è sempre random piena non simmetrica dunque il metodo di Arnoldi, le iterazioni arrivano a n dimensione della matrice. La matrice va da 3 a 100. Anche qui tempi e errori di calcolo, e infine la tabella riassuntiva dei tempi di calcolo.

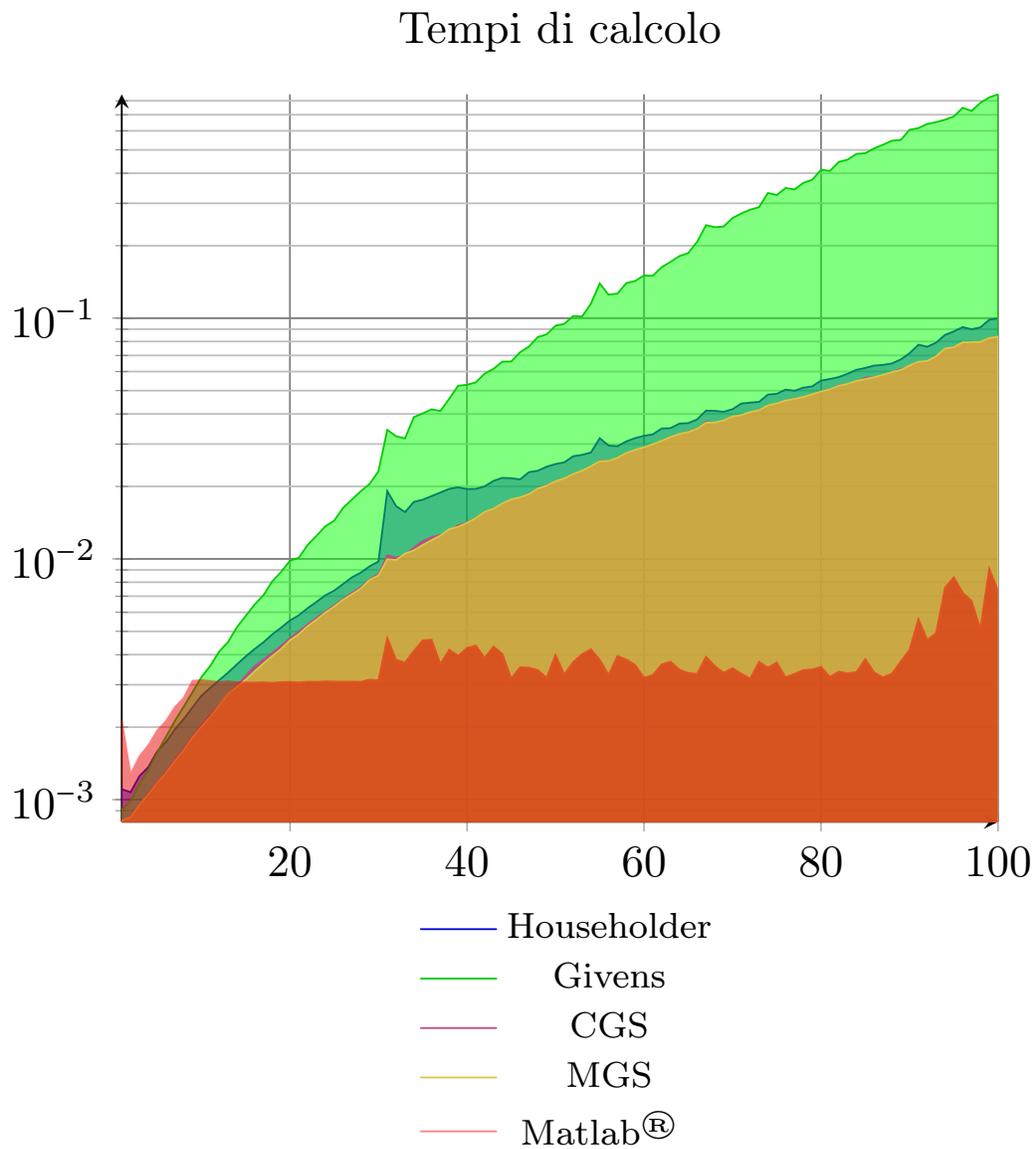


Figura 8.11: Tempo di calcolo della Soluzione Minimi Quadrati $A(m = n)$.

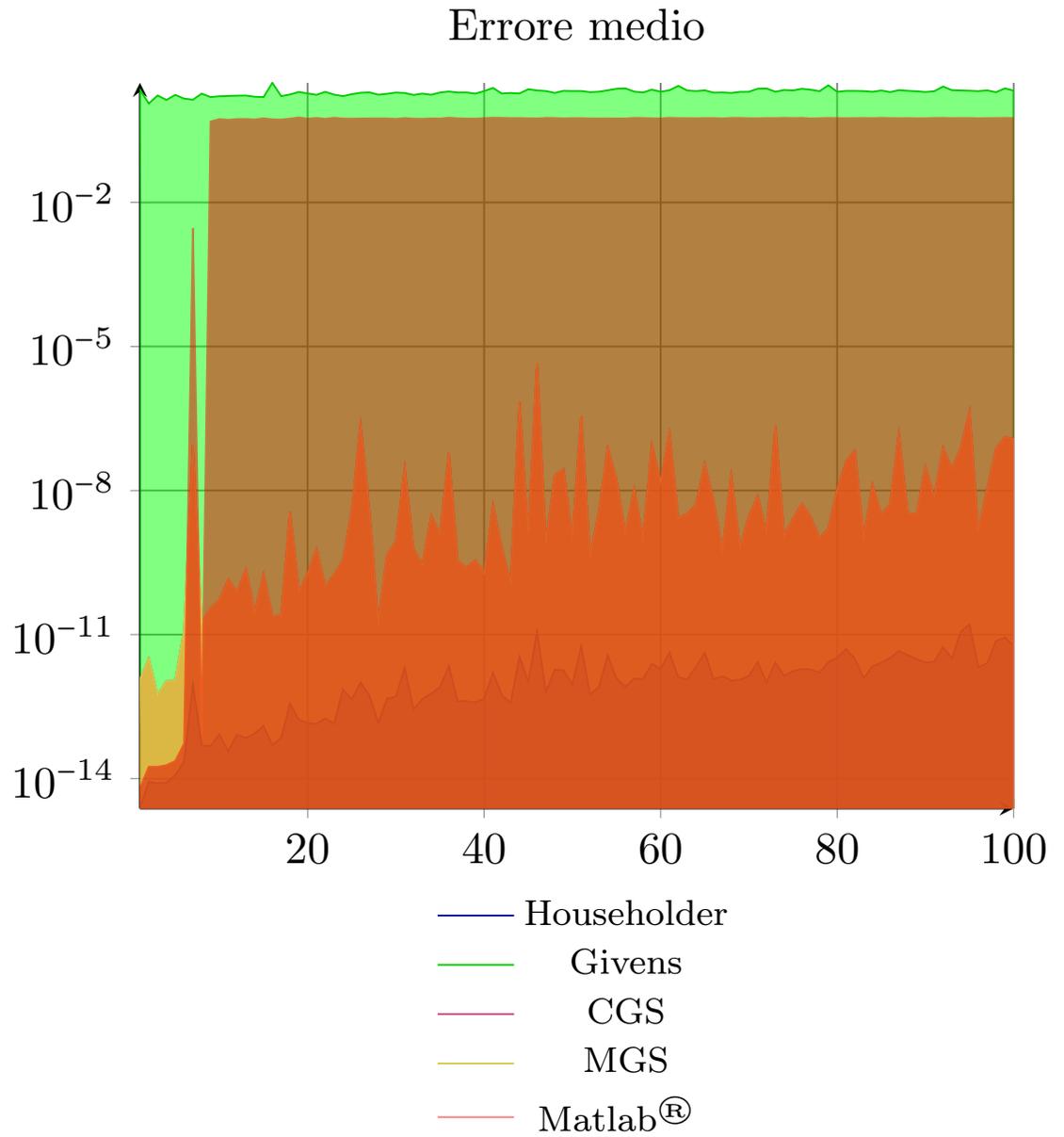


Figura 8.12: Calcolo dell'errore sulla soluzione di un problema ai Minimi Quadrati $A(m = n)$.

Tabelle riassuntive

Tabella 8.2: Tempo medio di calcolo [s]

Dimensione Matrice	Tempo medio [s]				
	Householder	Givens	CGS	MGS	Matlab
10	0.002699	0.003225	0.002036	0.002009	0.003136
20	0.005551	0.009822	0.004680	0.004587	0.003081
40	0.019506	0.052861	0.014151	0.014120	0.004268
60	0.032462	0.150477	0.028941	0.029117	0.003201
80	0.055053	0.413455	0.049643	0.049573	0.003562
100	0.099773	0.850902	0.082725	0.083850	0.007358

Riferimenti bibliografici

[Rod08] Giuseppe Rodriguez. *Algoritmi numerici*. Cagliari: Pitagora Editrice Bologna, 2008.

Siti web consultati

- [Wiki15a] *Oskar Perron*. 2015. URL: https://it.wikipedia.org/wiki/Oskar_Perron.
- [Wiki15b] *Issai Shur*. 2015. URL: https://it.wikipedia.org/wiki/Issai_Schur.
- [Wiki16a] *Walter Edwin Arnoldi*. 2016. URL: https://en.wikipedia.org/wiki/Walter_Edwin_Arnoldi.
- [Wiki16b] *André-Louis Cholesky*. 2016. URL: https://it.wikipedia.org/wiki/Andr f-Louis_Cholesky.
- [Wiki16c] *Jacques Solomon Hadamard*. 2016. URL: https://it.wikipedia.org/wiki/Jacques_Hadamard.
- [Wiki16d] *Charles Hermite*. 2016. URL: https://it.wikipedia.org/wiki/Charles_Hermite.
- [Wiki16e] *John William Strutt*. 2016. URL: https://it.wikipedia.org/wiki/John_William_Strutt_Rayleigh.
- [Wiki17a] *Friedrich Ludwig "Fritz" Bauer*. 2017. URL: https://en.wikipedia.org/wiki/Friedrich_L._Bauer.
- [Wiki17b] *Augustin-Louis Cauchy*. 2017. URL: https://it.wikipedia.org/wiki/Augustin-Louis_Cauchy.
- [Wiki17c] *Ferdinand Georg Frobenius*. 2017. URL: https://it.wikipedia.org/wiki/Ferdinand_Georg_Frobenius.
- [Wiki17d] *Johann Friedrich Carl Gauss*. 2017. URL: https://it.wikipedia.org/wiki/Carl_Friedrich_Gauss.
- [Wiki17e] *J rgen Pedersen Gram*. 2017. URL: https://it.wikipedia.org/wiki/J rgen_Pedersen_Gram.
- [Wiki17f] *David Hilbert*. 2017. URL: https://it.wikipedia.org/wiki/David_Hilbert.
- [Wiki17g] *Nikolay Mitrofanovich Krylov*. 2017. URL: https://en.wikipedia.org/wiki/Nikolay_Mitrofanovich_Krylov.
- [Wiki17h] *Cornelius (Cornel) Lanczos*. 2017. URL: https://en.wikipedia.org/wiki/Cornelius_Lanczos.
- [Wiki17i] *Pierre-Simon Laplace*. 2017. URL: https://it.wikipedia.org/wiki/Pierre_Simon_Laplace.
- [Wiki17j] *Lewis Fry Richardson*. 2017. URL: https://en.wikipedia.org/wiki/Lewis_Fry_Richardson.
- [Wiki17k] *Erhard Schmidt*. 2017. URL: https://it.wikipedia.org/wiki/Erhard_Schmidt.