

## NUMERICAL APPROXIMATION

Ponde chiave: Approssimazione, ottimizzazione, learning e statistica

Approssimazione:  $v \in V$ ,  $a_n \rightarrow v$ ,  $a_n \in \mathcal{Y} \subset V$  ( $\begin{array}{l} \text{es } v = c[a,b] \\ \mathcal{Y} = \mathbb{R}_n[x] \end{array}$ )

Ottimizzazione:  $f: X \rightarrow \mathbb{R}$ ,  $\min_{x \in X} f$  (ordinamento)

Learning:  $f: X \rightarrow \{0,1\}$

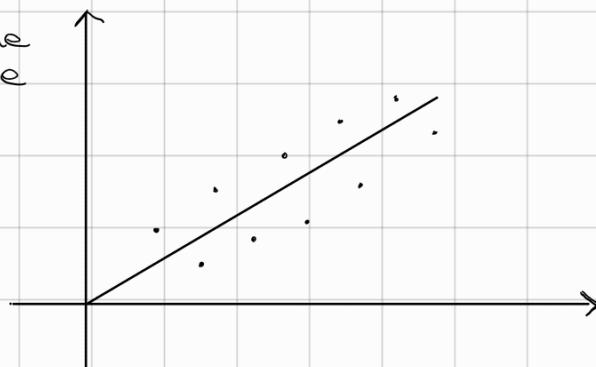
Meglio però ragionare con  $f: X \rightarrow [0,1]$

probabilità

Dal momento che è difficile trovare  $f$ , cerca una  $f_n$  che la approssima

Il problema più semplice è trovare una legge che rappresenti dei dati di un problema reale

Approssimazione: trovare una funzione che fa minimo errore



Ottimizzazione: trovare minimo

Esempio:

(SVD, approccio statistico)  $\rightarrow$  è un problema di approssimazione

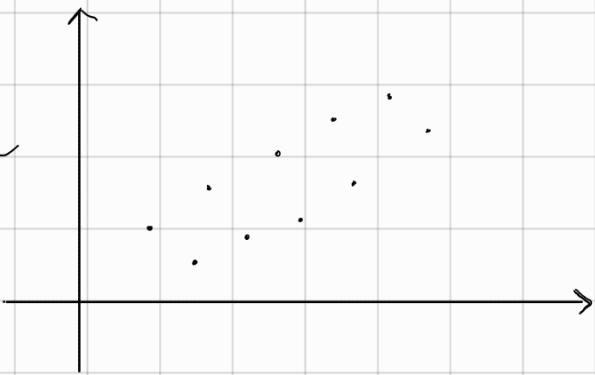
Voglio approssimare la matrice con una di range più basso.

## Minimi quadrati lineari

Ci chiediamo quale sia la migliore retta che genera questi punti.

Costruisca allora, scelta una retta  $y = mx + q$ ,

l'errore che questa retta commette in ogni punto.



Creata il vettore errore  $\epsilon$ , se ne minimizza la norma euclidea

$$\epsilon_i = mx_i + q - y_i \quad i = 1, \dots, n$$

$$\|\epsilon\| = \sqrt{\sum_{i=1}^n |\epsilon_i|^2}$$

Penso minimizzare  $\|\epsilon\|^2$  tante se  $f \geq 0$  i minimi coincidono.

$$f(m, q) = \sum_{i=1}^n |mx_i + q - y_i|^2$$

$$\frac{\partial f}{\partial m}(m, q) = \sum_{i=1}^n 2x_i(mx_i + q - y_i)$$

$$\frac{\partial f}{\partial q}(m, q) = \sum_{i=1}^n 2(mx_i + q - y_i)$$

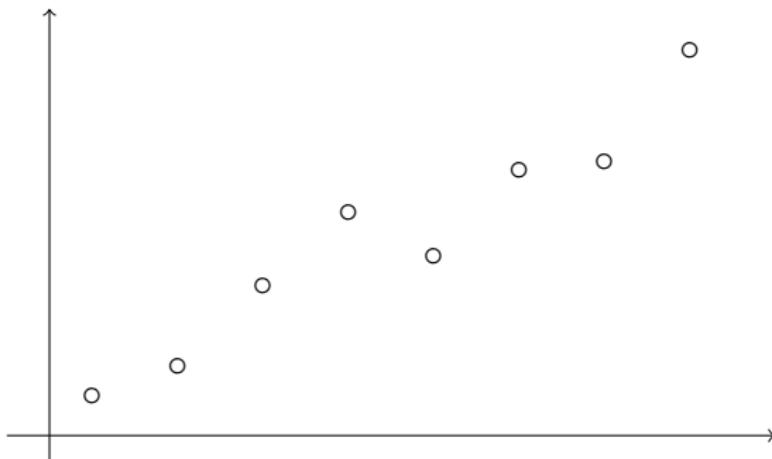
$$\text{Poniamo } \frac{\partial f}{\partial m} = \frac{\partial f}{\partial q} = 0$$

(i punti con  $\frac{\partial f}{\partial m} = 0 = \frac{\partial f}{\partial q}$  si dicono stazionari, non necessariamente minimi)

L'ho svolto sulle rette di Analisi 2.

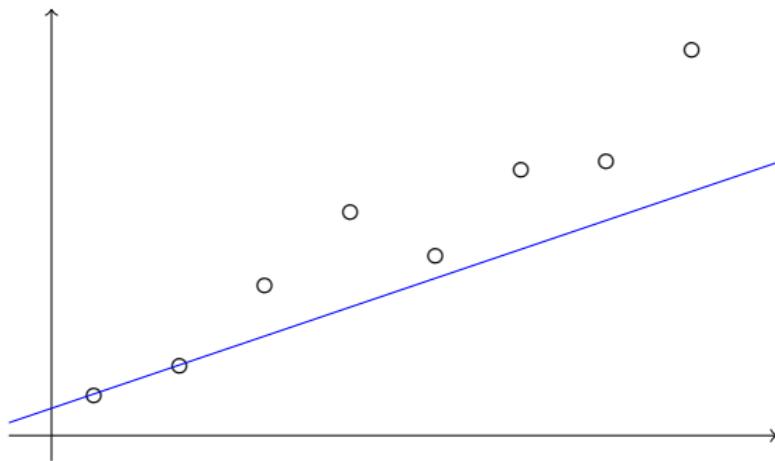
## An approximation problem

“Fit” the following data. Do data hide a linear law with errors?



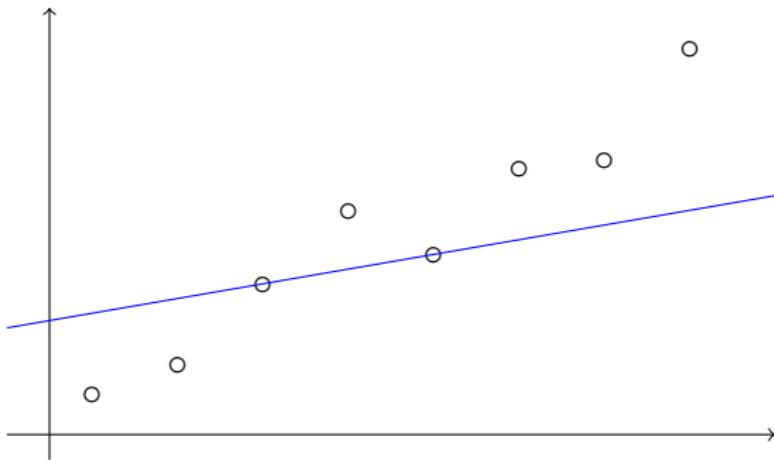
## An approximation problem

“Fit” the following data. Do data hide a linear law with errors?



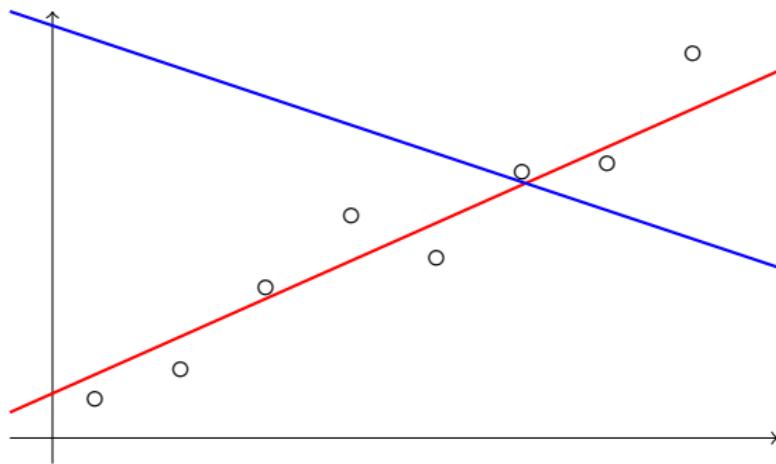
## An approximation problem

“Fit” the following data. Do data hide a linear law with errors?



## An approximation problem

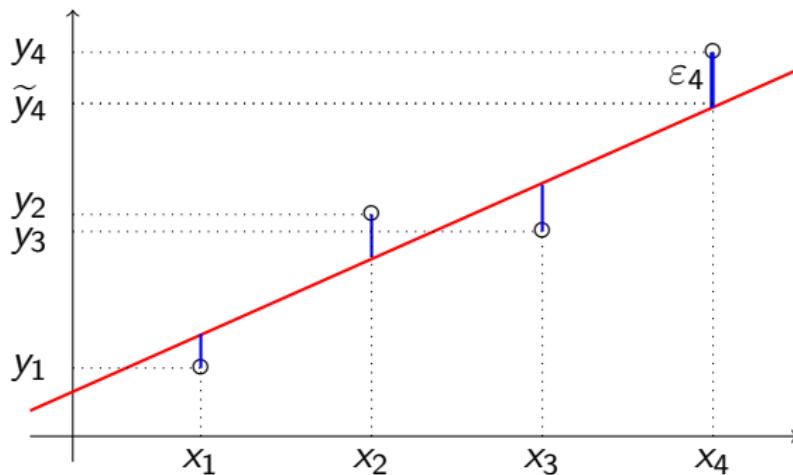
“Fit” the following data. Do data hide a linear law with errors?



## First approach

Search the “nearest” line to the set of points.

We look for the function  $y = mx + q$ . Idea: keep the errors low in approximating  $y_i$  with  $\tilde{y}_i = mx_i + q$ .



Define the (absolute) error on the data  $x_i$  as

$$\varepsilon_i(m, q) = \tilde{y}_i - y_i = mx_i + q - y_i.$$

## First approach

Define the (absolute) error on the data  $x_i$  as

$$\varepsilon_i(m, q) = \tilde{y}_i - y_i = mx_i + q - y_i.$$

We get the vector  $\varepsilon = [\varepsilon_1 \quad \dots \quad \varepsilon_n]^T$ .

We have the following problem

$$\operatorname{argmin}_{m,q \in \mathbb{R}} \|\varepsilon\| = \operatorname{argmin}_{m,q \in \mathbb{R}} \left( \sum_{i=1}^n |mx_i + q - y_i|^2 \right)^{\frac{1}{2}},$$

or, equivalently

$$\operatorname{argmin}_{m,q \in \mathbb{R}} \sum_{i=1}^n (mx_i + q - y_i)^2.$$

## First approach

Let us minimize  $f(m, q) := \sum_i (mx_i + q - y_i)^2$ .

The Hessian matrix

$$H = 2 \begin{bmatrix} \sum_i x_i^2 & \sum_i x_i \\ \sum_i x_i & n \end{bmatrix},$$

is positive semidefinite:  $h_{11} \geq 0$ ,  $h_{22} \geq 0$  and

$\det(H) = n \sum_i x_i^2 - (\sum_i x_i)^2 \geq 0$  follows from the Cauchy-Schwarz inequality

$$\sum_i x_i = \langle e, x \rangle \leq \|e\| \|x\| = \sqrt{n \sum_i x_i^2}.$$

It is positive definite if and only if  $x_i \neq x_j$  for some  $i$  and  $j$ .

$\forall A \in \mathbb{C}^{n \times n}$  es definita positiva se  $A^* = A \in \cup_{v \in \mathbb{C}^n \setminus \{0\}} \{v^* A v > 0\}$

The function  $f$  is convex.

Teorema:

Se  $A \in \mathbb{C}^{n \times n}$  è hermitiana, è definita positiva se e solo se vale uno dei seguenti

- ①  $\det A_l > 0, l = 1, \dots, n$  ( $A_l$  minore principale di testa)
- ② autovalori tutti positivi ( $\neq 0$ )
- ③  $\exists R$  triangolare superiore con elementi  $\mathbb{R}^+$  sulla diagonale tale che  $A = R^* R$  ( $\text{dim } A = \text{dim } R$ )  
(Fattorizzazione di Cycleski, simile a quella LU)
- ④  $\exists B \in \mathbb{C}^{n \times m}$  di range pieno ( $\text{rg } B = \min \{n, m\}$ ) t.c.  
 $A = B^* B$

Tornando a noi, usiamo il metodo dei minori  
 $2 \sum x_i^2 > 0 \Leftrightarrow x_i$  non sono tutti nulli (nel nostro problema  
è una cosa da considerare)

$$\det H = 2^n (n \sum x_i^2 - (\sum x_i)^2)$$

Usiamo la diseguaglianza di Cauchy-Schwarz:

$$\text{se scelgo } v = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, w = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \Rightarrow v^* w = \sum x_i$$

$$\Rightarrow |\sum x_i| \leq \sqrt{n} \sqrt{\sum x_i^2}$$

$$\Rightarrow (\sum x_i)^2 \leq n \sum x_i^2$$

$$\Rightarrow \text{effettivamente quel det è } > 0 \quad \left( \begin{array}{l} \text{se } \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \text{ non è multiplo} \\ \text{di } \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \text{ cioè se e solo se} \\ \text{gli } x_i \text{ non sono tutti coincidenti} \end{array} \right)$$

Quindi se i dati non sono tutti coincidenti allora  
 $H > 0$  (usiamo questo simbolo per def positiva)  $\Rightarrow$  le punti stazionari sono di minimo (e  $H > 0$  implica anche che la funzione è convessa strictamente)

$$\textcircled{2} \quad A \begin{bmatrix} x \\ y \end{bmatrix} = b$$

$$A = 2 \begin{bmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & n \end{bmatrix}, \quad b = \begin{bmatrix} 2 \sum x_i y_i \\ 2 \sum y_i \end{bmatrix}$$

Le soluzioni del sistema sono i punti stazionari dato che  $A = H$ .

Se  $x_i$  non sono tutti coincidenti allora la matrice è invertibile  $\Rightarrow$  c'è un solo punto stazionario e che è minimo poiché  $H > 0$

Se  $x_i$  sono coincidenti  $\sum x_i^2 \geq 0$  e

$$2^n (n \sum x_i^2 - (\sum x_i)^2) = 0$$

Criterio:

$A \in \mathbb{C}^{n \times n}$  hermitiana è semidefinita  $\Leftrightarrow$  vale uno o

$$\textcircled{1} \quad V^* A V \geq 0$$

\textcircled{2} Se  $\det A_k \geq 0$  per ogni minore principale

(non per forza  
di testa)

\textcircled{3} tutti gli autovalori sono  $\geq 0$

$$\textcircled{4} \quad \exists B \in \mathbb{C}^{n \times n} \text{ t.c. } A = B^* B$$

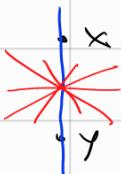
Nel nostro caso controllando anche  $n > 0$  abbiamo che  $H$  è semidefinita positiva  $\Rightarrow$  anche se  $x_i$  non sono tutti coincidenti  $H$  è semi-definita positiva  $\Rightarrow$  tutti i punti stazionari sono di minimo locale.

Che accade se invece della norma euclidea, che risente della grandezza dell'origine, prendessi un'altra norma?

→ stock: calcoli sono più  
difficili  
Si può usare qualsiasi norma, non sappiamo qual è  
la migliore

Potrei cercare le minime di  $\sum \delta^2([x_i], r)$  con  $r$  retta

Sono problemi diversi perché se  $x$  e  $y$  sono sulla stessa ascisse  
il secondo problema ha una una soluzione e l'altro  
problema ne trova infinte



## First approach

The gradient is zero if

$$\frac{\partial f}{\partial m} = \sum_{i=1}^n 2x_i(mx_i + q - y_i) = 0,$$

$$\frac{\partial f}{\partial q} = \sum_{i=1}^n 2(mx_i + q - y_i) = 0.$$

This leads to the linear system

$$M\tilde{x} = c, \quad M = \begin{bmatrix} 2\sum x_i^2 & 2\sum x_i \\ 2\sum x_i & 2n \end{bmatrix}, \quad \tilde{x} = \begin{bmatrix} m \\ q \end{bmatrix}, \quad c = \begin{bmatrix} 2\sum x_i y_i \\ 2\sum y_i \end{bmatrix}.$$

## Second approach

Se procede con la retta di regressione, l'approccio è diverso? no.

The result is obtained from a random vector

$$(X, Y) : \Omega \rightarrow \mathbb{R}^2$$

and we look for a regression between variables of the type  
 $Y = \beta_1 X + \beta_0$ .

We known that, estimators for  $\beta_1$  and  $\beta_0$  are (if  $\text{var}(X) \neq 0$ )

$$\beta_1 = \frac{\text{cov}(X, Y)}{\text{var}(X)}, \quad \beta_0 = \mathbb{E}[Y] - \beta_1 \mathbb{E}[X].$$

### Exercise

*Is this a different solution?* 

## Third approach

Let  $L^2$  be the space of functions on  $\{x_1, \dots, x_n\}$  with the measure that “counts the points” (isomorphic to  $\mathbb{R}^n$ )

$$\int f d\mu := \sum_i f(x_i).$$

We look for the best approximation with respect to  $\mathcal{F} = \{1, x\}$ .

The solution is of the type  $\alpha_0 + \alpha_1 x$ , where  $\alpha = [\alpha_0 \ \alpha_1]^T$  solves the linear system

$$A\alpha = b, \quad A = \begin{bmatrix} \langle \varphi_0, \varphi_0 \rangle & \langle \varphi_1, \varphi_0 \rangle \\ \langle \varphi_0, \varphi_1 \rangle & \langle \varphi_1, \varphi_1 \rangle \end{bmatrix}, \quad b = \begin{bmatrix} \langle \varphi_0, \psi \rangle \\ \langle \varphi_1, \psi \rangle \end{bmatrix},$$

where  $\varphi_0 = 1$ ,  $\varphi_1 = x$  and  $\psi = [y_1, \dots, y_n]$ .

### Exercise

*Is this a different solution?* 

Fitting Kurve



$$c_1 e^x + c_2 e^{2x}$$

Combination zweier Kurven der Form  
exponentielle Kurve  $\curvearrowleft$  - - -

## Linear fitting

Let  $\mathcal{F} = \text{span}\{\varphi_0, \dots, \varphi_\ell\}$ , with  $\varphi_0, \dots, \varphi_\ell$  linearly independent real functions.

Let

$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix}, \quad \begin{bmatrix} x_2 \\ y_2 \end{bmatrix}, \quad \dots, \quad \begin{bmatrix} x_n \\ y_n \end{bmatrix},$$

find  $f \in \mathcal{F}$  such that

$$\sum_{i=1}^n |f(x_i) - y_i|^2$$

is minimum.

(Problema di ottimizzazione)

## Linear fitting

Find  $f \in \mathcal{F}$  such that

$$\sum_{i=1}^n |f(x_i) - y_i|^2$$

is minimum.

Writing  $f = \alpha_0\varphi_0 + \cdots + \alpha_\ell\varphi_\ell$ , the problem can be stated as

Find the minimum of  $\|Ax - b\|$ , with ,  $A \in \mathbb{C}^{n \times \ell}$

$$A = \begin{bmatrix} \varphi_0(x_1) & \cdots & \varphi_\ell(x_1) \\ \vdots & \ddots & \vdots \\ \varphi_0(x_n) & \cdots & \varphi_\ell(x_n) \end{bmatrix}, \quad b = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad x = \begin{bmatrix} \alpha_0 \\ \vdots \\ \alpha_\ell \end{bmatrix}.$$

$$\|Ax - b\|_2^2 = \sum_i \left| \sum_j \varphi_j(x_i) \alpha_j - y_i \right|^2$$

→ il problema si pone come trovare  $\alpha_0, \dots, \alpha_n$  t.c.

$$\sum_i \left| \sum_j \varphi_j(x_i) \alpha_j - y_i \right|^2$$

30/09

Pseudouno dei dati  $\begin{bmatrix} x_1 \\ y_1 \end{bmatrix}, \dots, \begin{bmatrix} x_n \\ y_n \end{bmatrix} \in \mathbb{R}^2$  e

$$Y = \text{span} \{ \varphi_0, \dots, \varphi_n \} \subset C[a, b], \varphi_i \text{ l.i. e } x_i \in [a, b]$$

Trovare la soluzione di  $\underset{\alpha \in Y}{\arg \min} \| \epsilon \|_2$

$$\epsilon = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}, \quad \epsilon_i = f(x_i) - y_i$$

$f = \sum_{j=0}^n \alpha_j \varphi_j$  → il vero problema è trovare

$$\underset{\alpha_0, \dots, \alpha_n \in \mathbb{R}}{\arg \min} \left\| \sum_{j=0}^n \alpha_j \varphi_j(x_i) - y_i \right\|_2^2$$

Se raccolgo i dati noti in una matrice

$$A = \begin{bmatrix} \varphi_0(x_1) & \cdots & \varphi_m(x_1) \\ \vdots & \ddots & \vdots \\ \varphi_0(x_n) & \cdots & \varphi_m(x_n) \end{bmatrix} \quad b = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \quad x = \begin{bmatrix} \alpha_0 \\ \vdots \\ \alpha_n \end{bmatrix}$$

allora dobbiamo trovare  $\underset{x \in \mathbb{R}^{m+1}}{\arg \min} \|Ax - b\|^2$

$$\text{ove } (Ax - b)_i = \sum_{j=0}^m \alpha_j \varphi_j(x_i) - y_i = \epsilon_i$$

(il problema è lo stesso scritto però in forma compatta)

## Overdetermined systems

A classical problem is the solution of a linear system  $Ax = b$  where  $A \in \mathbb{R}^{m \times n}$ , with  $m \geq n$ ,  $b \in \mathbb{R}^m$  and  $x \in \mathbb{R}^n$ .

If  $m > n$  the system is said to be **overdetermined** and has no solution for almost every  $b \in \mathbb{R}^m$ .

When a system has no solution one might be interested in  $x$  such that  $\|Ax - b\|$  is minimus (hoping that  $Ax \approx b$ ).

Same problem as before.

# The linear least squares problem

## Problem

Let  $A \in \mathbb{C}^{m \times n}$  and  $b \in \mathbb{C}^m$ , find  $x \in \mathbb{C}^n$  such that  $\|Ax - b\|$  is minimum.

For any solution  $x^*$ , the quantity  $Ax^* - b$  is said to be **residual** and we are interested in its norm.

## Exercise

*If the linear system  $Ax = b$  has a solution, then the set of solutions of the linear least squares problem coincides with the set of solutions of the linear system.*

In the general case, we have a complete characterization about the existence and uniqueness of solutions.

# The linear least squares problem

## Theorem

Let  $A \in \mathbb{C}^{m \times n}$  with  $m \geq n$  and  $b \in \mathbb{C}^m$ . The linear least squares problem with data  $A$  and  $b$  has solutions and the set of solutions coincides with the solution of the square linear system  $n \times n$

$$A^*Ax = A^*b.$$

There exists a unique solution if and only if  $\text{rank}(A) = n$ .

There exists unique a solution  $x^*$  of minimum norm and it belongs to  $\ker(A^*A)^\perp$ .

- The problem has **always** a solution.
- It can be **reduced** to a square linear system with Hermitian positive semidefinite coefficient.
- **Uniqueness** is obtained by choosing the solution with minimum norm.
- We have both a linear algebra and an analysis proof. 

# The linear least squares problem

## Theorem

Let  $A \in \mathbb{C}^{m \times n}$  with  $m \geq n$  and  $b \in \mathbb{C}^m$ . The linear least squares problem with data  $A$  and  $b$  has solutions and the set of solutions coincides with the solution of the square linear system  $n \times n$

$$A^*Ax = A^*b.$$

## Proof.

Observe that  $\text{range}(A)^\perp = \ker A^*$

$$v \in \text{range}(A)^\perp \iff v^*Aw = 0, w \in \mathbb{C}^n$$

$$\iff \langle A^*v, w \rangle = 0, w \in \mathbb{C}^n \iff A^*v = 0 \iff v \in \ker A^*.$$

We know that  $\mathbb{C}^m = \text{range}(A) \oplus \text{range}(A)^\perp$ , so that there exist unique  $b_1 \in \text{range}(A)$  and  $b_2 \in \ker A^*$  such that  $b = b_1 + b_2$ .



# The linear least squares problem

## Theorem

Let  $A \in \mathbb{C}^{m \times n}$  with  $m \geq n$  and  $b \in \mathbb{C}^m$ . The linear least squares problem with data  $A$  and  $b$  has solutions and the set of solutions coincides with the solution of the square linear system  $n \times n$

$$A^* A x = A^* b.$$

## Proof.

There exist unique  $b_1 \in \text{range}(A)$  and  $b_2 \in \ker A^*$  such that  
 $b = b_1 + b_2$

$$\|Ax - b\|^2 = \langle Ax - b_1 - b_2, Ax - b_1 - b_2 \rangle$$

$$\langle Ax - b_1, Ax - b_1 \rangle + \langle b_2, b_2 \rangle = \|Ax - b_1\|^2 + \|b_2\|^2.$$

Since  $\langle Ax - b_1, b_2 \rangle = 0$ , because  $Ax - b_1 \in \text{range}(A)$  and  $b_2 \in \text{range}(A)^\perp$ .



# The linear least squares problem

## Theorem

Let  $A \in \mathbb{C}^{m \times n}$  with  $m \geq n$  and  $b \in \mathbb{C}^m$ . The linear least squares problem with data  $A$  and  $b$  has solutions and the set of solutions coincides with the solution of the square linear system  $n \times n$

$$A^* A x = A^* b.$$

## Proof.

There exist unique  $b_1 \in \text{range}(A)$  and  $b_2 \in \ker A^*$  such that  
 $b = b_1 + b_2$

$$\|Ax - b\|^2 = \|Ax - b_1\|^2 + \|b_2\|^2.$$

The minimum is attained when  $\|Ax - b_1\|^2$  is minimum, but since  $b_1 \in \text{range}(A)$ , the system  $Ax = b_1$  has a solution, and the solutions of the least squares problem are solution of the system

$$Ax = b_1.$$

# The linear least squares problem

## Theorem

Let  $A \in \mathbb{C}^{m \times n}$  with  $m \geq n$  and  $b \in \mathbb{C}^m$ . The linear least squares problem with data  $A$  and  $b$  has solutions and the set of solutions coincides with the solution of the square linear system  $n \times n$

$$A^*Ax = A^*b.$$

## Proof.

There exist unique  $b_1 \in \text{range}(A)$  and  $b_2 \in \ker A^*$  such that  $b = b_1 + b_2$ . We prove that  $Ax = b_1$  is equivalent to  $A^*Ax = A^*b$ .

$$A^*Ax = A^*b \stackrel{A^*(Ax - b) = 0}{\iff} Ax - b_1 - b_2 \in \ker A^*$$

$$\stackrel{b_2 \in \ker A^*}{\iff} Ax - b_1 \in \ker A^* \stackrel{\ker A^* = \text{range}(A)^\perp}{\iff} Ax - b_1 \in \text{range}(A)^\perp$$

$$\stackrel{Ax - b_1 \in \text{range}(A)}{\iff} Ax - b_1 = 0.$$

Since  $\text{range}(A) \cap \text{range}(A)^\perp = 0$ .

□

# The linear least squares problem

## Theorem

Let  $A \in \mathbb{C}^{m \times n}$  with  $m \geq n$  and  $b \in \mathbb{C}^m$ . The linear least squares problem with data  $A$  and  $b$  has solutions and the set of solutions coincides with the solution of the square linear system  $n \times n$

$$A^*Ax = A^*b.$$

There exists a unique solution if and only if  $\text{rank}(A) = n$ .

There exists unique a solution  $x^*$  of minimum norm and it belongs to  $\ker(A^*A)^\perp$ .

## Proof.

The linear system  $A^*Ax = A^*b$  has unique solution if and only if  $\text{rank}(A^*A) = n$ , but

$$\text{rank}(A^*A) = \text{rank}(A).$$

because <sup>a</sup>

$$A^*Av = 0 \Leftrightarrow v^*A^*Av = 0 \Leftrightarrow \|Av\|^2 = 0 \Leftrightarrow Av = 0.$$

---

<sup>a</sup>If  $M \in \mathbb{C}^{m \times n}$  with  $m \geq n$ ,  $\text{rank}(M) < n$  if and only if  $Mv = 0$ , for  $v \neq 0$ .



# The linear least squares problem

## Theorem

Let  $A \in \mathbb{C}^{m \times n}$  with  $m \geq n$  and  $b \in \mathbb{C}^m$ . The linear least squares problem with data  $A$  and  $b$  has solutions and the set of solutions coincides with the solution of the square linear system  $n \times n$

$$A^* A x = A^* b.$$

There exists a unique solution if and only if  $\text{rank}(A) = n$ .

There exists unique a solution  $x^*$  of minimum norm and it belongs to  $\ker(A^* A)^\perp$ .

## Proof.

The solutions of a linear system are an affine subspace, where there exists a point of minimum norm.

For the last statement see the references.



Esempio (si ottiene un sistema semidefinito)

$Ax = b$        $A \in \mathbb{C}^{m \times n}$        $b \in \mathbb{C}^m$       Se  $m > n$  es  
 sono determinate  
 spesso più sol  
 $\nexists$  (ma è ben posto)

Sono dei problemi motivi da sol. ma e' estremamente associato con la sol perch $\bar{e}$   $ng(A) \neq ng(Al)$

Allora in questo caso si ragiona così: si cerca di mettere le più vicine possibile a una sol.

Intuitiva data dalle informazioni che ho

Cerco la x per cui  $Ax \simeq b$

$\Rightarrow$  Re problema' creato

$$\underset{x \in \mathbb{C}^n}{\operatorname{argmin}} \|Ax - b\|$$

Dej

Se  $x^*$  è una soluzione allora chiamiamola  
 $Ax^* - b$  residuo

## Esercizio:

Se il sistema lineare  $Ax=b$  ha soluzione, allora l'insieme  $W$  delle soluzioni nel problema dei minimi quadrati (LS), con  $A, b$  come sopra, allora  $W$  coincide con l'insieme  $V$  delle sol. del sistema lineare.

Se  $x^* \in V$ ,  $\Lambda x^* - b = 0$  per ogni  $x \in Q^n$

$$\|Ax - b\| \geq 0 = \|Ax^* - b\| \Rightarrow x^* \in W$$

$$\text{Se } \hat{x} \in W \text{ e } x^* \in V \Rightarrow \|Ax^* - b\| \leq \|Ax^* - b\| = 0$$

$$\Rightarrow \|Ax^* - b\| = 0 \Rightarrow Ax^* = b \Rightarrow \hat{x} \in V$$

Teorema del problema dei minimi quadrati:

Sia  $A \in \mathbb{C}^{m \times n}$ ,  $m \geq n$  e  $b \in \mathbb{C}^m$ . Il problema dei minimi quadrati con dati  $A$  e  $b$  ha soluzione  
+ slide

OSS:

- Prob. ha sempre sol (non sempre unica)
- Il prob. si riduce a un sistema lineare

Richiami di analisi:

$f: \mathcal{S} \rightarrow \mathbb{R}$ ,  $\mathcal{S} \subset \mathbb{B}$ ,  $\mathbb{B}$  sp. di Banach (con norma  $\|\cdot\|$ )

$f$  è differenziabile in  $x \in \mathcal{S}$  se  $\exists$  un'applicazione lineare continua  $Df(x): \mathbb{B} \rightarrow \mathbb{R}$  t.c.

$$\lim_{h \rightarrow 0} \frac{|f(x+h) - f(x) - Df(x)[h]|}{\|h\|} = 0$$

$Df \Rightarrow$  chiamata derivata di Frechet di  $f$  in  $x$   
o anche

$$f(x+h) = f(x) + Df(x)[h] + o(\|h\|)$$

$$f(x+h) = f(x) + Df(x)[h] + \frac{1}{2} D^2 f(x)[h, h] + o(\|h\|^2)$$

$D^2 f(x): \mathbb{B} \times \mathbb{B} \rightarrow \mathbb{R}$

derivata  
seconda

lineare

$$x \in \mathbb{R}^n$$

$$x \in \mathbb{R}$$

$$\nabla f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$A \in \mathbb{R}^{1 \times n}$$

$$A = \left[ \frac{\partial}{\partial x_1}(x), \dots, \frac{\partial}{\partial x_n}(x) \right]$$

$$\nabla f(x)[h] = \langle \nabla f(x), h \rangle, h \in \mathbb{R}^n, \nabla f(x) \in \mathbb{R}^n \cong \mathbb{R}^{n \times 1}$$

$$\nabla^2 f(x)[h, h] = h^T H h$$

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \dots \\ \dots & \dots \end{bmatrix}$$

Dimostrazione <sup>1</sup> del Teorema dei minimi quadrati:

$$\langle Av, w \rangle = \langle v, A^* w \rangle$$

Vogliamo trovare i minimi di  $\|Ax - b\|^2 = \varphi(x)$

$$\varphi(x+h) - \varphi(x) = \langle A(x+h) - b, A(x+h) - b \rangle - \langle Ax - b, Ax - b \rangle =$$

$$= \cancel{\langle Ax - b, Ax - b \rangle} + \langle Ah, Ah \rangle + 2 \langle Ax - b, Ah \rangle - \cancel{\langle Ax - b, Ax - b \rangle}$$

$$= \underbrace{\langle 2A^*(Ax - b), h \rangle}_{\text{lineare}} + \underbrace{h^T A^* A h}_{\text{bilineare (O(C||h||)) posse}} \quad$$

$$\|h^T A^* A h\| \leq \|h^T\| \|A^* A\| \|h\| \leq \|h\| \leq O(C\|h\|)$$

$$\nabla \varphi(x) = 2(A^* Ax - A^* b) = 0 \Leftrightarrow A^* Ax = A^* b$$

$$\nabla^2 \varphi(x) = A^* A \geq 0 \quad (= \text{tutte le sottosassette sono > zero in loc})$$

## Dimostrazione ②:

$$\text{Im } A^+ = \ker A^* \quad (\text{dimostrazione questo})$$

$$r \stackrel{=Aw}{\in} \text{Im } A^+ \Leftrightarrow \langle v, Aw \rangle = 0 \quad w \in \mathbb{C}^n$$

$$0 = \langle A^* v, w \rangle = 0 \quad w \in \mathbb{C}^n \Leftrightarrow A^* v = 0 \Leftrightarrow v \in \ker A^*$$

$$\text{ma } \mathbb{C}^m = \text{Im } A \oplus \text{Im } A^+$$

$$\Rightarrow b = b_1 + b_2 \text{ con } b_1 \in \text{Im } A, b_2 \in \text{Im } A^+ = \ker A^*$$

$$\begin{aligned} &= \|Ax - b\|^2 - \|Ax - b_1 - b_2\|^2 = \|Ax - b_1\|^2 + \|b_2\|^2 - 2\operatorname{Re} \langle Ax - b_1, b_2 \rangle \\ &\quad \xrightarrow{\text{costante}} \in \text{Im } A \quad \in \text{Im } A^+ \\ &= \|Ax - b_1\|^2 + \|b_2\|^2 \quad (\operatorname{argmin}_f = \operatorname{argmin}_{f+c}) \end{aligned}$$

$$\Rightarrow \operatorname{argmin}_x \|Ax - b\|^2 = \operatorname{argmin}_x \|Ax - b_1\|^2$$

ma  $Ax = b_1$  ha soluzione perché  $b_1 \in \text{Im } A$

$$\Rightarrow \operatorname{argmin}_x \|Ax - b\|^2 = \left\{ x \mid Ax = b_1 \right\}$$

Dimostriamo che  $Ax = b_1 \Leftrightarrow A^* A x = A^* b$

$$A^* A x = A^* b \Leftrightarrow A^*(Ax - b) = 0 \Leftrightarrow Ax - b \in \ker A^*$$

$$\Leftrightarrow Ax - b_1 - b_2 \in \text{Im } A^+ \Leftrightarrow Ax - b_1 \in \text{Im } A^+ \quad \in \text{Im } A$$

$$\text{ma } \text{Im } A^+ \cap \text{Im } A = \{0\} \Leftrightarrow Ax - b_1 = 0 \Leftrightarrow Ax = b_1$$

Dimostriamo la parte dell' inverso  
e si può trovare la soluz. unica se e solo se  
 $A^*A$  è invertibile

$A^*A$  è semidefinita positiva.

Se  $\text{rg } A < n \Rightarrow Ax=0$  ha sol. non nulli  $\Leftrightarrow \|Ax\|^2 = 0$

$$\Leftrightarrow X^*A^*Ax = 0$$

Se  $A^*Av = 0 \Rightarrow v^*A^*Av = 0 \Rightarrow \|Av\|^2 = 0 \Rightarrow Av = 0$   
 $A^*A$  non invertibile

$\Rightarrow$  Se  $\text{rg } A = n$  allora  $A^*A$  è invertibile

Vediamo ora algoritmi per risolvere il sistema:

1 - Fattorizzazione di Cholesky ( $\text{rg } A = n$ )

2 - Fattorizzazione QR

3 - SVD

# The Choleksy approach

## Theorem

Let  $A \in \mathbb{C}^{n \times n}$  be positive definite, there exists unique  $R \in \mathbb{C}^{n \times n}$  upper triangular such that  $A = R^*R$  and the diagonal elements of  $R$  are positive.

## Proof.

The principal minor of a positive definite matrix are positive definite  $\Rightarrow A$  admits a unique LU factorization.

Write  $U = D\tilde{U}$ , where  $D$  is diagonal and  $\tilde{U}$  is upper triangular with 1 on the diagonal

$$\begin{bmatrix} u_{11} & u_{12} & \cdots & u_{1n} \\ 0 & u_{22} & \cdots & u_{2n} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \ddots & 0 & u_{nn} \end{bmatrix} = \begin{bmatrix} u_{11} & 0 & \cdots & 0 \\ 0 & u_{22} & \cdots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \ddots & 0 & u_{nn} \end{bmatrix} \begin{bmatrix} 1 & u_{12}/u_{11} & \cdots & u_{1n}/u_{11} \\ 0 & 1 & \cdots & u_{2n}/u_{22} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \ddots & 0 & 1 \end{bmatrix}$$



# The Choleksy approach

## Theorem

Let  $A \in \mathbb{C}^{n \times n}$  be positive definite, there exists unique  $R \in \mathbb{C}^{n \times n}$  upper triangular such that  $A = R^*R$  and the diagonal elements of  $R$  are positive.

## Proof.

Since

$$A = LU = L(D\tilde{U}), \quad A = A^* = \tilde{U}^*D^*L^*$$

$L(D\tilde{U})$  and  $\tilde{U}^*(D^*L^*)$  are LU factorizations of  $A$ , by uniqueness of LU factorization  $L = \tilde{U}^*$  and we can write  $A = LDL^*$ .



# The Choleksy approach

## Theorem

Let  $A \in \mathbb{C}^{n \times n}$  be positive definite, there exists unique  $R \in \mathbb{C}^{n \times n}$  upper triangular such that  $A = R^*R$  and the diagonal elements of  $R$  are positive.

## Proof.

We have  $A = LDL^*$ , we prove that  $D$  has positive diagonal elements. If  $v \neq 0$ , then

$$v^*Dv = v^*L^{-1}AL^{-*}v = w^*Aw > 0,$$

with  $w = L^{-*}v \neq 0$  since  $v \neq 0$ . This implies that  $D$  is definite positive and thus has a positive diagonal.



# The Choleksy approach

## Theorem

Let  $A \in \mathbb{C}^{n \times n}$  be positive definite, there exists unique  $R \in \mathbb{C}^{n \times n}$  upper triangular such that  $A = R^*R$  and the diagonal elements of  $R$  are positive.

## Proof.

$A = LDL^*$ ,  $D$  has positive diagonal, we can write

$$D = \begin{bmatrix} u_{11} & 0 & \cdots & 0 \\ 0 & u_{22} & \cdots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \ddots & 0 & u_{nn} \end{bmatrix} = \begin{bmatrix} \sqrt{u_{11}} & 0 & \cdots & 0 \\ 0 & \sqrt{u_{22}} & \cdots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \ddots & 0 & \sqrt{u_{nn}} \end{bmatrix} \begin{bmatrix} \sqrt{u_{11}} & 0 & \cdots & 0 \\ 0 & \sqrt{u_{22}} & \cdots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \ddots & 0 & \sqrt{u_{nn}} \end{bmatrix} = FF^*$$

and thus  $A = (LF)(F^*L)$ . Setting  $R = F^*L$  we have the proof since  $A = R^*R$ ,  $R$  is upper triangular and  $R$  has positive diagonal elements.



## The Cholesky approach

Features of the Cholesky factorization.

- Characterizes positive definite matrices;
- Parametrize (complex) positive definite matrices as  $\mathbb{R}^{n(n-1)} \times \mathbb{R}_+^n$ ;
- Can be computed by an algorithm requiring  $n^3/3 + o(n^3)$  arithmetic operations.

If  $A$  has rank  $n$ , the Cholesky factorization of  $A^*A$  allows one to solve the least squares problem

$$A^*Ab = A^*b,$$

Drawback:  $A^*A$  must be computed and can have a large condition number (larger than the one of  $A$ ).

# 1 - Fattorizzazione di Cholesky:

Teorema

Sei  $A \in \mathbb{C}^{n \times n}$  def positiva,  $\exists ! R \in \mathbb{C}^{n \times n}$  triangolare superiore con elementi positivi sulla diagonale t.c  $A = R^*R$  (fattorizzazione di Cholesky)

Dimo:

per ind. su  $n$ :

$n=1$ :  $A = a > 0$   $A = \sqrt{a}\sqrt{a}$  in modo unico

$$R = R^* = \sqrt{a}$$

$n-1=n$ :  $A \in \mathbb{C}^{n \times n}$

$$A = \begin{bmatrix} A_{n-1} & u \\ u^* & \alpha \end{bmatrix}$$

Mostriamo che  $A_{n-1}$  è def positiva poiché  $A$  è def positiva

e anche  $\lambda \in \mathbb{R} \setminus \{0\}$  risolve  $A + \lambda I$  definita positiva (minore > 0)

Inoltre  $A$  invertibile e anche  $A_{n-1}$

$$\begin{bmatrix} I & P \\ -u^* A_{n-1}^{-1} & 1 \end{bmatrix} \begin{bmatrix} A_{n-1} & u \\ u^* & \lambda \end{bmatrix} =$$

$$= \begin{bmatrix} A_{n-1} & u \\ -u^* \tilde{A}_{n-1}^{-1} + u^* & \lambda - u^* \tilde{A}_{n-1}^{-1} u \end{bmatrix}$$

$$\det = \det A_{n-1} \cdot \gamma > 0 \Rightarrow \gamma > 0$$

$$A = R^* R$$

$$R = \begin{bmatrix} R_{n-1} & Y \\ 0 & X \end{bmatrix}$$

$$= \begin{bmatrix} A_{n-1} & u \\ u^* & \lambda \end{bmatrix} = \begin{bmatrix} R_{n-1}^* & 0 \\ Y^* & \bar{X} \end{bmatrix} \begin{bmatrix} R_{n-1} & Y \\ 0 & X \end{bmatrix} =$$

$$= \begin{bmatrix} R_{n-1}^* R_{n-1} & R_{n-1}^* Y \\ Y^* R_{n-1} & Y^* Y + \bar{X} X \end{bmatrix}$$

$$\left\{ \begin{array}{l} A_{n-1} = R_{n-1}^* R_{n-1} \text{ per } N_p \text{ invertibile} \text{ la sol unica} \\ \lambda = R_{n-1}^* Y \end{array} \right.$$

$$\left\{ \begin{array}{l} \lambda = R_{n-1}^* Y \Leftrightarrow R_{n-1} \text{ ho gli positivi reale} \Rightarrow \det R_{n-1} > 0 \\ \Rightarrow \det R_{n-1}^* > 0 \Rightarrow \text{esiste unica la sol unica } Y \end{array} \right.$$

$$\lambda = Y^* Y + \bar{X} X, X > 0 \Leftrightarrow |X|^2 = \lambda - Y Y^*$$

$$\text{ma } Y = R_{n-1}^{-1} u \Rightarrow Y^* = u^* (R_{n-1}^{-1})^* = u^* R_{n-1}^{-1}$$

$$\Rightarrow \|x\|^2 = \alpha - u^* R_{n-1}^{-1} R_{n-1}^{*} u = \alpha - u^* A_{n-1}^{-1} u = \gamma$$

$\Rightarrow x \in$  sol per  $\gamma$  se e solo se  $x = \sqrt{\gamma}$

OSS:

$$\Rightarrow A_{n-1} = R_{n-1}^* R_{n-1} \quad A_{n-1}^{-1} = (R_{n-1}^* R_{n-1})^{-1} = R_{n-1}^{-1} (R_{n-1}^*)^{-1}$$

$$\Rightarrow u = R_{n-1}^* y \Rightarrow y = (R_{n-1}^*)^{-1} u$$

$$\Rightarrow u^* = (R_{n-1}^* y)^* = y^* R_{n-1} \Rightarrow y^* = u^* R_{n-1}^{-1}$$

Algoritmo che calcola fattorizzazione di Cholesky

$$A = R^* R \quad (\text{vedere su } \mathbb{C}^{3 \times 3})$$

$$\begin{bmatrix} R_{11} & 0 & 0 \\ \bar{R}_{12} & R_{22} & 0 \\ \bar{R}_{13} & \bar{R}_{23} & R_{33} \end{bmatrix} \begin{bmatrix} R_{11} & R_{12} & R_{13} \\ 0 & R_{22} & R_{23} \\ 0 & 0 & R_{33} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ \bar{a}_{12} & a_{22} & a_{23} \\ \bar{a}_{13} & \bar{a}_{23} & a_{33} \end{bmatrix}$$

$$\left\{ \begin{array}{l} R_{11}^2 = a_{11} \\ R_{11} R_{12} = a_{12} \\ R_{11} R_{13} = a_{13} \\ \bar{R}_{12} R_{12} + R_{22}^2 = a_{22} \\ \bar{R}_{12} R_{13} + R_{22} R_{23} = a_{23} \\ \bar{R}_{13} R_{13} + \bar{R}_{23} R_{23} + R_{33}^2 = a_{33} \end{array} \right.$$

$$\left\{ \begin{array}{l} r_{ii} = \sqrt{\alpha_{ii}} \\ r_{i2} = \alpha_{i2} / r_{ii} \\ r_{i3} = \alpha_{i3} / r_{ii} \\ r_{22} = \sqrt{\alpha_{22} - r_{12}^2} \\ r_{23} = - \\ r_{33} = - \end{array} \right.$$

Vediamo le similitudini

$$(R^* R)_{ij} = \alpha_{ij}, \quad i \leq j$$

$$\begin{aligned} (R^* R)_{ij} &= \sum_{l=1}^n (R^*)_{il} (R)_{lj} = \sum_{l=1}^n \overline{r}_{ei} \underbrace{r_{lj}}_{0 \text{ se } l > j} = \\ &\stackrel{\min(i,j)}{=} \sum_{l=1}^i \overline{r}_{ei} r_{lj} = \sum_{l=1}^i \overline{r}_{ei} r_{lj} = \alpha_{ij} \quad i=1, \dots, n \quad i \leq j \\ &\text{by?} \end{aligned}$$

$$\sum_{l \neq i} \overline{r}_{ei} r_{li} + \overline{r}_{ii} r_{ii} = \alpha_{ii} \quad i = j$$

$$\Rightarrow r_{ii} = \sqrt{\alpha_{ii} - \sum_{l \neq i} |r_{ei}|^2}$$

$$\text{Per } i \leq j \quad \sum_{l \neq i} \overline{r}_{ei} r_{ej} + r_{ii} r_{ij} = \alpha_{ij}$$

$$\Rightarrow r_{ij} = \frac{\alpha_{ij} - \sum_{l \neq i} \overline{r}_{ei} r_{el}}{r_{ii}}$$

(dove calcolare  
in ordine  $\alpha_{ij}$  per regole  
di sx o dx)

L'implementazione lo faccio

Per  $i = 1, \dots, n \Rightarrow r_{ii} = \dots$

per  $j = i+1, \dots, n \Rightarrow r_{ij} = \dots$

Ha senso se notiamo che  $r_{nn} \neq 0$   
e se  $r_{ni} \neq 0$

Costo complessivo

Per  $r_{ii}$ :  $2(i-1)$  ops + 1 sqrt  $\rightarrow \gamma = O(1)$  n di ops da sqrt  
Per  $r_{ij}$ :  $2(i-1) + 1$  ops  
Divisione

$$\Rightarrow \sum_{i=1}^n 2(i-1) + \gamma + \sum_{i=1}^n \sum_{j=i+1}^n 2(i-1) + 1$$

$$\text{per } n \rightarrow \infty \quad O(n^2) + \sum_{i=1}^n (n-i)(2(i-1)+1) \approx$$

$$\approx \sum_{i=1}^n 2i(n-i) = 2\left(n \sum i - \sum i^2\right) \approx 2\left(\frac{n^2}{2} - \frac{n^3}{3}\right) \approx \frac{n^3}{3} \text{ ops}$$

(La metà del tempo che ci vuole Gauss-Jordan  
usa metà parametri)

Dunque:

01 / 10

- calcolare  $A^*A$  e  $A^*b \rightarrow 2mn^2$
- Factorizzare  $A^*A$  come  $R^*R \rightarrow \frac{n^3}{3}$
- $R^*y = A^*b \rightarrow n^2$
- $Rx = y \rightarrow n^2$

Esempio:  $A = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$   $b = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$   $Ax = b$

$$\begin{cases} 2x = 1 \\ x = 1 \end{cases} \quad \text{Non ha sol. chiuso}$$

Nella vedremo le metodi per min  $\|Ax-b\|^2$

## The QR factorization approach

If  $A = QR = \begin{bmatrix} Q_1 & Q_2 \end{bmatrix} \begin{bmatrix} R_1 \\ 0 \end{bmatrix}$ , with  $R_1 \in \mathbb{C}^{n \times n}$

$$\|Ax - b\|^2 = \|QRx - b\|^2 = \|Q(Rx - Q^*b)\|^2 \stackrel{\|Qv\|=\|v\|}{=} \|Rx - Q^*b\|^2$$

$$\begin{aligned} &= \left\| \begin{bmatrix} R_1 \\ 0 \end{bmatrix} x - \begin{bmatrix} Q_1^* \\ Q_2^* \end{bmatrix} b \right\|^2 = \left\| \begin{bmatrix} R_1 x - Q_1^* b \\ Q_2^* b \end{bmatrix} \right\|^2 \\ &\stackrel{\left\| \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \right\|^2 = \|v_1\|^2 + \|v_2\|^2}{=} \|R_1 x - Q_1^* b\|^2 + \|Q_2^* b\|^2. \end{aligned}$$

The minimum of  $\|Ax - b\|^2$  is reached for  $x$  such that  $R_1 x = Q_1^* b$  and the norm of the residual is  $\|Q_2^* b\|$ .

Can deal also with the case  $\text{rank}(A) < n$  (QR factorization with column pivoting).

## QR factorization

A **QR factorization** of a matrix  $A \in \mathbb{C}^{n \times n}$  consists in writing the matrix  $A$  as the product  $QR$  where  $Q \in U(n)$  is unitary and  $R \in \mathbb{C}^{n \times n}$  upper triangular.

When  $A \in \mathbb{R}^{n \times n}$ , we may require that  $Q$  is orthogonal and  $R$  is real.

Can be extended to rectangular matrices, where  $R$  is upper triangular with the same size as  $A$ .

# Existence of the QR factorization

## Theorem

*Any square matrix admits a QR factorization.*

## Proof.

*Idea: use the fact that any unitary vector  $v \in \mathbb{C}^n$ , can be completed to a basis of  $\mathbb{C}^n$ .*



# Existence of the QR factorization

## Theorem

*Any square matrix admits a QR factorization.*

## Proof.

We proceed by induction on the size  $n$  of the matrix  $A$ . The case  $n = 1$  is given by  $A = QR$  with  $Q = 1$  and  $R = A$ . □

# Existence of the QR factorization

## Theorem

*Any square matrix admits a QR factorization.*

## Proof.

Assume that  $n > 1$  and let  $v \in \mathbb{C}^n$  be the first column of  $A$ . If  $v = 0$  set  $q_1 = [1 \ 0 \ \cdots \ 0]^T$ , elsewhere set  $q_1 = v/\|v\|$ .

There exist  $q_2, \dots, q_n$  spanning an  $n - 1$  dimensional subspace of  $\mathbb{C}^n$  orthogonal to  $\text{span}(q_1)$ .

Let  $U = [q_1 | Q_2]$ , where  $Q_2 = [q_2 | \cdots | q_n]$  and partition  $A = [v | B]$  with  $B \in \mathbb{C}^{n \times (n-1)}$ . Notice that  $U$  is unitary and

$$U^* A = \begin{bmatrix} q_1^* \\ Q_2^* \end{bmatrix} \begin{bmatrix} v & B \end{bmatrix} = \begin{bmatrix} \|v\| & q_1^* B \\ 0 & Q_2^* B \end{bmatrix}.$$



# Existence of the QR factorization

## Theorem

*Any square matrix admits a QR factorization.*

## Proof.

By induction we can write  $Q_2^*B = \tilde{Q}\tilde{R}$ , with  $\tilde{Q} \in U(n-1)$  and  $\tilde{R}$  upper triangular. Setting  $\hat{Q} = \begin{bmatrix} 1 & 0 \\ 0 & \tilde{Q}^* \end{bmatrix} \in U(n)$  we get that

$$R = \hat{Q}(U^*A) = \begin{bmatrix} 1 & 0 \\ 0 & \tilde{Q}^* \end{bmatrix} \begin{bmatrix} \|v\| & q_1^*B \\ 0 & Q_2^*B \end{bmatrix} = \begin{bmatrix} \|v\| & q_1^*B \\ 0 & \tilde{R} \end{bmatrix}$$

is upper triangular, thus  $A = QR$  with  $Q = U\hat{Q}^*$  is a QR factorization of  $A$ . □

## Uniqueness of the QR factorization

The QR factorization is **not unique**, think to the one-dimensional case  $a = e^{i\theta}(e^{-i\theta}a)$ .

It is easy to characterize all QR factorizations of a nonsingular matrix in terms of phase matrices.

A **phase matrix**  $F \in \mathbb{C}^{n \times n}$  is a diagonal unitary matrix

$$F = \begin{bmatrix} e^{i\theta_1} & & & \\ & e^{i\theta_2} & & \\ & & \ddots & \\ & & & e^{i\theta_n} \end{bmatrix}.$$

### Exercise

*Prove that a unitary and upper triangular matrix is a phase matrix.*

# Uniqueness of the QR factorization

## Theorem

Let  $A \in GL(n)$  and let  $QR$  be a factorization of  $A$ .

For any  $QR$  factorization of  $A$ , say  $A = \tilde{Q}\tilde{R}$  there exists a phase matrix  $F$  such that  $\tilde{Q} = QF^*$  and  $\tilde{R} = FR$ . Conversely, any phase matrix  $F$  yields a  $QR$  factorization  $A = (QF^*)(FR)$ .

In particular, there exists unique a  $QR$  factorization of  $A$  where  $R$  has positive diagonal entries.

## Proof.

Let  $A = Q_1R_1 = Q_2R_2$  be two  $QR$  factorizations, with  $A$  invertible.

Every matrix is invertible, and  $Q_1R_1 = Q_2R_2$  implies that  $F := Q_2^*Q_1 = R_2R_1^{-1}$ , the matrix  $F$  is unitary and upper triangular and thus it is a phase matrix.



# Uniqueness of the QR factorization

## Theorem

Let  $A \in GL(n)$  and let  $QR$  be a factorization of  $A$ .

For any  $QR$  factorization of  $A$ , say  $A = \tilde{Q}\tilde{R}$  there exists a phase matrix  $F$  such that  $\tilde{Q} = QF^*$  and  $\tilde{R} = FR$ . Conversely, any phase matrix  $F$  yields a  $QR$  factorization  $A = (QF^*)(FR)$ .

In particular, there exists unique a  $QR$  factorization of  $A$  where  $R$  has positive diagonal entries.

## Proof.

Conversely, if  $F$  is a phase matrix, then  $QF^*$  is unitary and  $FR$  is upper triangular, and since  $F^*F = I$ , we have that  $(QF^*)(FR)$  is a  $QR$  factorization of  $A$ .



# Uniqueness of the QR factorization

## Theorem

Let  $A \in GL(n)$  and let  $QR$  be a factorization of  $A$ .

For any QR factorization of  $A$ , say  $A = \tilde{Q}\tilde{R}$  there exists a phase matrix  $F$  such that  $\tilde{Q} = QF^*$  and  $\tilde{R} = FR$ . Conversely, any phase matrix  $F$  yields a QR factorization  $A = (QF^*)(FR)$ .

In particular, there exists unique a QR factorization of  $A$  where  $R$  has positive diagonal entries.

## Proof.

Consider the phase matrix  $F = \text{diag}(\bar{r}_{11}/|r_{11}|, \dots, \bar{r}_{nn}/|r_{nn}|)$ . We have that  $\tilde{R} := FR$  is upper triangular with positive diagonal entries and thus  $\tilde{Q}\tilde{R}$ , with  $\tilde{Q} = QF^*$ , is the required QR factorization.<sup>a</sup>

The uniqueness follows from the fact that any two QR factorization differs by a phase matrix and if not the identity, it would yield a triangular factor with some nonpositive diagonal entry. □

---

<sup>a</sup>Since  $\det(A) = \det(Q)\det(R)$ , we have that  $\det(R) \neq 0$  and  $r_{ii} \neq 0$ .

## Uniqueness of the QR factorization

In the real case, a similar result holds with  $Q$  orthogonal and  $R$  real.

When  $A$  is singular, there is no simple result on the uniqueness is not known. Think to the case  $A = 0$ , where every unitary matrix  $Q$  yields the QR factorization  $A = Q0$ .

## How to compute a QR factorization?

Recall the Gaussian elimination.

- 1 Apply plane rotations (Givens algorithm).
- 2 Apply reflections wrt hyperplanes (Householder algorithm).
- 3 Orthogonalize (Gram-Schmidt algorithm).

## Complex Givens matrices

A complex Givens matrix can be written as

$$G = \begin{bmatrix} c & -\bar{s} \\ s & c \end{bmatrix},$$

where  $c = \cos \varphi$ ,  $s = \psi \sin \varphi$  (with  $|\psi| = 1$ ).

If  $a, b \neq 0$ , then with

$$\psi = -\frac{b}{a} \left| \frac{a}{b} \right|, \quad \cos \varphi = \frac{|a|}{\sqrt{|a|^2 + |b|^2}}, \quad \sin \varphi = \frac{|b|}{\sqrt{|a|^2 + |b|^2}},$$

we have that  $G \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \alpha \\ 0 \end{bmatrix}$ . (The cases  $a = 0$  or  $b = 0$  can be treated directly.)

## Givens method

For any vector  $w \in \mathbb{C}^2$ , there exists a  $2 \times 2$  Givens  $G$  such that  
 $Gw$  is a multiple of  $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ .

Let  $v \in \mathbb{C}^n$ , there exists a Givens matrix  $G_{n-1,n} \in \mathbb{C}^{n \times n}$  acting on the plane  $(x_{n-1}, x_n)$  such that

$$v^{(2)} := G_{n-1,n}v = \begin{bmatrix} v_1^{(2)} \\ v_2^{(2)} \\ \vdots \\ v_{n-1}^{(2)} \\ 0 \end{bmatrix}.$$

Notice that  $v_1^{(2)} = v_1, \dots, v_{n-2}^{(2)} = v_{n-2}$ , that is, just the last two elements of  $v$  have been changed, by premultiplying it by  $G_{n-1,n}$ .

## Givens method

If  $n > 2$ , we continue by constructing a matrix  $G_{n-2,n-1}$  acting on the plane  $x_{n-2}, x_{n-1}$ , such that

$$v^{(3)} := G_{n-2,n-1} v^{(2)} = G_{n-2,n-1} G_{n-1,n} v = \begin{bmatrix} v_1^{(3)} \\ v_2^{(3)} \\ \vdots \\ v_{n-2}^{(3)} \\ 0 \\ 0 \end{bmatrix}.$$

## Givens method

And so on, we can construct  $G_{n-3,n-2}, \dots, G_{23}, G_{12}$  such that

$$G_{12} G_{23} \cdots G_{n-2,n-1} G_{n-1,n} v = \begin{bmatrix} \gamma_1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

Where  $|\gamma_1| = \|v\|$ , since rotations are isometries.

In the real case, the Givens matrices are plane rotations and we can choose the last rotation  $G_{12}$  so that  $\gamma_1 = \|v\|$ .

## Givens method

If  $v$  is the first column of a square matrix  $A$ , then we have

$$G_{12} G_{23} \cdots G_{n-2,n-1} G_{n-1,n} A = \begin{bmatrix} \gamma_1 & * \\ 0 & B_2 \end{bmatrix} := A_2.$$

We can produce a QR factorization of  $A$  by recursively repeating the same procedure for  $B_2$ , observing that there exist Givens matrices  $\tilde{G}_{12}, \dots, \tilde{G}_{n-2,n-1} \in U(n-1)$  such that

$$\tilde{G}_{12} \tilde{G}_{23} \cdots \tilde{G}_{n-2,n-1} B_2 = \begin{bmatrix} \gamma_2 & * \\ 0 & B_3 \end{bmatrix},$$

with  $B_3 \in \mathbb{C}^{(n-2) \times (n-2)}$ .

## Givens method

From a Givens matrix  $G \in U(n-1)$  acting on the plane  $x_i, x_j$ , we can construct a Givens matrix  $\widehat{G} \in U(n)$

$$\widehat{G} = \begin{bmatrix} 1 & 0 \\ 0 & G \end{bmatrix},$$

acting on the coordinates  $x_{i+1}, x_{j+1}$ .

Thus, by defining

$$G_{i+1,i+2}^{(2)} = \begin{bmatrix} 1 & 0 \\ 0 & \widetilde{G}_{i,i+1} \end{bmatrix}, \quad i = 1, \dots, n-2$$

we obtain

$$\begin{aligned} G_{23}^{(2)} G_{34}^{(2)} \cdots G_{n-1,n-2}^{(2)} A_2 &= \begin{bmatrix} \gamma_1 & * \\ 0 & \widetilde{G}_{12} \widetilde{G}_{23} \cdots \widetilde{G}_{n-2,n-1} B_2 \end{bmatrix} \\ &= \begin{bmatrix} \gamma_1 & * & * \\ 0 & \gamma_2 & * \\ 0 & 0 & B_3 \end{bmatrix} := A_3. \end{aligned}$$

## Givens method

The procedure can be repeated for  $n - 1$  steps, getting a **upper triangular** matrix  $R := A_n$ , such that

$$\underbrace{G_{n-1,n}^{(n-1)} G_{n-2,n-1}^{(n-2)} G_{n-1,n}^{(n-2)} \cdots G_{23}^{(2)} \cdots G_{n-1,n-2}^{(2)} G_{12}^{(1)} \cdots G_{n-1,n}^{(1)}}_{Q^*} A = R,$$

where we have denoted  $G_{ij}$  by  $G_{ij}^{(1)}$ .

The matrix  $Q^*$  is such that  $Q^*A = R$ , that is  $A = QR$  and we **have found a factorization** of  $A$  using  $\frac{1}{2}(n - 1)n$  Givens matrices.

## Example

Applying real Givens matrices to an orthogonal matrix  $A \in \mathcal{O}(n)$ , and choosing the last rotation at each step so that the diagonal element is positive, at the end we get a phase matrix of the type

$$R = \begin{bmatrix} 1 & 0 & \cdots & \cdots & 0 \\ 0 & 1 & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & 1 & 0 \\ 0 & \cdots & \cdots & 0 & \pm 1 \end{bmatrix},$$

where the bottom left element is 1 if  $A \in \mathcal{SO}(n)$ . This means that any  $n$  dimensional notation can be written as the product of  $\frac{1}{2}(n - 1)n$  plane rotations. The corresponding angles, provide a **local parametrization** of the differentiable manifold  $\mathcal{SO}(n)$ .

$U(n)$  requires  $n^2$  parameters.

## Computational cost

$$\underbrace{G_{n-1,n}^{(n-1)} G_{n-2,n-1}^{(n-2)} G_{n-1,n}^{(n-2)} \cdots G_{23}^{(2)} \cdots G_{n-1,n-2}^{(2)} G_{12}^{(1)} \cdots G_{n-1,n}^{(1)}}_{Q^*} A = R,$$

$\frac{1}{2}n(n - 1)$  matrix multiplications, but each multiplication modifies  
**just two rows.**

At the first step  $n - 1$  Givens rotations acting on two rows require  $6(n - 1)^2$  operations (excluding the ones needed to compute  $c$  and  $s$ )

At step  $k$  we have  $6(n - k)^2$  operations, for a total cost of

$$\sum_{k=1}^{n-1} 6(n - k)^2 = 2n^3 + o(n^3)$$

operations.

# Implementation

Apparently complicated but...

not really

```
for k = 1:n-1
    for h = n:-1:k+1
        A(h-1:h,:) = givens(A(h-1,k),A(h,k))*A(h-1:h,:);
    end
end
```

Only one instruction.

## Gram-Schmidt algorithm

Problem: find an orthogonal basis of a subspace  $\mathcal{W} \subset \mathbb{C}^n$ .

Let  $v_1, \dots, v_k \in \mathbb{C}^n$ . An orthonormal basis of  $\mathcal{W} = \text{span}\{v_1, \dots, v_k\}$  can be obtained by the following algorithm, for  $j = 1 : k$ ,

- Obtain  $u_j$  by removing from  $v_j$  the projection onto the span of the previous vectors  $w_1, \dots, w_{j-1}$ ;
- Normalize  $u_j$  to obtain  $w_j$ .

The set  $w_1, \dots, w_k$  is an orhtogonal basis of  $\mathcal{W}$ .

$$u_j = v_j - P_j v_j = v_j - \sum_{\ell=1}^{j-1} \frac{\langle w_\ell, v_j \rangle}{\langle w_\ell, w_\ell \rangle} w_\ell, \quad w_j = \frac{u_j}{\|u_j\|}.$$

## QR factorization and the Gram-Schmidt algorithm

$$u_j = v_j - P_j v_j = v_j - \sum_{\ell=1}^{j-1} \frac{\langle w_\ell, v_j \rangle}{\langle w_\ell, w_\ell \rangle} w_\ell, \quad w_j = \frac{u_j}{\|u_j\|}.$$

This simple algorithm can be seen as a matrix factorization.  
Indeed,

$$v_j = \sum_{\ell=1}^{j-1} (w_\ell^* v_j) w_\ell + \|u_j\| w_j = \sum_{\ell=1}^j r_{\ell j} w_\ell, \quad j = 1, \dots, k,$$

that, can be written in matrix form as

$$A = QR,$$

where  $A = [v_1 | \dots | v_k]$ ,  $Q = [w_1 | \dots | w_k]$  and

$$R = \begin{bmatrix} r_{11} & r_{12} & \cdots & r_{1k} \\ 0 & r_{22} & \ddots & r_{2k} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & r_{kk} \end{bmatrix}.$$

## QR factorization and the Gram-Schmidt algorithm

When  $k = n$ , we get a QR factorization of the matrix  $A$ .

The Gram-Schmidt algorithm can be extended to the case where the vectors  $v_1, \dots, v_k \in \mathbb{C}^n$  are not necessarily independent (but one of them is nonzero).

The difference in the step is that the vector  $u$  can be zero, and in this case, the step of the algorithm is just skipped. We will get a number of orthogonal vectors equal to the dimension of the  $\text{span}\{v_1, \dots, v_k\}$ .

### Exercise

*Explain how to get a QR factorization of a singular matrix using the Gram-Schmidt algorithm.*

## QR factorization and the Gram-Schmidt algorithm

In a finite arithmetic computation, the Gram-Schmidt algorithm is unstable, and a modified Gram-Schmidt algorithm is preferable

```
for i = 1 : k
    u = v(:, i);
    for j = 1 : i - 1 % this is empty if i = 1
        u = u - (u' * w(:,j)) * w(:, j);
    end
    w(:, i) = u / norm(u);
end
```

where the difference is that we project the current value of  $u$  on  $w_j$  and not  $v_j$ . The two algorithms are mathematically equivalent, but in finite arithmetic only the modified one is numerically stable.

### Exercise

*Determine the computational cost of the Gram-Schmidt method.*

## The Householder method

Householder matrices can be used to annihilate elements on a vector and this can be used to construct a QR factorization of a matrix.

Let  $v \in \mathbb{R}^n$ . There exist reflections  $P$ , with respect to a vector hyperplane, such that  $Pv = \alpha e_1$ . The vector hyperplane can be chosen orthogonal to  $w = v + \|v\|e_1$  (or  $w = v - \|v\|e_1$ ).

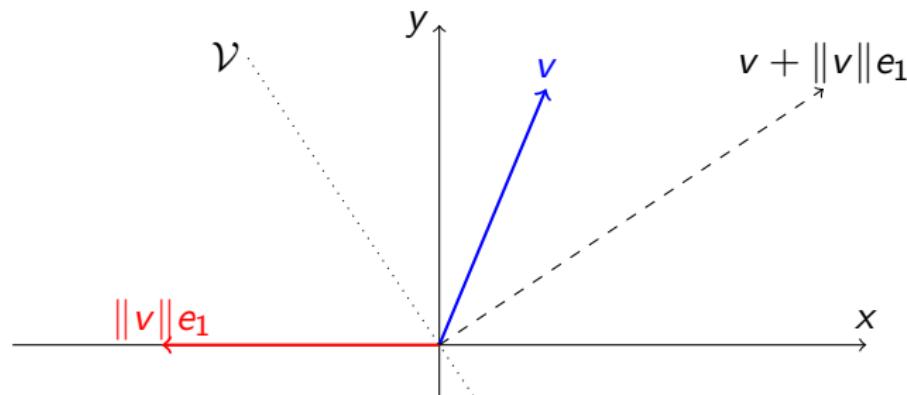


Figure: Moving the vector  $v \in \mathbb{R}^2$  to  $\|v\|e_1$  using a reflection with respect to a line  $\mathcal{V}$  orthogonal to  $v + \|v\|e_1$ .

## Householder method

In the complex case we choose the vector  $w$  such that  $w + \theta\|v\|e_1$ , where

$$\theta = \begin{cases} v_1/|v_1|, & v_1 \neq 0, \\ 1, & v_1 = 0. \end{cases}$$

(An equivalent choice is  $w - \theta\|v\|e_1$ .)

## Householder method

Using Householder matrices one can obtain a **QR factorization** of a matrix  $A \in \mathbb{C}^{n \times n}$  in  $n - 1$  steps.

At the first step, if the first column of  $A$ , say  $c_1$ , is not a multiple of  $e_1$ , then one constructs a Householder matrix  $P_1$  such that  $P_1 c_1 = \alpha_1 f_1$ , else set  $P_1 = I$ . We have

$$A_2 = P_1 A = \begin{bmatrix} \alpha_1 & * \\ 0 & B_2 \end{bmatrix}.$$

The procedure can be repeated finding an Householder matrix (or the identity)  $\tilde{P}_2 \in \mathbb{C}^{(n-1) \times (n-1)}$  such that  $\tilde{P}_2 c_2 = \alpha_2 e_1$ , where  $c_2$  is the first column of  $B_2$ .

## Householder method

Using Householder matrices one can obtain a **QR factorization** of a matrix  $A \in \mathbb{C}^{n \times n}$  in  $n - 1$  steps.

At the first step, if the first column of  $A$ , say  $c_1$ , is not a multiple of  $e_1$ , then one constructs a Householder matrix  $P_1$  such that  $P_1 c_1 = \alpha_1 f_1$ , else set  $P_1 = I$ . We have

$$A_2 = P_1 A = \begin{bmatrix} \alpha_1 & * \\ 0 & B_2 \end{bmatrix}.$$

The procedure can be repeated finding an Householder matrix (or the identity)  $\tilde{P}_2 \in \mathbb{C}^{(n-1) \times (n-1)}$  such that  $\tilde{P}_2 c_2 = \alpha_2 e_1$ , where  $c_2$  is the first column of  $B_2$ . Defining the matrix  $P_2 = \begin{bmatrix} 1 & 0 \\ 0 & \tilde{P}_2 \end{bmatrix}$  (that is still an Householder matrix or the identity), one has

$$A_3 = P_2 A_2 = \begin{bmatrix} 1 & 0 \\ 0 & \tilde{P}_2 \end{bmatrix} \begin{bmatrix} \alpha_1 & * \\ 0 & B_2 \end{bmatrix} = \begin{bmatrix} \alpha_1 & * & * \\ 0 & \alpha_2 & * \\ 0 & 0 & B_3 \end{bmatrix}.$$

## Householder method

The procedure can be repeated  $n - 1$  times obtaining at the end an **upper triangular matrix**

$$R = P_{n-1}P_{n-2} \cdots P_2P_1A,$$

from which we get the factorization

$$A = QR, \quad Q = (P_{n-1} \cdots P_1)^{-1} = P_1P_2 \cdots P_{n-1},$$

where we have used the fact that Householder matrices are unitary and Hermitian.

Since the vector parameter of the Householder matrices can be explicitly constructed, then the previous construction can be seen as an algorithm.

In the real case, at each step there are two choices, one can reflect using  $v + \|v\|f_1$  or  $v - \|v\|f_1$ . This is mathematically equivalent, and one can choose the reflection in order to get a positive diagonal entry on the matrix  $R$ . In practice, one chooses the sign that gives the better numerical stability, that is “plus” when  $v_1$  is positive and “minus” when  $v_1$  is negative.

# The Cartan-Dieudonné theorem.

The Householder method provide a constructive proof of the following.

## Theorem

*Any nontrivial linear isometry in  $\mathbb{R}^n$  is the composition of at most  $n$  reflections with respect to hyperplanes.*

## Proof.

Let  $Q$  be the orthogonal matrix associated with the isometry in an orthogonal basis of  $\mathbb{R}^n$ . Applying the Householder method to  $Q$ , where the reflections are chosen so that the first  $n - 1$  diagonal elements of  $R$  are positive, yields

$$Q = P_1 P_2 \cdots P_{n-1} R,$$

where  $P_1, \dots, P_{n-1}$  are reflection or the identity, while  $R$  is upper triangular and orthogonal with positive diagonal elements, except at most  $r_{nn}$ .



# The Cartan-Dieudonné theorem.

The Householder method provide a constructive proof of the following.

## Theorem

*Any nontrivial linear isometry in  $\mathbb{R}^n$  is the composition of at most  $n$  reflections with respect to hyperplanes.*

## Proof.

$$Q = P_1 P_2 \cdots P_{n-1} R,$$

The matrix  $R$  is a phase matrix and thus there are two cases:  
either  $R = I$  or  $R = \text{diag}(1, \dots, 1, -1)$ , that is a reflection.

Equation above is thus a decomposition of  $Q$  as the product of at most  $n$  reflections. □

# The Cartan-Dieudonné theorem.

The Householder method provide a constructive proof of the following.

## Theorem

*Any nontrivial linear isometry in  $\mathbb{R}^n$  is the composition of at most  $n$  reflections with respect to hyperplanes.*

## Proof.



The theorem is not true for the identity in dimension one since the isometry  $x \rightarrow x$ , corresponding to the matrix  $[1]$ , cannot be written as the product of one reflection.

## Implementation

The idea of the implementation is not to construct every reflection matrix, but to use just the parameter vectors.

The first step, with  $P_1 = I - \beta w_1 w_1^*$ , where  $\beta = 1/\|w_1\|^2$ , requires the product  $P_1 A$ . Without constructing  $P_1$ , one can compute the product as  $(I - \beta w_1 w_1^*)A = A - \beta w_1 (w_1^* A)$  and this requires

- one vector norm and one multiplication to get  $\beta$  ( $2n$  operations and one square root);
- one matrix multiplication  $u = A^* w_1$  ( $2n^2 - n$  operations);
- one scalar-vector product  $v = \beta w_1$  ( $n$  operations);
- the product  $vu^*$  ( $n^2$  operations);
- the sum  $A - vu^*$  ( $n^2$  operations).

The cost of the first step is then  $4n^2 + o(n^2)$  operations.

## Implementation

The  $k$ -th step require the same procedure in the right-lower  $(n - k + 1) \times (n - k + 1)$  block, say  $B_k$  and thus the cost is of about  $4(n - k + 1)^2$  operations.

The total asymptotic cost of the Householder algorithm for computing  $R$  is of  $\sum_{k=1}^{n-1} 4(n - k + 1)^2 \sim \frac{4}{3}n^3$  arithmetic operations. This allows one to solve a linear system, but **does not compare favorably with Gaussian elimination** whose cost is  $\frac{2}{3}n^3$  operations.

# Implementation

A pseudocode implementation is

```
for k = 1 : n-1
    B = A(k:n, k:n);
    w = B(:,k); % select the column k from the index k to the end
    w(1) = w(1) + sign(w(1)) * norm(w);
    beta = 2 / norm(w);
    u = w' * B;
    A(k:n, k:n) = B - (beta * w) * u;
end
```

If one needs the  $Q$  factor a further computation is needed.

## Exercise

*Describe a method to construct the factor  $Q$  of the QR factorization using Householder matrices and that require  $\frac{4}{3}n^3 + o(n^3)$  operations.*

## Thin QR factorization

If  $A \in \mathbb{C}^{m \times n}$  the algorithms for the QR factorization can be applied. If  $m > n$ ,  $n$  steps are needed, since also the last column should be processed.

We get  $A = QR$  with  $Q \in U(m)$  and  $R \in \mathbb{C}^{m \times n}$  upper triangular (a rectangular matrix  $R$  is upper triangular if  $r_{ij} = 0$  for  $i > j$ ).

When  $m > n$  we can block the factor as  $Q = [ Q_1 \quad Q_2 ]$  with  $Q_1 \in \mathbb{C}^{m \times n}$  and  $R = \begin{bmatrix} R_1 \\ 0 \end{bmatrix}$ , with  $R_1 \in \mathbb{C}^{n \times n}$ .

Using the block multiplication we discover that  $A = Q_1 R_1$ , where  $R_1$  is square and upper triangular. This factorization is said to be **thin QR factorization** and might be useful in some cases.

Notice that  $Q_1 \in \text{St}(m, n)$ .

## Fattorizzazione QR:

$A \in \mathbb{C}^{n \times n}$ ,  $A = QR$  fatt QR di A se Q è unitaria ( $Q^*Q = I \Leftrightarrow Q^{-1} = Q^* \Leftrightarrow QQ^* = I$ ) e R triangolare superiore

Ricordiamo che se Q unitaria  $\|Q\|_2 = 1$

$$\Rightarrow \sqrt{\rho(Q^*Q)} = 1 \Rightarrow \mu(A) = \|A'\| \|A\|$$

$$\mu(Q) = \|Q'\| \|Q\|$$

Condizionamento

$Q^*$  è unitaria poiché  $(Q^*)^* = Q^*Q = I$

Inoltre, A invertibile,

$$\|I\| = \|AA'\| \leq \|A\| \|A'\| = \mu(A)$$

Se la norma è indotta  $\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} \Rightarrow \|I\| = 1$   
 $\Rightarrow \mu(A) \geq 1$

OSS che il condizionamento dipende dalla norma, ed

esempio con  $\|\cdot\|_F$   $\|I\|_F = \sqrt{n} \Rightarrow \mu_F(A) \geq \sqrt{n}$

Condizionamento

Va bene comunque utilizzare qualsiasi norma

Se  $A \in \mathbb{R}^{n \times n}$ , allora  $Q \in O(n)$ ,  $R \in \mathbb{R}^{n \times n}$  triangolare sup  $\Rightarrow$  ortogonale ( $Q^TQ = I$ )

Se  $A \in \mathbb{C}^{n \times n}$   $A = QR$   $Q \in U(n)$   $R \in \mathbb{C}^{n \times n}$  unitaria  $r_{ij} = 0 \Leftrightarrow j > i$

Ora  $\uparrow$  scesa

$$\begin{pmatrix} 0 & \square \\ \square & 1 \end{pmatrix} \quad \begin{pmatrix} 0 & \square \\ \square & 0 \end{pmatrix}$$

Ci chiediamo se esiste lo fatt QR, e a differenza  
 dei fatt LU, lo QR esiste sempre, dimostrando  
 per induzione:

$$n=1 \quad A = QR = [1]A$$

$$n \geq 1 \quad (n-1 = \dim)$$

Sia  $v \in \mathbb{C}^n$  la prima colonna di  $A$

$$q_1 = \begin{cases} \begin{pmatrix} 1 \\ 0 \end{pmatrix} & \text{se } v=0 \\ \frac{v}{\|v\|} & \text{se } v \neq 0 \end{cases}$$

$\text{Span}^{\perp}\{q_1\}$  è un spazio di dimensione  $n-1$  e una sua base orthonormale è  $q_2, \dots, q_n$

$$\text{Costruiamo } U = \left[ q_1 \mid Q_2 \right], \quad Q_2 = \left( q_2 \mid \dots \mid q_n \right)$$

$U$  è unitaria (per verificare)

$$A = \left[ v \mid B \right]_P \in \mathbb{C}^{n \times (n-1)}$$

$$U^* A = \left[ \frac{q_1^*}{Q_2^*} \right] \left[ v \mid B \right] = \left[ \frac{q_1^* v}{Q_2^* v} \mid \frac{q_1^* B}{Q_2^* B} \right] = \left[ \frac{\|v\|}{0} \mid \frac{q_1^* B}{Q_2^* B} \right]$$

$$q_1^* v = \frac{v^* v}{\|v\|} = \frac{\|v\|^2}{\|v\|} = \|v\|$$

$$Q_2^* v = \left[ \frac{q_2^*}{\vdots} \right] v = \left[ \frac{q_2^* v}{q_n^* v} \right] = 0$$

Per ipotesi induttiva  $\exists \tilde{Q} \in U(n-1)$  e  $\tilde{R} \in \mathbb{C}^{(n-1) \times (n-1)}$  triangolare  
sup t.c.  $Q_2^* B = \tilde{Q} \tilde{R}$

Costruiamo

$$\hat{Q} = \begin{bmatrix} 1 & 0 \\ 0 & \tilde{Q}^* \end{bmatrix} \in U(n)$$

$$\hat{Q}(U^*A) = \begin{bmatrix} I & 0 \\ 0 & \tilde{Q} \end{bmatrix} \left[ \begin{array}{c|c} I_{n-e} & Q_1^*B \\ \hline 0 & Q_2^*B \end{array} \right] = \begin{bmatrix} I_{n-e} & Q_1^*B \\ 0 & \tilde{Q}^*(Q_2^*B) \end{bmatrix}.$$

ove  $\tilde{R} = \tilde{Q}^*Q_2^*B$  ( $\hat{Q}U^*A = R \Rightarrow A = \underbrace{U^*\tilde{Q}}_Q R$ )

$$\Rightarrow \hat{Q}(U^*A) = \begin{bmatrix} I_{n-e} & * \\ 0 & \tilde{R} \end{bmatrix} = R \quad \square$$

Dimostriamo che se  $U \in U(n)$   $\begin{bmatrix} I_{n-e} & 0 \\ 0 & U \end{bmatrix}$  è unitone

$$\begin{bmatrix} I_{n-e} & 0 \\ 0 & U \end{bmatrix}^* \begin{bmatrix} I_{n-e} & 0 \\ 0 & U \end{bmatrix} = \begin{bmatrix} I_{n-e} & 0 \\ 0 & U^* \end{bmatrix} \begin{bmatrix} I_{n-e} & 0 \\ 0 & U \end{bmatrix} =$$

$$= \begin{bmatrix} I_{n-e} & 0 \\ 0 & U^* \end{bmatrix} \begin{bmatrix} I_{n-e} & 0 \\ 0 & U \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & U^*U \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & I \end{bmatrix} = I$$

in  $\hat{Q}$  è unitone

Se  $U = [q_1 | \dots | q_n]$  unitone  $\Rightarrow q_1, \dots, q_n$  base orthonormata di  $\mathbb{C}^n$

Vediamo ora l'unicità della fattorizzazione QR:

$n=1$   $a = 1 \cdot a$  oppure  $a = (-1)(-a)$  oppure  
 $a = i(-ia)$

Le matrici unitonie  $1 \times 1$  sono:

$$\text{U}(1) = \{x^*x = 1 \Rightarrow \bar{x}x = 1 \Leftrightarrow |x|^2 = 1 \Rightarrow |x| = 1\}$$

$\Rightarrow U(1) = \{e^{i\theta} | \theta \in \mathbb{R}\}$  è un gruppo in  $\mathbb{C}$

Ci sono infinite soluzioni per  $n=1$

o ha elementi positivi nel caso IR

Se  $a \neq 0$   $a = \frac{a}{|a|} \cdot |a|$  questa è unica

$$\underbrace{a}_{Q} \quad \underbrace{|a|}_{R}$$

Allora supponiamo di avere

$$A = Q_1 R_1 = Q_2 R_2, \text{ se } A \in G \subset (n, \mathbb{C})$$
$$\Rightarrow R_1, R_2 \text{ invertibili } (r_{ii} \neq 0 \ \forall i = 1, \dots, n)$$

$$\Rightarrow Q_1^* Q_2 R_2 R_1^{-1} = Q_1^* Q_2 R_2 R_2^{-1} \Rightarrow R_2 R_2^{-1} = Q_1^* Q_2 = F$$

Se ci sono due fatti  $QR = F$  è contemporaneamente  
invertibile e triangolare superiore

$$\begin{bmatrix} r_{11} & r_{12} & \cdots & r_{1n} \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \ddots & \ddots & r_{nn} \end{bmatrix} = \begin{bmatrix} e^{i\theta_1} & 0 & \cdots & 0 \\ 0 & e^{i\theta_2} & \cdots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & e^{i\theta_n} \end{bmatrix}$$

Se fosse anche invertibile la somma delle vette (i)  
monino 1 è anche le righe

$$\Rightarrow \sqrt{\sum_{i=0}^n |r_{1i}|^2} = 1 \Rightarrow |r_{11}| = 1 \Rightarrow r_{11} = e^{i\theta_1}$$

$$\Rightarrow \sqrt{|e^{i\theta_1}|^2 + |r_{12}|^2 + \cdots + |r_{1n}|^2} = 1$$

" "

$$\Rightarrow |r_{12}|^2 + \cdots + |r_{1n}|^2 = 0$$

$\Rightarrow F$  è diagonale con elementi di moduli 1 sulla  
diagonale. Una matrice del genere è detta di fase

# Teorema (slide)

... e v

Se la matrice non è invertibile ⇒ + complicato

Come calcolare lo fatt QR?

(cerca un algoritmo migliore di Gauss inverso se lo trovi)

Faccia rotazioni e riflessioni ( $\Rightarrow$  da Gauss che è -cessazione di righe)

$$A = QR \Rightarrow Q^* A = R$$

$$Q^* \begin{bmatrix} c_1 \\ \vdots \\ c_n \end{bmatrix} = \begin{bmatrix} r_1 \\ \vdots \\ r_n \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} Q^* c_1 \\ \vdots \\ Q^* c_n \end{bmatrix} = \begin{bmatrix} r_1 \\ \vdots \\ r_n \end{bmatrix}$$

trovare rotazioni per cui  $\sum c_i = r_i \quad i=1, \dots, n$

Ci sono 3 modi per calcolare lo fatt QR

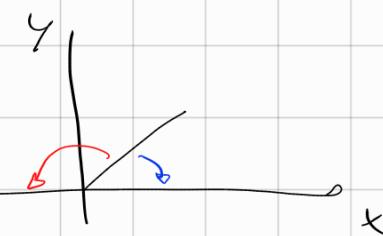
Usiamo le rotazioni planari:

Given's matrices

$$\text{Data } V = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \in \mathbb{R}^2$$

esiste una rotazione lineare

che manda  $V$  sull'asse  $x$ ? Sì,



$$\det Q = \pm 1$$

$$\mathcal{O}(n) = \{ Q \in \mathbb{Q}^{n \times n} \mid Q^T Q = I \}$$

$$1 = \det Q^T Q = \det Q^T \det Q = \det^2 Q$$

$\text{SO}(n) = \left\{ Q \in \text{O}(n), \det Q = 1 \right\}$  est un groupe  
groupe orthogonal spécial  
(rotations)

Le noyau  $\text{SO}(2)$  sera  $\left\{ \begin{bmatrix} \cos \theta & -\sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}, \theta \in [0, 2\pi] \right\}$ ,  
noyau d'angle  $\theta$  autour de l'origine

La matrice complexe di Givens è

$$G = \begin{bmatrix} c & -s \\ s & c \end{bmatrix} + \text{sliders.}$$

dove hanno  $Q$  per cui:

$$\begin{bmatrix} c & -s \\ s & c \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \alpha \\ 0 \end{bmatrix} \quad \dots$$

Ma questa rotazione mi permette di trovare  
la fattorizzazione QR.

8/10

$A = QR$ ,  $Q \in \text{U}(n)$  B triangolare superiore

$$\begin{bmatrix} a \\ b \end{bmatrix} \in \mathbb{R}^2 \rightarrow \begin{bmatrix} \alpha \\ 0 \end{bmatrix}$$


$$|\alpha| = \sqrt{a^2 + b^2}$$

Matrice di Givens.

$$\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} c & -s \\ s & c \end{bmatrix} \quad c^2 + s^2 = 1$$

$$\begin{bmatrix} c & -s \\ s & c \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \alpha \\ 0 \end{bmatrix} \quad \begin{cases} ca - sb = \alpha \\ sa + cb = 0 \end{cases}$$

Se  $b=0 \Rightarrow c=1$  e  $s=0$  ( $c=-1, s=0$ )

$$\text{Se } b \neq 0 \Rightarrow c = -\frac{sa}{b} = 0 - \frac{sa^2}{b} - sb = 2$$

$$-\frac{s(a^2+b^2)}{b} = \pm \sqrt{a^2+b^2} \Rightarrow s = \mp \frac{b}{\sqrt{a^2+b^2}}$$

$$c = \pm \frac{a}{\sqrt{a^2+b^2}}$$

CASO COMPLESSO:

$$\begin{bmatrix} c & -s \\ s & c \end{bmatrix}$$

$\leftarrow$  Nella stessa rotazione nel caso complesso  
se  $c = \cos \varphi = \frac{|a|}{\sqrt{|a|^2+|b|^2}}$ ,  $s = \operatorname{sen} \varphi \cdot e^{i\theta}$

$$a \neq 0 \neq b$$

$$\operatorname{sen} \varphi = \frac{|b|}{\sqrt{|a|^2+|b|^2}}, \varphi = \frac{a}{|a|} \cdot \frac{|b|}{b}$$

$$\text{Se } a=0 \quad \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}, c=0, s=\pm 1$$

Vediamo una rotazione autonoma su  $\mathbb{R}^3$

• Sul piano  $(x_1, x_2)$

$$\begin{bmatrix} \cos \varphi & -\operatorname{sen} \varphi & 0 \\ \operatorname{sen} \varphi & \cos \varphi & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} G \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ x_3 \end{bmatrix}$$

• Sul piano  $(x_2, x_3)$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \varphi & -\operatorname{sen} \varphi \\ 0 & \operatorname{sen} \varphi & \cos \varphi \end{bmatrix}$$

• Sul piano  $(x_1, x_3)$

$$\begin{bmatrix} \cos \varphi & 0 & -\operatorname{sen} \varphi \\ 0 & 1 & 0 \\ \operatorname{sen} \varphi & 0 & \cos \varphi \end{bmatrix}$$

Definiamo allora la matrice di Givens su  $\mathbb{C}^n$  nel caso dei complessi:

Una matrice di Givens e' del tipo

rotazione  $i, j$

$$G_{ij} = \begin{bmatrix} 1 & 0 & & 0 \\ 0 & 1 & -\bar{s} & \\ & s & c & \\ 0 & 0 & & 1 \end{bmatrix}$$

Agisce nel piano  
( $x_i, x_j$ )

$$G_{ij} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_{i-1} \\ cx_i - \bar{s}x_j \\ x_{i+1} \\ \vdots \\ x_{j-1} \\ sx_i + cx_j \\ x_{j+1} \\ \vdots \\ x_n \end{bmatrix}$$

$$(G_{ij})_{lkl} = \begin{cases} 1 & l=i, l=j \\ c & l=i, l=j \\ s & l=j, l=i \\ -\bar{s} & l=i, l=j \\ 0 & \text{altrimenti} \end{cases}$$

Metodo di Givens per la fattorizzazione QR:

- PASSO 1: Come estrarre gli elementi di un vettore tramite "rotazione" (per i complessi non e' rotazione)

$$v \in \mathbb{C}^n$$

$$v = \begin{bmatrix} v_1^{(1)} \\ \vdots \\ v_n^{(1)} \end{bmatrix}$$

Trovare matrici di Givens t.c.

$$Gv = \begin{bmatrix} x \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$Gv = G \begin{bmatrix} v_n^{(1)} \\ v^{(1)} \end{bmatrix} = \begin{bmatrix} v_{n-1}^{(2)} \\ 0 \end{bmatrix}$$

$$G_{n-1,n} = \begin{bmatrix} 1 & & & \\ & \ddots & & 0 \\ & & 1 & \\ 0 & & & G \end{bmatrix}$$

$$G_{n-1,n} \begin{bmatrix} v_1^{(1)} \\ \vdots \\ v_n^{(1)} \end{bmatrix} = \begin{bmatrix} v_1^{(1)} \\ \vdots \\ v_{n-2}^{(1)} \\ v_{n-1}^{(1)} \\ 0 \end{bmatrix} = \begin{bmatrix} v_1^{(2)} \\ \vdots \\ v_{n-2}^{(2)} \\ v_{n-1}^{(2)} \\ 0 \end{bmatrix}$$

PASSO 2:

$$G \begin{bmatrix} v_{n-2}^{(2)} \\ v_{n-1}^{(2)} \end{bmatrix} = \begin{bmatrix} v_{n-2}^{(3)} \\ 0 \end{bmatrix}$$

$$G_{n-2,n-1} = \begin{bmatrix} 1 & & & \\ & \ddots & & 0 \\ & & 1 & \\ 0 & & & G \end{bmatrix}$$

$$G_{n-2,n-1} \begin{bmatrix} v_1^{(2)} \\ \vdots \\ v_{n-1}^{(2)} \\ 0 \end{bmatrix} = \begin{bmatrix} v_1^{(2)} \\ \vdots \\ v_{n-3}^{(2)} \\ v_{n-2}^{(2)} \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} v_1^{(3)} \\ \vdots \\ v_{n-2}^{(3)} \\ 0 \\ 0 \end{bmatrix}$$

Dopo  $n-1$  passi  $\Rightarrow$  ottiene un vettore

$$G_{12} \cdot G_{23} \cdot G_{34} \cdots G_{n-2,n-1} \cdot G_{n-1,n} \cdot v = \begin{bmatrix} \lambda \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$$

Se  $v \in \mathbb{R}^n$  è possibile scegliere l'ultima  $G_{12}$  t.c.  $\lambda > 0$

Se poniamo ora da  $A_1 = [v | B_1]$  e applico  $G_{12}, \dots, G_{n-1,n}$

$$\Rightarrow \underbrace{G_{12} \cdots G_{n-1,n}}_Q \cdot A_1 = Q [v | B_1] = [Q_1 v | Q_1 B_1] =$$

$$= \left[ \begin{array}{c|c} v & * \\ \hline 0 & B_2 \end{array} \right] = A_2$$

Possiamo interpretare questo come la prima passo di un algoritmo per calcolare la fattorizzazione QR di A.

Sia  $v^{(1)}$  la prima colonna di A, si conoscano  $G_{12}^{(1)}, \dots, G_{n-1,n}^{(1)}$  e le applichiamo ad A

$$G_{12}^{(1)} \cdots G_{n-1,n}^{(1)} v^{(1)} = \left[ \begin{array}{c} v \\ 0 \end{array} \right]$$

$$\Rightarrow A_2 = G_{12}^{(1)} \cdots G_{n-1,n}^{(1)} A_1 = \left[ \begin{array}{c|c} v & * \\ \hline 0 & B_2 \end{array} \right] = A_2$$

$\uparrow C^{(n-1) \times (n-1)}$

$$A_2 = \left[ \begin{array}{c|c} v & * \\ \hline 0 & B_2 \end{array} \right], \text{ sia } v^{(2)}$$

la prima colonna di  $B_2$

$\exists \tilde{G}_{12}^{(2)}, \dots, \tilde{G}_{n-2,n-1}^{(2)}$  di Givens t.c.

$$\tilde{G}_{12}^{(2)} \cdots \tilde{G}_{n-2,n-1}^{(2)} v^{(2)} = \left[ \begin{array}{c} v \\ 0 \end{array} \right] \Rightarrow B_2 = \left[ \begin{array}{c|c} v & * \\ \hline 0 & B_3 \end{array} \right]$$

Per ogni  $i$  su  $A_2$  si conoscano delle Givens estese ad  $A_2$  che lasciano la prima componente invariata:

$$G_{23}^{(2)} = \begin{bmatrix} 1 & 0 \\ 0 & \tilde{G}_{22}^{(2)} \end{bmatrix}, \dots, G_{i,i+1}^{(2)} = \begin{bmatrix} 1 & 0 \\ 0 & \tilde{G}_{i-i,i}^{(2)} \end{bmatrix}$$

con  $i = 2, \dots, n-1$

Il secondo passo è dunque

$$G_{23}^{(2)} G_{34}^{(2)} \dots G_{n-1,n}^{(2)} A_2 = \left[ \begin{array}{c|c} \gamma_3 & * \\ \hline 0 & \left[ \begin{array}{c|c} \gamma_2 & * \\ \hline 0 & B_3 \end{array} \right] \end{array} \right] = A_3$$

Al passo l sono necessarie  $n-l$  rotazioni:

Sia  $v^{(l)}$  la prima colonna di  $B_l$ , allora

$$\tilde{G}_{22}^{(2)} \dots \tilde{G}_{n-l,n-l+1}^{(2)} v^{(l)} = \begin{bmatrix} \gamma_l \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$G_{l,l+1}^{(l)} = \begin{bmatrix} 1 & 0 \\ 0 & \tilde{G}_{l,l+1}^{(l)} \end{bmatrix} \dots$$

$$A_{l+1} = G_{l,l+1}^{(l)} \dots G_{n-1,n}^{(l)} A_l \quad l=1, \dots, n-1$$

$$R = A_n = G_{n-1,n}^{(n-1)} A_{n-1} = G_{n-1,n}^{(n-1)} G_{n-2,n-1}^{(n-2)} G_{n-1,n}^{(n-2)} A_{n-2} =$$

$$= G_{n-1,n}^{(n-1)} \cdot G_{n-2,n-1}^{(n-2)} \cdot G_{n-1,n}^{(n-2)} \cdots G_{l,l+1}^{(l)} \cdots G_{l-1,n}^{(l)} \cdots G_{22}^{(2)} \cdots G_{n-1,n}^{(1)} \cdot A$$

$$R = U A \Rightarrow A = U^* R, \quad A = QR$$

In tutto si hanno  $n-l$  al passo l

$$\Leftrightarrow \sum_{l=1}^{n-1} (n-l) = \sum_{l=1}^{n-1} l = \frac{n(n-1)}{2}$$

A differenza della fct LU, lo QR  $\nexists$  sempre perché R non ha algoritmo men ha condizioni (ad es Gauss aveva PIVOT  $\neq 0$ )

$G = \text{givens } (a, b)$

For  $l = 1 : n-1$

% passo l

For  $i = l : n-l$

% "costruiremo"  $G_{i,i+1}^{(l)}$  e "moltiplichiamo" per A

$G = \text{givens } (A(i, l), A(i, l+1))$ ;

$A(i:i+1, l:n) = G * A(i:i+1, l:n)$ ;

Costo computazionale:

$6(n-l)$  operazioni

$\curvearrowleft$  2 costruzioni lunghe di righe

$$\begin{aligned} &= \sum_{l=1}^{n-1} \sum_{i=l}^{n-1} 6(n-l+1) = \sum_{l=1}^{n-1} 6(n-l+1)(n-l) = \\ &= \sum_{l=1}^{n-1} 6(n-l)^2 + \sum_{l=1}^{n-1} 6(n-l) \sim \frac{6n^3}{3} + \Theta(n^3) \sim 2n^3 \end{aligned}$$

Questo costo computazionale è un po' più alto di quello dell'algoritmo di Gauss, ma ha motivi evidenti ed ha informazioni in più.

Se applica l'algoritmo su una matrice elettiva

$$Q \in \mathbb{R}^{3 \times 3} \quad Q^T Q = I$$

$$\begin{bmatrix} q_{11} & q_{12} & q_{13} \\ q_{21} & q_{22} & q_{23} \\ q_{31} & q_{32} & q_{33} \end{bmatrix}$$

per una rotazione  $(y, z) \rightarrow$

$$\begin{bmatrix} q_{11} & q_{12} & q_{13} \\ \tilde{q}_{21} & \tilde{q}_{22} & \tilde{q}_{23} \\ 0 & \tilde{q}_{32} & \tilde{q}_{33} \end{bmatrix}$$

rot  $(x, y)$

$$\begin{bmatrix} x & 0 & 0 \\ 0 & \hat{q}_{22} & \hat{q}_{23} \\ 0 & \tilde{q}_{32} & \tilde{q}_{33} \end{bmatrix}$$

rot  $(y, z)$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \tilde{q}_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \pm 1 \end{bmatrix}$$

$\Rightarrow$  con 3 rotazioni sono riusciti

pratica  
elettiva

ad ottenere una matrice elettiva con det  $\neq 1$

Se  $Q \in SO(3) \exists G_{23}, G_{12}, \tilde{G}_{23}$  t.c

$$G_{23} G_{12} \tilde{G}_{23}^T Q = I \Rightarrow Q = \tilde{G}_{23}^T G_{12}^T G_{23}^T$$

$SO(3)$  è varietà topologica di dim 3 (si dimostra che  
è differenziabile)

Possiamo fattorizzare una matrice  $SO(n)$  con  $n^2$  angoli  
(parametri)

Fattorizzazione QR magra (Thin QR factorization)

Se  $A \in \mathbb{C}^{m \times n}$ ,  $\exists U \in U(m)$ ,  $R \in \mathbb{C}^{m \times n}$  t.c

$$r_{i,j} = 0 \quad i > j, \quad A = QR$$

strong rep

$$A = QR = \left[ \underbrace{Q_1}_{n} \mid \underbrace{Q_2}_{m-n} \right] \left[ \begin{array}{c|c} R_1 & \\ \hline 0 & \end{array} \right]_{n \times n} = Q_1 R_1 + Q_2 0$$

Osserviamo che  $Q_1 \in \text{St}(m, n, \mathbb{C}) = \{U \in \mathbb{C}^{m \times n}, U^T U = I\}$   
In Vettori

Def

16/10

$A \in \mathbb{C}^{m \times n}$  è detta di Hessenberg superiore se  
 $a_{ij} = 0 \quad i > j + 1$

Metodo di Householder:

impiega  $\frac{4}{3}n^3$  operazioni per R e  $\frac{4}{3}n^3$  per Q

$A \in \mathbb{C}^{m \times n} \quad m \geq n$

Supponiamo di avere una fattorizzazione QR di A

Trovare le queste di minima norma  $\|Ax - b\|$

$$\|Ax - b\|^2 = \|QRx - \underbrace{QQ^*b}_{\text{"I" }}\|^2 = \|\cancel{Q}(Rx - Q^*b)\|^2 =$$

$$\begin{bmatrix} R_1 \\ 0 \end{bmatrix} \begin{bmatrix} Q_1^* \\ Q_2^* \end{bmatrix}$$

$$Q = \underbrace{\begin{bmatrix} Q_1 & Q_2 \end{bmatrix}}_m$$

$$= \left\| \begin{bmatrix} R_1 x - Q_1^* b \\ Q_2^* b \end{bmatrix} \right\|^2 = \|R_1 x - Q_1^* b\|^2 + \|Q_2^* b\|^2$$

Ma x è presente solo nel primo addendo

$$\Rightarrow \arg \min \|Ax - b\|^2 = \arg \min \|R_1 x - Q_1^* b\|^2$$

$$\text{Se } \text{rg } A = n, \quad A = Q R, \quad R = \begin{bmatrix} R_1 \\ 0 \end{bmatrix}$$

$R_1$  è invertibile

$\Rightarrow R_A^* b = Q^* b$  ha sol unica  $\Rightarrow$  c'è anche soluzione dei minimi quadrati.

Una base ortogonale del sottospazio è una base dello spazio ortogonale:

$$V \subset \mathbb{C}^n \\ V = \text{Span} \{ v_1, \dots, v_l \}, \quad v_1, \dots, v_l \perp \perp$$

(Si trova una base ortogonale con l'algoritmo V )  
 Gram-Schmidt in geometria

$$V = [v_1 | v_2 | \dots | v_l] \in \mathbb{C}^{n \times l}$$

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

↑  
fattorizzazione  
magna

$$\Rightarrow v_i = \sum_{j=1}^l q_j e_j$$

$q_1, \dots, q_l$  base di  $V$ ,  $q_j^* q_k = 0 \quad j \neq k$

$\Rightarrow q_{l+1}, \dots, q_n$  base di  $V^\perp$

Algoritmo di Gram-Schmidt:

$$v_1, \dots, v_l \quad (\perp \perp \Leftrightarrow v_i \neq 0 \forall i)$$

Meglio trovare una base ortonormale  $w_1, \dots, w_l$

$$w_1 = \frac{v_1}{\|v_1\|}$$

$$w_2 = \frac{v_2 - \langle v_2, w_1 \rangle w_1}{\|v_2 - \langle v_2, w_1 \rangle w_1\|}$$

$w \neq 0$  perché  $\perp \perp$

$$\Rightarrow w_i = \frac{v_i - \sum_{j=1}^i \langle v_i, w_j \rangle w_j}{\|v_i - \sum_{j=1}^i \langle v_i, w_j \rangle w_j\|}$$

OSS:

$$v_2 = \|v_2\|w_2 + \langle v_2, w_2 \rangle w_2 \quad (\text{è comb lineare di } w_1 \text{ e } w_2)$$

$$v_i = \sum_{j=1}^i \langle v_i, w_j \rangle w_j + \|v_i\| w_i \quad \begin{array}{l} \text{è comb lineare} \\ \text{e } w_i \text{ non normalizzata} \end{array}$$

$h_{j,i}$

$$W = [w_1 | \dots | w_l]$$

$$v_i = [w_1 | \dots | w_l] \begin{bmatrix} h_{1,i} \\ h_{2,i} \\ \vdots \\ h_{n,i} \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\Rightarrow V = [v_1 | \dots | v_e] = [w_1 | \dots | w_l] \begin{bmatrix} h_{1,1} & h_{1,2} & \dots & h_{1,l} \\ h_{2,1} & h_{2,2} & \dots & h_{2,l} \\ \vdots & \vdots & \ddots & \vdots \\ h_{e,1} & h_{e,2} & \dots & h_{e,l} \end{bmatrix}$$

factorizzazione QR magra

Costo computazionale

2n OPS per  $\langle \cdot, \cdot \rangle$

n OPS quando  $\langle \cdot, \cdot \rangle \cdot w_i$

n OPS con le somme

$$\Rightarrow \text{COSTO} \sum_{i=1}^l \sum_{j=1}^{i-1} 1 + \sum_{i=1}^l (2n + \gamma) =$$

$n$  OPS per la rotazione  
per moltiplicare e

## The SVD approach

The linear least squares can be solved with SVD. See the next lecture.

$$= 4n \sum_{i=1}^l i - 2n \ln \frac{e^2}{2} = 2nl^2$$

Se  $l=n \Rightarrow 2n^3$  (caso Givens)

L'algoritmo di Gram-Schmidt è però poco stabile

SVD

15/10

$$A \in \mathbb{C}^{m \times n}, A = U \Sigma V^*, U \in \mathbb{U}(n), V \in \mathbb{U}(m), \Sigma \in \mathbb{R}^{m \times n} \text{ t.c. } \sigma_{ij} = 0 \quad i \neq j \\ \sigma_{ii} = \sigma_i \quad i = 1, \dots, p, \quad p = \min\{m, n\} \\ \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$$

$$[A] = [U] \begin{bmatrix} \sigma_1 & & & 0 \\ & \ddots & & \\ 0 & & \sigma_p & \\ & & & \ddots & 0 \\ 0 & & & & \vdots \end{bmatrix} [V^*]$$

$$[A] = [U] \begin{bmatrix} \sigma_1 & & & 0 \\ & \ddots & & \\ 0 & & \sigma_n & \\ & & & \ddots & 0 \\ 0 & & & & \vdots \end{bmatrix} [V^*]$$

$\sigma_1, \dots, \sigma_p$  sono i valori singolari di  $A$   
le colonne di  $U$  e  $V$  sono i vettori singolari

PCA principal component analysis

Teorema:

SLIDE

(Teorema simile al Teorema spettrale)

Dimostrazione:

Per induzione su  $n$  e per ogni  $m$

• Se  $A=0 \Rightarrow A=UOV^*$  è una SVD

• Se  $A \neq 0 \Rightarrow A=U\Sigma V^*$

$$\Rightarrow A^* = V\Sigma^* U$$

$\Sigma = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix}$

è una SVD.  $\left( \begin{array}{l} \text{suppongo } m \geq n \\ \text{senza perdita di gen} \end{array} \right)$

Se  $n=1$ ,  $A \in \mathbb{C}^{n \times 1} \Rightarrow \exists$  fatt QR

$$\Rightarrow A = Q \begin{bmatrix} \alpha \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} 1 \end{bmatrix} = \alpha - \|A\| > 0 \quad A \neq 0$$
$$= Q \begin{bmatrix} |\alpha| \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} \frac{\alpha}{|\alpha|} \end{bmatrix}$$

$\alpha = |\alpha| \frac{\alpha}{|\alpha|} \rightarrow U(1)$

$\Rightarrow$  O.K.

$n-1 = DN$

$A \in \mathbb{C}^{m \times n}$ ,  $m \geq n$ ,  $A \neq 0$   $\left( \|A\|_2 = \max_{\|x\|=1} \|Ax\| \right)$

$\exists x \in \mathbb{C}^n$  s.t.  $\|Ax\| = \|A\| = \sigma_1$

$$y = \frac{Ax}{\sigma_1}, \|y\| = \frac{\|Ax\|}{\sigma_1} = 1$$

$\Rightarrow$  slide

## References

The material of this presentation (and the missing proofs) can be found in

- D. Bini, M. Capovani, O. Menchi, *Metodi Numerici per l'Algebra Lineare*, Zanichelli, Bologna, 1982 (Capitolo 7), in Italian.
- J. Stoer, R. Burlisch, *Introduction to Numerical Analysis, Third edition*, Texts in Applied Mathematics, 12. Springer-Verlag, New York, 2002, in English.

## Regression

$Y = \beta_1 X + \beta_0$ , with

$$\beta_1 = \frac{\text{cov}(X, Y)}{\text{var}(X)}, \quad \beta_0 = \mathbb{E}[Y] - \beta_1 \mathbb{E}[X].$$

## Regression

$Y = \beta_1 X + \beta_0$ , with

$$\beta_1 = \frac{\text{cov}(X, Y)}{\text{var}(X)}, \quad \beta_0 = \mathbb{E}[Y] - \beta_1 \mathbb{E}[X].$$

Observe that

$$M = \begin{bmatrix} 2n\mathbb{E}[X^2] & 2n\mathbb{E}[X] \\ 2n\mathbb{E}[X] & 2n \end{bmatrix}, \quad c = \begin{bmatrix} 2n\mathbb{E}[XY] \\ 2n\mathbb{E}[Y] \end{bmatrix}.$$

We claim that  $M \begin{bmatrix} \beta_1 \\ \beta_0 \end{bmatrix} = c$ :

## Regression

$$Y = \beta_1 X + \beta_0, \text{ with}$$

$$\beta_1 = \frac{\text{cov}(X, Y)}{\text{var}(X)}, \quad \beta_0 = \mathbb{E}[Y] - \beta_1 \mathbb{E}[X].$$

Observe that

$$M = \begin{bmatrix} 2n\mathbb{E}[X^2] & 2n\mathbb{E}[X] \\ 2n\mathbb{E}[X] & 2n \end{bmatrix}, \quad c = \begin{bmatrix} 2n\mathbb{E}[XY] \\ 2n\mathbb{E}[Y] \end{bmatrix}.$$

We claim that  $M \begin{bmatrix} \beta_1 \\ \beta_0 \end{bmatrix} = c$ :

$$2n\mathbb{E}[X]\beta_1 + 2n\mathbb{E}[Y] - 2n\mathbb{E}[X]\beta_1 = 2n\mathbb{E}[Y]$$

## Regression

$$Y = \beta_1 X + \beta_0, \text{ with}$$

$$\beta_1 = \frac{\text{cov}(X, Y)}{\text{var}(X)}, \quad \beta_0 = \mathbb{E}[Y] - \beta_1 \mathbb{E}[X].$$

Observe that

$$M = \begin{bmatrix} 2n\mathbb{E}[X^2] & 2n\mathbb{E}[X] \\ 2n\mathbb{E}[X] & 2n \end{bmatrix}, \quad c = \begin{bmatrix} 2n\mathbb{E}[XY] \\ 2n\mathbb{E}[Y] \end{bmatrix}.$$

We claim that  $M \begin{bmatrix} \beta_1 \\ \beta_0 \end{bmatrix} = c$ :

$$2n\mathbb{E}[X]\beta_1 + 2n\mathbb{E}[Y] - 2n\mathbb{E}[X]\beta_1 = 2n\mathbb{E}[Y]$$

$$\begin{aligned} 2n\mathbb{E}[X^2]\beta_1 + 2n\mathbb{E}[X]\mathbb{E}[Y] - 2n\mathbb{E}[X]^2\beta_1 \\ = 2n \left( \text{var}(X) \frac{\text{cov}(X, Y)}{\text{var}(X)} + \mathbb{E}[X]\mathbb{E}[Y] \right) = 2n\mathbb{E}[XY]. \end{aligned}$$



## Hilbert spaces approximation

$\mathcal{H} = L^2_\mu$ , with  $\int f d\mu = \sum_i f(x_i)$ .

Find the best approximation on  $\mathcal{F} = \text{span}\{\varphi_0, \varphi_1\}$ ,  $\varphi_0 = 1$ ,  $\varphi_1 = x$  of  $f$  such that  $f(x_i) = y_i$ .

$$A = \begin{bmatrix} \langle \varphi_1, \varphi_1 \rangle & \langle \varphi_0, \varphi_1 \rangle \\ \langle \varphi_1, \varphi_0 \rangle & \langle \varphi_0, \varphi_0 \rangle \end{bmatrix}, \quad b = \begin{bmatrix} \langle f, \varphi_1 \rangle \\ \langle f, \varphi_0 \rangle \end{bmatrix}.$$

Where, for instance,  $\langle \varphi_1, \varphi_1 \rangle = \int x^2 d\mu = \sum x_i^2$ , and  
 $\langle f, \varphi_1 \rangle = \int f x d\mu = \sum x_i y_i$  

## Minimization by derivation

In the real case, with  $f(x) = \|Ax - b\|^2$ , we have

$$\begin{aligned}f(x + h) &= \\ \|A(x + h) - b\|^2 &= \langle Ax - b + Ah, Ax - b + Ah \rangle \\ &= \langle Ax - b, Ax - b \rangle + 2\langle Ax - b, Ah \rangle + \langle Ah, Ah \rangle \\ &= f(x) + \langle 2(A^T Ax - A^T b), h \rangle + hA^T Ah\end{aligned}$$

## Minimization by derivation

In the real case, with  $f(x) = \|Ax - b\|^2$ , we have

$$\begin{aligned}f(x + h) &= \\ \|A(x + h) - b\|^2 &= \langle Ax - b + Ah, Ax - b + Ah \rangle \\ &= \langle Ax - b, Ax - b \rangle + 2\langle Ax - b, Ah \rangle + \langle Ah, Ah \rangle \\ &= f(x) + \langle 2(A^T Ax - A^T b), h \rangle + h A^T A h\end{aligned}$$

and we discover the Taylor series

$$Df(x)[h] = \underbrace{\langle 2(A^T Ax - A^T b), h \rangle}_{\text{gradient}}, \quad D^2f(x)(h, h) = h \underbrace{A^T A}_{\text{Hessian}} h,$$



Sia  $A \in \mathbb{C}^{m \times n}$  con SVD  $A = U \Sigma V^*$

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_l > \sigma_{l+1} = 0 \quad (l \leq p = r, \sigma_{l+1} = 0)$$

$$A = [U_1 | U_2] \left[ \begin{array}{c|c} \sigma_1 & 0 \\ \hline 0 & 0 \end{array} \right] \left[ \begin{array}{c} v_1^* \\ \vdots \\ v_n^* \\ \hline v_2^* \end{array} \right] ]_{n-l} =$$

$$= [U_1 | U_2] \left[ \begin{array}{c} \Sigma v_i^* \\ 0 \end{array} \right] = U_1 \Sigma v_1^* \text{ SVD magne}$$


---

$$A = \sum_{i=1}^{\min\{n, m\}} \sigma_i u_i v_i^* = U \Sigma V^*$$

$$\text{Im}(A) = \text{Span}\{u_1, \dots, u_l\} \quad \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_l > \sigma_{l+1} = 0$$

$$\text{rg } A = l$$

$$\text{Ker } A = \text{Span}\{v_{l+1}, \dots, v_n\}$$

$\Rightarrow$  algoritmi per calcolare queste 3 cose

+ Lezione 21/10

Ripasso di SVD

22/10

$$A \in \mathbb{C}^{m \times n} \quad \exists U \in \mathbb{U}(m), V \in \mathbb{U}(n)$$

$$\Sigma \in \mathbb{R}^{m \times n} \quad \Sigma = \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_n \end{bmatrix}, \quad \sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq \sigma_n \geq 0$$

$$I.c \quad A = U \Sigma V^*$$

E poi c'è il modo magro considerando solo i valori non nulli:

+ SLIDE + REC

$$\sigma_1 = \|A\|_2$$

+ chiedi

Proviamo a risolvere il problema dei minimi quadrati con la SVD:

la SVD

$$\|Ax - b\|_2^2 \stackrel{!}{=} \|U\Sigma V^* x - U\Sigma V^* b\|_2^2 = \|U(\Sigma V^* x - V^* b)\|_2^2 =$$

Possiamo scrivere  $y = V^* x$ , allora,  $m \geq n$

$$\begin{aligned} &= \left\| \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_n \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} U_1^* b \\ U_2^* b \\ \vdots \\ U_m^* b \end{bmatrix} \right\|_2^2 \\ &= \left\| \begin{bmatrix} \sigma_1 y_1 - U_1^* b \\ \sigma_n y_n - U_n^* b \\ \vdots \\ -U_m^* b \end{bmatrix} \right\|_2^2 \\ &= \sum_{i=1}^l (\sigma_i y_i - u_i^* b)^2 + \sum_{i=l}^m |u_i^* b|^2 \end{aligned}$$

Trovare il minimo ns<sub>p</sub> a x equivale a trovare min ns<sub>p</sub> y  
del problema appena trovato

I punti di minimo sono quelli che annullano

$$y_i = \frac{u_i^* b}{\sigma_i} \quad i = 1, \dots, l$$

$$y_i \in \mathbb{C} \quad i = l+1, \dots, n$$

$$x = Vy = [v_1, \dots, v_n] \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \sum y_i v_i = \sum_{i=1}^l \frac{u_i^* b}{\sigma_i} v_i + \sum_{i>l} y_i v_i$$

+ REC

$$\|x\|_2 = \|Vy\|_2 = \|y\|_2 = |y_1|^2 + \dots + |y_l|^2 + |y_{l+1}|^2 + \dots + |y_n|^2$$

e quindi  $y_i = \frac{u_i^* b}{\sigma_i} \quad i = 1, \dots, l$

$$y_i = 0 \quad i > l$$

$$\text{da } x = \sum_{i=1}^l \frac{u_i^* b}{\sigma_i} v_i$$

1. Come si calcola la SVD?
2. Come si calcolano gli autovettori?

$$A \in \mathbb{C}^{m \times n} \text{ è inv se } m=n, \det A \neq 0$$

$$AB = BA = I$$

Def

Data A SVD  $A = U \Sigma V^*$ , la pseudoinversa è

$$A^+ = V \Sigma^+ U^* \text{ dove } \Sigma^+ \in \mathbb{R}^{n \times m} \text{ è } \Sigma^+ = \begin{bmatrix} 1/\sigma_1 & & 0 \\ & \ddots & \\ 0 & & 1/\sigma_n & \dots & 0 \end{bmatrix}$$

$$(\Sigma^+)_{ij} \begin{cases} 1/\delta_i & i=j < l+1 \\ 0 & \text{otherwise} \end{cases}$$

ES

$$V = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} \quad Q = \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix} \quad I \Rightarrow V^+ = \begin{bmatrix} 1/3 \\ 0 \\ 0 \end{bmatrix} Q^*$$

PROPRIETÀ DI  $A^+$ :

- Se  $A$  è invertibile  $\tilde{A} = A^+$

Dim:

$$A = U \Sigma V^* = U \begin{bmatrix} \delta_1 & & 0 \\ & \ddots & \\ 0 & & \delta_n \end{bmatrix} V^*$$

$$\tilde{A} = (V^*)^{-1} \begin{bmatrix} \delta_1 & & 0 \\ & \ddots & \\ 0 & & \delta_n \end{bmatrix} U^{-1} = V \underbrace{\begin{bmatrix} 1/\delta_1 & & \\ & \ddots & \\ & & 1/\delta_n \end{bmatrix}}_{\Sigma^+} U^* = A^+$$

- $AXA = A$ ,  $XAX = X$

- Se  $m=n$  e  $\det A = n \Rightarrow A^+ = (A^* A)^{-1} A^* \quad A^* A = I$

Dim:

$$A = U \begin{bmatrix} \Sigma^+ \\ 0 \end{bmatrix} V^* \quad \Sigma^+ \text{ quadrata e invertibile}$$

$$A^+ = V \begin{bmatrix} \Sigma^+ \\ 0 \end{bmatrix} V^*$$

$$A^* A = V \begin{bmatrix} \Sigma^+ \\ 0 \end{bmatrix} V^* \cancel{U^*} \begin{bmatrix} \Sigma^+ \\ 0 \end{bmatrix} V^* = V \Sigma^2 V^*$$

$$(A^* A)^\dagger = V (\Sigma^2)^\dagger V^*$$

$$(A^* A)^\dagger A^* = V (\Sigma_1^2)^\dagger V^* V [\Sigma_1 \ 0] U^* = V [\Sigma_1 \ 0] U^* = A^+$$

$\|Ax - b\|^2$ , le soluzioni del problema dei minimi quadrati

la minima normina è  $x = A^+ b = + cose$

+ SLIDE + REC

Dove che  $\varrho = A (A^* A)^\dagger A^* =$   
 $\underbrace{A^+}$

$$= U \begin{bmatrix} \Sigma^2 \\ 0 \end{bmatrix} V^* V \begin{bmatrix} \Sigma^2 \\ 0 \end{bmatrix} U^* = U \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix} U^*$$

$$\Rightarrow \text{Span} \{u_1, \dots, u_n\} = V$$

Def:

$$z = \rho e^{i\theta}, z \neq 0$$

↑ introduce unitaria

introduce def pos.

def positivo

$A \in \mathbb{C}^{n \times n}$   $A = U \tilde{U}^\dagger \rightarrow$  decomposizione polare  
 $\tilde{U}$  unitaria

+ SLIDE

$$\text{Se } A = \tilde{U} \tilde{\Sigma} \tilde{V}^* =_{\text{p}} U \tilde{\Sigma} V^*$$

$$A = \underbrace{\tilde{U} \tilde{V}^*}_{\tilde{U}} \left( \underbrace{\tilde{V} \Sigma \tilde{V}^*}_{H} \right)$$

+ SLIDES

28/10

## Trasformata discreta di Fourier

$$\omega_n = e^{i \frac{2\pi}{n}} \quad z^n = 1 \quad e^{i \frac{2\pi}{n} j} \quad j = 0, \dots, n-1$$

È radice principale dell'unità

$$\left( e^{i \frac{2\pi}{n}} \right)^l = e^{i \frac{2\pi}{n} l}$$

$$G = \left\{ e^{i \frac{2\pi}{n} j}, \quad j=0, \dots, n-1 \right\}$$

$$1, \omega_n, \omega_n^2, \dots, \omega_n^{n-1}$$

Una radice dell'unità  
è detta positiva se le  
sue potenze generano  
tutte le radici

+ SLIDES

29/10

SLIDES su FFT

12/11

$$F_n \times (F_n)_{i,j} = \omega_n^{ij} \quad i, j = 0, \dots, n-1$$

$$\omega_n = e^{-i \frac{2\pi}{n}}$$

IDFT(x) valutazione  $f(z) = x_0 + x_1 z + \dots + x_{n-1} z^{n-1}$

$$\tilde{F}_n^{-1}(x) = \frac{1}{n} F_n^* x$$

$$FFT \quad \left( \frac{3}{2} n^2 \log_2 n \right)$$

+ slide ( $n \leq 4$  del gruppo FFT)

$$a(x) = \sum_{i=0}^{n-1} a_i x^i, \quad b(x) = \sum_{i=0}^{n-1} b_i x^i$$

$$c(x) = \sum_{i=0}^{2n-2} c_i x^i = ab$$

Algoritmo standard lo interpretiamo come un prodotto  
matrice vettore

$$\begin{bmatrix} c_0 \\ c_1 \\ \vdots \\ c_{2n-2} \end{bmatrix} = \begin{pmatrix} a_0 & & & & \\ a_1 & a_0 & & & \\ \vdots & \ddots & \ddots & & \\ & \ddots & \ddots & a_0 & \\ & & \ddots & \ddots & a_1 \\ & & & \ddots & \ddots \\ & & & & \ddots & a_{n-1} \\ & & & & & \ddots & a_n \\ & & & & & & \ddots \\ & & & & & & & a_{n-1} \end{pmatrix} \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_{n-1} \end{bmatrix}$$

$$n = 2^p \quad a(x) = \sum_{i=0}^{2^p-1} a_i x^i \quad (+\text{slide})$$

$$2^{p-1} < 2n \leq 2^p$$

$$[\log_2 n] = p, \quad a_n = \dots = a_{2^p-1} = 0$$

Prodotto tra polinomi tramite FFT

$$\text{Se } p + c \quad 2^p \geq 2n - 2 \quad e \quad 2^{p-1} < 2n - 2$$

Svolgendo i  $n$  complessi  $\omega_n^i = a(\omega_n^i) \quad i = 0, \dots, m-1$

$$\text{con } m = 2^p$$

$$\beta_i = b(\omega_n^i) \quad i = 0, \dots, m-1$$

Si calcola poi  $y_i = \alpha_i \beta_i$   $i=0, \dots, m-1$   
 Il prodotto di  $n$  complessi

Si haora il polinomio  $f(x) = p(\omega_n^i) = y_i$   $i=0, \dots, m-1$

$$\Rightarrow p = a \cdot b$$

$p$  è il polinomio d'interpolazione sulle radici dell'unità  
 e valo  $y_i$ , ma

$$ab(\omega_n^i) = a(\omega_n^i)b(\omega_n^i) = \alpha_i \beta_i = y_i$$

minore e uguale  
 a  $m$

$\Rightarrow ab = p$  poiché il polinomio di interpolazione è unico

Perciò:

- 2IDFT
  - $n$  prodotti di numeri complessi  $\rightarrow m$  ops
  - DFT
- $\frac{9}{2} m \log_2 m$  ops

$$\Rightarrow O(m \log_2 m)$$

Quel è il max  $n$  di cifre in base  $B$  della mantissa  
 di un coefficiente di  $\tilde{c}$ ? ( $\leq (B-1)^2 N$ )

$$\Rightarrow \log_B ((B-1)^2 N) =$$

$$= \lg_B N + 2 \lg_B (B-1) = c \lg_2 N + \gamma = \Theta(\lg_2 N)$$

Ho due mob per vedere le foto sulle

$$\frac{n!}{n(n-1)!} \xrightarrow{n \rightarrow \infty} 1$$

$$F_2 x = \begin{bmatrix} F_{\frac{1}{2}} x_p + D_{\frac{1}{2}} T_{\frac{1}{2}} x_s \\ F_{\frac{1}{2}} x_p - D_{\frac{1}{2}} T_{\frac{1}{2}} x_s \end{bmatrix} = \begin{bmatrix} F_{\frac{1}{2}} \\ F_{\frac{1}{2}} \end{bmatrix} + \begin{bmatrix} D_{\frac{1}{2}} T_{\frac{1}{2}} & 0 \\ 0 & D_{\frac{1}{2}} T_{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} x_p \\ x_s \end{bmatrix}$$

$$= \begin{bmatrix} I & D_{\frac{n}{2}} \\ I & -D_{\frac{n}{2}} \end{bmatrix} \begin{bmatrix} F_{n/2} & 0 \\ F_{n/2} & 0 \end{bmatrix} \begin{bmatrix} X_F \\ X_D \end{bmatrix} =$$

$$= \begin{bmatrix} I & D_{n_2} \\ I & -D_{n_2} \\ I & -D_{n_2} \end{bmatrix} \begin{bmatrix} I & D_{n_2} \\ I & -D_{n_2} \\ I & -D_{n_2} \end{bmatrix}^T = \begin{bmatrix} O & F_{n_1 n_2} \\ F_{n_1 n_2}^T & F_{n_1 n_2} F_{n_1 n_2}^T \\ O & F_{n_1 n_2} F_{n_1 n_2}^T \end{bmatrix} = \begin{bmatrix} x_{i=0}(\zeta) \\ x_{i=2}(\zeta) \\ x_{i=1}(\zeta) \\ x_{i=3}(\zeta) \end{bmatrix}$$

el 1 pass della FFT si deve p

1 blocchi si devono  $2^{f-l}$

Dimesione dei blocchi è  $2^l$

Dimesione di  $D_q$  è  $2^{f-l-1}$

→ Implementazione con ciclo for.

For  $l=2:p$  % Calcolo  $x = \Delta_l x$

For  $j=0: 2^{f-l}$  % siamo i blocchi

For  $q=0: 2^{l-1}$

$$\gamma = \omega * x (q + 2^{l-1})$$

$$x(q + 2^{l-1}) = x(i) - \gamma$$

$$x(q) = x(i) + \gamma$$

$$Ax = b \quad A \in \mathbb{C}^{n \times n} \quad \det A \neq 0 \quad b \neq 0$$

(A di grande dimensione e spesso)

Metodi dei sottospazi di Krylov:

$$\begin{aligned} k, L \subset \mathbb{C}^n, \dim k = \dim L = n \\ \text{approssimazione dello spazio} \\ \text{del sistema lineare} \\ x_m = x_0 + k \\ b - Ax_n \perp L \end{aligned}$$

} metodo di  
proiezione

$$1. \quad k = L = k_n(A, r_0) \quad r(x) = b - Ax$$

FOR

(full orthogonalization )  
method

$$r_0 = b - Ax_0$$

$$r_n = b - Ax_n$$

CG (con iugate gradient se A è def positiva)

$$k_n(Av) = \text{Span} \{ v, Av, A^2v, \dots, A^{m-1}v \}$$

$$2. \quad k = k_n(\overset{\text{invertibile}}{A}, r_0), \quad L = AR \quad \text{GMRES}$$

$$3. \quad k = k_n(A, r_0), \quad k_n(A^\top, b) \quad \text{BICG-STAB}$$

Se  $v \neq 0$

e  $d = n$  immagine  $\left[ \begin{array}{c} d < n \\ \text{dimostra per } m \in \mathbb{C}^{n \times n} \\ m \geq d \end{array} \right]$

$$k_1(Av) \subset k_2(Av) \subset \dots \subset k_d(Av) = \dots = k_m(Av)$$

ove  $d$  è il grado di  $v$  rispetto ad  $A$ , grado minimo di un polinomio  $h$  c.  $p(A)v = 0$

$$K_m(A, v) = \{ p(A)v, \deg p < m \}$$

$$\tilde{A} = p(A), \quad x_* = \tilde{A}b = p(A)b \in K_m(A, b) \quad (\text{per } m \leq d)$$

$$x_* = \tilde{A}b = q(A)b \quad \text{con } \deg q < n$$

La sol di 1. si puo' scrivere esplicitamente come:

$$V_m \in \mathbb{C}^{m \times n}, \quad \text{Span } V_m = K_m$$

$$K_m = K_m(A, r_0), \quad \text{space } M = \text{Span}\{\text{colonne}\}$$

$$x_n = x_0 + V_m (V_m^* A V_m)^{-1} V_m^* r_0$$

Le cose da fare ora sono

- trovare una base
- e' vero che  $V_m^* A V_m$  e' invertibile?
- e' necessaria l'invertibile  $V_m^* A V_m$ ?
- Se  $m$  e' piccolo, quanto e' buona l'approssimazione?
- Come scegliere  $m$ ?

Se  $m \leq d$   $V, AV, \dots, A^{m-1}V$  e' una base, questa base per e' mal condizionata ( $\|M\|_2, \|M'\|_2$  grande)

Infatti se ad esempio  $A$  e' diagonale allora sono

$$A = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_n \end{bmatrix}, \quad \text{se c'e' un autovalore piu' grande}$$

degli altri quando  $\mu_2(M)$  per le molte volte  
tendono ad allontanarsi dal prodotto dell'autovettore  
massimo con le rette

$\Rightarrow$  Ci serve un'altra base (ben diversa da  $\{v\}$ )  
con  $\mu_2(M) \approx \|M\|_2 \|\tilde{M}\|_2$  piccolo)

$\Rightarrow$  Ci servono le medie: esistono per cui sappiamo  
che  $\mu_2(M) = 1 \Leftrightarrow M \in N(0)$

Trovare una base ortogonale dello  $\mathcal{K}_n(A, v)$

Usiamo allora i metodi di Arnoldi:  
(Pertanto  $A$  e  $v$  e costretti  $v_1, \dots, v_m$  base ortogonale di  $\mathcal{K}_n(A, v)$ )

$$(0) \quad v_1 = \frac{v}{\|v\|}$$

$$(1) \quad Av_1 = u_1, \quad \tilde{u}_1 = u_1 - \underbrace{\langle u_1, v_1 \rangle v_1}_{h_{11}} \quad (\tilde{u}_1 \perp v_1)$$

$$v_2 = \frac{\tilde{u}_1}{\|\tilde{u}_1\|} \quad (\tilde{u}_1 \neq 0) \quad \text{se } h_{11} \neq 0$$

$$(2) \quad Av_2 = u_2, \quad \tilde{u}_2 = u_2 - \underbrace{\langle u_2, v_1 \rangle v_1}_{h_{21}} - \underbrace{\langle u_2, v_2 \rangle v_2}_{h_{22}}$$

$$v_3 = \frac{\tilde{u}_2}{\|\tilde{u}_2\|} \quad h_{21} \quad h_{22}$$

$$(l) \quad \tilde{v}_{m+1} = \tilde{u}_{m+1}, \quad \tilde{u}_{m+1} = u_{m+1} - \sum_{j \leq m} \langle u_{m+1}, u_j \rangle v_j$$

$\underbrace{h_{m+1}}_{\|u_{m+1}\|}$

ottimizzazione

metodo matrice vettore

$\sim O(mn) \text{ ops} + mn$

La base che ho ottenuto è base di  $K_n$ ?

Teorema:

Se il metodo non si amesta fino al passo l'induso allora  $\text{span}\{v_1, \dots, v_l\} = K_l(A, v)$ ,  $i \leq l+1$

Se il metodo si amesta al passo l allora

$\text{span}\{v_1, \dots, v_{l-1}\} = K_m(A, v) = K_l(A, v) \quad m \geq l$  e  
 $A^l v \in K_l(A, v)$  (LUCKY BREAK DOWN)

→ perchē nel caso di amestō va tutto bene (è la svariazione migliore)

Teorema:

se il metodo  $\mathcal{H}$  di Arnoldi per  $K_{m+1}(A, v)$  non si amesta fino al passo m e sia

$$V_l = [v_1 | \dots | v_l] \in \mathbb{C}^{l \times n} \quad \text{dove } v_1, \dots, v_{m+1} \text{ sono i vettori}$$

ottenuti dal metodo

$$\bar{H}_m = \begin{bmatrix} h_{11} & \cdots & h_{1m} \\ h_{21} & \ddots & \vdots \\ \vdots & \ddots & \vdots \\ h_{m+1,1} & \cdots & h_{mm} \end{bmatrix} \quad \begin{array}{l} (m+1) \times m \\ \in \mathbb{C} \text{ con gli } h_{ij} \text{ sono i} \\ \text{coeff del metodo di Arnoldi} \end{array}$$

$H_m \in \mathbb{C}^{m \times m}$  la matrice ottimale selezionando le  
minime righe di  $\bar{H}_m$  allora si ha

$$2. AV_m = V_{m+1} \bar{H}_m$$

$$2. AV_m = V_m H_m + h_{m+1,m} V_{m+1} J_m^T$$

$$3. V_m^* AV_m = H_m$$

Dimostrazione

$$\textcircled{1} \quad 1 \leq l \leq m$$

$$\begin{aligned} AV_l &= u_l = \tilde{u}_l + \sum_{j \leq l} h_{jl} v_j = h_{l+1,l} v_{l+1} + \sum_{j \leq l} h_{jl} v_j = \\ &= \sum_{j \leq l+1} h_{jl} v_j = \begin{bmatrix} v_1 | \dots | v_{m+1} \end{bmatrix} \begin{bmatrix} h_{1l} \\ \vdots \\ h_{l+1,l} \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \\ &= V_m \begin{bmatrix} h_{1,l} \\ \vdots \\ h_{l+1,l} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \end{aligned}$$

$$AV_m = A \begin{bmatrix} v_1 | v_2 | \dots | v_{m+1} \end{bmatrix} = \begin{bmatrix} AV_1 | AV_2 | \dots | AV_{m+1} \end{bmatrix} =$$

ove

$$AV_1 = V_{m+1} \begin{bmatrix} h_{11} \\ h_{21} \\ \vdots \\ 0 \end{bmatrix}, \quad AV_2 = V_{m+1} \begin{bmatrix} h_{12} \\ h_{22} \\ h_{32} \\ ? \\ \vdots \\ 0 \end{bmatrix}, \quad \dots$$

$$\Rightarrow AV_m = \begin{bmatrix} h_{11} & h_{12} & \dots & h_{1m+1} \\ h_{21} & h_{22} & \dots & h_{2m+1} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & h_{m+1,m+1} \end{bmatrix} = V_{m+1} \bar{H}_m$$

(2)

$$\left[ \begin{array}{c|c} V_m & V_{m+1} \\ \hline \underbrace{\qquad}_{m \times m} & \underbrace{(m+1) \times 1} \end{array} \right] \left[ \begin{array}{c} H_m \\ \hline \underbrace{0 \cdots 0}_{m+1, m} \end{array} \right] =$$

$\xrightarrow{m \times m}$

$\xrightarrow{2 \times (m+1)}$

$$= V_m H_m + V_{m+1} [0, \dots, 0] = V_m H_m + h_{m+1,m} V_{m+1} f_m^T$$

(3)

$$V_m^* A V_m = \underbrace{V_m^* V_m H_m}_{\substack{\parallel \\ \text{rechte Seite}}} + \underbrace{V_m^* h_{m+1,m} V_{m+1} f_m^T}_{\substack{\parallel \\ \text{rechte Seite}}} = H_m$$

$\xrightarrow{0}$

$\xrightarrow{V_m^* V_{m+1} f_m^T}$

Vediamo i metodi:

FOM (se  $A$  è def positiva ha CG)

$$K = L = K_m(A, r_0)$$

$$\left\{ \begin{array}{l} X_m = X_0 + V_m Y_m \\ Y_m = (V_m^* A V_m)^{-1} V_m^* r_0 \end{array} \right.$$

sol del  
problema

$$\left\{ \begin{array}{l} x_n = x_0 + k \\ b - Ax_0 \perp L \end{array} \right.$$

Se  $V_m$  è la matrice le cui colonne sono le base ottenuta tramite il metodo I Arnoldi

$$-V_m^* A V_m = H_m$$

$v_1 = \frac{r_0}{\|r_0\|_2}$  sol del problema in  $K_m(A, r_0)$

$$-V_m^* r_0 = \underbrace{V_m^* V_1 \|r_0\|_2}_L = \|r_0\|_2 f_1$$

$(10. -0)$

$\approx f_1$

$\Rightarrow$  Primo  $\Rightarrow$  calcolo  $V_m$  con Arnoldi, poi  $\Rightarrow$  valore  
 Restante linea  $H_m Y_m = \|\tau_{\text{roll}}\|_2 f_1$ , poi l'approssimazione  
 ne è  $X_m = X_0 + V_m Y_m$

Come lo scegli per questo  $m$ ?

Se  $A$  è def positiva  $\Rightarrow H_m$  è simmetrica  $\Rightarrow$  è tridiagonale

$$\|X_m - X^*\|_A = \min_{X \in X_0 + k} \|X - X_m\|_A$$

si accelera  
 Arnoldi e  
 ottiene lanczos

① Posa  $X_m$  come a priori l'errore

② Oppure lo trasforma in un metodo iterativo con  
 si rinnova a posteriori dell'errore

Teorema (per FOM)

$$\|b - Ax_m\|_2 = h_{m+1,m} \left\| \int_m^\top Y_m \right\| \quad (\text{succ } \rightarrow m \rightarrow n)$$

Dunque

$$\begin{aligned}
 b - Ax_m &= b - Ax_0 - AV_m Y_m = r_0 - AV_m Y_m = \\
 &= r_0 - V_m H_m Y_m - V_m h_{m+1,m} \underbrace{V_{m+1} \int_m^\top Y_m}_m = \\
 &= r_0 - V_m \|\tau_{\text{roll}}\|_2 f_1 - V_m V_{m+1} \int_m^\top h_{m+1,m} Y_m = \\
 &= \underbrace{r_0 - V_m \|\tau_{\text{roll}}\|_2}_{\tau_0} - V_m V_{m+1} \int_m^\top h_{m+1,m} Y_m =
 \end{aligned}$$

$$\Rightarrow \|b - Ax_m\|_2 = \|Y_m v_{m+1} f_m^T h_{m+1,m}\| = h_{m+1,m} |f_m^T Y_m|$$

Terremo. \*

Se  $A$  è definitiva ( $x_* = \tilde{A}^{-1}b$ ,  $b \neq 0$ ),  $m=0, 1, \dots, l$   
dunque  $\Delta(A, b)$

$$\|x_m - x_*\|_A \leq 2 \left( \frac{\sqrt{\mu} - 1}{\sqrt{\mu} + 1} \right)^n \|x_m - x_0\|_A$$

con  $\mu = n$  e condizionalità,  $\mu = \|A\|_2 \|\tilde{A}^{-1}\|_2$

26/11

Coste computazionale del metodo di Arnoldi

$$\tilde{u}_i = u_i - \sum_{j \leq i} \langle u_i, v_j \rangle v_j = 4n$$

↑  
n              ↑  
2n              n

$$\frac{\tilde{u}_i}{\|\tilde{u}_i\|} \rightarrow 2n+1 \text{ sqrt ops}$$

$$\Rightarrow \sum_{i=1}^n \sum_{j \leq i} 4n \sim \frac{4n^2 m^2}{2} \sim 2nm^2$$

$\Rightarrow$  Arnoldi "costa"  $2nm^2$  ops se  $A$  è spesso

(CAUSS)

Per risolvere il sistema lineare pesso avere  $\frac{2}{3}n^3$  ops

Per il passo 3 ( $x_n = x_0 + v_m y_m$ ) ha  $2nm$

Si implementa così:

For  $i = 1, 2, \dots$

% calcolo  $v_i$  con Arnoldi e calcolo  $y_i, x_i$

if ( $\|b - Ax\|_m < tol$ )

stop

else

gradiente  
conjugato

Se  $A$  definita positiva FOM si chiama CG e si semplifica di molto

$$H_m = V_m^* A V_m, H_m^* = V_m^* A^* (V_m^*)^* = V_m^* A V_m = H_m$$

$\Rightarrow H_m$  hermitiano + hessenberg  $\Rightarrow$  tridiagonale

metodo di Lanczos

$$\tilde{u}_i = u_i - \langle u_i, v_i \rangle v_i - \langle u_i, v_{i-1} \rangle v_{i-1}$$

$$\begin{bmatrix} \cdot & & \\ & \ddots & 0 \\ \cdot & 0 & \ddots \end{bmatrix}$$

$\nearrow$  Il metodo di Arnoldi si semplifica di molto

$$x_{m+1} = x_m + d_m d_m$$

Gli autovettori di  $H_m$ , i valori di Ritz,  
approssimano gli autovettori di  $A$

$$2. \begin{cases} x_m = x_0 + k_m \\ b - Ax_m \perp Ak_m \end{cases}$$

A invertibile

Metodo GMRES

(generalized minimal residual)

Si dimostra che  $b - Ax_m + Akm$  equivale a  
dove  $\|b - Ax_m\|_2 = \min_{x \in x_0 + km} \|b - Ax\|_2$

Calcolo  $V_m$  con Arnoldi senza troncatura

$$x_m = x_0 + V_m y_m \quad r_0 = V_0 \|r_0\| = \|r_0 - V_{m+1} f_m\|_2 \Leftrightarrow V_m = \frac{r_0}{\|r_0\|}$$

$$b - Ax_m = b - Ax_0 - A V_m y_m = r_0 - V_{m+1} \bar{H}_{m+1} y_m \quad \downarrow$$

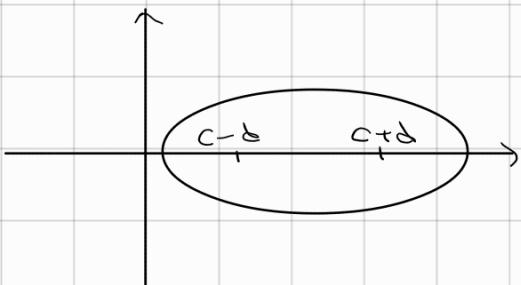
$$= V_{m+1} (\|r_0\| f_m - \bar{H}_m y_m)$$

$$\|b - Ax_m\|_2 = \|V_{m+1} (\|r_0\| f_m - \bar{H}_m y_m)\|_2$$

$x_m$  che minimizza il residuo è del tipo  $x_0 + V_m y_m$

dove  $y_m$  minimizza  $\|r_0\| f_m - \bar{H}_m y_m\|_2$

$\Rightarrow y_m$  è la soluzione del problema dei minimi quadrati con dati  $\bar{H}_m$  e  $\|r_0\| f_m$



Se gli autovalori di  $A$  si trovano nell'ellisse di centro  $c$ , fuochi  $c-d$  e  $c+d$  e semiasse maggiore  $a$

$$\text{e } A \text{ è diagonalizzabile } X^T A X = D = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_n \end{bmatrix}$$

$$\|r_m\|_2 \leq \mu_2(X) \frac{T_m(\frac{a}{d})}{|T_m(\frac{c}{d})|} \|r_0\|$$

Tm per il Chebyshev non è una conseguenza

# Problemi di ottimizzazione

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  continua

trovare min  $f(x) \in \arg \min_{x \in \mathbb{R}^n} f(x)$

ottimizzazione  
non vincolata

$f: X \rightarrow \mathbb{R}$        $X \subseteq \mathbb{R}^n$

min  $f(x)$   
 $x \in X$

ottimizzazione  
vincolata

OSS

trovare min e max e' indifferente poiche'

$$\max_{x \in X} f = -\min_{x \in X} (-f)$$

e inoltre  $\min_{x \in X} (f(x) + c) = (\min_{x \in X} f(x)) + c$

$x^*$  min globale per  $f$  in  $X$  se  $f(x) \geq f(x^*) \quad x \in X$   
 " locale "  $f(x) \geq f(x^*)$  in  $B(x^*, r)$

ed o' sette se vale la maggiore delle sette.

L'ottimizzazione e' divisa in:

- Analisi ( $\exists$  del minimo etc. - )

- Metodi ottimizzazione

Esempi:

- $f(x,y) = x^2 + y^2$   $x^2 + y^2 \geq 0$   $x^2 + y^2 = 0 \Leftrightarrow (x,y) = (0,0)$

$\Rightarrow \min_{(x,y) \in \mathbb{R}^2} f(x,y) = 0$  e argmini  $f(x,y) = \{(0,0)\}$

Nessun ha massimo

- $f(x,y) = -x^2 - y^2$  ha massimo assoluto e nessun minimo

- $f(x,y) = x^2 - y^2$  non ha né max né min globali  
(local max less per uno)

- $\left\{ \begin{array}{l} x_0 \in \mathbb{R}^{n_0} \\ y_i = A_{i-1} x_{i-1} - b_{i-1}, \quad i = 1, \dots, l \\ x_i = \phi_i(y_i) \end{array} \right.$  noti numeri definiti  
in modo iterativo

ove  $A_i \in \mathbb{R}^{n_i \times n_{i-1}}$ ,  $b_i \in \mathbb{R}^{n_i}$  (pesi)

$$\phi: \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{n_i}, (\phi_i(x_i))_i = \varphi_i(x_i), \varphi_i: \mathbb{R} \rightarrow \mathbb{R}$$

$$x_0 \xrightarrow{f} x_l$$

$\uparrow$   
funzione di  
attivazione  
(deve essere monotonica)

Ottimizzazione non vincolata:

$f \in C^1(\mathbb{R}^n)$ ,  $x_0 \in \mathbb{R}^n$  esiste  $L(x_0) : \mathbb{R}^n \rightarrow \mathbb{R}$

$$f(x_0 + h) = f(x_0) + L(x_0)[h] + o(h) \quad h \in \mathbb{R}^n$$

$h \rightarrow 0$

$L(x_0)$  se esiste  $\Rightarrow$  derivata di Frechet  
in  $x_0$  ( $Df(x_0)$ )

Se  $Df(x_0)$  esiste  $\forall x_0 \in \mathbb{R}^n \Rightarrow f$  è derivabile in  $\mathbb{R}^n$

$$\begin{aligned} Df : \mathbb{R}^n &\longrightarrow \text{Hom}(\mathbb{R}^n, \mathbb{R}) \\ x_0 &\mapsto Df(x_0) \end{aligned}$$

$$\lim_{\substack{\text{per} \\ h \rightarrow 0}} \frac{|f(x_0 + h) - f(x_0) - L(x_0 + h)|}{\|h\|} = 0$$

## Proposizioni

Per ogni  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R} \quad \exists \mathbf{v} \in \mathbb{R}^n$  t.c.  $\mathbf{g}(w) = \langle \mathbf{v}, w \rangle$   
 $\forall w \in \mathbb{R}^n$

OSS

$f$  derivabile  $\overset{x_0}{\underset{\mathbb{R}^n}{f}} : \mathbb{R}^n \rightarrow \mathbb{R} \quad Df(x_0) : \mathbb{R}^n \rightarrow \mathbb{R}$

$$Df(x_0)[h] = \langle \nabla f(x_0), h \rangle, \quad \nabla f(x_0) \in \mathbb{R}^n$$

$A \in \mathbb{R}^{1 \times n}$  associato a  $Df(x_0)$

$$A = \left[ \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right] = \nabla f(x_0)^T$$

$Df(x_0)$  è detto derivata nella direzione  $h$

||

$$\nabla f(x_0)^T h = \lim_{t \rightarrow 0} \frac{f(x_0 + th) - f(x_0)}{t} = \frac{\partial f}{\partial h}$$

Tegoreo (cond. necessarie del I ordine per l'ottimale)

Se  $f \in C^1(\mathbb{R}^n)$  e  $x^*$  è minima  $\Rightarrow \nabla f(x^*) = 0$

Dmo:

Se  $f(x) \geq f(x^*)$  in  $\in B(x^*, r)$

$$\Rightarrow f(x^* + th) \geq f(x^*)$$

$\Rightarrow$  per ogni  $h \in \mathbb{R}^n$   $x^* + th \in B(x^*, r)$ ,  $t \in [0, t_0]$

$$\lim_{t \rightarrow 0} \frac{f(x^* + th) - f(x^*)}{t} = \nabla f(x^*)^T h \geq 0$$

ma è anche  $\leq 0$

$$\Rightarrow \nabla f(x^*)^T h = 0 \quad \forall h \in \mathbb{R}^n \Rightarrow \langle \nabla f(x^*), h \rangle = 0$$

$$\Rightarrow \nabla f(x^*) = 0$$

$f(x) \in C^2(\Omega)$

02/12  
03/12

$\nabla f(x^*) = 0 \quad \nabla^2 f(x) \geq 0 \quad \Leftarrow x^* \text{ minimum}$

Termin:

$\nabla f(x^*) = 0, \nabla^2 f(x) > 0 \Rightarrow x^* \text{ e- minimum}$

Durch:  
 $h \in S^{n-1}$

$$f(x^* + th) - f(x^*) = \nabla f(x^*)[th] + \frac{1}{2} \nabla^2 f(x^*)[th, th] + o(\|th\|^2)$$

$$\lim_{t \rightarrow 0} \frac{f(x^* + th) - f(x^*)}{t^2} = \frac{1}{2} h^T \nabla^2 f(x^*) h \geq \frac{1}{2} \varepsilon \|h\|^2 = \frac{1}{2} \varepsilon$$

$$\exists t_0(h) \in [0, t_0[h]], f(x^* + th) > f(x^*)$$

$(t, h) \rightarrow f(x^* + th) - f(x^*)$  es iste  $M(h)$  innerhalb von  $S^{n-1}$  stetig

$(t, h) \in [0, t_s] \times M(h), t_s > 0$

$$S^{n-1} = \bigcup_{h \in S^{n-1}} M(h) = \bigcup_{s=1}^n \overbrace{M(h_s)}^{t_s}$$

$$h^T A h \geq \varepsilon h^T h = \varepsilon \|h\|^2$$

$t \in [0, \max_{s=1, \dots, n} t_s] \Rightarrow f(x^* + th) > f(x^*), B(x^*, \max t_s) \setminus \{x^*\}$

Esempio:

$A \in \mathbb{R}^{n \times m}$  simmetrica,  $b \in \mathbb{R}^n$

Consideriamo  $\varphi(x) = \frac{1}{2} x^T A x - b^T x$ ,  $\nabla \varphi(x) = Ax - b$

$$\begin{aligned}\varphi(x+h) &= \frac{1}{2} (x+h)^T A (x+h) - b^T (x+h) - \frac{1}{2} x^T A x - b^T x = \\ &= (x^T A - b^T) h + \frac{1}{2} h^T A h\end{aligned}$$

I punti stazionari sono i punti in cui il gradiente è 0

→ I punti stazionari sono le soluzioni di  $Ax - b = 0$

Se  $A \in GL(n) \Rightarrow \exists 1.$  punto stazionario

Se  $A \succ 0 \Rightarrow \exists$  minima globale

Se  $A$  ha un autovalore negativo  $\Rightarrow A \not\succ 0 \Rightarrow$  non è una minima

$$\varphi(x) = \|Ax - b\|_2^2 = (Ax - b)^T (Ax - b) = x^T A^T A x - 2x^T A^T b - b^T b$$

$$2A^T A x - 2A^T b = 0$$

Proprietà del gradiente:

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  allora  $f(x) = b^T x$ ,  $b \in \mathbb{R}^n$

Se  $b = 0$   $f(x) \equiv 0$

Se  $b \neq 0$   $f(x) = 0 \Leftrightarrow x \perp b \Rightarrow \text{ker } f = \{b^\perp\} = \text{span}\{b\}^\perp$   
ha dimensione  $n-1$

$\text{ker } f$  divide  $\mathbb{R}^n$  in due parti in cui  $f(x)$  ha segno opposto

$f \in C^1(\mathbb{R})$  nell'intorno di  $x_0 \Rightarrow$  compatta come

$$f(x) - f(x_0) \approx f(x_0) + \nabla f(x_0)^T \underbrace{(x - x_0)}_h, \quad x \in M(x_0)$$

Teorema:

Sia  $f \in C^1(\mathbb{R}^n)$  essa  $x_0 \in \mathbb{R}^n$  t.c.  $\nabla f(x_0) \neq 0$   
 allora se  $h \in \mathbb{R}^n$  t.c.  $\nabla f(x_0)^T h < 0$ , esiste  $t_0 > 0$  t.c.  
 $\exists t \in [0, t_0] \quad f(x_0 + th) < f(x_0)$

Dimo:

$$f(x_0 + th) - f(x_0) = \underbrace{\langle \nabla f(x_0), th \rangle}_{\hookrightarrow = t \nabla f(x_0)^T h} + o(t) \quad \text{plicc } f \in C^1(\mathbb{R}^n)$$

$$\lim_{t \rightarrow 0} \frac{f(x_0 + th) - f(x_0)}{t} = \nabla f(x_0)^T h < 0 \Rightarrow \exists t_0 > 0 \text{ t.c.} \\ \text{per } t \in [0, t_0], \quad f(x_0 + th) < f(x_0)$$

Metodi di ottimizzazione:

- Valore iniziale  $x_0 \in \mathbb{R}^n$
- iterazione: sceglie la direzione  $d_l \in \mathbb{R}^n$   
 $\text{e possa } \lambda_l \in \mathbb{R}^+$

Se  $\nabla f(x_l) = 0 \Rightarrow$  min ferme

Se  $\nabla f(x_l) \neq 0$ :

$$x_{l+1} = x_l + \lambda_l d_l$$

Dal questo metodo non spesso converge a un minimo o a un massimo ( $\nabla f(x_e) = 0$  potrebbe essere solo zero).  
Spesso solo che converga a un punto stazionario, anziché che sia vera ~~successione~~ a convergere a un punto stazionario.

Sceglieremo di tale che  $\nabla f(x_e)^T d \leq 0$   
⇒ direzione di discesa

Quale è la direzione di massima discesa?

- Una scelta potrebbe essere  $-\nabla f(x_e)$
- Un'altra potrebbe essere  $-A_e \nabla f(x_e)$  con  $A_e$  definita positiva ( $-\nabla f(x_e)^T A \nabla f(x_e) \leq 0$ )

(Questi metodi si chiamano metodi del gradiente)

Metodo di Newton:

I CASO

$$\underbrace{\nabla f(x) = 0}_{\text{Sistema non lineare}}$$

$$x_{e+1} = x_e - \frac{\hat{f}(x_e)}{\hat{f}'(x_e)}$$

$\hat{f}: \mathbb{R}^n \rightarrow \mathbb{R}^n$   
 $\hat{f}(x) = \circ$

$x \in \mathbb{R}^n$  Immagine

$\nabla f(x) \in \mathbb{R}^n$  equazioni

$$x_{e+1} = x_e - \hat{J}(x_e)^{-1} [\hat{f}(x_e)] =$$

$$= x_e - \hat{J}'(x_e) \hat{f}(x_e)$$

↑ Jacobiano

II CASO

Se  $\hat{f} = \nabla f$ ,  $x_{e+1} = x_e - H(x_e)^{-1} \nabla f(x_e)$  metodo di Newton

oppure considero  $x_{e+1} = x_e + \alpha_e \underbrace{(-H(x_e)^{-1} \nabla f(x_e))}_{\text{de}}$

$H(x_e)$  def pos  $\Rightarrow H(x_e)^{-1}$  e def pos

Metodi che usano solo derivate prim'ordine del I ordine  
 Metodi che usano anche derivate seconda solo del II ordine

Fissato  $\alpha_e$ , cerchiamo  $\lambda_e$

Siamo nel caso in cui siamo d'ov'e' il minimo  
 restrizione di  $f$  allo retta  $x_e + \lambda \alpha_e$ ,  $\lambda \in \mathbb{R}$

$\varphi(\lambda) = f(x_e + \lambda \alpha_e)$ , cerca argmin  $f(x_e + \lambda \alpha_e) \geq \alpha_e$   
 exact line search

e un'altro di questi è  
 metodo del gradiente  
 +  
 exact line search  $\Rightarrow$  Steepest descent

Idea alternativa: inexact line search  
 cioè scelgo  $\lambda_e$  in modo che garantisco che sia  
 facile da trovare e che garantisca "decrescita  
 sufficiente" per la convergenza globale.

Un esempio è:

• Regola di ARMJO (backtracking)

È possibile,  $\beta \in [0, 1]$  parametri che variano  
 ad esempio  $\beta = 1$  & quanto più meno

$\sigma \in [0, 1[$

Al posso il scegliere le più piccole intere non  
 negative ma tale che

La decrescita è utile  
 per garantire la  
 convergenza  
 perché

$$f(x_e + \bar{\lambda} \beta^m \alpha_e) \leq f(x_e) + \bar{\lambda} \beta^m \nabla f(x_e)^T \alpha_e$$

$$\frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{x}^T \mathbf{b}$$

- Calcolo del gradiente (subito corrente)

↳ differenziazione numerica

↳ differenziazione analitica

- Valutare il gradiente

(Nelle reti neurali si nasce e calcola e valutare)

$$\begin{cases} \mathbf{x}_0 \\ \mathbf{y}_i = \mathbf{A}_i \mathbf{x}_{i-1} + \mathbf{b}_i, \quad i = 1, \dots, l \\ \mathbf{x}_i = \sigma(\mathbf{y}_i) \end{cases}$$

↳ σ non lineare  
però

$$\mathbf{x}_l = f(\mathbf{x}_0; \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_l, \mathbf{b}_1, \dots, \mathbf{b}_l)$$

LOSS FUNCTION

$$\tilde{\mathcal{L}} = \sum_{i=1}^N \| f(\tilde{\mathbf{x}}_i; \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_l, \mathbf{b}_1, \dots, \mathbf{b}_l) - \tilde{\mathbf{y}}_i \|^2$$

Trovare  $\mathbf{A}_1, \dots, \mathbf{A}_l, \mathbf{b}_1, \dots, \mathbf{b}_l$  tali che  $\tilde{\mathcal{L}}$  sia minima

$$f_1(x) = \|x\|^2 = \mathbf{x}^T \mathbf{x}, \quad \nabla f_1(x)[h] = 2 \mathbf{x}^T h$$

$$f_2(x) = \sigma(x), \quad f_2(x+h) - f_2(x) = \left[ \begin{array}{c} \sigma(x_1+h_1) - \sigma(x_1) \\ \sigma(x_2+h_2) - \sigma(x_2) \\ \vdots \\ \sigma(x_n+h_n) - \sigma(x_n) \end{array} \right] =$$

$$= \left[ \begin{array}{c} \sigma'(x_1) h_1 \\ \vdots \\ \sigma'(x_n) h_n \end{array} \right] + O(\|h\|)$$

$$\Rightarrow Df_2(x)[h] = D\sigma(x)[h] = \sigma'(x) \odot h$$

metodo punto per punto

$$f_3(A) = Ax + b, \quad Df_3(A)[H] = Hx$$

$\mathbb{R}^{n \times n} \rightarrow \mathbb{R}^n$

$$f_4(b) = Ax + b \quad Df_4(b)[h] = h$$

$$f(A, b) = Ax - b, \quad Df(A, b)[H, h] = Hx + h =$$

$$= [x^T \otimes I \quad I] \begin{bmatrix} \text{vec}(H) \\ h \end{bmatrix}$$

09/12

$$D(f \circ g)(x)[h] = Df(g(x)) [Dg(x)[h]]$$

$$\text{Quindi } \varphi \in C^1 \Rightarrow \varphi'(\alpha) = Df(x + \alpha d)[d] = \nabla f(x + \alpha d)^T d$$

$$\Rightarrow \varphi'(0) = \nabla f(x)^T d$$

$$\|d\| = 1, \quad \nabla f(x)^T d = (\cos \theta) \| \nabla f(x) \| \quad \begin{pmatrix} \text{e max per } \theta = 0 \\ \text{e min per } \theta = \pi \end{pmatrix}$$

NOTAZIONE:

$$d_e = -\nabla f(x_e) = -g(x_e)$$

$$d_e = A_e g(x_e), \quad A_e \text{ def positiva}$$

$$d_e = -g(\hat{x}_e) + \beta_e g(x_{e-}) \quad \leftarrow \text{In questa categoria c'è il metodo del CG (metodo simile e re mancatum)}$$

Si dimostra che applicando CG alla funzione quadratica

$$p(x) = \frac{1}{2} x^T A x - x^T b, \quad A > 0, \text{ allora ottengo il FOM.}$$

Poss scegliere  $\lambda_e$  attraverso:

- PASSO FISSO: quando  $\beta I \geq H(x) \geq \alpha I$  ci sono dei passi fissi che garantiscono la convergenza globale

- ARMIJO'S RULE:

$$\bar{\lambda}, \beta \in [0, 1], \sigma \in ]0, 1[$$

ad ogni passo  $\lambda_e = \bar{\lambda} \beta^{m_e}$  dove  $m_e \in \mathbb{N}$  è il più piccolo intero non negativo tale che

$$f(x_e + \lambda_e d_e) \leq f(x_e) + \sigma g(x_e)^T d_e, \quad \lambda_e > \lambda_e \quad (\text{sufficiente decrescita})$$

Lemma:

Sia  $\gamma: I \rightarrow \mathbb{R}$ ,  $I$  intorno di 0 in  $\mathbb{R}$ , tale che

$\gamma(t) = \gamma_0 + \psi(t)$  con  $\lim_{t \rightarrow 0} \psi(t) = 0$ , allora per ogni

$0 < \gamma < 1$  esiste  $t_0 > 0$  t.c.  $t \in [0, t_0] \Rightarrow \gamma(t) > \sigma \gamma_0$ .

[ $\square$ ]

Se  $\gamma_0 > 0$   
Se  $\gamma_0 < 0$

Teorema:

Dato  $f \in C^1(\mathbb{R}^n)$ ,  $x \in \mathbb{R}^n$  e  $d \in \mathbb{R}^n$  direzione di decrescita per  $f$  in  $x$ ,  $0 < \sigma < 1$ , esiste  $t_0 > 0$  t.c.  $t \in [0, t_0]$

$$f(x + t d) \leq f(x) + \sigma \nabla f(x)^T d t \quad (\nabla f(x) \neq 0)$$

Dur:

$$f(x+td) = f(x) + t \nabla f(x)^T d + \underbrace{\psi(t)t}_{\text{resto dove } \psi(t) \rightarrow 0}$$

per le Lemma  $\underbrace{\nabla f(x)^T d + \psi(t)}_{< 0} \leq \sigma \nabla f(x)^T d$  per  $t \in [0, t_0]$

Tesima:

$\nabla f(x)$  non è mai zero

Sia  $\{x_e\}_e$  una successione infinita ottenuta dal metodo

$x_0 \in \mathbb{R}^n$ ,  $x_{e+1} = x_e + \lambda_e d_e$  per ottimizzazione  $f \in C^2(\mathbb{R}^n)$ , dove la scelta di  $d_e$  è ottenuta tramite la regola di Armijo con parametri  $\bar{\alpha} > 0$ ,  $\beta \in ]0, 1[$ ,  $\sigma \in ]0, 1[$ .

Se  $\{d_e\}_e$  è gradient-related allora ogni punto limite è stazionario

Def

$\{d_e\}_e$  successione di direzioni è gradient-related se per ogni sottosequenza di  $\{x_e\}_e$  che converge a un punto non stazionario

$$\liminf_{e \in \mathcal{L}_1} g(x_e), d_e < 0$$

$\{\|d_e\|\}_{e \in \mathcal{L}_1}$  è limitata,  $\mathcal{L}_1 \subset \{0, 1, \dots\}$  sottosequenze infinite

Dur:

$$\mathcal{L}_1 \subseteq \mathcal{L}_0 = \{0, 1, 2, \dots\}$$

P.A

- 1) Supponiamo che esiste  $\{x_e\}_{e \in \mathcal{L}_1}$  t.c  $\lim_{e \in \mathcal{L}_1} x_e = x^*$   $g(x^*) \neq 0$
- 2) Dimostriamo che  $f(x_{e+1}) - f(x) \rightarrow 0$

Supponiamo  $f$  è continua allora  $\lim_{e \in \mathcal{L}_1} f(x_e) = f(\lim_{e \in \mathcal{L}_1} x_e) = f(x^*)$

$$\{f(x_e)\}_{e \in \mathcal{L}_1}$$

3)  
<)

5)  $\tilde{d}_e = \frac{d_e}{\|d_e\|}$ ,  $\tilde{d}_e = \frac{d_e}{\|d_e\|}$ ,  $d_e \neq 0$  perché  $e \in \{\|d_e\|\}$  lembite ...  
gradient related  
ed è una direzione  
 $\rightarrow$  bisca a

E' stata una sottosequenza  $\tilde{d}_e$  che converge

6)  $\exists \theta_e \in [0, 1]$   
 $f(x_e + \tilde{d}_e \tilde{d}_e) - f(x_e) = Df(x_e + \theta_e \tilde{d}_e)(\tilde{d}_e) =$

Tangente al valore medio

$$\langle g(x^*), d^* \rangle \geq \theta_e \langle g(x^*), d^* \rangle$$

$$A(x_{e+1} - x_e) = g(x_{e+1}) - g(x_e)$$

n° incognite, n° equazioni  $\Rightarrow$  sottosistema nato

Ogni soluzione da cui mettere fatto gressi-Newton

Reti neurali

$$d_{e+1} = -A g_e$$

$$\begin{cases} x_0 \in \mathbb{R}^n \\ y_e = A_e x_{e-1} + b_e & e = 1, \dots, L \\ x_e = \sigma(y_e) \\ f(x_0) = y = x_L \end{cases}$$

$$\begin{matrix} x_0^{(1)}, \dots, x_0^{(n)} \\ y_0^{(1)}, \dots, y_0^{(n)} \end{matrix}$$

$$\mathcal{L} = \sum_{i=1}^n \|y_i^{(i)} - f(x_i^{(i)}, A_s, b_s, A_e, b_e)\|^2$$

$$\mathcal{L}(A_s, b_s, A_e, b_e) : (\mathbb{R}^{n_0 \times n_s} \times \mathbb{R}^{n_s} \times \mathbb{R}^{n_s \times n_e} \times \mathbb{R}^{n_e}) \rightarrow \mathbb{R}$$

$$x_3 = \sigma(y_3) = \sigma(A_3 x_2 + b_3) = \sigma(A_3 \sigma(y_2) + b_3) =$$

$$= \sigma(A_3 + \sigma(A_2 x_1 + b_2) + b_3) = \sigma(A_3 + \sigma(A_2 \sigma(y_1) + b_2) + b_3) =$$

$$= \sigma(A_3 + \sigma(A_2 + \sigma(A_1 x_0 + b_1) + b_2) + b_3)$$

$$DF(H, h) = D\sigma(y_3)[Hx_2 + h] = \sigma'(y_3) \odot [x_2^\top \otimes I \cdot I] \begin{bmatrix} \text{vec}(H) \\ h \end{bmatrix}$$

back propagation

costo computazionale : n di OPS (non conta + delle valutazioni delle reti) neuroni)