

# Il metodo del Gradiente Coniugato

Gerardo Toraldo

Università di Napoli Federico II

A.A. 2014-2015

## Minimizzazione di una funzione quadratica convessa

$$\begin{aligned} \min f(x) &= \frac{1}{2}x^T Ax - b^T x \quad \text{s.t.} \quad x \in \mathbb{R}^n \\ \text{con } A &\in \mathbb{R}^{n \times n} \text{ simmetrica definita positiva e } b \in \mathbb{R}^n \end{aligned} \quad (1)$$

Risolvere il problema (1) è equivalente a risolvere uno dei problemi seguenti:

$$Ax = b,$$

$$\min_{x \in \mathbb{R}} \|Ax - b\|.$$

## Metodo Steepest Descent (SD)

$x_0 \in \mathbb{R}^n; k \leftarrow 0;$

$g_0 \leftarrow Ax_0 - b;$

**while** *not\_stopcondition*

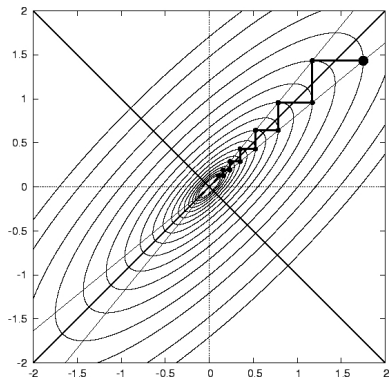
$$\alpha_k = \frac{g_k^T g_k}{g_k^T A g_k};$$

$$x_{k+1} \leftarrow x_k - \alpha_k g_k;$$

$$g_{k+1} \leftarrow g_k - \alpha_k A g_k;$$

$$k \leftarrow k + 1;$$

**endwhile**



Notiamo che  $x_{k+1} = x_k + v^*$  dove

$$v^* = \operatorname{argmin} f(x_k + v) \text{ con } v \in \operatorname{Span}\{g_k\}$$

## Algoritmo di discesa con "passo ottimo"

$x_0 \in \mathbb{R}^n$ ;  $g_0 \leftarrow Ax_0 - b$ ;  $k \leftarrow 0$ ;

$g_0 \leftarrow Ax_0 - b$ ;

**while** *not\_stopcondition*

scegli  $v_k$ ;

$$\alpha_k = -\frac{v_k^T g_k}{v_k^T A v_k};$$

$$x_{k+1} \leftarrow x_k + \alpha_k v_k;$$

$$g_{k+1} \leftarrow g_k + \alpha_k A v_k;$$

$$k \leftarrow k + 1;$$

**endwhile**

Notiamo che

$$x_{k+1} = x_k + v^*, \text{ con}$$

$$v^* = \operatorname{argmin}_{v \in \operatorname{Span}\{v_k\}} f(x_k + v)$$

$$\text{e } v \in \operatorname{Span}\{v_k\}$$

Ad ogni passo la funzione obiettivo  $f(x)$  viene **minimizzata in uno spazio affine unidimensionale (o, equivalentemente,  $f(x_0 + v)$  viene minimizzata in uno spazio vettoriale unidimensionale)**

E' possibile minimizzare  $f(x_0 + v)$  in **sottospazi di dimensione via via crescente?**

$$x_1 = \operatorname{argmin}_{x = x_0 + v, v \in \operatorname{Span}\{v_0\}} f(x)$$

$$x_2 = \operatorname{argmin}_{x = x_0 + v, v \in \operatorname{Span}\{v_0, v_1\}} f(x)$$

.....

$$x_{k+1} = \operatorname{argmin}_{x = x_0 + v, v \in \operatorname{Span}\{v_0, v_1, \dots, v_k\}} f(x)$$

## Algoritmo di discesa con "passo ottimo"

$x_0 \in \mathbb{R}^n$ ;  $g_0 \leftarrow Ax_0 - b$ ;  $k \leftarrow 0$ ;

$g_0 \leftarrow Ax_0 - b$ ;

**while** *not\_stopcondition*

scegli  $v_k$ ;

$$\alpha_k = -\frac{v_k^T g_k}{v_k^T A v_k};$$

$$x_{k+1} \leftarrow x_k + \alpha_k v_k;$$

$$g_{k+1} \leftarrow g_k + \alpha_k A v_k;$$

$$k \leftarrow k + 1;$$

**endwhile**

Notiamo che

$$x_{k+1} = x_k + v^*, \text{ con}$$

$$v^* = \operatorname{argmin}_v f(x_k + v)$$

$$\text{e } v \in \operatorname{Span}\{v_k\}$$

Ad ogni passo la funzione obiettivo  $f(x)$  viene **minimizzata in uno spazio affine unidimensionale (o, equivalentemente,  $f(x_0 + v)$  viene minimizzata in uno spazio vettoriale unidimensionale)**

E' possibile minimizzare  $f(x_0 + v)$  in **sottospazi di dimensione via via crescente?**

$$x_1 = \operatorname{argmin}_x f(x) \text{ con } x = x_0 + v, v \in \operatorname{Span}\{v_0\}$$

$$x_2 = \operatorname{argmin}_x f(x) \text{ con } x = x_0 + v, v \in \operatorname{Span}\{v_0, v_1\}$$

.....

$$x_{k+1} = \operatorname{argmin}_x f(x) \text{ con } x = x_0 + v, v \in \operatorname{Span}\{v_0, v_1, \dots, v_k\}$$

## Minimizzazione in sottospazi

### Il problema

$$\min f(x) = \frac{1}{2}x^T Ax - b^T x \quad \text{s.t.} \quad x = x_0 + v \quad \text{con} \quad v \in S = \text{Span}\{v_0, v_1, \dots, v_k\}$$

può essere formulato come segue:

$$\begin{aligned} \min f(x_0 + \alpha_0 v_0 + \alpha_1 v_1 + \dots + \alpha_k v_k) &\Leftrightarrow \\ \min f(x_0 + V\alpha) \quad \text{con} \quad \alpha \in \mathbb{R}^{k+1}, \quad V = [v_0, v_1, \dots, v_k] \in \mathbb{R}^{n \times k+1} \end{aligned}$$

Quindi

$$\begin{aligned} \bar{\alpha} = \operatorname{argmin} f(x_0 + V\alpha) &\Leftrightarrow \nabla_{\alpha} f(x_0 + V\bar{\alpha}) = 0 \\ &\Leftrightarrow V^T \nabla f(x_0 + V\bar{\alpha}) = 0 \end{aligned}$$

### Condizioni di ortogonalità

$$v_i^T \nabla f(x_0 + V\bar{\alpha}) = 0, \quad i = 0, 1, \dots, k \quad \Leftrightarrow \quad \nabla f(x_0 + V\bar{\alpha}) \in S^{\perp} = \text{Span}\{v_0, v_1, \dots, v_k\}^{\perp}$$

## Minimizzazione in sottospazi

### Il problema

$$\min f(x) = \frac{1}{2}x^T Ax - b^T x \quad \text{s.t.} \quad x = x_0 + v \quad \text{con} \quad v \in S = \text{Span}\{v_0, v_1, \dots, v_k\}$$

può essere formulato come segue:

$$\begin{aligned} \min f(x_0 + \alpha_0 v_0 + \alpha_1 v_1 + \dots + \alpha_k v_k) &\Leftrightarrow \\ \min f(x_0 + V\alpha) \quad \text{con} \quad \alpha \in \mathbb{R}^{k+1}, \quad V = [v_0, v_1, \dots, v_k] \in \mathbb{R}^{n \times k+1} \end{aligned}$$

Quindi

$$\begin{aligned} \bar{\alpha} = \operatorname{argmin} f(x_0 + V\alpha) &\Leftrightarrow \nabla_{\alpha} f(x_0 + V\bar{\alpha}) = 0 \\ &\Leftrightarrow V^T \nabla f(x_0 + V\bar{\alpha}) = 0 \end{aligned}$$

### Condizioni di ortogonalità

$$v_i^T \nabla f(x_0 + V\bar{\alpha}) = 0, \quad i = 0, 1, \dots, k \quad \Leftrightarrow \quad \nabla f(x_0 + V\bar{\alpha}) \in S^{\perp} = \text{Span}\{v_0, v_1, \dots, v_k\}^{\perp}$$

## Minimizzazione in un sottospazio bidimensionale: perché scegliere direzioni coniugate?

Supponiamo di aver eseguito le prime due iterazioni:

$$x_1 = x_0 + \alpha_0 v_0; \quad x_2 = x_1 + \alpha_1 v_1.$$

La scelta del passo ottimo garantisce che

$$\nabla f(x_1)^T v_0 = 0 \quad \text{e} \quad \nabla f(x_2)^T v_1 = 0.$$

Richiedere che  $x_2$  **minimizzi la funzione in  $\text{Span}\{v_0, v_1\}$**  significa richiedere che valga anche la condizione di ortogonalità

$$\nabla f(x_2)^T v_0 = 0.$$

Ricordando che  $\nabla f(x_2) = \nabla f(x_1) + \alpha_1 A v_1$ , abbiamo

$$\nabla f(x_2)^T v_0 = \nabla f(x_1)^T v_0 + \alpha_1 (A v_1)^T v_0 = \alpha_1 v_0^T A v_1$$

e, dunque,  $x_2$  minimizza la funzione nello spazio affine

$$A = \{x : x = x_0 + \alpha_0 v_0 + \alpha_1 v_1, \quad \alpha_0, \alpha_1 \in \mathbb{R}\}$$

se e solo se le direzioni di ricerca  $v_0$  e  $v_1$  sono  **$A$ -coniugate**, ovvero  $v_0^T A v_1 = 0$ .



## Minimizzazione in un sottospazio bidimensionale: perché scegliere direzioni coniugate?

Supponiamo di aver eseguito le prime due iterazioni:

$$x_1 = x_0 + \alpha_0 v_0; \quad x_2 = x_1 + \alpha_1 v_1.$$

La scelta del passo ottimo garantisce che

$$\nabla f(x_1)^T v_0 = 0 \quad \text{e} \quad \nabla f(x_2)^T v_1 = 0.$$

Richiedere che  $x_2$  **minimizzi la funzione in  $\text{Span}\{v_0, v_1\}$**  significa richiedere che valga anche la condizione di ortogonalità

$$\nabla f(x_2)^T v_0 = 0.$$

Ricordando che  $\nabla f(x_2) = \nabla f(x_1) + \alpha_1 A v_1$ , abbiamo

$$\nabla f(x_2)^T v_0 = \nabla f(x_1)^T v_0 + \alpha_1 (A v_1)^T v_0 = \alpha_1 v_0^T A v_1$$

e, dunque,  $x_2$  minimizza la funzione nello spazio affine

$$A = \{x : x = x_0 + \alpha_0 v_0 + \alpha_1 v_1, \quad \alpha_0, \alpha_1 \in \mathbb{R}\}$$

se e solo se **le direzioni di ricerca  $v_0$  e  $v_1$  sono  $A$ -coniugate, ovvero  $v_0^T A v_1 = 0$ .**

## Il metodo delle direzioni coniugate

### Proposizione [Coniugatezza e indipendenza lineare]

Se i vettori  $v_1, v_2, \dots, v_k$  sono **A-coniugati**, allora essi sono anche **linearmente indipendenti**

**Metodo delle direzioni coniugate:** metodo di discesa con "passo ottimo" in cui le direzioni di ricerca sono mutuamente coniugate.

- :-( Necessità di determinare delle direzioni coniugate
- :-) Il metodo termina dopo un numero finito di passi (al più  $n$ )

## Metodo delle direzioni coniugate

$v_0, v_1, \dots, v_{n-1}$  vettori non nulli  $A$ -coniugati

$x_0 \in \mathbb{R}^n$ ;  $g_0 \leftarrow Ax_0 - b$ ;  $k \leftarrow 0$ ;

**while**  $k < n$  **and**  $\|g_k\| > tol$

$$\alpha_k = -\frac{v_k^T g_k}{v_k^T A v_k};$$

$$x_{k+1} \leftarrow x_k + \alpha_k v_k;$$

$$g_{k+1} \leftarrow g_k + \alpha_k A v_k;$$

$$k \leftarrow k + 1;$$

**endwhile**

Come scegliere le direzioni  $v_k$ ?

## Metodo del Gradiente Coniugato

$x_0 \in \mathbb{R}^n$ ;  $g_0 \leftarrow Ax_0 - b$ ;  $v_0 = -g_0$ ;  $k \leftarrow 0$ ;

**while**  $k < n$  **and**  $\|g_k\| > tol$

$$\alpha_k = -\frac{v_k^T g_k}{v_k^T A v_k};$$

$$x_{k+1} \leftarrow x_k + \alpha_k v_k;$$

$$g_{k+1} \leftarrow g_k + \alpha_k A v_k;$$

$$\beta_k \leftarrow ?;$$

$$v_{k+1} = -g_{k+1} + \beta_k v_k;$$

$$k \leftarrow k + 1;$$

**endwhile**

La direzione di ricerca all'iterazione  $k$  è costruita come combinazione lineare dell'antigradiente in tale iterazione e della direzione di ricerca all'iterazione precedente.

## Il Metodo del gradiente coniugato: proprietà

Siano  $v_0, v_1, \dots, v_k$  le direzioni  $A$ -coniugate generate dall'algoritmo CG e siano  $x_i = x_0 + v^*$ , con  $v^* = \operatorname{argmin} f(x_0 + v)$ ,  $v \in \operatorname{Span}\{v_0, v_1, \dots, v_{i-1}\}$ ,  $i \leq k+1$ , le corrispondenti approssimazioni della soluzione. Quale deve essere il valore di  $\beta_k$  affinché la direzione

$$v_{k+1} = -g_{k+1} + \beta_k v_k$$

sia  $A$ -coniugata alle precedenti?

Essendo  $g_i = \beta_{i-1} v_{i-1} - v_i$ , si ha

$$\operatorname{Span}\{g_0, g_1, \dots, g_i\} = \operatorname{Span}\{v_0, v_1, \dots, v_i\} \quad \forall i \leq k$$

e, ricordando la formula ricorsiva che lega i gradienti, si ha anche

$$\operatorname{Span}\{v_0, v_1, \dots, v_i\} = \operatorname{Span}\{g_0, g_1, \dots, g_i\} = \operatorname{Span}\{g_0, Ag_0, \dots, A^i g_0\} \quad \forall i \leq k.$$

Infine, in virtù della coniugatezza delle direzioni e della ottimalità del passo, risulta

$$g_i \perp v_j, \quad j = 1, \dots, i-1 \quad \text{e quindi} \quad g_i \perp g_j, \quad j = 0, \dots, i-1.$$

## Il Metodo del gradiente coniugato: proprietà

Siano  $v_0, v_1, \dots, v_k$  le direzioni  $A$ -coniugate generate dall'algoritmo CG e siano

$x_i = x_0 + v^*$ , con  $v^* = \operatorname{argmin} f(x_0 + v)$ ,  $v \in \operatorname{Span}\{v_0, v_1, \dots, v_{i-1}\}$ ,  $i \leq k+1$ ,  
le corrispondenti approssimazioni della soluzione. Quale deve essere il valore di  $\beta_k$   
affinché la direzione

$$v_{k+1} = -g_{k+1} + \beta_k v_k$$

sia  $A$ -coniugata alle precedenti?

Essendo  $g_i = \beta_{i-1}v_{i-1} - v_i$ , si ha

$$\operatorname{Span}\{g_0, g_1, \dots, g_i\} = \operatorname{Span}\{v_0, v_1, \dots, v_i\} \quad \forall i \leq k$$

e, ricordando la formula ricorsiva che lega i gradienti, si ha anche

$$\operatorname{Span}\{v_0, v_1, \dots, v_i\} = \operatorname{Span}\{g_0, g_1, \dots, g_i\} = \operatorname{Span}\{g_0, Ag_0, \dots, A^i g_0\} \quad \forall i \leq k.$$

Infine, in virtù della coniugatezza delle direzioni e della ottimalità del passo, risulta

$$g_i \perp v_j, \quad j = 1, \dots, i-1 \quad \text{e quindi} \quad g_i \perp g_j, \quad j = 0, \dots, i-1.$$

## Il metodo del Gradiente Coniugato: proprietà

Da  $v_{k+1} = -g_{k+1} + \beta_k v_k$  discende

$$v_{k+1}^T A v_i = -g_{k+1}^T A v_i + \beta_k v_k^T A v_i, \quad i = 1, \dots, k.$$

Inoltre, per  $i \leq k$ ,

$$0 = g_{k+1}^T g_i = g_{k+1}^T (g_{i-1} + \alpha_{i-1} A v_{i-1}) \Rightarrow g_{k+1}^T A v_{i-1} = 0$$

Sfruttando la coniugatezza di  $v_0, v_1, \dots, v_k$  si ha quindi

$$v_{k+1}^T A v_i = 0, \quad i = 0, 1, \dots, k-1, \quad (2)$$

$$v_{k+1}^T A v_k = -g_{k+1}^T A v_k + \beta_k v_k^T A v_k. \quad (3)$$

Da (2) segue che, per qualsiasi scelta di  $\beta_k$ , la direzione  $v_{k+1}$  è coniugata a  $v_0, v_1, \dots, v_{k-1}$ , mentre da (3) segue che affinché  $v_{k+1}$  sia coniugata anche a  $v_k$  deve essere

$$\beta_k = \frac{g_{k+1}^T A v_k}{v_k^T A v_k}.$$

figuracg.jpg



## Il metodo del Gradiente Coniugato: proprietà (riepilogo)

Abbiamo dimostrato che il metodo del Gradiente Coniugato

- genera una sequenza di direzioni coniugate

$$v_k = -g_k + \beta_{k-1}v_{k-1},$$

$$\text{con } \beta_k = (g_{k+1}^T A v_k) / (v_k^T A v_k);$$

- costruisce una sequenza di punti

$$x_k = x_0 + v^*, \text{ con } v^* = \operatorname{argmin} f(x_0 + v), v \in \operatorname{Span}\{v_0, v_1, \dots, v_{k-1}\};$$

- converge in al più  $n$  passi;
- è tale che

$$\operatorname{Span}\{v_0, v_1, \dots, v_k\} = \operatorname{Span}\{g_0, g_1, \dots, g_k\} = \operatorname{Span}\{g_0, A g_0, \dots, A^k g_0\}$$

## Il metodo del Gradiente Coniugato: espressioni di $\alpha_k$ e $\beta_k$

- Da  $v_k = -g_k + \beta_{k-1}v_{k-1}$  e  $v_{k-1}^T g_k = 0$  discende

$$\alpha_k = -\frac{v_k^T g_k}{v_k^T A v_k} = \frac{g_k^T g_k}{v_k^T A v_k}.$$

- Da  $A v_k = (g_{k+1} - g_k)/\alpha_k$  discende

$$\beta_k = \frac{g_{k+1}^T A v_k}{v_k^T A v_k} = \frac{g_{k+1}^T g_{k+1}}{g_k^T g_k}.$$

## Metodo del Gradiente Coniugato

$x_0 \in \mathbb{R}^n$ ;  $g_0 \leftarrow Ax_0 - b$ ;  $v_0 = -g_0$ ;  $k \leftarrow 0$ ;

**while**  $k < n$  **and**  $\|g_k\| > tol$

$$\alpha_k \leftarrow \frac{g_k^T g_k}{v_k^T A v_k} \quad \boxed{-\frac{v_k^T g_k}{v_k^T A v_k}};$$

$$x_{k+1} \leftarrow x_k + \alpha_k v_k;$$

$$g_{k+1} \leftarrow g_k + \alpha_k A v_k;$$

$$\beta_k \leftarrow \frac{g_{k+1}^T g_{k+1}}{g_k^T g_k} \quad \boxed{\frac{g_{k+1}^T A v_k}{v_k^T A v_k}};$$

$$v_{k+1} \leftarrow -g_{k+1} + \beta_k v_k;$$

$$k \leftarrow k + 1;$$

**endwhile**

### Complessità computazionale per iterazione

- Un solo prodotto matrice vettore per iterazione ( $Av_k$ );
- Il prodotto  $Ax_k$  non viene mai eseguito;
- memoria di lavoro: 4 vettori ( $x$ ,  $v$ ,  $Av$ , e  $r$ )

## Il metodo del Gradiente Coniugato: convergenza

### Convergenza e autovalori distinti

Se  $m$  è il numero di autovalori distinti di  $A$ , allora il metodo CG converge alla soluzione in un numero di passi  $k \leq m$ .

**Dimostrazione.** Siano  $\lambda_1, \lambda_2, \dots, \lambda_m$  gli autovalori distinti di  $A = UDU^T$ . Supponiamo per assurdo che i vettori  $g_0, g_1, \dots, g_m$  siano tutti non nulli. Osserviamo inoltre che

$$\prod_{i=1}^m (A - \lambda_i I) = U \left( \prod_{i=1}^m (D - \lambda_i I) \right) U^T = 0.$$

Allora le matrici

$$I = A^0, A, A^2, \dots, A^m$$

sono linearmente dipendenti, e tali sono i vettori

$$g_0, Ag_0, A^2 g_0, \dots, A^m g_0.$$

Ma

$$\text{Span}\{g_0, Ag_0, A^2 g_0, \dots, A^m g_0\} = \text{Span}\{g_0, g_1, g_2, \dots, g_m\},$$

e ciò è in contrasto l'ipotesi di assurdo.

## Il metodo del Gradiente Coniugato: convergenza

### Convergenza e autovalori distinti

Se  $m$  è il numero di autovalori distinti di  $A$ , allora il metodo CG converge alla soluzione in un numero di passi  $k \leq m$ .

**Dimostrazione.** Siano  $\lambda_1, \lambda_2, \dots, \lambda_m$  gli autovalori distinti di  $A = UDU^T$ . Supponiamo per assurdo che i vettori  $g_0, g_1, \dots, g_m$  siano tutti non nulli. Osserviamo inoltre che

$$\prod_{i=1}^m (A - \lambda_i I) = U \left( \prod_{i=1}^m (D - \lambda_i I) \right) U^T = 0.$$

Allora le matrici

$$I = A^0, A, A^2, \dots, A^m$$

sono linearmente dipendenti, e tali sono i vettori

$$g_0, Ag_0, A^2g_0, \dots, A^mg_0.$$

Ma

$$\text{Span}\{g_0, Ag_0, A^2g_0, \dots, A^mg_0\} = \text{Span}\{g_0, g_1, g_2, \dots, g_m\},$$

e ciò è in contrasto l'ipotesi di assurdo.

## il Metodo del gradiente coniugato: errore assoluto

Definiamo errore assoluto al passo  $k$  la quantità

$$e_k = x_k - x^*.$$

Essendo  $Ax^* = b$ , si ha

$$e_k^T A v_i = (Ax_k - Ax^*)^T v_i = g_k^T v_i = 0, \quad i = 0, 1, \dots, k-1$$

( $e_k$  è  $A$ -ortogonale (coniugato) alle prime  $k$  direzioni di ricerca)

Inoltre, essendo  $Ae_0 = g_0$ , si ha

$$\text{Span}\{g_0, g_1, \dots, g_k\} = \text{Span}\{v_0, v_1, \dots, v_k\} = \text{Span}\{Ae_0, A^2e_0, \dots, A^{k+1}e_0\}$$

## Il metodo del gradiente coniugato: errore assoluto

Si osservi che  $x_{k+1} \in X_k$ , con

$$X_k = \{x : x = x_0 + v\}, \quad v \in \text{Span}\{v_0, v_1, \dots, v_k\},$$

e quindi  $e_{k+1} \in E_k$ , con

$$E_k = \{e : e = e_0 + v\}, \quad v \in \text{Span}\{v_0, v_1, \dots, v_k\} = \text{Span}\{Ae_0, A^2e_0, \dots, A^{k+1}e_0\}.$$

Esistono quindi dei coefficienti reali  $\gamma_1, \gamma_2, \dots, \gamma_{k+1}$  tali che

$$\begin{aligned} e_{k+1} &= e_0 + \sum_{i=1}^{k+1} \gamma_i A^i e_0 \\ &= \left( 1 + \sum_{i=1}^{k+1} \gamma_i A^i \right) e_0 = p_{k+1}(A) e_0, \end{aligned}$$

dove  $p_{k+1} \in \Pi_{k+1}^1$  e  $\Pi_{k+1}^1$  è l'insieme dei polinomi di grado  $k+1$  che valgono 1 nell'origine.

## Il metodo del gradiente coniugato: errore assoluto

### Caratterizzazione dell'errore assoluto

L'errore assoluto  $e_{k+1}$  risolve il problema

$$\begin{aligned} \min \quad & u^T A u \\ \text{s.t.} \quad & u \in U = \{u : u = e_0 + s, s \in S\} \\ \text{con } S = & \text{Span}\{v_0, v_1, \dots, v_k\} = \text{Span}\{Ae_0, A^2e_0, \dots, A^{k+1}e_0\} \end{aligned} \quad (4)$$

**Dimostrazione.** Il problema (4) può essere riformulato come

$$\begin{aligned} \min \quad & (e_0 + V\alpha)^T A(e_0 + V\alpha) \\ \text{con } V = & [v_0, v_1, \dots, v_k] \text{ e } \alpha \in \mathbb{R}^{k+1} \end{aligned}$$

o, equivalentemente,

$$\min \alpha^T V^T A V \alpha + 2e_0^T A V \alpha$$

(minimizzazione di una funzione convessa).



## Il metodo del gradiente coniugato: errore assoluto

### Caratterizzazione dell'errore assoluto

L'errore assoluto  $e_{k+1}$  risolve il problema

$$\begin{aligned} \min \quad & u^T A u \\ \text{s.t.} \quad & u \in U = \{u : u = e_0 + s, s \in S\} \\ \text{con } S = & \text{Span}\{v_0, v_1, \dots, v_k\} = \text{Span}\{Ae_0, A^2e_0, \dots, A^{k+1}e_0\} \end{aligned} \tag{4}$$

**Dimostrazione.** Il problema (4) può essere riformulato come

$$\begin{aligned} \min \quad & (e_0 + V\alpha)^T A(e_0 + V\alpha) \\ \text{con } V = & [v_0, v_1, \dots, v_k] \text{ e } \alpha \in \mathbb{R}^{k+1} \end{aligned}$$

o, equivalentemente,

$$\min \alpha^T V^T A V \alpha + 2e_0^T A V \alpha$$

(minimizzazione di una funzione convessa).

## Il metodo del gradiente coniugato: errore assoluto

Scrivendo le condizioni di ottimo si ha:

$$\begin{aligned} V^T A V \alpha + V^T A e_0 &= 0 & \Leftrightarrow & V^T A (V \alpha + e_0) = 0 & \Leftrightarrow \\ V^T A u &= 0 & \Leftrightarrow & v_i^T A u = 0, \quad i = 0, 1, \dots, k \end{aligned}$$

L'ultima relazione è proprio la condizione di  $A$ -ortogonalità tra  $e_{k+1}$  e le direzioni di ricerca  $v_0, \dots, v_k$  e quindi vale la tesi.

La proposizione precedente può essere anche formulata come segue:

### Caratterizzazione dell'errore assoluto

$$\begin{aligned} \|e_{k+1}\|_A &= \min_{u \in U} \|u\|_A = \min_{\gamma \in \mathbb{R}^{k+1}} \|e_0 + \sum_{i=1}^{k+1} \gamma_i A^i e_0\| \\ &= \min \|p(A)e_0\|_A, \quad p \in \Pi_{k+1}^1. \end{aligned}$$

## Il metodo del gradiente coniugato: errore e spettro di $A$

Si ha  $(u_1, u_2, \dots, u_n)$  base ortonormale di autovettori di  $A$

$$\begin{aligned} \|p(A)(x_0 - x^*)\|_A^2 &= (x_0 - x^*)^T A (p(A))^2 (x_0 - x^*) = \\ &= \left( \sum_1^n \gamma_i u_i \right)^T A (p(A))^2 \left( \sum_1^n \gamma_i u_i \right) = \left( \sum_1^n \gamma_i u_i \right)^T \sum_1^n \gamma_i A (p(A))^2 u_i = \\ &= \left( \sum_1^n \gamma_i u_i \right)^T \sum_1^n \gamma_i A (p(\lambda_i))^2 u_i = \left( \sum_1^n \gamma_i u_i \right)^T \sum_1^n \gamma_i \lambda_i (p(\lambda_i))^2 u_i = \\ &= \sum_1^n \gamma_i^2 \lambda_i (p(\lambda_i))^2 \leq \max_i (p(\lambda_i))^2 \sum_1^n \gamma_i^2 \lambda_i = \max_i (p(\lambda_i))^2 \|x_0 - x^*\|_A^2 \end{aligned}$$

Stima dell'errore relativo

$$\frac{\|e_{k+1}\|_A^2}{\|e_0\|_A^2} \leq \min_{p \in \Pi_{k+1}^1} \max_{1 \leq i \leq n} (p(\lambda_i))^2$$

## Il metodo del gradiente coniugato: errore e spettro di $A$

Si ha  $(u_1, u_2, \dots, u_n)$  base ortonormale di autovettori di  $A$

$$\begin{aligned} \|p(A)(x_0 - x^*)\|_A^2 &= (x_0 - x^*)^T A (p(A))^2 (x_0 - x^*) = \\ &= \left( \sum_1^n \gamma_i u_i \right)^T A (p(A))^2 \left( \sum_1^n \gamma_i u_i \right) = \left( \sum_1^n \gamma_i u_i \right)^T \sum_1^n \gamma_i A (p(A))^2 u_i = \\ &= \left( \sum_1^n \gamma_i u_i \right)^T \sum_1^n \gamma_i A (p(\lambda_i))^2 u_i = \left( \sum_1^n \gamma_i u_i \right)^T \sum_1^n \gamma_i \lambda_i (p(\lambda_i))^2 u_i = \\ &= \sum_1^n \gamma_i^2 \lambda_i (p(\lambda_i))^2 \leq \max_i (p(\lambda_i))^2 \sum_1^n \gamma_i^2 \lambda_i = \max_i (p(\lambda_i))^2 \|x_0 - x^*\|_A^2 \end{aligned}$$

### Stima dell'errore relativo

$$\frac{\|e_{k+1}\|_A^2}{\|e_0\|_A^2} \leq \min_{p \in \Pi_{k+1}^1} \max_{1 \leq i \leq n} (p(\lambda_i))^2$$

## I polinomi di Chebyshev

$$\cos[(k+1)\theta] = 2 \cos \theta \cos k\theta - \cos[(k-1)\theta].$$

$$\cosh[(k+1)\theta] = 2 \cosh \theta \cosh k\theta - \cosh[(k-1)\theta],$$

$$\cosh x = \cos(ix) = \frac{e^x + e^{-x}}{2}, \quad \cosh^{-1} x = \ln(x + \sqrt{x^2 - 1}), \quad i = \sqrt{-1}.$$

Il polinomio di Chebyshev di prima specie di grado  $k$  si può definire come segue:

$$T_k(\cosh \theta) = \cosh(k\theta), \quad \theta = \cosh^{-1} x$$

$$\Leftrightarrow T_k(x) = \cosh(k \cosh^{-1} x) = \cosh(k \ln(x + \sqrt{x^2 - 1}))$$

$$\Leftrightarrow T_k(x) = \frac{1}{2} \left( (x + \sqrt{x^2 - 1})^k + (x - \sqrt{x^2 - 1})^k \right).$$

$$T_0(x) = 1, \quad T_1(x) = x \quad T_{k+1}(x) = 2xT_k(x) - T_{k-1}(x)$$

### Teorema

Sia

$$Q(x) = T_k \left( \frac{2x - (b+a)}{b-a} \right) / T_k \left( \frac{-(b+a)}{b-a} \right).$$

Allora

$$\min_{p \in \Pi_k^1} \max_{x \in [a,b]} |p(x)| = \max_{x \in [a,b]} |Q(x)| = \left| T_k \left( \frac{-(b+a)}{b-a} \right) \right|^{-1}.$$

## I polinomi di Chebyshev

$$\cos[(k+1)\theta] = 2 \cos \theta \cos k\theta - \cos[(k-1)\theta].$$

$$\cosh[(k+1)\theta] = 2 \cosh \theta \cosh k\theta - \cosh[(k-1)\theta],$$

$$\cosh x = \cos(ix) = \frac{e^x + e^{-x}}{2}, \quad \cosh^{-1} x = \ln(x + \sqrt{x^2 - 1}), \quad i = \sqrt{-1}.$$

Il polinomio di Chebyshev di prima specie di grado  $k$  si può definire come segue:

$$T_k(\cosh \theta) = \cosh(k\theta), \quad \theta = \cosh^{-1} x$$

$$\Leftrightarrow T_k(x) = \cosh(k \cosh^{-1} x) = \cosh(k \ln(x + \sqrt{x^2 - 1}))$$

$$\Leftrightarrow T_k(x) = \frac{1}{2} \left( (x + \sqrt{x^2 - 1})^k + (x - \sqrt{x^2 - 1})^k \right).$$

$$T_0(x) = 1, \quad T_1(x) = x \quad T_{k+1}(x) = 2xT_k(x) - T_{k-1}(x)$$

### Teorema

Sia

$$Q(x) = T_k \left( \frac{2x - (b+a)}{b-a} \right) / T_k \left( \frac{-(b+a)}{b-a} \right).$$

Allora

$$\min_{p \in \Pi_k^1} \max_{x \in [a,b]} |p(x)| = \max_{x \in [a,b]} |Q(x)| = \left| T_k \left( \frac{-(b+a)}{b-a} \right) \right|^{-1}.$$

## I polinomi di Chebyshev

$$\cos[(k+1)\theta] = 2 \cos \theta \cos k\theta - \cos[(k-1)\theta].$$

$$\cosh[(k+1)\theta] = 2 \cosh \theta \cosh k\theta - \cosh[(k-1)\theta],$$

$$\cosh x = \cos(ix) = \frac{e^x + e^{-x}}{2}, \quad \cosh^{-1} x = \ln(x + \sqrt{x^2 - 1}), \quad i = \sqrt{-1}.$$

Il polinomio di Chebyshev di prima specie di grado  $k$  si può definire come segue:

$$T_k(\cosh \theta) = \cosh(k\theta), \quad \theta = \cosh^{-1} x$$

$$\Leftrightarrow T_k(x) = \cosh(k \cosh^{-1} x) = \cosh(k \ln(x + \sqrt{x^2 - 1}))$$

$$\Leftrightarrow T_k(x) = \frac{1}{2} \left( (x + \sqrt{x^2 - 1})^k + (x - \sqrt{x^2 - 1})^k \right).$$

$$T_0(x) = 1, \quad T_1(x) = x \quad T_{k+1}(x) = 2xT_k(x) - T_{k-1}(x)$$

### Teorema

Sia

$$Q(x) = T_k \left( \frac{2x - (b+a)}{b-a} \right) / T_k \left( \frac{-(b+a)}{b-a} \right).$$

Allora

$$\min_{p \in \Pi_k^1} \max_{x \in [a,b]} |p(x)| = \max_{x \in [a,b]} |Q(x)| = \left| T_k \left( \frac{-(b+a)}{b-a} \right) \right|^{-1}.$$

## I polinomi di Chebyshev

$$\cos[(k+1)\theta] = 2 \cos \theta \cos k\theta - \cos[(k-1)\theta].$$

$$\cosh[(k+1)\theta] = 2 \cosh \theta \cosh k\theta - \cosh[(k-1)\theta],$$

$$\cosh x = \cos(ix) = \frac{e^x + e^{-x}}{2}, \quad \cosh^{-1} x = \ln(x + \sqrt{x^2 - 1}), \quad i = \sqrt{-1}.$$

Il polinomio di Chebyshev di prima specie di grado  $k$  si può definire come segue:

$$T_k(\cosh \theta) = \cosh(k\theta), \quad \theta = \cosh^{-1} x$$

$$\Leftrightarrow T_k(x) = \cosh(k \cosh^{-1} x) = \cosh(k \ln(x + \sqrt{x^2 - 1}))$$

$$\Leftrightarrow T_k(x) = \frac{1}{2} \left( (x + \sqrt{x^2 - 1})^k + (x - \sqrt{x^2 - 1})^k \right).$$

$$T_0(x) = 1, \quad T_1(x) = x \quad T_{k+1}(x) = 2xT_k(x) - T_{k-1}(x)$$

### Teorema

Sia

$$Q(x) = T_k \left( \frac{2x - (b+a)}{b-a} \right) / T_k \left( \frac{-(b+a)}{b-a} \right).$$

Allora

$$\min_{p \in \Pi_k^1} \max_{x \in [a,b]} |p(x)| = \max_{x \in [a,b]} |Q(x)| = \left| T_k \left( \frac{-(b+a)}{b-a} \right) \right|^{-1}.$$



## Il metodo del gradiente coniugato: errore relativo e condizionamento

Consideriamo  $T_k(s(\tau))$ , con

$$s(\tau) = \frac{2\tau - (\lambda_1 + \lambda_n)}{\lambda_n - \lambda_1}$$

Osserviamo che  $-1 \leq s(\tau) \leq 1 \Leftrightarrow \lambda_1 \leq \tau \leq \lambda_n$ . Inoltre, considerato

$$s_0 = s(0) = -\frac{\lambda_1 + \lambda_n}{\lambda_n - \lambda_1} = -\frac{\kappa + 1}{\kappa - 1}, \quad \text{con } \kappa = \frac{\lambda_n}{\lambda_1},$$

risulta

$$T_k(s_0) = \frac{(-1)^k}{2} \left( \left( \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^k + \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \right).$$

# Il metodo del gradiente coniugato: errore relativo e condizionamento (cont.)

## Teorema

$$\frac{\|e_k\|_A}{\|e_0\|_A} \leq \frac{2}{\left(\frac{a-1}{a+1}\right)^k + \left(\frac{a+1}{a-1}\right)^k} \leq 2 \left(\frac{a-1}{a+1}\right)^k, \quad \text{con } a = \sqrt{\kappa(A)}$$

**Dimostrazione.** Consideriamo il polinomio

$$Q(\tau) = \frac{T_k(s(\tau))}{T_k(s_0)}.$$

Si ha:

$$\frac{\|e_k\|_A^2}{\|e_0\|_A^2} \leq \max_{1 \leq i \leq n} |Q(\lambda_i)| \leq \max_{\lambda_1 \leq \tau \leq \lambda_n} |Q(\tau)| = \frac{1}{|T_k(s_0)|},$$

da cui segue la prima disuguaglianza del teorema. Per dimostrare la seconda disuguaglianza basta osservare che

$$\frac{2}{\left(\frac{a-1}{a+1}\right)^k + \left(\frac{a+1}{a-1}\right)^k} = \frac{2 \left(\frac{a-1}{a+1}\right)^k}{\left(\frac{a-1}{a+1}\right)^{2k} + 1}.$$

## Il metodo del gradiente coniugato: cluster di autovalori

Il risultato del teorema precedente può essere raffinato, per tenere conto, in particolare, del caso in cui gli autovalori siano concentrati intorno a  $\lambda_1$ .

### Teorema

Sia  $m$  un intero minore di  $n$ . Per il metodo CG vale la seguente maggiorazione per gli errori<sup>a</sup>:

$$\frac{\|e_k\|_A}{\|e_0\|_A} \leq 2 \left( \frac{a-1}{a+1} \right)^{k-m}, \quad \text{con } a = \sqrt{\kappa_m(A)} \text{ dove } \kappa_m(A) = \sqrt{\frac{\lambda_{n-m}}{\lambda_1}}$$

<sup>a</sup>Composite convergence bounds based on Chebyshev polynomials and finite precision conjugate gradient computations, T. Gergelits, · Z. Strakos

Esempio:

$\Lambda = \{1, 1.1, 1.2, 1.3, 1.4, 1.5, 1.6, 1.7, 1.8, 1000\}$ , se consideriamo  $m = 1$  avremo  $\kappa_m(A) = \lambda_9/\lambda_1 = 1.8$  e quindi  $\frac{a-1}{a+1} \simeq 0,146$

$\Lambda = \{1, 999.1, 999.2, 999.3, 999.4, 999.5, 999.6, 999.7, 999.8, 1000\}$ , se consideriamo  $m = 1$  avremo  $\kappa_m(A) = \lambda_9/\lambda_1 = 999.8$  e quindi  $\frac{a-1}{a+1} \simeq 0,939$

Si osservi che in entrambi gli esempi è  $\kappa(A) = 1000$

## Metodi di Krylov

$$Ax = b \quad (A \text{ non necessariamente sdp}) \quad (5)$$

Se  $A$  è non singolare la soluzione di (5) è  $x^* = A^{-1}b$ . A partire da un punto iniziale  $x_0$  un **metodo iterativo di Krylov** calcola la generica iterazione  $h$ —ma minimizzando **un certo errore** nello spazio affine  $h$ —dimensionale:

$$x_0 + K_h, \quad \text{dove } K_h = \text{Span}\{r_0, Ar_0, \dots, A^{h-1}r_0\}$$

dove  $r_0 = Ax_0 - b$  (residuo iniziale)

- $K_h = K_h(A, r_0)$   $k$ —mo spazio di Krylov;
- $r(x) = Ax - b$  residuo in  $x$  ( $r_h = b - Ax_h$ );
- Se  $A$  è non singolare, allora  $K_n = \mathbb{R}^n$ , e quindi un metodo di Krylov termina dopo (al più)  $n$  passi.

## CG

Il metodo CG é un metodo di Krylov. L'iterazione  $x_k$  é soluzione del problema

$$\min_{x \in x_0 + K_k} \|x^* - x\|_A, \text{ dove } K_k = \text{Span}\{r_0, Ar_0, \dots, A^{k-1}r_0\} \quad (6)$$

### Minimizzazione dell'errore

L'iterazione  $k$ -ma determina in  $x_0 + K_k$  il punto che ha distanza minima dalla soluzione  $x^*$  secondo la norma indotta da  $A$

## GMRES

Nel metodo GMRES, l'iterazione  $x_k$  è soluzione del problema

$$\min_{x \in x_0 + K_k} \|b - Ax\|, \text{ dove } K_k = \text{Span}\{r_0, Ar_0, \dots, A^{k-1}r_0\}$$

### Minimizzazione del residuo

Notiamo che, se  $x \in x_0 + K_k$ , allora

$$r(x) = -Ax + b = r_0 - \sum_0^k \gamma_i A^i r_0 = p(A)r_0, \text{ con } p \in \pi_k^1$$

### Teorema (dimostrazione: kelley)

Sia  $A$  non singolare, e sia  $x_k$  la  $k$ -ma iterazione generata da GMRES. Allora

$$\|r_k\| = \min_{p \in \pi_k^1} \|p(A)r_0\|$$

$$\|r_k\| = \min_{p \in \pi_k^1} \|p(A)r_0\| \Rightarrow \|r_k\| \leq \|r_0\| \min_{p \in \pi_k^1} \|p(A)\|$$

e quindi vale la seguente

### Maggiorazione del $k$ -mo residuo

$$\frac{\|r_k\|}{\|r_0\|} \leq \|p(A)\| \text{ per ogni } p \in \pi_k^1 \quad (7)$$

Necessità di stime dell'errore più "operative" di (7) (ad es., collegate allo spettro o al condizionamento di  $A$ ). Occorre fare ulteriori ipotesi su  $A$

### Maggiorazione del $k$ -mo residuo

$$\text{Se } \|I - A\| \leq \rho < 1, \text{ allora } \|g_k\| \leq \rho^k \|r_0\|$$

Dimostrazione: in (7) si prenda  $p(A) = (I - A)^k$ .

## Matrici Diagonalizzabili

Una matrice  $A$  si dice **diagonalizzabile** se esiste una matrice non singolare  $V \in \mathbb{C}^{n \times n}$  tale che

$$A = V^{-1} \Lambda V \Leftrightarrow V A V^{-1} = \Lambda$$

$$\Lambda = \begin{pmatrix} \lambda_1 & 0 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 & 0 \\ 0 & 0 & \cdot & 0 & 0 \\ 0 & 0 & 0 & \cdot & 0 \\ 0 & 0 & 0 & 0 & \lambda_1 \end{pmatrix}$$

con  $V$  **matrice diagonale** costituita dagli autovalori di  $A$ .

- ➊  $V A = \Lambda V$  (le colonne  $k$ -di  $V$  è autovettore destro di  $A$  di autovalore  $\lambda_k$ . Analogamente le righe di  $V^{-1}$  sono autovettori sinistri di  $A$ ).
- ➋ Dalla invertibilità di  $V$  segue che esiste una base di autovettori d/s (**condizione necessaria e sufficiente per la diagonalizzazione**)
- ➌ Se  $V$  è **ortogonale**, allora  $V$  ha come colonne gli autovettori di  $A$ ,  $V^{-1} = V^H$ , con  $V^H$  complessa coniugata di  $V$ .
- ➍ Se  $A$  è Hermitiana (simmetrica), gli autovettori possono essere scelti in maniera tale da formare una base ortonormale di  $\mathbb{C}^n$  ( $\mathbb{R}^n$ )



## Matrici Diagonalizzabili: stima del residuo

Se  $A = V^{-1}\Lambda V$ , allora, per ogni polinomio  $p(\cdot)$  si ha che

$$p(A) = V^{-1}p(\Lambda)V. \quad (8)$$

Da (7-8) segue che

$$\|r_k\| = \min_{p \in \pi_k^1} \|p(A)r_0\| \leq \|V\| \|p(\Lambda)\| \|V^{-1}\| \|r_0\| \quad \forall p \in \pi_k^1 \quad (9)$$

$$\leq \kappa(V) \max_{i \leq n} |p(\lambda_i)| \|r_0\| \quad \forall p \in \pi_k^1 \quad (10)$$

### Teorema (terminazione finita per GMRES ( $A$ diagonalizzabile))

Se  $A$  possiede solo  $k$  autovalori distinti  $\lambda_1, \lambda_2, \dots, \lambda_k$ , allora GMRES termina dopo (al più)  $k$  passi.

Dimostrazione Il polinomio

$$p(x) = \prod_{i=1}^k \left( \frac{\lambda_i - x}{\lambda_i} \right)$$

è tale che  $p(0) = 1$ ,  $p(\Lambda) = 0$ , e quindi, dalla (9) segue la tesi

## GMRES (Generalized Minimal RESidual method)

Nel metodo GMRES, l'iterazione  $x_k$  é soluzione del problema

$$\min_{x \in x_0 + K_k} \|b - Ax\|, \text{ dove } K_k = \text{Span}\{r_0, Ar_0, \dots, A^{k-1}r_0\} \quad (11)$$

Supponendo di conoscere una base ortonormale  $\{u_1, u_2, \dots, u_k\}$  di  $K_k$ , allora il problema (11) potrebbe essere formulato come

$$\min \|b - A(x_0 + y_1u_1 + y_2u_2 + \dots + y_ku_k)\| \quad (12)$$

o, equivalentemente,

$$\min \|b - A(x_0 + U_k y)\|, \quad y \in \mathbb{R}^k \iff \min \|r_0 - A(U_k y)\|, \quad y \in \mathbb{R}^k \quad (13)$$

dove  $U = [u_1, u_2, \dots, u_k]$

## GMRES: Costruzione di $U_{k+1}$

Costruzione iterativa dei vettori  $u_i$ .

Alla prima iterazione:

$$r_0 = Ax_0 - b; u_1 = r_0 / \|r_0\| \quad (14)$$

Supponiamo di voler costruire  $u_2$  tale che  $\{u_1, u_2\}$  sia una base ortonormale di  $\text{Span}\{r_0, Ar_0\}$ . Possiamo scegliere:  $u_2 = Au_1 + \beta u_1$ . Imponendo la condizione di ortonormalità:

$$u_2^T u_1 = 0 \Leftrightarrow (Au_1 + \beta u_1)^T u_1 = 0 \Leftrightarrow \beta = -h^{11} \quad \text{dove } h^{11} = (Au_1)^T u_1 = u_1^T Au_1;$$

$$\text{e quindi } u_2 = \frac{Au_1 - h^{11}u_1}{\|Au_1 - h^{11}u_1\|}$$

Analogamente, volendo costruire  $u_3$  tale che  $\{u_1, u_2, u_3\}$  sia una base ortonormale di  $\text{Span}\{g_0, Ag_0, A^2g_0\}$ :

$$u_3 = Au_2 + \alpha u_1 + \beta u_2;$$

$$u_1^T u_3 = 0 \Leftrightarrow \alpha = -u_1^T Au_2 = -h^{12};$$

$$u_2^T u_3 = 0 \Leftrightarrow \beta = -u_2^T Au_2 = -h^{22}; \text{ e quindi } u_3 = \frac{Au_2 - h^{12}u_1 - h^{22}u_2}{\|Au_2 - h^{12}u_1 - h^{22}u_2\|}$$

## GMRES: Algoritmo di Arnoldi

### Algoritmo Arnoldi

$$\textcircled{1} \quad g_0 = Ax_0 - b; u_1 = g_0 / \|g_0\|$$

$\textcircled{2} \quad \text{for } i=1:k$

$$u_{i+1} = \frac{Au_i - \sum_{j=1}^i h^{ji} u_j}{\|Au_i - \sum_{j=1}^i h^{ji} u_j\|} \quad \text{dove } h^{ji} = (Au_j)^T u_i \quad (15)$$

È immediato osservare che  $u_{i+1}$  è linearmente indipendente dai vettori  $u_1, \dots, u_i$ , appartiene a  $K_{i+1}$ ; inoltre, per  $k \leq i$

$$u_k^T \left( Au_i - \sum_{j=1}^i h^{ji} u_j \right) = u_k^T Au_i - \sum_{j=1}^i h^{ji} u_k^T u_j = u_k^T Au_i - h^{ki} u_k^T u_k = 0$$

e quindi i vettori  $u_1, \dots, u_i, u_{i+1}$ , rappresentano una base ortonormale di  $K_{i+1}$

## Algoritmo di Arnoldi: *The happy breakdown theorem*

Sia  $A$  non singolare, e supponiamo che

$$Au_i - \sum_{j=1}^i h^{ji} u_j = 0$$

Allora

$$x^* = A^{-1}b \in x_0 + K_i$$

Dimostrazione: Kelley (lemma 3.4.1)

Posto  $\gamma_i = \|Au_i - \sum_{j=1}^i h^{ji}u_j\|$ , considerata la matrice

$$P = \begin{bmatrix} h^{11} & \gamma_1 & & & & \\ h^{21} & h^{22} & \gamma_2 & & & \\ & \ddots & \ddots & \ddots & & \\ h^{n-1,1} & h^{n-1,2} & \ddots & \ddots & h^{n-1,n-1} & \gamma_{n-1} \end{bmatrix}, \quad (16)$$

se  $H = P^T$ , e  $H_k$  è la sottomatrice di  $H$  ottenuta prendendo le prime  $k$  colonne e  $k+1$  righe,

$$H_k = \begin{bmatrix} h^{11} & h^{21} & \ddots & \ddots & h^{k1} \\ \gamma_1 & h^{22} & & & h^{k2} \\ & \gamma_2 & \ddots & \ddots & \ddots \\ & & \ddots & \ddots & h^{kk} \\ & & & \ddots & \gamma_k \end{bmatrix}, \quad (17)$$

allora, per le direzioni di ricerca generate dal metodo di Arnoldi vale la formula

$$\boxed{AU_k = U_{k+1}H_k}$$

Osserviamo che la matrice  $H_k$  é **Hessemberg superiore**.

## GMRES: Algoritmo di Arnoldi

Supponiamo che l'algoritmo proceda senza *breakdown*. Ricordiamo che alla generica iterazione occorre risolvere il problema

$$\min \|g_0 - A(U_k)y\|$$

Essendo  $U_{k+1}H_k = AU_k$  e osservando che  $g_0/\|g_0\| = U_{k+1}e_1$ , tale problema può essere scritto come

$$\min \|g_0 - (AU_k)y\| \Leftrightarrow \min \|g_0 - (U_{k+1}H_k)y\| \Leftrightarrow \min \|U_{k+1}(\|g_0\|e_1 - H_k y)\|$$

Quindi, la generica iterazione  $k$ -ma può essere calcolata come  $x_k = x_0 + H_k y_k$ , con

$$y_k = \operatorname{argmin} \|\|g_0\|e_1 - H_k y\|$$

## GMRES: Algoritmo di Arnoldi

$$y_k = \operatorname{argmin} \|\|g_0\|e_1 - H_k y\|^2 = \operatorname{argmin} \frac{1}{2} y^T A_k y - b_k^T y$$

dove  $A_k = H_k^T H_k$  e  $b_k = H_k^T \|g_0\|e_1$ , e quindi ad ogni passo del metodo di Arnoldi occorrerà risolvere un sistema della forma  $y = (H_k^T H_k)^{-1} H_k^T \|g_0\|e_1$

### GMRES (Arnoldi)

$x_0 \in \mathbb{R}^n$ ;  $k \leftarrow 0$ ;

$r_0 \leftarrow b - Ax_0$ ,  $u_1 = r_0 / \|r_0\|$ ,

**while** *not\_stopcondition*

$k = k + 1$ ;

**for**  $j = 1 \dots k$

$h_{jk} \leftarrow (Au_k)^T u_j$

**endfor**

$u_{k+1} \leftarrow Au_k - \sum_{j=1}^k h_{jk} u_j$ ,  $h_{k+1,k} \leftarrow \|u_{k+1}\|$

$u_{k+1} \leftarrow u_{k+1} / \|u_{k+1}\|$

$y_k \leftarrow (H_k^T H_k)^{-1} H_k^T \|g_0\|e_1$

**endwhile**

$x_{sol} = x_0 + U_k y_k$



## GMRES: Algoritmo di Arnoldi

Complessità di tempo: all'iterazione  $k$ —ma

- Calcolo dei coefficienti  $h_{jk}$ : un **prodotto matvet**  $(Av_k) + k\text{vetvet}$ ;
- Risoluzione di un sistema  $k \times k$  (**la matrice è fattorizzata come prodotto di matrici di Hessemberg**)
- altri due prodotti **vetvet**

Al termine un **prodotto matvet** per ricostruire la soluzione

Complessità di spazio:

**Occorre memorizzare la matrice  $H$  e tutte le direzioni di ricerca  $u_i$**

Inconvenienti numerici:

- Il fenomeno del *fill-in* rende il metodo poco adatto ai problemi sparsi; in particolare, dover conservare tutte le direzioni di ricerca pone seri problemi di occupazione di memoria (possibile rimedio: restart)
- L'accumulo dell'errore di roundoff può produrre la perdita di ortogonalità fra le direzioni di ricerca (e il conseguente deterioramento delle proprietà di convergenza del metodo)(possibile rimedio: ortogonalizzazione di Gram-Schmidt modificata)

per i rimedi: vedi Kelley.

## Precondizionamento - CG

**Obiettivo:** risolvere un sistema equivalente a quello di partenza, con un indice di condizionamento piú favorevole.

**Strategia** Determinare una matrice non singolare  $M$ , che sia una "approssimazione di  $A^{-1}$ ", in modo tale che lo spettro di  $MA$  "si stringa" intorno a 1, riconducendo la risoluzione di  $Ax = b$  alla risoluzione di

$$MAx = Mb \quad (18)$$

Nota: se  $A$  è spd, non è detto che  $AM$  lo sia (quindi (18) non è applicabile a CG). Si considera quindi il piú generale schema di precondizionamento

$$MANy = Mb \quad \text{con} \quad y = N^{-1}x \quad (19)$$

con  $M$  e  $N$  precondizionatori sinistro e destro rispettivamente.

**Nota:** Se  $A$  è sdp e  $N = M^T$ , allora la matrice del sistema (19) sarà sdp

## Precondizionamento CG - Jacobi

Metodo di Jacobi per  $Ax = b$ :

$$\begin{aligned} A &= L + D + U \\ Ax = b &\Leftrightarrow (L + D + U)x = b; \quad Dx = b - (L + U)x \\ x_{k+1} &= D^{-1}b - D^{-1}(L + U)x_k \end{aligned}$$

Precondizionatore sinistro

$$M = (m_{ij}) : m_{ij} = \begin{cases} a_{ii} & \text{se } i = j \\ 0 & \text{altrimenti} \end{cases}$$

## Precondizionamento CG - Cholesky incompleto

Si consideri una matrice  $B$  spd e sia  $B = U^T U$  la sua fattorizzazione di Cholesky. Nel sistema preconditionato  $MANy = Mb$  in (19) possiamo scegliere  $M = U^{-T}$ ,  $N = U^{-1}$ .

Si osservi che, se  $B = A$ ,

$$\begin{aligned} U^{-T}AU^{-1}y = U^{-T}b &\Leftrightarrow U^{-T}U^TUU^{-1}y = U^{-T}b \Leftrightarrow \\ U^{-T}U^TUU^{-1}y = U^{-T}b &\Leftrightarrow Iy = U^{-T}b \end{aligned} \quad (20)$$

La matrice  $M$  dovrà essere scelta

- ❶ in modo da "approssimare bene"  $A$  ed avere condizionamento più favorevole  $(\kappa(U^{-T}AU^{-1}) \simeq 1)$ ;
- ❷ tale che il costo computazionale per calcolare  $U$  sia accettabile;
- ❸  $U$  sia facilmente "invertibile"

## Teorema

Sia

$$Q(x) = T_k \left( \frac{2x - (b+a)}{b-a} \right) / T_k \left( \frac{-(b+a)}{b-a} \right)$$

Allora  $\max |Q(x)|_{x \in [a,b]} = \min_{p \in \pi_k^1} \max |p(x)|_{x \in [a,b]}$

## Metodi di Newton inesatti

Il metodo di Newton per il sistema di equazioni non lineare  $F(x) = 0$

$$x_{k+1} = x_k + s_k \text{ dove}$$

$$F'(x_k)s_k = -F(x_k) \Leftrightarrow r_k = F'(x_k)s_k + F(x_k) = 0$$

**Metodi Newton-inesatti:** (44) é risolto in maniera approssimata con un **residuo**

$$r_k = F'(x_k)s_k + F(x_k)$$

(possibilmente) **non nullo** (e passo  $s_k = F'(x_k)^{-1}r_k - F'(x_k)^{-1}F(x_k)$ ). Che accuratezza richiedere su  $r_k$ ? Si considera, come "misura di accuratezza" il rapporto

$$\eta_k = \frac{\|r_k\|}{\|F(x_k)\|} \simeq \frac{\|F(x_{k+1})\|}{\|F(x_k)\|}$$

### Metodo di Newton

$x_0 \in \mathbb{R}^n$ ;  $k \leftarrow 0$ ;

**while** *not\_stopcondition*

    solve  $F'(x_k)s = F(x_k)$ ;

$x_{k+1} \leftarrow x_k + s$ ;  $k \leftarrow k + 1$ ;

**endwhile**

### Metodo di Newton Inesatto

$x_0 \in \mathbb{R}^n$ ;  $k \leftarrow 0$ ;

**while** *not\_stopcondition*

**approx**solve  $F'(x_k)s = F(x_k)$ ;

$x_{k+1} \leftarrow x_k + s$ ;  $k \leftarrow k + 1$ ;

**endwhile**



## Teorema: Convergenza metodi NI

Supponiamo che  $\{\eta_k\}$  sia tale che  $\eta_k < t < 1$ , allora esiste  $\varepsilon > 0$  tale che, se  $\|x_0 - x^*\| < \varepsilon$  la successione  $\{x_k\}$  converge a  $x^*$  **q-linearmente** e, in particolare,

$$\|x_{k+1} - x^*\|_* < t\|x_k - x^*\|_* \quad \text{dove } \|y\|_* = \|F'(x^*)y\| \quad (21)$$

Osserviamo che, poiché  $F'(x_*)$  è non singolare, per ogni  $\gamma > 0$  esiste  $\rho > 0$  tale che, se  $\|x - x^*\| < \rho$  allora

$$\begin{aligned} \|F'(x) - F'(x^*)\| &< \gamma \\ \|F'(x)^{-1} - F'(x^*)^{-1}\| &< \gamma \\ \|F(x) - F'(x^*)(x - x_*)\| &< \gamma\|x - x_*\| \end{aligned}$$

e, di conseguenza

$$\|F(x_k)\| = \| [F(x_k) - F'(x^*)e_k] + F'(x^*)e_k \| \leq 2\gamma\|e_k\|$$

Inoltre,

$$\frac{1}{\mu}\|y\| \leq \|y\|_* \leq \mu\|y\| \quad \text{dove } \mu = \max\{\|F'(x^*)\|, \|F'(x^*)^{-1}\|\}$$



$$\begin{aligned}
 F'(x^*)e_{k+1} &= F'(x^*)(e_k + s_k) = F'(x^*) [e_k + F'(x_k)^{-1}r_k - F'(x_k)^{-1}F(x_k)] \\
 &= F'(x^*)F'(x_k)^{-1} [F'(x_k)e_k + r_k - F(x_k)] = \\
 &= F'(x^*)F'(x_k)^{-1} \left[ r_k + \underbrace{(F'(x_k) - F'(x^*))e_k}_{\text{}} + \underbrace{(-F(x_k) + F'(x^*)e_k)}_{\text{}} \right];
 \end{aligned}$$

$$\begin{aligned}
 \|F'(x^*)F'(x_k)^{-1}\| &= \|F'(x^*)(F'(x_k)^{-1} - F'(x^*)^{-1}) + I\| \\
 &\leq 1 + \|F'(x^*)\| \cdot \|F'(x_k)^{-1} - F'(x^*)^{-1}\| \leq 1 + \gamma\mu
 \end{aligned}$$

$$\begin{aligned}
 \|e_{k+1}\|_* &\leq (1 + \gamma\mu)(\|r_k\| + \|F'(x_k) - F'(x^*)\|\|e_k\| + \|F(x_k) - F'(x^*)e_k\|) \\
 &\leq (1 + \gamma\mu)(\eta_k\|F(x_k)\| + \gamma\|e_k\| + \gamma\|e_k\|) \leq (1 + \gamma\mu)(2\gamma\eta_k + 2\gamma)\|e_k\| \\
 &\leq (1 + \gamma\mu)(2\gamma\eta_k + 2\gamma)\mu\|e_k\|_*
 \end{aligned} \tag{22}$$

Lemma: (maggiorazione di  $\|F\|$ ).

Se

$$\alpha = \max\{\|F'(x_*)\| + \frac{1}{2\|F(x_*)^{-1}\|}, 2\|F(x_*)^{-1}\|\}$$

allora

$$\frac{1}{\alpha}\|x - x^*\| \leq \|F(x)\| \leq \alpha\|x - x^*\|$$

## Teorema: Velocità di Convergenza metodi NI

Supponiamo che la successione  $\{x_k\}$  generata da un metodo NI converga a  $x^*$ . La convergenza è superlineare se e solo se

$$\|r_k\| = o(\|F(x_k)\|) \text{ per } k \rightarrow \infty$$

**Dimostrazione:** supponiamo  $\|e_{k+1}\| = o(\|e_k\|)$

$$\begin{aligned} r_k &= \underbrace{F'(x_k)s_k + F(x_k)}_{= F'(x_k)e_{k+1} - F'(x_k)e_k + F(x_k) - F'(x^*)e_k + F'(x^*)e_k} = \\ &= [F(x_k) - F'(x^*)e_k] + [(F'(x_k) - F'(x^*))e_k] + [(F'(x^*) + F'(x_k) - F'(x^*))e_{k+1}] \end{aligned}$$

e quindi

$$\|r_k\| \leq o(\|e_k\|) + o(1)\|e_k\| + [\|F'(x^*)\| + o(1)] o(\|e_k\|)$$

e quindi  $\|r_k\| = o(\|e_k\|) = o(\|F(x_k)\|) \quad k \rightarrow \infty$

Se  $\|r_k\| = o(\|F(x_k)\|)$ , ricordando ( 22) avremo

$$\begin{aligned} \|e_{k+1}\|_* &\leq (1 + \gamma\mu)(\|r_k\| + \|F'(x_k) - F'(x^*)\|\|e_k\| + \|F(x_k) - F'(x^*)e_k\|) \\ &= o(\|F(x_k)\|) + o(1)(\|e_k\|) + o(\|e_k\|) = o(\|e_k\|) \end{aligned}$$

## Teorema: Velocità di Convergenza metodi NI

Supponiamo che la successione  $\{x_k\}$  generata da un metodo NI converga a  $x^*$ . Se  $\eta_k \rightarrow 0$ , allora la convergenza è **q-superlineare**.

Il generico metodo di Newton inesatto per il sistema di equazioni non lineare  $F(x) = 0$  può essere scritto nella forma

$$x_{k+1} = x_k + s_k \quad \text{dove} \quad (23)$$

$$H_k s_k = -F(x_k) \quad (24)$$

In questo caso

$$r_k = F'(x_k) (-H_k F(x_k)) + F(x_k) \quad \text{e quindi}$$

$$\|r_k\| = \|F'(x_k) (-H_k F(x_k)) + F(x_k)\| \leq \|I - F'(x_k) H_k\| \cdot \|F(x_k)\|$$

In particolare, la successione  $\eta_k$  è maggiorata dalla successione  $\rho_k$  dove

$$\rho_k = \|I - F'(x_k) H_k\|$$

### Teorema: Velocità di Convergenza metodi NI

Supponiamo che la successione  $\{x_k\}$  generata da un metodo NI converga a  $x^*$ . Se  $F'(x)$ , è Lipschitz continua in  $x^*$  allora la convergenza è **q-quadratica**, cioè esiste  $C > 0$  tale che

$$\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2.$$

**Metodo di Newton Inesatto globalizzato**

$x_0 \in \mathbb{R}^n$ ;  $k \leftarrow 0$ ;  $t \in (0, 1)$

**while**  $\|F(x_k)\| > tol$

trova  $\eta_k \in [0, 1]$  e  $s_k$  tali che

$$\|F(x_k) + F'(x_k)s_k\| \leq \eta_k \|F(x_k)\| \quad (C1) \ \&\&$$

$$\|F(x_k + s_k)\| \leq [1 - t(1 - \eta_k)] \|F(x_k)\| \quad (C2)$$

$k \leftarrow k + 1$ ;  $x_{k+1} = x_k + s_k$

**endwhile**

Notiamo che, posto

$$\text{pred}_k = \|F(x_k)\| - \|F(x_k) + F'(x_k)s_k\|, \quad \text{ared}_k = \|F(x_k)\| - \|F(x_k + s_k)\|$$

$$(C1) \Rightarrow \text{pred}_k \geq (1 - \eta_k) \|F(x_k)\| \quad (25)$$

$$(C2) \Rightarrow \text{ared}_k \geq t(1 - \eta_k) \|F(x_k)\| \quad (26)$$

e quindi la condizione su  $s_k$  implica

$$\text{ared}_k \geq t(\text{pred}_k)$$

che può essere interpretata come una condizione di fedeltà del modello lineare di  $F$  ripetuto ad  $F$ .

**Metodo di Newton Inesatto con *Backtracking***

$x_0 \in \mathbb{R}^n$ ;  $k \leftarrow 0$ ;  $\eta_{max} \in [0, 1)$ ,  $0 < \theta_{min} < \theta_{max} < 1$

**while**  $\|F(x_k)\| > tol$

$\eta_k \in [0, \eta_{max}]$  calcola  $s_k$  tale che

$$\|F(x_k) + F'(x_k)s_k\| \leq \eta_k \|F(x_k)\|$$

**while**  $\|F(x_k + s_k)\| > [1 - t(1 - \eta_k)]\|F(x_k)\|$

scegli  $\Theta \in [\theta_{min}, \theta_{max}]$

$$s_k = \Theta s_k, \quad \eta_k = 1 - \Theta(1 - \eta_k)$$

**endwhile**

$$k \leftarrow k + 1; \quad x_{k+1} = x_k + s_k$$

**endwhile**

**Teorema: Convergenza NIB**

Supponiamo che la successione  $\{x_k\}$  generata da un metodo NIB abbia un punto di accumulazione  $x^*$ , con  $F'(x_*)$  invertibile, allora  $F(x_*) = 0$  e  $x_k \rightarrow x_*$ . Inoltre, per  $k$  sufficientemente grande, i valori iniziali di  $s_k$  ed  $\eta_k$  saranno accettabili.

Se la successione è limitata allora (in alternativa):

- Converge ad una soluzione  $x_*$
- Ammette punti di accumulazione con Jacobiano singolare.

## Implementazione di NIB: scelta dei parametri

$$\eta_{max} = 0.9; \quad \theta_{min} = 0.1; \quad \theta_{max} = 0.5; \quad t = 1e - 4$$

Termini forzanti e convergenza:

- $\eta_k = \eta < 1 \Rightarrow$  convergenza **lineare**
- $\eta_k \rightarrow 0 \Rightarrow$  convergenza **superlineare**
- $\eta_k = O(\|F(x_k)\|) \Rightarrow$  convergenza **quadratica**

Una possibile scelta per i termini forzanti:

$$\eta_k = \min\{\eta_{max}, \tilde{\eta}_k\} \text{ dove}$$

$$\tilde{\eta}_k = \frac{|||F(x_k)|| - ||F(x_{k-1}) + F'(x_{k-1})s_{k-1}|||}{||F(x_{k-1})||}$$



## Metodo delle secanti

Il metodo di Newton per il sistema di equazioni non lineare  $F(x) = 0$

$$x_{k+1} = x_k - (F'(x_k))^{-1} F(x_k) \quad (27)$$

è una immediata generalizzazione del metodo delle tangenti (Newton-Raphson) a funzioni di più variabili (formalmente ottenuta sostituendo la derivata della funzione con lo Jacobiano del sistema). Nel caso di funzioni di una variabile, le varianti del metodo di Newton che non utilizzano le derivate (secanti, falsa posizione, ...)

$$x_{k+1} = x_k - \frac{b - a}{F(b) - F(a)} F(x_k) \quad (28)$$

sono formalmente ottenute da (27) sostituendo la derivata con una approssimazione alle differenze finite  $\frac{F(b)-F(a)}{b-a}$  (con  $a, b$  opportunamente scelti).

**Un tale procedimento non è generalizzabile a funzioni di più variabili.** Occorre stabilire delle opportune condizioni da imporre alla approssimazione  $B_k$  della matrice Jacobiana  $F'(x_k)$

## Condizione secante

$$B_k s_k = y_k, \text{ dove } s_k = x_k - x_{k-1}, y_k = F(x_k) - F(x_{k-1}) \quad (29)$$

La condizione (29) ammette infinite soluzioni (che costituiscono uno spazio affine di dimensione  $n^2 - n$ )

## Condizione secante: Broyden (good)

$$B_{k+1} = \operatorname{argmin}_{X \in \mathbb{R}^{n \times n}} \|X - B\|_F \quad \text{s.t. } Xs_k = y_k \quad \text{con } B = B_k \quad (30)$$

Il problema (30) nelle  $n^2$  variabili  $x_{ij}$  può essere formulato come <sup>1</sup>

$$\min \Phi(X) = XX - 2BX \quad \text{s.t. } Xs = y, \quad (31)$$

o, equivalentemente,

$$\min \Phi(X) = 0.5(XX - 2BX) \quad \text{s.t. } \Theta_i(X) = X^i s = y_i, \quad i = 1..n, \quad (32)$$

dove  $X^i$  è la  $i$ -ma riga di  $X$  e gli indici di  $s_k, y_k$  sono omessi per semplificare le notazioni. Le condizioni di Lagrange per (32) sono

$$\nabla \Phi(X) = X - B = \sum_i^n \lambda_i \nabla \Theta_i(X) \Leftrightarrow X^i - B^i = \lambda_i s \Leftrightarrow X - B = \lambda s^T \quad (33)$$

Per ricavare il vettore dei parametri di Lagrange in (33), moltiplicando ambo i membri per  $s$  e ricordando il vincolo  $Xs = y$  si ottiene

$$Xs - Bs = \lambda s^T s \Leftrightarrow y - Bs = \lambda s^T s \quad \text{da cui} \quad (34)$$

$$\lambda = \frac{y - Bs}{s^T s} \quad \text{e quindi } X = B + \lambda s^T = B + \frac{ys^T - (Bs)s^T}{s^T s} \quad (35)$$

<sup>1</sup>il prodotto matrice matrice va inteso come prodotto componente per componente:  
 $C = AB, \quad c_{ij} = a_{ij}b_{ij}$

## Condizione secante: Broyden (good)

$$B_k = B_{k-1} + \frac{y_k s_k^T - (B_{k-1} s_k) s_k^T}{s_k^T s_k} = B_{k-1} + p_k s_k^T \quad \text{dove} \quad p_k = \frac{y_k - B_{k-1} s_k}{s_k^T s_k} \quad (36)$$

- Con la (36) all'iterazione  $k$  viene calcolata una approssimazione della matrice Jacobiana  $F'(x_k)$  a partire da quella calcolata all'iterazione  $k - 1$  mediante un *update* di rango 1.
- La risoluzione di un sistema lineare è l'operazione più onerosa per il calcolo del passo ad ogni iterazione;
- In generale la successione  $\{B_k\}$  non converge a  $F'(x^*)$ ;
- La formula (36) non preserva simmetria e pd.
- Convergenza superlineare:  $\lim_{k \rightarrow \infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} = 0$

**Metodo di Broyden (diretto) Good**

$$x_0 \in \mathbb{R}^n; \quad B_0 \in \mathbb{R}^{n \times n}; \quad k \leftarrow 1;$$

$$F_0 = F(x_0); \quad s_1 \leftarrow -B_0^{-1} F_0;$$

$$x_1 \leftarrow x_0 + s_1, \quad F_1 \leftarrow F(x_1);$$

$$y_1 = F_1 - F_0$$
**while**  $\|F(x_k)\| > tol$ 

$$B_k \leftarrow B_{k-1} + \frac{y_k s_k^T - (B_{k-1} s_k) s_k^T}{s_k^T s_k}$$

$$s_{k+1} \leftarrow -B_k^{-1} F_k;$$

$$x_{k+1} \leftarrow x_k + s_{k+1}; \quad F_{k+1} = F(x_{k+1})$$

$$k \leftarrow k + 1; \quad y_k = F_k - F_{k-1}$$
**endwhile**

## Condizione secante: Broyden (bad)

Idea: fare una update di una **matrice** che approssimi l'inversa della matrice Jacobiana.

## Formula di Sherman Morrison Woodbury

$$(A + vw^T)^{-1} = A^{-1} - \frac{A^{-1}vw^T A^{-1}}{1 + w^T A^{-1}v}$$

$$\begin{aligned} \left( A^{-1} - \frac{A^{-1}vw^T A^{-1}}{1 + w^T A^{-1}v} \right) (A + vw^T) &= I + A^{-1}vw^T - \frac{A^{-1}vw^T + A^{-1}vw^T A^{-1}vw^T}{1 + w^T A^{-1}v} = \\ I + \frac{A^{-1}vw^T + A^{-1}vw^T \overbrace{w^T A^{-1}v} - A^{-1}vw^T - A^{-1}v \overbrace{w^T A^{-1}v} w^T}{1 + w^T A^{-1}v} &= I \end{aligned}$$

Utilizzando SMW con  $B_{k-1} \rightarrow A$ ,  $p_k \rightarrow v$ ,  $s_k \rightarrow w$ ,  $H = (B)^{-1}$  si ha

$$H_k = (B_k)^{-1} = (B_{k-1} + p_k s_k^T)^{-1} = H_{k-1} - \frac{H_{k-1} p_k s_k^T H_{k-1}}{1 + s_k^T H_{k-1} p_k} = \quad (37)$$

$$H_{k-1} - \frac{(H_{k-1} y_k - s_k) s_k^T H_{k-1} / s_k^T s_k}{1 + \frac{s_k^T H_{k-1} y_k - s_k^T s_k}{s_k^T s_k}} = H_{k-1} + \frac{(s_k - H_{k-1} y_k) s_k^T H_{k-1}}{s_k^T H_{k-1} y_k} \quad (38)$$

## Condizione secante: Broyden (bad)

$$H_k = H_{k-1} + \frac{(s_k - H_{k-1}y_k)s_k^T H_{k-1}}{s_k^T H_{k-1}y_k} \quad (39)$$

- Con la (39) all'iterazione  $k$  viene calcolata una approssimazione dell'inversa della matrice Jacobiana a partire da quella calcolata all'iterazione  $k - 1$  mediante un *update* di rango 1.
- Un prodotto matrice vettore è il calcolo più oneroso richiesto per il calcolo del passo ad ogni iterazione;
- In generale la successione  $\{H_k\}$  non converge a  $F'(x^*)^{-1}$ ;
- La formula (39) non preserva simmetria e pd.
- Convergenza superlineare:  $\lim_{k \rightarrow \infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} = 0$

## Matrici di Broyden: invertibilità

### Matrix Determinant Lemma (MDL)

Siano  $A \in \mathbb{R}^{n \times n}$ ,  $u, v \in \mathbb{R}^n$ , con  $A$  non singolare. Allora

$$\det(A + uv^T) = (1 + v^T A^{-1} u) \det(A)$$

### Matrici di Broyden: invertibilità

$$B_{k-1} \text{ invertibile} \Rightarrow B_k \text{ invertibile sse } s_k^T B_{k-1}^{-1} y_k \neq 0$$

$$H_{k-1} \text{ invertibile} \Rightarrow H_k \text{ invertibile sse } s_k^T H_{k-1} y_k \neq 0$$

**Dimostrazione:**

$$\det(B_k) = \det\left(B_{k-1} + \frac{y_k s_k^T - (B_{k-1} s_k) s_k^T}{s_k^T s_k}\right) = \det\left(B_{k-1} + p_k s_k^T\right) \left[p_k = \frac{y_k - B_{k-1} s_k}{s_k^T s_k}\right]$$

$$\text{Applicando MDL si ha } \det(B_k) = \det(B_{k-1}) \left(1 + s_k^T B_{k-1}^{-1} \frac{y_k - B_{k-1} s_k}{s_k^T s_k}\right)$$

$$= \det(B_{k-1}) \left(1 + \frac{s_k^T B_{k-1}^{-1} y_k}{s_k^T s_k} - 1\right) = \det(B_{k-1}) \frac{s_k^T B_{k-1}^{-1} y_k}{s_k^T s_k}$$



Broyden per  $F(x) = Ax - b$ 

$$B_k = B_{k-1} + \frac{y_k s_k^T - (B_{k-1} s_k) s_k^T}{s_k^T s_k} = B_{k-1} + p_k s_k^T \quad \text{dove} \quad p_k = \frac{y_k - B_{k-1} s_k}{s_k^T s_k}$$

Se  $F(x) = Ax - b$ , allora  $y_k = A s_k$  e, posto  $E_{k-1} = B_{k-1} - A$ , si ha

$$p_k = \frac{(A - B_{k-1}) s_k}{s_k^T s_k} = -\frac{E_{k-1} s_k}{s_k^T s_k}, \text{ e quindi } B_k = B_{k-1} - E_{k-1} \frac{s_k s_k^T}{s_k^T s_k}$$

$$E_k = E_{k-1} \left( I - \frac{s_k s_k^T}{s_k^T s_k} \right) \Rightarrow \|E_k\|_2 \leq \|E_{k-1}\|_2 \text{ essendo } \left\| I - \frac{s_k s_k^T}{s_k^T s_k} \right\|_2 = 1$$

Al crescere delle iterazioni, le matrici di Broyden "si avvicinano" ad  $A$ .

## Convergenza dei Metodi di Broyden

Broyden: convergenza finita per  $Ax = b$  (Gay 1979)

Nel caso in cui sia  $F(x) = Ax - b$ , il metodo di Broyden *Good* [Bad] **converge in al più  $2n$  iterazioni**, purché la matrice  $B_0$  [ $H_0$ ] sia non singolare e durante le iterazioni sia soddisfatta la condizione  $s_k^T B_{k-1}^{-1} y_k \neq 0$  [ $s_k^T H_{k-1} y_k \neq 0$ ]

Broyden: convergenza per  $F \in C^1$

Nel caso in cui  $F(x)$  sia di classe  $C^1$ ,  $x_*$  sia tale che  $F(x_*) = 0$ ,  $\det(F'(x_*)) \neq 0$  con  $F'(x)$  Lipschitz continua in un intorno di  $x_*$ , e inoltre si verifichi  $s_k^T B_{k-1}^{-1} y_k \neq 0$  [ $s_k^T H_{k-1} y_k \neq 0$ ]

il metodo di Broyden *Good* [Bad] è localmente convergente intorno a  $x_*$ , con velocità di convergenza superlineare.