

## Lezione 25

In questa lezione si conclude la descrizione della procedura per la ricerca di una fattorizzazione QR di una matrice e si discute un procedimento che utilizza tale procedura per la ricerca della soluzione di un sistema di equazioni lineari. Infine si approfondisce la nozione di *costo* di una procedura e lo si utilizza per confrontare i due procedimenti per la ricerca della soluzione di un sistema di equazioni lineari introdotti. Contestualmente si descrive un modello più realistico dell'insieme dei numeri di macchina. Le parti di testo in **magenta** sono asserti *omessi in classe*.

Riprendiamo la descrizione del procedimento di ricerca di una *fattorizzazione QR* di un'assegnata matrice iniziato nella Lezione scorsa.

- *Esempio* (continua).

– *Secondo passo*

Siano  $\Omega$  e  $\Theta$  le matrici determinate dal *Primo passo* e:

$$\Delta = \text{diag}(\|\omega_1\|, \|\omega_2\|, \|\omega_3\|) \in \mathbb{R}^{3 \times 3}$$

Si ricordi che  $\omega_1 \neq 0$  e  $\omega_2 \neq 0$ . Se anche  $\omega_3 \neq 0$  allora  $\Delta$  è invertibile e la coppia:

$$U = \Omega \Delta^{-1} \quad , \quad T = \Delta \Theta$$

è una fattorizzazione QR di  $A$ .

La procedura seguente formalizza il procedimento descritto:

- $[U, T] = \text{GS}(A)$

//  $A$  matrice  $n \times n$  ad elementi reali.

$\omega_1 = a_1$ ;

per  $k = 1, \dots, n - 1$  **ripeti**:

se  $\omega_k = 0$  allora: interrompi la costruzione;

altrimenti:

$$\quad d_k = \omega_k \cdot \omega_k$$

$$\quad \text{per } j = k + 1, \dots, n \text{ **ripeti**: } \theta_{kj} = \frac{a_j \cdot \omega_k}{d_k};$$

$$\quad \omega_{k+1} = a_{k+1} - (\omega_1 \theta_{1,k+1} + \dots + \omega_k \theta_{k,k+1});$$

se  $\omega_n = 0$  allora: interrompi la costruzione;

altrimenti:

$$\quad d_n = \omega_n \cdot \omega_n;$$

$$\quad \Delta = \text{diag}(\sqrt{d_1}, \dots, \sqrt{d_n});$$

$$U = (\omega_1, \dots, \omega_n) \Delta^{-1}; T = \Delta \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & \theta_{ij} & \\ & & & \ddots \\ & & 0 & & 1 \end{bmatrix}$$

La procedura GS termina correttamente *se e solo se* applicata ad una matrice *invertibile*.

- *Osservazione*

- (1) La procedura GS (che realizza il *metodo di ortonormalizzazione di Gram-Schmidt*, da cui il nome) mostra che *qualunque matrice invertibile ammette fattorizzazione QR*. Si osservi che se  $U, T$  è una fattorizzazione QR di  $A$  ed  $E \in \mathbb{R}^{n \times n}$  è una matrice diagonale tale che  $|e_{11}| = 1, \dots, |e_{nn}| = 1$  allora la coppia:

$$U' = UE \quad , \quad T' = ET$$

è a sua volta una fattorizzazione QR di  $A$ . Dunque: *la fattorizzazione QR non è unica*.

- (2) Esistono altre procedure per la ricerca di una fattorizzazione QR di una matrice, più generali di GS e preferibili ad essa da un punto di vista numerico. Queste procedure terminano correttamente *in ogni caso* e dunque mostrano che *qualunque matrice*  $n \times n$  *ammette fattorizzazione QR*. Scilab ha una funzione predefinita per la ricerca di una fattorizzazione QR di una matrice che utilizza il *metodo di Householder*:<sup>1</sup>

\* `qr`

Questa *funzione predefinita* restituisce una *fattorizzazione QR* di un'assegnata matrice  $A$ . Precisamente, se  $A$  è una matrice  $n \times n$ :

$$[U, T] = \text{qr}(A)$$

restituisce la coppia  $U, T$  di matrici  $n \times n$  che *approssima* una fattorizzazione QR di  $A$ . Come già osservato  $A$  può non essere invertibile (vedere l'Esercizio 4).

- *Procedimento di ricerca della soluzione di un sistema di equazioni lineari*

Siano  $A \in \mathbb{R}^{n \times n}$  e  $b \in \mathbb{R}^n$ . Per verificare l'invertibilità di  $A$  e cercare la soluzione del sistema  $Ax = b$  si utilizza il procedimento seguente:

- (1) Si applica ad  $A$  una procedura di calcolo di una fattorizzazione QR;
- (2) Dette  $U, T$  le matrici generate dalla procedura:
  - (2a) Si calcola  $c = U^T b$  soluzione del sistema  $Ux = b$ ;
  - (2b) Se  $T$  non invertibile allora il procedimento si arresta *altrimenti* si determina la soluzione  $x^*$  del sistema  $Tx = c$  utilizzando la procedura di sostituzione all'indietro.

Il procedimento è *soddisfacente*: Se  $A$  non è invertibile allora la matrice  $T$  è non invertibile ed il procedimento si arresta prematuramente, *altrimenti* il vettore finale  $x^*$  è la soluzione del sistema dato.

- *Uso del calcolatore*

Utilizzando il calcolatore (che, si ricordi, opera in  $F(\beta, m)$ ), il procedimento si modifica come segue: (0) I dati  $A$  e  $b$  del sistema si trasformano in  $\hat{A}$  e  $\hat{b}$ , matrice e colonna di elementi  $\text{rd}(a_{ij}), \text{rd}(b_i), i, j = 1, \dots, n$ ; (1) Si applica ad  $\hat{A}$  la procedura `qr`; (2) Dette  $\hat{U}, \hat{T}$  le matrici determinate al passo precedente: si calcola il vettore  $\hat{c} = \hat{U}^T * \hat{b}$ ; (3) Se  $\hat{T}$  invertibile: si applica la procedura `SI`, realizzazione di SI, alla coppia  $\hat{T}, \hat{c}$  per il calcolo del vettore  $\hat{x}$ . Il vettore finale è utilizzato per approssimare la soluzione  $x^*$  del sistema  $Ax = b$ .

Anche in questo caso per giudicare l'accuratezza di  $\hat{x}$  come approssimazione di  $x^*$  occorre studiare il *condizionamento* del calcolo della soluzione del sistema  $Tx = c$ , ovvero indagare il *numero di condizionamento* della matrice  $T$ . Scelta in  $\mathbb{R}^n$  la norma euclidea  $N_2$ , si studia:

$$c_2(T) = \|T^{-1}\|_2 \|T\|_2$$

Si ha:

$$A = UT \quad \Rightarrow \quad T = U^T A \quad \text{e} \quad A^{-1} = T^{-1} U^T \quad \Rightarrow \quad T^{-1} = A^{-1} U$$

Dunque:

$$\|T\|_2 \leq \|U^T\|_2 \|A\|_2 \quad \text{e} \quad \|T^{-1}\|_2 \leq \|A^{-1}\|_2 \|U\|_2$$

Ma: per ogni matrice  $M$  ortogonale si ha  $\|M\|_2 = 1$ . Perciò:

$$\|T\|_2 \leq \|A\|_2 \quad \text{e} \quad \|T^{-1}\|_2 \leq \|A^{-1}\|_2 \quad \Rightarrow \quad c_2(T) \leq c_2(A)$$

In questo caso la *procedura di fattorizzazione QR produce un fattore destro con proprietà di condizionamento non peggiori di quelle della matrice iniziale*. Il procedimento di ricerca della soluzione del sistema di equazioni lineari esposto sopra è *soddisfacente anche* quando realizzato al calcolatore.

<sup>1</sup>Si veda, ad esempio: [https://en.wikipedia.org/wiki/QR\\_decomposition#Using\\_Householder\\_reflections](https://en.wikipedia.org/wiki/QR_decomposition#Using_Householder_reflections).

## (G) Costo

La nozione di *costo*, già implicitamente adottata in precedenza, è quella “aritmetica:”

- *Definizione* (costo aritmetico).

Sia  $\phi$  un algoritmo, ovvero una sequenza *finita* di composizioni di *funzioni predefinite*. Si chiama *costo aritmetico* di  $\phi$  il numero  $C(\phi)$  di *pseudo-operazioni aritmetiche eseguite per calcolare  $\phi$* .<sup>2</sup>

Sia  $\phi$  un algoritmo. Il costo di  $\phi$  deve essere una quantità almeno approssimativamente proporzionale al *tempo* necessario al calcolatore per portare a termine il calcolo. La definizione adottata riesce nell'intento se:

- (a) La maggior parte del tempo impiegato dal calcolatore per calcolare  $\phi$  è speso nell'esecuzione di pseudo-operazioni aritmetiche.
- (b) Il tempo necessario per eseguire ciascuna pseudo-operazione aritmetica è lo stesso, in particolare *non dipende dagli operandi*.

Il sussistere della condizione (a) *dipende da  $\phi$* . Ad esempio, se  $\phi$  è la realizzazione di una procedura che calcola la norma infinito di un vettore il costo aritmetico è *zero* – per questo algoritmo *nessuna* delle funzioni predefinite calcolate è una pseudo-operazione aritmetica – ma non è vero che il calcolatore impiega tempo zero a calcolare  $\phi$ . Nel seguito la definizione di costo sarà applicata solo ad algoritmi che calcolano *quasi esclusivamente* pseudo-operazioni aritmetiche.

La condizione (b), invece, *non dipende da  $\phi$*  e secondo il modello di calcolatore presentato finora *non può essere verificata*. Si consideri, ad esempio, il calcolo in  $F(2, 53)$  dello pseudo-prodotto  $2^{b_1} \otimes 2^{b_2} = 2^{b_1+b_2}$ . *Non è ragionevole* supporre che il tempo necessario per il calcolo sia indipendente da quali numeri interi  $b_1$  e  $b_2$  occorre sommare (calcolare la somma di due numeri interi *non può* richiedere un tempo *indipendente* dal numero di cifre necessario per rappresentare gli addendi – si pensi al solo tempo necessario per *leggere gli addendi e scrivere la somma*). In un modello *più realistico* di calcolatore, però, l'insieme dei numeri di macchina è il seguente:

- *Definizione* (numeri in virgola mobile con esponente limitato ed elementi denormalizzati)

Siano  $\beta$  un numero intero maggiore di uno,  $m$  un numero intero positivo,  $b_{min}$  e  $b_{max}$  numeri interi tali che  $b_{min} < b_{max}$ . Si indica con:

$$F(\beta, m, b_{min}, b_{max})$$

l'insieme di tutti i numeri reali  $x$  tali che:

$$x = (-1)^s \beta^b 0.c_1 \cdots c_m$$

con  $s \in \{0, 1\}$ ,  $b$  numero intero tale che  $b_{min} \leq b \leq b_{max}$  e  $c_1, \dots, c_m$  cifre in base  $\beta$  *arbitrarie*. Gli elementi dell'insieme che si ottengono scegliendo  $c_1 \neq 0$  si dicono *normalizzati*. Gli elementi non nulli che si ottengono scegliendo  $b = b_{min}$  e  $c_1 = 0$  si dicono *denormalizzati*.  $F(\beta, m, b_{min}, b_{max})$  è un insieme di *numeri in virgola mobile con esponente limitato ed elementi denormalizzati*.

Adottando questo modello, l'ipotesi (b) è verificata. Nell'Esercitazione vedremo come l'insieme dei numeri di macchina in *Scilab* sia fedelmente descritto dal nuovo modello.

Vediamo alcuni esempi di determinazione del costo di una procedura e poi confrontiamo i due procedimenti proposti per la ricerca della soluzione di un sistema di equazioni lineari.

- *Esempio*

– *Prodotto riga per colonna*

Sia  $\text{prc}_n$  la realizzazione del *prodotto riga per colonna*  $a^T b$ , con  $a, b \in \mathbb{R}^n$ , ottenuta sostituendo ciascuna operazione aritmetica con la corrispondente pseudo-operazione e specificando l'ordine di composizione delle pseudo-operazioni. La corrispondenza biunivoca

---

<sup>2</sup>Il costo va valutato *nel caso peggiore*: il numero di pseudo-operazioni aritmetiche eseguite per calcolare  $\phi(x)$  potrebbe dipendere da  $x$ .

tra operazioni aritmetiche e pseudo-operazioni così stabilita rende possibile determinare il costo di  $\text{prc}_n(a, b)$  calcolando il numero di operazioni aritmetiche in  $a^\top b$ . Dette  $a_i, b_i$  le componenti di  $a, b$  si ha:

$$a^\top b = a_1 b_1 + \dots + a_n b_n$$

Il calcolo di  $a^\top b$  richiede  $n$  prodotti e  $n - 1$  somme. Allora:

$$C(\text{prc}_n) = 2n - 1$$

– *Prodotto matrice per colonna*

Sia  $\text{pmc}_n$  la realizzazione del *prodotto matrice per colonna*  $Ab$ , con  $A \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n$ . Si ragiona come nel caso del prodotto riga per colonna, ovvero si calcola il numero di operazioni aritmetiche in  $Ab$ . Dette  $r_1, \dots, r_n$  le righe di  $A$  si ha:

$$Ab = \begin{bmatrix} r_1 b \\ \vdots \\ r_n b \end{bmatrix}$$

Ciascuna delle  $n$  componenti del vettore  $Ab$  è un prodotto riga per colonna. Il calcolo di  $Ab$  richiede quindi  $n^2$  prodotti e  $n(n - 1)$  somme. Allora:

$$C(\text{pmc}_n) = n C(\text{prc}_n) = 2n^2 - n$$

– *Prodotto matrice triangolare per colonna*

Sia  $\text{pmtc}_n$  la realizzazione del *prodotto matrice triangolare per colonna*  $Tc$ , con  $T \in \mathbb{R}^{n \times n}$  matrice triangolare,  $c \in \mathbb{R}^n$ . Si ragiona come nel caso del prodotto matrice per colonna. Dette  $r_1, \dots, r_n$  le righe di  $T$  si ha:

$$Tc = \begin{bmatrix} r_1 c \\ \vdots \\ r_n c \end{bmatrix}$$

Ciascuna delle  $n$  componenti del vettore  $Tc$  è un prodotto riga per colonna. Questa volta, però, ciascuna componente ha un costo diverso. Supponendo, ad esempio,  $T$  triangolare superiore ed evitando di calcolare operazioni con risultato noto a priori (per ogni  $x \in \mathbb{R}$  si ha  $0x = 0$  e  $0 + x = x$ ) si ottiene:

$$C(\text{pmtc}_n) = C(\text{prc}_n) + C(\text{prc}_{n-1}) + \dots + C(\text{prc}_1) = 2(1 + 2 + \dots + n) - n = n^2$$

Il costo di  $\text{pmtc}_n$  è circa *metà* del costo di  $\text{pmc}_n$ .

– *Procedura SI*

Come già sappiamo (Lezione 17) il calcolo di  $\text{SI}(T, c)$  per  $T \in \mathbb{R}^{n \times n}$  e  $c \in \mathbb{R}^n$  richiede  $n$  divisioni e  $\frac{1}{2} n(n - 1)$  prodotti e somme. Allora:

$$C(\text{SI}) = n^2$$

Il costo è uguale a quello di  $\text{pmtc}_n$ . Gli stessi risultati si ottengono per  $\text{SA}$ .

– *Procedura EGPP*

Il calcolo di  $\text{EGPP}(A)$  per  $A \in \mathbb{R}^{n \times n}$  risulta richiedere:

$$\sum_{k=1}^{n-1} k = \frac{1}{2} n(n - 1) \text{ divisioni e } \sum_{k=1}^{n-1} k^2 = \frac{1}{6} (n - 1)n(2n - 1) \text{ prodotti e somme}$$

In totale:

$$C(\text{EGPP}) = \frac{2}{3} n^3 - \frac{1}{2} n^2 - \frac{1}{6} n$$

Si osservi che il calcolo di  $\text{EGPP}(A)$  richiede anche *confronti*, in numero trascurabile rispetto alle pseudo-operazioni aritmetiche.

Il procedimento per per la ricerca della soluzione di un sistema di equazioni lineari che utilizza la fattorizzazione LR ha dunque costo:

$$C(\text{EGPP}) + C(\text{SA}) + C(\text{SI}) = \frac{2}{3} n^3 + \text{termini di grado inferiore in } n$$

Un calcolo analogo per il procedimento per per la ricerca della soluzione di un sistema di equazioni lineari che utilizza la fattorizzazione QR porta ad un costo:<sup>3</sup>

$$C(\text{qr}) + C(\text{pmc}_n) + C(\text{SI}) = \frac{4}{3} n^3 + \text{termini di grado inferiore in } n$$

che è circa *il doppio* del precedente. Si osservi che si è scelto di esprimere il costo dei procedimenti mostrando esplicitamente solo *il termine dominante* al crescere di  $n$ . In entrambi i casi *il termine dominante è generato dalla procedura di ricerca della fattorizzazione*.

• *Osservazione (calcolo della matrice inversa)*

La funzione predefinita `inv` di *Scilab* determina un'approssimazione della matrice inversa di una data matrice invertibile. Sia  $A \in \mathbb{R}^{n \times n}$  una matrice invertibile. Dette  $e_1, \dots, e_n$  le colonne della matrice identica, il calcolo di  $Y = \text{inv}(A)$  avviene, sostanzialmente, in questo modo:

- $[S, D, P] = \text{EGPP}(A)$
- per  $k = 1, \dots, n$  ripeti:
  - $c = \text{SA}(S, P e_k);$
  - $y_k = \text{SI}(D, c);$
- $Y = (y_1 \dots, y_n).$

Dunque si calcolano le  $n$  colonne della matrice inversa come soluzione degli  $n$  sistemi lineari  $Ay = e_1, \dots, Ay = e_n$ . Tutti i sistemi hanno *la stessa matrice* e per la soluzione degli  $n$  sistemi è sufficiente calcolare *una sola volta* la fattorizzazione di  $A$ . Questo fa sì che il termine dominante del costo del procedimento sia ancora un multiplo di  $n^3$ .

### Esercizi

1. Dimostrare che la coppia  $U, T$  definita nel *Secondo passo* è una fattorizzazione QR di  $A$ .
2. Sia  $A \in \mathbb{R}^{3 \times 3}$ . Dimostrare che se  $\text{GS}(A)$  termina prematuramente allora  $A$  non è invertibile.
3. Sia:

$$A = \begin{bmatrix} 0 & 2 \\ 1 & -1 \end{bmatrix}$$

Determinare  $\text{GS}(A)$  e dedurre dal risultato due diverse fattorizzazioni QR di  $A$ .

4. Sia  $A \in \mathbb{R}^{n \times n}$  la matrice nulla. Per  $n = 5$  verificare che la procedura  $\text{GS}$  applicata ad  $A$  termina prematuramente mentre  $\text{qr}(A)$  determina una fattorizzazione QR *esatta*.
5. Dimostrare, utilizzando la definizione di norma indotta, che: *se  $M$  è una matrice ortogonale allora  $\|M\|_2 = 1$ .*
6. Verificare che il calcolo della matrice  $H_k$  nel passo  $k$ -esimo di  $\text{EGPP}(A)$  richiede  $n - k$  divisioni e che il calcolo del prodotto  $H_k P_k A^{(k)}$  richiede  $(n - k)^2$  prodotti e somme.
7. Verificare che il calcolo di  $\text{EGPP}(A)$  richiede non più di  $\frac{1}{2} n(n - 1)$  confronti. Determinare l'errore relativo commesso approssimando il numero di pseudo-operazioni aritmetiche e confronti richiesto dal calcolo di  $\text{EGPP}(A)$  con il numero delle sole pseudo-operazioni aritmetiche.
8. Sia  $\text{ata}_n$  la realizzazione del prodotto  $A^T A$  con  $A \in \mathbb{R}^{n \times n}$ . Determinare  $C(\text{ata}_n)$ . Si tenga conto che la matrice  $A^T A$  è *simmetrica*.

<sup>3</sup>Si veda la pagina di Wikipedia citata in precedenza. Si osservi che il calcolo della fattorizzazione QR richiede anche  $n - 1$  radici quadrate.