

Esercitazione 3

Approssimazione di Autovalori

Esercizio 1

1. Costruire una matrice di connessione $A \in \mathbb{R}^{100 \times 100}$ (dove $N = 100$ rappresenta il numero di pagine) con i seguenti comandi:

```
>> A=randi([0, 1],100,100);  
>> s=sum(A);  
>> for j=1:size(A,1)  
A(j,:)=A(j,:)./s;  
end
```

Tale matrice soddisfa la proprietà (2) nell'approfondimento a fine documento.

(Suggerimento per l'interpretazione: scrivere una matrice di elementi $a_{i,j} = 0, 1$ dove uno indica l'esistenza di un link dalla pagina p_j alla pagina p_i e zero l'assenza di link. Se non si vuole studiare un caso reale, questa matrice può essere generata in modo casuale tramite la funzione Matlab/Octave `randi`. Successivamente modificare A (normalizzando le colonne) in modo che rispetti la (2)).

2. Costruire una matrice di connessione $B \in \mathbb{R}^{5 \times 5}$ con i comandi

```
>> B=[  
0 0 0 1 0  
1 0 0 0 0  
0 1 0 0 0  
0 1 0 0 1  
1 1 1 1 0];  
>> s=sum(B);  
>> for j=1:size(B,1)  
B(j,:)=B(j,:)./s;  
end
```

La matrice B è relativa a 5 pagine web, come rappresentato in Fig. 1, secondo l'interpretazione data dall'approfondimento.

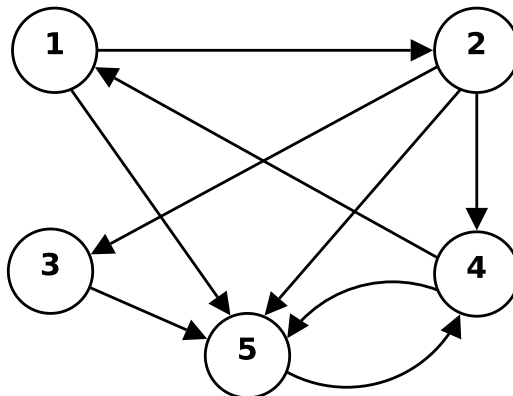


Figura 1: Schema di un web di 5 pagine. Ogni cerchio rappresenta una pagina web, ogni freccia un link.

3. Metodo delle potenze

Consideriamo una matrice $A \in \mathbb{C}^{n \times n}$ e supponiamo che i suoi autovalori siano così ordinati:

$$|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq |\lambda_4| \geq \dots \geq |\lambda_n|.$$

Denotiamo con \mathbf{x}_1 l'autovettore di norma unitaria associato all'autovalore dominante λ_1 . Se gli autovettori di A sono linearmente indipendenti, λ_1 e \mathbf{x}_1 possono essere calcolati tramite la seguente procedura, nota come *metodo delle potenze*:

- (a) si considera un vettore arbitrario $\mathbf{x}^{(0)} \in \mathbb{C}^n$ e si pone

$$\mathbf{y}^{(0)} = \mathbf{x}^{(0)} / \|\mathbf{x}^{(0)}\|, \quad \lambda^{(0)} = (\mathbf{y}^{(0)})^H A \mathbf{y}^{(0)};$$

- (b) si calcola:

per $k = 1, 2, \dots$

$$\mathbf{x}^{(k)} = A \mathbf{y}^{(k-1)}, \quad \mathbf{y}^{(k)} = \frac{\mathbf{x}^{(k)}}{\|\mathbf{x}^{(k)}\|}, \quad \lambda^{(k)} = (\mathbf{y}^{(k)})^H A \mathbf{y}^{(k)};$$

dove $\mathbf{y}^H = \bar{\mathbf{y}}^T$.

- (c) si termina l'algoritmo quando $|\lambda^{(k)} - \lambda^{(k-1)}| < \epsilon |\lambda^{(k)}|$, dove ϵ è una tolleranza assegnata.

Questo metodo genera una successione di vettori $\{\mathbf{y}^{(k)}\}$ di norma unitaria tali da allinearsi alla direzione dell'autovettore \mathbf{x}_1 , per $k \rightarrow \infty$. Si dimostra inoltre che $\lambda^{(k)} \rightarrow \lambda_1$, per $k \rightarrow \infty$. Si noti che il calcolo di $\lambda^{(0)}$ serve soltanto per testare la condizioni di uscita al passo $k = 1$.

Scrivere una funzione `eigpower` che implementa il metodo delle potenze per la ricerca dell'autovalore di modulo massimo e dell'autovettore associato. Una possibile intestazione potrebbe essere:

```
function [lambda,x,iter]=eigpower(A,tol,nmax,x0)
```

dove i parametri di ingresso sono:

- la matrice A ;
- la tolleranza `tol`;
- il numero massimo di iterazioni `nmax`;
- il vettore iniziale `x0`.

I parametri di uscita associati sono l'autovalore di modulo massimo `lambda`, l'autovettore associato `x` ed il numero di iterazioni effettuate `iter`.

4. Tramite la funzione `eigpower` si verifichi che A e B ammettono 1 come autovalore di modulo massimo e si calcoli il corrispondente autovettore (*PageRank*). Si assuma una tolleranza di 10^{-6} , un numero massimo di iterazioni pari a 100 ed un vettore iniziale $x_0 = [1/N, 1/N, \dots, 1/N]^T$.

Esercizio 2

Metodo delle potenze inverse

Dato che gli autovalori della matrice A^{-1} sono i reciproci di quelli di A , è possibile utilizzare il metodo delle potenze per approssimare anche l'autovalore di A di modulo minimo λ_n ($0 < |\lambda_n| < |\lambda_{n-1}| \leq |\lambda_{n-2}| \leq \dots \leq |\lambda_1|$). Si tratta del *metodo delle potenze inverse*:

1. si considera un vettore arbitrario $\mathbf{x}^{(0)} \in \mathbb{C}^n$ e si pone

$$\mathbf{y}^{(0)} = \mathbf{x}^{(0)} / \|\mathbf{x}^{(0)}\|, \quad \mu^{(0)} = (\mathbf{y}^{(0)})^H A^{-1} \mathbf{y}^{(0)};$$

2. si calcola:

per $k = 1, 2 \dots$

$$\mathbf{x}^{(k)} = A^{-1} \mathbf{y}^{(k-1)}, \quad \mathbf{y}^{(k)} = \frac{\mathbf{x}^{(k)}}{\|\mathbf{x}^{(k)}\|}, \quad \mu^{(k)} = (\mathbf{y}^{(k)})^H A^{-1} \mathbf{y}^{(k)};$$

dove $\mathbf{y}^H = \bar{\mathbf{y}}^T$.

3. si termina l'algoritmo quando $|\mu^{(k)} - \mu^{(k-1)}| < \epsilon |\mu^{(k)}|$, dove ϵ è una tolleranza assegnata.

Il metodo fornisce un'approssimazione di $\mu = 1/\lambda_n$, cioè $1/\mu^{(k)} \rightarrow \lambda_n$ per $k \rightarrow \infty$.

La matrice A^{-1} non viene realmente calcolata. Ad ogni passo si risolve invece il sistema lineare $A\mathbf{x}^{(k)} = \mathbf{y}^{(k-1)}$. Potrebbe quindi essere conveniente fattorizzare una volta sola, inizialmente, la matrice A con una opportuna fattorizzazione, così che ad ogni iterazione vengano risolti sistemi più semplici. Si noti inoltre che gli autovettori di A ed A^{-1} sono gli stessi ed è quindi possibile scrivere il quoziente di Rayleigh che definisce $\mu^{(k)}$ come $(\mathbf{y}^{(k)})^H A \mathbf{y}^{(k)}$, cioè con la matrice A al posto di A^{-1} ; in questo modo $\mu^{(k)}$ convergerà direttamente a λ_n , per $k \rightarrow \infty$, e non al suo reciproco. Il calcolo di $\mu^{(0)}$ serve soltanto per testare la condizioni di uscita al passo $k = 1$.

1. Scrivere una funzione `invpower` che implementa il metodo delle potenze per la ricerca dell'autovalore di modulo minimo e dell'autovettore associato. Una possibile intestazione potrebbe essere:

```
function [lambda,x,iter]=invpower(A,tol,nmax,x0)
```

dove `lambda` è l'autovalore di modulo minimo, mentre gli altri parametri di ingresso e di uscita hanno lo stesso significato attribuito dalla funzione `eigpower`. Si risolva il sistema lineare $A\mathbf{x}^{(k)} = \mathbf{y}^{(k-1)}$ tramite la fattorizzazione LU della matrice A e il metodo di sostituzione in avanti e all'indietro (utilizzare le funzioni `fwsb.m` e `bksb.m`).

2. Applicare la funzione scritta alla matrice ottenuta con il comando `toeplitz(1:4)`, usando tolleranza pari a 10^{-6} , numero massimo di iterazioni pari a 100 e vettore iniziale $x_0 = (1, 2, 3, 4)^T$. Provare poi ad usare la funzione `invpower` ponendo $x_0 = (1, 1, 1, 1)^T$ e commentare il risultato.

Approfondimento: Google PageRank

Dalla sua nascita, Google ha impiegato brevissimo tempo per imporsi come il più utilizzato motore di ricerca rispetto a tutti i concorrenti. Il successo è dipeso dalla accuratezza delle risposte, molto superiore a quelle di altri motori di ricerca. Google infatti stila una propria classifica (*ranking*) dell'importanza delle pagine web, per ordinare i risultati proposti; molto spesso questa classifica è proprio quella desiderata dall'utente.

L'algoritmo che ordina le pagine web per importanza è denominato *PageRank* ed è stato sviluppato da S. Brin e L. Page presso la Stanford University [1]. Il principio su cui si basa è il seguente:

- se una pagina web A ha un collegamento (*link*) verso una pagina B, questo è interpretato come un voto di A in favore di B (cioè alza B in graduatoria);
- i votanti non sono tutti uguali: il voto di chi è alto in classifica (cioè ha ricevuto molti link) vale maggiormente di quello di chi è in basso.

Questo principio è descritto schematicamente in Fig. 2. La pagina B riceve molti link e quindi è considerata molto importante. La pagina C è considerata più importante di E, anche se riceve meno link, perchè il voto di B conta molto di più dei tanti "piccoli" voti che riceve E (ad es. se il "Washington Post" o il "New York Times" pubblicassero un link al sito web di "Calcolo Numerico ed Elementi di Analisi", questo varrebbe molto di più dei sei link dai siti web personali di docenti ed esercitatori del corso!).

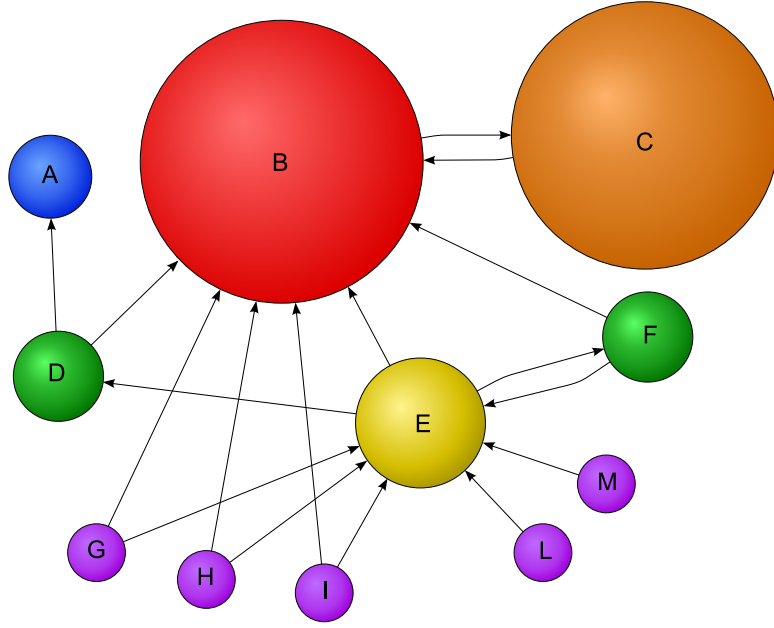


Figura 2: Diagramma schematico del *ranking*: ogni sfera rappresenta una pagina web, ogni freccia un link, la dimensione delle sfere è l'importanza attribuita alla pagina.

Dal punto di vista formale, si definisce il *ranking* (importanza relativa) $r(p)$ della pagina p come la somma di tutti i *ranking* $r(q)$ di tutte le pagine q che puntano p :

$$r(p) = \sum_{q \rightarrow p} \frac{r(q)}{\#q}; \quad (1)$$

la somma è ponderata da $\#q$, numero di link presenti nella pagina q (se una pagina q ha un solo link verso p è probabile che un lettore interessato lo segua, viceversa se la pagina q ha moltissimi link è poco probabile che il lettore scelga proprio quello verso p). È facile notare che si tratta di un problema in forma implicita: per conoscere il ranking di p si devono conoscere quelli delle altre pagine, che a loro volta si basano su quello di p .

Fortunatamente il problema è affrontabile in modo relativamente semplice, una volta scritto in forma matriciale. Siano $\{p_1, p_2, \dots, p_N\}$ tutte le pagine web censite e sia $A \in \mathbb{R}^{N \times N}$ la matrice di connessione, il cui elemento $a_{i,j}$ è dato da:

$$a_{i,j} = \begin{cases} \frac{1}{\#p_j} & \text{se esiste un link da } p_j \text{ a } p_i \\ 0 & \text{altrimenti.} \end{cases}$$

dove $\#p_j$ è il numero di link presenti sulla pagina p_j .

Si noti che gli $a_{i,j}$ possono essere interpretati come una distribuzione di probabilità, in particolare rappresentano la probabilità che un persona che clicca a caso giunga alla pagina p_i a partire dalla pagina p_j . La matrice gode della proprietà:

$$\sum_{i=1}^N a_{i,j} = 1. \quad (2)$$

Se il *ranking* delle pagine web p_i è rappresentato dal vettore colonna *PageRank* $\mathbf{r} = [r_1, r_2, \dots, r_n]^T$, allora l'equazione (1) equivale al seguente problema:

$$\mathbf{r} = A\mathbf{r}.$$

Il *PageRank* è quindi l'autovettore corrispondente all'autovalore 1 del problema agli autovalori associato. Si può dimostrare che se i λ_i sono gli autovalori di A allora $|\lambda_i| \leq 1$. Inoltre $\lambda_1 = 1$ ha molteplicità uno¹.

Dato che il numero di pagine censite N è dell'ordine di grandezza delle decine di miliardi, il calcolo con un metodo diretto dell'autovettore *PageRank* è computazionalmente troppo oneroso, anche per chi dispone di risorse di calcolo eccezionali. Si utilizza quindi un metodo iterativo, che restituisce una soluzione approssimata, basato sul metodo delle potenze. In questo caso particolare, arrestare il numero di iterazioni al passo k significa aver considerato solo le pagine p_j che distano da p_i al più k click².

Riferimenti bibliografici

- [1] S. Brin and L. Page, "The anatomy of a large-scale hypertextual web search engine", *Computer Networks and ISDN Systems*, vol. 30, no. 1-7, pp. 107 – 117, 1998. Proceedings of the Seventh International World Wide Web Conference.

¹Ciò vale sotto l'ipotesi che ogni pagina web abbia almeno un link. Un trucco per soddisfare l'ipotesi per le pagine senza link potrebbe essere quello di attribuire loro un link ad ogni pagina esistente. In questo laboratorio non si affronta questo caso.

²Per ragioni di semplicità di esposizione tralasciamo la possibilità di introdurre un fattore di rilassamento (*damping factor*) [1], che comunque non modifica il principio base dell'algoritmo.