

Progetto di Calcolo Scientifico: Spectral Clustering

Valeria Tateo

Ottobre 2024

L'obiettivo del progetto è suddividere un insieme di dati in k gruppi, detti cluster. Definita una metrica, i membri di uno stesso cluster sono "vicini", mentre vettori in cluster differenti sono ben separati.

1 Introduzione

Il problema del clustering può essere studiato introducendo il concetto di matrice Laplaciana associata ad un grafo. Dato un insieme di dati, si definisce una funzione

$$s(x_i, x_j) \geq 0$$

simmetrica e non negativa, che misura quanto due dati sono simili. Consideriamo il grafo pesato non orientato $G(V, E)$ e i pesi w_{ij} per gli archi (x_i, x_j) che siano funzione di $s(x_i, x_j)$. La matrice Laplaciana si definisce come

$$L = D - W$$

dove $W = (w_{ij})_{ij}$ è la matrice di adiacenza pesata e D è la matrice diagonale che ha come entrate le somme sulle righe di W .

2 Grafi di Somiglianza

In base alla funzione s scelta, si possono determinare Grafi di Somiglianza differenti. Di seguito si riportano gli algoritmi utilizzati per calcolare la matrice Laplaciana associata a tre tipi diversi di Grafi di Somiglianza.

Algorithm 1 Si connettono tutti i punti la cui distanza è minore di ϵ

```
1: procedure  $\epsilon$ -NEIGHBORHOOD( $v, \epsilon$ )
2:    $n \leftarrow \text{length}(v)$ 
3:    $W \leftarrow \text{zeros}(n, n)$ 
4:    $D \leftarrow \text{zeros}(n, n)$ 
5:
6:   for  $i = 1$  to  $n$  do
7:     for  $j = 1$  to  $n$  do
8:       if  $(i \neq j)$  and  $\|v(i) - v(j)\|_2 < \epsilon$  then
9:          $W(i, j) \leftarrow \|v(i) - v(j)\|_2$        $\triangleright$  In alternativa:  $W(i, j) \leftarrow \exp\left(-\frac{\|v(i) - v(j)\|_2^2}{2}\right)$ 
10:         $D(i, i) \leftarrow D(i, i) + W(i, j)$ 
11:       end if
12:     end for
13:   end for
14:
15:    $L \leftarrow D - W$ 
16:
17: return  $(L, D, W)$ 
18: end procedure
```

Algorithm 2 Si connette ogni vertice con i k-vertici più vicini

```
procedure K_NEAREST_NEIGHBOR_GRAPH( $v, k$ )
2:    $n \leftarrow \text{length}(v)$ 
    $W \leftarrow \text{zeros}(n, n)$ 
4:    $D \leftarrow \text{zeros}(n, n)$ 

6:   for  $i = 1$  to  $n$  do

8:      $\text{dists} \leftarrow \exp\left(-\frac{(v-v(i))^2}{2}\right)$             $\triangleright$  In alternativa:  $\text{dists} \leftarrow \|v(i) - v(j)\|_2$ 
        $\text{dists}(i) \leftarrow -\infty$                                 $\triangleright$  Esclude l'elemento stesso
10:     $[\sim, \text{sorted\_indices}] \leftarrow \text{sort}(\text{dists}, \text{'descend'})$ 
        $\text{nearest\_indices} \leftarrow \text{sorted\_indices}(1 : k)$         $\triangleright$  Seleziona i k vicini più vicini
12:     $W(i, \text{nearest\_indices}) \leftarrow 1$                     $\triangleright$  Aggiorna la matrice dei pesi

14:   end for

16:   for  $i \leftarrow 1$  to  $n$  do

18:     for  $j \leftarrow 1$  to  $n$  do

20:       if  $i \neq j$  then

22:          $D(i, i) \leftarrow D(i, i) + W(i, j)$ 

24:       end if

26:     end for

28:   end for

30:    $L \leftarrow D - W$ 

32: return (L,D,W)
end procedure
```

Algorithm 3 Si connettono tutti i punti e si pone il peso degli archi uguale alla somiglianza, dove la funzione di somiglianza è la densità di una gaussiana standard

```
procedure FULLY_CONNECTED_GRAPH( $v, \sigma$ )
    $n \leftarrow \text{length}(v)$ 
3:    $W \leftarrow \text{zeros}(n, n)$ 
    $D \leftarrow \text{zeros}(n, n)$ 

6:   for  $i = 1$  to  $n$  do
     for  $j = 1$  to  $n$  do
       if ( $i \neq j$ ) then
9:          $W(i, j) \leftarrow \exp\left(-\frac{\|v(i)-v(j)\|_2^2}{2}\right)$     $\triangleright$  In alternativa:  $W(i, j) \leftarrow \|v(i) - v(j)\|_2$ 
            $D(i, i) \leftarrow D(i, i) + W(i, j)$ 
       end if
12:    end for
  end for

15:    $L \leftarrow D - W$ 

   return (L,D,W)
18: end procedure
```

3 Inverse subspace iteration con shift negativo

Al fine di implementare una strategia di clustering, è necessario studiare un algoritmo che ci permetta di calcolare gli autovettori associati agli autovalori più piccoli di una matrice semidefinita positiva. Il seguente algoritmo ne è un esempio: sia A una matrice, $num_eigenvectors$ il numero di autovettori che si vuole calcolare, $num_iterations$ il numero di iterazioni (ad esempio pari a 100) che l'algoritmo esegue e $shift$ lo shift negativo utilizzato (ad esempio pari a -1).

Algorithm 4

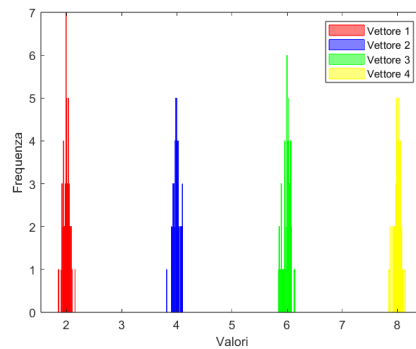
```

procedure INVERSE_SUBSPACE_ITERATION( $A, num\_eigenvectors, num\_iterations, shift$ )
   $n \leftarrow size(A, 1)$ 
   $X \leftarrow rand(n, num\_eigenvectors)$ 
  4: for  $iter = 1$  to  $num\_iterations$  do
     $X \leftarrow (A - shift \times eye(n)) \setminus X$ 
     $[X, \sim] \leftarrow qr(X, 0)$ 
  end for
  8:  $eigenvalues \leftarrow diag(X^T A X)$  ▷ Calcola gli autovalori approssimati
     $eigenvectors \leftarrow X$ 
    return ( $eigenvectors, eigenvalues$ )
end procedure

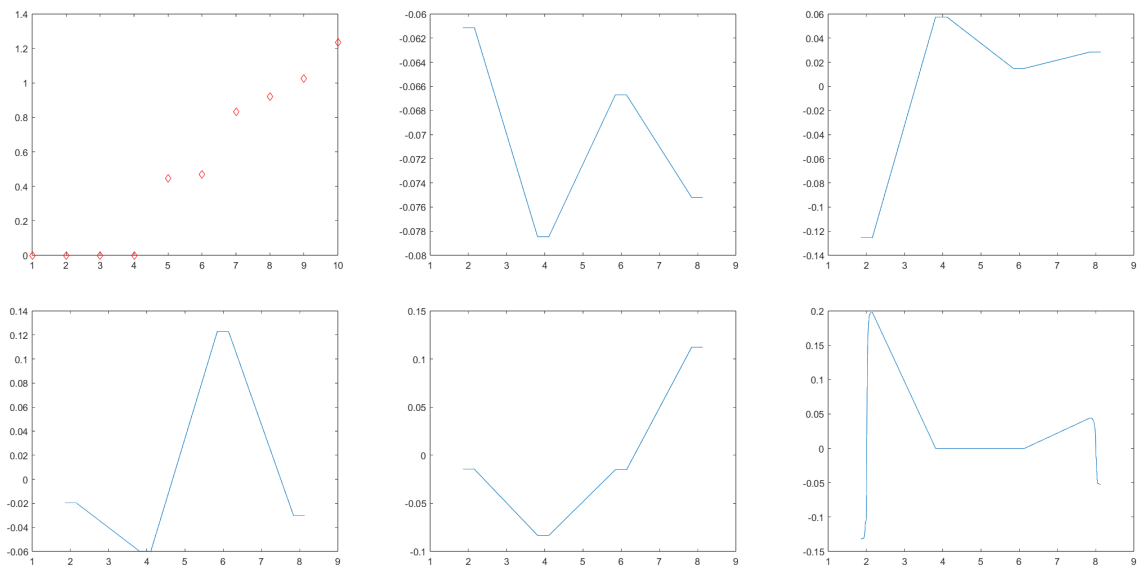
```

Per verificare la correttezza dell'algoritmo, lo testiamo sulla matrice Laplaciana calcolata utilizzando l'Algoritmo del 10-Nearest Neighbor Graph e su quella calcolata con l'Algoritmo del Fully Connected Graph, partendo da un insieme di dati costituito da 200 punti reali presi da 4 distribuzioni Gaussianne (50 per distribuzione) centrate in 2, 4, 6, 8 con varianza 16.

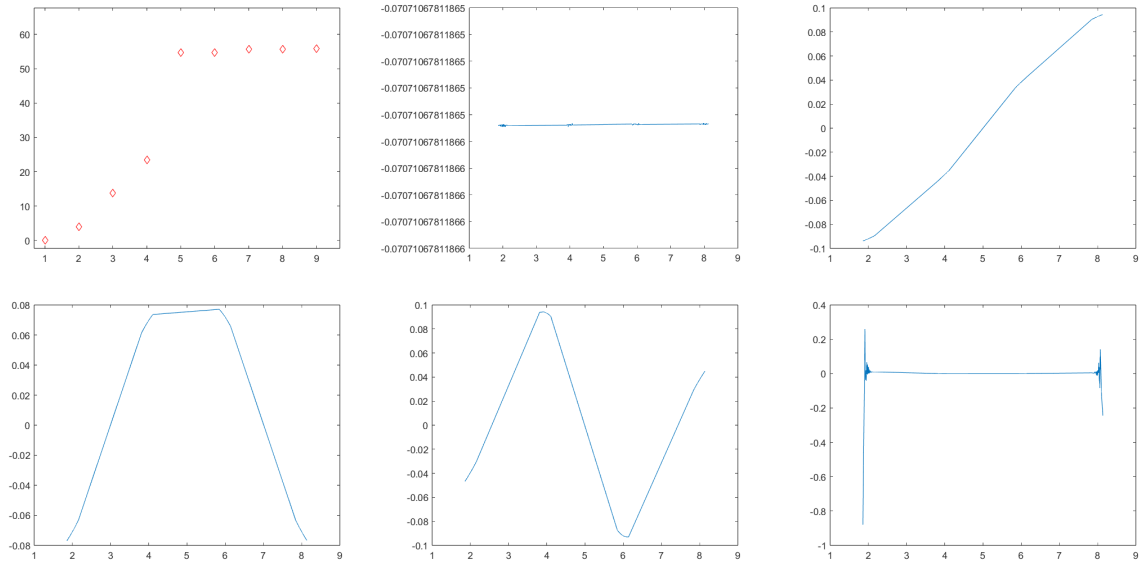
Di seguito è riportato l'istogramma relativo ai dati utilizzati e i primi 10 autovalori e 5 autovettori di L , ottenuti utilizzando gli algoritmi sopra citati.



Algoritmo del 10-Nearest Neighbor Graph:



Algoritmo del Fully Connected Graph:



Osserviamo che nel primo caso ci sono 4 autovalori nulli, dunque gli autovettori corrispondenti individuano 4 cluster, in accordo con il fatto che i dati sono stati presi da 4 gaussiane diverse. Nel secondo caso invece c'è un unico autovalore nullo: ciò deriva dal fatto che il grafo è completamente connesso.

4 Strategia di clustering

Utilizzando l'algoritmo kmeans di Matlab e le funzioni sopra riportate, possiamo eseguire l'algoritmo **Unnormalized Spectral Clustering**. L'obiettivo è quello di ottenere una rappresentazione dei dati tale che i punti di colore diversi appartengano a cluster diversi. Sia x il vettore dei dati, ϵ la soglia per l'algoritmo ϵ -neighborhood, k il numero di vicini dell'algoritmo k -nearest_neighbor_graph, c il numero di cluster cercati e $num_eigenvectors$ il numero di autovettori della matrice Laplaciana.

Algorithm 5

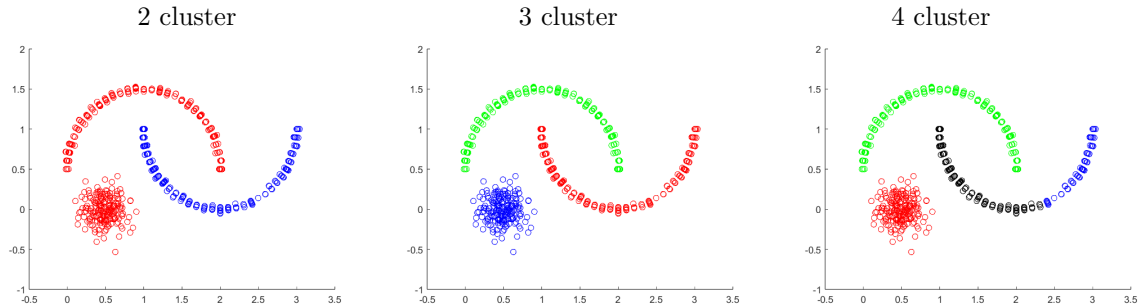
```

1: function SPECTRAL_CLUSTERING( $x, \epsilon, k, c, num\_eigenvectors$ )
2:    $[L, \sim, \sim] \leftarrow \epsilon\_neighborhood(x, \epsilon)$ 
3:      $\triangleright$  In alternativa:  $k\_nearest\_neighbor\_graph(x, k)$  o  $fully\_connected\_graph(x)$ 
4:    $[U, \sim] \leftarrow inverse\_subspace\_iteration(L, num\_eigenvectors, 100, -1)$ 
5:    $y \leftarrow kmeans(U(:, 1:c), c)$ 
6:    $clust \leftarrow cell(1, c)$ 
7:   for  $i = 1$  to  $length(x)$  do
8:     for  $j = 1$  to  $c$  do
9:       if  $y(i) == j$  then
10:         $z \leftarrow i$ 
11:         $clust\{j\} \leftarrow [clust\{j\}, z]$ 
12:       end if
13:     end for
14:   end for
15:    $col \leftarrow \{'ro', 'bo', 'go', 'ko'\}$ 
16:   for  $c$  do
17:      $a \leftarrow find(y == i)$ 
18:     hold on
19:     for  $j = 1$  to  $size(a, 1)$  do
20:        $plot(x(a(j)), col\{i\})$ 
21:     end for
22:   end for
23:   return  $clust$ 
24: end function

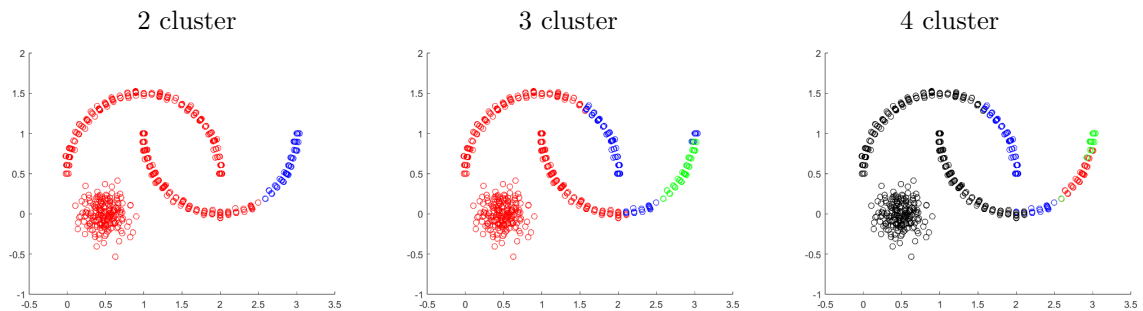
```

Siamo pronti per testare la strategia di clustering descritta su un file di dati in \mathbb{R}^2 al variare di alcuni parametri. Varieranno: l'algoritmo per la costruzione del grafo utilizzato; la funzione di somiglianza (distanza euclidea e densità gaussiana standard); numero di cluster cercati (2, 3, 4); numero di autovettori della matrice Laplaciana.

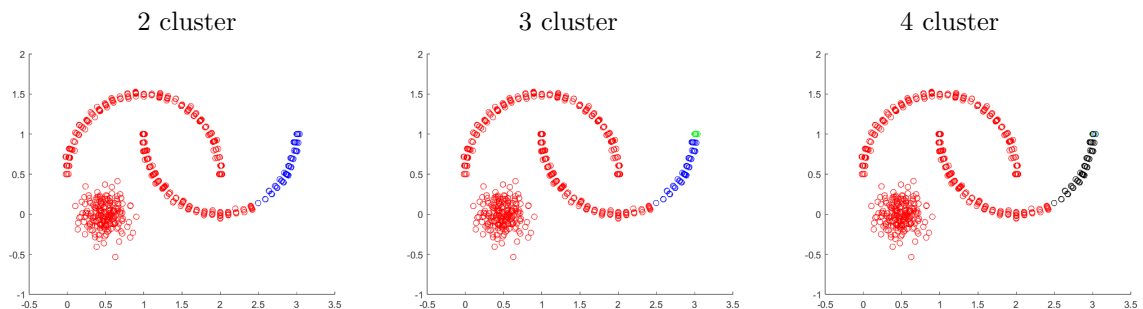
Caso 1: Algoritmo: ϵ -neighborhood ($\epsilon = 0.25$); Funzione di somiglianza: distanza euclidea; numero di autovettori della matrice Laplaciana pari al numero di cluster cercati.



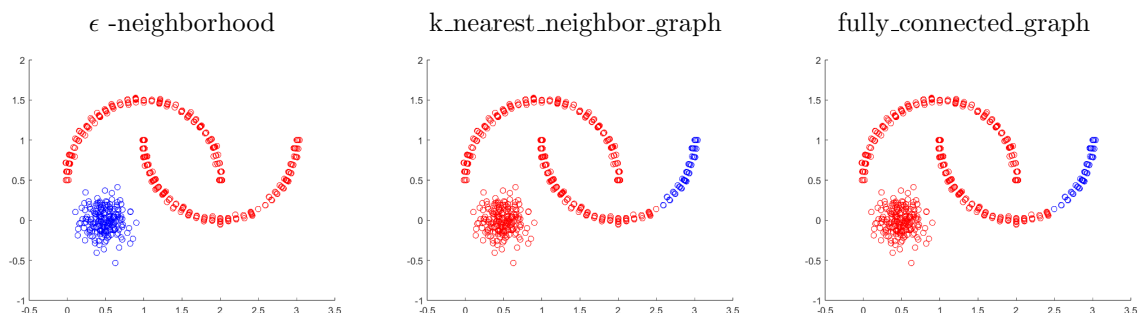
Caso 2: Algoritmo: k_nearest_neighbor_graph (k=12); Funzione di somiglianza: densità gaussiana; numero di autovettori della matrice Laplaciana pari al numero di cluster cercati.



Caso 3: Algoritmo: fully_connected_graph; Funzione di somiglianza: densità gaussiana; Numero di autovettori della matrice Laplaciana pari al numero di cluster cercati.



Caso 4: Algoritmi: ϵ -neighborhood ($\epsilon = 0.25$), k_nearest_neighbor_graph (k=12) e fully_connected_graph; Funzione di somiglianza: densità gaussiana; Caso del vettore di Fiedler.



5 Conclusioni

Dai risultati ottenuti è evidente che la strategia di clustering migliore è quella del caso 1.3. In generale, l'algoritmo ϵ -neighborhood funziona meglio rispetto agli altri. Notiamo infatti che:

- utilizzando l'algoritmo `k_nearest_neighbor_graph`, i dati vengono suddivisi in maniera casuale all'aumentare del numero di cluster;
- utilizzando l'algoritmo `fully_connected_graph` e aumentando il numero di autovettori, vengono individuati 2 cluster principali e i restanti cluster includono pochi dati.

Infine, nel caso 4 abbiamo analizzato il metodo del vettore di Fiedler che funziona bene per ogni grafo considerato.

6 Bibliografia

- [1] Luxburg. *A Tutorial on Spectral Clustering*