



A.D. 1308
unipg
DIPARTIMENTO
DI INGEGNERIA

Tesina di progetto
Signal Processing and Optimization for Big Data

Corso di Laurea Magistrale in Ingegneria Informatica e Robotica

Curriculum Data Science

A.A. 2024-2025

DIPARTIMENTO DI INGEGNERIA

docente

Prof. Paolo BANELLI

**IMPLEMENTAZIONE ALGORITMO PER IL CLUSTERING
SPETTRALE DI UN GRAFO**

studente

363551 **Andrea Tomassoni** andrea.tomassoni@studenti.unipg.it

0. Indice

1	Introduzione	2
2	Analisi Teorica del problema	3
2.1	Eigengap Heuristic	3
2.2	Spectral Graph Clustering	3
2.2.1	k-Cut	3
2.2.2	RatioCut	4
2.2.3	NormalizedCut	5
3	Implementazione	7
3.1	Eigengap Heuristic	7
3.1.1	Algoritmo	7
3.1.2	Codice Matlab	7
3.2	Spectral Normalized k-Cut	8
3.2.1	Codice Matlab	8
4	Dataset	9
4.1	Dati Sintetici	9
4.2	Dataset GARR	10
5	Prova Sperimentale	11
5.1	Rete Sintetica	11
5.2	Rete GARR	14

1. Introduzione

Lo scopo di questa tesina è l'implementazione di un'algoritmo per il clustering spettrale di un grafo. L'idea di fondo è quella di trovare uno strumento che tramite la clusterizzazione del grafo permetta di analizzare una rete internet per andare ad identificare eventuali porzioni della rete poco connesse.

L'algoritmo implementato è composto da due parti: nella prima parte si va ad identificare il numero di cluster ottimale per il grafo analizzato tramite un euristica (Eigengap Heuristic). Mentre in una seconda fase si sfrutta il dato precedente ottenuto per andare a clusterizzare il grafo.

L'algoritmo è stato implementato in Matlab ed è stato testato sia su dati sintetici, generati tramite un apposito script, sia su un grafo che rappresenta la rete del Consorzio GARR.

2. Analisi Teorica del problema

2.1 Eigengap Heuristic

Questa euristica viene utilizzata per calcolare il numero ottimale di cluster in cui dividere il grafo. Si basa sull'idea che i primi k autovalori del Laplaciano di un grafo siano piccoli, mentre al $k+1$ esimo autovalore si verifichi un salto.

Nello specifico, il procedimento è il seguente:

- calcolare e ordinare gli autovalori del Laplaciano

$$\lambda_1 \leq \lambda_2 \leq \lambda_3 \leq \dots \leq \lambda_n$$

- calcolare l'ampiezza di tutti i gap tra gli autovalori.

$$gap_i = \lambda_{i+1} - \lambda_i \quad (2.1)$$

- trovare il gap massimo

$$k = \arg \max_i (gap_i) \quad (2.2)$$

2.2 Spectral Graph Clustering

L'obiettivo di questo procedimento è quello di dividere il grafo in k cluster tramite minimizzazione del problema di ottimizzazione k -cut.

2.2.1 k-Cut

Il problema è formalmente definito come segue.

Posta la dimensione totale del taglio (*total cut size*) come

$$C(A_1, A_2, \dots, A_k) = \frac{1}{2} \sum_{i=1}^k W(A, \bar{A}_i) \quad (2.3)$$

si vuole andare a minimizzare la dimensione del taglio, ovvero

$$\min_{A_1, \dots, A_k} C(A_1, \dots, A_k) \quad (2.4)$$

2.2.2 RatioCut

Per prevenire soluzioni che vanno ad associare ad un cluster un singolo nodo, è possibile introdurre il *RatioCut* R , il quale va a penalizzare i cluster a bassa cardinalità.

$$R(A_1, \dots, A_k) = \frac{1}{2} \sum_{i=1}^k \frac{W(A_i, \bar{A}_i)}{|A_i|} \quad (2.5)$$

Nota

Il problema di minimizzazione di R è di tipo *NP-Hard*, ma si può ricorrere ad una versione rilassata del problema ponendo:

$$R(A_1, \dots, A_k) = \frac{1}{2} \sum_{j=1}^k \underline{h}_j^\top L \underline{h}_j \quad (2.6)$$

con

$$(\underline{h}_j)_i = \begin{cases} \frac{1}{|A_j|^{1/2}}, & \text{if } v_i \in A_j \\ 0, & \text{otherwise} \end{cases} \quad (2.7)$$

$$i = 1, \dots, n, \quad j = 1, \dots, k$$

Il problema precedente si può riscrivere come

$$\min_{\{A_1, \dots, A_k\}} R(A_1, \dots, A_k) = \frac{1}{2} \text{trace} (H^\top L H) \quad (2.8)$$

con

$$H = [\underline{h}_1, \dots, \underline{h}_k]$$

Il problema rilassato da risolvere diventa:

$$\begin{aligned} \arg \min_{H=\{h_1, \dots, h_n\}} \quad & \frac{1}{2} \text{trace} \{H^\top L H\} = \frac{1}{2} \sum_{i=1}^k \underline{h}_i^\top L \underline{h}_i \\ \text{s.t.} \quad & H^\top H = I \end{aligned} \quad (2.9)$$

Tramite il teorema di *Rayleigh-Ritz* è possibile dimostrare che la soluzione di questo problema sono i primi k autovettori del Laplaciano L .

La soluzione ottenuto va ora riportata ad una versione intera, in modo da assegnare ogni nodo ad un singolo cluster. Questo è possibile tramite l'utilizzo dell'algoritmo *k-means* sulle componenti spettrali. Nello specifico, posto:

$$\tilde{H} = [\underline{u}_1, \underline{u}_2, \dots, \underline{u}_k] = \begin{bmatrix} \underline{y}_1^\top \\ \vdots \\ \underline{y}_n^\top \end{bmatrix} \quad (2.10)$$

andando ad eseguire l'algoritmo *k-means* sui vettori $\{\underline{y}_i\}$ vengono identificati i centri (cluster) alla distanza minima.

2.2.3 NormalizedCut

Nel caso di grafi non regolari, in cui il grado di alcuni nodi può differire molto da quello degli altri, è preferibile utilizzare il *NormalizedCut*, una versione modificata del *RatioCut* in cui si va normalizzare tenendo conto non solo dei nodi dell'insieme preso in considerazione, ma anche del grado di questi.

Il *NormalizedCut* è definito come:

$$N(A_1, \dots, A_k) = \sum_{i=1}^k \frac{\text{cut}(A_i, \bar{A}_i)}{\text{vol}(A_i)} \quad (2.11)$$

dove

$$\text{vol}(A_i) = \sum_{v \in A_i} \deg(v)$$

Quindi, minimizzare N equivale a risolvere:

$$\begin{aligned} \arg \min_{H=\{h_1, \dots, h_n\}} \quad & \frac{1}{2} \text{trace} \{H^\top L_{sym} H\} \\ \text{s.t.} \quad & H^\top H = I \end{aligned} \quad (2.12)$$

dove

$$L_{sym} = I - D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$$

Il problema risulta quindi analogo al caso del *RatioCut* e di conseguenza anche in questo caso si procederà prima con una soluzione rilassata del problema attraverso i primi k autovettori del Laplaciano (L_{sym}) e successivamente si andrà ad eseguire l'algoritmo *k-means* sui vettori $\{\underline{y}_i\}$

Nota

Nel *NormalizedCut*, gli autovettori $[\underline{u}_1, \underline{u}_2, \dots, \underline{u}_k]$ vengono normalizzati come segue

$$\tilde{\mathbf{u}}_i = \frac{\mathbf{u}_i}{\|\mathbf{u}_i\|}$$

prima che venga eseguito *k-means*. Lo scopo di questo passaggio è quello di riportare ogni autovettore sulla sfera unitaria, in modo da confrontare solo la direzione del vettore e non la sua lunghezza. In questo modo si evita che i vettori la cui norma è maggiore dominino sugli altri, evitando così il rischio di creare cluster errati.

Il *NormalizedCut* quindi è preferibile nei casi reali in cui la distribuzione del grado dei nodi potrebbe non essere uniforme. Il *RatioCut* infatti non distigue se un nodo che viene tagliato è "importante" oppure no, con il rischio di creare cluster sbilanciati.

Nello specifico, *NormalizedCut* utilizza il Laplaciano L_{sym} invece di L in quanto il Laplaciano L nelle porzioni di grafo in cui la distribuzione del grado dei nodi è molto disomogenea tenderà a favorire tagli che separano nodi di grado basso anche se questi non rappresentano i veri cluster.

E' stato inoltre dimostrato¹ che, dal punto di vista della teoria della convergenza, il Laplaciano L non è consistente al crescere della mole di dati.

¹Von Luxburg, U. (2007): "A Tutorial on Spectral Clustering"

3. Implementazione

3.1 Eigengap Heuristic

3.1.1 Algoritmo

L'algoritmo utilizzato per il calcolo del numero di cluster ottimale prende in input la matrice di adiacenza A ed esegue i seguenti passaggi:

- Calcolo matrice dei Gradi D
- Calcolo della matrice $D^{-\frac{1}{2}}$
- Calcolo del Laplaciano L_{sym}
- Calcolo degli autovalori di L_{sym}
- Ordinamento degli autovalori calcolati
- Calcolo dei gap
- Selezione del gap e del relativo autovalore

In output viene restituito il k ottimo trovato tramite questa euristica.

3.1.2 Codice Matlab

```
function k_opt = k_eigengap_heuristic(A, max_k)
    D = diag(sum(A, 2));
    D_inv_sqrt = diag(1 ./ sqrt(diag(D) + eps));
    L_sym = eye(size(A)) - D_inv_sqrt * A * D_inv_sqrt;
    [~, Lambda_vals] = eigs(L_sym, max_k + 1, 'smallestreal');
    eigvals = sort(diag(Lambda_vals));
    gaps = diff(eigvals);
    [~, k_opt] = max(gaps);
end
```

Nota: Nel codice matlab è presente un secondo parametro in ingresso alla funzione che deve essere oggetto di tuning in base alla cardinalità del problema.

3.2 Spectral Normalized k-Cut

L'algoritmo per la clusterizzazione del grafo prende in input la matrice di adiacenza A del grafo e il numero di cluster k ottimale calcolato precedentemente. I passaggi sono i seguenti:

- Calcolo matrice dei Gradi D
- Calcolo della matrice $D^{-\frac{1}{2}}$
- Calcolo del Laplaciano L_{sym}
- Calcolo dei k autovettori più piccoli di L_{sym}
- Normalizzazione degli autovettori $\{\underline{u}_i\}$ (righe della matrice U) a norma 1
- Calcolo dei cluster tramite l'algoritmo k-means che prende in input i vettori $\{\underline{y}_i\}$ (colonne della matrice U)

In output vengono restuiti i k cluster calcolati.

3.2.1 Codice Matlab

```
function clusters = spectral_k_norm_cut(A, k)
    D = diag(sum(A, 2));
    D_inv_sqrt = diag(1 ./ sqrt(diag(D) + eps));
    L_sym = eye(size(A)) - D_inv_sqrt * A * D_inv_sqrt;
    [U, ~] = eigs(L_sym, k, 'smallestreal');
    U_norm = U ./ vecnorm(U, 2, 2);

    labels = kmeans(U_norm, k, 'Replicates', 10);
    clusters = cell(1, k);
    for i = 1:k
        clusters{i} = find(labels == i);
    end
end
```

Nota

Il codice completo dell'implementazione è disponibile alla repository Github del progetto

4. Dataset

L'algoritmo implementato è stato testato su due diversi dataset: uno sintetico, utilizzato principalmente in fase di sviluppo e debug del codice. Il secondo dataset invece rappresenta una rete internet reale.

4.1 Dati Sintetici

Il dataset sintetico è stato costruito dinamicamente tramite un'apposita funzione Matlab¹ che permette di generare un grafo indicando i seguenti parametri:

- Numero di cluster,
- Numero di nodi per cluster,
- Probabilità di connessione dei nodi intra-cluster (2 nodi appartenenti allo stesso cluster),
- Probabilità di connessione dei nodi inter-cluster (2 nodi appartenenti a cluster differenti),
- Pesi degli archi intra-cluster,
- Pesi degli archi inter-cluster.

Questa funzione permette di definire un grafo come l'unione di più cluster. L'idea di fondo era quella di riprodurre una rete internet e quindi diversi AS (Autonomous System) collegati tra loro.

Nota

Questa tipologia di struttura del grafo permette, tarando i parametri, di andare a controllare la corretta esecuzione e funzionamento dell'algoritmo implementato.

¹graphN.m - disponibile nella repository del progetto

4.2 Dataset GARR

Il secondo dataset utilizzato rappresenta l'infrastruttura di rete del consorzio GARR, il quale si occupa di fornire connessioni a banda ultralarga a università, centri di ricerca e musei. La scelta è ricaduta su questa rete in quanto rende accessibile a tutti le informazioni riguardanti l'infrastruttura, la posizione dei nodi e le specifiche di connettività delle rete². E' stato quindi possibile scaricare l'elenco dei nodi della rete, le connessioni tra i nodi e anche la banda di queste connessioni.

I dati estratti ottenuto sono stati manualmente corretti ed elaborati in modo da renderli comprensibili da matlab, come nell'esempio sottostante:

Nodes.txt

```
r11.fi04 [1625,1775]
r11.fi05 [1675,1825]
...
r11.pg00 [525,1375]
r11.rm02 [625,1325]
```

Links.txt

```
r11.fi04 r11.fi02 10000
r11.fi05 r11.fi02 10000
...
r11.pg00 r12.pg00 20000
r11.pg00 r11.rm02 10000
```

Nota

Per quanto riguarda questo dataset, ne sono state testate due versioni differenti: la prima contenente solamente i nodi della rete GARR; mentre una seconda è comprensiva anche degli Endpoint³ esterni alla rete GARR.

²GARR Weather Map

³Es. Google, Microsoft GEANT, ...

5. Prova Sperimentale

Di seguito vengono riportati i risultati di un'esecuzione dell'algoritmo implementato, sia nel caso di rete generata da dataset sintetico che della rete GARR.

Nota

Nelle seguenti visualizzazioni, le etichette degli archi che riportano i pesi sono state rimosse per una migliore visibilità della figura statica.

5.1 Rete Sintetica

Per la generazione della rete sintetica sono stati utilizzati i seguenti parametri:

Parametri

```
% Numero di cluster
clusters_num = 3;

% Numero di nodi per singolo cluster
nodes = [15, 25, 20];

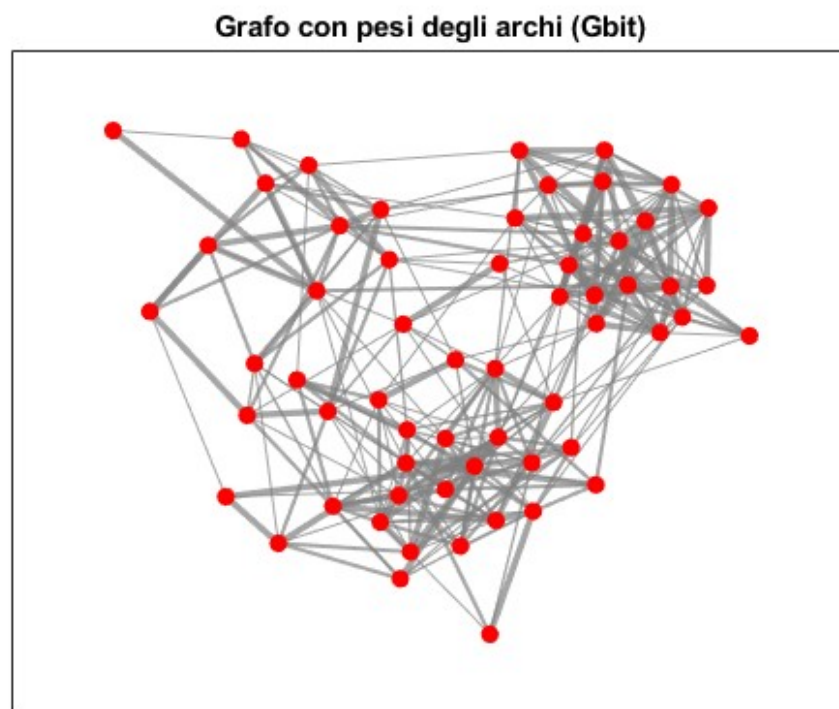
% Probabilità di connessione dei nodi intra-cluster
p_intra = [0.4, 0.3, 0.5];

% Probabilità di connessione dei nodi inter-cluster
p_inter = [0 0.05 0.05; 0.05 0 0.05; 0.05 0.05 0];

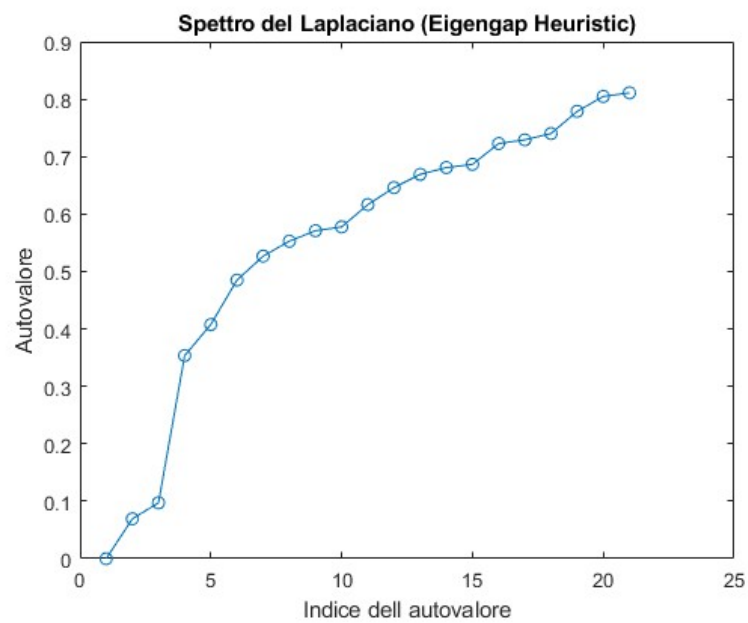
% Pesi degli archi intra-cluster
weights_intra = {[10, 25, 50, 100]};

% Pesi degli archi inter-cluster
weights_inter = {[1, 5, 10, 25, 50]}
```

Ed è stata ottenuta la rete in figura

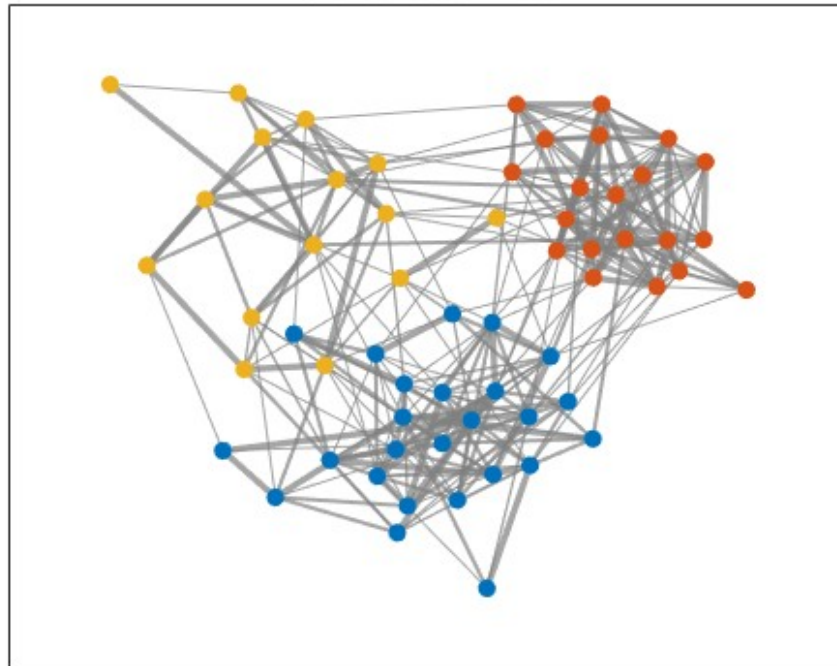


Successivamente sono state applicate la Eigengap Heuristic e la funzione di clusterizzazione.

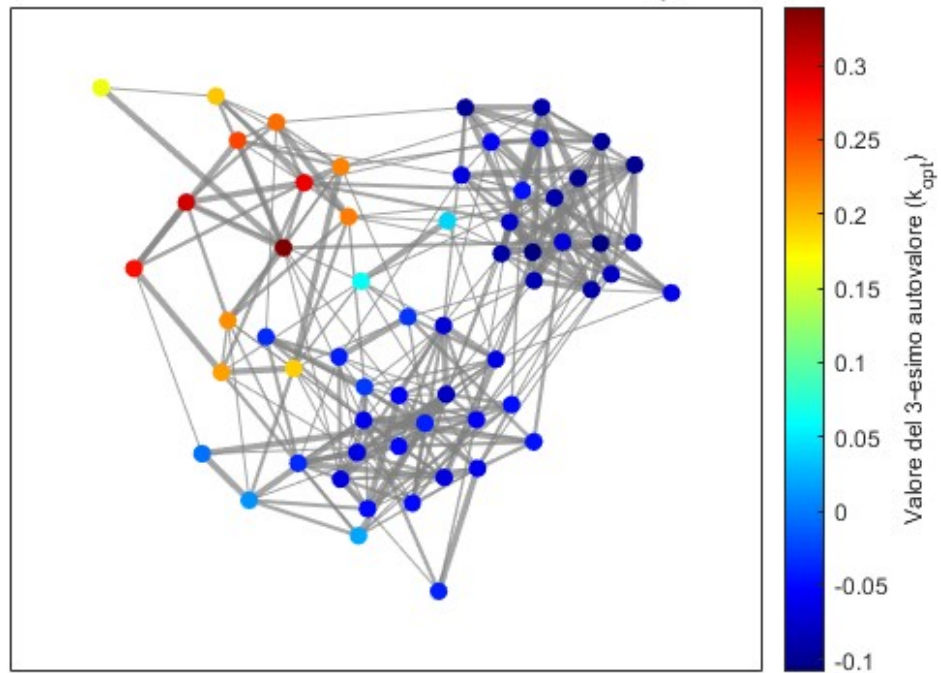


Numero di cluster: 3

Grafo clusterizzato

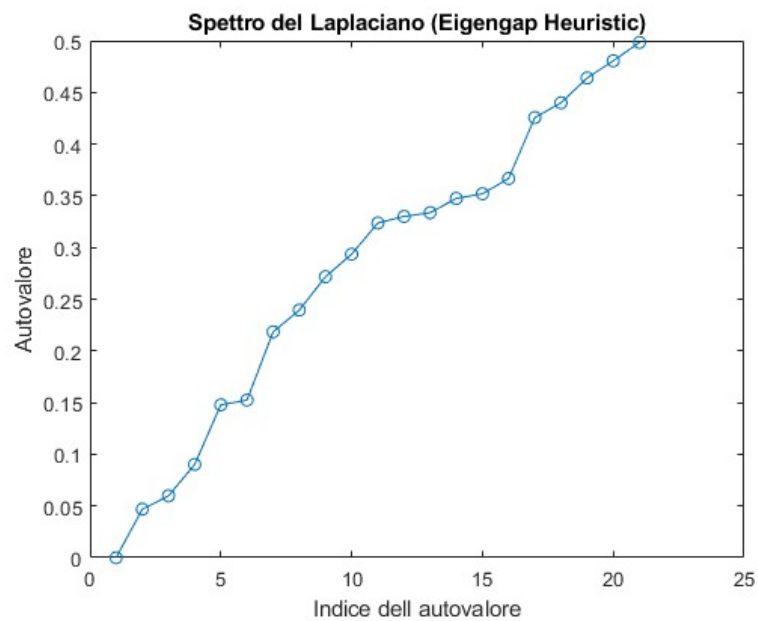
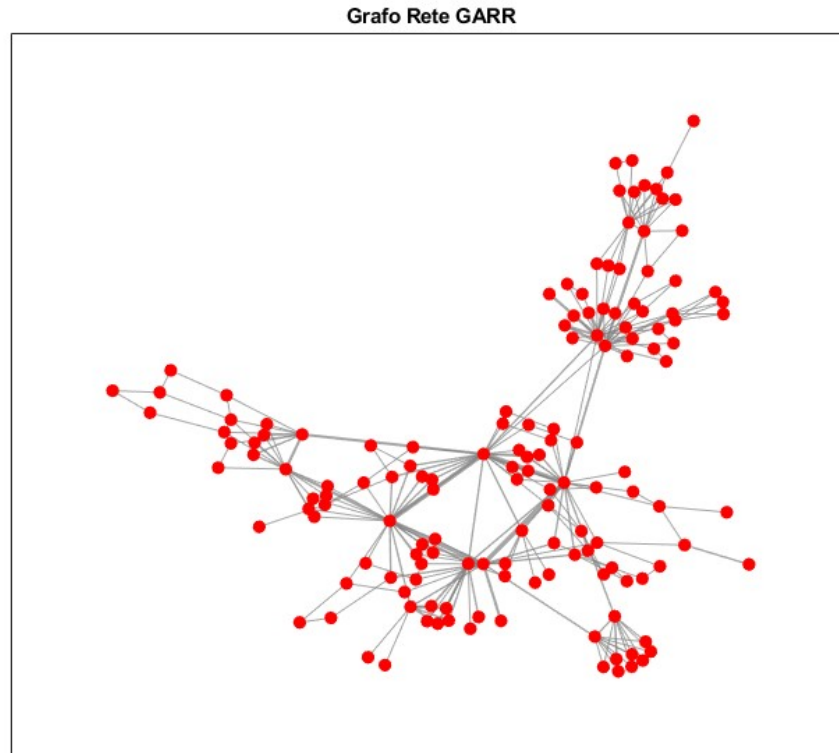


Grafo con valori del 3-esimo autovalore (k_{opt})



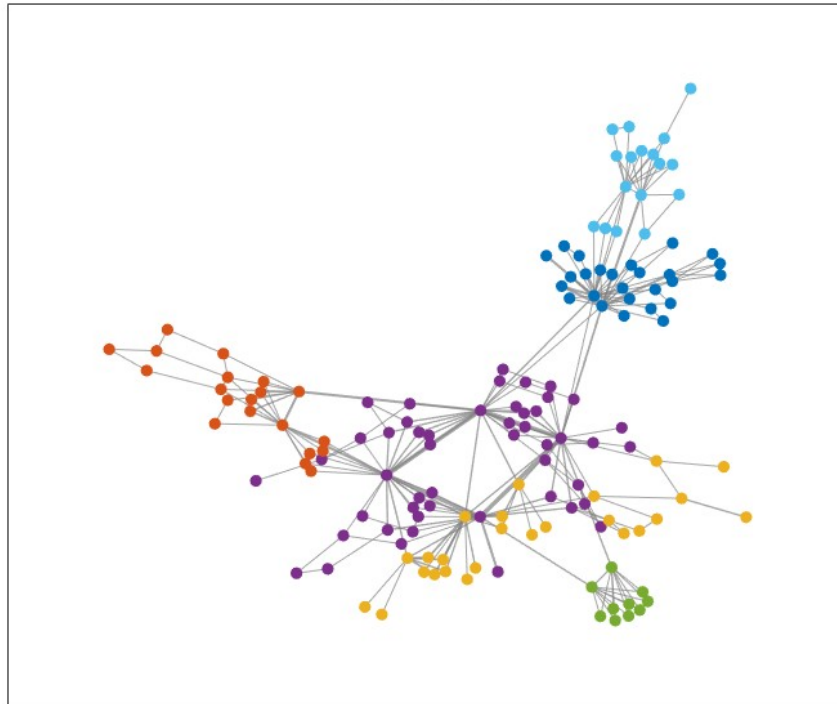
5.2 Rete GARR

Per quanto riguarda il dataset della rete GARR verrà mostrata l'esecuzione dell'algoritmo utilizzando il dataset che comprende anche gli Endpoint esterni.



Numero di cluster: 6

Grafo clusterizzato



Grafo con valori del 6-esimo autovalore (k_{opt})

