

# Spectral Clustering Analysis



Network Statistics for Data Science (2AMS30)

Andrea Mangrella  
Maarten van Sluijs  
Roëlle Bänffer

# Community detection

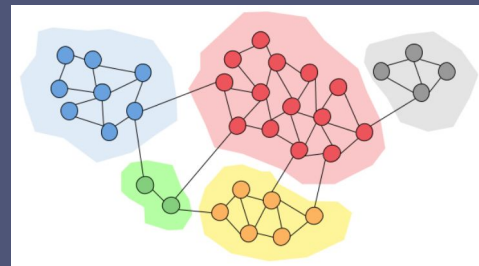
---

- Communities: groups of vertices having high probability of being connected to each other
- Goal: predict community labels based on the edges
- Internal and external degree
- Artificial benchmarks
- Different types of uses: Predicting behavior, Anomaly detection, Recommendation systems, Improving tasks efficiency

Let  $G = (V, E)$  be a graph with  $m$  communities  $C_1, \dots, C_m \subset V$ .

$$d_i^{int}(C) = \sum_{j \in C} 1[(i, j) \in E]$$

$$d_i^{ext}(C) = d_i - d_i^{int}$$



# Technical details

---

A set of objects  $\{x_1, x_2, \dots, x_n\}$

Pairwise similarity  $s_{i,j}$

- Generally Gaussian:  $s(x_i, x_j) = \exp\left(-\|x_i - x_j\|^2 / (2\sigma^2)\right)$

Construct a similarity graph  $G = (V, E)$

- $\epsilon$ -neighbourhood graph
- K-nearest neighbour graph
- Mutual k-nearest neighbour

Properties of  $G$ :

- Weighted adjacency matrix  $W$
- Degree matrix  $D$  (diagonal)

# Technical details

---

A graph's Laplacian matrix is needed for clustering:

- Unnormalized Laplacian:

$$L = D - W$$

- Normalized Laplacians:

Symmetric laplacian

$$L_{\text{sym}} = I - D^{-0.5} W D^{0.5}$$

Random walk laplacian

$$L_{\text{rw}} = I - D^{-1} W$$

# of connected components = multiplicity of eigenvalue 0

# Unnormalized Clustering algorithm

---

Given a similarity matrix  $S = \mathbb{R}^{n \times n}$  and  $k$  clusters do:

1. Construct one of the similarity graphs
2. Compute laplacian  $L$
3. Compute first  $k$  eigenvectors  $u_1, \dots, u_k$  from  $L$
4. Construct  $U$  as:
5. Represent point  $y_i$  as row  $i$  of  $U$
6. Cluster points  $y$  using kmeans

$$U = \begin{pmatrix} u_{11} & . & . & . & u_{k1} \\ . & & & & . \\ . & & & & . \\ u_{1i} & . & . & . & u_{ki} \\ . & & & & . \\ . & & & & . \\ u_{1n} & . & . & . & u_{kn} \end{pmatrix}$$

# Normalized Clustering algorithm

---

Given a similarity matrix  $S = \mathbb{R}^{n \times n}$  and  $k$  clusters do:

- |   |   |
|---|---|
| 1. Construct one of the similarity graphs   | 1. Construct one of the similarity graphs                                 |
| 2. Compute laplacian $L$  | 2. Compute laplacian $L$  |
| 3. Compute first $k$ eigenvectors $u_1, \dots, u_k$ from the problem $L^*u = \lambda^*D^*u$ | 3. Compute first $k$ eigenvectors $u_1, \dots, u_k$ from $L_{\text{sym}}$ |
| 4. Construct $U$ as before  | 4. Form $U$ as before   |
| 5. Represent point $y_i$ as row $i$ of $U$  | 5. Form matrix $T$ by normalizing the columns of $U$                      |
| 6. Cluster points $y$ using kmeans  | 6. Represent point $y_i$ as row $i$ of $T$                                |
|   | 7. Cluster points $y$ using kmeans  |

# Pros/Cons of Spectral Clustering:

---

## PROS:

In this presentation I decided to focus on proving empirically that Spectral Clustering is heavily influenced by **Sparsity**

- Easy to compute thanks to advanced linear algebra libraries
- Often outperform standard methods
- Simple to implement



## CONS:

- No proof we are going to compare the results of graphs
- Choice of clustering algorithm for eigenvectors can introduce issues related to initialization and local optima.
- Exploiting the **Stochastic Block Method** random graph properties to generate **dense** and **sparse** graphs with deducible community **ground truths** for all nodes!

# Metrics used to analyze the Results:

---

## Mutual Information:

For both scores we are mainly going to refer to their adjusted version, as it's more reliable in a graph setting. Mutual information is how (on average) the change we see in one variable is related to the change in another.

$$\text{Mutual Information} = \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \left[ \frac{p(x, y)}{p(x)p(y)} \right]$$

## Rand Index:

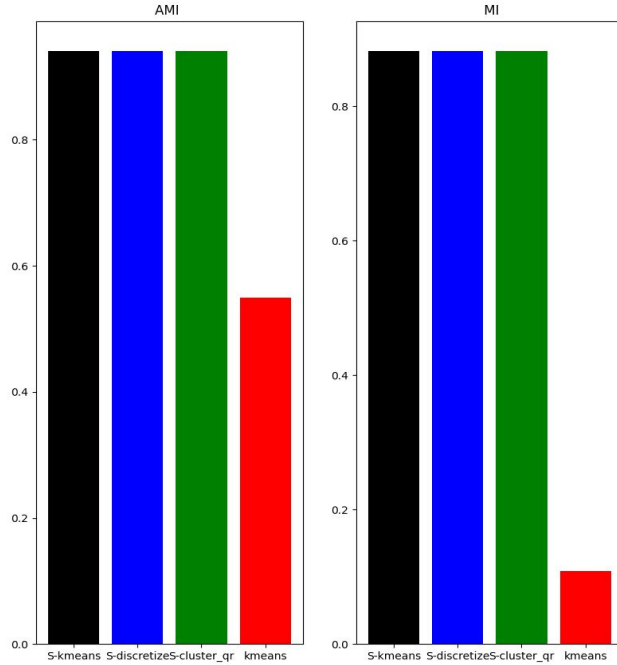
The Rand Index computes a similarity measure between two clusterings by considering all pairs of samples and counting pairs that are correctly assigned (in respect to the ground truths).

$$\text{Rand Index} = (\text{number of agreeing pairs}) / (\text{number of pairs})$$

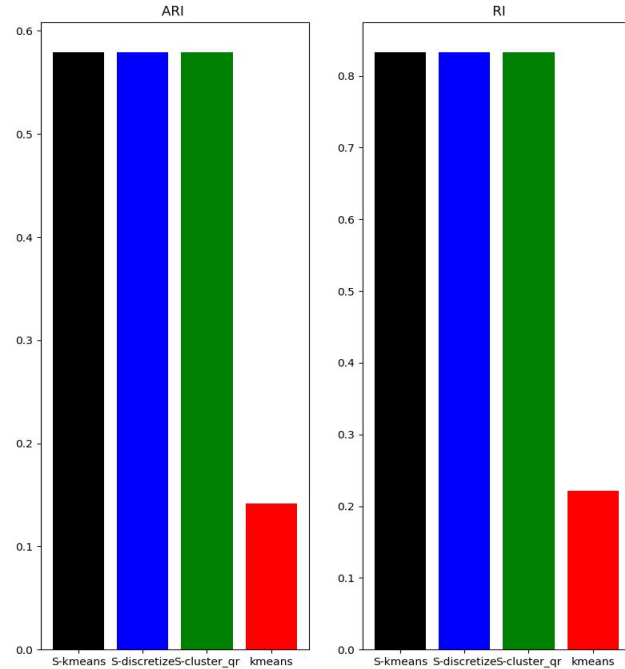


# Karate Club Analysis:

Karate Club  
AMI and MI

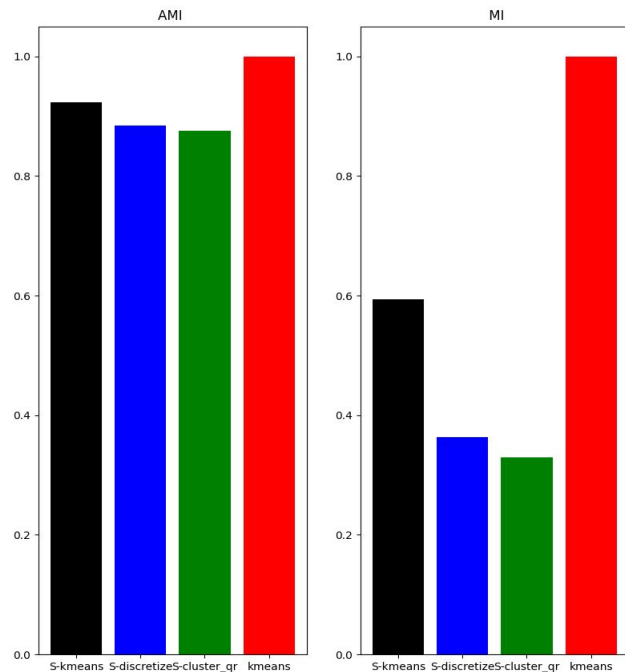


Karate Club  
ARI and RI

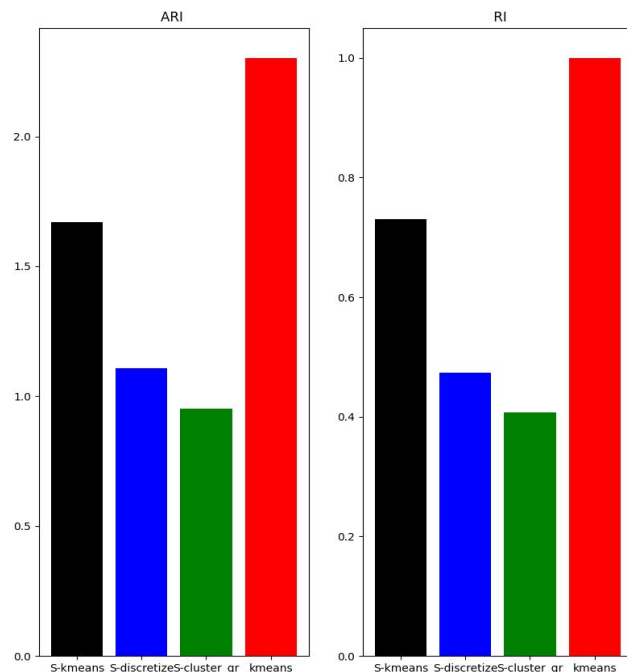


# SBM Analysis (n=10):

SBM Sparse with n=10  
AMI and MI

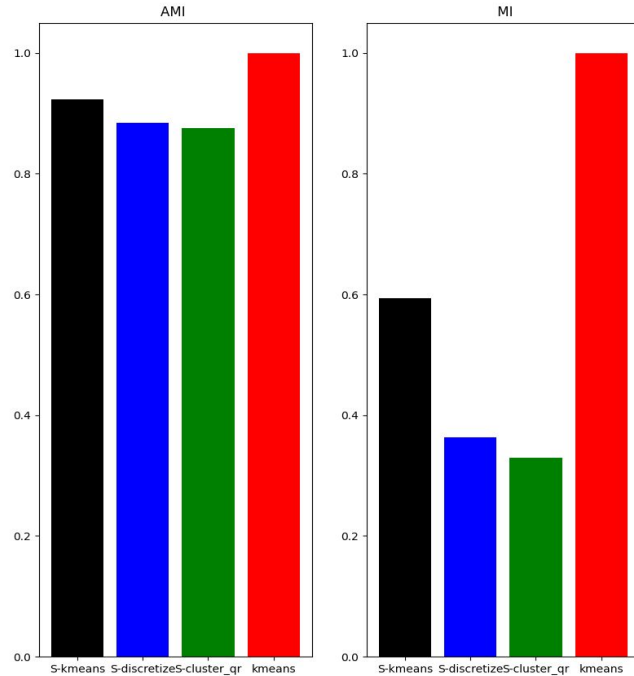


SBM Sparse with n=10  
ARI and RI

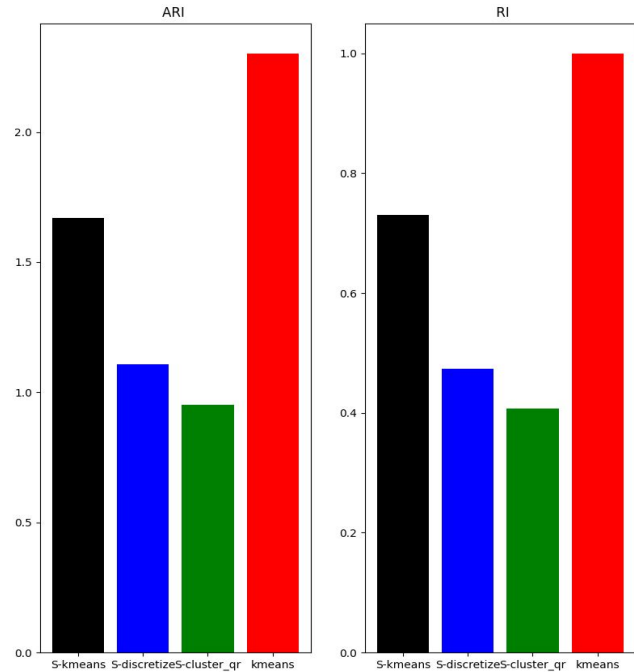


# SBM Analysis (n=1000):

SBM Sparse with n=1000  
AMI and MI



SBM Sparse with n=1000  
ARI and RI



# Bibliography:

---

- ❑ Fortunato, Santo, and Darko Hric. "**Community detection in networks: A user guide.**" Physics reports 659 (2016): 1-44. Section 4.3
- ❑ Von Luxburg, Ulrike. "**A tutorial on spectral clustering.**" Statistics and computing 17 (2007): 395-416.
- ❑ Karataş, A., & Şahin, S. (2018). Application Areas of Community Detection: A Review.
- ❑ [Code](#) for plotting the results shown in the presentation