

Characterizing measures for the assessment of cluster analysis and community detection

Nejat Arinik

Vincent Labatut

Rosa Figueiredo

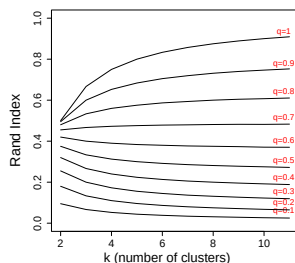
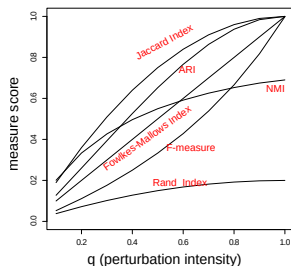
Laboratoire Informatique d'Avignon, University of Avignon, France.

MARAMI

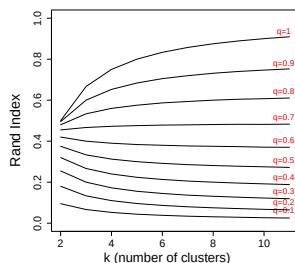
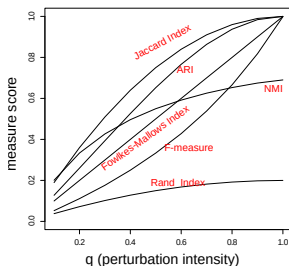
Oct 15, 2020

- **Context:** comparing two non-overlapping partitions through an external measure → cluster analysis, graph partitioning
 - ground-truth vs. estimated partition, 2 estimated partitions

- **Context:** comparing two non-overlapping partitions through an external measure → cluster analysis, graph partitioning
 - ground-truth vs. estimated partition, 2 estimated partitions
- **Issues:**
 - profusion of available measures → trend to follow popular measures
 - lack of comprehensive comparison



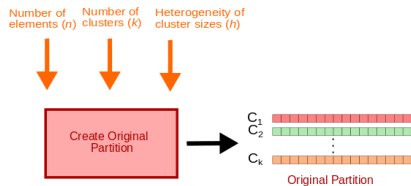
- **Context:** comparing two non-overlapping partitions through an external measure → cluster analysis, graph partitioning
 - ground-truth vs. estimated partition, 2 estimated partitions
- **Issues:**
 - profusion of available measures → trend to follow popular measures
 - lack of comprehensive comparison



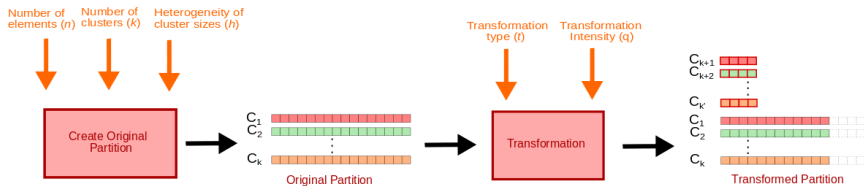
- **Our solution:** a new framework of evaluation

- 1 Our framework
 - Characterizing of the measures
 - Regression Analysis
- 2 Experiments
- 3 Results
- 4 Practical case
- 5 Conclusion

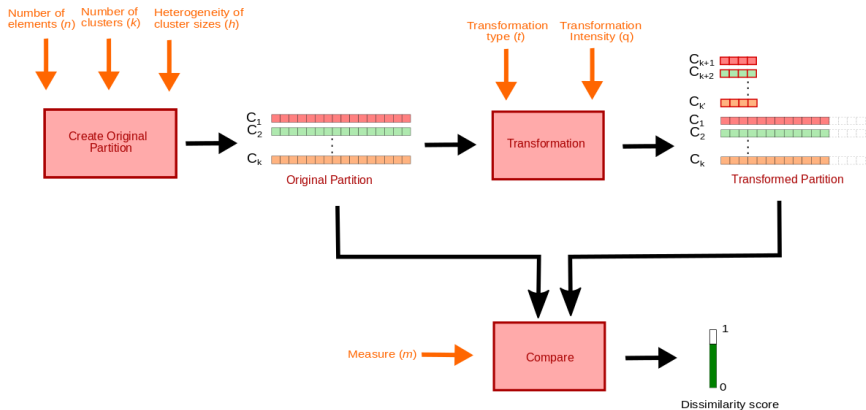
Our framework: Characterizing of the measures



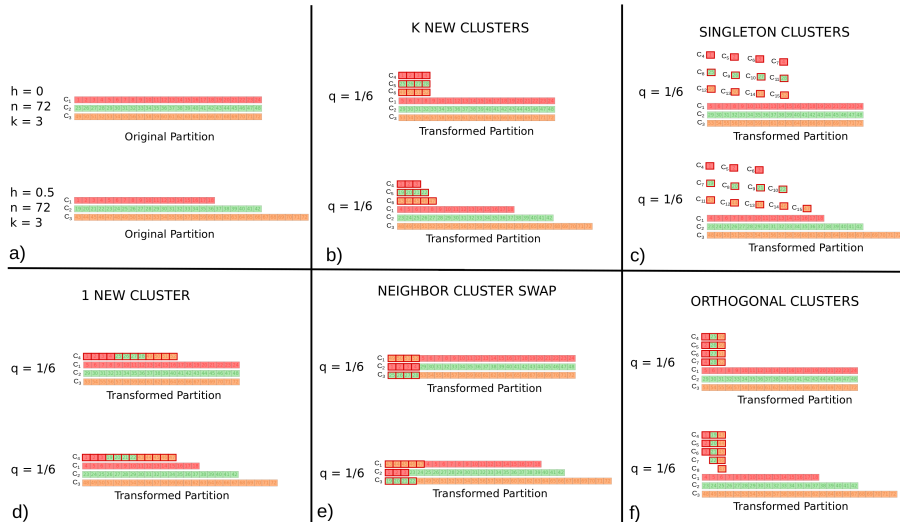
Our framework: Characterizing of the measures



Our framework: Characterizing of the measures



Our framework: proposed deterministic transformations



Our framework: Regression Analysis

Our multiple linear regression model:

$$y = \sum_i \sum_j \left(\beta_{0ij} t_i m_j \right. \\ \left. + \beta_{1ij} n t_i m_j + \beta_{2ij} k t_i m_j + \beta_{3ij} p t_i m_j + \beta_{4ij} h t_i m_j \right. \\ \left. + \beta_{5ij} n k t_i m_j + \beta_{6ij} n h t_i m_j + \beta_{7ij} n p t_i m_j + \beta_{8ij} k h t_i m_j + \beta_{9ij} k p t_i m_j + \beta_{10ij} h p t_i m_j \right) \\ + \epsilon, \tag{1}$$

- $\beta_{.ij}$: regression coefficients
- t_i ($1 \leq i \leq T$) and m_j ($1 \leq j \leq M$): binary dummy variables, where T = number of transformations and M = number of measures
- ϵ : common error

Our framework: Regression Analysis

Our multiple linear regression model:

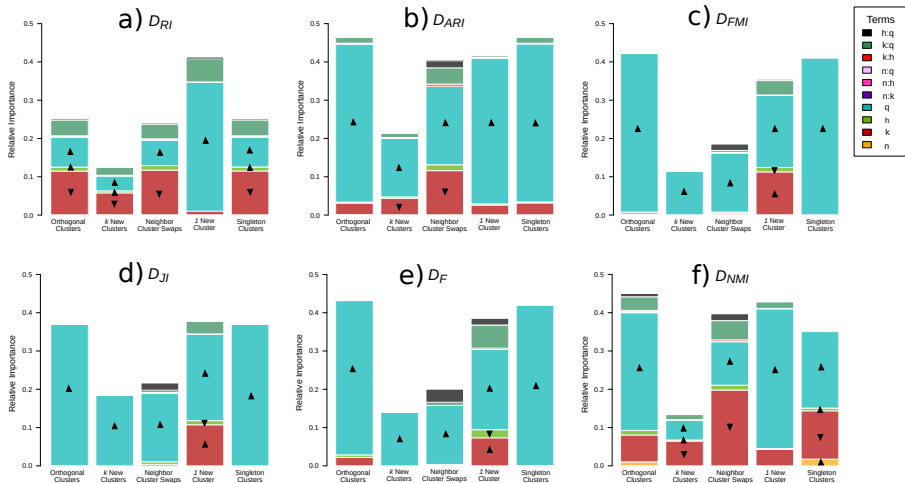
$$y = \sum_i \sum_j \left(\beta_{0ij} t_i m_j \right. \\ \left. + \beta_{1ij} n t_i m_j + \beta_{2ij} k t_i m_j + \beta_{3ij} p t_i m_j + \beta_{4ij} h t_i m_j \right. \\ \left. + \beta_{5ij} n k t_i m_j + \beta_{6ij} n h t_i m_j + \beta_{7ij} n p t_i m_j + \beta_{8ij} k h t_i m_j + \beta_{9ij} k p t_i m_j + \beta_{10ij} h p t_i m_j \right) \\ + \epsilon, \tag{1}$$

- $\beta_{.ij}$: regression coefficients
- t_i ($1 \leq i \leq T$) and m_j ($1 \leq j \leq M$): binary dummy variables, where T = number of transformations and M = number of measures
- ϵ : common error

- **Relative importance analysis** → squared beta weights

- Our data: 50,000 pairs of partitions
 - measures (\mathbf{m}) = {Rand Index (D_{RI}), Adjusted Rand Index (D_{ARI}), Fowlkes-Mallows Index (D_{FMI}), Jaccard Index (D_{JI}), F-measure (D_F), Normalized Mutual Information (D_{NMI})}
 - transformations (\mathbf{t}) = {K New Cluster, Singleton Clusters, 1 New Cluster, Neighbor Cluster Swap, Orthogonal Clusters}
 - number of elements (\mathbf{n}) = 3240, 4320, ..., 12960
 - number of clusters (\mathbf{k}) = 2, 3, ..., 11
 - heterogeneity of cluster sizes (\mathbf{h}) = 0, 0.1, ..., 0.9
 - transformation intensity (\mathbf{q}) = 0.1, 0.2, ..., 1
- Regression assumptions
 - by design \rightarrow no collinearity between the quantitative variables
 - large dataset & central limit theorem \rightarrow no issue with the residuals
 - heteroscedasticity \rightarrow increase of the variance in y with parameter q

Results



Practical case: vote application in European Parliament

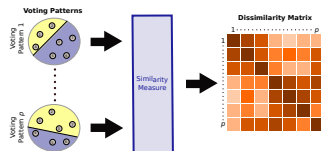
N. Arinik & R. Figueiredo & V. Labatut. Multiple Partitioning of Multiplex Signed Networks. Social Networks, 2020, 60, 83–102.

Requirements:

- up to 3 clusters
- n fixed

Expectations:

- detecting an extra cluster, or a missing one, is an important change \rightarrow the difference of k between the original and transformed partitions
- the effect of k should be stronger than that of h
- a dissimilarity score should decrease, when h increases



Practical case: vote application in European Parliament

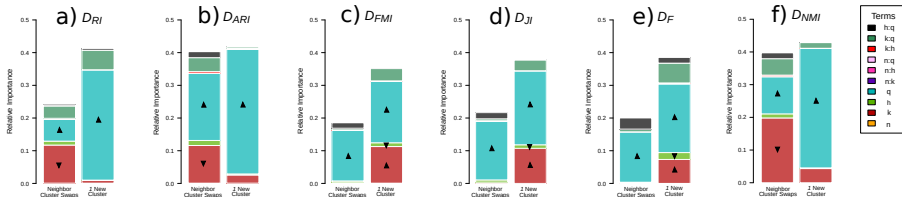
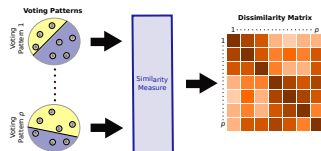
N. Arinik & R. Figueiredo & V. Labatut. Multiple Partitioning of Multiplex Signed Networks. Social Networks, 2020, 60, 83–102.

Requirements:

- up to 3 clusters
- n fixed

Expectations:

- detecting an extra cluster, or a missing one, is an important change \rightarrow the difference of k between the original and transformed partitions
- the effect of k should be stronger than that of h
- a dissimilarity score should decrease, when h increases



Practical case: vote application in European Parliament

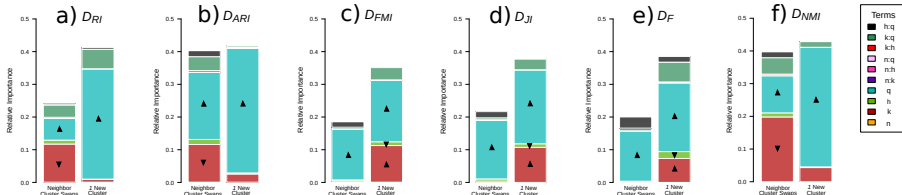
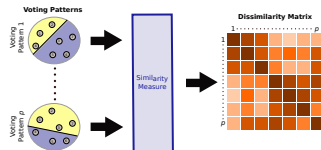
N. Arinik & R. Figueiredo & V. Labatut. Multiple Partitioning of Multiplex Signed Networks. Social Networks, 2020, 60, 83–102.

Requirements:

- up to 3 clusters
- n fixed

Expectations:

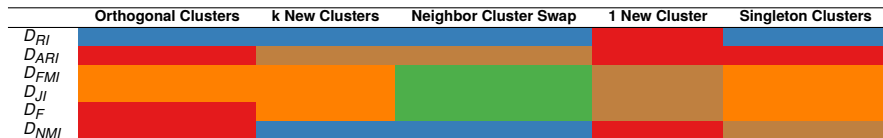
- detecting an extra cluster, or a missing one, is an important change \rightarrow the difference of k between the original and transformed partitions
- the effect of k should be stronger than that of h
- a dissimilarity score should decrease, when h increases



selected measure: D_F (F-measure)

Conclusion & Further research

- a new generic framework of evaluation
- ease of interpretation for the results
- typology of measures based on their performances



Conclusion & Further research

- a new generic framework of evaluation
- ease of interpretation for the results
- typology of measures based on their performances

	Orthogonal Clusters	k New Clusters	Neighbor Cluster Swap	1 New Cluster	Singleton Clusters
D_{RI}	Blue	Blue	Blue	Red	Blue
D_{ARI}	Red	Brown	Brown	Red	Red
D_{FMI}	Orange	Orange	Green	Brown	Orange
D_{JI}	Orange	Orange	Green	Brown	Orange
D_F	Red	Blue	Blue	Red	Brown
D_{NMI}	Red	Blue	Blue	Red	Brown

- evaluation with more measures and transformations
- designing the framework for graph similarity measures
- designing the framework for overlapping partitions

Thank you for your attention!

Contact Information:

Nejat ARINIK

nejat.arinik@univ-avignon.fr