

Relative Validity Criteria for Community Mining Algorithms

Reihaneh Rabbany*, Mansoreh Takaffoli*, Justin Fagnan*, Osmar R. Zaiane* and Ricardo J. G. B. Campello*†

*Computing Science Department, University of Alberta, Edmonton, Canada T6G-2E5

Email: {rabbanyk,takaffol,fagnan,zaiane,rcampell}@ualberta.ca

†Department of Computer Science, University of São Paulo, São Carlos, SP, Brazil, C.P. 668

Email: campello@icmc.usp.br

Abstract—Grouping data points is one of the fundamental tasks in data mining, which is commonly known as clustering if data points are described by attributes. When dealing with interrelated data that does not have any attributes and is represented in the form of nodes and their relationships, this task is also referred to as community mining. There has been a considerable number of approaches proposed in recent years for mining communities in a given network. But little work has been done on how to evaluate community mining results. The common practice is to use an agreement measure to compare the mining result against a ground truth, however, the ground truth is not known in most of the real world applications. In this paper, we investigate relative clustering quality measures defined for evaluation of clustering data points with attributes and propose proper adaptations to make them applicable in the context of social networks. Not only these relative criteria could be used as metrics for evaluating quality of the groupings but also they could be used as objectives for designing new community mining algorithms.

I. INTRODUCTION

Data Mining is the analysis of large scale data to discover meaningful patterns such as groups of data records (cluster analysis), unusual records (anomaly detection) or dependencies (association rule mining) which are crucial in a very broad range of applications. It is a multidisciplinary field that involves methods at the intersection of artificial intelligence, machine learning, statistics and database systems. The recent growing trend in the Data Mining field is the analysis of structured/interrelated data, motivated by the natural presence of relationships between data points in a variety of the present-day applications. The structures in these inter-related data are usually represented using networks, known as complex networks or information networks; examples are the hyperlink networks of web pages, citation or collaboration networks of scholars, biological networks of genes or proteins, trust and social networks of humans and much more.

All these networks exhibit common statistical properties, such as power law degree distribution, small-world phenomenon, relatively high transitivity, shrinking diameter, and densification power laws [1], [2]. Network clustering, a.k.a. community mining, is one of the principal tasks in the analysis of complex networks. Many community mining algorithms have been proposed in recent years (for a complete survey refer to [3]). These algorithms evolved very quickly from simple heuristic approaches to more sophisticated optimization based methods that are explicitly or implicitly trying to maximize

the goodness of the discovered communities. The broadly used explicit maximization objective is the modularity, first introduced by Newman et al. [1].

Although there have been many methods presented for detecting communities, very little work has been done on how to evaluate the results and validate these methods. The difficulties of evaluation is due to the fact that the interesting communities that have to be discovered are hidden in the structure of the network, thus, the true results is not known for comparison. Furthermore, there are no other means to measure the goodness of the discovered communities in a real network. We also do not have any large enough dataset with known communities, often called ground truth, to use as a benchmark to generally test and validate the algorithms. The common practice is to use synthetic benchmark networks and compare the discovered communities with the built-in ground truth. However, it is shown that the networks generated with the current benchmarks disagree with some of the characteristics of real networks. These facts motivate investigating a proper objective for evaluation of community mining results.

Defining an objective function to evaluate community mining is non-trivial. Aside from the subjective nature of the community mining task, there is no formal definition on the term community. Consequently, there is no consensus on how to measure “goodness” of the discovered communities by a mining algorithm. However, the well-studied clustering methods in the Machine Learning field are subject to similar issues and yet there exists an extensive set of validity criteria defined for clustering evaluation, such as Davies-Bouldin index [4], Dunn index [5], and Silhouette [6] (for a complete survey refer to [7]). In this paper, we describe how these criteria could be adapted to the context of community mining in order to compare results of different community mining algorithms. Also, these criteria can be used as alternatives to modularity to design novel community mining algorithms.

In the following, we first briefly introduce well-known community mining algorithms, and different evaluation/comparison methods including available benchmarks. Next, different ways to adapt clustering validity criteria to handle comparison of community mining results is proposed. Then, we extensively compare and discuss the adapted criteria on real and synthetic networks. Finally, we conclude with a brief analysis of these results.

II. RELATED WORK

A community is roughly defined as “densely connected” individuals that are “loosely connected” to others. A great number of community mining algorithms have been developed in the last few years having different interpretations of this definition. Basic heuristic approaches mine communities by assuming that the network of interest divides naturally into some subgroups, determined by the network itself. For instance, the Clique Percolation Method [8] finds groups of nodes that can be reached via chains of k -cliques. The common optimization approaches mine communities by maximizing the overall “goodness” of the result. The most credible “goodness” objective is known as modularity Q , proposed in [1], which considers the difference between the fraction of edges that are within the communities and the expected such fraction if the edges are randomly distributed. Several community mining algorithms for optimizing modularity Q have been proposed such as fast modularity [9], and Max-Min modularity [10]. Local modularity M [11], and local modularity L [12], are local variants of modularity Q , where the ratio of internal and external edges is calculated by identifying boundary nodes of a detected local community. Although many mining algorithms are based on the concept of modularity, Fortunato et al. [13] have shown that that modularity cannot accurately evaluate small communities due to its resolution limit. Hence, any algorithm based on modularity is biased against small communities. As an alternative to optimizing modularity Q , we previously proposed TopLeaders community mining approach [14], which implicitly maximizes the overall closeness of followers and leaders, assuming that a community is a set of followers congregating around a potential leader. There are many other alternative methods. One notable family of approaches mine communities by utilizing information theory concepts such as compression (e.g. Infomod [15], Infomap [16]), and entropy (e.g. entropy-base [17]). For a survey on different community mining techniques refer to [3].

The standard procedure for evaluating the results of a community mining algorithm is the external evaluation of results, particularly when comparing accuracy of different algorithms. Which is assessing the agreement between the results and the ground truth that is known for the benchmark datasets. These benchmarks are typically small real world datasets or synthetic networks, [1], [10], [15]. On the other hand, there is no well-defined criterion for evaluating the resulted communities for networks without any ground truth, which is the case in most of real world applications. The common practice is to validate the results partly by a human expert, for example see [11]. However, the community mining problem is NP-complete; thus, the human expert validation is limited and rather based on narrow intuition than on an exhaustive examination of the relations in the given network. Alternatively, modularity Q is sometimes reported to show the quality of discovered communities. In this paper, we investigate other potential measures for comparing different (non-overlapping) community mining results and compare their

assessment with the modularity Q . All these new measures are adapted from well-grounded traditional clustering criteria for evaluating data points with attributes. Recently Vendramin et al. comprehensively compared their performances in [7], based on the idea that the better a criterion the more correlative is its ranking of different partitions to ranking of external index.

The external evaluation, requires knowing the true communities. For this purpose, several generators are proposed for synthesizing networks with built-in ground truth. GN benchmark [18] is the first synthetic network generator. This benchmark is a graph with 128 nodes, with expected degree 16, and is divided into four groups of equal sizes; where the probabilities of the existence of a link between a pair of nodes of the same group and of different groups are z_{in} and $1 - z_{in}$, respectively. However, the same expected degree for all the node, and equal-size communities are not accordant to real social network properties. LFR benchmark [19] amends GN benchmark by considering power law distributions for degrees and community sizes. Similar to GN benchmark, each node shares a fraction $1 - \mu$ of its links with the other nodes of its community and a fraction μ with the other nodes of the network. In this paper, we generate our synthetic networks using LFR benchmark, due to its more realistic structure.

There are recent studies on the comparison of different community mining algorithms in terms of evaluating their performance on synthetic and real networks. For example refer to studies by Danon et al. [20], Gustafsson et al. [21], and Lancichinetti and Fortunato [22]. All these studies are based on the agreements of the generated communities with the true one in the ground truth and are using GN and/or LFR benchmarks. Orman et al. [23] further performed a qualitative analysis of the identified communities by comparing distribution of resulted communities with the community size distribution of the ground truth. None of these studies, however, considers any different validity criteria other than modularity to evaluate the goodness of the detected communities. In this paper, we plan to examine potential validity criteria specifically defined for evaluation of community mining results. In the future, these criteria not only can be used as a mean to measure the goodness of discovered communities, but also as an objective function to detect communities.

III. COMMUNITY QUALITY CRITERIA

In this section, we overview several validity criteria that could be used as relative indexes for comparing and evaluating different partitionings of a given network. All of these criteria are generalized from well-known clustering criteria. The clustering quality criteria are defined with the implicit assumption that data points consist of vectors of attributes. Consequently their definition is mostly integrated or mixed with the definition of the distance measure between data points. The commonly used distance measure is the Euclidean distance, which cannot be defined for graphs. Therefore, we first review different possible distance measures that could be used in graphs. Then, we present generalizations of criteria that could use any notion of distance.

Distance Between Nodes

Let A denotes the adjacency matrix of the graph, and A_{ij} is the weight of the edge between nodes n_i and n_j . The distance $d(i, j)$ denotes the distance between n_i and n_j , which can be computed by one of the following measures.

1) *Edge Path (ED)*: The distance between two nodes is the inverse of their incident edge weight ¹:

$$d_{ED}(i, j) = \frac{1}{A_{ij}}$$

2) *Shortest Path Distance (SPD)*: The distance between two nodes is the length of the shortest path between them, which could be computed using the well-known Dijkstras Shortest Path algorithm [24].

3) *Adjacency Relation Distance (ARD)*: The distance between two nodes is the structural dissimilarity of them, that is computed by the difference between their immediate neighbourhood [24]:

$$d_{ARD}(i, j) = \sqrt{\sum_{k \neq j, i} (A_{ik} - A_{jk})^2}$$

4) *Neighbour Overlap Distance (NOD)*: The distance between two nodes is the ratio of the unshared neighbours between them [3]:

$$d_{NOD}(i, j) = 1 - \frac{|\aleph_i \cap \aleph_j|}{|\aleph_i \cup \aleph_j|}$$

where \aleph_i is the set of nodes directly connected to n_i . Note that there is a close relation between this measure and the previous one, since similarly d_{NOD} could be re-written as:

$$d_{NOD}(i, j) = 1 - \frac{\sqrt{\sum_{k \neq j, i} (A_{ik} + A_{jk})^2} - \sqrt{\sum_{k \neq j, i} (A_{ik} - A_{jk})^2}}{\sqrt{\sum_{k \neq j, i} (A_{ik} + A_{jk})^2} + \sqrt{\sum_{k \neq j, i} (A_{ik} - A_{jk})^2}}$$

The latter formulation of d_{NOD} in terms of the adjacency matrix can be straightforwardly generalized for weighted graphs.

5) *Pearson Correlation Distance (PCD)*: The Pearson correlation coefficient between two nodes is the correlation between their corresponding rows of the adjacency matrix [3]:

$$C(i, j) = \frac{\sum_k (A_{ik} - \mu_i)(A_{jk} - \mu_j)}{N\sigma_i\sigma_j}$$

where N is the number of nodes, the average $\mu_i = (\sum_k A_{ik})/N$ and the variance $\sigma_i = \sqrt{\sum_k (A_{ik} - \mu_i)^2/N}$. Then the distance between two nodes is computed as $d_{PCD}(i, j) = 1 - C(i, j)$ which lies between 0 (when the two nodes are most similar) and 2 (when the two nodes are most dissimilar).

¹For avoiding division by zero when A_{ij} is zero, $1/\epsilon$ is returned where ϵ is a very small number; same is true for all other formula whenever a division by zero may occur.

6) *ICloseness Distance (ICD)*: The distance between two nodes is computed as the inverse of the connectivity between their common neighbourhood:

$$d_{ICD}(i, j) = \frac{1}{\sum_{k \in \aleph_i \cap \aleph_j} ns(k, i)ns(k, j)}$$

where $ns(k, i)$ denotes the neighbouring score between node k and i that is computed iteratively (for complete formulation refer to [25]).

Community Centroid

In addition to the notion of distance measure, most of the cluster validity criteria use averaging between the numerical data points to determine the centroid of a cluster. The averaging is not defined for nodes in a graph, therefore we modify the criteria definitions to use a generalized centroid notion, in a way that, if the centroid is set as averaging, we would obtain the original criteria definitions, but we could also use other alternative notions for centroid of a group of data points.

Averaging data points results in a point with the least average distance to the other points. When averaging is not possible, using medoid is the natural option, which is perfectly compatible with graphs. More formally, the centroid of a community could be obtained as:

$$\bar{C} = \arg \min_{m \in C} \sum_{i \in C} d(i, m)$$

Relative Validity Criteria

Here we present our generalization of the well-known clustering validity criteria defined as quality measures for internal evaluation of clustering results. All these criteria are originally defined based on distances between data points, which is in all cases the Euclidean or norm of difference between their vectors of attributes, refer to [7] for comparative analysis of these criteria in the clustering context. Here we alter the formulae to use a generalized distance metric, so that we could plug in our graph distance measure. The second alteration is generalizing the mean over data points to a general centroid notion, which would be set as averaging in the presence of attributes and *medoid* in our case of dealing with graphs and in the absence of attributes.

In a nutshell, in every criterion, the average of points in a cluster is replaced with a generalized notion of centroid, and distances between data points are generalized from Euclidean/norm to a generic distance. Consider a clustering $C = \{C_1 \cup C_2 \cup \dots \cup C_k\}$ of N data points, where \bar{C} denotes the centroid of data points belonging to C . The quality of C could be defined using one of the following criteria.

7) *Variance Ratio Criterion*: This criterion measures the ratio of the between-community distances to within-community distances which could be generalized as follows:

$$VRC = \frac{\sum_{l=1}^k |C_l| d(\bar{C}_l, \bar{C})}{\sum_{l=1}^k \sum_{i \in C_l} d(i, \bar{C}_l)} \times \frac{N - k}{k - 1}$$

where \bar{C}_l is the centroid of the community C_l , and \bar{C} is the centroid of the entire network. The original clustering formula proposed in [26] for attributes vectors, is obtained if the centroid is fixed to averaging of attributes vectors and distance to (square of) Euclidean distance.

8) *Davies-Bouldin index*: This minimization criterion calculates the worst-case within-community to between-community distances ratio averaged over all communities [4]:

$$DB = \frac{1}{k} \sum_{l=1}^k \max_{m \neq l} ((\bar{d}_l + \bar{d}_m) / d(\bar{C}_l, \bar{C}_m))$$

$$\bar{d}_l = \frac{1}{|C_l|} \sum_{i \in C_l} d(i, \bar{C}_l)$$

9) *Dunn index*: : This criterion considers both the minimum distance between any two communities and the length of the largest community diameter (i.e. the maximum or the average distance between all the pairs in the community) [5]:

$$Dunn = \min_{l \neq m} \left\{ \frac{\delta(C_l, C_m)}{\max_p \Delta(C_p)} \right\}$$

where δ denotes distance between two communities and Δ is the diameter of a community. Different variations of calculating δ and Δ are available; δ could be single, complete or average linkage, or only the difference between the two centroids. Moreover, Δ could be maximum or average distance between all pairs of nodes, or the average distance of all nodes to the centroid. For example the single linkage for δ and maximum distance for Δ are $\delta(C_l, C_m) = \min_{i \in C_l, j \in C_m} d(i, j)$ and $\Delta(C_p) = \max_{i, j \in C_p} d(i, j)$. Therefore we have different variations of Dunn index in our experiments, each indicated by two indexes for different methods to calculate δ (i.e. single(0), complete(1), average(2), and centroid(4)) and different methods to calculate Δ (i.e. maximum(0), average(1), average to centroid(4)).

10) *Silhouette Width Criterion*: This criterion measures the average of silhouette score for each data point. The silhouette score of a point shows the goodness of the community it belongs to by calculating the difference between the distance to its nearest neighbouring community and its own community [6]:

$$SWC = \frac{1}{N} \sum_{l=1}^k \sum_{i \in C_l} \frac{\min_{m \neq l} d(i, C_m) - d(i, C_l)}{\max \{ \min_{m \neq l} d(i, C_m), d(i, C_l) \}}$$

where $d(i, C_l)$ is the distance of point i to community C_l , which is originally set to be the average distance (called SWC2) (i.e. $1/|C_l| \sum_{j \in C_l} d(i, j)$) or could be the distance to its centroid (called SWC4) [27] (i.e. $d(i, \bar{C}_l)$). An alternative formula for Silhouette is proposed in [28] :

$$ASWC = \frac{1}{N} \sum_{l=1}^k \sum_{i \in C_l} \frac{\min_{m \neq l} d(i, C_m)}{d(i, C_l)}$$

11) *PBM*: This criterion is based on the within-community distances and the maximum distance between centroids of communities [29]:

$$PBM = \frac{1}{k} \times \frac{\max_{l,m} d(\bar{C}_l, \bar{C}_m)}{\sum_{l=1}^k \sum_{i \in C_l} d(i, \bar{C}_l)}$$

12) *C-Index*: This criterion compared the sum of the within-community distances to the worst and best case scenarios [30]. The best case scenario is where the within-community distances are the shortest distances in the graph, and the worst case scenario is where the within-community distances are the longest distances in the graph.

$$\theta = \frac{1}{2} \sum_{l=1}^k \sum_{i,j \in C_l} d(i, j)$$

$$CIndex = \frac{\theta - \min \theta}{\max \theta - \min \theta}$$

The $\min \theta / \max \theta$ is computed by summing the m_1 smallest/largest distances between every two points, where $m_1 = \sum_{l=1}^k \frac{|C_l|(|C_l|-1)}{2}$.

13) *Z-Statistics*: This criterion is similar to C-Index, however with different formulation [31]:

$$ZIndex = \frac{\theta - E(\theta)}{\sqrt{var(\theta)}}$$

$$E(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N d(i, j)$$

$$Var(\theta) = \frac{\left(\sum_{i=1}^N \sum_{j=1}^N d(i, j) \right)^2 - 2 \sum_{i=1}^N \left(\sum_{j=1}^N d(i, j) \right)^2}{N(N-1)}$$

$$- \frac{\left(\sum_{i=1}^N \sum_{j=1}^N d(i, j) \right)^2}{N^2} + \frac{\sum_{i=1}^N \sum_{j=1}^N d(i, j)^2}{N}$$

14) *Point-Biserial*: This criterion computes the correlation of the distances between any node and the members of its community [32]. Intuitively, nodes that are in the same community should be separated by shorter distances than those which are not.

$$PB = \frac{M_1 - M_0}{S} \sqrt{\frac{m_1 m_0}{m^2}}$$

where m is the total number of distances i.e. $N(N-1)/2$ and S is the standard deviation of all pairwise distances i.e. $\sqrt{\frac{1}{m} \sum_{i,j} (d(i, j) - \frac{1}{m} \sum_{i,j} d(i, j))^2}$, while M_1 , M_0 are respectively the average of within and between-community distances, and m_1 and m_0 represent the number of within and between community distances. More formally

$$m_1 = \sum_{l=1}^k \frac{N_l(N_l - 1)}{2} \quad m_0 = \sum_{l=1}^k \frac{N_l(N - N_l)}{2}$$

$$M_1 = 1/2 \sum_{l=1}^k \sum_{i,j \in C_l} d(i,j) \quad M_0 = 1/2 \sum_{l=1}^k \sum_{\substack{i \in C_l \\ j \notin C_l}} d(i,j)$$

15) *Modularity*: is the well-known criterion proposed by Newman et al. [1] specifically for the context of community mining. This criterion considers the difference between the fraction of edges that are within the community and the expected such fraction if the edges were randomly distributed. Let E denotes the number of edges in the network i.e. $E = \frac{1}{2} \sum_{ij} A_{ij}$, then Q-modularity is defined as:

$$Q = \frac{1}{2E} \sum_{l=1}^k \sum_{i,j \in C_l} [A_{ij} - \frac{\sum_j A_{ij} \sum_i A_{ij}}{2E}]$$

The computational complexity of different validity criteria is provided in the previous work by Vendramin et al. [7].

IV. COMPARISON METHODOLOGY AND RESULTS

In this section, we compare the proposed relative community criteria. First, we describe the approach we have used for the comparison. Then, we report the criteria performances in different settings.

The performance of a criterion could be examined by how well it could rank different partitioning of a given dataset. More formally, consider we have a dataset d and a set of m different possible partitionings, i.e. $P(d) = \{p_{d1}, p_{d2}, \dots, p_{dm}\}$. Then, the performance of criterion c on dataset d could be determined by how much its values, $I_c(d) = \{c(p_{d1}), c(p_{d2}), \dots, c(p_{dm})\}$, correlate with the "goodness" of these partitionings. Assuming that the true partitioning (i.e. ground truth) p_d^* is known for dataset d , the "goodness" of partitioning p_{di} could be determined using partitioning agreement measure a , a.k.a. external evaluation. Hence, for dataset d with set of possible partitionings $P(d)$, the external evaluation provides $E(d) = \{a(p_{d1}, p_d^*), a(p_{d2}, p_d^*), \dots, a(p_{dm}, p_d^*)\}$, where (p_{d1}, p_d^*) denotes the "goodness" of partitioning p_{d1} comparing to the ground truth. Then, the performance score of criterion c on dataset d could be examined by the correlation of its values $I_c(d)$ and the values obtained from the external evaluation $E(d)$ on different possible partitionings. Finally, the criteria are ranked based on their average performance score over a set of datasets. The following procedure summarizes our comparison approach.

```

D ← {d1, d2, ..., dn}
for all dataset d ∈ D do
  {generate m possible partitionings}
  P(d) ← {pd1, pd2, ..., pdm}
  E(d) ← {a(pd1, pd*), a(pd2, pd*), ..., a(pdm, pd*)}
  for all c ∈ Criteria do
    Ic(d) ← {c(pd1), c(pd2), ..., c(pdm)}
    {compute the correlation}
    scorec(d) ← correlation(E, I)
  end for
end for
scorec ←  $\frac{1}{n} \sum_{d=1}^n \text{score}_c(d)$ 

```

Dataset	K^*	#	\bar{K}	\overline{AMI}
strike	3	60	3.17±1∈[2,5]	0.59±0.27∈[-0.04,1]
polboks	3	60	3.17±1.13∈[2,6]	0.44±0.25∈[0.04,1]
karate	2	60	3.57±1.23∈[2,6]	0.46±0.27∈[-0.02,1]
football	11	60	10.17±4.55∈[4,19]	0.68±0.16∈[0.4,1]

TABLE I
STATISTICS FOR SAMPLE PARTITIONINGS OF EACH REAL WORLD DATASET. FOR EXAMPLE FOR THE KARATE CLUB DATASET WHICH HAS 2 COMMUNITIES IN ITS GROUND TRUTH, WE HAVE GENERATED 60 DIFFERENT PARTITIONINGS WITH AVERAGE 3.63±1.38 CLUSTERS RANGING FROM 2 TO 7 AND THE "GOODNESS" OF SAMPLES IS ON AVERAGE 0.46±0.24 AND RANGES FROM 0.01 TO 1, IN TERMS OF THEIR AMI AGREEMENT WITH THE TRUE PARTITIONING OF THE GROUND TRUTH.

External evaluation is done with an agreement measure, which computes the agreement between two given partitionings or between a partitioning and the ground truth. There are several choices for the partitioning agreement measure. The commonly used ones are pair counting based, such as Adjusted Rank Index (ARI) [33] and Jaccard Coefficient [34], and the information theoretic-based, such as Normalized Mutual Information (NMI) [35], [20] and the Adjusted Mutual Information(AMI) [36].

There are also different ways to compute the correlation between two vectors. The classic options are Pearson Product Moment coefficient or the Spearman's Rank correlation coefficient. The reported results in our experiments are based the Spearman's Correlation, since we are interested on the correlation of rankings that a criterion provides for different partitionings and not the actual values of that criterion. However, the reported results mostly agree with the results obtained by using Pearson correlation, which are reported in the supplementary materials².

Sampling the Partitioning Space

In our comparison, we generate different partitionings for each dataset d to sample the space of all possible partitionings. For doing so, given the perfect partitioning, p_d^* , we randomized different versions of p_d^* by randomly merging and splitting communities and swapping nodes between them. The detailed procedure of sampling is described more in the supplementary materials.

A. Results on Real World Datasets

We first compare performance of different criteria on five well-known real-world benchmarks: Karate Club (weighted) by Zachary [37], Sawmill Strike data-set [38], NCAA Football Bowl Subdivision [18], and Politician Books from Amazon [39]. Table I shows general statistics about the datasets and their generated samples. We could see that the randomized samples cover the space of partitionings according to their external index range.

²Available here:
<http://webdocs.cs.ualberta.ca/~rabbanyk/criteriaComparison>

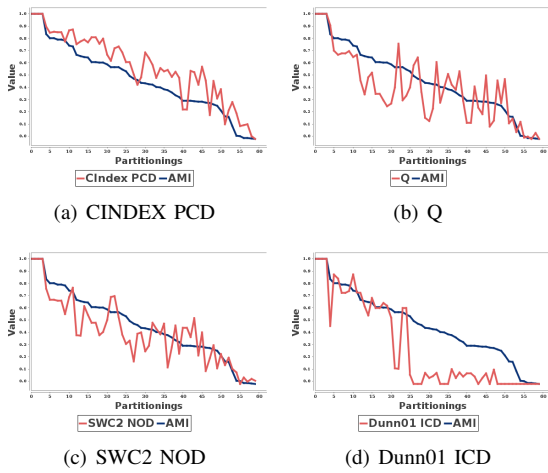


Fig. 1. Visualization of correlation between an external agreement measure and the relative quality criteria for Karate dataset. The x axis indicates different random partitionings, and the y axis indicates the value of the index. While, the blue/darker line represents the value of the external index for the given partitioning and the red/lighter line represents the value that the criterion gives for the partitioning. Please note that the value of criteria are not generally normalized and in the same range as the external indexes, in this figure AMI. For the sake of illustration therefore, each criterion’s values are scaled to be in the same range as of the external index.

Figure 1 exemplifies how different criteria exhibit different correlations with the external index. It visualizes the correlation between few selected internal indexes and an external index for one of our datasets listed in Table I.

Similar analysis is done for all 4 datasets \times 19 criteria \times 7 distances \times 4 external indexes, which produced over 2000 of such correlations. The top ranked criteria based on their average performance over these datasets are summarized in Table II. Based on these results, *CIndex* when used with *PCD* distance has a higher correlation with the external index comparing to the modularity *Q*. And this is true regardless of the choice of *AMI* as the external index, since it is ranked above *Q* also by *ARI* and *NMI*.

The correlation between a criterion and an external index depends on how close the randomized partitionings are from the true partitioning of the ground truth. This can be seen in the Figure 1. For example, the *Dunn01* (single linkage network diameter and average linkage within community scores) with the *ICD* distance agrees strongly with the external index in samples with higher external index value, i.e. closer to the ground truth, but not on further samples. On the other hand, *Q* is very well matched for the samples too far or too close to the ground truth but is not doing as well as others in the middle. With this in mind, we have divided the generated clustering samples into three sets of easy, medium and hard samples and re-ranked the criteria in each of these settings. Since the external index determines how far a sample is from the optimal result, the samples are divided into three equal length intervals according to the range of the external index. Table III, reports the rankings in each of these three settings. We could see that these average results supports our earlier hypothesis, i.e. when considering partitionings medium far from the true

Rank	Criterion	AMI	ARI	Jaccard	NMI
1	CIndex PCD	0.907±0.058	1	1	1
2	SWC2 NOD	0.857±0.031	4	4	2
3	Q	0.85±0.083	2	2	3
4	CIndex ARD	0.826±0.162	6	15	5
5	CIndex SPD	0.811±0.126	3	10	4
6	ASWC2 NOD	0.809±0.043	5	11	6
7	CIndex NOD	0.794±0.096	12	3	9
8	SWC2 PCD	0.789±0.103	7	7	8
9	SWC4 NOD	0.778±0.075	9	5	7
10	ASWC2 PCD	0.772±0.088	10	9	10
11	SWC2 SPD	0.751±0.121	8	6	11
12	Dunn01 ICD	0.742±0.111	18	24	12
13	ASWC2 SPD	0.733±0.116	11	8	13
14	Dunn00 PCD	0.721±0.1	21	30	14
15	DB ICD	0.712±0.063	24	22	16
16	Dunn00 ICD	0.707±0.133	28	28	15
17	Dunn04 ICD	0.703±0.055	25	23	17
18	SWC4 PCD	0.7±0.072	14	12	21
19	SWC4 SPD	0.681±0.081	15	13	23
20	SWC2 ARD	0.681±0.302	17	26	19

TABLE II
OVERALL RANKING OF CRITERIA ON THE REAL WORLD DATASETS, BASED ON AMI AND USING THE SPEARMAN’S CORRELATION. THE FIRST/THIRD COLUMN SHOWS THE AVERAGE RANKING/CORRELATION OF CRITERIA WITH THE AMI EXTERNAL INDEX. RANKING BASED ON OTHER EXTERNAL INDEXES IS ALSO PROVIDED IN SEPARATE COLUMNS.

partitioning, *CIndex PCD* performs significantly better than the modularity *Q*, which is true in neither the near optimal samples nor the samples very far from the ground truth. One may conclude based on this experiment that, *CIndex PCD* is a more accurate evaluation criterion comparing to *Q* especially when the results might not be very perfect or very poor.

Near Optimal Samples					
Rank	Criterion	AMI	ARI	Jaccard	NMI
1	Q	0.736±0.266	5	5	2
2	CIndex PCD	0.72±0.326	1	1	3
3	SWC2 SPD	0.718±0.389	3	3	4
4	CIndex SPD	0.716±0.14	4	4	1
5	SWC2 ICD	0.713±0.396	2	2	5
Medium Far Samples					
Rank	Criterion	AMI	ARI	Jaccard	NMI
1	CIndex PCD	0.608±0.202	8	18	1
2	CIndex NOD	0.58±0.053	39	13	2
3	CIndex ARD	0.513±0.313	26	62	5
4	Dunn01 ICD	0.457±0.173	58	83	8
5	SWC2 NOD	0.447±0.19	5	9	3
6	ASWC2 PCD	0.446±0.191	7	3	9
7	SWC2 PCD	0.446±0.19	6	2	10
8	Dunn04 ICD	0.439±0.109	43	37	11
9	Dunn41 SPD	0.437±0.177	56	47	15
10	Dunn01 SPD	0.434±0.205	29	67	7
11	Q	0.409±0.353	4	7	16
Far Far Samples					
Rank	Criterion	AMI	ARI	Jaccard	NMI
1	SWC2 NOD	0.634±0.217	3	13	1
2	ASWC2 NOD	0.583±0.191	5	21	2
3	Q	0.498±0.179	4	38	5
4	CIndex PCD	0.493±0.282	2	4	13
5	CIndex SPD	0.437±0.291	1	11	4

TABLE III
DIFFICULTY ANALYSIS OF THE RESULTS: CONSIDERING RANKING FOR PARTITIONINGS NEAR OPTIMAL GROUND TRUTH, MEDIUM FAR AND VERY FAR. REPORTED RESULT ARE BASED ON AMI AND THE SPEARMAN’S CORRELATION.

Dataset	K^*	#	\bar{K}	\overline{AMI}
network5	2	60	$3.5 \pm 1.36 \in [2, 8]$	$0.44 \pm 0.22 \in [0.11, 1]$
network9	5	60	$5.5 \pm 2.05 \in [2, 10]$	$0.69 \pm 0.19 \in [0.37, 1]$
network7	4	60	$5.2 \pm 2.65 \in [2, 12]$	$0.47 \pm 0.19 \in [0.13, 1]$
network3	2	60	$3.3 \pm 1.13 \in [2, 6]$	$0.47 \pm 0.23 \in [0.11, 1]$
network6	5	60	$5.8 \pm 2.55 \in [2, 12]$	$0.68 \pm 0.2 \in [0.27, 1]$
network8	5	60	$5.37 \pm 2.04 \in [2, 10]$	$0.67 \pm 0.21 \in [0.32, 1]$
network10	6	60	$5.33 \pm 2.51 \in [2, 11]$	$0.63 \pm 0.19 \in [0.24, 1]$
network1	4	60	$3.4 \pm 1.17 \in [2, 6]$	$0.46 \pm 0.23 \in [-0, 1]$
network2	3	60	$3.1 \pm 1.27 \in [2, 7]$	$0.49 \pm 0.22 \in [0.13, 1]$
network4	7	60	$5.17 \pm 2.49 \in [2, 12]$	$0.57 \pm 0.2 \in [0.18, 1]$

TABLE IV

STATISTICS FOR SAMPLE PARTITIONINGS OF EACH SYNTHETIC DATASET. THE BENCHMARK GENERATION PARAMETERS: 100 NODES WITH AVERAGE DEGREE 5 AND MAXIMUM DEGREE 50 WHERE SIZE OF EACH CLUSTER COULD BE BETWEEN FROM 5 TO 50 WHERE MIXING PARAMETER IS .1.

B. Synthetic Benchmarks Datasets

Lastly, we compare the criteria on a larger set of synthetic benchmarks. We have generated our dataset using the LFR benchmarks [19] which are the generators widely in use for community mining evaluation. Similar to the last experiment, Table V reports the ranking of the criteria according to their average performance on synthesized datasets of Table IV.

Based on datasets of Table IV, modularity Q overall outperforms other criteria especially in ranking poor partitionings; while, *CIndex PCD* performs better in ranking finer results.

Overall Results					
Rank	Criterion	AMI	ARI	Jaccard	NMI
1	Q	0.894 ± 0.018	1	2	1
2	ASWC2 NOD	0.854 ± 0.056	3	4	2
3	SWC2 NOD	0.854 ± 0.051	4	3	3
4	CIndex PCD	0.826 ± 0.07	2	1	4
5	CIndex SPD	0.746 ± 0.137	8	24	5
6	SWC2 PCD	0.743 ± 0.047	5	5	6
7	ASWC2 PCD	0.739 ± 0.048	6	6	7
8	Dunn00 PCD	0.707 ± 0.11	11	26	8
Near Optimal Results					
Rank	Criterion	AMI	ARI	Jaccard	NMI
1	CIndex PCD	0.729 ± 0.17	1	1	1
2	Q	0.722 ± 0.111	6	5	5
3	SWC2 SPD	0.717 ± 0.185	18	18	2
4	SWC4 NOD	0.709 ± 0.201	5	6	4
5	SWC2 ICD	0.704 ± 0.216	15	15	3
Medium Far Results					
Rank	Criterion	AMI	ARI	Jaccard	NMI
1	SWC2 NOD	0.455 ± 0.191	5	11	3
2	CIndex PCD	0.453 ± 0.245	1	2	5
3	Q	0.45 ± 0.236	2	9	2
4	ASWC2 NOD	0.435 ± 0.187	4	14	1
5	Dunn00 ARD	0.386 ± 0.243	119	111	7
Far Far Results					
Rank	Criterion	AMI	ARI	Jaccard	NMI
1	Q	0.63 ± 0.139	1	4	2
2	ASWC2 NOD	0.596 ± 0.164	2	2	3
3	SWC2 NOD	0.57 ± 0.159	3	3	5
4	CIndex SPD	0.565 ± 0.132	4	25	1
5	CIndex PCD	0.446 ± 0.142	5	1	21

TABLE V

OVERALL RANKING AND DIFFICULTY ANALYSIS OF THE SYNTHETIC RESULTS. HERE COMMUNITIES ARE WELL-SEPARATED WITH MIXING PARAMETER OF .1. SIMILAR TO THE LAST EXPERIMENT, REPORTED RESULT ARE BASED ON AMI AND THE SPEARMAN'S CORRELATION.

The LFR generator, can generate networks with different level of difficulty for the partitioning task, by changing how well separated the communities are in the ground truth. To examine the effect of this difficulty parameter, we have ranked the criteria for different values of this parameter. We observed that indeed the modularity Q is the superior criterion for these synthetic benchmarks up to some level of how mixed are the communities, but this changes in more difficult settings. ³ Table VI reports overall ranking of the criteria for a difficult set of datasets that have high mixing parameter. We could see that in this setting *PB* index used with *PCD*, *NOD*, *SPD* or *ARD* distance, is a significantly more reliable criterion comparing to modularity Q considering its much higher performance variance.

Rank	Criterion	AMI	ARI	Jaccard	NMI
1	PB PCD	0.454 ± 0.15	1	1	1
2	PB NOD	0.448 ± 0.146	2	2	2
3	PB SPD	0.445 ± 0.144	3	3	4
4	PB ARD	0.44 ± 0.149	4	4	5
5	VRC ICD	0.424 ± 0.117	5	5	3
6	Q	0.391 ± 0.381	17	6	12
7	CIndex ARD	0.365 ± 0.173	6	7	6
8	ASWC4 SPD	0.358 ± 0.101	12	12	7
9	DB PCD	0.358 ± 0.108	15	9	10
10	ASWC4 NOD	0.357 ± 0.114	10	10	8

TABLE VI

OVERALL RANKING OF CRITERIA BASED ON AMI & SPEARMAN'S CORRELATION ON THE SYNTHETIC BENCHMARKS WITH SAME PARAMETERS BUT MUCH HIGHER MIXING PARAMETER, .7. WE COULD SEE THAT IN THESE SETTINGS, OTHER INDEXES OUTPERFORMED Q MODULARITY.

In short, the ranking of criteria also strongly depends on the difficulty of the network itself, apart from how far we are sampling from the ground truth. Altogether, choosing the right criterion for evaluating different community mining results, depends both on the application, i.e. how well-separated communities might be in the given network, and also on the algorithm that produces these results, i.e. how fine the results might be. For example if the problem is hard and communities mixed heavily, modularity Q might not distinguish the good and bad partitionings very well. While if we are choosing between fine and well separated clusterings, it indeed is the superior criterion.

V. CONCLUSION

In this paper, we generalized well-known clustering validity criteria originally used as quantitative measures for evaluating quality of clusters of data points represented by attributes. The first reason of this generalization is to adapt these criteria in the context of community mining of interrelated data. The only commonly used criterion to evaluate the goodness of detected communities in a network is the modularity Q. Providing more validity criteria can help researchers to better evaluate and compare community mining results in different settings. Also,

³Results for other settings available here: <http://webdocs.cs.ualberta.ca/~rabbanyk/criteriaComparison>

these adapted validity criteria can be further used as objectives to design new community mining algorithms. Our generalized formulation is independent of any particular distance measure unlike most of the original validity criteria that are defined based on the Euclidean distance. The adopted versions therefore could be used as community criteria when plugged in with different graph distances. In our experiments, several of these adopted criteria exhibit high performances on ranking different partitionings of a given dataset. Which makes them possible alternatives for the Q modularity. However a more careful examination is needed as the rankings depends significantly on the experiment's settings and the criteria should be chosen based on the application.

ACKNOWLEDGMENT

Ricardo Campello acknowledges Fapesp and CNPq for financial support.

REFERENCES

- [1] M. E. J. Newman and M. Girvan, "Finding and evaluating community structure in networks," *Physical Review E*, vol. 69, no. 2, p. 026113, 2004.
- [2] J. Leskovec, J. Kleinberg, and C. Faloutsos, "Graphs over time: densification laws, shrinking diameters and possible explanations," in *ACM SIGKDD international conference on Knowledge discovery in data mining*, 2005, pp. 177–187.
- [3] S. Fortunato, "Community detection in graphs," *Physics Reports*, vol. 486, no. 35, pp. 75–174, 2010.
- [4] D. L. Davies and D. W. Bouldin, "A cluster separation measure," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. PAMI-1, no. 2, pp. 224–227, 1979.
- [5] J. C. Dunn, "Well-separated clusters and optimal fuzzy partitions," *Journal of Cybernetics*, vol. 4, no. 1, pp. 95–104, 1974.
- [6] P. Rousseeuw, "Silhouettes: a graphical aid to the interpretation and validation of cluster analysis," *Journal of Computational and Applied Mathematics*, vol. 20, no. 1, pp. 53–65, 1987.
- [7] L. Vendramin, R. J. G. B. Campello, and E. R. Hruschka, "Relative clustering validity criteria: A comparative overview," *Statistical Analysis and Data Mining*, vol. 3, no. 4, pp. 209–235, 2010.
- [8] G. Palla, I. Derenyi, I. Farkas, and T. Vicsek, "Uncovering the overlapping community structure of complex networks in nature and society," *Nature*, vol. 435, no. 7043, pp. 814–818, 2005.
- [9] M. E. J. Newman, "Modularity and community structure in networks," *Proceedings of the National Academy of Sciences*, vol. 103, no. 23, pp. 8577–8582, 2006.
- [10] J. Chen, O. R. Zaiane, and R. Goebel, "Detecting communities in social networks using max-min modularity," in *SIAM International Conference on Data Mining*, 2009, pp. 978–989.
- [11] F. Luo, J. Z. Wang, and E. Promislow, "Exploring local community structures in large networks," *Web Intelligence and Agent Systems*, vol. 6, pp. 387–400, 2008.
- [12] J. Chen, O. R. Zaiane, and R. Goebel, "Detecting communities in large networks by iterative local expansion," in *International Conference on Computational Aspects of Social Networks*, 2009, pp. 105–112.
- [13] S. Fortunato and M. Barthélemy, "Resolution limit in community detection," *Proceedings of the National Academy of Sciences*, vol. 104, no. 1, pp. 36–41, 2007.
- [14] R. Rabbany, J. Chen, and O. R. Zaiane, "Top leaders community detection approach in information networks," in *SNA-KDD Workshop on Social Network Mining and Analysis*, 2010.
- [15] M. Rosvall and C. T. Bergstrom, "An information-theoretic framework for resolving community structure in complex networks," *Proceedings of the National Academy of Sciences*, vol. 104, no. 18, pp. 7327–7331, 2007.
- [16] —, "Maps of random walks on complex networks reveal community structure," *Proceedings of the National Academy of Sciences*, vol. 105, no. 4, pp. 1118–1123, 2008.
- [17] E. C. Kenley and Y.-R. Cho, "Entropy-based graph clustering: Application to biological and social networks," in *IEEE International Conference on Data Mining*, 2011.
- [18] M. Girvan and M. E. J. Newman, "Community structure in social and biological networks," *Proceedings of the National Academy of Sciences*, vol. 99, no. 12, pp. 7821–7826, 2002.
- [19] A. Lancichinetti, S. Fortunato, and F. Radicchi, "Benchmark graphs for testing community detection algorithms," *Physical Review E*, vol. 78, no. 4, p. 046110, 2008.
- [20] L. Danon, A. D. Guílerá, J. Duch, and A. Arenas, "Comparing community structure identification," *Journal of Statistical Mechanics: Theory and Experiment*, no. 09, p. 09008, 2005.
- [21] M. Gustafsson, M. Hrnquist, and A. Lombardi, "Comparison and validation of community structures in complex networks," *Physica A*, vol. 367, no. 0, pp. 559–576, 2006.
- [22] A. Lancichinetti and S. Fortunato, "Community detection algorithms: A comparative analysis," *Physical Review E*, vol. 80, no. 5, p. 056117, 2009.
- [23] G. K. Orman, V. Labatut, and H. Cherifi, "Qualitative comparison of community detection algorithms," in *International Conference on Digital Information and Communication Technology and Its Applications*, vol. 167, 2011, pp. 265–279.
- [24] S. Wasserman and K. Faust, *Social Network Analysis: Methods and Applications*. Cambridge University Press, 1994.
- [25] R. Rabbany and O. R. Zaiane, "A diffusion of innovation-based closeness measure for network associations," in *IEEE International Conference on Data Mining Workshops*, 2011, pp. 381–388.
- [26] T. Calinski and J. Harabasz, "A dendrite method for cluster analysis," *Communications in Statistics-theory and Methods*, vol. 3, pp. 1–27, 1974.
- [27] E. R. Hruschka, R. J. Campello, and L. N. de Castro, "Evolving clusters in gene-expression data," *Information Sciences*, vol. 176, no. 13, pp. 1898–1927, 2006.
- [28] E. Hruschka, R. Campello, and L. de Castro, "Improving the efficiency of a clustering genetic algorithm," in *Ibero-American Conference on Artificial Intelligence*, 2004, pp. 861–870.
- [29] M. Pakhira and A. Dutta, "Computing approximate value of the pbm index for counting number of clusters using genetic algorithm," in *International Conference on Recent Trends in Information Systems*, 2011, pp. 241–245.
- [30] E. C. Dalrymple-Alford, "Measurement of clustering in free recall," *Psychological Bulletin*, vol. 74, pp. 32–34, 1970.
- [31] L. J. Hubert and J. R. Levin, "A general statistical framework for assessing categorical clustering in free recall," *Psychological Bulletin*, vol. 83, pp. 1072–1080, 1976.
- [32] G. Milligan and M. Cooper, "An examination of procedures for determining the number of clusters in a data set," *Psychometrika*, vol. 50, no. 2, pp. 159–179, 1985.
- [33] L. Hubert and P. Arabie, "Comparing partitions," *Journal of Classification*, vol. 2, pp. 193–218, 1985.
- [34] P. Jaccard, "Étude comparative de la distribution florale dans une portion des alpes et des jura," *Bulletin del la Société Vaudoise des Sciences Naturelles*, vol. 37, pp. 547–579, 1901.
- [35] T. O. Kvalseth, "Entropy and correlation: Some comments," *Systems, Man and Cybernetics, IEEE Transactions on*, vol. 17, no. 3, pp. 517–519, may 1987.
- [36] N. X. Vinh, J. Epps, and J. Bailey, "Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance," *Journal of Machine Learning Research*, vol. 11, pp. 2837–2854, 2010.
- [37] W. W. Zachary, "An information flow model for conflict and fission in small groups," *Journal of Anthropological Research*, vol. 33, pp. 452–473, 1977.
- [38] W. d. Nooy, A. Mrvar, and V. Batagelj, *Exploratory Social Network Analysis with Pajek*. Cambridge University Press, 2004.
- [39] V. Krebs, "Books about us politics," <http://www.orgnet.com/>.