

Bridge Bounding: A Local Approach for Efficient Community Discovery in Complex Networks

Symeon Papadopoulos^{*}
Informatics and Telematics
Institute
P.O.Box 60361, 57001, Themi
Thessaloniki, Greece
papadop@iti.gr

Andre Skusa
Lycos Europe GmbH
Carl-Bertelsmann-Str. 29
33311 Gütersloh, Germany
andre.skusa@lycos-
europe.com

Athena Vakali
Department of Informatics
Aristotle University of
Thessaloniki
54124 Thessaloniki, Greece
avakali@csd.auth.gr

Yiannis Kompatsiaris
Informatics and Telematics
Institute
P.O.Box 60361, 57001, Themi
Thessaloniki, Greece
ikom@iti.gr

Nadine Wagner
Lycos Europe GmbH
Carl-Bertelsmann-Str. 29
33311 Gütersloh, Germany
nadine.wagner@lycos-
europe.com

ABSTRACT

The increasing importance of Web 2.0 applications during the last years has created significant interest in tools for analyzing and describing collective user activities and emerging phenomena within the Web. Network structures have been widely employed in this context for modeling users, web resources and relations between them. However, the amount of data produced by modern web systems results in networks that are of unprecedented size and complexity, and are thus hard to interpret. To this end, *community detection* methods attempt to uncover natural groupings of web objects by analyzing the topology of their containing network.

There are numerous techniques adopting a *global* perspective to the community detection problem, i.e. they operate on the complete network structure, thus being computationally expensive and hard to apply in a streaming manner. In order to add a *local* perspective to the study of the problem, we present *Bridge Bounding*, a local methodology for community detection, which explores the local network topology around a seed node in order to identify edges that act as boundaries to the local community. The proposed method can be integrated in an efficient global community detection scheme that compares favorably to the state of the art. As a case study, we apply the method to explore the topic structure of the *LYCOS iQ* collaborative question/answering application by detecting communities in the networks created from the collective tagging activity of users.

^{*}This author is also affiliated with the Aristotle University of Thessaloniki.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 200X ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

Categories and Subject Descriptors

H.3.3 [Information Systems]: Information Storage and Retrieval—*clustering, information filtering*; G.2.2 [Discrete Mathematics]: Graph Theory

General Terms

Algorithms, Experimentation

Keywords

Graph partitioning, local community detection, tag network

1. INTRODUCTION

Network structures (also called *graphs* in mathematical literature) and the associated analysis methods have long emerged as a valuable tool for modeling and analyzing the relations among objects in a variety of established scientific disciplines, e.g. social sciences and biology [21]. Recent years however have witnessed a substantial adoption of network analysis techniques in the field of computer science, and more specifically, in the modeling and analysis of massive data sets produced by online information systems, such as Web 2.0 applications.

In the field of network research, the problem of *community detection* has lately attracted significant interest since identifying the community structure of large networks can improve our understanding of the complex relations that exist among their elements. The origins of this problem can be traced in the fields of citation study [13], bibliometrics [28] and social network analysis [25]. More recently, this problem has been restated in the context of web graphs, i.e. the networks created from mapping the web hyperlink structure to the directed network model. Two seminal web community definitions were formulated by Kumar et al. [18] and Flake et al. [11]. According to the first, a community is a *dense directed bipartite subgraph* of the web graph [18]. The latter definition states that a community is a *vertex subset of the graph such that each of its members has at least as many edges to other members of the community as it does to*

non-member vertices [11]. Although these two community definitions are different, they both result in the formulation of community detection as a problem of finding a partition of a graph into subgraphs that maximizes some measure of within-subgraph density.

Due to the extremely high complexity of providing an exact solution to the community detection problem for the complete network¹, several attempts have been made to derive approximate solutions at reduced computational costs, with some of the most efficient techniques having a complexity of $O(n \log^2 n)$ [22] and $O(m + n)$ [29] for networks of n nodes and m edges. Despite being very efficient, most of the existing approaches adopt a *global* perspective, i.e. they operate on the full network, in order to output the detected community structure. In practice, however, there is frequently a need to explore the network structure at a *local* level, e.g. in interactive network visualization [26] and information retrieval applications [27]. Such applications impose severe constraints on the response time of the underlying network analysis processes, thus prohibiting the use of global community detection methods. To date, only few methods have been proposed that can be used for community detection at a local level [2, 29]. However, they are either unsuitable for networks of scale-free topology (frequently emerging in practice) [2] or are not local by design, thus not achieving maximum efficiency when applied as local [29].

The situation described above motivated us to introduce a methodology for performing community detection at a local level; we call the proposed methodology *Bridge Bounding*. Bridge Bounding initiates the community detection process from a seed node in the network and progressively attaches neighboring nodes to the community as long as the edges connecting these nodes do not act as *boundaries*. Thus, community detection is formulated as a problem of identifying edges that act as community boundaries, (which we also call *bridges*, since they connect communities of the network to each other [8, p.140]). This problem is tackled by means of *local network topology functions*, i.e. functions that examine the network structure around an edge (local network topology) and produce a measure of the extent that these edges act as bridges. An example of such a function is the *edge clustering coefficient* [24]. In that way, we ensure that the proposed approach has low complexity and at the same time is capable of precisely identifying community boundaries.

In order to demonstrate the benefits of our approach, we applied it to both synthetic and real networks. As a first step, we validated Bridge Bounding by testing its performance on the known community structure of synthetic networks and comparing it with the widely cited approach of Girvan and Newman [15]. The proposed method could successfully detect the synthetic communities across a variety of network generation parameters and achieved equivalent or better performance than the competing method, while being computationally much more efficient. Subsequently, we employed Bridge Bounding to explore the community structure of two tag networks, English and German, created by tags used to annotate questions in the LYCOS iQ question/answering application². A set of tag communities consisting of semantically related tags were extracted, thus

revealing the structure of topics associated with the collective question-answering activity of users. The extracted tag community structure can be exploited for improved topic interest monitoring and automatic tag recommendation to users of the application.

The rest of the paper is structured as follows. Section 2 presents an overview of existing work in the field of community detection in complex networks. Subsequently, the formal description of the proposed community detection methodology is presented in Section 3. Section 4 presents the results and insights we obtained by applying Bridge Bounding both to synthetic networks with known community structure and to the LYCOS iQ tag network. Finally, Section 5 summarizes the basic contributions of the paper and delineates our future work.

2. RELATED WORK

The problem of community finding in large complex networks has attracted considerable research interest for some time now. Its origins can be traced back to the first studies of the hyperlink structure of the web, e.g. to the observation of Gibson et al. [14] that communities emerge spontaneously around authoritative web pages which are identified by means of hub pages. Then, the works by Kumar et al. [18] and Flake et al. [11] formally defined and systematically tackled the problem of community detection. In the following, we provide a list of existing methods for community detection classified according to the approach they adopt. A more detailed discussion of existing community detection methods is contained in the survey by Danon et al. [19].

Subgraph enumeration. Kumar et al. [18] consider communities as dense bipartite subgraphs of the web (seen as a directed graph). A natural way to identify dense subgraph structures is by means of graph partition enumeration. In order to drastically reduce the vast number of subgraphs that are possible by complete enumeration, the authors employ a series of heuristic pruning techniques. An extension of this definition led to the notion of γ -dense communities [9], which can be efficiently discovered based on more sophisticated subgraph enumeration and pruning criteria.

Maximum flow. Flake et al. [11] define communities as subsets of vertices that have more links (undirected) to each other than to the rest of the network nodes. To detect such communities on the web, they integrate a maximum flow strategy with an iterative crawling process. A stricter community definition was considered by Ino et al. [17] and a technique was devised to detect them that was based on both the maximum flow algorithm and an iterative graph partitioning and contraction process.

Divisive-Agglomerative methods. According to Girvan and Newman [15], the community structure of a large network should be revealed by progressively removing edges with high edge betweenness, i.e. by following a divisive approach. Following the same approach but with the use of different measures, namely the edge *clustering coefficient* and the *bridging centrality*, Radicchi et al. [24] and Hwang et al. [16], respectively, could uncover the underlying community structure of complex networks. Later, the measure of *modularity* was defined by Newman and Girvan, as a means to quantify the quality of a network partition into communities [23]. More specifically, modularity reflects the extent to which a given network partition is characterized by higher intra-community density in comparison to the one

¹The problem is believed to be NP-complete [23].

²We collected data from both the German (<http://iq.lycos.de>) and the English (<http://iq.lycos.co.uk>) version of the application

that would be observed in a random partition of the same network. Building upon this measure, the methods by Newman [22] and Clauset et al. [6] describe efficient implementations of community detection by means of agglomerative strategies.

Seed-based Flooding. An alternative approach to assigning the nodes of a network to communities was presented by da Fontoura Costa in [7]. There, the community detection process starts from a set of *hub nodes* and is implemented as a parallel flooding process emanating from the hubs. Although being seed-based, the technique in [7] is not local since it requires simultaneous discovery of all communities in a network. Thus, a local method for community finding was described by Bagrow and Boltt [2]. The authors consider an expanding neighborhood around the starting node (which they call *l-shell*) to constitute the community around it. In order to finish the expansion process, the authors employ a criterion quantifying the change in the *total emerging degree* of the community [2].

Hybrid. A combined strategy for community detection is provided by Du et al. [10]. The authors consider a three-step community detection process: (a) detection of maximal cliques (subgraph enumeration), (b) initial network partition by progressive expansion of the maximal cliques (flooding) and (c) adjustment of the original partition in order to maximize modularity.

Most of the methods presented above are global, meaning that they need to process the whole network in order to output the identified community structure. Even though some of these methods achieve low complexity (linear to the size of the network), their use is still prohibitive, when there is need for extremely responsive community detection, e.g. in interactive exploration of large networks, which can be only feasible by means of local processing of the network. We could only find two local methods [2, 29] that are suitable for identifying communities within such applications. However, we consider the first of those [2] as unsuitable for graphs of scale-free nature (since the *l-shell* would contain the whole graph after just few expansion steps), and the second [29] as not achieving maximum efficiency, since it is not local by design (i.e. there are redundant computation steps when applying the method locally). We consider that our proposed methodology addresses the community detection problem from a local perspective in a more intuitive and efficient way.

Most existing community detection methods, to date, have been applied to two types of networks: (a) networks created from crawling part of the web and (b) networks reflecting the social relations and/or interactions among people. Recently, there has also been some work highlighting the value of community detection in tagging systems.³ Part of the case study in [4], which mainly deals with the evaluation of the effectiveness of tags as a means to annotate blog articles, describes the induction of a tag hierarchy by means of a standard hierarchical clustering scheme based on cosine similarity. In another study [27], a method based on Spectral Recursive Embedding is proposed to carry out multi-clustering on the two bipartite graphs formed by the documents-words and documents-tags interrelationships in order to improve the precision of tag recommendation. Finally, Cattuto et al. [5] exploit the tag overlap between online resources in order to

³Community detection is frequently termed *clustering* in the respective literature.

identify resource communities by means of spectral methods. In this work, we apply our proposed methodology to the tag network created from the collective tagging activity of the LYCOS iQ users. In that way, we show that the topological properties of tag networks can be exploited to extract tag groups that are semantically related to each other.

3. METHODOLOGY

In this section, we will first (Section 3.1) present the basic notations and definitions from graph theory that are necessary to formalize the problem of community detection. Then, we will introduce the Bridge Bounding community detection methodology in Section 3.2.

3.1 Basic notation and definitions

We consider undirected graphs $G = (V, E)$, where V is the set of vertices and $E \subseteq V \times V$ is the set of edges connecting the vertices. An edge connecting nodes $i, j \in V$ is denoted as e_{ij} . For a vertex s of the graph, we consider its neighborhood $N(s)$ consisting of all vertices which are directly connected to s , i.e. $\forall n \in N(s) : e_{sn} \in E$. We define the degree of vertex v as $d(v) = |N(v)|$. In a similar way, the neighborhood of an *edge* e_{st} consists of all edges that share at least one endpoint with e_{st} , $N(e_{st}) = \{e_{xy} | \{x, y\} \cap \{s, t\} \neq \emptyset\}$.

Global community detection algorithms process a graph G in order to partition the graph into a set of communities, $\mathbf{P} \equiv \{C_0, C_1, \dots, C_K\}$, where $C_i \subseteq V$. When the communities produced by a method are mutually exclusive, then $C_i \cap C_j = \emptyset, \forall i, j \in \{1, 2, \dots, K\}$, with $i \neq j$. During the community detection process, we consider the set of nodes $C_U \in \mathbf{P}$ comprising all nodes that have not been assigned to any community until that moment. For convenience, we also employ the mapping $g_C : V \rightarrow \mathbf{P}$, which returns the community a vertex is assigned to (or C_U if the vertex has not been assigned to any community yet).

Local methods for community detection adopt a seed-based approach, i.e. given G and a node s in the graph, a local method will produce a community C_s around the node. It is possible to induce a global community detection method based on a local one by repeatedly applying the local community detection method to randomly selected nodes from C_U until this set is empty (i.e. all nodes of the graph have been assigned to some community). In the context of our evaluation (Section 4), we are going to induce such a global community detection scheme by employing the local Bridge Bounding method, which we describe in Section 3.2. We will refer to this scheme as *progressive* community detection.

3.2 Community detection by Bridge Bounding

Bridge Bounding is based on a simple strategy in order to identify the community C_s surrounding a seed node s . A formal description of this strategy is presented below, in Algorithm 1. Starting from s , each node n belonging to the neighborhood of s is considered a member of C_s as long as it meets two conditions (line 8 of the algorithm): (a) it is not already member of another community and (b) the edge connecting it to s is not a community boundary, i.e. not a *bridge* (in the sense of [8, p.140]). Then, all neighbors of the newly assigned nodes (the frontier set F) are checked against the same conditions and are attached to C_s (line 9, lines 5-6) if they meet them. This process is repeated un-

til it is not possible to attach additional nodes to C_s (line 3). Thus, Bridge Bounding is equivalent to a flooding process, similar to the one described in [7], which stops when all nodes belonging to its frontier are adjacent to a bridge (community boundary).

Algorithm 1 LocalCommunityDetection

Require: Seed node $s \in G = (V, E)$

Require: Community mapping $g_C : V \rightarrow \mathbf{P}$

Require: Bridge function $b : E \rightarrow [0.0, 1.0]$

```

1:  $C_s = \emptyset$ 
2: Frontier set  $F = \{s\}$ 
3: while  $|F| > 0$  do  $\{F$  is non-empty $\}$ 
4:    $c \leftarrow F.pop()$ 
5:    $C_s \leftarrow C_s \cup \{c\}$ 
6:    $C_U \leftarrow C_U \setminus \{c\}$ 
7:   for all  $n \in N(c)$  such that  $e_{cn} = (c, n) \in E$  do
8:     if  $g_C(n) = C_U$  and  $b(e_{cn}) \leq B_L$  then
9:        $F.push(n)$ 
10:    end if
11:  end for
12: end while
13:  $\mathbf{P} \leftarrow \mathbf{P} \cup C_s$ 

```

The quality of the community structure output by Bridge Bounding is entwined with the success of quantifying the bridging behavior of edges. Let us consider the function $b : E \rightarrow [0, 1]$, which maps edges to real numbers in the given interval, to quantify the extent to which they act as bridges. In order for Bridge Bounding to make a binary decision on whether an edge e is a bridge or not (in order to stop or continue the community flooding process along this edge), the output of the bridging function, $b(e)$, is compared against some threshold B_L (which can be derived by analysis of the distribution of $b(e)$ values as will be shown later).

The problem of quantifying the bridging behavior of edges on a graph has been already studied and several measures based on graph topology have been developed with the goal of capturing the extent to which an edge acts as a bridge between different communities. One of the first attempts to define $b(e)$ was by means of its *betweenness centrality* as described in [23]. For a given edge e , its betweenness centrality is defined as the fraction of shortest paths running along the edge, $\sigma_{st}(e)$ to the number of all possible shortest paths σ_{st} between s and t .

$$b_{\Phi}(e_{st}) = \Phi(e_{st}) = \sum_{s \neq t \in V} \frac{\sigma_{st}(e)}{\sigma_{st}} \quad (1)$$

An extension to this measure, called *bridging centrality*, appeared in [16]. Bridging centrality was defined as the rank product of the edge betweenness (Equation 1) and the edge *bridging coefficient*, which made use of the local network topology to quantify the extent to which an edge acts as a bridge.

The measures of betweenness and bridging centrality are global bridging measures, i.e. they are computed by processing the whole graph. To reduce the computational requirements, one may consider local bridging measures, e.g. the *edge-clustering coefficient* [24]:

$$C_{st}^{(3)} = \frac{z_{st}^{(3)}}{\min[(d(s) - 1), (d(t) - 1)]} \quad (2)$$

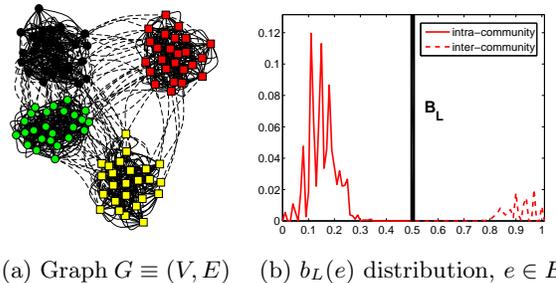


Figure 1: Relation of edge position in the graph and local bridging b_L probability distribution function (pdf). Edges drawn with dashed lines on the network of Figure 1(a) are also the ones with the highest local bridging values (the part of the distribution in Figure 1(b) plotted in dashed line).

where $z_{st}^{(3)}$ is the number of triangles containing that edge. Note that the larger the clustering coefficient is, the less the edge acts like a bridge. Hence, we define the *local bridging* of an edge as:

$$b_L(e_{st}) = 1 - C_{st}^{(3)} = 1 - \frac{|N(s) \cap N(t)|}{\min[(d(s) - 1), (d(t) - 1)]} \quad (3)$$

In order for $b_L(e)$ to have a low value, the two endpoints of e need to have a lot of common neighbors (relative to their degree). Effectively, this means that in order to move from one of the endpoints to the other, one has multiple options in addition to e . Thus, e is considered as an intra- (or within-) community edge. In the opposite case, when the two endpoints of a bridge have very few or no neighbors in common, then this edge is crucial for the connection between its endpoints. For that reason, we consider in the latter case, where $b_L(e)$ has a high value, that e is an inter-community edge or bridge.

In order to derive a decision threshold B_L for identifying the bridge edges of the graph (see line 8 of Algorithm 1), one needs to inspect the distribution of b_L values among the edges of the graph. Figure 1 illustrates how the position of edges on a graph with community structure affects their local bridging values. The graph of Figure 1(a) was generated to comprise a synthetic four-community structure. Edges that link different communities with each other, i.e. *inter-community* edges, are drawn in dashed line. According to the distribution of Figure 1(b), these edges are characterized by high b_L values, therefore they can be separated by means of thresholding from the intra-community edges.

The exact probability distribution function of b_L for a given graph is available only after computing the local bridging function for each edge of the graph, introducing in that way a global graph processing step in the Bridge Bounding methodology. However, this step does not impose severe restrictions on the computational process. First, according to Equation 3, the computation of b_L can be carried out in a streaming fashion, since only the neighborhoods of the two endpoints of each edge are required during the computation. To further reduce the computational requirements, it is possible to derive an approximation of the b_L probability distribution by computing the local bridging values of a small random subset of the network edges. Finally, one

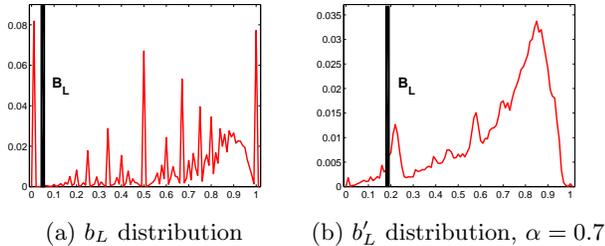


Figure 2: Distributions of first-order (b_L) and second-order (b'_L) local bridging on the English tag network of Section 4.2. Note that due to the b_L distribution shape, it is impossible to select a value for B_L such that less than 8% of the network edges are considered intra-community.

could even completely skip the distribution estimation step if it has been already performed for a graph of similar nature in the past (in which case one could reuse the previously estimated threshold B_L).

The simple measure of local bridging (Equation 3) employed by Bridge Bounding is ideal for networks with very clear community structure (such as the one of Figure 1(a)). However, the measure is often not well-suited for detecting communities in real networks. In particular, when a network is characterized by scale-free topology, the distribution of b_L values will have a *spiky* shape, similar to the one in Figure 2(a), where the depicted b_L distribution comes from the English LYCOS iQ tag network of Section 4.2. In such cases, it is hard to differentiate between bridge and non-bridge edges. For instance, according to Figure 2(a), 8% of the network edges have local bridging $b_L = 0$, thus $\forall B_L \geq 0$, Algorithm 1 will always consider 8% of the network edges as non-bridges. In networks with scale-free topology (which commonly emerge in practice), such a decision would cause Bridge Bounding to detect a community structure that consists of one large community and many *singleton* communities, i.e. communities comprising just one member. The reason for such an outcome is that scale-free networks maintain a large connected component even when a large fraction of their edges are removed⁴ [1]. Figure 3 illustrates the output of Bridge Bounding on a scale-free graph generated by the *preferential attachment* model of Barabási-Albert [3].

In order to alleviate this problem, we consider the 2^{nd} order local bridging of an edge e , $b'_L(e)$, by computing the weighted sum (with a mixing parameter α) of its local bridging, $b_L(e)$ and the mean local bridging of the edges constituting its neighborhood:

$$b'_L(e_{st}) = \alpha \cdot b_L(e_{st}) + (1 - \alpha) \frac{1}{|N(e_{st})|} \sum_{e \in N(e_{st})} b_L(e) \quad (4)$$

By applying Equation 4, we carry out a smoothing of the local bridging function by taking into account the values of the function in the neighborhood around a given edge. The α parameter defines the extent to which the values of the neighboring edges are taken into account in the compu-

⁴Although Bridge Bounding does not explicitly remove edges from the underlying network, it treats bridging edges as bounds, i.e. as non-existent.

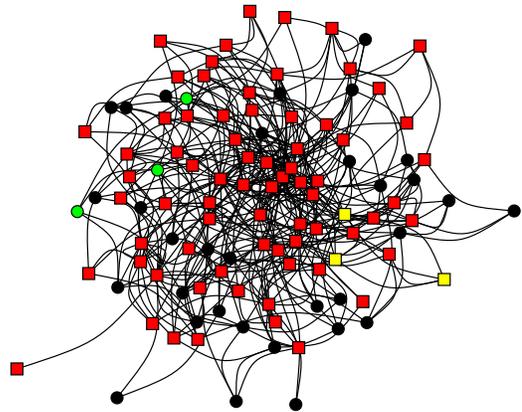


Figure 3: Community structure found by Bridge Bounding on a 100-node scale-free graph. The structure consists of one large (red squares), two small (green circles, yellow squares) and 32 singleton communities (black circles).

tation of b'_L . Figure 2(b) illustrates the distribution of b'_L (using $\alpha = 0.7$) for the LYCOS iQ English tag network of Section 4.2. Low b'_L values are distributed more evenly in comparison to the b_L ones. Hence, it is possible to select a value for B_L such that only a very-low fraction of edges are considered as intra-community ($\simeq 1\%$ in this example).

Effectively, the computation of 2^{nd} order local bridging makes use of topological information from a wider neighborhood around a given edge in comparison to local bridging. Following this, one could consider the ν^{th} order local bridging, $b_L^{(\nu)}$, which for sufficiently high values of ν , would utilize topological information from the whole graph. Obviously, since the computation of $b_L^{(\nu)}$ is carried out in an iterative manner, the complexity of computing the measure increases with its order ν .

In terms of complexity, a progressive global community detection scheme based on Bridge Bounding is decomposed in two steps: (a) local network topology function computation and (b) community detection. Computing the basic local bridging measure for a graph of n nodes and m edges with average node degree \bar{d} has a complexity of $O(\bar{d}^2 \cdot m)$ since for each edge, we need to find the intersection of two sets of average size \bar{d} .⁵ The community detection step has a complexity of $O(\bar{d} \cdot n)$, when Algorithm 1 is used in the global community detection scheme described in 3.1, since for each node of the network \bar{d} candidate nodes are considered as candidates for assignment to the community that is currently being created. Thus, in total, Bridge Bounding scales with $O(\bar{d}^2 \cdot m + \bar{d} \cdot n)$.

4. EVALUATION

In this section, we present a series of experiments we carried out in order to gain insights into the performance of the proposed approach. The first part of the experiments

⁵For the computation of higher-order local bridging, the complexity raises to $O(\nu \cdot \bar{d}^2 \cdot m)$. However, we consider that most applications of Bridge Bounding will make use of second- or at most third-order local bridging functions.

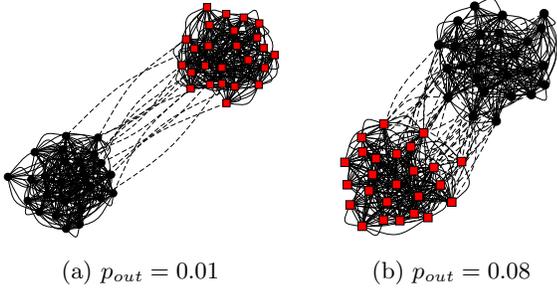


Figure 4: Sample synthetic mixtures of communities generated using the same set of parameters $\{N = 50, K = 2, z_{tot} = 18\}$ but different values for p_{out} .

compares the performance of progressive global community detection based on Bridge Bounding with the one achieved by the community detection method of Girvan-Newman [15]. This comparison is carried out on synthetic networks with known (predefined) community structure, thus giving the possibility for objective measurement of the method performance. In the second part of the experiments, we aim at gaining insights into real-world complex networks. Therefore, we apply our community detection technique on two networks created from the user tagging activities in the English and German version of the LYCOS iQ question / answering application. Since there is no ground truth concerning the community structure of the LYCOS iQ tag network, we use our subjective judgement in order to draw conclusions on the performance of the proposed method.

4.1 Synthetic networks

We created a parameterized community mixture generator following the strategy described in [23] and [20]. According to this, the generation process results in a network with N nodes which consists of K communities. We control the average degree z_{tot} of the network nodes, as well as the probability p_{out} that a node's edge will connect to a node of a different community. Thus, out of the z_{tot} edges of each node (on average), $z_{out} = p_{out} \cdot z_{tot}$ edges connect the node to nodes of different communities. Obviously, higher values of p_{out} will lead to networks with less profound community structure. Figure 4 depicts the difference in the conspicuity of community structure in relation to the fraction p_{out} of inter-community edges. This network generation process can be described by a four-element parameter set comprising N , K , z_{tot} , and p_{out} . We also consider a fifth parameter, namely the *community size variation* s_{var} , which is calculated as the ratio of the biggest community size to the size of the smallest one. In this case, each community C_i will have a different average node degree z_{tot}^i and therefore we define $z_{tot} = \frac{1}{K} \cdot \sum_{i=1}^K z_{tot}^i$. In the end, we consider the five-element parameter set:

$$S_{PAR} = \{N, K, z_{tot}, p_{out}, s_{var}\} \quad (5)$$

Two widely used measures to evaluate the effectiveness of data partitioning methods, e.g. community detection, when the *true* partition structure is known (which is the case when testing with synthetic networks) are (a) the fraction F_C of correctly classified instances [23] and (b) the *Normalized*

Mutual Information (NMI) introduced in [12] and applied for the evaluation of community detection in [20]. Consider two partitions of the n -node graph, $\mathbf{P}^a = \{C_0^a, C_1^a, \dots, C_{K_a}^a\}$ (true community structure) and $\mathbf{P}^b = \{C_0^b, C_1^b, \dots, C_{K_b}^b\}$ (community structure found by algorithm). The fraction F_C of correctly classified instances is straightforward to compute only when $K_a = K_b = K$. When the true number of communities K_a differs from the number of communities K_b found by the algorithm, we need to first identify a subset of the found communities $\mathbf{P}_c^b \subseteq \mathbf{P}^b$, that can be matched to a subset of the true communities, $\mathbf{P}_c^a \subseteq \mathbf{P}^a$. We consider two communities as matching if they present overlap of more than 50%. Then, assuming that community $C_x^a \in \mathbf{P}_c^a$ is the matching community of $C_i^b \in \mathbf{P}_c^b$, F_C is computed by the following equation.

$$F_C = \frac{1}{n} \cdot \sum_{C_i^b \in \mathbf{P}_c^b} |C_x^a \cap C_i^b| \quad (6)$$

The Normalized Mutual Information between the true partition, \mathbf{P}^a , and the one found by the algorithm, \mathbf{P}^b , quantifies the extent to which they are similar to each other from an information-theoretic point of view [12].

$$\text{NMI}(\mathbf{P}^a, \mathbf{P}^b) = \frac{-2 \cdot \sum_{i=1}^{K_a} \sum_{j=1}^{K_b} n_{ij}^{ab} \log\left(\frac{n_{ij}^{ab} \cdot n}{n_i^a \cdot n_j^b}\right)}{\sum_{i=1}^{K_a} n_i^a \log\left(\frac{n_i^a}{n}\right) + \sum_{j=1}^{K_b} n_j^b \log\left(\frac{n_j^b}{n}\right)} \quad (7)$$

In Equation 7, n_i^a and n_j^b denote the number of nodes in communities C_i^a and C_j^b respectively, and n_{ij}^{ab} denotes the number of shared nodes between communities $C_i^a \in \mathbf{P}^a$ and $C_j^b \in \mathbf{P}^b$. In general, NMI is preferred to the simplistic F_C measure, since it handles gracefully the cases where $K_a \neq K_b$. F_C is presented here together with NMI mainly due to the ease in its interpretation.

To demonstrate the effectiveness of Bridge Bounding in detecting the underlying community structure of networks, we compare the performance of the progressive global community detection scheme (see Section 3.1) based on Bridge Bounding in terms of both F_C and NMI to the performance of the community detection method by Girvan and Newman (GN) [15] on a multitude of synthetic networks. Since the GN method employs a divisive approach, it results in a hierarchical community structure, which contains multiple graph partitions to communities. Therefore, we needed to select a single partition from the hierarchy, which we would use to evaluate the performance of the method. The strategy used by Newman and Girvan in [23] to make this selection is to calculate the *modularity* Q of each partition and select the partition which maximizes it.

The modularity of a network partition into K communities is calculated from the $K \times K$ symmetric matrix \mathbf{e} whose element e_{ij} is the fraction of all edges in the network that link vertices in community i to vertices in community j . Further, we define the row (or column) sums $\alpha_i = \sum_j e_{ij}$ which represent the fraction of edges that connect to vertices in community i . Based on the above definitions, the measure of modularity is defined as:

$$Q = \sum_i (e_{ii} - \alpha_i^2) \quad (8)$$

This quantity measures the fraction of edges in the net-

Table 1: Comparison of performance between a global scheme based on Bridge Bounding with local bridging (BB), Bridge Bounding with 2nd order local bridging (BB’) and the method of Girvan and Newman (GN) [15]. The performance is measured on synthetic networks generated using the set $S_{PAR}^1 = \{200, 4, 40, p_{out}, 1.0\}$ of parameters, with p_{out} being the free parameter.

p_{out}	F_C			NMI		
	BB	BB’	GN	BB	BB’	GN
0.01	100	100	100	1.0	1.0	1.0
0.05	100	100	100	1.0	1.0	1.0
0.1	100	100	50	1.0	1.0	0.86
0.15	100	99	50	1.0	.98	0.86
0.20	99	74	50	0.98	0.84	0.86
0.25	24	24	0	0.54	0.56	0.02

Table 2: Similar comparison of performance as in Table 1, but on synthetic networks that were generated using the set $S_{PAR}^2 = \{200, 4, 40, 0.01, s_{var}\}$ of parameters, with s_{var} being the free parameter.

s_{var}	F_C			NMI		
	BB	BB’	GN	BB	BB’	GN
1.1	100	100	100	1.0	1.0	1.0
1.5	100	100	100	1.0	1.0	1.0
1.6	99.5	100	100	0.99	1.0	1.0
1.7	88	98	100	0.82	0.96	1.0
1.8	85.5	97	100	0.79	0.95	1.0
1.9	58.5	87	90	0.68	0.82	0.88
2.0	12.5	80	82	0.45	0.73	0.81
2.5	0	62	75	0.45	0.63	0.72

work that connect vertices of the same community (i.e. intra-community edges) minus the expected value of the same quantity in a network with the same community partition but random connections between the vertices. If the number of intra-community edges is no better than random, we would get $Q = 0$. For perfect separation to communities (i.e. communities that are completely disconnected from each other on the graph), we get $Q = 1$. In practice, modularity values in the range from 0.3 to 0.7 indicate significant community structure.

We created two sets of networks containing synthetic communities. The first set of such networks was generated holding the four network generation parameters of Equation 5 constant and varying p_{out} . This is a widely adopted test process [23, 24, 20] to test the performance of a community detection method as the communities of the synthetic graph gradually become less well-separated. Table 1 presents the comparison between the performance of Bridge Bounding (by use of both first- and second-order local bridging) and the GN method [15]. Both Bridge Bounding methods present equally good or better performance than GN across the range of p_{out} values that were used for testing.

A further test involved the generation of an additional set of networks by varying the s_{var} parameter in order to end up with networks comprising communities of unequal sizes. Table 2 provides an overview of the results obtained from the three methods of our study. Apparently, the use of the local bridging function (Equation 3) becomes problematic

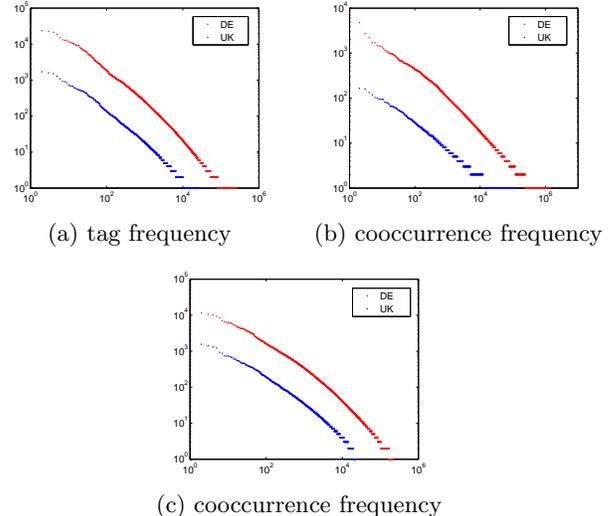


Figure 5: Rank plots of tag, cooccurrence frequencies and node degrees for the German and English LYCOS iQ tag networks.

for Bridge Bounding as soon as the size variation among the underlying communities exceeds a certain value (e.g. for $s_{var} \geq 2$, we measured $NMI(BB) < 0.5$). In contrast, Bridge Bounding with the use of 2nd order local bridging as well as the GN method yielded consistently better results in this series of tests. Hence, it becomes clear that the use of more sophisticated local topology measures, such as the 2nd order local bridging, could be crucial for the success of the proposed method.

4.2 LYCOS iQ tag network

LYCOS iQ is a collaborative question/answering application where people ask and answer questions on any topic. The application is available in six languages, German, English, French, Danish, Swedish and Dutch with German attracting the largest community of users. In order to support the users’ efforts of searching for relevant questions, the application incorporates a tagging functionality, similar to the one used in typical *social tagging systems* such as delicious⁶ and flickr⁷. There are no static categories and tags are not predefined by the system, but the users’ inputs are checked against tags existing in the system database to prevent duplicates.

Question submitters have the possibility of attaching more than one tag to each of their questions. Therefore, it is possible to create a tag network from the collaborative tagging activities of users. In this network, the vertex set comprises the tags chosen by users to tag their questions and the edge set contains the co-occurrences between tags in the users’ questions. When a question is tagged with more than two tags, then all possible pairwise co-occurrences are added to the network. For each tag of the network, its frequency (tf) is available. Further, the co-occurrence frequency (cf) between each pair of tags is available.

Figure 5 illustrates the rank plots of tag and cooccurrence frequencies as well as of the node degrees observed

⁶<http://delicious.com>

⁷<http://flickr.com>

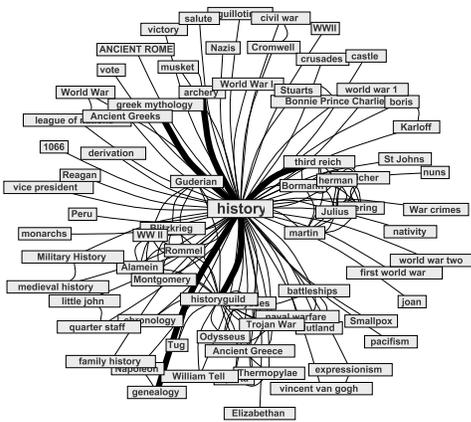


Figure 9: Community around tag “history”.

ure 10(b)), (b) the questions submitted by LYCOS iQ users (and consequently the tags used to describe them) are more focused to particular aspects of music, e.g. pop music artists.

Further, a noteworthy observation regarding the structure of the communities around “film” (Figure 10(c)) and “animals” (Figure 10(d)) is the existence of small cliques (between 3 and 5 members) within them. Those correspond to tags related to particular films in the “film” community (e.g. “batman”-“Christian Bale”-“comic”) or tags related to groups of animals (e.g. “leopards”-“panthers”-“mammals”) in the “animals” community. This indicates the existence of semantic hierarchies within topics (e.g. “mammals” are a subclass of “animals”; “leopards”, “panthers” are a subclass of “mammals”), which could be further validated by means of machine learning techniques [30].

As stated earlier, detecting the topic communities within a tag network, similar to the one created from LYCOS iQ application (nowadays, there are plenty of Web 2.0 applications incorporating collaborative tagging characteristics), can be beneficial for both the users and the administrators of the application. Users can be provided with a community view of the tags that are related to their context. For instance, when a LYCOS iQ user submits a question to the system, the text of her question can be parsed and matched against the tags already available in the system. Then, by identifying the community (or communities) that her question belongs to, it is possible to recommend relevant tags for use as descriptors of the question or relevant questions that have been tagged with tags belonging to the respective community. Further, administrators of such applications could use community detection in the context of a content monitoring and trend tracking framework for supporting the operation of important administrative tasks, e.g. online ad targeting or content moderation (which is most frequently synonymous to spam detection).

5. CONCLUSIONS

We introduced Bridge Bounding, a local methodology for community detection in large networks. The methodology is based on the notion of local network topology functions to quantify the extent to which edges act as community boundaries, i.e. bridges. We showed that use of local bridging, a topology function based on the widely used edge cluster-

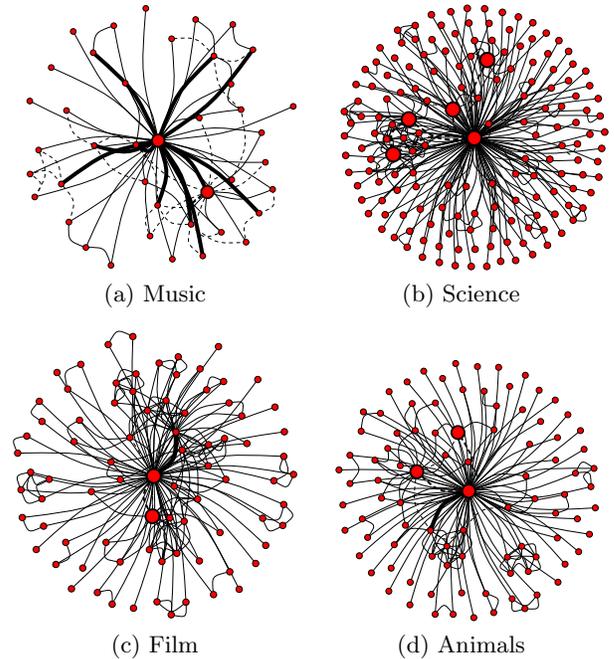


Figure 10: Further examples of community shapes. The presented communities were created using “music”, “science”, “film” and “animals” as seed nodes.

ing coefficient, resulted in successful discovery of existing community structure in synthetic networks, but failed to do so in networks of scale-free topology. For that reason, we employed the second- and higher-order local bridging functions to derive smoother estimates of the bridging properties of edges. The proposed methodology is extremely efficient, scaling with $O(\bar{d}^2 \cdot m + \bar{d} \cdot n)$ for networks of n nodes and m edges with average node degree \bar{d} .

A series of tests on synthetic networks with controlled community structure provides evidence that the Bridge Bounding method (with use of the 2^{nd} order local bridging function) performs equally well or better than the widely used method of Girvan and Newman. Moreover, application of our method on two large tag networks coming from the LYCOS iQ question/answering application proved beneficial in studying the underlying topic structure and can benefit both users and administrators of Web 2.0 applications with social tagging features.

In the future, we plan to carry out more thorough evaluation tests on the tag communities produced by Bridge Bounding. Specifically, we plan to conduct a user study among selected LYCOS iQ users in order to derive manual judgements on the quality of the detected communities. Subsequently, we are going to consider the potential of new edge bridging functions and of more sophisticated strategies for community detection based on Bridge Bounding. Instead of the currently employed fixed-threshold strategy for deciding whether an edge is intra- or inter-community, we will test the potential of adaptive threshold strategies. Finally, we intend to look into extensions that will endow the method with capabilities for uncovering hierarchical relations within the community structure.

6. ACKNOWLEDGMENTS

This work was supported by the WeKnowIt project, partially funded by the European Commission, under contract number FP7-215453. We would also like to acknowledge the use of the JUNG⁸ framework for parts of our implementation. Finally, we would like to acknowledge the use of the English and the German tag data sets from the LYCOS iQ application operated by LYCOS Europe.

7. REFERENCES

- [1] R. Albert, H. Jeong, and A.-L. Barabasi. Error and attack tolerance of complex networks. *Nature*, 406(6794):378–382, July 2000.
- [2] J. Bagrow and E. Boltt. A local method for detecting communities. *Physical Review E*, 72:046108, 2005.
- [3] A.-L. Barabasi and R. Albert. Emergence of scaling in random networks. *Science*, 286(5439):509–512, 1999.
- [4] C. H. Brooks and N. Montanez. Improved annotation of the blogosphere via autotagging and hierarchical clustering. In *WWW '06: Proceedings of the 15th international conference on World Wide Web*, pages 625–632, New York, NY, USA, 2006. ACM.
- [5] C. Cattuto, A. Baldassarri, V. D. P. Servedio, and V. Loreto. Investigating community structure in social tagging systems. *Advances in Complex Systems (ACS)*, 11(4):597–608, August 2008.
- [6] A. Clauset, M. E. J. Newman, and C. Moore. Finding community structure in very large networks. *Physical Review E*, 70:066111, 2004.
- [7] L. da F. Costa. Hub-based community finding, 2004.
- [8] W. de Nooy, A. Mrvar, and V. Batagelj. *Exploratory Social Network Analysis with Pajek (Structural Analysis in the Social Sciences)*. Cambridge University Press, 2005.
- [9] Y. Dourisboure, F. Geraci, and M. Pellegrini. Extraction and classification of dense communities in the web. In *WWW '07: Proceedings of the 16th international conference on World Wide Web*, pages 461–470, New York, NY, USA, 2007. ACM.
- [10] N. Du, B. Wu, X. Pei, B. Wang, and L. Xu. Community detection in large-scale social networks. In *Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 workshop on Web mining and social network analysis*, pages 16–25, New York, USA, 2007. ACM.
- [11] G. W. Flake, S. Lawrence, and C. L. Giles. Efficient identification of web communities. In *KDD '00: Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 150–160, New York, NY, USA, 2000. ACM.
- [12] A. L. N. Fred and A. K. Jain. Robust data clustering. *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2:128–136, 2003.
- [13] E. Garfield. *Citation Indexing: Its Theory and Application in Science, Technology, and Humanities*. John Wiley & Sons, New York, NY, USA, 1979.
- [14] D. Gibson, J. Kleinberg, and P. Raghavan. Inferring web communities from link topology. In *HYPERTEXT '98: Proceedings of the ninth ACM conference on Hypertext and hypermedia*, pages 225–234, New York, NY, USA, 1998. ACM.
- [15] M. Girvan and M. E. J. Newman. Community structure in social and biological networks. *Proceedings of the National Academy of Sciences*, 99(12):7821–7826, June 2002.
- [16] W. Hwang, T. Kim, M. Ramanathan, and A. Zhang. Bridging centrality: graph mining from element level to group level. In *KDD '08: Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 336–344, New York, NY, USA, 2008. ACM.
- [17] H. Ino, M. Kudo, and A. Nakamura. Partitioning of web graphs by community topology. In *WWW '05: Proceedings of the 14th international conference on World Wide Web*, pages 661–669, New York, NY, USA, 2005. ACM.
- [18] R. Kumar, P. Raghavan, S. Rajagopalan, and A. Tomkins. Trawling the web for emerging cyber-communities. *Computer Networks*, 31(11-16):1481–1493, 1999.
- [19] J. D. L. Danon, A. Diaz-Guilera and A. Arenas. Community structure identification. *arXiv*, cond-mat/0505245v1, 2005.
- [20] J. D. L. Danon, A. Diaz-Guilera and A. Arenas. Comparing community structure identification. *Journal of Statistical Mechanics*, P09008, 2005.
- [21] M. E. J. Newman. The structure and function of complex networks. *SIAM Review*, 45(2):167–256, 2003.
- [22] M. E. J. Newman. Fast algorithm for detecting community structure in networks. *Physical Review E*, 69:066133, 2004.
- [23] M. E. J. Newman and M. Girvan. Finding and evaluating community structure in networks. *Physical Review E*, 69:026113, 2004.
- [24] F. Radicchi, C. Castellano, F. Cecconi, V. Loreto, and D. Parisi. Defining and identifying communities in networks. In *Proceedings of the National Academy of Science of the United States of America*, volume 101, pages 2658–2663, March 2004.
- [25] J. P. Scott. *Social Network Analysis: A Handbook*. SAGE Publications, 2000.
- [26] B. Shneiderman and A. Aris. Network visualization by semantic substrates. *IEEE Transactions on Visualization and Computer Graphics*, 12(5):733–740, 2006.
- [27] Y. Song, Z. Zhuang, H. Li, Q. Zhao, J. Li, W. C. Lee, and C. L. Giles. Real-time automatic tag recommendation. In *SIGIR '08: Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 515–522, New York, NY, USA, 2008. ACM.
- [28] H. White and K. McCain. Bibliometrics. *Annual Review of Information Science and Technology*, 24:119–186, 1989.
- [29] F. Wu and B. A. Huberman. Finding communities in linear time: A physics approach. *arXiv*, cond-mat/0310600, 2003.
- [30] M. Zhou, S. Bao, X. Wu, and Y. Yu. An unsupervised model for exploring hierarchical semantics from social annotations. In *ISWC/ASWC2007: Proceedings of the 6th International Semantic Web Conference and 2nd Asian Semantic Web Conference*, pages 673–686, Heidelberg, 2007. Springer Verlag.

⁸<http://jung.sourceforge.net>