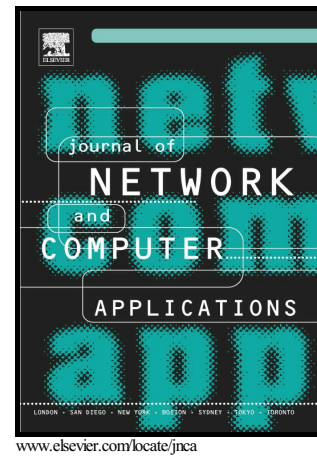


Hypergraph partitioning for social networks based on information entropy modularity

Wenyin Yang, Guojun Wang, Md Zakirul Alam Bhuiyan, Kim-Kwang Raymond Choo



PII: S1084-8045(16)30232-6
DOI: <http://dx.doi.org/10.1016/j.jnca.2016.10.002>
Reference: YJNCA1729

To appear in: *Journal of Network and Computer Applications*

Received date: 16 May 2016
Revised date: 25 September 2016
Accepted date: 4 October 2016

Cite this article as: Wenyin Yang, Guojun Wang, Md Zakirul Alam Bhuiyan and Kim-Kwang Raymond Choo, Hypergraph partitioning for social networks based on information entropy modularity, *Journal of Network and Computer Applications*, <http://dx.doi.org/10.1016/j.jnca.2016.10.002>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting galley proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



Hypergraph partitioning for social networks based on information entropy modularity¹

Wenyin Yang^{a,b}, Guojun Wang^{c,a,*}, Md Zakirul Alam Bhuiyan^{d,a}, Kim-Kwang Raymond Choo^{e,f}

^a*School of Information Science and Engineering, Central South University, Changsha 410083, China*

^b*School of Electronic and Information Engineering, Foshan University, Foshan 528000, China*

^c*School of Computer Science and Educational Software, Guangzhou University, Guangzhou 510006, China*

^d*Department of Computer and Information Sciences, Temple University, Philadelphia, PA 19122, USA*

^e*Department of Information Systems and Cyber Security, University of Texas at San Antonio, San Antonio, TX 78249, USA*

^f*School of Information Technology & Mathematical Sciences, University of South Australia, Adelaide, SA 5095, Australia*

Abstract

A social network is a typical scale-free network with power-law degree distribution characteristics. It demonstrates several natural imbalanced clusters when it is abstracted as a graph, and expands quickly under its generative mechanism. Hypergraph is superior for modeling multi-user operations in social networks, and partitioning the hypergraph modeled social networks could ease the scaling problems. However, today's popular hypergraph partitioning tools are not sufficiently scalable; thus, unable to achieve high partitioning quality for naturally imbalanced datasets. Recently proposed hypergraph partitioner, hyperpart, replaces the balance constraint with an entropy constraint to achieve high-fidelity partitioning solutions, but it is not tailored for scale-free networks, like social networks. In order to achieve scalable and high quality partitioning results for hypergraph modeled social networks, we propose a partitioning method, EQHyperpart, which utilizes information-Entropy-based modularity Q value (EQ) to direct the hypergraph partitioning process. This EQ considers power-law degree distribution while describing the "natural" structure of scale-free networks. We then apply simulated annealing and introduce a new definition of hyperedge cut, micro cut, to avoid the local minima in convergence of partitioning, developing EQHyperpart into two specific partitioners, namely: EQHyperpart-SA and EQHyperpart-MC. Finally, we evaluate the utility of our proposed method using classical social network datasets, including Facebook dataset. Findings show that EQHyperpart partitioners are more scalable than competing approaches, achieving a tradeoff between modularity retaining and cut size minimizing under balance constraints, and an auto-tradeoff without balance constraints for hypergraph modeled social networks.

© 2011 Published by Elsevier Ltd.

Keywords: Scale-free network, Social network partitioning, Hypergraph partitioning, Information entropy, Modularity

¹This is an extended version of the conference paper [1], with more than 50% new content.

1. Introduction

Social networks, a source of big data [2], face complicated scalability challenges in data management, due to the large numbers of data nodes and complex data relationships. Such challenges are compounded in online social networks. **Scaling out** is generally employed to address scaling problems by deploying data over several servers, such as virtual machines in the cloud, and partitioning the query loads between these servers [3]. Hash-based partitioning approach is widely used in commodity systems, such as dynamo and Cassandra. However, neglecting data relations still leads to heavy query loads [4, 5]. To improve partitioning quality, approaches based on modeling social network structure and user interactions have been proposed [3, 6, 7]. Generally, social networks are abstracted by graphs, and the query loads partitioning problem can be reduced to graph partitioning problems.

Graph partitioning comprises **dyadic** graph partitioning [8, 9] and hypergraph partitioning [10, 11, 12]. While there exist numerous partitioning solutions based on usual graphs, or dyadic graphs, there has been limited focus on the use of hypergraph partitioning. Recent studies have argued that hypergraphs are more powerful than dyadic graphs for modeling groups in many domains [7, 13, 14, 15], including partitioning and other operations on complex networks [7, 14]. Therefore, in this paper, we mainly focus on the partitioning of hypergraph modeled social networks. For instance, social network users could be denoted by vertices of a hypergraph, and the multi-user operation, such as multi-gathering of data from multiple social network (e.g., Facebook) friends in a single operation, could be modeled via a hyperedge in the hypergraph. Inter-server query costs of multi-way interactions are defined as the cut cost of hypergraphs. The hypergraph partitioning problem is to partition the vertices of a hypergraph into k disjoint nonempty equal-size partitions, such that the number of the hyperedges connecting vertices in different partitions (called the cut) or the cut size is minimized. An example of undirected hypergraph partitioning is depicted in Fig. 1. Here, a hypergraph $H = (V, N)$ with vertex set $V = \{v_1, v_2, \dots, v_{14}\}$ and hyperedge set $N = \{n_1, n_2, \dots, n_5\}$ is partitioned into four parts, namely: $T = \{T_1, T_2, \dots, T_4\}$, where the cut size is 4 since hyperedges n_1, n_2, n_3 , and n_4 are in cut state. The values in vertex circles denote the weight of vertices, assigned as the degree of vertex in this hypergraph.

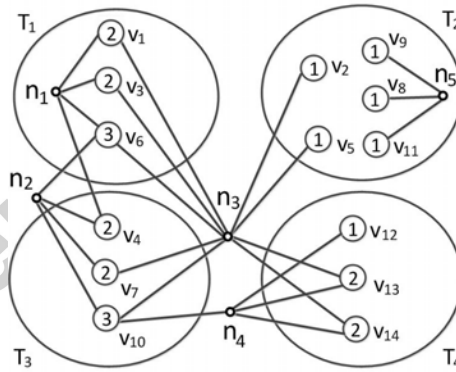


Fig. 1: An example of hypergraph partitioning. A hypergraph with 14 vertices and 5 hyperedges is partitioned into 4 parts, leading to a cutsize of 4.

Min-cut graph partitioning [9, 10, 11] and modularity optimization [16, 17, 18, 19] are popular network partitioning solutions. Min-cut graph partitioning algorithms are generally designed to produce equal sized partitions such that the number of inter-partition edges or hyperedge is minimized, and modularity optimization algorithms aim to maximize the modularity (see [3, 18]), a.k.a. Q value, which is the most

*Corresponding author

Email addresses: cswyyang@163.com (Wenyan Yang), csgjwang@gmail.com (Guojun Wang), zakirulalam@gmail.com (Md Zakirul Alam Bhuiyan), raymond.choo@fulbrightmail.org (Kim-Kwang Raymond Choo)

common method and measure designed to estimate the strength of dividing a network into modules. Networks with a high Q value have dense connections between the nodes within modules, but sparse connections between nodes in different modules. However, the social network is a typical scale-free network characterized by power-law distribution of degrees, community structures and quick expansion under its generative mechanism [20], and existing approaches are generally inadequate for the hypergraph partitioning of it. For example,

- Community structures of scale-free networks are not fully considered. Popular min-cut hypergraph partitioning tools, such as hMETIS [10] and khMETIS [11], provide load balancing while performing min-cut partitioning. However, min-cut and load balance are not necessarily the only constraints for all applications [12], including complex networks partitioning. Social networks possess typical “quasi-balanced” datasets with natural clusters, and the cluster structures are foundation for the expansion of networks; thus, it is important to also consider community structures to ensure scalability.
- Modularity optimization algorithms are not well adapted to changing graph structures when social networks grow rapidly. Specifically, existing partitioning algorithms based on community detection are sensitive to graph structure [17], and the node placement solutions are unstable. The statistical feature of complex systems could be used to model the modularity; thus, mitigating this sensitivity [21].
- Most hypergraph partitioners do not consider the power-law distribution of degrees. For example, the partitioners in hyperpart, Yaros2013hyperpart use information entropy as a constraint for hypergraph partitioning. Compared to equal-sized partitioners, such as hMETIS and khMETIS, hyperpart is able to produce high-fidelity partitions for imbalanced datasets. In other words, if the dataset is composed of different sizes of clusters due to natural characteristics, then the fixed-part partitioning result could reserve the natural group members within a part to the utmost extent, even though it is imbalanced in size. However, hyperpart is not suitable for scale-free networks, including social networks, because it uses an information entropy that treats every node equally. This does not allow one to reflect the degree distribution feature of scale-free networks.

To address the above gap, we proposed the use of a scale-free network featured information entropy in our previous work [1], which considers the vertex degrees rather than the vertex amounts of communities in the hyperpart to reflect the energy distribution in the network system. Since the distribution rules of entropy coincide with the vertex degree distribution, we used Q values computed based on the information entropy to depict the community structures of scale-free networks, and to direct the min-cut hypergraph partitioning process, producing high-fidelity partition results. The statistic characteristics of this information entropy could overcome the drawback due to unstable assignment. The preliminary experiment findings detailed in [1] demonstrate that this partitioning can achieve tradeoff between modularity retaining and cut size minimizing of hypergraph modeled complex networks.

We also remark that a partition method designed to get tradeoff among load balance, minimum cut cost, and modularity, is more scalable in the long run than existing min-cut hypergraph partitioning algorithms (e.g., hMETIS and khMETIS). Therefore, while partitioning, we seek to address the three above mentioned factors to ensure scalability. In addition, to avoid the local minima in the search for optimum values of Q value and cut cost, simulated annealing algorithm and micro cut are adopted for optimization. We propose *micro cut* as a measure for hypergraph cut state, which provides micro hints to hypergraph partitioners for facilitating a view into the future moves.

The contributions of this paper are two-fold.

- (1) We put forward a social network partitioner, EQHyperpart, which utilizes information Entropy based modularity Q (EQ) to direct the low cost partitioning process. It considers modularity maintaining, cut size minimizing and balance factors during partitioning to improve scalability.
- (2) We propose a new hyperedge cut metric, micro cut, and optimize the partitioning quality of EQHyperpart by adopting simulated annealing and micro cut, forming two variant partitioners, namely: EQHyperpart-SA and EQHyperpart-MC.

We then evaluate the partition quality of our proposed partitioning algorithm with competing hypergraph partitioners, such as hMETIS, khMETIS, and hyperpart. Findings demonstrate that EQHyperpart is more scalable, and EQHyperpart achieves an auto-tradeoff between cut size, modularity, and balance level.

The organization of the remainder of this paper is as follows. In Section 2, we introduce background knowledge and related work. Next, we present the definition of scale-free information entropy and EQ in Section 3. Section 4 describes our proposed EQHyperpart algorithm and micro cut. Section 5 presents our evaluation findings. Finally, in Section 6, we conclude this paper and outline future work.

2. Background and related work

2.1. Social networks modeling

A social network structure comprises a set of social actors, sets of dyadic ties, and other social interactions between actors. Existing related research on social networks could be classified into two categories, namely: one is to analyze the social network structures and study the theories explaining the patterns observed in these structures [22, 23, 24]; and the other is to solve existing problems based on social networks [25, 26, 27, 28, 29, 30, 31]. Modeling the social network as a graph, namely social graph, is an effective and widely-used mathematical tool in existing studies.

Generally, social networks could be modeled as undirected graphs or directed graphs according to their properties [22]. Classical undirected social graphs consist of friendship and interaction graphs [23]. The former indicates logical friendship between social network users, and the latter reflects visible interaction in social network, such as wall posts and other messages. Typical directed social graphs include latent graph [24] and following graph. The former models passive actions of social network users (e.g., profile browsing), while following graph reveals the subscription relationship in social networking applications such as Twitter.

In addition, social networks could be categorized into uniform and weighted models. Uniform social network model treats every edge equally, i.e. users access any of their friends' data in equal probability. Weighted social network models assign diverse meanings to edges based on the underlying motivations.

Moreover, social graph could be divided into dyadic social graph and social hypergraph. Dyadic social graphs only reflect the two-way relations of social network users via edges, while social hypergraphs can capture multi-way relations, such as multi-way interactions, multi-user static relationships, and multi-user dynamic operations, via hyperedges [7]. Hypergraphs are generally more effective in modeling groups, but the concept of a group depends on the social networking context [13]. For example, one hyperedge in Fig.1 could represent one group of social network users who interact closely via recent wall postings, or those have the same group label, like “colleague”, “classmate”, etc.

2.2. Information entropy

Entropy was first introduced as a thermodynamic concept in 1872 [32], and subsequently used in information theory [33]. The macro significance of entropy is a measure of the uniformity of system energy distribution, representing the object state as stable or not stable. The information entropy (also known as Shannon entropy) is used to characterize the uncertainty about the source of information, and increases with more sources of greater randomness. The source can be characterized by the probability distribution of its samples. When taken from a finite sample, Shannon defined the entropy H of a discrete random variable X with possible values $\{x_i | i = 1, \dots, N\}$ and probability mass function $P(x_i)$ as (1), where b is the base of the logarithm, and the unit of entropy is Shannon (defined by IEC 80000-13) for $b = 2$. Shannon entropy has been applied in modeling of complex systems [34] and Big Data applications [35], etc.

$$H(X) = - \sum_{i=1}^N P(x_i) \log_b P(x_i) \quad (1)$$

2.3. Hypergraph partitioning

A hypergraph is the generalization of a graph, where a hyperedge, a.k.a net, can connect a group of vertices. The concept of undirected hypergraph was first given by C. Berge in 1970s [36], and directed hypergraph theories were subsequently developed. Partitioning is an important concept in both graph and hypergraph theories. Given a hypergraph H , k -way partitioning assigns vertex set V of H to k disjoint nonempty partitions, aiming to minimize a given cost function of such an assignment (see Fig.1).

Pins allocation in very-large-scale integration (VLSI) design and computation load distribution in parallel computing are typical applications of hypergraph partitioning, in which the number of parts are known *a priori*, and balance constraint are preferred [10]. Hypergraph partitioning of social network could be adopted for user data allocation among fixed number of servers to improve inter-user data access performance [7, 14].

With certain constraints such as balance, the problem of optimally partitioning a hypergraph is known to be NP-hard [37]. There has been a number of heuristic algorithms with near-linear runtime presented in the literature. Kernighan-Lin (KL) algorithm [38] is the pioneer of both graph partitioning and community structure, which enables vertex swaps in an equal-sized bisection to reduce its cut. The strict equipartition requirement is relaxed to a balance constraint by the Fiduccia-Mattheyses (FM) algorithm [39], which introduces single-vertex move rules and enables k -way partitioning based on recursively bi-partitioning. Later, a direct k -way partitioning based on FM algorithm is put forward [39]. However, execution time increases as graph size grows. Karypis [9] proposed a multilevel framework for hypergraph partitioning, including coarsening, partitioning and uncoarsening phases. FM and / or its derivatives are applied in the middle phase. Partitioning tools, hMETIS [10] and khMETIS [11], were then developed to implement the multilevel framework via recursive bisection and direct k -way approaches, respectively. Following KL algorithm [38], several direct k -way partitioning methods based on FM algorithm [39] and a multilevel framework [9] for hypergraph partitioning are proposed. For example, hMETIS and khMETIS are designed to implement the multilevel framework via recursive bisection and direct k -way approaches, respectively. UMPa [40] is a multi-objective hypergraph partitioner also using multilevel. This framework is effective in reducing both execution time and cut size, but is limited to low levels of imbalance.

In hMETIS, unbalance factor (UBfactor) denoted by b is used to control deviation of each part size, within the upper bounce $n \cdot ((50 + b)/100)^{\log_2 k}$ and lower bounce $n \cdot ((50 - b)/100)^{\log_2 k}$, where n is the number of vertices and k is the number of parts. While in khMETIS, the heaviest part size is up to $(1 + b/100) \cdot (n/k)$, where b is also the unbalance factor, n is the number of vertices and k is the number of parts, indicating that the weight of the heaviest partition should not be more than $b\%$ greater than the average weight. For example, in Fig.1, the unbalance factor is $b = 10$ under hMETIS metric while $b = 43$ under khMETIS metric, assuming that each having unit vertex weight.

Yaros J. R. proposed the use of information entropy in (2) as an imbalance constraint [12], such that it enables the proposed partitioner, hyperpart, to find high-fidelity solutions for given levels of imbalance. In other words, the partitioning result is more similar to the actual partition result. However, $P(x_i)$ in (1) is assigned as $|V_i|/|V|$ in (2), where k denotes the number of parts, and $|V_i|$ denotes the number of vertex in part i . This implies every node is treated equally and the node importance is ignored, which can be reflected by node degrees. Therefore, such an approach cannot be applied to a scale-free complex network directly, which is characterized by the power-law degree distribution.

$$E_u = - \sum_i^k \frac{|V_i|}{|V|} \log_b \left(\frac{|V_i|}{|V|} \right) \quad (2)$$

Hyperedge cut is a standard cost function measuring the partition quality. Cut size refers to the number of nets in cut state, that is spanning more than one partition. In addition, the number of parts a cut net spans is referred to as the Sum of External Degrees (SOED). The similar $K - 1$ measure has penalty of parts spanned outside the base part. For example, in Fig.1, the cut sizes of net n_2 and n_3 are both 1, but SOED of n_2 and n_3 are 2 and 4, respectively. Thus, the $K - 1$ cut sizes of them are 1 and 3, respectively. hMETIS aims to minimize cut size, while khMETIS is designed to optimize $K - 1$ cut size or SOED. The typical objective

of hypergraph partitioning is to minimize inter-partition communication, and in this context, $K - 1$ cut size is the most suitable metric.

2.4. Community detection

Complex networks are naturally divided into subgroups. Originating from the small-world concept, the mutuality of ties, the frequency of ties among subgroup members, or the closeness of subgroup members could lead to a community within complex network. The community structure detector seeks to find those subgroups of complex networks. Modularity is a widely used metric to evaluate the detection capabilities of complex networks. The value of modularity, i.e. Q value, is defined in (3), where k denotes the number of communities, e_{ii} is the ratio of the number of edges inside community i to the total number of edges in the whole network, and e_{ij} is the ratio of the number of edges between community i and j to the total number of edges in the whole network. To increase Q value, inner-connection should be higher and inter-connection should be lower.

$$Q = \sum_i^k (e_{ii} - (\sum_j^k e_{ij})^2) \quad (3)$$

Current detection methods include minimum-betweenness [41], clustering centrality [42], random walk [43], eigenvectors of matrices [44], target function optimization [45], and methods based on density of edge connection [46]. Detection speed [47] and detecting overlapping community [48] are the focus of recent research in recent times, and a number of information-theoretic methods have been introduced. Examples include detection methods based on information-theoretic entropy [21], minimum description length (MDL) mutual information [49], and information bottle [50]. These methods are used in dyadic graph modeled complex networks. A review of the literature suggests that community structure detection in hypergraph modeled networks remains a topic that is understudied. One of the few work on hypergraph modeled networks is that of Xie [51], who proposed to view hyperedges as vertices and establish a network for hyperedges by their similarity. This modularity method is then applied to find communities for document words association.

A key limitation of community structure detection methods for both usual graphs [16, 21, 45, 47, 48] and hypergraph [51] is the possibility that no good division of the network exists. However, the goal of graph partitioning is to find the best division of the network regardless of whether a good division exists [18]. Despite this limitation, community structure detection methods could be adapted to find a good division.

3. EQ in hypergraph modeled social networks

In this section, we describe the social networks modeled using hypergraph in Section 3.1 prior to introducing the definitions of E (i.e. information entropy reflecting characteristics of scale-free networks) and EQ (i.e. modularity based on the Q value in hypergraph) in Section 3.2 and Section 3.3, respectively.

3.1. Hypergraph modeled social networks

In our work, social networks are modeled using hypergraph, i.e. vertices in hypergraph may represent the social network users, and hyperedges may indicate the relationships among users. Diverse context could be assigned to weights of vertices or hyperedges, such as workload and access cost. Hub nodes with higher degree in social network may indicate owning more friends or being more active. New nodes with little degree tend to link the hub nodes according to the generation mechanism of social networks. Interconnections within a group become closer than those among groups, forming community structures gradually. Therefore, the degree of nodes, or generally, the weight of vertices, implies the probability of joining a community. Hypergraph partitioning is one of the key approaches with extensive application in mining of subgroups. Minimum cut of hypergraph is more intuitive and effective than that of dyadic in some cases [14].

3.2. Definition of entropy for social networks

In scale-free complex networks, including social networks, the degree distribution follows the power-law, indicating that the energy is unevenly distributed in these networks, which is also known as *nonhomogeneity*. The latter property demonstrates that a scale-free network is a kind of *ordered network*, while scaled network, such as random network, belongs to *disordered network*. Entropy is used to quantify the property of *order*.

Existing studies on *Barabási – Albert* (BA) model [52] suggest that the scale-free network is caused by the growing and preferential attachment mechanisms that the new nodes preference to connect with hub nodes. To be specific, when a new node joins a scale-free network, the probability of node i chosen to be connected by the new node is decided by the degree of node i . Therefore, the degree of a node can be a baseline metric to reflect the importance of a vertex in a network. Therefore, the importance of vertex i can be defined by (4):

$$I_i = \frac{d_i}{\sum_{i=1}^N d_i}, \quad (4)$$

where N is the number of vertices in the network, and d_i is the degree of vertex i . We assume that $d_i > 0$; thus, $I_i > 0$. The importance of vertices is different in an ordered network. In a scaled network, however, the importance of vertices is roughly equivalent, and this is why a scaled network is called a disordered network. To quantitatively measure the order, network structure entropy is defined by (5):

$$E = - \sum_{i=1}^N I_i \log I_i \quad (5)$$

It is trivial to prove that when the network is completely uniform (i.e. $I_i = 1/N$), E reaches the peak. When all the vertices connect to one hub vertex, say the first vertex (i.e. $d_1 = N - 1$, $d_j = 1$ ($j > 1$)), E falls to the bottom, because the network is the most nonuniform.

3.3. Definition of Q value based on entropy

As mentioned in Section 3.1, energy distributes unevenly in an ordered scale-free network, which is associated with community structures. It can be inferred that the energy concentrates inside the community, leading to uneven distribution of the energy and resulting in the (obvious) community structure. We suggest using Entropy-based Q value (EQ) to describe the community structure property for a scale-free network. The denser within communities and the sparser among communities, the greater Q value will become and the more obvious the community structure will be. Thus, we define EQ as the difference between Community Structure Entropy (CSE) and Inter-Community Entropy (ICE).

3.3.1. Community Structure Entropy (CSE)

Based on the preferential attachment mechanisms of scale-free network, it can be inferred that when the community structure is formed, the new node chooses a community to join that also follows the preferential mechanisms. In other words, the importance, say the quantity, of a community determines the probability of new node's accession.

Let $Y = (y_1, y_2, \dots, y_{|Y|})$ be a variable of community, and $X = (x_1, x_2, \dots, x_{|X|})$ be a variable of node. We define CSE as a conditional entropy $H(Y | X)$ by (6) to measure the uncertainty or the disorder situation of X , on the condition of already existing Y , where N is the total number of vertices, M is the total number of communities, and $P(x_i, y_j) = P(y_j|x_i)P(x_i)$ is the joint probability.

$$H(Y|X) = - \sum_{i=1}^N \sum_{j=1}^M P(x_i, y_j) \log P(y_j|x_i) = - \sum_{i=1}^N \sum_{j=1}^M P(y_j|x_i)P(x_i) \log P(y_j|x_i) \quad (6)$$

In the community structure entropy, $P(x_i)$ denotes the importance of a node i , which is evaluated by (4), and $P(y_j|x_i)$ represents the probability that community j contains node i . As we have previously discussed,

the probability that a node joins a community is determined by the importance of the community. Therefore, $P(y_j|x_i)$ is defined as m_j/N , the proportion of node numbers of community j to the total node numbers of the whole network. Thus, CSE is calculated by (7), where Z_{ij} is an assignment matrix. If node i is assigned to community j , then $Z_{ij} = 1$, otherwise 0.

$$E_{CS} = - \frac{\sum_{i=1}^N \sum_{j=1}^M Z_{ij} d_i * (\frac{m_j}{N}) * \log(\frac{m_j}{N})}{\sum_{i=1}^N d_i} \quad (7)$$

3.3.2. Inter-Community Entropy (ICE)

ICE refers to the uncertainty among communities in scale-free networks. We focus on the ICE of hypergraph modeled scale-free network in this paper. It is well known that multilevel partitioning framework is popular in many hypergraph partitioners, which comprises three phases, namely: coarsening, partitioning, and uncoarsening.

Hypergraph is coarsened by merging vertices and / or edges using some heuristics algorithms. In an extreme case, a vertex-level hypergraph could be coarsened to a community-level one. Inspired by this, the communities can be viewed as super-vertices after coarsening, and the association among communities can be considered to be the cut hyperedges among communities. Hence, ICE can be defined as (8), where C is the current $K - 1$ cut size of the hypergraph modeled network and M is the total number of communities.

$$E_{IC} = -C * (\frac{1}{M}) * \log \frac{1}{M} \quad (8)$$

3.3.3. Entropy-based Q value

According to the idea of modularity introduced in Section 2.4, based on the definition of CSE and ICE, the entropy-based Q value EQ is defined by (9):

$$EQ = E_{CS} - (E_{IC})^2 \quad (9)$$

4. Hypergraph partitioning based on EQ

In this section, we first discuss the basic idea of EQHyperpart algorithm. Then, we present two optimization methods for the local minima avoidance. Here, Section 4.1 introduces the optimized algorithm by adopting Simulated Annealing (SA), named EQHyperpart-SA, in detail. The definition of micro cut and the optimization approach based on micro cut, forming EQHyperpart-MC, are discussed in Section 4.2.

4.1. Algorithm

Note that, in our previous work [1], we presented the basic idea of a hypergraph partitioning algorithm based on EQ, named as EQHyperpart, which is designed based on the idea of khMETIS, the k -way counterpart of hMETIS, but the minimum cuts size metric is replaced by the maximum EQ.

According to the single-vertex move rules, in each iteration, EQHyperpart selects the vertex with the highest gain and the highest delta-EQ-value to move and freeze, until all the vertices are frozen. Then, EQHyperpart rollbacks to the point with the highest EQ value and unfreezes all vertices and starts the new move operations again. The best partition solution search terminates when the EQ does not increase in the last round or the iteration number exceeds a threshold.

This basic partitioning solution, however, is easy to converge to a local optima in the solution space. We choose SA, a probabilistic technique for approximating the global optimum of a given function, to solve this problem. SA starts from an admissible solution of the problem (denoted by S in (10)). The search strategy in the neighbourhood of such a solution will be more intensive in the more promising regions, penalizing the searches that move far from these regions but accepting with certain probability P_s , defined in (10), and searches that worsen the solution (denoted by S' in (10)). Here, the function $y \equiv E(x)$ refers to the

energy level of the given solution, t is the temperature and k is a constant. The temperature t , which is the acceptance criterion governed by a random number generator and a control parameter (a.k.a cooling ratio r), slowly modifying its value by $t = r * t$ ($0 < r < 1$), drives the system toward the final solution, which corresponds to a local minimum of the objective function [53].

$$P_s = k * \exp\left(-\frac{E(S') - E(S)}{t}\right) \quad (10)$$

We adapt SA in EQHyperpart to search for the highest EQ value. When EQ value does not increase in the last round, the algorithm rolls back to the point where the second highest EQ value is seen with the probability defined by (10). The EQHyperpart based on SA (EQHyperpart-SA) comprises a sequence of operations depicted in Algorithm 1.

Firstly, after randomly distributing the pins of vertices, we compute the EQ value of the network, and the possible gains for each vertex (line 1). At the beginning of each iteration, we unfreeze all vertices to be ready for move (line 4). Then, the algorithm enters the inner while loop (lines 5-10). In this loop, we select the best move according to the compound conditions, including gain value, incremental EQ value, and the unbalance ratio (lines 6-7). After performing the move of vertex v from FP part to the TP part and locking vertex v (line 8), we update the gain values and pin distributions in an incremental manner (line 9) and compute the new EQ value of the whole network (line 10). Then, the EQ_{new} is assigned the largest EQ within this loop (line 11) and the corresponding move index is recorded. An iteration ends when the whole hypergraph is frozen, and we compare the highest EQ value from this iteration with the EQ value from the last one. If EQ_{new} is no higher than EQ_{old} , we unwind sequences of executed moves back to the point where the partition with EQ_{new} is seen, in certain probability, and assign the EQ_{new} to EQ_{old} (lines 13-17). In other case, we rollback to the point when the highest EQ value occurred (lines 18-24). The temperature is updated after the rollback (line 25).

If the number of iteration exceeds the predetermined value, or we achieve the same highest EQ value for certain times, indicating that no increase of EQ is possible for any further move, the outer loop terminates (line 26). Finally, the partitioned result is output.

Besides partitioning under certain balance constraint, EQHyperpart-SA enables maximum modularity partitioning without any constraint. In other words, we can ignore the unbalance constraints in line 6 of Algorithm 1, and achieve a relative reasonable partitioning result, which automatically obtains a tradeoff among the lowest cut, modularity, and balance level. Findings outlined in Section 5 verify the correctness of this function.

4.2. Micro cut

To speed up the partition of hypergraph, k -way FM-based partitioners, such as khMETIS, move vertex based on the gain which is defined by the cut cost benefits affected by critical nets. Usually, the cut size or $K - 1$ cut size is used as the gain measure. It leads to the well-known issue of convergence to local optima for this kind of partitioners.

Take net n_3 in Fig. 1 as an example. This net owns 9 vertex pins, and spans four parts initially; thus, the $K - 1$ cut size of n_3 is 3. It is clear that net n_3 is not critical to any part, because no move operation performed on any pin of net n_3 can change its cut-state (i.e. the $K - 1$ cut size). Consider all vertex pins of n_3 , the move gains measured by $K - 1$ cut size are zero (the best gain possible), from the located part to any other parts. In other words, the partitioner cannot see any benefit in moving any pin. In this case, the partitioner is stuck at this local optima, and makes the best move randomly.

To reduce the nearsightedness of partitioners, we define a measure for cut state, namely *micro cut* (see 11). This metric provides partitioners with some micro hints to inform future moves.

$$\sum_{e \in E} \left[1 - \sum_{\pi \in \Pi} \left(\frac{|\pi \cap e|}{\sum_{\pi \in \Pi} |\pi \cap e|} \right)^{(1+\alpha)} \right] \quad (11)$$

Here, E is the set of nets, Π is the set of partitions, and $\alpha > 0$. $|\pi \cap e|$ denotes the number of pins located in part π owned by the net e .

Algorithm 1 Algorithm of EQHyperpart-SA

Input: hypergraph $HG = (V, N)$, part number K , unbalance factor ε , temperature t , cooling ration r , iteration control parameter c , l

Output: partitioned result $P = (P_1, P_2, \dots, P_K)$

- 1: Initialize pin distributions, compute EQ value and gains for all possible moves from each vertex's current part to $(K - 1)$ other parts.
- 2: Set $EQ_{new} = EQ$, $EQ_{old} = EQ$
- 3: **repeat**
- 4: Unfreeze all vertices.
- 5: **while** there is any valid move **do**
- 6: $HGainList \leftarrow$ Select the highest gain moves that do not violate unbalance constraints.
- 7: $BestMove(v, FP, TP) \leftarrow$ Select the highest delta-EQ-value move in $HGainList$.
- 8: Move vertex v from FP part to TP part, and freeze v .
- 9: Update the gains of unfreezed neighbours of v and the pin distributions.
- 10: Update EQ .
- 11: **end while**
- 12: **if** $EQ > EQ_{new}$ **then**
- 13: Draw a random number $y \in (0, 1)$
- 14: **if** $y < \exp(-(E(EQ_{new}) - E(EQ_{old}))/t)$ **then**
- 15: Rollback to the point when the non-highest EQ value EQ_{new} is seen.
- 16: Set $EQ_{old} = EQ_{new}$
- 17: **else**
- 18: Rollback to the point when the highest EQ value EQ_{old} is seen.
- 19: **end if**
- 20: **else**
- 21: Rollback to the point when the highest EQ value EQ_{new} is seen.
- 22: Set $EQ_{old} = EQ_{new}$
- 23: **end if**
- 24: $t = r * t$
- 25: **until** EQ does not increase for c times or iteration number exceeds threshold l .

Fig. 2 illustrates the move of net n_3 under the indication by micro cut. With the micro cut, there is a positive gain in moving v_2 . The micro cut sizes of n_3 listed in Table 1 show that moving v_2 from parts T_2 to T_1 produces the highest gain, and guide the partitioner to take this move. With this move, net n_3 becomes critical to part T_2 . At this point, the partitioner is “enlightened”, and moves v_5 to part T_1 , reducing the $K - 1$ cut size from 3 to 2.

In essence, the micro cut metric encourages moving pins from a part with less pins to one with more pins owned by the same net, leading to the decrease of final gains in the $K - 1$ cut size metric. In view of this, the micro cut could be applied when there is a tie in highest gain move candidates, or take the place of $K - 1$ cut size when calculating the move gains. On the basis of EQHyperpart-SA, we choose the former usage of micro cut, forming a new version of EQHyperpart, named EQHyperpart-MC (EQHyperpart based on Micro Cut).

As shown in Algorithm 2, EQHyperpart-MC combines two cut size metrics to choose proper moves (line 6-7) for further local optima avoidance. In the following sections, the term *EQHyperpart* refers to two variant partitioners, namely: EQHyperpart-SA and EQHyperpart-MC.

5. Experiments

In this section, we evaluate the partitioning quality of EQHyperpart partitioners by evaluating the algorithms using several real-world datasets, in comparison to hMETIS, khMETIS, and hyperpart. Note

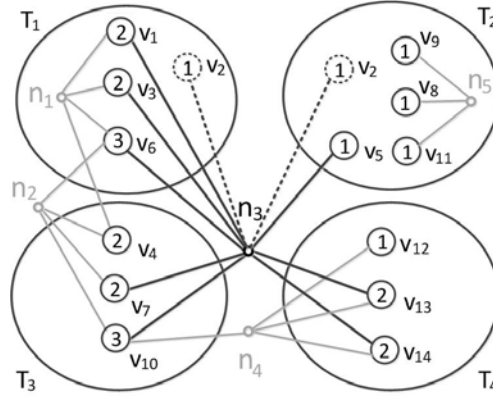


Fig. 2: Move of net n_3 by the hints of move gains under micro cut metric.

Table 1: Move Gain Calculation Example. (Take vertex v_2 as an example, with $\alpha = 1.1$.)

Type of Cut Size	T_2	$T_2 \rightarrow T_1$	$T_2 \rightarrow T_3$	$T_2 \rightarrow T_4$			
	Before Move	After Move	Move Gain	After Move	Move Gain	After Move	Move Gain
$K - 1$ Cut	3	3	0	3	0	3	0
Micro Cut	0.1278	0.1186	0.0092	0.1223	0.0055	0.1223	0.0055

that experiments are performed by EQHyperpart-SA and EQHyperpart-MC unless stated otherwise.

5.1. Experiments setup

The evaluations are three-fold. Firstly, we evaluate the scalability of EQHyperpart and traditional graph partitioning methods by measuring the increment speed of cross-partition query cost with the growth of network size. Secondly, several metrics are adopted to measure the partitioning quality of EQHyperpart, including the $K - 1$ cut size reducing extend, modularity retainment, and their tradeoff findings. Finally, the auto-tradeoff ability between cut size, modularity retainment, and unbalance level are evaluated without balance constraints in EQHyperpart.

We consider four datasets in the experiment. Three classical real-world social network datasets (i.e. Karate Club, Dolphin social network, and American College Football) are obtained from Mark Newman's personal website [54] and UCI Network Data Repository [55]. These three datasets offer the actual community numbers and partition details. Another dataset is part of Facebook data offered by SNAP (Stanford Network Analysis Project) [56, 57]. It reflects the earlier stage of a large social network. The Fast Girvan-Newman (FGN) algorithm [45] is performed on the Facebook dataset to detect communities, which act as the actual community structure for evaluation of hypergraph partitioning qualities.

These four datasets are modeled using hypergraph prior to partitioning execution, according to the inherent characteristic or inner-association of each dataset. For example, hyperedges are formed by the same associations, like friendship relationship, game relationship, and so forth. Moreover, the hypergraphs can be either unweighted or weighted, which take the degree as vertex weight. However, the hyperpart partitioner only supports unweighted hypergraph partitioning. Table 2 outlines the basic information of these datasets. Note that the unbalance level (UB-level) here uses the unbalance metric of khMETIS.

EQHyperpart is implemented in C++ using STL and Boost libraries, which is capable of processing weighted or unweighted hypergraph. It also allows the use of a hypergraph input file compatible with hMETIS. We perform the experiments on a Windows 2008 server, which has dual 2.0 GHz Intel Xeon processors with 8GB of RAM.

Algorithm 2 Algorithm of EQHyperpart-MC

Input: hypergraph $HG = (V, N)$, part number K , unbalance factor ε , temperature t , cooling ration r , iteration control parameter c, l

Output: partitioned result $P = (P_1, P_2, \dots, P_K)$

- 1: Initialize pin distributions, compute EQ value and gains based on $K - 1$ cut size for all possible moves from each vertex's current part to $(K - 1)$ other parts.
- 2: Set $EQ_{new} = EQ$, $EQ_{old} = EQ$
- 3: **repeat**
- 4: Unfreeze all vertices.
- 5: **while** there is any valid move **do**
- 6: $HGainList \leftarrow$ Select moves with the highest gain based on $K - 1$ cut size under unbalance constraints.
- 7: $BestMove(v, FP, TP) \leftarrow$ Select move with the highest gain based on micro cut size in $HGainList$.
- 8: Move vertex v from FP part to TP part, and freeze v .
- 9: Update the gains of unfreezed neighbours of v and the pin distributions.
- 10: Update EQ .
- 11: **end while**
- 12: **if** $EQ > EQ_{new}$ **then**
- 13: Draw a random number $y \in (0, 1)$
- 14: **if** $y < \exp(-(E(EQ_{new}) - E(EQ_{old}))/t)$ **then**
- 15: Rollback to the point when the non-highest EQ value EQ_{new} is seen.
- 16: Set $EQ_{old} = EQ_{new}$
- 17: **else**
- 18: Rollback to the point when the highest EQ value EQ_{old} is seen.
- 19: **end if**
- 20: **else**
- 21: Rollback to the point when the highest EQ value EQ_{new} is seen.
- 22: Set $EQ_{old} = EQ_{new}$
- 23: **end if**
- 24: $t = r * t$
- 25: **until** EQ does not increase for c times or iteration number exceeds threshold l .

5.2. Partitioner scalability experiments

As discussed in Session 3.2, the growing and preferential attachment mechanisms guide the generation of scale-free network. Before the next repartitioning point, the scale-free complex networks increase their sizes based on the latest partition results with the generative mechanisms during the interval. We seek to evaluate the scalability of partitioners under these mechanisms, that is, to measure the reduction rate of inter-partition query cost based on different initial placement configurations produced by different partitioners.

User nodes in Facebook form a typical scale-free network, such that they obey the generative mechanisms. In the experiments conducted on the Facebook dataset, the first 350 user nodes in this dataset are partitioned by respective partitioners to obtain the initial vertices placement layouts. Assuming that every node accesses its directed neighbour nodes only once, which are linked by the same hyperedge, followed by the calculation of average query cost of the whole network, denoted by C_{init} in (12). In the next placement phase, the successive nodes are placed in the proper partition according to the growing and preferential attachment mechanisms until the size of network reaches 1000. Finally, the average query cost is recalculated, denoted by C_{grown} in (12). The query cost saving rate R_{CS} is computed by (12), indicating the higher saving rate, the more scalable the partitioner.

$$R_{CS} = \frac{C_{init} - C_{grown}}{C_{init}} \quad (12)$$

Table 2: Real-world Dataset Characteristics.

Dataset	Vertices	Edges	Pins	Communities	$K - 1$ cut size	UB-level
Karate Club	34	78	190	2	11	12 16 (w)
Dolphin	62	159	380	2	10	33 41 (w)
American Football	115	616	1347	12	347	36 33 (w)
Facebook	348 ~ 4039	11422 ~ 88234	23192 ~ 180507	Unidentified	Unidentified	Unidentified

Fig. 3 depicts the query cost saving rate in the unbalance levels produced by hMETIS, khMETIS, and EQHyperpart-SA, respectively. It can be observed that EQHyperpart-SA is more scalable, and this is mainly due to the modularity retaining ability of EQHyperpart-SA and the generative mechanism which impels the node to join the partition with more related nodes.

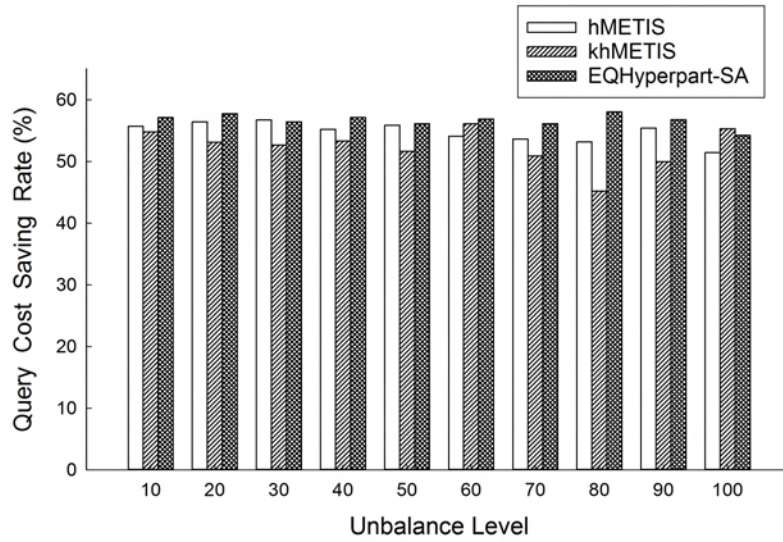


Fig. 3: Scalability comparison for partitioners.

5.3. Partitioning quality experiments

To evaluate the partitioning quality at different balance levels, five partitioners are performed on the unweighted and weighted social network hypergraphs, respectively, with different unbalance factors. For hMETIS, its unbalance factor (UBFactor) ranges from 1 to 50 in steps of 0.5. UBFactor of khMETIS, EQHyperpart-SA and EQHyperpart-MC ranges from 5 to 100 in steps of 1. The low entropy of hyperpart ranges from the high entropy to half the high entropy, in steps of 0.005. Finally, for each of these partitioners in each test, the partition findings selected for comparison are those with balance levels closest to the EQHyperpart.

5.3.1. $K-1$ cut size within balance constraints

Statistical results of two datasets, including American Football and Facebook, are shown in Fig. 4, where we display $K - 1$ cut size at different balance levels. We observe that EQHyperpart partitioners outperform other competing algorithms in the weighted network, because entropy based Q value is designed for nonuniform distributed scale-free networks. Specifically, EQHyperpart-MC significantly reduces $K - 1$ cut size, with the help of micro cut. Note that there are some data missed in findings of hMETIS and hyperpart, because they are not able to produce partitioning results at certain balance level requirements.

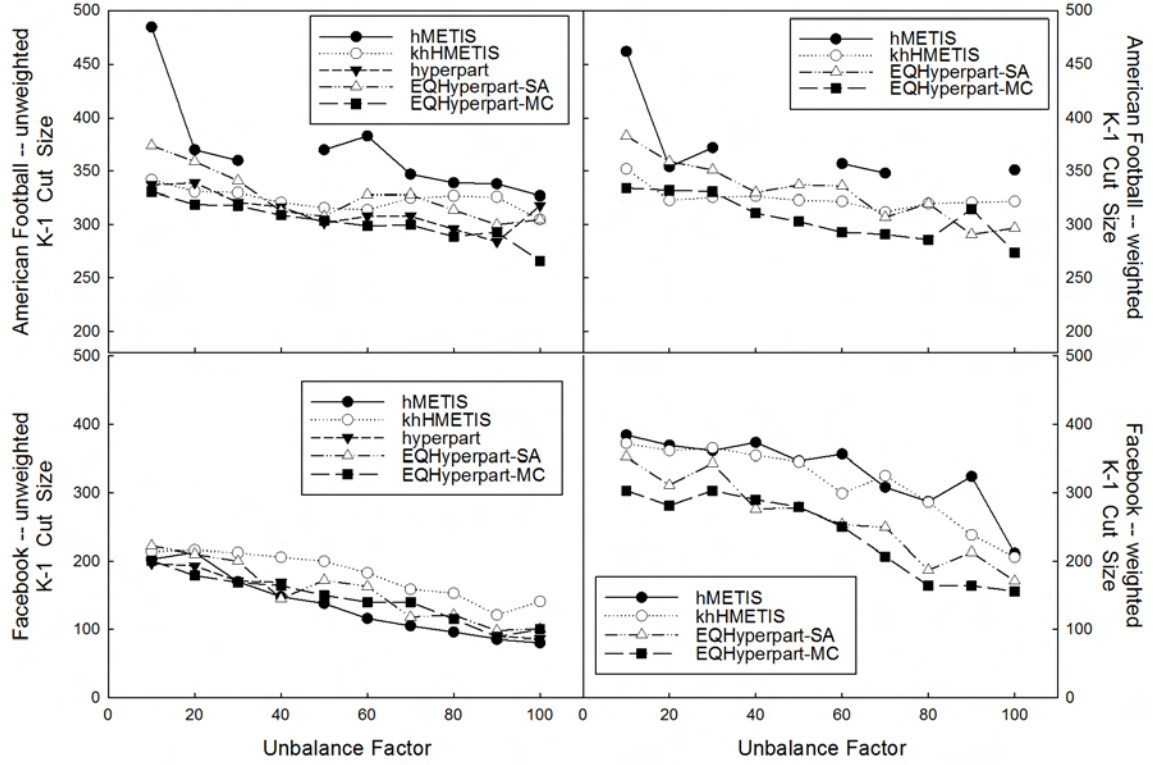


Fig. 4: Comparison of $K - 1$ cut size with balance constraints.

5.3.2. Modularity features retainment ability

To evaluate the retainment ability of modularity features, we adopt the typical metrics in information retrieval system, including recall rate, precision rate, and F-score, as evaluation measures.

First, we match the partition $t \in T$ to the natural community $c \in C$ with the maximum degree of overlapping. Then, a confusion matrix is built based on the part-community-match results. Given the confusion matrix, we estimate for each community $c \in C$ the following quantities, namely: $\alpha(c, T)$ is the number of vertices correctly assigned to c , $\beta(c, T)$ is the number of vertices incorrectly assigned to c , and $\gamma(c, T)$ is the number of vertices incorrectly not assigned to c . Then, the averaged recall is defined by (13) and the averaged precision is defined by (14). In order to consider these two measures, the weighted harmonic mean of them, named F-Score, or F-Measure is widely used. Wherein the $F1$ is the most common, calculated as (15).

$$R(T) = \frac{\sum_c \alpha(c, T)}{\sum_c \alpha(c, T) + \gamma(c, T)} \quad (13)$$

$$P(T) = \frac{\sum_c \alpha(c, T)}{\sum_c \alpha(c, T) + \beta(c, T)} \quad (14)$$

$$F1 = \frac{2 * P(T) * R(T)}{P(T) + R(T)} \quad (15)$$

The real world partitioning results (i.e. natural communities) of Karate Club, Dolphin, and American Football Team datasets are known beforehand. We achieve the natural communities of Facebook dataset by performing the FGN algorithm [45]. The part numbers correspond to the real community numbers, and the Facebook dataset is divided into 4 parts in this experiment.

Fig. 5 and Fig. 6 depict the F-Measure findings at different designated balance levels, produced by five partitioners (i.e. hMETIS, khMETIS, hyperpart, EQHyperpart-SA, and EQHyperpart-MC) on four unweighted hypergraph modeled datasets and four weighted datasets, respectively.

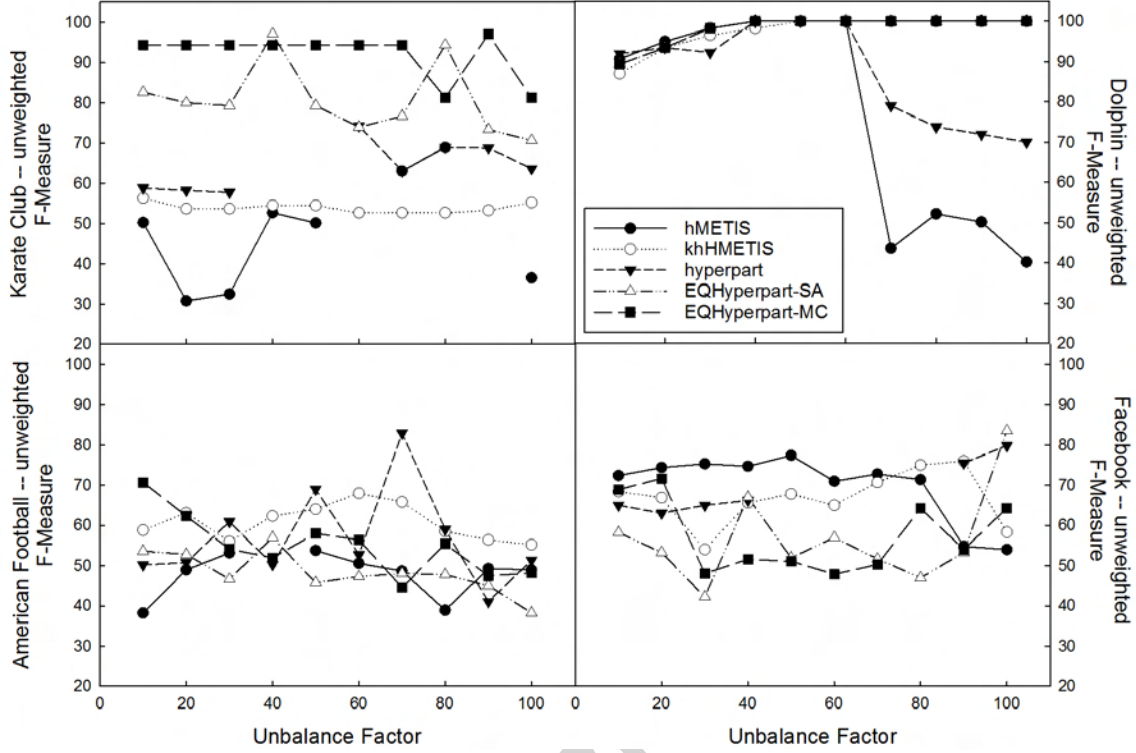


Fig. 5: F-Measure findings for unweighted datasets.

We also take the worse percentage of $K - 1$ cut size than true partition into consideration. This represents the deviation percentage of calculated $K - 1$ cut size to the one in real world, and provides another perspective for similarity evaluation. The lower the deviation percentage, the more similar the partitioning result is to the real world situation. Fig. 7 and Fig. 8 show the $K - 1$ cut size deviation relative to true partitions for unweighted and weighted datasets, respectively.

The findings for each dataset can be summarized as follows.

- In terms of Karate Club dataset, EQHyperpart partitioners, including both EQHyperpart-SA and EQHyperpart-MC, outperform in terms of modularity retaining ability, with the highest F-measure values and the closest $K - 1$ cut size, especially on unweighted dataset.
- In terms of Dolphin dataset, EQHyperpart partitioners achieve full F-score on stricter unbalance level demand than other partitioners, and maintain the community characteristics when the unbalance upper limitation requirement increases.
- In terms of American Football dataset, the performance of EQHyperpart partitioners is modest. According to the natural partitioning of American Football dataset, there exists overlapping communities. The EQ design in EQHyperpart does not consider this particular characteristic of datasets. This may result in the deterioration of partitioning quality on such datasets.
- In terms of Facebook dataset, EQHyperpart partitioners has a better performance on weighted dataset, with a higher F-measure and a lower $K - 1$ cut size deviation. Specifically, EQHyperpart-MC performs better than EQHyperpart-SA in general.

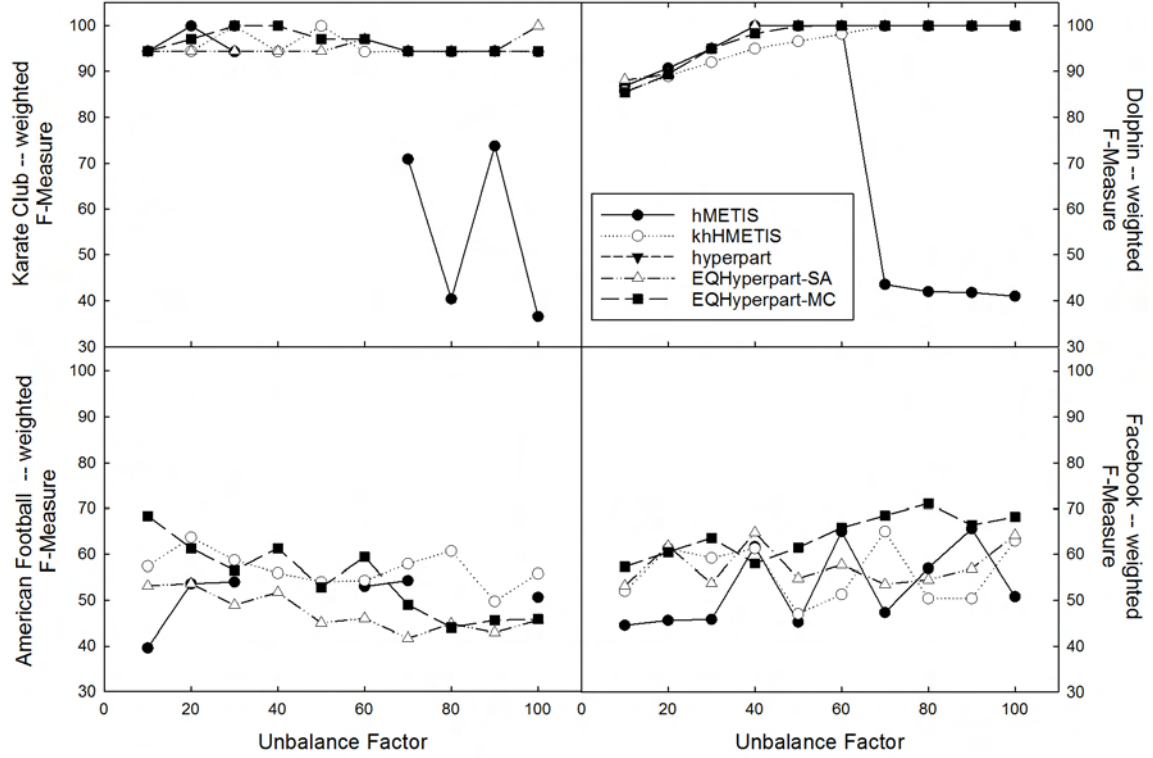


Fig. 6: F-Measure findings for weighted datasets.

We also noted from the findings that when the unbalance levels are close to the points of true partitions (e.g., 16 in Karate Club, 41 in Dolphin, and 33 in American Football weighted dataset), EQHyperpart is extremely close to the best partitioner in F-Measure, and the $K - 1$ cut sizes is extremely close to the real-world values. This implies that EQHyperpart partitioners are capable of maintaining the modularity, according to the nature characteristics of these networks. In addition, partition results produced by EQHyperpart on weighted datasets are better than those on unweighted datasets. This is particularly evident in Karate Club and Facebook datasets. The reason for this is that the weight reflects the node importance, which affects the generation of communities; thus, playing an important role in partitioning results of social networks.

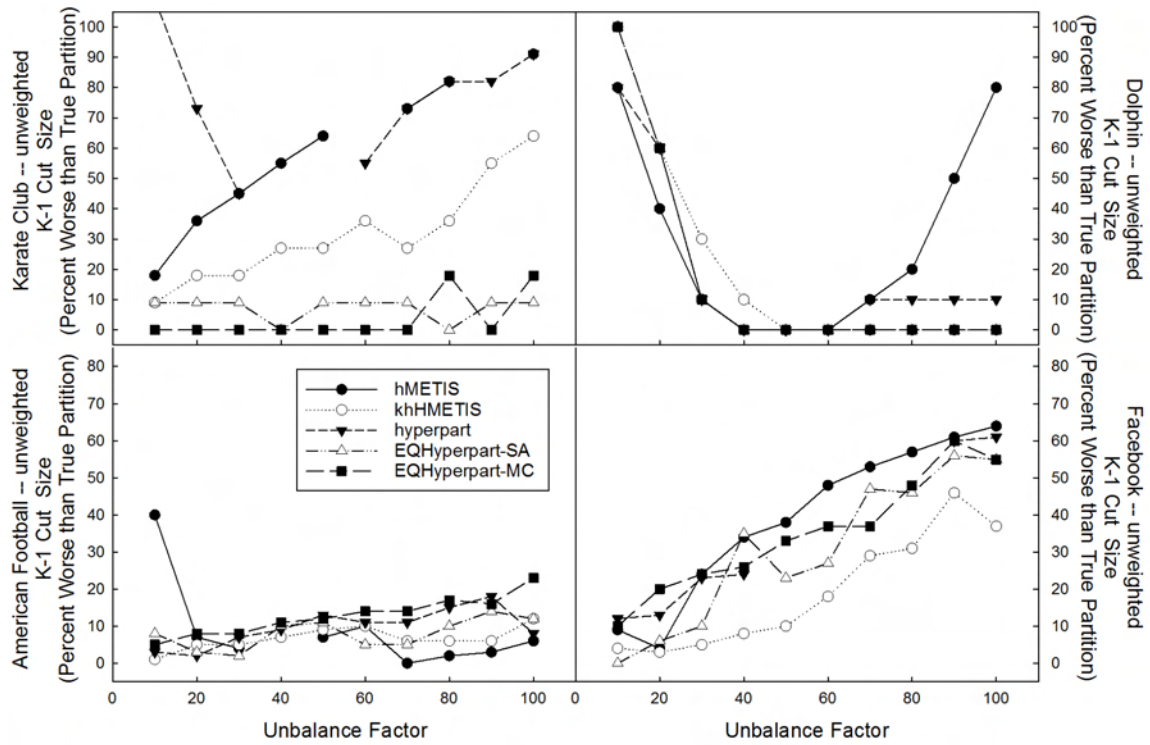
In summary, we have demonstrated that EQHyperpart retains the modularity characteristics of social networks under balance constraints.

5.3.3. Tradeoff between modularity and cut size

Because EQHyperpart partitioning favors the increment of modularity retaining ability, and the decrement of $K - 1$ cut size, we quantify the tradeoff between them under balance constraints using (16). Here, FS_i , CS_i and TO_i indicate the F-Score value, $K - 1$ cut size and the tradeoff value under unbalance level i , respectively. $|E|$ denotes the number of edges of the dataset, and a and b are weighting coefficients for the two factors, subjected to $a + b = 1$.

$$TO_i = a * FS_i * 100 + b * \left(1 - \frac{CS_i}{|E|}\right) * 100. \quad (16)$$

Table 3 shows the average tradeoff of each weighted dataset of the unbalance levels, ranging from 0 to 100 by a factor of 10, and the coefficients a and b are both assigned 0.5. It can be observed that the average

Fig. 7: $K - 1$ cut size deviation to true partition for unweighted datasets.

tradeoff value of EQHyperpart partitioners exceeds the value of the other three partitioners; thus, validating its tradeoff ability between modularity retaining and $K - 1$ cut size minimizing.

Table 3: Tradeoff findings.(Findings are rounded. (w) indicates the dataset is modeled using weighted hypergraph, otherwise, unweighted.)

	hMETIS	khMETIS	hyperpart	EQHyperpart-SA	EQHyperpart-MC
Karate Club	71	72	77	83	89
Dolphin	85	95	90	95	95
Football	44	54	53	48	53
Facebook	84	83	70	78	78
Karate Club(w)	81	91		90	91
Dolphin(w)	84	94		95	94
Football(w)	45	52		47	53
Facebook(w)	75	77		78	81

5.4. Auto-tradeoff partitioning experiments

Both EQHyperpart-SA and EQHyperpart-MC enable partitioning without balance constraint. In other words, EQHyperpart partitioners could terminate automatically at a point that results in an appropriate tradeoff between modularity, $K - 1$ cut size and the balance level.

General partitioning favors the decrement of $K - 1$ cut size and unbalance factor to obtain low query cost and load balancing. We argue that the modularity retaining ability is also important in social network

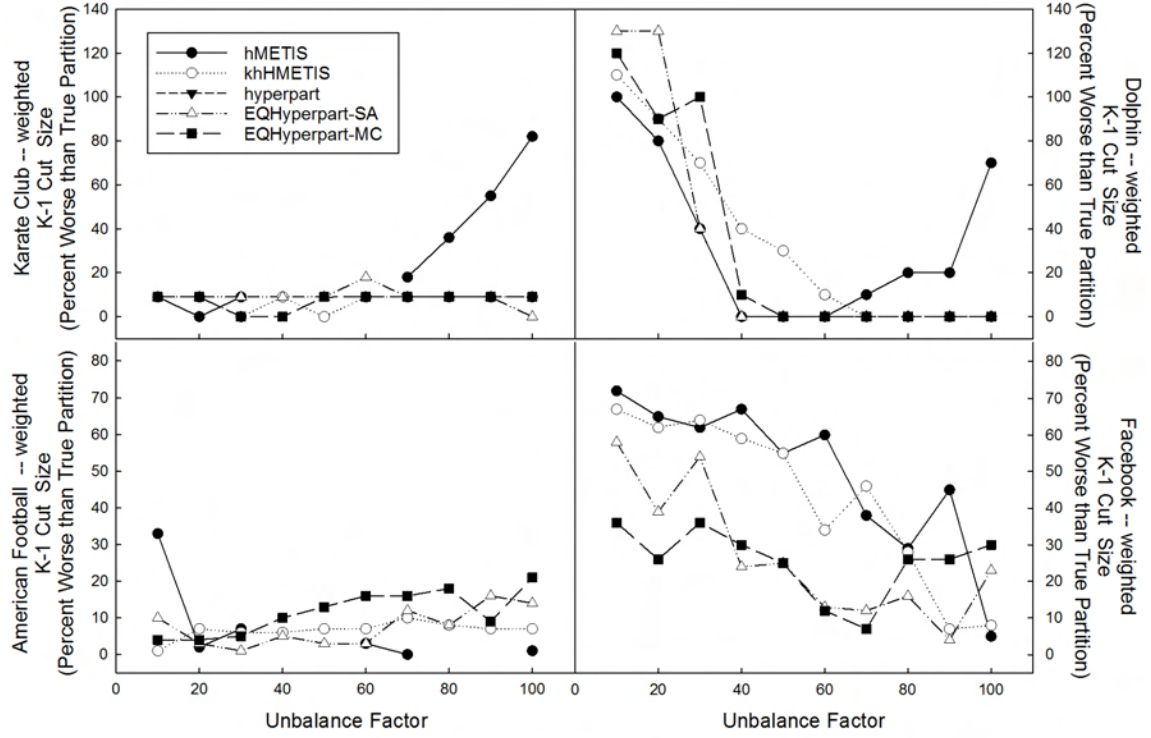


Fig. 8: $K - 1$ cut size deviation to true partition for weighted datasets.

partitioning, so we quantify the auto-tradeoff ability among these three elements using (17). Here, FS_{auto} , CS_{auto} and UB_{auto} indicate the F-Score value, $K - 1$ cut size and the unbalance factor value produced by auto-tradeoff partitioner, respectively. $|E|$ denotes the number of edges of the dataset, and UB_{max} is the max unbalance factor. Similarly, a , b , and c are weighting coefficients, subjected to $a + b + c = 1$.

$$TO_{auto} = a * FS_{auto} * 100 + b * (1 - \frac{CS_{auto}}{|E|}) * 100 + c * (1 - \frac{UB_{auto}}{UB_{max}}) * 100. \quad (17)$$

Table 4 illustrates the auto-tradeoff scores of the four unweighted and weighted datasets produced by EQHyperpart-SA. The unbalance factor adopts the metric of hMETIS; thus, $UB_{max} = 50$. Auto-tradeoff score is calculated by setting the coefficients $a = 0.2$, $b = 0.5$ and $c = 0.3$. EQHyperpart-SA was performed without any balance constraint for many times, and the most appropriate findings under the coefficient setting are listed in Table 4. In addition, the tradeoff score of real-world partition result for each dataset is included for comparison.

- In the datasets of Karate Club and Dolphin, the tradeoff score of auto-tuning EQHyperpart-SA are approximated to that of real-world placement. The performance on Facebook is modest, and the worst performance is observed on American Football Team.
- On the condition of this coefficient setting, the real-world placement is not necessarily the one with highest score. Some partitioning result produced by EQHyperpart-SA may have better overall performance. In fact, this is the universal phenomenon under other coefficient settings, but other auto-tuning partitioning findings could not be displayed due to space limitation. This is in accordance with that under balance constraints.
- The auto-tradeoff scores do not vary between unweighted and weighted datasets, with the exception of Facebook dataset. Specifically, EQHyperpart-SA achieves a better tradeoff on the weighted Facebook

Table 4: Auto-tradeoff findings. (Findings are rounded. (R) denotes the real-world partition result, and (w) denotes the weighted hypergraph modeled dataset, otherwise, unweighted.)

	$K - 1$ Cut Size	UBFactor	F-Score	Tradeoff Score
Karate Club(R)	11	6	100	89
Karate Club	9	9	81	85
Karate Club(w)	12	0	94	91
Dolphin(R)	10	17	100	87
Dolphin	12	15	98	87
Dolphin(w)	10	17	100	87
Football(R)	347	8	100	67
Football	181	33	31	52
Football(w)	181	33	31	52
Facebook(R)	223	12	100	92
Facebook	23	31	53	72
Facebook(w)	198	15	62	83

dataset than on the unweighted dataset, as the model of hypergraph affects the effect of auto-tuning partitioning results. We remark that Facebook is a classical scale-free network whose vertex degree obeys the power law, which makes it more suitable to be modeled as a weighted hypergraph than an unweighted hypergraph. The larger the social network dataset, the stronger this effect.

In summary, the EQHyperpart produces satisfactory results without any balance constraints. In other words, we have validated the effectiveness of the auto-tradeoff partitioning capability of EQHyperpart.

5.5. Performance experiments

To evaluate the performance, we employ the execution time of five partitioners over the four unweighted datasets for comparison. Because the partitioning times on different unbalance levels over the same dataset are similar, we calculated the mean value of them for each partitioner on each dataset, and normalized the runtime relative to that of khMETIS. As shown in Fig. 9, khMETIS and hMETIS are effective in reducing execution time. However, this comes at the cost of unsatisfactory partitioning qualities as mentioned above. We also noted from the findings that, EQHyperpart-SA runs faster than hyperpart, but EQHyperpart-MC runs slower than hyperpart. The reason for this is that EQHyperpart-MC utilizes a more elaborate cut size based on EQHyperpart-SA, which needs a little increase on calculation cost in exchange for the decrease of communication cost. Comprehensively, the running performance of EQHyperpart is modest and acceptable.

6. Conclusion

In this paper, we studied hypergraph partitioning method for social networks, and presented a hypergraph partitioner, EQHyperpart, which utilizes modularity Q based on the scale-free featured information Entropy (EQ) to guide the low cost partitioning process.

We then demonstrated that EQHyperpart achieves low cut size while retaining the modularity characteristics at certain balance levels, and has an effective auto-tradeoff partitioning capability. We also proposed two variant partitioners EQHyperpart-SA and EQHyperpart-MC that utilize simulated annealing and micro cut heuristic to optimize the partition quality. We compared our approach with one state-of-the-art and two popular hypergraph partitioners, and demonstrated that EQHyperpart is more scalable and suitable for partitioning social networks, especially on the weighted hypergraph modeled ones.

Future work includes extending this research by expanding our partitioning method and validating the approach using other widely used social network datasets.

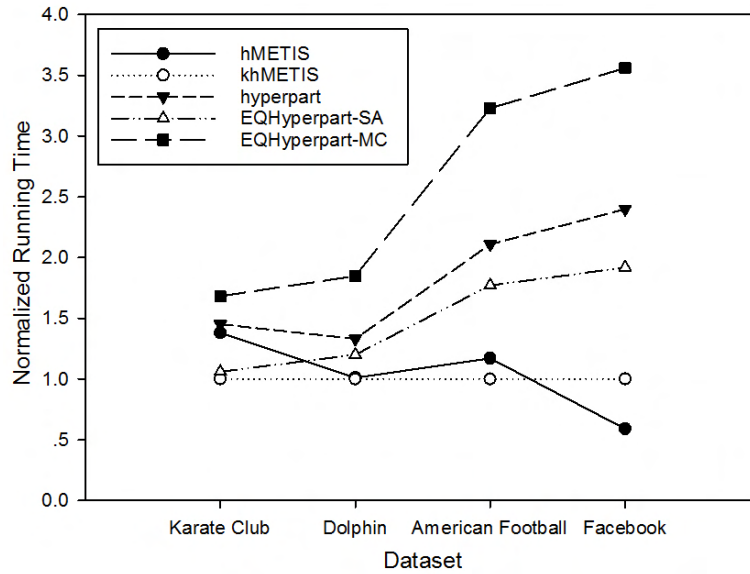


Fig. 9: Normalized running time comparison.

Acknowledgments

This work is supported in part by the National Natural Science Foundation of China under Grant Numbers 61632009, 61472451, 61272151, and 61402543, the High Level Talents Program of Higher Education in Guangdong Province under Funding Support Number 2016ZJ01, the China Guangdong Provincial Science & Technology Program under Grant Number 2015A030313638, and the Youth Creative Talents Program of Department of Education in Guangdong Province under Grant Number 2015KQNCX179.

References

- [1] W. Yang, G. Wang, M. Z. A. Bhuiyan, Partitioning of hypergraph modeled complex networks based on information entropy, in: Proceedings of the 15th International Conference on Algorithms and Architectures for Parallel Processing, 2015, pp. 678–690.
- [2] I. Saleh, T. Wei, M. B. Blake, Social-network-sourced big data analytics, *Internet Computing IEEE* 17 (5) (2013) 62–69.
- [3] J. M. Pujol, V. Erramilli, G. Siganos, X. Yang, N. Laoutaris, P. Chhabra, P. Rodriguez, The little engine(s) that could: scaling online social networks, *ACM Sigcomm Computer Communication Review* 40 (4) (2010) 1162–1175.
- [4] J. M. Pujol, G. Siganos, V. Erramilli, P. Rodriguez, J. M. Pujol, G. Siganos, Scaling online social networks without pains, *Proc of Netdb*.
- [5] M. Yuan, D. Stein, B. Carrasco, J. M. F. Trindade, Y. Lu, Partitioning social networks for fast retrieval of time-dependent queries, in: Data Engineering Workshops (ICDEW), 2012 IEEE 28th International Conference on, 2012, pp. 205–212.
- [6] C. Curino, E. Jones, Y. Zhang, S. Madden, Schism: a workload-driven approach to database replication and partitioning, *Proceedings of the Vldb Endowment* 3 (1-2) (2010) 48–57.
- [7] A. Turk, R. Oguz Selvitopi, H. Ferhatosmanoglu, C. Aykanat, Temporal workload-aware replicated partitioning for social networks, *IEEE Transactions on Knowledge & Data Engineering* 26 (11) (2014) 2832–2845.
- [8] S. Arora, S. Rao, U. Vazirani, Expander flows, geometric embeddings and graph partitioning, *Journal of the ACM* 56 (2) (2009) 1–37.
- [9] G. Karypis, V. Kumar, A fast and high quality multilevel scheme for partitioning irregular graphs, *Siam Journal on Scientific Computing* 20 (1) (1998) 359–392.
- [10] G. Karypis, R. Aggarwal, V. Kumar, S. Shekhar, Multilevel hypergraph partitioning: applications in VLSI domain, *Very Large Scale Integration Systems IEEE Transactions on* 7 (1) (1999) 69–79.
- [11] G. Karypis, V. Kumar, Multilevel k-way hypergraph partitioning, in: Proceedings of the 36th annual ACM/IEEE Design Automation Conference, 1999, pp. 343–348.
- [12] J. R. Yaros, T. Imielinski, Imbalanced hypergraph partitioning and improvements for consensus clustering, in: the 25th IEEE International Conference on Tools with Artificial Intelligence (ICTAI), IEEE Press, 2013, pp. 358–365.

- [13] B. Heintz, A. Chandra, Beyond graphs: toward scalable hypergraph analysis systems, *ACM SIGMETRICS Performance Evaluation Review* 41 (4) (2014) 94–97.
- [14] W. Yang, G. Wang, Directed social hypergraph data allocation strategy in online socail networks, *Journal of Chinese Computer Systems* 36 (7) (2015) 1559–1564.
- [15] A. Guzzo, A. Pugliese, A. Rullo, D. Saccà, Intrusion detection with hypergraph-based attack models, in: *Graph Structures for Knowledge Representation and Reasoning*, Springer, 2014, pp. 58–73.
- [16] C. Aaron, M. E. J. Newman, M. Cristopher, Finding community structure in very large networks, *Physical Review E Statistical Nonlinear & Soft Matter Physics* 70 (6) (2004) 264–277.
- [17] H. Kwak, Y. Choi, Y. H. Eom, H. Jeong, S. Moon, Mining communities in networks: a solution for consistency and its evaluation, in: *Proceedings of the 9th ACM SIGCOMM Conference on Internet Measurement Conference*, ACM, 2009, pp. 301–314.
- [18] M. E. J. Newman, Modularity and community structure in networks, *Proceedings of the National Academy of Sciences* 103 (23) (2006) 8577–8582.
- [19] V. D. Blondel, J. L. Guillaume, R. Lambiotte, E. Lefebvre, Fast unfolding of communities in large networks, *Journal of Statistical Mechanics Theory & Experiment* 30 (2) (2008) 155–168.
- [20] A. L. Barabási, E. Bonabeau, Scale-free networks, *Scientific American* 288 (3) (2003) 60–69.
- [21] X. Deng, B. Wang, B. Wu, S. Yang, Research and evaluation on modularity modeling in community detecting of complex network based on information entropy, in: *the 3rd IEEE International Conference on Secure Software Integration & Reliability Improvement*, 2009, pp. 297–302.
- [22] L. Jin, Y. Chen, T. Wang, P. Hui, A. V. Vasilakos, Understanding user behavior in online social networks: a survey, *Communications Magazine IEEE* 51 (9) (2013) 144–150.
- [23] C. Wilson, B. Boe, A. Sala, K. P. N. Puttaswamy, B. Y. Zhao, User interactions in social networks and their implications, *Natural Product Communications* 6 (1) (2011) 137–40.
- [24] J. Jiang, C. Wilson, X. Wang, P. Huang, W. Sha, Y. Dai, B. Y. Zhao, Understanding latent interactions in online social networks, *ACM Transactions on the Web* 7 (4) (2010) 369–382.
- [25] T. M. Wang, W. T. Lee, T. Y. Wu, H. W. Wei, Y. S. Lin, New P2P sharing incentive mechanism based on social network and game theory, *Journal of Network & Computer Applications* 41 (3) (2012) 47–55.
- [26] T. Ma, J. Zhou, M. Tang, Y. Tian, A. Al-Dhelaan, M. Al-Rodhaan, S. Lee, Social network and tag sources based augmenting collaborative recommender system, *Ieice Transactions on Information & Systems* 98 (4) (2015) 902–910.
- [27] G. Carullo, A. Castiglione, A. De Santis, F. Palmieri, A triadic closure and homophily-based recommendation system for online social networks, *World Wide Web* 18 (6) (2015) 1579–1601.
- [28] N. Zaerpour, M. Rabbani, A. H. Gharehgozli, R. Tavakkoli-Moghaddam, A comprehensive decision making structure for partitioning of make-to-order, make-to-stock and hybrid products, *Soft Computing* 13 (2009) 1035–1054.
- [29] G. Carullo, A. Castiglione, G. Cattaneo, A. De Santis, U. Fiore, F. Palmieri, Feeltrust: Providing trustworthy communications in ubiquitous mobile environment, in: *IEEE 27th International Conference on Advanced Information Networking and Applications*, 2013, pp. 1113–1120.
- [30] G. Carullo, A. Castiglione, A. De Santis, F. Palmieri, A triadic closure and homophily-based recommendation system for online social networks, *World Wide Web* 18 (2015) 15791601.
- [31] T. Ma, J. Zhou, M. Tang, Y. Tian, A. Al-Dhelaan, M. Al-Rodhaan, S. Lee, Social network and tag sources based augmenting collaborative recommender system, *IEICE transactions on Information and Systems* E98-D (4) (2015) 902–910.
- [32] L. Boltzmann, Further studies on the thermal equilibrium of gas molecules, *The Kinetic Theory of Gases. Series: History of Modern Physical Sciences* 1 (2003) 262–349.
- [33] C. E. Shannon, A mathematical theory of communication, *Bell System Technical Journal* 27 (3) (1948) 379–423.
- [34] P. G. Popescu, F. Pop, A. Herişanu, N. Țăpuș, New inequalities between information measures of network information content, *Mathematical Problems in Engineering* 2013 (1) (2013) 151–164.
- [35] P. G. Popescu, E.-I. Slușanschi, V. Iancu, F. Pop, A new upper bound for shannon entropy. A novel approach in modeling of big data applications, *Concurrency and Computation: Practice and Experience* 28 (2) (2016) 351–359.
- [36] C. Berge, *Graphs and Hypergraphs*, Elsevier Science Ltd., 1985.
- [37] R. Borndörfer, O. Heismann, The hypergraph assignment problem, *Discrete Optimization* 15 (4) (2015) 15–25.
- [38] B. W. Kernighan, S. Lin, An efficient heuristic procedure for partitioning graphs, *Bell System Technical Journal* 49 (2) (1970) 291–307.
- [39] C. M. Fiduccia, R. M. Mattheyses, A linear-time heuristic for improving network partitions, in: *Proceedings of the 19th Design Automation Conference*, 1982, pp. 175–181.
- [40] M. Deveci, K. Kaya, B. Uçar, Ümit V. Çatalyürek, Hypergraph partitioning for multiple communication cost metrics: Model and methods, *Journal of Parallel & Distributed Computing* 77 (2015) 69–83.
- [41] M. E. J. Newman, M. Girvan, Finding and evaluating community structure in networks, *Physical Review E Statistical Nonlinear & Soft Matter Physics* 69 (2) (2004) 026113.
- [42] F. Santo, L. Vito, M. Massimo, Method to find community structures based on information centrality, *Physical Review E Statistical Nonlinear & Soft Matter Physics* 70 (5) (2004) 148–168.
- [43] P. Pons, M. Latapy, Computing communities in large networks using random walks, in: *Proceedings of the 20th International Conference on Computer and Information Sciences*, Springer Berlin Heidelberg, 2005, pp. 284–293.
- [44] M. E. Newman, Finding community structure in networks using the eigenvectors of matrices, *Physical Review E Statistical Nonlinear & Soft Matter Physics* 74 (3) (2006) 92–100.
- [45] M. E. Newman, Fast algorithm for detecting community structure in networks, *Physical Review E Statistical Nonlinear & Soft Matter Physics* 69 (6) (2004) 066133–066133.
- [46] G. Palla, I. Derényi, I. Farkas, T. Vicsek, Uncovering the overlapping community structure of complex networks in nature and

- society, *Nature* 435 (7043) (2005) 814–818.
- [47] J. He, D. Chen, A fast algorithm for community detection in temporal network, *Physica A: Statistical Mechanics & its Applications* 429 (2015) 87–94.
- [48] J. Eustace, X. Wang, Y. Cui, Overlapping community detection using neighborhood ratio matrix, *Physica A: Statistical Mechanics & its Applications* 421 (2015) 510–521.
- [49] R. Martin, C. T. Bergstrom, An information-theoretic framework for resolving community structure in complex networks, *Proceedings of the National Academy of Sciences* 104 (18) (2007) 7327–7331.
- [50] H. Shen, X. Cheng, H. Chen, Y. Liu, Information bottleneck based community detection in network, *Chinese Journal of Computers* 31 (4) (2008) 677–686.
- [51] X. Zheng, D. Y. Yi, Z. Z. Ouyang, L. Dong, Hyperedge communities and modularity reveal structure for documents, *Chinese Physics Letters* 29 (3) (2012) 038902.
- [52] A.-L. Barabási, R. Albert, H. Jeong, Scale-free characteristics of random networks: the topology of the world-wide web, *Physica A: Statistical Mechanics & its Applications* 281 (1-4) (2000) 69–77.
- [53] S. Z. Selim, K. Alsultan, A simulated annealing algorithm for the clustering problem, *Pattern Recognition* 24 (91) (1991) 1003–1008.
- [54] M. E. J. Newman, Website of Mark Newman, <http://www-personal.umich.edu/~mejn/netdata/> (2015).
- [55] UCI, UCI network data repository, <http://networkdata.ics.uci.edu/> (2015).
- [56] SNAP, Stanford large network dataset collection, <http://snap.stanford.edu/data/index.html> (2015).
- [57] J. J. McAuley, J. Leskovec, Learning to discover social circles in ego networks, *Advances in Neural Information Processing Systems* (2012) 539–547.