



A new clustering method and its application in social networks

Peixin Zhao^{a,*}, Cun-Quan Zhang^{b,1}

^aSchool of Management, Shandong University, #27 Shanda Nan Rd., Jinan 250100, PR China

^bDepartment of Mathematics, West Virginia University, Morgantown, WV 26506-6310, USA

ARTICLE INFO

Article history:

Received 27 April 2010

Available online 30 June 2011

Communicated by L. Heutte

Keywords:

Clustering

Graph theory

Hierarchical tree

Social network

ABSTRACT

In a graph theory model, clustering is the process of division of vertices into groups, with a higher density of edges within groups than between them. In this paper, we introduce a new clustering method for detecting such groups and use it to analyse some classic social networks. The new method has two distinguished features: non-binary hierarchical tree and the feature of overlapping clustering. A non-binary hierarchical tree is much smaller than the binary-trees constructed by most traditional methods and, therefore, it clearly highlights meaningful clusters which significantly reduces further manual efforts for cluster selections. The present method is tested by several bench mark data sets for which the community structure was known beforehand and the results indicate that it is a sensitive and accurate method for extracting community structure from social networks.

© 2011 Elsevier B.V. All rights reserved.

1. Introduction

Clustering is an important task for the discovery of community structures in networks. Its goal is to sort cases (people, things, events, etc.) into clusters so that the degree of association is relatively strong between members of the same cluster and relatively weak between members of different clusters. Merriam-Webster (2008) defines cluster analysis as “a statistical classification technique for discovering whether the individuals of a population fall into different groups by making quantitative comparisons of multiple characteristics.” Various clustering algorithms have been proposed in the literature in many different scientific disciplines. Jain (2009) broadly divided these algorithms into two groups: (i) hierarchical method and (ii) partitional method. Hierarchical clustering algorithms recursively find nested clusters either in agglomerative mode or in divisive mode. The most well-known hierarchical algorithms are single-link and complete-link; in single-link hierarchical clustering, the two clusters whose two closest members have the smallest distance are merged in each step; in complete-link case, the two clusters whose merger has the smallest diameter are merged in each step. Compared to hierarchical clustering algorithms, partitional clustering algorithms find all the clusters simultaneously as a partition of the data and do not impose a hierarchical structure. The most popular and the simplest partitional algorithm is *K*-means (Steinhaus, 1956). Berkhin (2009) listed another six classifications besides the above two main groups: (iii) grid-based methods, (iv) methods based on co-occur-

rence of categorical data, (v) constraint-based clustering, (vi) clustering algorithms used in machine learning, (vii) scalable clustering algorithms, (viii) algorithms for high dimensional data.

In recent years, a growing number of clustering algorithms for categorical data have been proposed based on various centrality measures. For instance, vertex betweenness has been studied by Freeman (1977) as a measure of the centrality to detect communities in a network; Girvan and Newman (2002) and Newman and Girvan (2004) generalize vertex betweenness centrality to edge in order to discover community structures; Frey and Dueck (2007) devised a method called affinity propagation, which takes as input measures of similarity between pairs of data points. Newman (2006) utilizes the eigenvectors of matrices to find community structure in networks; Rosvall and Bergstrom (2008) use the probability flow of random walks on a network as a proxy for information flows in the real system and decompose the network into modules. Wu and Huberman (2004) propose an approach for discovering the communities based on the property of resistor networks. For reviews see Refs. Newman (2004a) and Danon et al. (2005).

A social network (Wasserman and Faust, 1994; Scott, 2000) is a set of people or groups each of which has connections of some kind to some or all of the others. A cluster is a collection of individuals with dense relationship internally and sparse relationships externally. Based on this criterion, we introduce a new clustering method for detecting community structures. In this method, individuals and their relationships are denoted by weighted graphs, then the graph density we defined gives a better quantity depict of whole correlation among individuals in a community, so that a reasonable clustering output can be presented. Compared with other methods, this method has two important features:

* Corresponding author. Tel.: +86 531 88363169.

E-mail addresses: pxzhao@sdu.edu.cn (P. Zhao), cqzhang@math.wvu.edu (C.-Q. Zhang).

¹ Tel.: +1 304 293 2011x2332.

- (1) A much smaller hierarchical trees that clearly highlight meaningful clusters.
- (2) Overlapping clusters.

To evaluate the effectiveness of our method, we applied it to analyse some classic bench mark data sets whose clusters are already known. These data sets include Karate Club, Davis southern club women, Dolphin, Books about US politics, American College Football. The accuracy of the outputs in those classical benchmark data sets is a supporting evidence of futhre applicability of the new method.

The rest of the paper is organized as follows. In Section 2 we introduce the details of the new dense subgraph clustering method. In Section 3 we apply it to some classic social networks and compare its results with that of known clusters. Finally, a summary and conclusions are given in Section 4.

2. A new clustering method

A graph or network is one of the most commonly used models to present real-valued relationships of a set of input items. Let $G = (V, E)$ be a graph with the vertex set V and the edge set E with weight $w(e)$ on every edge e . Models with un-weighted graphs (the weight of every edge is set to 1) have been extensively studied in graph theory. In an un-weighted graph G , a subgraph H of G is defined as a *clique* if every pair of vertices of H is joined by one edge (Bondy and Murty, 1976; Diestel, 2005; West, 1996). It is well-known that the search of maximum cliques in graphs is an NP-complete problem (Gary and Johnson, 1979). Therefore, it is not practical to define cliques as clusters. Furthermore, there is no appropriate definition for a clique in a weighted graph. However, in order to closely represent the nature and the real situation of the inputs in most applications (different degrees of similarity for clustering problems), we should use weighted graph models which are much more appropriate than un-weighted models. For simplification or other practical reasons, many designers of clustering methods may set a specific threshold, such as that any edge with weight below the threshold is deleted and the remaining ones have no associated weight. However, one may not be able to expect an accurate output since the cut-off (by threshold) may cause a loss of important information.

For a subgraph $C(|V(C)| > 1)$, we define the density of C by

$$d(C) = \frac{2 \sum_{e \in E(C)} w(e)}{|V(C)|(|V(C)| - 1)}. \quad (1)$$

As seen above, if $w(e) = 1$ for every edge e in C and $d(C) = 1$, then the subgraph C induces a clique. For a weighted graph, a subgraph C is called a Δ -quasi-clique if $d(C) \geq \Delta$ for some positive real number Δ .

Since clustering is a process that detects all dense subgraphs in G and construct a hierarchically nested system to illustrate their inclusion relation, a heuristic process is applied here for finding all quasi-cliques with density in various levels. The core of the algorithm is deciding whether or not to add a vertex to an already selected dense subgraph C . For a vertex $v \notin V(C)$, we define the contribution of v to C by

$$c(v, C) = \frac{\sum_{u \in V(C)} w(uv)}{|V(C)|}. \quad (2)$$

A vertex v is added into C if $c(v, C) > \alpha d(C)$ where α is a function of some user specified parameters.

Instance: $G = (V, E)$ is a graph with edge weights $w: E(G) \mapsto \mathbb{R}^+$.

Question: Detects Δ -quasi-cliques in G with various levels of Δ , and construct a hierarchically nested system to illustrate their inclusion relation.

Sub-Algorithm Growing(C, G):

(Grow a Community C in G)
 while $V(G) - V(C) \neq \emptyset$
 begin
 pick $v \in V(G) - V(C)$ such that $c(v, C)$ is a maximum
 if $c(v, C) > \alpha_n d(C)$ then add v to C (where $n = |V(C)|$,
 $\alpha_n = 1 - \frac{1}{2\lambda(n+t)}$
 with $\lambda \geq 1$ and $t \geq 1$ as user specified parameters)
 else return
 end

Sub-Algorithm Decompose(G, w_0):

(decompose a graph G into communities using edges with weights at least w_0).
 Let $E_0 = \{e \in E(G) : w(e) \geq w_0\}$
 For each $e = uv \in E_0$ in decreasing order of $w(e)$
 begin
 if either u or v is not in any community
 then
 begin
 create a new empty community C and add u, v in it
 Growing(C, G)
 end
 end
 end

Sub-Algorithm Merging(G):

For any two communities C_i and C_j in G , if $|C_i \cap C_j| > \beta \min(|C_i|, |C_j|)$ then merge C_i and C_j into a new community $C = C_i \cup C_j$ (where β is a user specified parameter).
 Contract each community to a vertex. The weight of an edge is defined by

$$w(C_i, C_j) = \frac{\sum_{e \in E_{ij}} w(e)}{|E_{ij}|}$$

where the set of crossing edges $E_{ij} = \{v_i v_j : v_i \in C_i, v_j \in C_j, v_i \neq v_j\}$

Main-Algorithm

(generate hierarchic clustering tree or a graph G)
 while $E(G) \neq \emptyset$
 begin
 Choose w_0 according to some criterion
 Decompose(G, w_0)
 Merging(G)
 Store the resulted graph to G
 end
 Trace the movement of each vertex and generate the hierarchic tree.

3. Applications of the method in social networks

In this section we present a number of applications of our method to some classic social networks for which the community structure is already known and compare its results with that of preceded methods.

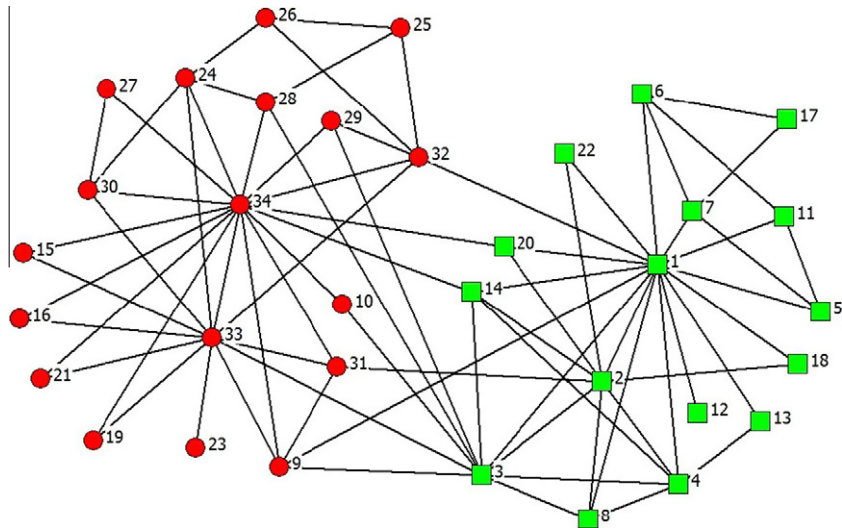


Fig. 1. The network of friendships in the karate club study in Zachary (1977).

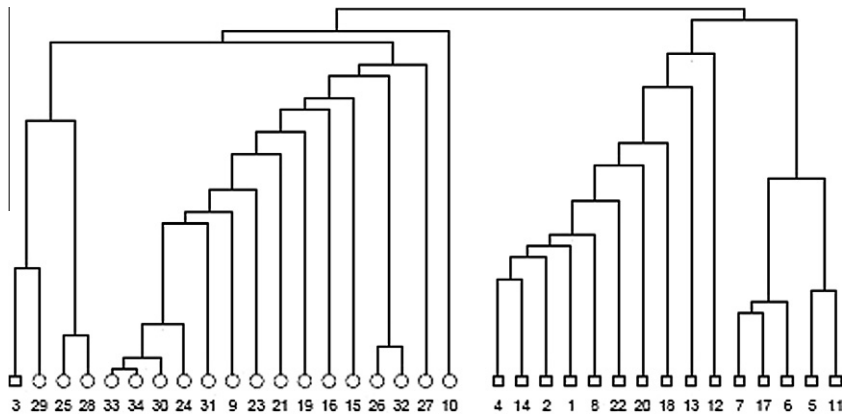


Fig. 2. Result of Girvan and Newman algorithm (Girvan and Newman, 2002).

3.1. Zachary's karate club study

The first social network is the well known “karate club” of Zachary (1977). He observed 34 members of a karate club over two years. During the course of observation, the club members split into two groups because of the disagreement between the administrator of the club and the club's instructor, the members of one group left to start their own club. Zachary constructed a simple unweighted graph to show the friendships between two

members of the club, each member in the club is represented by a node, and edge is drawn if the two members are friends outside the club activities. Fig. 1 shows the network, with the administrator and instructor were respected by node 1 and 34 respectively.

Table 1
Matrix A: Attendance records of eighteen women.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14
Eleanor	0	1	0	1	0	0	0	1	0	0	0	1	0	0
Brenda	0	1	0	1	0	0	1	1	0	1	0	1	1	0
Dorothy	0	0	0	0	0	1	0	0	0	0	0	1	0	0
Verne	0	0	0	1	1	1	0	0	0	0	0	1	0	0
Flora	1	0	0	0	0	1	0	0	0	0	0	0	0	0
Olivia	1	0	0	0	0	1	0	0	0	0	0	0	0	0
Laura	0	1	1	1	0	0	1	1	0	1	0	1	0	0
Evelyn	0	1	1	0	0	1	1	1	0	1	0	1	1	0
Pearl	0	0	0	0	0	1	0	1	0	0	0	1	0	0
Ruth	0	1	0	1	0	1	0	0	0	0	0	1	0	0
Sylvia	0	0	0	1	1	1	0	0	1	0	1	0	0	1
Katherine	0	0	0	1	1	1	0	0	1	0	1	1	0	1
Myrna	0	0	0	0	1	1	0	0	1	0	0	1	0	0
Theresa	0	1	1	1	0	1	1	1	0	0	0	1	1	0
Charlotte	0	1	0	1	0	0	1	0	0	0	0	0	1	0
Frances	0	1	0	0	0	0	1	1	0	0	0	1	0	0
Helen	1	0	0	1	1	0	0	0	1	0	0	1	0	0
Nora	1	0	0	1	1	1	0	1	1	0	1	0	0	1

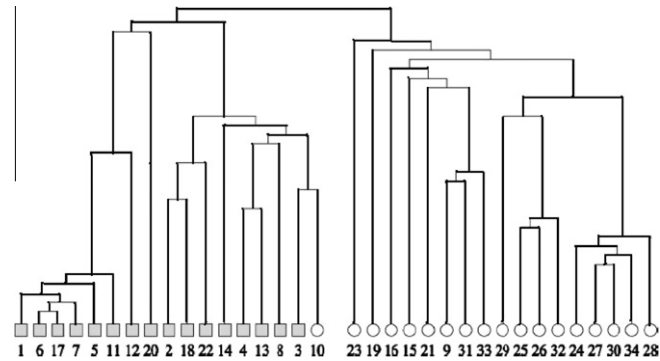


Fig. 3. Result of modularity method (Newman, 2004b).

Green square represent individuals associated with the administrator and red circle represent those associated with the instructor.

The algorithm of Girvan and Newman (2002) and modularity method (Newman, 2004b) have been respectively used for detecting the communities in this network. Figs. 2 and 3 show the hierarchical trees of their results. Both of their divisions are almost perfect except one node is classified incorrectly (see Node 3 in

Fig. 2 and Node 10 in Fig. 3). As a completely different approach, the Kernighan–Lin algorithm (Kernighan and Lin, 1970) was also applied to this network, it detects the two factions perfectly—every vertex is correctly classified. Pointed out by Newman (2004b) that an algorithm should look for groups of size 16 and 18, which are the known size of the groups into which the network split. Giving any other sizes will lead to “wrong” results.

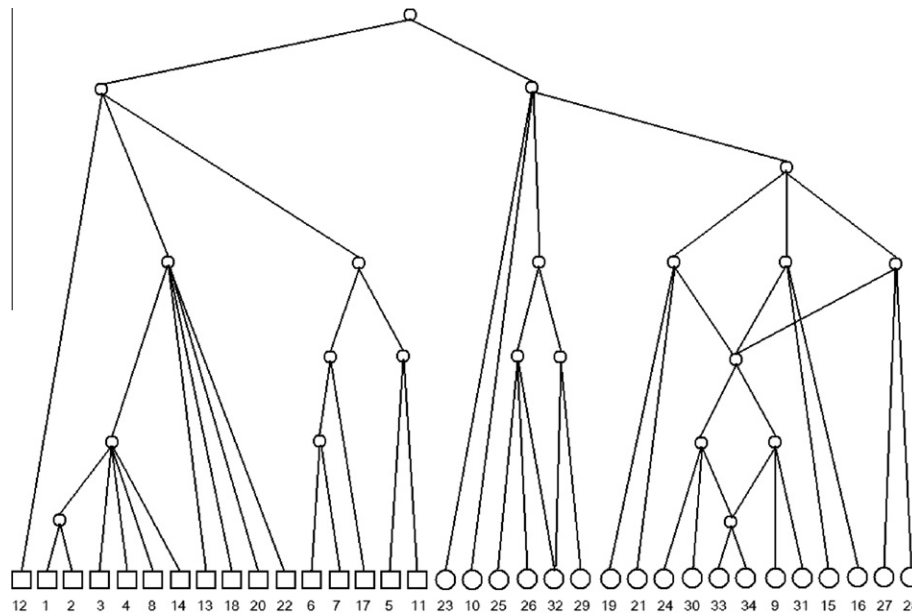


Fig. 4. Result of our method.

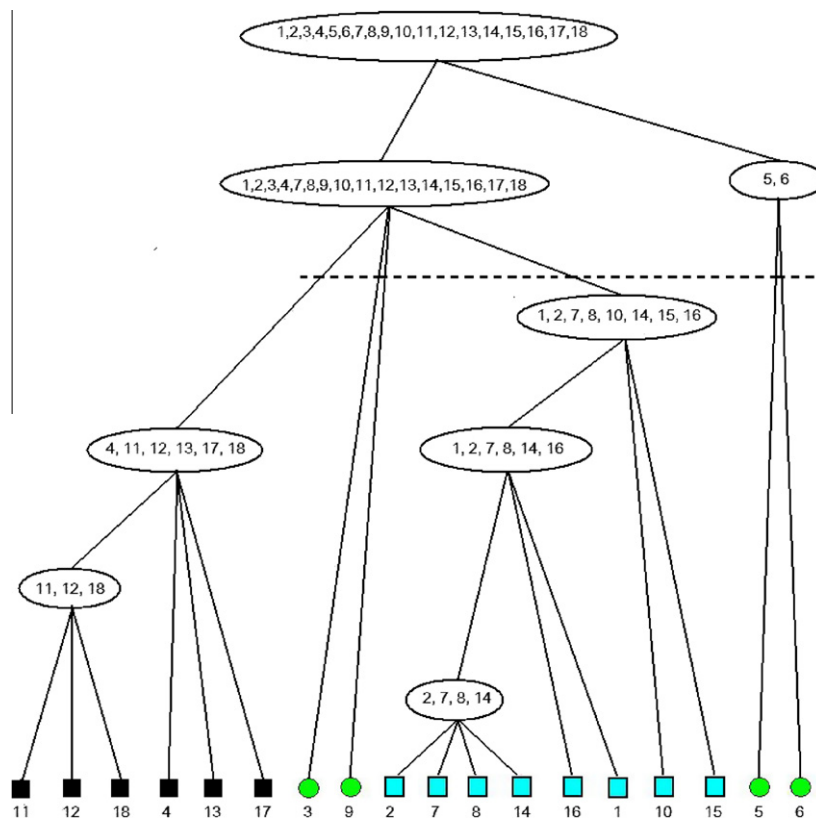


Fig. 5. Dendrogram of the groups found by our method in Southern Club Women Network.

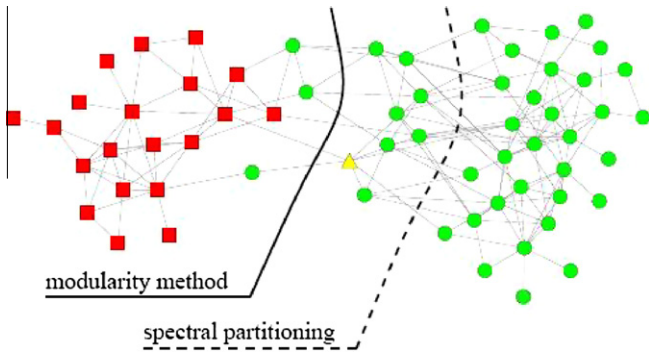


Fig. 6. Division of dolphins network in (Newman, 2006).

Fig. 4 illustrates the output dendrogram derived by applying our method to this network. One may notice immediately that the biggest two groups corresponds perfectly with the actual factions observed by Zachary.

3.2. Davis southern club women

The data of the social participation of eighteen women in “Old City” was collected by Davis et al. (1941). The data (see Table 1) is a table with 18 rows—one for each woman—and 14 columns, one for each “event” (such as a club meeting, a church supper, a card party, etc.), held during the course of a year. For the simplicity, we use number 1–18 denote the 18 women, then a matrix A is generated to record their attendances of events: $A(i, j)$ is 1 if woman i

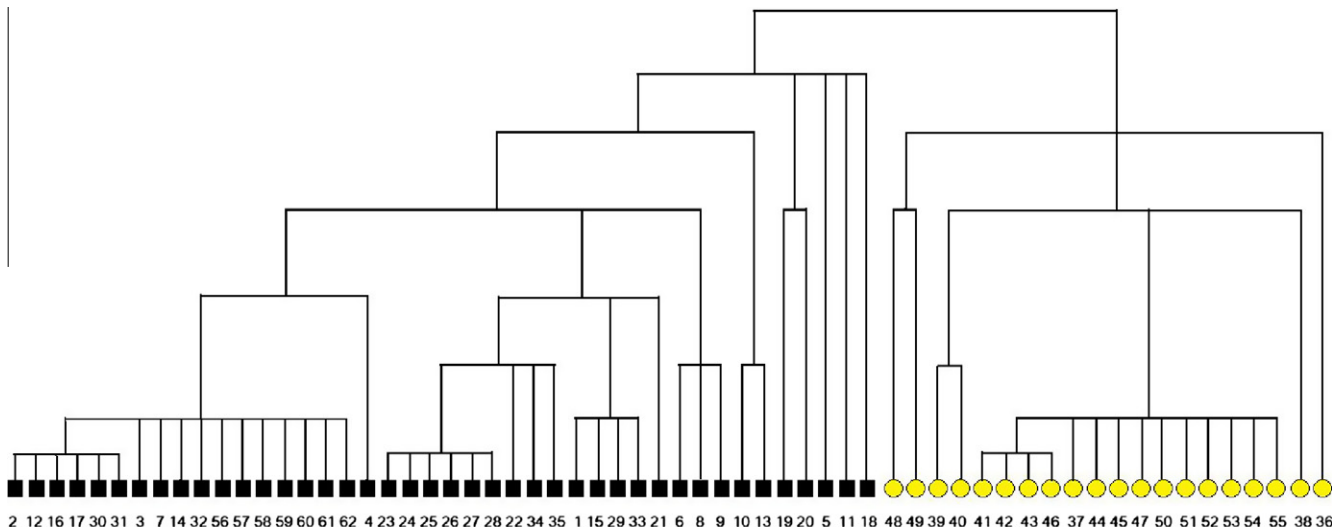


Fig. 7. The dendrogram of dolphins network by our method.

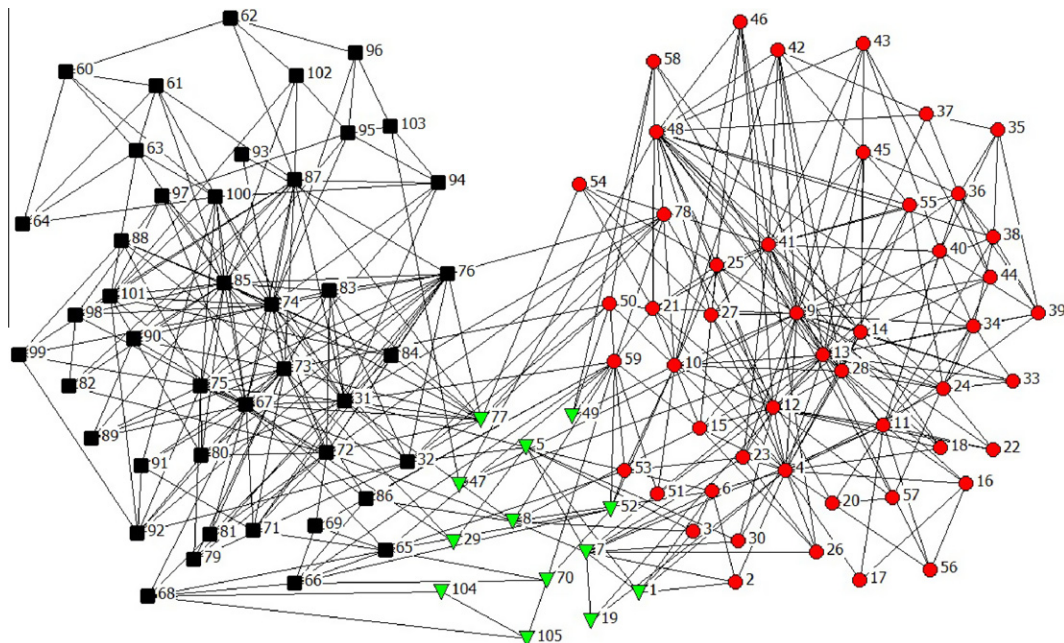


Fig. 8. Correct division of political books.

attended social event j , and 0 otherwise. The goal of the study was to determine the clique structure according to their records of attendances. The clique membership reported by Homans (1950) is as follows. Group 1: 1, 2, 7, 8, 14, 15, 16; Group 2: 11, 12, 13, 17, 18; women not clearly belonging to either groups: 3, 4, 5, 6, 9, 10.

For the using of our method, we define cell (i,j) of the symmetric matrix is the inner product of the i th and the j th row in matrix A . Then the dendrogram of the output is illustrated in Fig. 5. It is easy to know from the figure that the 18 women are mainly divided into 2 groups, one includes 4, 11, 12, 13, 17, 18, the other includes 1, 2, 7, 8, 10, 14, 15, 16, the rest 4 women (3,5,6,9) not clearly belong to either groups. This result is identical with the result of duality method in (Breiger, 1974). Both of the two divisions correspond almost perfectly to the alignments of Homans except tow women (4 and 10) are misclassified.

3.3. Dolphin's network

The next social network was constructed from observations of a bottlenose dolphin community (Lusseau et al., 2003; Lusseau, 2003; Lusseau and Newman, 2004). There are 62 nodes and 159 edges in this network: nodes represent the dolphins, edges between nodes represent associations between dolphin pairs occurring more often than expected by chance. This network is interest because, during the course of the study, the dolphin group split into two smaller subgroups following the departure of a key member of the population. The subgroups of the actual division in (Newman, 2006) are represented by the shapes of the vertices (see Fig. 6), the squares and circles represent the actual division of the network observed when the dolphin community split into two as a result of the departure of a keystone individual. The individual who departed is represented by the yellow triangle.

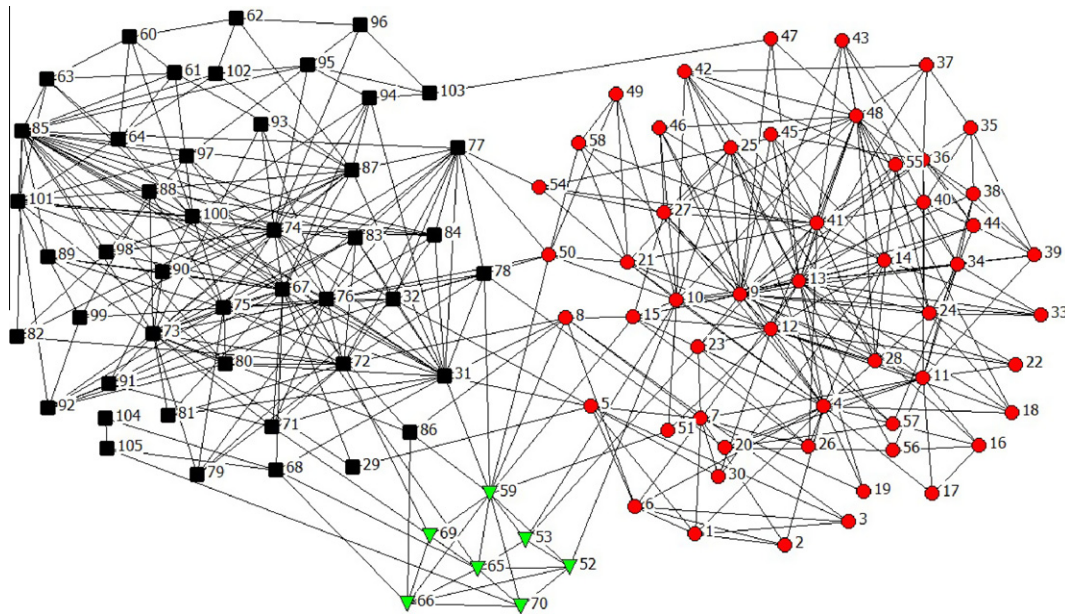


Fig. 9. Result of algorithm of Girvan and Newman (2002) in political books.

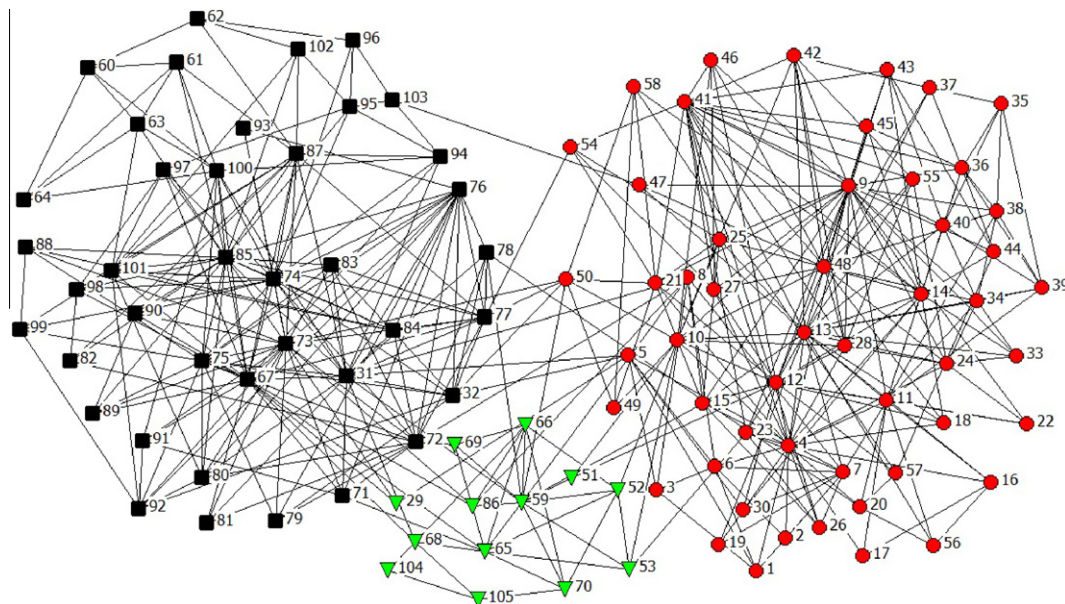


Fig. 10. Result of our method in political books.

The dotted line denotes the division of the network into two equal-sized groups found by the standard spectral partitioning method. The solid curve represents the division found by the method based on the leading eigenvector of the modularity matrix in (Newman, 2006). Its result corresponds quite closely to the actual split—all but 3 of the 62 dolphins are misclassified.

We have also applied our method to this network, the result is shown in Fig. 7 which corresponds perfectly with the actual division.

The algorithms of (Newman and Girvan, 2004) and Newman (2004a) also give precisely the same result.

3.4. Political books network

The dataset of books about US politics compiled by Krebs (2009). The 105 nodes represent 105 books about US politics sold by the online bookseller Amazon.com. The 441 edges represent

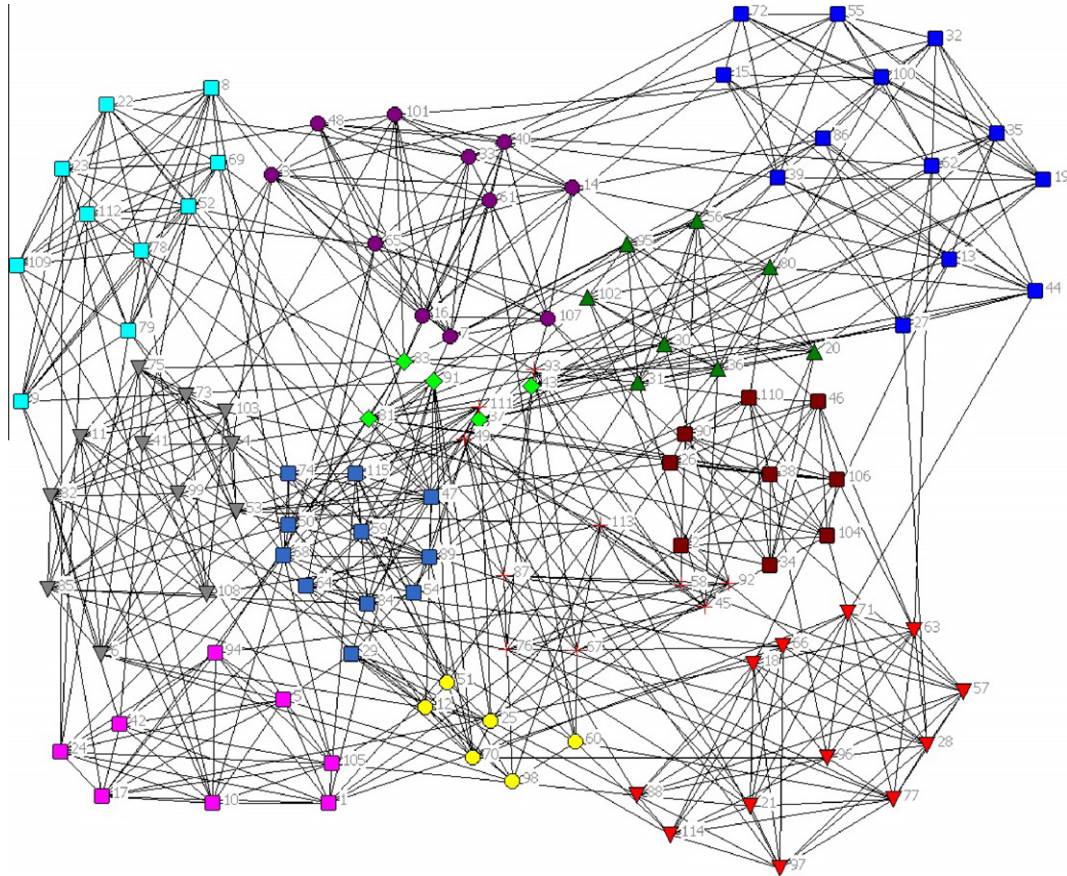


Fig. 11. Actual communities of US college football teams.

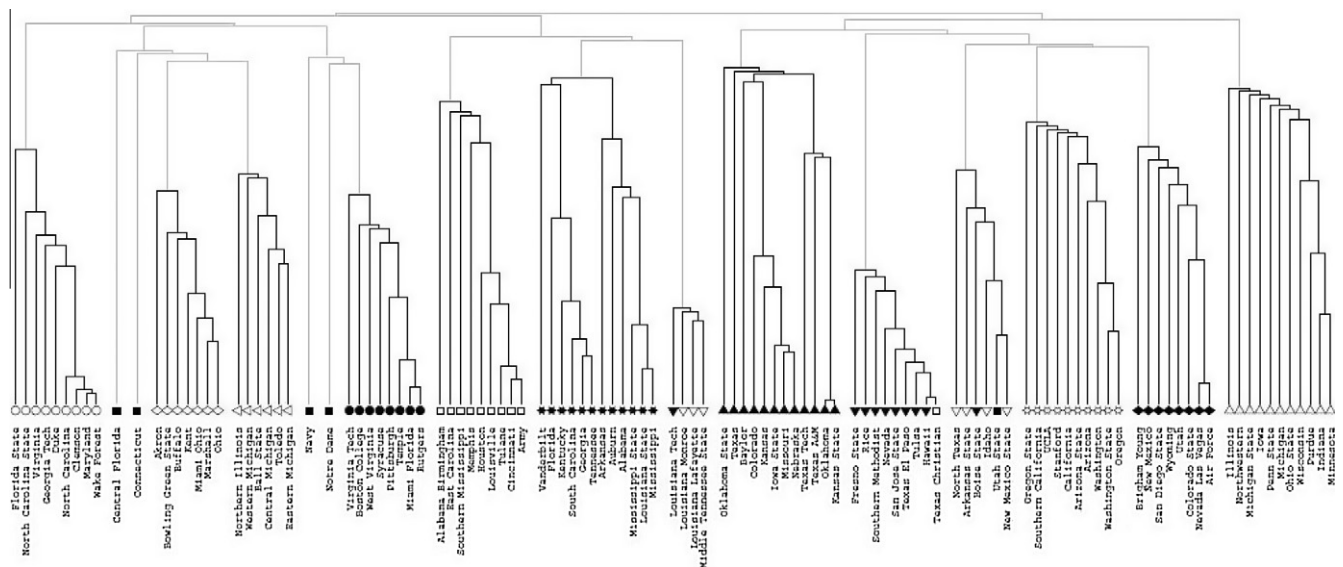


Fig. 12. The dendrogram generated by the algorithm of Girvan and Newman (2002).

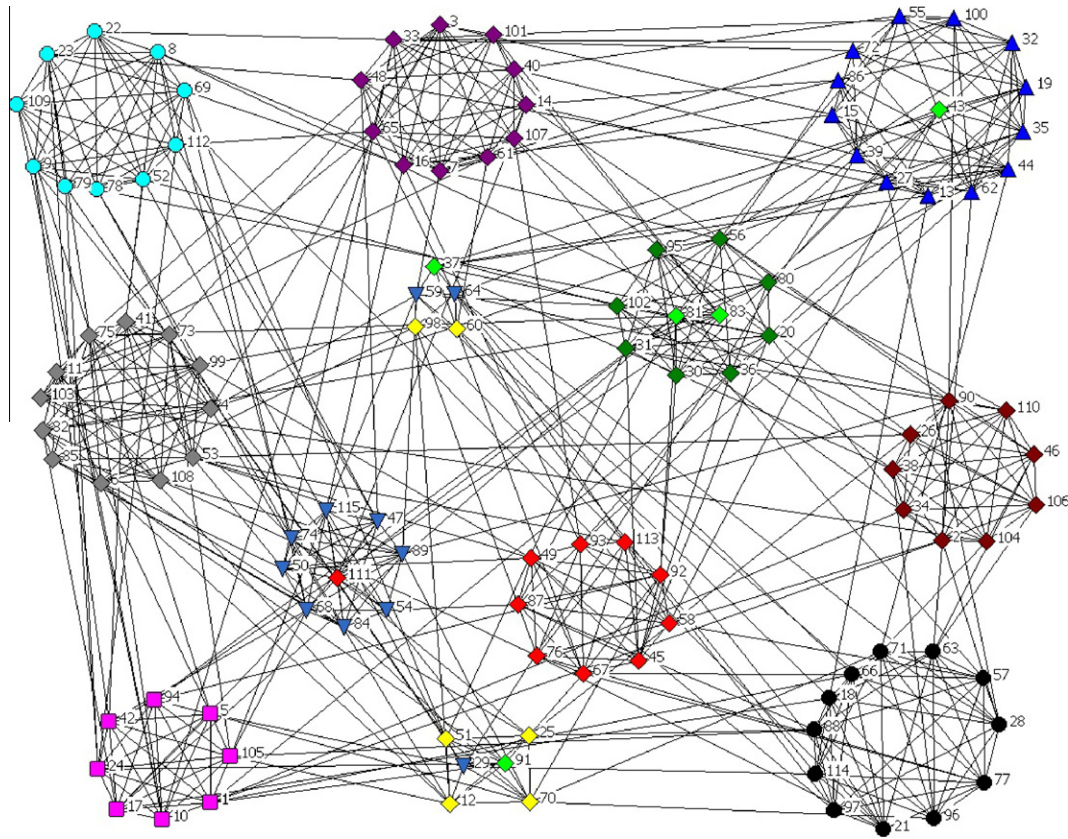


Fig. 13. Division of US college football team network by our method.

frequent co-purchasing of books by the same buyers, as indicated by the “customers who bought this book also bought these other books” feature on Amazon. The nodes have been given values “l”, “n”, or “c” to indicate whether they are “liberal”, “neutral”, or “conservative”. These alignments were assigned separately by Newman (2009) based on a reading of the descriptions and reviews of the books posted on Amazon (see Fig. 8, the three clusters are illustrated using three different shapes and colors: black square for “liberal” books, red circle for “conservative” books and green triangle for “neutral” books). The goal is to detect these clusters that represent the different political orientations of the books.

If the algorithm of Girvan–Newman is applied to this dataset, the books will be classified into 3 clusters (see Fig. 9). Compared with the correct classification in Fig. 8, it has 17 books are classified incorrectly: (1, 5, 7, 8, 19, 29, 47, 49, 53, 59, 65, 66, 69, 77, 78, 104, 105).

Feeding this network into our method, the results are shown in Fig. 10. The books are classified into 3 clusters; the number of incorrectly classified books is also 17: (1, 5, 7, 8, 19, 47, 49, 51, 53, 59, 65, 66, 68, 69, 77, 78, 86).

3.5. US college football

The next dataset is taken from another classic studies in social networks: US college football (Girvan and Newman, 2002). This network representing the schedule of games between American college football teams in 2000 season. 115 teams are divided into “conferences” containing about 8 to 12 teams each. They play an average of about 7 intra-conference games and 4 inter-conference games in a season. Fig. 11 shows the actual conferences. The dendrogram generated by the edge betweenness algorithm (Girvan and Newman, 2002) is shown in Fig. 12, the number of misclassified teams is 11.

The result of our method is illustrated in Fig. 13, which indicates some improvement of accuracy: the number of misclassified teams is 10.

3.6. An example of overlapping

Most of clustering methods generate non-overlapping groups, which are useful for graph drawing but are not as good for group analysis, since real social groups are usually more complex, involving different degrees of overlap among groups. In this section we use an example in (Santamaria and Roberto, 2008) to illustrate the overlapping feature of our method.

Fig. 14(a) shows three groups are represented as complete subgraphs. Fig. 14(b) shows edges are hidden and replaced by transparent hulls wrapping the elements in each group. The elements

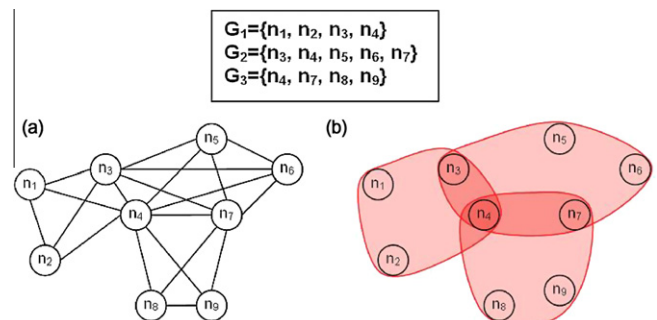


Fig. 14. The relationships between nodes and their clusters in Santamaria and Roberto (2008).

like n_4 , present in the three groups are highlighted by hull overlapping. Our method also give the same result precisely.

4. Conclusion

In this paper, we introduce a new clustering method for detecting community structure in network and use it to analyse some classic social networks. Compared with the existing methods, the new method has two distinguished features:

- (1) A much smaller hierarchical trees that clearly highlight meaningful clusters. It was pointed out in SAS/STAT User's Guide (SAS Institute Inc, 2003), "there are no completely satisfactory methods for determining the number of population clusters for any type of cluster analysis". Hence, a relatively small hierarchical tree in an output will significantly reduce the human involvement in the final selection of clusters. In Figs. 2 and 3, one may notice that each hierarchical tree (for karate club) with 34 leafs (inputs) has 33 internal nodes. With the new method, the hierarchical tree (see Fig. 4) contains only 22 internal vertices. Similarly, a hierarchical tree for women club with 18 leafs usually have 17 internal nodes, while the tree in Fig. 5 has only 8 internal vertices.
- (2) The feature of overlapping clustering does reflect the complexity of our real world. One may notice the overlapping clustering feature in Fig. 4: At the right end of the graph, a small group of club members {24,30,33,34,9,31} hold multi-memberships in three internal clusters (internal nodes on the third level, right end). The overlapping clustering is a concept that has recently received increased attention in (Palla et al., 2005; Pereira-Leal et al., 2004; Futschik and Carlisle, 2005), etc. One may also notice this feature in Fig. 14: The element n_3 and n_7 present in two groups and the node n_4 presents in three groups. According to mathematical definition in graph theory, the "hierarchical trees" in Figs. 4 and 14 are not really trees – they are *hierarchical networks* in which the relations of clusters are hierarchically nested.

For further research, we will consider to develop an automated value selection method for each parameter: Determine a function for each parameter in terms of some structural information of the input graph, so that graphs with different structures (density, connectivity, locally or globally) will be automatically assigned proper values.

Acknowledgements

The research of the first author is partially supported by Independent Innovation Foundation of Shandong University, IIFSDU(2009TS014). The research of the second author is partially supported by WV EPSCoR grant.

Appendix A. The main algorithm in detail

- Step 0.** $w_0 \leftarrow \gamma \max\{w(e) : \forall e \in E(G)\}$ where γ ($0 < \gamma < 1$) is a user specified parameter.
- Step 1.** (The initial step) Sort the edge set $\{e \in E(G) : w(e) \geq w_0\}$ as a sequence $S = e_1, \dots, e_m$ such that $w(e_1) \geq w(e_2) \geq \dots \geq w(e_m)$. $\mu \leftarrow 1$ where μ is the indicator of edges, $p \leftarrow 0$ where p is the indicator of communities, $\ell \leftarrow 1$ where ℓ is the indicator of the levels in the hierarchical system, and $\mathcal{L}_\ell \leftarrow \emptyset$ where \mathcal{L}_ℓ is the community sets in the ℓ th hierarchical level.
- Step 2.** (Starting a new search). $p \leftarrow p + 1$, $C_p \leftarrow V(e_\mu)$. $\mathcal{L}_\ell \leftarrow \mathcal{L}_\ell \cup \{C_p\}$.
- Step 3.** (Growing)

Substep 3.1. If $V(G) - V(C_p) = \emptyset$, go to Step 4, otherwise choose $v \in V(G) - V(C_p)$ such that $c(v, C_p)$ is a maximum. If

$$c(v, C_p) \geq \alpha_n d(C_p), \quad (3)$$

where $n = |V(C_p)|$ and $\alpha_n = 1 - \frac{1}{2\lambda(n+t)}$ with $\lambda \geq 1$ and $t \geq 1$ as user specified parameters, then $C_p \leftarrow C_p \cup \{v\}$ and go back to Substep 3.1.

Substep 3.2. $\mu \leftarrow \mu + 1$. If $\mu > m$ go to Step 4. Otherwise continue.

Substep 3.3. Suppose $e_\mu = xy$. If at least one of $x, y \notin \bigcup_{i=1}^{p-1} V(C_i)$ then go to Step 2, otherwise go to Substep 3.2.

Step 4. (Merging)

Substep 4.1. List all members of \mathcal{L}_ℓ as a sequence C_1, \dots, C_s such that

$$|V(C_1)| \geq |V(C_2)| \geq \dots \geq |V(C_s)|, \quad (4)$$

where $s = |\mathcal{L}_\ell|$, $h \leftarrow 2$, $j \leftarrow 1$.

Substep 4.2. If $|C_j \cap C_h| > \beta \min(|C_j|, |C_h|)$ (where β ($0 < \beta < 1$) is a user specified parameter), then $C_{s+1} \leftarrow C_j \cup C_h$ and the sequence \mathcal{L}_ℓ is rearranged as follows:

$$C_1, \dots, C_{s-1} \leftarrow \text{deleting } C_j, C_h \text{ from } C_1, \dots, C_{s+1} \quad (5)$$

$s \leftarrow s - 1$, $h \leftarrow \max\{h - 2, 1\}$, and go to Substep 4.4.

Substep 4.3. $j \leftarrow j + 1$. If $j < h$ go to Substep 4.2.

Substep 4.4. $h \leftarrow h + 1$ and $j \leftarrow 1$. If $h \leq s$ go to Substep 4.2.

Step 5. Contract each $C_p \in \mathcal{L}_\ell$ as a vertex:

$$V(G) \leftarrow \left[V(G) - \bigcup_{p=1}^s V(C_p) \right] \cup \{C_1, \dots, C_s\}, \quad (6)$$

$$w(uv) \leftarrow w(C_i, C_j) = \frac{\sum_{e \in E_{ij}} w(e)}{|E_{ij}|}, \quad (7)$$

if the vertex u is obtained by contracting C_i and v is obtained by contracting C_j where the set of crossing edges $E_{ij} = \{xy : x \in C_i, y \in C_j, x \neq y\}$. For $t \in V(G) - \{C_1, \dots, C_s\}$, define $w(t, C_i) = w(\{t\}, C_i)$. Other cases are defined similarly.

If $|V(G)| \geq 2$ then go to Step 6, otherwise go to END.

Step 6. $\ell \leftarrow \ell + 1$, $\mathcal{L}_\ell \leftarrow \emptyset$, $w_0 \leftarrow \gamma \max\{w(e) : \forall e \in E(G)\}$, and go to Step 1 (to start a new search in a higher level of the hierarchical system).

END.

If the input data is an unweighted graph G , the adjacency information are used for establishing the similarity matrix of G . Let $A = (a_{ij})$ be the adjacency matrix of G where

$$a_{ij} = \begin{cases} 1, & \text{there is an edge between node } i \text{ and node } j, \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

then the inner product of the i th and the j th row of A is used to describe the similarity between node i and j and stored as $G(i, j)$ in the similarity matrix G .

References

- Berkhin, P., 2009. Survey of clustering data mining techniques, <<http://www.ee.ucr.edu/barth/EE242/>>.
- Bondy, J.A., Murty, U.S.R., 1976. Graph theory with applications. Macmillan, London.
- Breiger, R.L., 1974. The duality of persons and groups. Social Forces 53 (2), 181–190.
- Danon, L., Duch, J., Diaz-Guilera, A., Arenas, A., 2005. Comparing community structure identification. J. Stat. Mech., P09008.
- Davis, A., Gardner, B.B., Gardner, M.R., 1941. Deep South: A Social Anthropological Study of Caste and Class. University of Chicago Press, Chicago.
- Diestel, R., 2005. Graph theory, . Graduate Texts in Mathematics, third ed., 173. Springer, heidelberg.
- Freeman, L., 1977. A set of measures of centrality based upon betweenness. Sociometry 40, 35–41.
- Frey, Brendan J., Dueck, Delbert, 2007. Clustering by passing messages between data points. Science 315, 972–976.

- Futschik, M.E., Carlisle, B., 2005. Noise-robust soft clustering of gene expression time course. *J. Bioinform. Comput. Biol.* 3, 965–988.
- Gary, M.R., Johnson, D.S., 1979. *Computers and Intractability*. Freeman, NY.
- Girvan, M., Newman, M.E.J., 2002. Community structure in social and biological networks. *Proc. Natl. Acad. Sci.* 99 (12), 7821–7826.
- Homans, G.C., 1950. *The human group*. Harcourt, Brace and World, New York.
- Jain, A.K., 2009. Data clustering: 50 years beyond *K*-Means. <<http://dataclustering.cse.msu.edu/papers/JainDataClusteringPRL09.pdf>>.
- Kernighan, B.W., Lin, S., 1970. An efficient heuristic procedure for partitioning graphs. *Bell Syst. Tech. J.* 49, 291–307.
- Krebs, V., 2009. <<http://www.orgnet.com/>>.
- Lusseau, D., 2003. emergent properties of a dolphin social network. *Proc. R. Soc. Lond. B* 270, S186–S188.
- Lusseau, D., Newman, M.E.J., 2004. Identifying the role that individual animals play in their social network. *Proc. R. Soc. Lond. B (Suppl.)* 271, S377–S481.
- Lusseau, D., Schneider, K., Boisseau, O.J., Haase, P., Slooten, E., Dawson, S.M., 2003. The bottlenose dolphin community of doubtful sound features a large proportion of long-lasting associations. Can geographic isolation explain this unique trait? *Behavioral Ecology and Sociobiology* 54, 396–405.
- Merriam-Webster Online Dictionary. 2008. Cluster analysis. <<http://www.merriam-webster-online.com>>.
- Newman, M.E.J., 2004a. Detecting community structure in networks. *Eur. Phys. J. B* 38, 321–330.
- Newman, M.E.J., 2004b. Fast algorithm for detecting community structure in networks. *Phys. Rev. E* 69, 066133.
- Newman, M.E.J., 2006. Finding community structure in networks using the eigenvectors of matrices. *Phys. Rev. E* 74, 036104.
- Newman, M.E.J., Girvan, M., 2004. Finding and evaluating community structure in networks. *Phys. Rev. E* 69, 026113.
- Newman, M.E.J., 2009. <<http://www-personal.umich.edu/mejn/netdata/>>.
- Palla, G., Derenyi, I., Farkas, I., Vicsek, T., 2005. Uncovering the overlapping community structure of complex networks in nature and society. *Nature*, Vol. 435 (7043), 814–818.
- Pereira-Leal, J.B., Enright, A.J., Ouzounis, C.A., 2004. Detection of functional modules from protein interaction networks. *PROTEINS: Struct. Func. Bioinform.* 54, 49–57.
- Rosvall, M., Bergstrom, C.T., 2008. Maps of random walks on complex networks reveal community structure. *Proc. Natl. Acad. Sci. USA* 105 (4), 1118–1123.
- Santamaria, Rodrigo, Roberto, Theron, 2008. Overlapping clustered graphs: Co-authorship networks visualization. *Lect. Notes Comput. Sci.* 5166, 190–199.
- SAS Institute Inc., 2003. Introduction to Clustering Procedures, Chapter 8 of SAS/STAT User's Guide. (SAS OnlineDocTM: Version 8) <http://www.math.wpi.edu/saspdf/stat/pdfdx.htm>.
- Scott, J., 2000. *Social Network Analysis: A handbook*. second ed. Sage Publications, London.
- Steinhaus, H., 1956. Sur la division des corp materiels en parties. *Bull. Acad. Polon. Sci.*, C1. III IV, 801–804.
- Wasserman, S., Faust, K., 1994. *Social Network Analysis*. Cambridge University Press, Cambridge.
- West, W., 1996. *Introduction to Graph Theory*. Prentice Hall, Upper Saddle River, NJ.
- Wu, F., Huberman, B.A., 2004. Finding communities in linear time: A physics approach. *Eur. Phys. J. B* 38, 331–338.
- Zachary, W.W., 1977. An information flow model for conflict and fission in small groups. *J. Anthropol. Res.* 33 (4), 452–473.