# Community Extraction in Hypergraphs Based on Adjacent Numbers

Hiroyuki Miyagawa      Maiko Shigeno      Satoshi Takahashi
Mingchao Zhang

Graduate School of Systems and Information Engineering, University of Tsukuba,
Ibaraki, 305-8573, Japan. *

**Abstract**    Community extraction plays an important role in network analysis. This paper discusses communities in hypergraphs. Since hypergraphs can represent multi-relational networks, they are important structures in many domains. In this study, the definition of communities based on the number of adjacent links/nodes is extended to a hypergraph. Four extended models are proposed. The adequacy of these models for hypergraph communities is investigated by using a hypergraph obtained from joint work relationships.

**Keywords**    Community; Hypergraph; Minimum cut algorithm

## 1   Introduction

Community extraction plays an important role in the analysis of networks such as the WWW and social networks. For example, communities in a web graph may correspond to sets of websites dealing with the related topics. Since web links are created to guide users to the related pages, subgraphs with dense connections can be regarded as communities. Thus, the process of finding a community is related to that of finding a subgraph so that connections are dense within and are sparse outside.

Many researchers have given different notions of communities and have developed various algorithms to detect communities. There are two types of community extraction approaches: one involves finding one or more communities at a time, and the other involves dividing a network into communities such that each node belongs to one of the communities. The latter is a general approach in community extraction. However, complete decomposition is not appropriate in many cases. Some objects may be classified as outliers or as objects that are not strongly connected to any specific group. Thus, this study focuses on the former approach. In a classical method of extracting one or more communities, one finds subgraphs relaxing a clique condition (See, for example [6] ). Flake, Lawrence and Giles [4] defined a community by using the number of adjacent links and proposed a flow-based algorithm.

Recently, hypergraphs have gained considerable attention, since they can represent complex situations for a network which has some attributions of elements. Community

extraction in hypergraphs has also been widely researched. Brinkmeier, Recknagel and Werner [3] defined a community by a minimum cut in a hypergraph, and extracted communities that included specified elements. Zhang, Takahashi and Shigeno [7] extended a maximum density subgraph to a hypergraph, and proposed an efficient extraction algorithm. Barber [2] proposed a community extraction method using modularity in a bipartite graph representing a hypergraph. His concept has been used to derive several variants for extraction algorithms.

In this study, we extend the definition of a community provided by Flake, Lawrence and Giles [4] to hypergraphs, and propose community extraction methods for hypergraphs. We also design four extended models. The extended community for each model can be identified efficiently by a flow-based algorithm. We compare the four proposed models by applying our method to real-world problems. In this study, we assume that an entire hypergraph is known and static.

## 2  Preliminaries

Let $G = (N, E)$ be an undirected graph with node set $N$ and edge set $E$. For a node $v \in N$ and a node set $C \subseteq N$, let $\delta_G(v, C) = \{(v, w) \in E \mid w \in C\}$. We call the cardinality of $\delta_G(v, C)$ the *adjacent number* of $v$ in $C$. In terms of undirected graphs, Flake, Lawrence and Giles [4] defined communities characterized of having more relations inside the community than outside.

**Definition 1.**
*A nonempty proper subset of nodes $C(\subset N)$ is a community, if each node $v$ in $C$ satisfies $|\delta_G(v, C)| \geq |\delta_G(v, N \setminus C)|$.*

Hereafter, the above community is referred to as an *adjacent number community*, or an *adj-community*, for short. For nodes $s$ and $t$, an *s-t cut* is a partition of the node set $(X, N \setminus X)$ such that $s \in X$ and $t \in N \setminus X$. A *minimum s-t cut* is an *s-t* cut that minimizes the number of edges connecting $X$ and $N \setminus X$, i.e., minimizes $|\{(u, v) \in E \mid u \in X, v \in N \setminus X\}|$. When graph $G$ has an edge capacity $c : E \to \mathbb{R}$, a minimum *s-t* cut implies an *s-t* cut $(X, N \setminus X)$ that minimizes the cut capacity $c(X, N \setminus X) = \sum\{c(u, v) \mid (u, v) \in E, u \in X, v \in N \setminus X\}$. When a directed graph is being considered, $c(X, N \setminus X)$ is defined by the sum of capacities of edges leaving from $X$ and entering to $N \setminus X$. Many efficient polynomial-time algorithms have been proposed for finding the minimum *s-t* cuts [1]. The following property, which is a slight variation of Theorem 1 in [4], shows that an adj-community can be found efficiently when using algorithms for an undirected graph with unit capacity.

**Theorem 1.**
*For specified nodes s and t, let $(X, N \setminus X)$ be a minimum s-t cut in G. If*

$$|\delta_G(s, X)| \geq |\delta_G(s, N \setminus X)|, \tag{1}$$

*then X is an adj-community.*

When the entire graph can be used for the calculation, adj-communities are found by repeating the following steps: **(Step 1)** In an appropriate manner, choose nodes $s$ and $t$; **(Step 2)** Find a minimum *s-t* cut $(X, N \setminus X)$; **(Step 3)** Check whether $X$ satisfies the condition (1); If yes, return $X$ as an adj-community. Flake, Lawrence and Giles [4]

proposed a method for choosing $s$ and $t$ for web graphs on the basis of the web properties. As pointed out by [5], note that their method cannot find all adj-communities that include $s$ and exclude $t$. Moreover, even if such communities exist, the method may find no communities. However, their method can be extended to one that operates over a graph induced by a crawl of fixed depth. The effectiveness of the extended method for a real-world case was reported in [4]. In addition, several variations of the adj-community have been discussed [5]. Since an adj-community is a basic idea for community extraction, the concept is worth focusing on. In this study, we extend this adj-community to hypergraphs.

At the last of this section, we introduce notations of hypergraphs. Given a finite set $N$, let $\mathscr{P}^*(N) = \{X \subseteq N \mid |X| \geq 2\}$ be the set of all subsets whose cardinalities are at least two. A hypergraph $\Gamma = (N, \mathscr{H})$ consists of a finite set $N$ of nodes and finite multi-set $\mathscr{H}$ of sets in $\mathscr{P}^*(N)$. Elements of $\mathscr{H}$ are referred to as *hyperedges*. A hypergraph is an extension of a graph in the sense that each hyperedge can connect more than two nodes. For each node $v \in N$ of hypergraph $\Gamma$, we define $\delta(v) = \{h \in \mathscr{H} \mid v \in h\}$.

# 3   Communities in hypergraphs

This section generalizes adj-communities to a hypergraph $\Gamma = (N, \mathscr{H})$. To do this, we generalize the notion of the adjacent number $|\delta_G(v, C)|$ of $v \in N$ in $C \subseteq N$. Since $\delta_G(v, C)$ is the set of edges, we introduce a set of adjacent hyperedges $\delta_\Gamma(v, C) = \{h \in \delta(v) \mid h \subseteq C \cup \{v\}\}$ for $v \in N$ and $C \subseteq N$. When $\Gamma$ is a graph $G$, i.e., the cardinality of each hyperedge is exactly two, it is clear that $\delta_G(v, C) = \delta_\Gamma(v, C)$ for any $v \in N$ and $C \subseteq N$. We then define communities based on the number of adjacent hyperedges as follows.

**Definition 2.**
*In a hypergraph $\Gamma = (N, \mathscr{H})$, a nonempty proper subset of nodes $C(\subset N)$ is a community, if each node $v \in C$ satisfies $|\delta_\Gamma(v, C)| \geq |\delta_\Gamma(v, N \setminus C)|$.*

Since the definition is based on the number of hyperedges, we call this community a *hyperedge-based-community*, or an *h-community* for short. We make a directed graph from $\Gamma$, in order to extract h-communities by using minimum *s-t* cut algorithms such as those in Flake, Lawrence and Giles's method [4]. Let $\mathscr{H}^+$ and $\mathscr{H}^-$ be two copies of $\mathscr{H}$. Denote by $h^+$ (resp. $h^-$) the copy of a hyperedge $h(\in \mathscr{H})$ belonging to $\mathscr{H}^+$ (resp. $\mathscr{H}^-$) is. We then construct a directed graph $\tilde{D}_\Gamma = (\tilde{N}, \tilde{A} \cup \tilde{A}^\pm)$, where $\tilde{N} = N \cup \mathscr{H}^+ \cup \mathscr{H}^-$, $\tilde{A} = \{(v, h^+)(h^-, v) \mid h \in \delta(v), v \in N\}$, and $\tilde{A}^\pm = \{(h^+, h^-)(h^-, h^+) \mid h \in \mathscr{H}\}$. An edge capacity $\tilde{c}: \tilde{A} \cup \tilde{A}^\pm \to \mathbb{R}$ is defined as

$$\tilde{c}(a) = \begin{cases} \infty & (a \in \tilde{A}), \\ 1 & (a \in \tilde{A}^\pm). \end{cases}$$

**Theorem 2.**
*For a pair of $s, t \in N$, suppose that $(Y, \tilde{N} \setminus Y)$ is a minimum s-t cut in $(\tilde{D}_\Gamma, \tilde{c})$. Let $C = Y \cap N$. If $|\delta_\Gamma(s, C)| \geq |\delta_\Gamma(s, N \setminus C)|$, then $C$ is an h-community.*

**Proof.** Note that $\tilde{c}(Y, \tilde{N} \setminus Y) < \infty$ since there exists an *s-t* cut with finite capacity, for example, $\tilde{c}(N \cup \mathscr{H}^+, \mathscr{H}^-) = |\mathscr{H}|$. Suppose $\hat{v} \in C$. If $\hat{v} \in h(\in \mathscr{H})$ and $h^+ \notin Y$, then the cut capacity $\tilde{c}(Y, \tilde{N} \setminus Y) \geq \tilde{c}(v, h^+) = \infty$. Thus, we obtain $\{h^+ \mid h \in \delta(\hat{v})\} \subseteq Y$, which implies that both of sets $\{h^+ \mid h \in \delta_\Gamma(\hat{v}, C)\}$ and $\{h^+ \mid \delta_\Gamma(\hat{v}, N \setminus C)\}$ are contained in

$Y$. Similarly, due to the edge set $\tilde{A}$, if $h^-$ belongs to $Y$, then $v$ contained in $h$ is in $C$. Thus, $\{h^- \mid h \in \delta_\Gamma(\hat{v}, N \setminus C)\} \subseteq \tilde{N} \setminus Y$ holds. In addition, if $h^- \notin Y$ for $h \in \delta_\Gamma(\hat{v}, C)$, we obtain $\tilde{c}(Y, \tilde{N} \setminus Y) = \tilde{c}(Y \setminus \{h^-\}, \tilde{N} \setminus (Y \setminus \{h^-\})) + 1$, which contradicts the minimality of $\tilde{c}(Y, \tilde{N} \setminus Y)$. Hence, we have $\{h^- \mid h \in \delta_\Gamma(\hat{v}, C)\} \subseteq Y$.

Assume that there exists $v \in C \setminus \{s\}$ with $|\delta_\Gamma(v, C)| < |\delta_\Gamma(v, N \setminus C)|$. By setting $Y' = Y \setminus (\{v\} \cup \{h^- \mid h \in \delta_\Gamma(v, C)\} \cup \{h^+ \mid h \in \delta_\Gamma(v, N \setminus C)\})$, we have

$$\tilde{c}(Y', \tilde{N} \setminus Y') = \tilde{c}(Y, \tilde{N} \setminus Y) + |\delta_\Gamma(v, C)| - |\delta_\Gamma(v, N \setminus C)| < \tilde{c}(Y, \tilde{N} \setminus Y),$$

which contradicts the minimality of $\tilde{c}(Y, \tilde{N} \setminus Y)$ since $Y'$ is an $s$-$t$ cut. $\square$
Thus, we can find an h-community by finding a minimum $s$-$t$ cut in $(\tilde{D}_\Gamma, \tilde{c})$ similar to the manner in which adj-communities are found in a graph.

The definition of h-communities appears to be a natural extension of adj-communities, since both definitions are based on the number of adjacent hyperedges or edges. However, h-communities do not consider all of the hyperedges that contain a node $v$ in a community. Therefore, for a node subset $C$, we classify all hyperedges by their contribution to relations in $C$. Let $\mathcal{H}(C)^> = \{h \in \mathcal{H} \mid |h \cap C| > |h \setminus C|\}$ and $\mathcal{H}(C)^\leq = \{h \in \mathcal{H} \mid |h \cap C| \leq |h \setminus C|\}$. When $\Gamma$ is a graph $G$, it is obvious that $\delta_G(v, C) = \delta(v) \cap \mathcal{H}(C)^>$ and $\delta_G(v, N \setminus C) = \delta(v) \cap \mathcal{H}(C)^\leq$ for any $C \subseteq N$ and $v \in C$. We then introduce a definition which uses all of the hyperedges in $\delta(v)$ for $v \in C$.

**Definition 3.**
*In a hypergraph $\Gamma = (N, \mathcal{H})$, a nonempty proper subset of nodes $C(\subset N)$ is a community, if each node $v \in C$ satisfies $|\delta(v) \cap \mathcal{H}(C)^>| \geq |\delta(v) \cap \mathcal{H}(C)^\leq|$.*

We call this community a *classified-hyperedges community*, or a *c-community* for short. Let $\tilde{G}_\Gamma = (N \cup \mathcal{H}, \tilde{E})$ be a bipartite graph, where $\tilde{E} = \{(v, h) \mid h \in \delta(v), v \in N\}$. This bipartite graph is widely used to represent $\Gamma$.

**Theorem 3.**
*For a pair of $s, t \in N$, suppose that $(Y, (N \cup \mathcal{H}) \setminus Y)$ is a minimum $s$-$t$ cut in $\tilde{G}_\Gamma$. Let $C = Y \cap N$. If $|\delta(s) \cap \mathcal{H}(C)^>| \geq |\delta(s) \cap \mathcal{H}(C)^\leq|$, then $C$ is a c-community.*

**Proof.** Let $\mathcal{H}(C)^= = \{h \in \mathcal{H} \mid |h \cap C| = |h \setminus C|\}$, and $\Delta(Y, (N \cup \mathcal{H}) \setminus Y)$ be the number of edges between $Y$ and $(N \cup \mathcal{H}) \setminus Y$. From the minimality of $\Delta(Y, (N \cup \mathcal{H}) \setminus Y)$, a hyperedge in $\mathcal{H} \cap Y$ belongs to $\mathcal{H}(C)^>$ or $\mathcal{H}(C)^=$. Similarly, a hyperedge in $\mathcal{H} \setminus Y$ belongs to $\mathcal{H}(C)^\leq$.

Suppose that there exists $v \in C \setminus \{s\}$ such that $|\delta(v) \cap \mathcal{H}(C)^>| < |\delta(v) \cap \mathcal{H}(C)^\leq|$. By setting $Y' = Y \setminus (\{v\} \cup (\delta(v) \cap \mathcal{H}(C)^=))$, we have

$$\begin{aligned}
&\Delta(Y', (N \cup \mathcal{H}) \setminus Y') \\
=\ &\Delta(Y, (N \cup \mathcal{H}) \setminus Y) + |\delta(v) \cap \mathcal{H}(C)^>| - |\delta(v) \setminus Y| - |\delta(v) \cap \mathcal{H}(C)^= \cap Y| \\
=\ &\Delta(Y, (N \cup \mathcal{H}) \setminus Y) + |\delta(v) \cap \mathcal{H}(C)^>| - |\delta(v) \cap \mathcal{H}(C)^\leq| < \Delta(Y, (N \cup \mathcal{H}) \setminus Y),
\end{aligned}$$

which contradicts the minimality of $\Delta(Y, (N \cup \mathcal{H}) \setminus Y)$, since $Y'$ is an $s$-$t$ cut. $\square$
Thus, we can also find a c-community by using a minimum $s$-$t$ cut algorithm for $\tilde{G}_\Gamma$.

We now return to the notion of the adjacent number in a graph. We can interpret the adjacent number as the total number of nodes including repeated memberships, each of

which belongs to $C$ and is jointed with $v$ by an edge. In other words, in a graph, the adjacent number $|\delta_G(v,C)|$ is equivalent to the cardinality of multiset $\{w \in N \mid (v,w) \in E, w \in C\}$. From this point of view, we extend the notion of adjacent numbers to the cardinality of a multiset of nodes, each of which belongs to $C$ and is contained in a hyperedge together with $v$. This multiset is denoted by $\hat{d}(v,C) = \{w \in C \mid h \in \delta(v), w \in h\}$.

**Definition 4.**
*In a hypergraph $\Gamma = (N, \mathscr{H})$, a nonempty proper subset of nodes $C(\subset N)$ is a community, if each node $v \in C$ satisfies $|\hat{d}(v,C)| \geq |\hat{d}(v, N \setminus C)|$.*

We call this community a *node-based-community*, or an *n-community* for short. It is clear that $|\hat{d}(v,C)| = \sum_{h \in \delta(v)} |h \cap (C \setminus \{v\})|$ holds. Hence, we can observe that an n-community corresponds to an adj-community in the undirected graph constructed by representing each hyperedge in $\Gamma$ by a clique. As described in the next section, n-communities do not seem to be so appropriate. One reason may be that n-communities ignore connections in one hyperedge.

We now give a definition combining the above three types of communities. For the strength of relations of $C(\subseteq N)$ and $v \in C$, we define a multiset as a collection of nodes, excluding $v$, that constitutes hyperedges covered by $C$. In other words, we consider a multiset $d(v,C)$ as defined by a collection of node sets $h \setminus \{v\}$, where $h \in \delta_\Gamma(v,C)$. In addition, to measure the relations between $v \in C$ and $N \setminus C$, we count the number of hyperedges not covered by $C$. Thus, we evaluate the relationships within and without a community by different measurements.

**Definition 5.**
*In a hypergraph $\Gamma = (N, \mathscr{H})$, a nonempty proper subset of nodes $C(\subset N)$ is a community, if each node $v \in C$ satisfies $|d(v,C)| \geq |\{h \in \delta(v) \mid h \nsubseteq C\}|$.*

We call this community a *mixed-criterion-community*, or an *mc-community* for short. When $\Gamma$ is a graph, an mc-community coincides with an adj-community. This mc-community can be extracted by a minimum *s-t* cut algorithm for a directed bipartite graph $D_\Gamma = (N \cup \mathscr{H}, A^F \cup A^B)$ with an edge capacity $c : A^F \cup A^B \to \mathbb{R}$, where $A^F = \{(v,h) \mid h \in \delta(v), v \in N\}$, $A^B = \{(h,v) \mid h \in \delta(v), v \in N\}$, and

$$c(e) = \begin{cases} 1 & (e \in A^F), \\ \infty & (e \in A^B). \end{cases}$$

**Theorem 4.**
*For a pair of $s,t \in N$, suppose that $(Y, (N \cup \mathscr{H}) \setminus Y)$ is a minimum s-t cut in $(D_\Gamma, c)$. Let $C = Y \cap N$. If $|d(s,C)| \geq |\{h \in \delta(s) \mid h \nsubseteq C\}|$, then $C$ is an mc-community.*

**Proof.** From the definition of the capacity, a hyperedge $h$ covered by $C$, if and only if $h \in Y$. If there exists a node $v \in C \setminus \{s\}$ with $|d(v,C)| < |\{h \in \delta(v) \mid h \nsubseteq C\}|$, we obtain, for an *s-t* cut $Y' = Y \setminus (\{v\} \cup \delta_\Gamma(v,C))$,

$$c(Y', (N \cup \mathscr{H}) \setminus Y') = c(Y, (N \cup \mathscr{H}) \setminus Y) + \sum_{h \in \delta_\Gamma(v,C)} |h \setminus \{v\}| - |\{h \in \delta(v) \mid h \nsubseteq C\}|$$

$$< c(Y, (N \cup \mathscr{H}) \setminus Y),$$

since $\sum_{h \in \delta_\Gamma(v,C)} |h \setminus \{v\}| = |d(v,C)|$. This inequality contradicts the minimality of $c(Y, (N \cup \mathscr{H}))$. $\square$

# 4  Experimental results

To investigate the adequacy of our communities, we performed experiments on our community extraction in a hypergraph representing joint work relationships. For a node set of a hypergraph, we collected authors of articles that appear in a specific journal. Authors for whom there were no joint works in the journal, were excluded. Each hyperedge was given by a set of co-authors for an article, while an article written by a single author was ignored.

Our experiments found a smaller community that contained a specified author under a required condition. The following is the procedure for finding an h-community.

**Step 1**  Set a specified author as $s$.

**Step 2**  For each node $v \in N \setminus \{s\}$, find a minimum $s$-$v$ cut $(Y_v, \tilde{N} \setminus Y_v)$ in $(\tilde{D}_\Gamma, \tilde{c})$. Let $(Y, \tilde{N} \setminus Y)$ be a cut attaining the minimum capacity among the obtained cuts satisfying the condition $|\delta_\Gamma(s, Y_v \cap N)| \geq |\delta_\Gamma(s, N \setminus Y_v)|$. (If there is no cut satisfying the condition, the procedure fails.)

**Step 3**  Choose a node $v$ attaining $\min\{\tilde{c}(Y_v, \tilde{N} \setminus Y_v) \mid v \in (Y \cap N) \setminus \{s\}\}$.

**Step 4**  If $C = (Y \cap Y_v) \cap N$ satisfies the condition $|\delta_\Gamma(s, C)| \geq |\delta_\Gamma(s, N \setminus C)|$, then update $(Y, \tilde{N} \setminus Y)$ by $(Y \cap Y_v, \tilde{N} \setminus (Y \cap Y_v))$ and go to Step 3. Otherwise, return $Y \cap N$ as an h-community.

We can find a c-community, n-community, and mc-community in a similar manner.

In the first experiment, we obtained a hypergraph from the Journal of the Operations Research Society of Japan, using vols. 38–56 (1995–2009). The obtained hypergraph had 570 nodes and 344 hyperedges. We selected two authors, Author1 and Author2, as the central researchers. Author1 had nine hyperedges, and Author2 had six hyperedges in this hypergraph. For Author1, the obtained h-, c-, n-, and mc-communities were composed of 18, 25, 22 and 29 authors, respectively. Figure 1 (a) shows the number of authors belonging to individual relations among the obtained communities as a Venn diagram. Also, Figure 1 (b) shows the result for Author2. We also show the result for Author2 on a subgraph of $\tilde{G}_\Gamma = (N \cup \mathcal{H}, \tilde{E})$ in Figure 2.



(a) for Author1

(b) for Author2

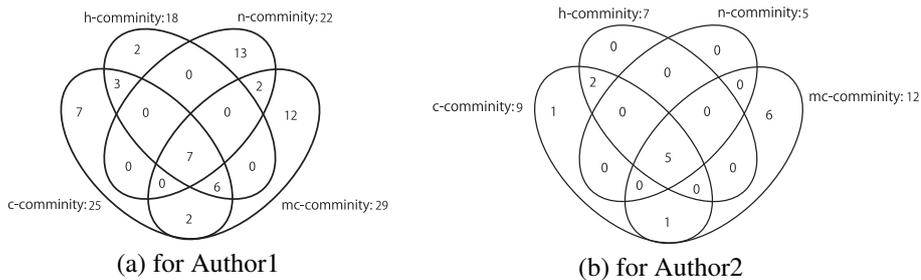Figure 1: The numbers of authors belonging to the obtained communities

In the second experiment, we used articles from Mathematical Programming Journal, Series A and Series B from vol.43 No.1(1989) to vol.115, No.10 (2008) in order to construct a hypergraph, which contained 1800 nodes and 1215 hyperedges. We selected the same author Author2 in the first experience. Author2 had 17 hyperedges in

Figure 2: The result for Author2 in the first experiment on a bipartite subgraph of $\tilde{G}_\Gamma = (N \cup \mathscr{H}, \tilde{E})$. Circles and rectangles represent authors and articles, respectively. Signs of "h-", "c-", "n-" and "mc-" beside each author node express types of the communities in which the author belongs to.

this hypergraph. We also selected Author3, who had 17 hyperedges. Figure 3 shows the numbers of authors belonging to individual relations among the obtained h-, c-, n-, and mc-communities for Author2. We also show the result for Author2 on a subgraph of $\tilde{G}_\Gamma = (N \cup \mathscr{H}, \tilde{E})$ in Figure 4.



(a) for Author2                         (b) for Author3

Figure 3: The numbers of authors belonging to the obtained communities

In our experiments, the obtained communities could extract groups of the authors carrying out active research in the fields of authors selected as a source. C-communities tend to be influenced by large-size hyperedges. Indeed, for Author1, a hyperedge consisting of six authors was related. As a result, the obtained c-community contained these six authors. The result for Authors3 was similar. In contrast, for Author2, the size of each related hyperedge was less than four, and the obtained c-communities were included in other communities. In our experiments, h-communities tended to be similar to c-communities. The result shown in Figure 2 illustrates that we could not obtain the minimal communities.

## 5   Conclusion

We proposed four definitions for hypergraph communities on the basis of adjacent numbers, not only an extension to a bipartite graph that represents a hypergraph. We
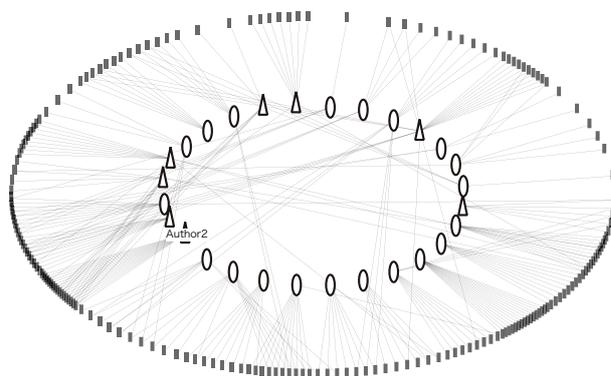
Figure 4: The result for Author2 in the second experiment on a bipartite subgraph of $\tilde{G}_\Gamma = (N \cup \mathcal{H}, \tilde{E})$. Each circle and triangle shows an author and a rectangle shows an article. All authors in this figure are in h-community. Authors that are contained in all obtained communities are illustrated by triangles.

also present the experimental results of extracting these communities in hypergraphs that represent joint work relationships. The experimental results suggest that the appropriate community must be chosen on the basis of the hypergraph characteristics, especially, the size of the hyperedges.

## Acknowledgements

# References

[1] R. K. Ahuja, T. L. Magnanti and J. B. Orlin. Network Flows: Theory, Algorithms, and Applications. Prentice Hall, 1993.

[2] M. J. Barber. Modularity and community detection in bipartite networks. Physical Review E, 76, 066102, 2007.

[3] M. Brinkmeier, S. Recknagel, and J. Werner. Communities in graphs and hypergraphs. In Proceedings of the sixteenth ACM conference on Conference on information and knowledge management, 869-872, 2007.

[4] G. Flake, S. Lawrence, and C. Giles. Efficient identification of web communities. In Proceedings of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 150-160, 2000.

[5] H. Ino, M. Kudo and A. Nakamura. Partitioning of web graphs by community topology. In Proceedings of the 14th international conference on World Wide Web, 661-669, 2005.

[6] S. Wasserman and K. Faust. Social Network Analysis: Methods and Applications. Cambridge University Press, 1994.

[7] M. Zhang, S. Takakashi, and M. Shigeno. A maximum density subset problem and its algorithm with approximate binary search (in Japanese). Transactions of the Operations Research Society of Japan, Vol. 53, 1-13, 2010.