

## یک راه حل افزایشی جهت خوشه بندی محتوایی - ساختاری یک گراف

سامان کشوری<sup>۱</sup>، محمدعلی جوادزاده<sup>۲\*</sup>، حسن نادری<sup>۳</sup>

۱- دانشجوی کارشناسی ارشد، ۲- استادیار، دانشگاه جامع امام حسین (ع) ۳- استادیار، دانشگاه علم و صنعت ایران  
(دریافت: ۹۸/۰۱/۲۱، پذیرش: ۹۸/۰۵/۱۵)

## چکیده

خوشه بندی گره های گراف از جنبه ساختاری یا محتوایی، همواره مورد توجه پژوهشگران حوزه داده کاوی بوده است؛ اما به خوشه بندی گراف بر مبنای ساختار و محتوا به طور ترکیبی کمتر توجه شده است. با توجه به نیاز خوشه بندی ساختاری-محتوایی در شبکه های اطلاعاتی که شبکه های اجتماعی نمونه ای از آن هاست، در این مقاله الگوریتم خوشه بندی ICS-Cluster ارائه شده که هر دو جنبه ساختار و محتوا را به صورت هم زمان در نظر می گیرد. هدف این روش، رسیدن به خوشه هایی با ساختار درونی منسجم (ساختاری) و مقادیر ویژگی (محتوایی) همگن در گراف است. در این روش ابتدا گراف اولیه به یک گراف ساختاری-محتوایی تبدیل می شود که در آن وزن هر یال (ارتباط) بیانگر شباهت ساختاری-محتوایی دو گره (موجودیت) است. خوشه بندی با توجه به وزن یال ها به صورت افزایشی انجام می شود بدین معنا که گره های یال با وزن بالا به عنوان خوشه در نظر گرفته می شوند و وزن یال های متصل به خوشه با یکدیگر ادغام شده و به صورت یک یال متصل به خوشه در نظر گرفته می شوند، این مراحل تا زمانی که الگوریتم به تعداد خوشه مورد نظر کاربر برسد، ادامه خواهد یافت. الگوریتم ICS-Cluster به هر تعداد خوشه که مدنظر کاربر است، گراف را خوشه بندی می کند. مقایسه الگوریتم مطرح شده با سه الگوریتم خوشه بندی ساختاری-محتوایی ارائه شده، بر اساس معیارهای شش گانه سنجش کیفیت خوشه، بیانگر عملکرد مناسب روش ICS-Cluster است. این معیارها معیارهای ساختاری تراکم خوشه، خطای یال و پیمانی، معیار محتوایی میانگین شباهت، معیار ساختاری-محتوایی CS-Measure و زمان اجرای روش ها است.

**کلیدواژه ها:** خوشه بندی گراف، خوشه بندی ساختاری-محتوایی، گراف خصوصیت، استخراج خوشه

## An Incremental Solution for Content-Structural Graph Clustering

S. Keshvari, M. A. Javadzadeh\*, H. Naderi

Imam Hossein University

(Received: 10/04/2019; Accepted: 06/08/2019)

## Abstract

Researchers have always been interested in graph nodes clustering based on content or structure. But less attention has been paid to clustering based on both structure and content. But a content-structural clustering is needed in information networks like social networks. In this paper, the ICS-Cluster algorithm is proposed which takes into consideration both the structure and content aspects of the nodes. The purpose of this approach is to gain a coherent internal structure (structural aspect) and homogeneous attribute values (content aspect) in the graph. In this approach firstly the graph is converted into a content-structural graph which edges' weight show similarity between the connected nodes. Incremental clustering is done based on edges' weight in this process the edges with the most weight is considered as clusters then the weight of connected edge to the cluster is aggregated and they'll be one edge, the process is repeated until the algorithm reaches the number of clusters that indicated by the user. ICS-Cluster algorithm number of cluster is indicated by the user. Comparing ICS-Cluster with other content structural algorithm based on six criteria for measuring cluster quality shows that ICS-Cluster has good performance. These criteria contain structural criteria (Modularity, Error Link, and Density), content criterion (Average Similarity), content-structural criterion (CS-Measure) and the run time.

**Keywords:** Graph Clustering, Content-Structural Clustering, Attributed Graph, Cluster Extraction

\*Corresponding Author E-mail: javadzadeh@ihu.ac.ir

## ۱. مقدمه

معنایی شبیه به هم هستند. در این حالت باید ارتباط ساختاری و معنایی بین خوشه‌های مجزا به حداقل خود برسد. در این مقاله یک روش خوشه‌بندی سلسله‌مراتبی ارائه شده است که ابتدا گراف اولیه - که در آن یال نشان‌دهنده ارتباط ساختاری و وزن آن بیانگر شباهت محتوایی است - به گرافی ساختاری-محتوایی تبدیل می‌شود. در این گراف دو گره که شباهت ساختاری-محتوایی آن‌ها بیش از حد آستانه است به یکدیگر متصل شده و وزن آن میزان شباهت ساختاری-محتوایی خواهد بود.

روش خوشه‌بندی ICS-Cluster<sup>۲</sup> بر روی گراف مذکور انجام می‌شود. در این روش ابتدا گره‌های اطراف یال با بیش‌ترین وزن به‌عنوان خوشه در نظر گرفته می‌شود و گره‌های دیگر به‌صورت مرحله‌به‌مرحله به خوشه‌های تشکیل شده اضافه می‌شود.

در ادامه مقاله، در بخش ۲، کارهای انجام شده در زمینه خوشه‌بندی گراف به‌صورت مختصر مطرح شده است. در بخش ۳، روش تحقیق و روش ارائه شده در این مقاله شرح داده شده است. در بخش ۴، با پیاده‌سازی روش‌های خوشه‌بندی بر روی مجموعه داده، عملکرد روش پیشنهادی با سایر روش‌های خوشه‌بندی ساختاری-محتوایی ارزیابی و مقایسه شده است. در پایان و در بخش ۵، نتیجه‌گیری کلی مقاله مطرح می‌شود.

## ۲. پیشینه تحقیق

خوشه‌بندی گراف یکی از راه‌های تحلیل گراف‌های حجیم و پیچیده است که تاکنون روش‌های خوشه‌بندی مختلفی برای آن ارائه و استفاده شده است [۱۳-۱۴]. بیشتر این روش‌ها برای عمل خوشه‌بندی تنها جنبه ساختاری گراف را در نظر می‌گیرند [۱۵-۱۶]، که از این دسته می‌توان TopGC و Louvain را نام برد [۱۷-۱۸].

برخی روش‌ها نیز مطرح شده‌اند که عمل خوشه‌بندی را بر اساس محتویات گره‌ها انجام می‌دهند، از جمله این دست روش‌ها می‌توان [۱۹] CAMAS k-Means و k-Medoids [۲۰-۲۱] را نام برد که تحقیقات برای ارائه روش خوشه‌بندی محتوایی همچنان در حال انجام است [۲۲-۲۳]. در بسیاری از کاربردهای دنیای واقعی هر دو جنبه ساختار و محتوا در کنار یکدیگر و به‌طور هم‌زمان موردنظر هستند [۲۴-۲۵].

به‌عنوان مثال این مسئله در کشف جوامع در شبکه‌های اجتماعی از اهمیت بالایی برخوردار است [۲۶-۲۸]، بنابراین، بهتر است در خوشه‌بندی شباهت ساختاری و محتوایی گره‌ها در کنار یکدیگر در نظر گرفته شود. در ادامه برخی از روش‌های خوشه‌بندی ساختاری-محتوایی بیان شده است.

گراف کاربرد بسیار وسیعی در مباحث علمی دارد زیرا بسیاری از مسائل واقعی علمی جهت تحلیل بر روی گراف مدل می‌شوند [۲، ۳]. از گراف به‌منظور تحلیل مسائل مختلف استفاده می‌شود [۴-۵]، شبکه‌های اجتماعی نیز بر روی گراف مدل می‌شوند و از آنجایی که استفاده از آن‌ها در جامعه در حال افزایش است حجم گراف حاصل از آن‌ها افزایش می‌یابد بنابراین، پرداختن به خوشه‌بندی آن‌ها بیش‌ازپیش اهمیت دارد [۶-۷]. محققان همواره سعی بر بهبود الگوریتم‌های خوشه‌بندی به‌منظور افزایش کیفیت خوشه‌های به‌دست‌آمده دارند [۸-۹].

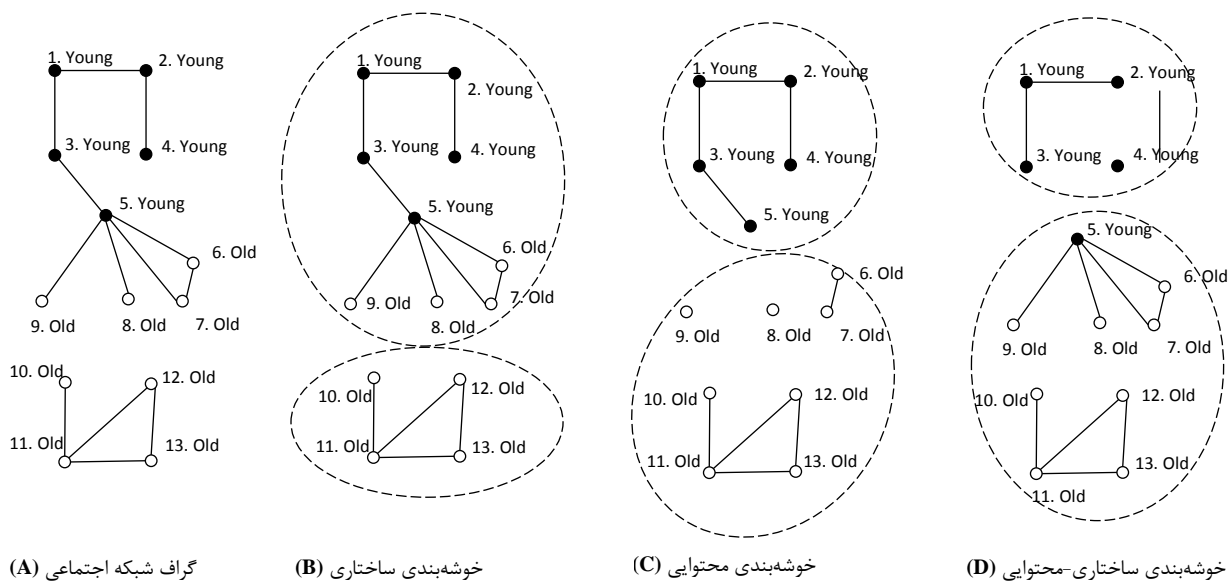
خوشه‌بندی مسئله‌ی مهم در تحلیل داده‌های گراف و یک روش دسته‌بندی بدون ناظر برای داده‌های گراف است. خوشه‌بندی شامل تقسیم گره‌های گراف به گروه‌هایی با گره‌های مشابه است. شباهت بین گره‌ها معمولاً توسط یک تابع هدف ریاضی تعریف می‌شود [۱۰]. معمولاً مسئله خوشه‌بندی گراف در رده مسائل NP-Hard قرار می‌گیرد.

حل این مسائل به صورت کلی از طریق روش‌های مکاشفه‌ای<sup>۱</sup> و تقریبی حاصل می‌شود. این توابع هدف در روش‌های مکاشفه‌ای به دو شیوه محلی و سراسری تعریف می‌شوند [۱۱-۱۲]. اغلب روش‌های خوشه‌بندی که تاکنون ارائه شده است فقط یک جنبه ساختاری یا محتوای گراف را در نظر گرفته‌اند و روش‌های خوشه‌بندی که هم ساختار و هم محتوا را در نظر بگیرند کمتر ارائه شده است.

به‌عنوان مثال خوشه‌بندی‌های مختلف افراد در شبکه‌های اجتماعی که در شکل (۱) نشان داده شده است. در این خوشه‌بندی اگر نیاز باشد افرادی که با یکدیگر به‌واسطه ویژگی‌های شخصی (در این مثال سن) در ارتباط بوده را خوشه‌بندی کرد، نمی‌توان تنها ویژگی‌های محتوایی یا ساختاری را در نظر گرفت. همان‌طور که در شکل نشان داده شده، اگر خوشه‌بندی بر اساس ساختار گراف تشکیل شده انجام شود، ویژگی‌های افراد نادیده گرفته خواهد شد، اگر ویژگی افراد مد نظر قرار گیرد، ارتباط بین آن‌ها نادیده گرفته خواهد شد؛ بنابراین، بهتر است ویژگی‌های ساختاری و محتوایی گراف به‌صورت هم‌زمان در نظر گرفته شود. مطابق با شکل (۱) برای فرد شماره ۵ باینکه از لحاظ محتوایی جوان است ولی چون به لحاظ ساختاری با چهار فرد میان‌سال در ارتباط است، در خوشه مربوط به میان‌سال‌ها قرار خواهد گرفت.

هدف کلی، تشخیص خوشه‌هایی از گراف است که گره‌های درونی هر خوشه به‌طور متراکم به یکدیگر مرتبط بوده و از لحاظ

<sup>۲</sup> Incremental Content-Structural Clustering<sup>۱</sup> Heuristic Methods



شکل ۱. نمونه شبکه اجتماعی با محتوای "سن"

همچنین خوشه‌های متراکم‌تری به دست می‌آورد.

DCM<sup>۴</sup>: این الگوریتم از مجموعه‌ای از جوامع کاندید شروع کرده و دو گام اصلی را به‌طور متناوب تا رسیدن به همگرایی تکرار می‌کند. در گام اول جوامع با بهترین امتیاز جامعه (از نظر ساختاری) را به دست آورده و در گام بعد سعی می‌کند یک توصیف مناسب برای این جوامع به دست آورد. برای تطابق بیشتر توصیف با جوامع، ممکن است در صورت لزوم جامعه را مقداری تغییر داده و گام‌های الگوریتم را تا همگرا شدن تکرار کند [۳۱]. این روش لزوماً همه گره‌ها را پوشش نمی‌دهد و برای مجموعه داده‌های کوچک کند و برای مجموعه داده‌های بزرگ سرعت بالا دارد.

CS-Cluster<sup>۵</sup>: در این الگوریتم، ابتدا گراف وزن دار ورودی -که در آن وزن بیانگر شباهت محتوایی است- به یک گراف وزن دار که وزن آن بیانگر فاصله ساختاری و یال آن نشان‌دهنده شباهت محتوایی دو گره است تبدیل می‌شود. سپس گره‌های مرکزی به کمک رابطه (۱) به دست می‌آیند.

$$C_v = \frac{(D_v)^2}{\sum_{e \in E_v} w_e} \quad (1)$$

که در آن  $C_v$  امتیاز مرکزیت گره  $v$ ،  $E_v$  مجموعه یال‌های متصل به  $v$  در گراف دوم،  $w_e$  وزن یال  $e$  در گراف دوم و  $D_v$  درجه گره  $v$  در گراف دوم است.

پس از انتخاب گره‌های مرکزی و در نظر گرفتن آن‌ها به عنوان مرکز خوشه، بر اساس فاصله موجود میان این گره‌ها

SA-Cluster<sup>۱</sup>: در این الگوریتم به ازای هر خصوصیت موجود در هر گره، یک گره خصوصیت به گراف اضافه می‌شود و گره‌هایی که دارای آن خصوصیت هستند از طریق یک یال خصوصیت به آن گره وصل می‌شوند. این الگوریتم یک الگوریتم تکرارشونده است که در آن ابتدا مرکزی‌ترین گره‌ها (گره‌ها با تراکم بالا) به عنوان مرکز خوشه انتخاب شده و سپس باقی گره‌ها به نزدیک‌ترین مراکز، تخصیص داده می‌شوند. سپس در هر خوشه مراکز به هنگام رسانی می‌شوند (همچنین وزن یال‌ها که بر اساس ویژگی‌های محتوایی و ساختاری تنظیم شده‌اند و سایر پارامترهای الگوریتم به روز می‌شود) و عملیات تا همگرا شدن تکرار می‌شود [۲۹]. از آنجاکه این الگوریتم به صورت تکرارشونده است، زمان اجرای بالایی دارد.

SANS<sup>۲</sup>: در این الگوریتم، ابتدا در گراف وزن دار ورودی شاخص وزن هر گره -که به صورت مجموع وزن یال‌های متصل به گره است- محاسبه شده و گره‌ای که بیش‌ترین شاخص وزن را دارد به عنوان مرکز<sup>۳</sup> خوشه انتخاب می‌شود؛ سپس گره‌های همسایه گره مرکزی به خوشه مربوط به آن اضافه می‌شوند و ادامه گره‌هایی که با گره‌های درون خوشه شباهت (شباهت محتوایی) بیش از حد آستانه دارند نیز به خوشه اضافه می‌شوند. پس از طی این مراحل، مجدداً از بین گره‌های باقیمانده گره‌ای که بیش‌ترین شاخص وزن را دارد به عنوان مرکز بعد انتخاب می‌شود و عملیات تا خوشه‌بندی همه گره‌ها تکرار می‌شود [۳۰]. زمان اجرای این الگوریتم نسبت به SA-Cluster بهتر است و

<sup>۴</sup> Description-Driven Community Detection  
<sup>۵</sup> Content-Structural Clustering

<sup>۱</sup> Structural Attribute Cluster  
<sup>۲</sup> Structural Attribute Neighbourhood Similarity  
<sup>۳</sup> Centroid

بین دو گره  $i$  و  $j$  در نظر گرفته خواهد شد.

$$W_{(i,j)} = a(sim_{(i,j)}) + b\left(\frac{1}{D_{(i,j)}}\right) \quad (2)$$

که در آن  $i \neq j$  بوده و  $a$  میزان اهمیت (وزن) شباهت محتوایی و  $b$  اهمیت (وزن) شباهت ساختاری میان دو گره  $i$  و  $j$  است. شباهت ساختاری بین دو گره، عکس فاصله آن دو از یکدیگر در نظر گرفته شده است. در این رابطه  $D(i,j)$  مطابق با رابطه (۳) برابر با تعداد یال (فاصله) بین دو گره  $i$  و  $j$  است. در واقع هرچه این مقدار بیشتر باشد، شباهت ساختاری کمتر خواهد بود. این فاصله از طریق نسخه Modified الگوریتم دایجسترا<sup>۴</sup> محاسبه می‌شود که pseudocode آن در شکل (۲) مشاهده می‌شود [۳۳].

$$D(i,j) = \sum_{e=i}^j E_e \quad (3)$$

**Input:** A weighted graph  $G(V, E, c)$ , where  $V$  is a set of vertices (represent nodes),  $E$  is a set of edges connecting the vertices in  $V$  (they represent the connectivity between nodes), and  $c$  is a set of cost values (measured in ETX) associated with the edges in  $E$ . The SN id is also given as an input.

**Output:** A shortest path tree rooted at the SN ( $s$ ) and that reaches each vertex in  $V$  through a shortest path.

**BEGIN**

```

1:  $N \leftarrow s$  //  $N$  contains only the source node
2:  $d(s) \leftarrow 0$ 
3:  $r \leftarrow s$ 
4: for each node  $n$  in  $V$  (except the SN  $s$ ) not in  $N$  do
5:    $d(n) \leftarrow \infty$ 
6: end for
7: repeat
8:   for each  $v$  directly connected to  $r$ , and not in  $N$  do
9:     if  $p(r)$  not NULL then
10:       $c(r, v) = c(r, v) - prob[p(r)][v] \times c(p(r), r)$ 
11:     end if
12:      $d(v) = \min\{d(v), d(r) + c(r, v)\}$ 
13:   end for
14:   find node  $x$  not in  $N$  and  $d(x)$  is minimum
15:   insert  $x$  into  $N$ 
16:    $p(x) = r$ 
17:    $r = x$ 
18: until  $N = V$ 
END

```

شکل ۲. الگوریتم دایجسترا استفاده شده [۳۳]

مقدار  $sim_{(i,j)}$  نیز از طریق شباهت جاکارد طبق رابطه (۴) محاسبه می‌شود [۳۴].

$$sim_{(i,j)} = J(C, C') = \frac{|C \cap C'|}{|C \cup C'|} \quad (4)$$

که در آن  $J$  شباهت جاکارد است که برابر است با میزان اشتراک موجود بین دو گره بر روی کل معیارهای موجود در گره، که در بهترین حالت، برابر با یک است.

به منظور ساخت گراف ساختاری-محتوایی، اگر  $W_{(i,j)}$  از مقداری که توسط کاربر مشخص می‌شود بیشتر بود، بین آن دو یالی را رسم کرده و وزن آن یال برابر با مقدار  $W_{(i,j)}$  خواهد بود. در غیر این صورت یالی رسم نمی‌شود. لازم به ذکر است که حد آستانه مورد نظر به صورت پیش فرض برابر با میانگین وزن یال‌ها

خوشه‌ها تا زمانی که تمامی گره‌ها پوشش داده شوند گسترش می‌یابند. نتایج حاصل از این پژوهش نشان می‌دهد از آنجاکه در طول انجام مراحل انتخاب گره‌های مرکزی تا گسترش خوشه در این روش هم‌زمان جنبه ساختار و محتوا دیده می‌شود، این روش به لحاظ ساختاری-محتوایی دارای عملکرد بهتری نسبت به سایر روش‌ها است [۳۲]. از نقاط ضعف‌های این روش می‌توان سرعت نامناسب، تراکم<sup>۱</sup> و پیمانگی<sup>۲</sup> به نسبت پایین اشاره کرد.

در روش‌های ارائه شده تاکنون میزان اهمیت شباهت محتوایی با ساختاری یکسان در نظر گرفته شده است، ولی در برخی از خوشه‌بندی‌ها میزان اهمیت شباهت ساختاری نسبت به شباهت محتوایی بیشتر بوده و بالعکس. به عنوان مثال در شبکه اجتماعی که در آن اطلاعات پروفایل کاربری زیادی وجود ندارد، اهمیت ساختاری گراف بیش از اهمیت محتوایی است، در حالی که در خوشه‌بندی صفحات وب، محتوا از اهمیت بالاتری نسبت به ساختار برخوردار است. در ادامه روشی ارائه شده که ضمن پارامتریک کردن اهمیت هر کدام، عملکرد مناسبی به لحاظ معیارهای مختلف از خود نشان می‌دهد.

### ۳. روش پیشنهادی

یک گراف ویژگی را به صورت  $G=(V,E,A)$  نشان داده می‌شود، به طوری که در آن  $V$  مجموعه گره‌ها،  $E$  مجموعه یال‌ها و  $A$  مجموعه خصوصیات هر گره گراف را نشان می‌دهد. به هر  $v_i \in V$  یک بردار ویژگی  $[a_1(v_i), \dots, a_m(v_i)]$  تخصیص داده می‌شود که در آن  $a_j(v_i)$  مقدار خصوصیت  $a_j$  در گره  $v_i$  است. سعی بر آن است گراف ویژگی  $G$  به صورت خوشه‌بندی شود که هر کدام از خوشه‌ها علاوه بر نزدیکی ساختاری، از لحاظ محتوایی نیز مشابه باشند. در ادامه روش ارائه شده تشریح می‌شود.

#### ۳-۱. روش ارائه شده (ICS-Cluster<sup>۳</sup>)

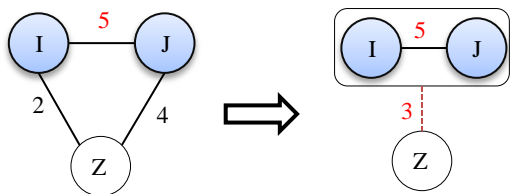
در این روش ابتدا با ترکیب ساختار و محتوای گراف اولیه گرافی ساختاری-محتوایی تشکیل می‌شود. همان طور که اشاره شد، در روش‌های قبلی میزان اهمیت شباهت محتوایی با ساختاری یکسان در نظر گرفته شده است ولی در روش ارائه شده کاربر می‌تواند با دادن ضریب (وزن)، اهمیت هر کدام را نسبت به دیگری مشخص نماید. اگر شباهت محتوایی بین دو گره  $i$  و  $j$  از  $sim_{ij}$  و فاصله ساختاری را  $D$  در نظر گرفته شود شباهت ساختاری-محتوایی طبق رابطه (۲) محاسبه می‌شود. در گراف ساختاری-محتوایی که ایجاد می‌شود این مقدار به عنوان وزن یال

<sup>1</sup> Density

<sup>2</sup> Modularity

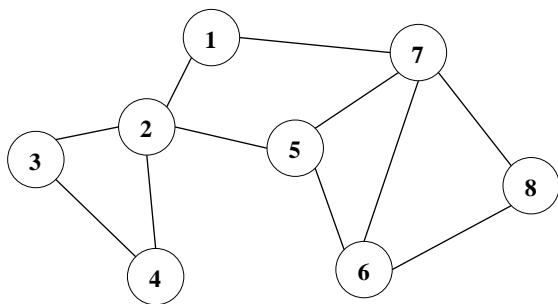
<sup>3</sup> Incremental Structural Content Clustering

<sup>4</sup> Dijkstra



شکل ۳. یک نمونه ادغام گره در روش ICS-Cluster

برای ادامه خوشه‌بندی مجدداً بیشترین مقدار وزن گره را در نظر گرفته و تا زمانی که تعداد خوشه‌ها به خوشه‌های مدنظر کاربر برسد، عمل خوشه‌بندی انجام خواهد شد. در ادامه مثالی از نحوه خوشه‌بندی ارائه می‌شود که در آن از مثال ارائه‌شده در [۳۲] استفاده شده است تا تفاوت‌های روش ارائه‌شده با آن روش به خوبی مشخص شود. در این مثال تعداد خوشه موردنظر کاربر ۲ در نظر گرفته شده است. در ابتدا گراف ساختاری را که در شکل (۴) داده شده است برای شروع کار در نظر گرفته می‌شود.



شکل ۴. ساختار گراف نمونه

ماتریس ارائه‌شده در زیر بیانگر میزان شباهت دو گره از لحاظ محتوایی به یکدیگر هستند.

$$sim = \begin{bmatrix} - & 0.40 & 0.52 & 0.60 & 0.32 & 0.30 & 0.30 & 0.15 \\ - & - & 0.60 & 0.60 & 0.22 & 0.48 & 0.19 & 0.36 \\ - & - & - & 0.70 & 0.42 & 0.52 & 0.21 & 0.32 \\ - & - & - & - & 0.30 & 0.23 & 0.09 & 0.16 \\ - & - & - & - & - & 0.40 & 0.70 & 0.50 \\ - & - & - & - & - & - & 0.45 & 0.47 \\ - & - & - & - & - & - & - & 0.48 \\ - & - & - & - & - & - & - & - \end{bmatrix}$$

درواقع با مشخص شدن شباهت بین دو گره در هر یال، گراف به صورت شکل (۵) در خواهد آمد که در آن وزن هر یال بیانگر شباهت محتوایی دو گره که ارتباط ساختاری دارند است.

در نظر گرفته شده است. بنابراین، گراف تشکیل شده یک گراف وزن دار است که وزن هر یال نشان‌دهنده شباهت ساختاری-محتوایی دو گره است. حین تشکیل گراف ساختاری-محتوایی از روی گراف قبلی، یال‌های (ارتباط‌های) حساس به دست می‌آیند. این یال‌ها شامل یال‌هایی است که بیشترین وزن را دارند. در این شرایط یال‌های حساس طبق رابطه (۵) به دست می‌آیند.

$$K\text{-Sensitive-edge} = K\text{-MAX}(w_i) \quad (5)$$

که منظور از K-Sensitive-edge ترتیب یال‌هایی است که دارای بیشترین مقدار وزن هستند و به عنوان یال‌های حساس محسوب می‌شوند. در این مرحله نه تنها یال‌های حساس به دست آمده بلکه گره‌های حساس نیز از طریق رابطه (۶) به دست می‌آید.

$$K\text{-Sensitive-node} = K\text{-MAX}(\sum_{i \in n} w_i) \quad (6)$$

که در آن، منظور از K-Sensitive-node تعداد K گره حساس است که برابر است با گره‌هایی که در آن جمع وزن یال‌های متصل به آن حداکثر هستند. در این رابطه i رئوس موردنظر، n تمامی رئوس و w وزن یال‌های متصل به گره است.

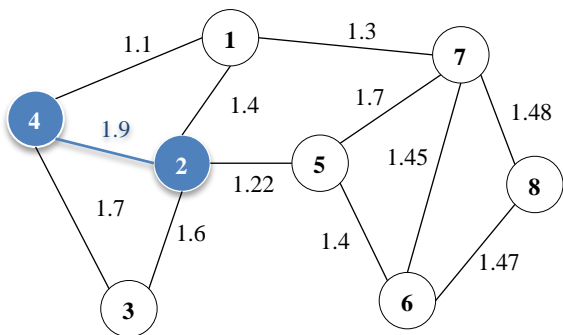
تا این مرحله گراف وزن‌داری وجود دارد که وزن هر یال میزان شباهت ساختاری و محتوایی بین دو گره آن را نشان می‌دهد؛ در ادامه و به منظور خوشه‌بندی ابتدا هر گره در گراف به عنوان یک خوشه در نظر گرفته می‌شود، سپس به ترتیب گره‌های اطراف یال‌های با وزن بالا با یکدیگر ادغام می‌شوند. با این کار این دو گره به عنوان یک مجموعه در نظر گرفته می‌شود. سپس این مجموعه به عنوان یک خوشه در نظر گرفته می‌شود و اگر گره دیگری به هر دو گره ادغام‌شده در ارتباط بوده، تنها با یک یال به این مجموعه متصل شده و وزن آن مطابق با رابطه (۷) برابر با میانگین حسابی وزن دو یال قبلی می‌شود.

$$W_{new} = \frac{1}{n} \sum_{i=1}^n W_i \quad (7)$$

که در آن، n تعداد یال‌های موجود در مجموعه (خوشه) بوده که به گره خارج از خوشه متصل هستند. به همین ترتیب به صورت سلسله مراتبی تا جایی که تعداد خوشه‌ها برابر با تعدادی باشد که مدنظر کاربر است. یعنی با توجه به شکل (۳) اگر فرض شود یال بین دو گره i و z دارای بیشترین وزن باشد، پس از ادغام آن‌ها گره z که قبلاً با وزن ۴ و ۲ به این گره‌ها وصل بود، با یک یال از طریق رابطه (۷) محاسبه و برابر با ۳ شده به روز می‌شود.

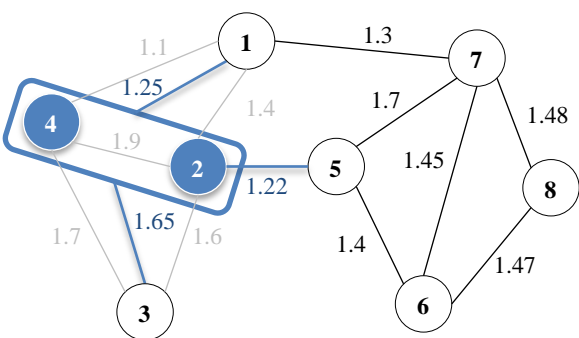
همان‌طور که مشاهده می‌شود، در گراف ساختاری یالی بین گره‌های ۱ و ۴ وجود نداشت که با اعمال اثر محتوا و ساختار کنار یکدیگر، به دلیل شباهت محتوایی بالای این دو گره، بین آن‌ها یالی رسم شده است. در مرحله بعدی بیش‌ترین شباهت بین دو گره (حساس‌ترین یال‌ها) نقش مهمی در ایجاد خوشه‌ها ایفا می‌کنند، با توجه به اینکه حین ساخت گراف ساختاری-محتوایی، این مقادیر به‌صورت نزولی ذخیره می‌شوند، دسترسی به آن‌ها ساده است.

مطابق با شکل (۷) بیش‌ترین مقدار در گراف فرض شده  $1/9$  بوده که مربوط به یال بین ۲ و ۴ است، در این مرحله این دو گره به‌عنوان یک خوشه در نظر گرفته می‌شوند.



شکل ۷. مرحله اول خوشه‌بندی روش ارائه‌شده

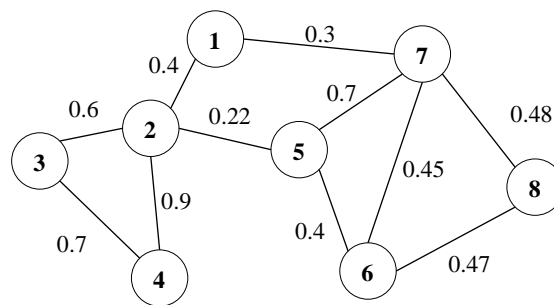
از آنجاکه دو گره ۲ و ۴ به‌عنوان یک خوشه در نظر گرفته‌شده‌اند، وزن یال گره‌های متصل به این مجموعه با توجه به رابطه (۷) به‌روز می‌شوند. نمایی از به‌روز شدن ارتباط در شکل (۸) نشان داده شده است.



شکل ۸. در نظر گرفتن گره‌های ۲ و ۴ به‌عنوان یک گره و ادغام گره‌های متصل به آن

به‌عنوان مثال با توجه به این‌که گره ۱ با دو گره ۲ و ۴ در ارتباط بوده وزن ارتباط این گره با خوشه موردنظر به‌صورت زیر محاسبه می‌شود.

$$w_{new(1C_1)} = \frac{1}{2}(1.1 + 1.4) = 1.25$$

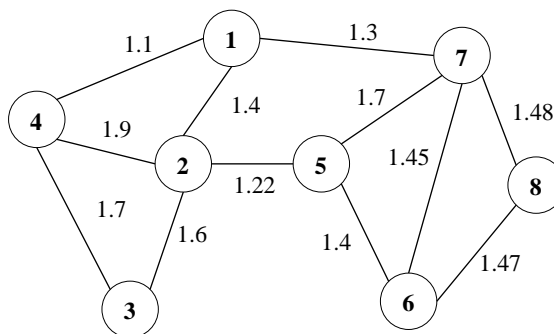


شکل ۵. گراف وزن‌دار نمونه

به‌منظور ساخت گراف جدید با توجه به راه‌حل ارائه‌شده برای هر دو گره موجود در گراف طبق رابطه (۲) وزن یال‌ها محاسبه می‌شود. به این منظور در این مثال میزان اهمیت شباهت ساختاری و محتوایی (یعنی مقدار  $a$  و  $b$  در این رابطه) هردو برابر با ۱ در نظر گرفته شده است. لازم به ذکر است که در این روش مقادیر  $a$  و  $b$  عددی بین ۰ تا ۱ باید در نظر گرفته شود تا از بروز پاسخ‌های غیرواقعی جلوگیری شود. ماتریس زیر نشان‌دهنده شباهت ساختاری-محتوایی دو گره است.

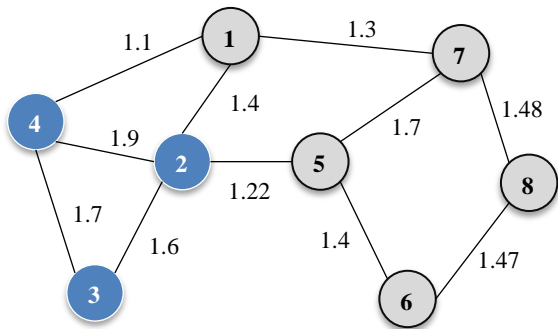
$$sim = \begin{bmatrix} - & 1.40 & 1.02 & 1.10 & 0.82 & 0.80 & 1.30 & 0.65 \\ - & - & 1.60 & 1.90 & 1.22 & 0.98 & 0.79 & 0.69 \\ - & - & - & 1.70 & 0.92 & 0.85 & 0.54 & 0.57 \\ - & - & - & - & 0.80 & 0.56 & 0.42 & 0.41 \\ - & - & - & - & - & 1.40 & 1.70 & 1.00 \\ - & - & - & - & - & - & 1.45 & 1.47 \\ - & - & - & - & - & - & - & 1.48 \\ - & - & - & - & - & - & - & - \end{bmatrix}$$

در ادامه باید میزان حد آستانه شباهت ساختاری-محتوایی دو گره تعیین شود، همان‌طور که گفته شد این شباهت به‌صورت پیش‌فرض برابر با میانگین وزن یال‌ها خواهد بود. میانگین شباهت دو گره به یکدیگر در این مثال برابر با  $1/0.5$  است که یال‌هایی که کمتر از این مقدار شباهت دارند رسم نمی‌شوند؛ در نتیجه گراف شکل (۶) نشان‌دهنده گراف ساختاری-محتوایی است که باید خوشه‌بندی شود.



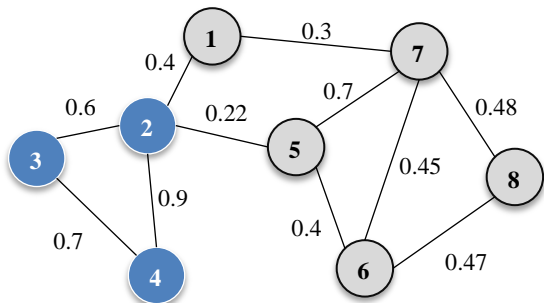
شکل ۶. گراف ساختاری محتوایی نمونه

تعداد خوشه‌ها به تعدادی مدنظر کاربر برسد. از آنجاکه گره ۱ بین دو خوشه قرار دارد، به خوشه‌ای تعلق خواهد داشت که با یال با وزن بیشتر به آن متصل است. خوشه‌بندی نهایی گراف ساختاری-محتوایی تشکیل شده از گراف مثال موردنظر را شکل (۱۱) نشان می‌دهد.



شکل ۱۱. خوشه‌های نهایی گراف ساختاری-محتوایی

در نهایت خوشه‌بندی ساختاری-محتوایی گراف اولیه به صورت شکل (۱۲) انجام شده است. به منظور ارزیابی روش ارائه شده معیارهای موجود بر روی خوشه‌های گراف اولیه محاسبه می‌شود تا بتوان روش ارائه شده را با سایر روش‌های خوشه‌بندی ساختاری-محتوایی مقایسه نمود.



شکل ۱۲. خوشه‌های نهایی در گراف وزن دار اولیه

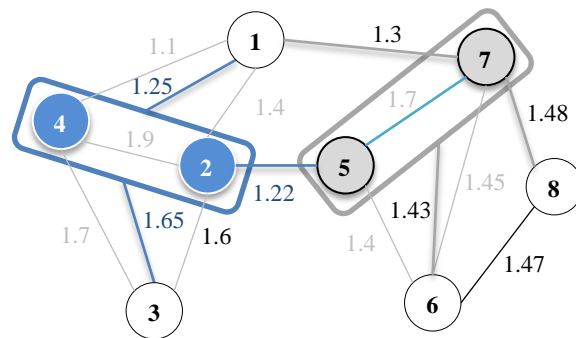
در شکل (۱۳) فلوجارت الگوریتم ارائه شده نشان داده شده است. ورودی الگوریتم ICS-Cluster شامل گراف ساختاری و ماتریس شباهت محتوایی یال‌های آن است. در ادامه گراف ساختاری-محتوایی بر اساس رابطه که در روش ارائه شده تولید می‌شود، که در آن وزن گره‌ها، بیانگر شباهت ساختاری-محتوایی بیش از حد آستانه مشخصی است.

در ادامه بیشترین وزن یال که نشان‌دهنده حداکثر شباهت ساختاری-محتوایی بین دو گره است را یافته و آن دو به عنوان یک خوشه در نظر گرفته می‌شوند. پس از آن وزن‌های متصل به گره‌های موجود در خوشه به‌روزرسانی می‌شوند، بدین معنا که

که در آن،  $C_1$  نشان‌دهنده خوشه اول (متشکل از گره‌های ۲ و ۴) است. این مقدار برای گره ۳ به صورت زیر محاسبه می‌شود.

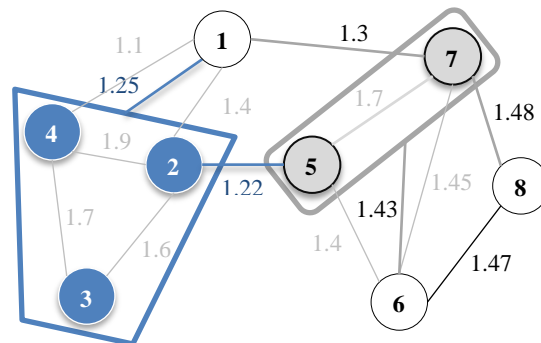
$$w_{new(3C_1)} = \frac{1}{2}(1.7 + 1.6) = 1.65$$

با توجه به این که گره ۵ تنها با گره ۲ از خوشه تشکیل شده در ارتباط است، مقدار موردنظر برای آن همان مقدار قبلی خواهد بود. در ادامه و پس از به‌روز شدن تمام ارتباطات با خوشه موردنظر، به مقدار بعدی بزرگ‌ترین  $w$  مراجعه می‌شود. این مقدار با توجه به شکل (۹)،  $1/7$  بوده که بین دو گره ۵ و ۷ وجود دارد. با توجه به اینکه این دو گره عضو هیچ خوشه‌ای نیستند، به عنوان خوشه جدید در نظر گرفته می‌شوند. برای ادامه کار مشابه قبل وزن یال مربوط به گره‌هایی که به این خوشه متصل هستند، با توجه به رابطه (۷) به‌روز می‌شوند. نمایی از آن در شکل (۹) مشاهده می‌شود.



شکل ۹. در نظر گرفتن گره‌های ۵ و ۷ به عنوان یک گره و ادغام گره‌های متصل به آن

در ادامه بیشترین مقدار،  $1/65$  بوده که مربوط به ارتباط بین گره ۳ با خوشه اول است. بنابراین، مطابق با شکل (۱۰) گره ۳ به خوشه یک اضافه می‌شود.



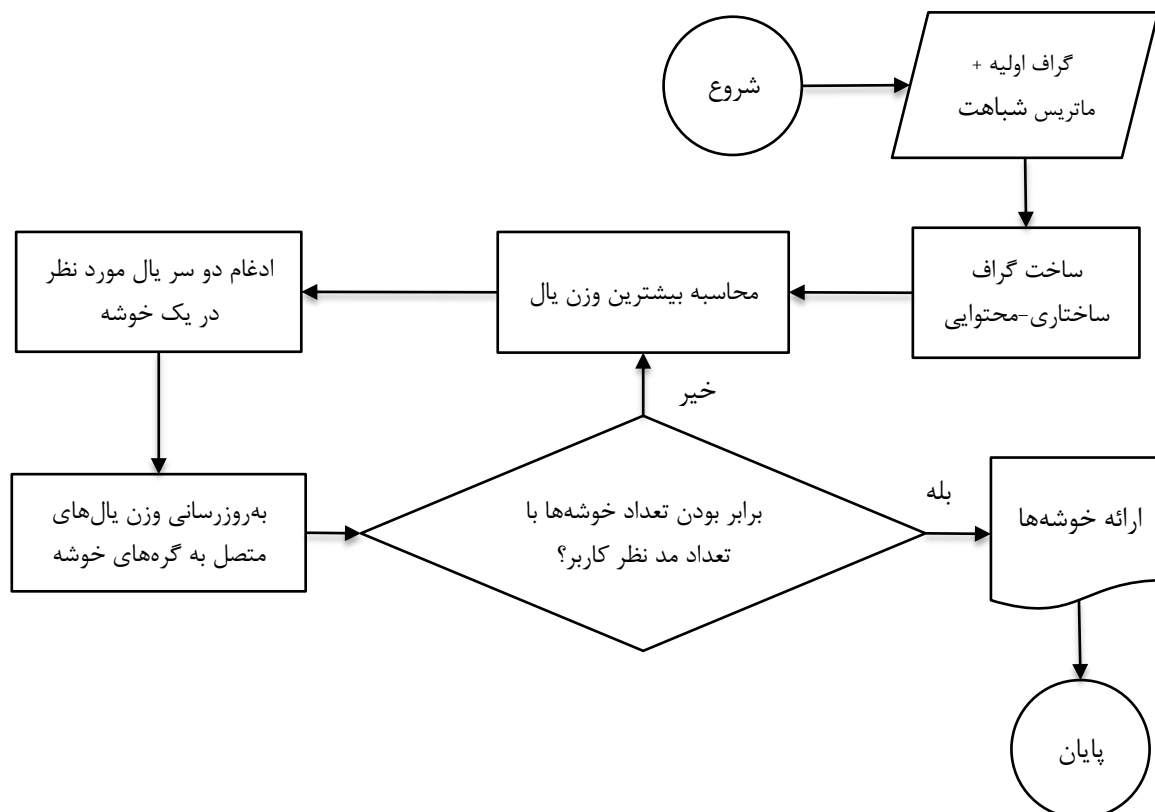
شکل ۱۰. اضافه شدن گره ۳۳ به خوشه اول (متشکل از گره‌های ۲ و ۴)

در ادامه نیز به همین ترتیب به صورت سلسله مراتبی یال‌های با وزن بالا محاسبه شده و به خوشه‌ها اضافه می‌شوند تا زمانی که

زمانی که تعداد خوشه‌ها به تعداد مدنظر کاربر برسد، ادامه میابد. در ادامه کیفیت خوشه‌های تولیدشده در روش ارائه شده (ICS-Cluster) با روش‌های خوشه‌بندی ساختاری-محتوایی با پیاده‌سازی این روش‌ها بر روی یک مجموعه داده، مورد ارزیابی قرار گرفته است.

تمام گره‌هایی که از یک یال به خوشه‌ای متصل هستند، به یک یال تبدیل شده و از وزن آن‌ها میانگین گرفته می‌شود و به‌عنوان وزن ارتباط آن گره با خوشه در نظر گرفته می‌شود.

پس از طی مراحل فوق، بررسی می‌شود که تعداد خوشه موجود به تعداد مدنظر کاربر رسیده یا خیر، در صورتی که تعداد برابر نبود، مجدداً حداکثر وزن یال محاسبه شده و خوشه‌بندی تا



شکل ۱۳. فلوجارت روش ICS-Cluster

روش‌های مورد مقایسه دارای پارامترهای متفاوتی هستند. برای مقایسه این روش‌ها محدودیت‌هایی وجود دارد. به‌عنوان مثال در روش SA-Cluster و روش ICS-Cluster تعداد خوشه‌ها از ابتدا مشخص می‌شود، ولی در روش CS-Cluster مشخص نیست و تنها دو پارامتر حد آستانه شباهت و فاصله مراکز خوشه‌ها در اختیار است؛ همچنین در روش SANS نیز فقط پارامتر حد آستانه شباهت در دسترس است.

#### ۴-۱. مجموعه داده<sup>۲</sup>

در آزمایش‌های انجام شده از مجموعه داده Delicious استفاده شده است.<sup>۳</sup> این مجموعه داده دارای ۱۸۶۱ گره و ۷۶۶۴ یال است. خصوصیت‌های هر گره (که تعداد آن‌ها ۱۳۵۰ عدد است)

#### ۴. ارزیابی

جهت ارزیابی روش ارائه شده و مقایسه آن با سایر روش‌های خوشه‌بندی ساختاری-محتوایی آزمایش‌های مختلفی بر اساس معیارهای مختلف انجام شده است. این آزمایش‌ها در یک سیستم کامپیوتری با پردازنده ۴ هسته‌ای GH۲،۲۰ و حافظه اصلی ۸ گیگابایتی انجام شده است. به‌منظور پیاده‌سازی الگوریتم‌ها از زبان جاوا در محیط اکتلیپس<sup>۱</sup> استفاده شده است.

در ادامه، نتایج مقایسه آزمایش‌های انجام شده بر روی چهار روش SANS، SA-Cluster، CS-Cluster و RLS-Cluster آمده است. از آنجاکه خوشه‌بندی DCM غیر همپوشان بوده و تمامی گره‌ها را پوشش نمی‌دهد، در مقایسه در نظر گرفته نشده است.

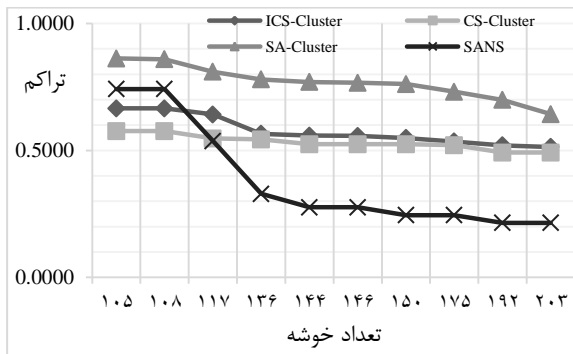
<sup>۲</sup> Dataset

<sup>۳</sup> <https://rdrr.io/cran/mlr.datasets/man/delicious.html>

<sup>۱</sup> Eclipse



گره‌های مشابه در هر نقطه‌ای از گراف به خوشه اضافه می‌شوند و این باعث می‌شود در نهایت تراکم کلی خوشه‌ها پایین باشد.



شکل ۱۴. مقایسه تراکم روش ICS-Cluster با سایر روش‌ها

#### ۲-۲-۴. خطای یال

با توجه به اینکه از جنبه ساختاری، خوشه‌های ایدئال است که تعداد یال‌های داخلی زیاد و تعداد یال‌های خروجی (بین خوشه‌های) کم دارد و معیار تراکم تنها تعداد یال‌های داخلی را در نظر می‌گیرد، بنابراین، حتی از لحاظ ساختاری معیار کاملی نبوده و همه جنبه‌ها را در نظر نمی‌گیرد. معیار خطای یال که در ادامه شرح داده می‌شود از لحاظ ساختاری کامل‌تر است، چون هم یال‌های داخلی و هم یال‌های خروجی را در نظر می‌گیرد. این معیار مجموع خطای یال خوشه‌های حاصل از عمل خوشه‌بندی را محاسبه می‌کند [۳۱]. این معیار نیز یک معیار ساختاری است و خوشه‌بندی را از نظر ساختاری مورد ارزیابی قرار می‌دهد. خطای یال شامل تعداد یال‌هایی است که به اشتباه حذف شده‌اند، به تعبیری دیگر تعداد یال‌هایی که وجود دارند در صورتی که بهتر بود حذف شوند. برای ارزیابی این معیار دو نوع خطا تعریف می‌شود، اولین نوع خطا که از طریق رابطه (۹) محاسبه می‌شود، تعداد یال‌های ناموجود بین گره‌ها در یک جامعه  $C$  را تعیین می‌کند [۳۱].

$$\epsilon_{\text{within}}(C, E) = \{(v, w) | v, w \in C \wedge v \neq w \wedge (v, w) \notin E\} \quad (9)$$

خطای نوع دوم تعداد یال‌های بین جامعه  $C$  و جوامع دیگر را بر اساس رابطه (۱۰) تعیین می‌کند.

$$\epsilon_{\text{between}}(C, E) = \{|(v, w) \in E | v \in C \wedge w \notin C|\} \quad (10)$$

این دو خطا روی هم رفته یک معیار خطای کلی را ارائه می‌دهند که به صورت رابطه (۱۱) تعریف می‌شود:

$$\epsilon(C, G) = \sum_{C \in \mathcal{C}} \epsilon_{\text{within}}(C, E) + \frac{\epsilon_{\text{between}}(C, E)}{2} \quad (11)$$

نیز به وسیله یک آرایه دودویی مشخص شده است. خصوصیات با استفاده از مقادیر ۰ و ۱ در آرایه مربوط به هر گره مشخص شده که برای تعیین وجود یا عدم وجود خصوصیت در یک گره استفاده می‌شوند. در واقع برای هر گره یک آرایه ۱۳۵۰ خانه‌ای وجود دارد که برای هر ویژگی یک خانه در نظر گرفته شده است. اگر گره دارای ویژگی مورد نظر بود، خانه مربوطه برای آن ویژگی مقدار ۱ خواهد داشت، در غیر این صورت مقدار خانه برابر با ۰ خواهد بود. این کار باعث می‌شود که خصوصیات گره‌ها به صورت دودویی ذخیره شده تا هنگام محاسبه شباهت، بتوان از شباهت جاکارد به راحتی استفاده نمود. لازم به ذکر است که در روش ICS-Cluster پس از تبدیل گراف اولیه به گراف ساختاری-محتوایی تعداد یال‌های گراف جدید از ۷۶۶۴ به ۸۷۳۹۹۴ یال می‌رسد که دلیل آن تأثیر شباهت محتوایی در کنار شباهت ساختاری است.

#### ۲-۴. معیارهای بررسی شده

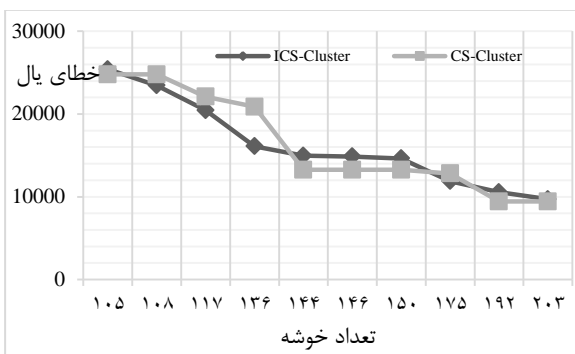
مقایسه‌های انجام شده بر اساس معیارهای مختلفی صورت پذیرفته که در ادامه تشریح می‌شوند.

##### ۱-۲-۴. تراکم

این معیار که خوشه‌بندی را از نظر ساختاری ارزیابی می‌کند، از طریق رابطه (۸) محاسبه می‌شود [۱۸]:

$$Density = \frac{\sum_{c \in C} E_c}{E} \quad (8)$$

که در آن،  $C$  مجموعه خوشه‌های حاصل از خوشه‌بندی،  $E_c$  یال‌های موجود در خوشه  $C$  و  $E$  مجموعه کل یال‌های گراف را نشان می‌دهد. شکل (۱۴) مقایسه تراکم خوشه‌های حاصل از خوشه‌بندی چهار روش را نشان می‌دهد. در روش SA-Cluster مراکز خوشه‌ها در نقاط متراکم گراف انتخاب می‌شود و سپس گره‌های همسایه این گره‌های مرکزی به هر خوشه تخصیص داده می‌شود. این عمل باعث می‌شود در نهایت خوشه‌ها تراکم بالایی داشته باشند. در روش ICS-Cluster نیز با توجه به اینکه ابتدا گراف اولیه به گرافی با تعداد یال بیشتر که بیانگر شباهت محتوایی و ساختاری است تبدیل شده و همچنین خوشه‌ها به صورت سلسله مراتبی از پایین به بالا تشکیل شده و در نهایت ارزیابی بر روی خوشه نهایی بر روی گراف اولیه انجام می‌شود، این معیار ساختاری برای این روش دارای مقدار مناسبی است. روش SANS تراکم کمی را به دست می‌آورد. علت این تراکم پایین در این روش این است که بعد از تعیین مراکز خوشه‌ها،



شکل ۱۶. مقایسه خطای یال روش‌های ICS-Cluster و CS-Cluster

### ۴-۲-۳. میانگین شباهت<sup>۱</sup>

این معیار خوشه‌بندی را از نظر محتوایی مورد ارزیابی قرار می‌دهد که میانگین شباهت خوشه‌های گراف را مطابق با رابطه (۱۲) محاسبه می‌کند.

$$Avg\_Sim\_Graph = \frac{\sum_{C \in C} C_{avgsim}}{|C|} \quad (12)$$

که در آن C مجموعه خوشه‌های حاصل از خوشه‌بندی،  $C_{avgsim}$  میانگین شباهت درون خوشه‌های خوشه C و |C| تعداد خوشه‌های حاصل از خوشه‌بندی را نشان می‌دهد. در واقع میانگین شباهت، مجموع میانگین شباهت داخلی خوشه‌ها به روی تعداد خوشه‌ها است. همان‌طور که در روش ICS-Cluster توضیح داده شد، در مرحله اولیه این الگوریتم، ابتدا گراف به یک گراف ساختاری-محتوایی تبدیل می‌شود که در آن نقش شباهت محتوایی دیده شده است. در روش CS-Cluster در هنگام تعیین مراکز خوشه‌ها و همچنین در زمان توسعه خوشه‌ها میزان شباهت گره‌ها در نظر گرفته می‌شود. در نتیجه همان‌طور که در نمودار شکل (۱۷) مشاهده می‌شود این دو روش بیشترین میانگین شباهت را دارند.



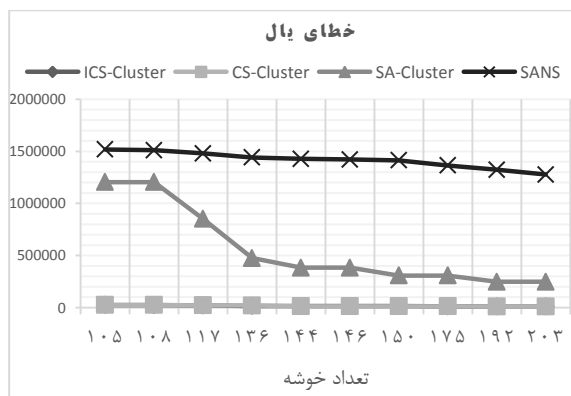
شکل ۱۷. مقایسه میانگین شباهت روش ICS-Cluster با سایر روش‌ها

### ۴-۲-۴. معیار ارزیابی CS-Measure

در اکثر روش‌های خوشه‌بندی، عمل خوشه‌بندی بر اساس ساختار گراف انجام می‌شود و معیارهای ارزیابی جوامع مورد استفاده نیز اکثراً بر اساس ویژگی‌های ساختاری گراف هستند، از جمله این

افراز پایه گراف حالتی است که در آن هر گره به صورت جداگانه به عنوان یک خوشه در نظر گرفته می‌شود. در روش ICS-Cluster، ابتدا یال با وزن بالا (مهم) انتخاب شده، سپس دوسر آن به عنوان یک خوشه در نظر گرفته می‌شود، در توسعه آن نیز گره‌هایی که با وزن بالا به خوشه متصل هستند نیز به آن اضافه می‌شوند، بنابراین، در این روش، خوشه‌ها دارای تعداد اتصالات داخلی بالا و تعداد اتصال خارجی کم هستند. در روش CS-Cluster نیز در هنگام انتخاب مراکز خوشه‌ها و هنگام توسعه هر خوشه وضعیت اتصال گره‌ها در نظر گرفته می‌شود. در نتیجه در نهایت خوشه‌ها در این دو حالت دارای کمترین خطای یال هستند. همان‌طور که در شکل (۱۵) مشاهده می‌شود، بعد از آن‌ها روش SANS، به نسبت خطای یال مناسبی دارد، ولی SA-Cluster بیشترین خطای یال را دارد. در روش SA-Cluster چون بعد از انتخاب مراکز خوشه‌ها، گره‌های مشابه مرکز خوشه در هر قسمت از گراف به خوشه مربوط به آن مرکز خوشه تخصیص داده می‌شوند و این گره‌ها گاهی در فاصله خیلی دور و دارای اتصال کمی است، در نتیجه در نهایت میزان خطای یال بسیار زیادی خواهد داشت.

در معیار خطای یال وزن یال‌ها در نظر گرفته نشده است؛ یعنی تأثیر وجود یک یال خارجی با وزن کم برابر تأثیر یک یال خارجی با وزن بسیار بالاست. در صورتی که تأثیر وجود این دو یال باید متفاوت باشد.



شکل ۱۵. مقایسه خطای یال روش ICS-Cluster با سایر روش‌ها

در شکل (۱۵)، به دلیل بودن خطای یال روش‌های ICS-Cluster و CS-Cluster، خطای یال روش‌های SA-Cluster و SANS Cluster به خوبی نشان داده نشده است، به همین دلیل در شکل (۱۶) تنها خطای یال این دو روش آمده است تا اختلاف آن‌ها به خوبی نشان داده شود.

<sup>1</sup> Average Similarity

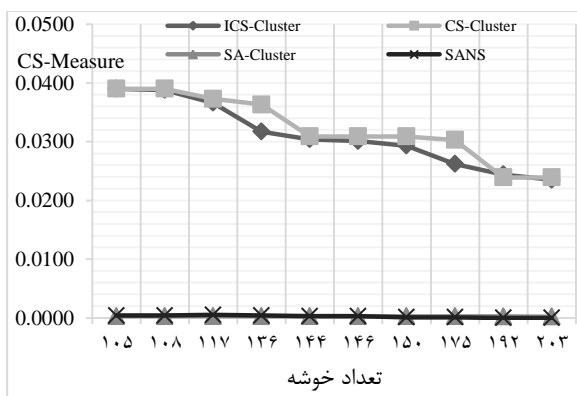
می‌شود. در ادامه میانگین وزن یال‌ها (شامل یال‌های موجود در خوشه، یال‌های ناموجود و یال‌های خروجی خوشه) محاسبه شده و به‌عنوان امتیاز خوشه در نظر گرفته می‌شود. پس از محاسبه امتیاز هر خوشه، میانگین امتیاز خوشه‌ها به‌عنوان امتیاز کلی خوشه‌بندی مطابق با رابطه‌های (۱۳) و (۱۴) محاسبه می‌شود [۳۲].

$$EScore(C) = \frac{\sum_{e \in E_{in,C}} We + \frac{\sum_{e \in E_{ex,C}} (-we)}{2}}{|E_{in,C}| + |E_{ex,C}|} \quad (13)$$

$$CS - Measure = \frac{\sum_{C \in C} EScore(C)}{|C|} \quad (14)$$

که در آن  $EScore(C)$  امتیاز جامعه  $C$ ،  $E_{in,C}$  مجموعه یال‌های داخلی جامعه  $C$ ،  $E_{ex,C}$  مجموعه یال‌های خروجی از جامعه  $C$ ،  $We$  وزن یال  $e$ ،  $|C|$  تعداد خوشه‌های حاصل از خوشه‌بندی موردنظر را نشان می‌دهد.

همان‌طور که توضیح داده شد در این معیار، وزن یال‌های خارجی در (۱-) ضرب شده و در محاسبه میانگین شباهت نظر گرفته می‌شود. در نتیجه تأثیر وجود یک یال خارجی با وزن بالا در کم کردن میانگین شباهت خوشه، خیلی بیشتر از وجود یک یال خارجی با وزن کم است. در نتیجه نسبت به معیار خطای یال دقیق‌تر عمل می‌کند. با توجه به اینکه این معیار، شباهت ساختاری و محتوایی را به‌طور هم‌زمان در نظر می‌گیرد، معیار مطلوب‌تری برای ارزیابی کیفیت خوشه‌های نهایی است؛ بنابراین، با توجه به شکل (۱۸)، می‌توان گفت روش CS-Measure نسبت به سایر روش‌ها در مجموع بهتر عمل می‌کند. با توجه به اینکه در مرحله اول الگوریتم روش ICS-Cluster گراف ایجاد شده گراف ساختاری-محتوایی بوده و در ایجاد و توسعه خوشه‌ها هم‌زمان هر دو جنبه در نظر گرفته می‌شود، این روش نیز دارای مقادیر مناسبی بر اساس این معیار است، در حالی که روش SA-Cluster و SANS به لحاظ این معیار دارای مقدار کمی هستند.



شکل ۱۸. مقایسه CS-Measure روش ICS-Cluster با سایر روش‌ها

معیارهای ساختاری معیار تراکم، پیمانی، میانگین درجه، نسبت برش، هدایت، خطای یال را می‌توان نام برد [۱۳]. برخی معیارهای کارایی نیز وجود دارد که محتوای گره‌ها را مدنظر دارند از جمله فراخوانی، دقت و صحت که البته این دو معیار در مورد خوشه‌بندی‌هایی که برای هر خوشه یک توصیف مشخص و خاص مدنظر است کاربرد دارد [۳۵]. معیار CS-Measure هر دو جنبه ساختاری و محتوایی را با یکدیگر ترکیب می‌کند.

همان‌طور که گفته شد، خطای یال خوشه به‌صورت رابطه‌های (۹-۱۱) تعریف می‌شود. چون خطای یال‌های خارجی برای هر دو خوشه‌ای که شامل گره‌های مبدأ و مقصد یال هستند محاسبه می‌شود، خطای یال خارجی تقسیم بر ۲ می‌شود [۳۴]. معیار فوق فقط ویژگی‌های ساختاری گراف را در محاسبات خود در نظر می‌گیرد. به عبارتی هنگام محاسبه خطای یال خارجی به این نکته که یال موردنظر بین دو گره با شباهت محتوایی بالا وجود دارد یا بین دو گره با شباهت محتوایی کم، توجهی نمی‌شود و تأثیر وجود این دو یال یکسان است. به‌طوری‌که مثلاً اگر بین دو خوشه یک یال با وزن وجود داشته باشد، امتیاز خوشه‌بندی باید کمتر از زمانی باشد که بین دو خوشه یک یال با وزن خیلی کم وجود دارد. در واقع حالت دوم وضعیتی بهتر را نشان می‌دهد. در صورتی که در معیار خطای یال ذکر شده این مورد رعایت نشده است و تأثیر وجود یا عدم وجود دو یال با وزن‌های مختلف بین دو خوشه یکسان است؛ بنابراین، با در نظر گرفتن مقدار شباهت گره‌ها و ترکیب این مقادیر شباهت با معیار خطای یال می‌توان به معیاری مناسب جهت ارزیابی خوشه‌بندی‌های ساختاری-محتوایی دست‌یافت.

در معیار CS-Measure امتیاز خوشه بر اساس شباهت درون خوشه‌ای و خطای یال خوشه محاسبه می‌شود. برای این کار در این معیار وزن یال‌های درون خوشه‌ای و بین خوشه‌ای، متناسب با وضعیتی که دارند به‌روز می‌شود. سپس، میانگین وزن یال‌های هر خوشه محاسبه می‌شود. به همین منظور، چون حالت ایدئال به این صورت است که خوشه‌ها یال‌های خروجی نداشته باشند و یا در صورت وجود این یال‌ها وزن کمی داشته باشند، وزن یال‌های خروجی به‌صورت منفی در نظر گرفته شده (w-) که این عمل باعث می‌شود در صورت وجود یال خروجی در خوشه، امتیاز کلی خوشه‌بندی متناسب با وزنی که یال خروجی دارد کاهش یابد.

اگر وزن یال‌های خروجی زیاد باشد این کاهش امتیاز زیاد خواهد بود و اگر وزن یال‌های خروجی کم باشد این کاهش امتیاز کم‌تر خواهد بود. وزن یال‌های ناموجود برابر صفر در نظر گرفته

۴-۲-۵. پیمانگی

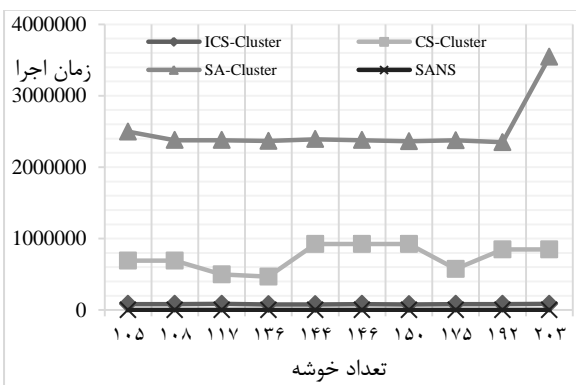
این معیار که در واقع خوشه‌بندی را از نظر ساختاری ارزیابی می‌کند به صورت رابطه (۱۵) محاسبه می‌شود [۳۶]:

$$Q = \frac{1}{4m} \sum_S (m_S - E(m_S)) \quad (15)$$

که در آن،  $m$  تعداد یال‌های گراف،  $m_S$  تعداد یال‌های موجود در خوشه  $S$ ،  $E(m_S)$  تعداد یال‌های مورد انتظار در خوشه  $S$  را نشان می‌دهد این معیار از طریق مجموع تفاضل یال‌های موجود در خوشه‌ها به یال‌های ممکن آن به روی مجموع یال‌های موجود در گراف محاسبه می‌شود.

همان‌طور که در شکل (۱۹) مشاهده می‌شود، با توجه به اینکه در روش SA-Cluster مراکز خوشه‌ها در نقاط متراکم گراف انتخاب می‌شود و سپس گره‌های همسایه این گره‌های مرکزی به هر خوشه تخصیص داده می‌شود. این عمل باعث می‌شود در نهایت خوشه‌ها پیمانگی بالایی داشته باشند.

روش ICS-Cluster نسبت به روش CS-Cluster دارای پیمانگی بالاتری بوده و روش SANS نیز کمترین پیمانگی را دارد. علت پیمانگی پایین در روش SANS این است که بعد از تعیین مراکز خوشه‌ها، گره‌های مشابه در هر نقطه‌ای از گراف به خوشه اضافه می‌شوند و این باعث می‌شود در نهایت پیمانگی کلی خوشه‌ها پایین باشد.

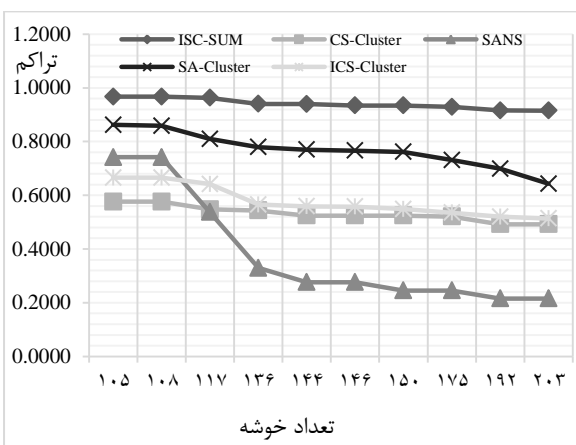


شکل ۲۰. مقایسه زمان اجرای روش ICS-Cluster با سایر روش‌ها (میلی‌ثانیه)

با توجه به تحقیقات صورت گرفته، با تغییر رابطه (۷) در بخش چهارم، برای برخی از معیارهای محاسبه شده بهبودهایی حاصل شد. در رابطه (۷) از میانگین حسابی برای ادغام گره‌ها استفاده شده است، در حالی که اگر به جای آن از جمع یال‌های ادغام شونده استفاده شود، نتایج پیمانگی و تراکم نسبت به محاسبه میانگین حسابی بهبود خواهند یافت. بدین معنا که مثلاً اگر دو رأس قرار است با یکدیگر ادغام شوند، وزن یال جدید مطابق با رابطه (۱۶) حاصل جمع وزن دو یال قبلی است.

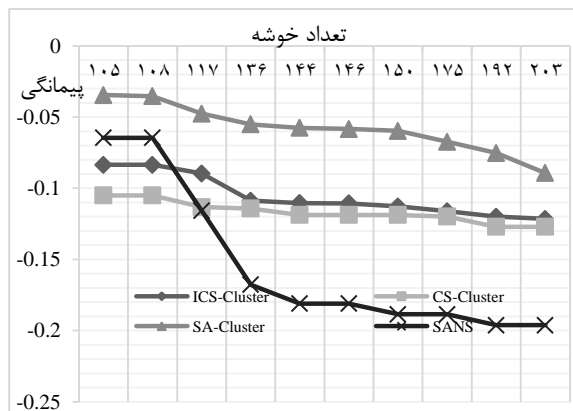
$$W_{new} = \sum_{i=1}^n W_i \quad (16)$$

که در آن،  $W_{new}$  وزن یال ادغام شده است،  $n$  تعداد یال‌هایی است که ادغام شده و  $W_i$  وزن هر کدام از آن‌ها است. همان‌طور که در شکل (۲۱) مشاهده می‌شود، این روش با نام ICS-SUM نشان داده شده و دارای تراکم بهتری نسبت به سایر روش‌ها است.



شکل ۲۱. مقایسه تراکم روش‌های خوشه‌بندی بر اساس رابطه (۱۶) (ICS-SUM)

مقایسه خاصیت پیمانگی برای الگوریتم‌ها در شکل (۲۲) نشان داده شده است که در آن روش ICS-SUM نسبت به سایرین دارای خاصیت پیمانگی بالاتری است. در واقع نتایج نشان می‌دهد که اگر در خوشه‌بندی نیاز به بالا بودن دو معیار تراکم و



شکل ۱۸. مقایسه پیمانگی روش ICS-Cluster با سایر روش‌ها

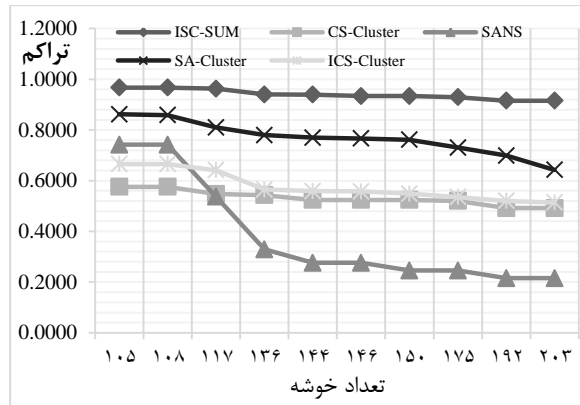
۴-۲-۶. زمان اجرا

با استفاده از این معیار مدت زمان خوشه‌بندی محاسبه می‌شود. با توجه به نمودار حاصل از آزمایش‌های انجام شده که در شکل (۲۰) مشاهده می‌شود؛ از بین روش‌های مقایسه شده، روش‌های ICS-Cluster و SANS نسبت به سایر روش‌ها سریع‌تر هستند. پس از آن‌ها، روش CS-Cluster دارای زمان اجرای قابل قبولی است. روش SA-Cluster به دلیل تکرار شونده‌گی، نسبت به بقیه کندتر عمل خوشه‌بندی را انجام می‌دهد.

بیشتر روش‌های خوشه‌بندی که تاکنون ارائه شده است فقط یک جنبه ساختاری یا محتوایی را در نظر گرفته‌اند و روش‌های خوشه‌بندی که ساختار و محتوا را باهم نظر بگیرند کمتر ارائه شده است. در این مقاله روش ICS-Cluster ارائه شد که خوشه‌بندی را بر اساس ساختار و محتوا انجام می‌دهد. در این روش ابتدا گراف اولیه به یک گراف ساختاری-محتوایی تبدیل شده، سپس عمل خوشه‌بندی بر اساس یال‌های مهم آن گراف انجام می‌شود. کارایی روش ارائه شده در مقاله با سایر روش‌های خوشه‌بندی ساختاری-محتوایی ارائه شده مورد ارزیابی قرار گرفت که نتایج کلی آن در جدول (۱) آمده است. همان‌طور که گفته شد خوشه‌بندی که بین ساختار و محتوا تعادل خوبی برقرار کند، به لحاظ خوشه‌بندی ساختاری-محتوایی عملکرد خوبی از خود نشان داده است. لذا روش خوشه‌بندی که بتواند به لحاظ تمامی معیارها (ساختاری، محتوایی و ساختاری-محتوایی) نتایج مناسبی ارائه دهد، می‌توان گفت که بهترین روش است.

پیمانی بود، بهتر است به جای میانگین حسابی از رابطه (۱۶) استفاده کرد.

کاربرد زیادی در عصر ارتباط دارد که با افزایش استفاده از اینترنت حجم گراف‌های تشکیل شده روبه افزایش است، به همین دلیل خوشه‌بندی این گراف‌ها از اهمیت بالایی برخوردار است.



شکل ۲۲. مقایسه پیمانی روش‌های خوشه‌بندی بر اساس رابطه (۱۶) (ICS-SUM)

جدول ۱. مقایسه کلی روش ICS-Cluster با سایر روش‌ها بر اساس معیارهای مختلف ارائه شده در این مقاله

رتبه	CS-Measure	میانگین شباهت	خطای یال	تراکم	زمان اجرا	پیمانی
۱	CS-Cluster	CS-Cluster	ICS-Cluster	SA-Cluster	SANS	SA-Cluster
۲	ICS-Cluster	ICS-Cluster	CS-Cluster	ICS-Cluster	ICS-Cluster	ICS-Cluster
۳	SA-Cluster	SANS	SANS	CS-Cluster	CS-Cluster	CS-Cluster
۴	SANS	SA-Cluster	SA-Cluster	SANS	SA-Cluster	SANS

داده است. لذا روش خوشه‌بندی که بتواند به لحاظ تمامی معیارها (ساختاری، محتوایی و ساختاری-محتوایی) نتایج مناسبی ارائه دهد، می‌توان گفت که بهترین روش است.

### ۵. نتیجه‌گیری

نتایج ارائه شده در مقایسه‌های انجام شده نشان می‌دهد که الگوریتم ICS-Cluster به لحاظ ساختاری و محتوایی دارای نتایج متعادلی است، بدین معنا که برخلاف روش‌های SANS و SA-Cluster که در آن‌ها در برخی از معیارها جواب مناسبی ارائه شده و در برخی دیگر جواب نامناسبی ارائه می‌شود، در روش ICS-Cluster پاسخ‌ها به نسبت خوب و متعادل است. روش SA-Cluster بر اساس معیارهای ساختاری عملکرد خوبی دارد برخلاف روش CS-Cluster که به لحاظ محتوایی عملکرد مناسبی از خود نشان می‌دهد. در کل تنها مزیت روش SANS سرعت اجرایی آن است. روش ICS-Cluster به لحاظ معیارهای ساختاری

گراف کاربرد زیادی در عصر ارتباط دارد که با افزایش استفاده از اینترنت حجم گراف‌های تشکیل شده روبه افزایش است، به همین دلیل خوشه‌بندی این گراف‌ها از اهمیت بالایی برخوردار است. بیشتر روش‌های خوشه‌بندی که تاکنون ارائه شده است فقط یک جنبه ساختاری یا محتوایی را در نظر گرفته‌اند و روش‌های خوشه‌بندی که ساختار و محتوا را باهم نظر بگیرند کمتر ارائه شده است. در این مقاله روش ICS-Cluster ارائه شد که خوشه‌بندی را بر اساس ساختار و محتوا انجام می‌دهد. در این روش ابتدا گراف اولیه به یک گراف ساختاری-محتوایی تبدیل شده، سپس عمل خوشه‌بندی بر اساس یال‌های مهم آن گراف انجام می‌شود. کارایی روش ارائه شده در مقاله با سایر روش‌های خوشه‌بندی ساختاری-محتوایی ارائه شده مورد ارزیابی قرار گرفت که نتایج کلی آن در جدول (۱) آمده است. همان‌طور که گفته شد خوشه‌بندی که بین ساختار و محتوا تعادل خوبی برقرار کند، به لحاظ خوشه‌بندی ساختاری-محتوایی عملکرد خوبی از خود نشان

- [9] Opsahl, T.; Panzarasa, P. "Clustering in Weighted Networks"; Soc. Networks 2009, 31, 155-163.
- [10] Aggarwal, C.; Wang, H. "Managing and Mining Graph Data"; Springer US, 2010.
- [11] Patkar, S. B.; Narayanan, H. "An Efficient Practical Heuristic for Good Ratio-Cut Partitioning"; 16<sup>th</sup> Int. Conf. VLSI Design 2003, 1-6.
- [12] Feldmann, A. E.; Foschini, L. "Balanced Partitions of Trees and Applications"; Algorithmica 2015, 71, 354-376.
- [13] Newman, M. "Community Detection in Networks: Modularity Optimization and Maximum Likelihood Are Equivalent"; Phys. Rev. 2016, 94, 1-8.
- [14] Bhatia, V.; Rani, R. "A Parallel Fuzzy Clustering Algorithm for Large Graphs Using PREGEL"; Expert Syst. Appl. 2017, 78, 135-144.
- [15] Fortunato, S.; Hricb, D. "Community Detection in Networks: A User Guide"; Phys. Rep. 2016, 659, 1-44.
- [16] Khatoun, M.; Banu, A. "A Survey on Community Detection Methods In Social Networks"; IJEME. 2015, 1, 8-18.
- [17] Elhadi, H.; Agam, G. "Structure and Attributes Community Detection: Comparative Analysis of Composite Ensemble and Selection Methods"; Proc. 7<sup>th</sup> Workshop on Social Network Mining and Analysis 2013, 1-7.
- [18] Harenberg, S.; Bello, G.; Gjeltema, L.; Ranshous, S.; Harlalka, J.; Seay, R.; Padmanabhan, K.; Samatova, N. "Community Detection in Large-Scale Networks: A Survey and Empirical Evaluation"; Computation Stat. 2014, 6, 426-439.
- [19] Bu, Z.; Gao, G.; Li, H. J.; Cao, J. "CAMAS: A Cluster-Aware Multiagent System for Attributed Graph Clustering"; Inform. Fusion 2017, 37, 10-21.
- [20] Weber L. M.; Robinson M. D. "Comparison of Clustering Methods for High-Dimensional Single-Cell Flow and Mass Cytometry Data"; Cold Spring Harbor Labs Journals 2016, 047613.
- [21] Shchukin, V.; Khristich, D.; Galinskaya, I. "Word Clustering Approach to Bilingual Document Alignment"; First Conf. Machine Translation 2016, 2, 953-994.
- [22] Skabar, A. "Clustering Mixed-Attribute Data Using Random Walk"; Procedia Comput. Sci. 2017, 108, 988-997.
- [23] Xu, Z.; Cheng, J.; Xiao, X.; Fujimaki, R.; Muraoka, Y. "Efficient Nonparametric and Asymptotic Bayesian Model Selection Methods for Attributed Graph Clustering"; Knowl. Inform. Syst. 2017, 53, 239-268.
- [24] Gross, G.; Nagi, R.; Sambhoos, K. "A Fuzzy Graph Matching Approach in Intelligence Analysis and Maintenance of Continuous Situational Awareness"; Inform. Fusion 2014, 18, 43-61.
- [25] Boobalan, M. P.; Lopez, D.; Gao, X. Z. "Graph Clustering Using K-Neighbourhood Attribute Structural Similarity"; Appl. Soft Comput. 2016, 47, 216-223.
- [26] Zhou, H.; Li, J.; Li, J.; Zhang, F.; Cui, Y. "A Graph Clustering Method for Community Detection in Complex Networks"; Physica A 2017, 469, 551-562.
- [27] Bai, L.; Cheng, X.; Liang, J.; Guo, Y. "Fast Graph Clustering with a New Description Model for Community Detection"; Inform. Sci. 2017, 388, 37-47.

و محتوایی پاسخ مناسبی ارائه می‌دهد. این تحقیق ضمن ارائه روش خوشه‌بندی ساختاری-محتوایی جدید، منجر به ارائه نکاتی شد که حائز اهمیت هستند:

- تبدیل گراف اولیه به یک گراف ساختاری-محتوایی که وزن یال‌های آن بیانگر شباهت ساختار و محتوایی است.
- دادن وزن موردنظر کاربر به ساختار یا محتوا با توجه به اهمیت هر کدام نسبت به دیگری که منجر به افزایش کیفیت خوشه‌های نهایی مدنظر کاربر می‌شود.
- خوشه‌بندی بر اساس یال‌های حساس (با توجه به اینکه در سایر روش‌ها بر اساس گره حساس خوشه‌بندی انجام می‌شود).
- در تحقیقات معیارهای جدیدی به‌منظور ارزیابی کیفیت خوشه‌ها در گراف ارائه می‌شود [۳۷]، اما تنها یک معیار ساختاری-محتوایی برای ارزیابی خوشه‌ها وجود دارد [۳۲]، برای کارهای آینده پیشنهاد می‌شود معیارهای ارزیابی بیشتر و بهتری برای ارزیابی کیفیت و کارایی این روش‌های خوشه‌بندی پیشنهاد داد تا ارزیابی‌های دقیق‌تر و کامل‌تری بر روی روش‌های خوشه‌بندی انجام شود. در زمینه خوشه‌بندی ساختاری-محتوایی گراف می‌توان روش‌هایی را پیشنهاد داد که این عمل را با سرعت و دقت بالاتری انجام دهند.

## ۶. مراجع‌ها

- [1] Sambhoos, K.; Nagi, R.; Sudit, M.; Stotz, A. "Enhancements to High Level Data Fusion Using Graph Matching and State Space Search"; Inform. Fusion 2010, 11, 4, 351-364.
- [2] Zhou, L.; Wang, Q.; Fujita, H. "One Versus One Multi-Class Classification Fusion Using Optimizing Decision Directed Acyclic Graph for Predicting Listing Status of Companies"; Inform. Fusion 2017, 36, 80-89.
- [3] Gross, G. A.; Nagi, R. "Precedence Tree Guided Search for the Efficient Identification of Multiple Situations of Interest – AND/OR Graph Matching"; Inform. Fusion 2016, 27, 240-254.
- [4] Guedes, G. P.; Ogasawara, E.; Bezerra, E.; Xexeo, G. "Discovering Top-Non-Redundant Clustering in Attributed Graphs"; Neurocomputing 2016, 210, 45-54.
- [5] Schoch, D.; Valente, T. W.; Brandes, U. "Correlations Among Centrality Indices and a Class of Uniquely Ranked Graphs"; Soc. Networks 50, 46-54.
- [6] Johnson, M.; Paulusma, D.; Leeuwen, E. J. V. "Algorithms for Diversity and Clustering in Social Networks Through Dot Product Graphs"; Soc. Networks 2015, 41, 48-55.
- [7] Ertem, Z.; Veremyev, A.; Butenko, S. "Detecting Large Cohesive Subgroups with High Clustering Coefficients in Social Networks"; Soc. Networks 2016, 46, 1-10.
- [8] Vörös, A.; Snijders, T. A. B. "Cluster Analysis of Multiplex Networks: Defining Composite Network Measures"; Soc. Networks 2017, 49, 93-112.

- [33] Halloush, R. "Overhearing-Aware Modified Dijkstra's Algorithm for Multicasting over Multi-Hop Wireless Networks"; *Int. J. Commun. Netw. Distrib. Syst.* 2016, 16, 240-260.
- [34] Pool, S.; Bonchi, F.; Leeuwen, M. "Description-Driven Community Detection"; *ACM Trans. Intel. Syst. Technol.* 2014, 5, 28.
- [35] Wang, M.; Wang, Ch.; Xu, Y. J.; Zhang, J. "Community Detection in Social Networks: An In-Depth Benchmarking Study with a Procedure-Oriented Framework"; *Proceedings of the VLDB Endowment* 2015, 8, 998-1009.
- [36] Yang, J.; Leskovec, J. "Defining and Evaluating Network Communities Based on Ground-Truth"; *Knowl. Inf. Syst.* 2015, 42, 181-213.
- [37] Biswas, A.; Biswas, B. "Defining Quality Metrics for Graph Clustering Evaluation"; *Expert Syst. Appl.* 2017, 71, 1-17.
- [28] Ruan, Y.; Fuhry, D.; Parthasarathy, S. "Efficient Community Detection in Large Networks Using Content and Links"; *Proc. 22<sup>nd</sup> Int. Conf. World Wide Web WWW '13*. ACM, New York, NY, USA, 2013, 1089-1098.
- [29] Zhou, Y.; Cheng, H.; Xu, Y. J. "Graph Clustering Based on Structural/Attribute Similarities"; *Proc. VLDB Endowment* 2009, 2, 718-729.
- [30] Parimala, M.; Daphne, L. "Graph Clustering Based on Structural Attribute Neighborhood Similarity (SANS)"; *IEEE Int. Conf. Elec. Comput. Commun. Technol.* 2015, 1-5.
- [31] Pool, S.; Bonchi, F.; Leeuwen, M. "Description-Driven Community Detection"; *ACM Trans. Intel. Syst. Technol.* 2014, 03, 201-210.
- [32] Rahmati, K.; Naderi, H.; Keshvari, S. "Content-Structural Graph Clustering and a New Measure For Its Evaluation"; *Adv. Defence Sci. Technol.* 2018, 03, 201-210. (In Persian)