



## روشی نوین برای خوشه‌بندی نیمه نظارتی شبکه‌های پیچیده مبتنی بر معیار پیمانگی

محمد قدیریان<sup>۱</sup>، نوشین بیگدلی<sup>۲\*</sup>

<sup>۱</sup> گروه مهندسی برق کنترل- دانشکده فنی و مهندسی- دانشگاه بین‌المللی امام خمینی (ره)- قزوین- ایران

<sup>۲</sup> گروه مهندسی برق کنترل- دانشکده فنی و مهندسی- دانشگاه بین‌المللی امام خمینی (ره)- قزوین- ایران

### چکیده

### مقاله پژوهشی

#### تاریخ دریافت:

۱۴۰۱/۰۸/۲۴

#### تاریخ پذیرش:

۱۴۰۱/۱۰/۱۹

#### کلیدواژه‌ها:

تجزیه نامنفی ماتریسی،  
خوشه‌بندی گراف، خوشه‌بندی  
نیمه نظارتی، معیار پیمانگی

#### نویسنده مسئول:

n.bigdeli@eng.ikiu.ac.ir

خوشه‌بندی، ابزاری پرکاربرد جهت تحلیل اطلاعات شبکه‌های پیچیده است که برای مدل‌سازی سامانه‌های پیچیده بکار می‌رود. پیمانگی، معیاری پایه و فراگیر جهت ارزیابی و صحت‌سنجی خوشه‌بندی شبکه‌ها است که دارای چالش‌هایی چون ان‌پی-سخت بودن مسئله و عدم امکان استفاده از دانش اولیه در خوشه‌بندی است. لذا، خوشه‌بندی مبتنی بر معیار پیمانگی، قابلیت تعمیم به خوشه‌بندی‌های نیمه نظارتی را ندارد. از طرفی، یکی از روش‌های خوشه‌بندی نیمه نظارتی، روش خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی (NMF) است. اما این روش، ویژگی‌های خاص شبکه‌ها را در نظر نمی‌گیرد. در این مقاله، برای غلبه بر چالش‌های نام‌برده و با ارائه اثباتی جدید، برای خوشه‌بندی مبتنی بر معیار پیمانگی، ساختاری مشابه با خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی نامتقارن ارائه می‌شود که در آن، امکان بهره‌گیری از دانش اولیه و حل به روش تکراری میسر می‌گردد. سپس، روش خوشه‌بندی نیمه نظارتی نوینی به نام تجزیه نیمه نظارتی نامنفی ماتریس‌های متقارن مبتنی بر معیار پیمانگی (SSNMF-Q) با بهره‌گیری از مزیت دانش اولیه و روش حل تکراری، به جای حل مسئله ان‌پی-سخت ارائه می‌گردد. برای ارزیابی روش پیشنهادی، از پنج مجموعه داده واقعی استفاده شده که نتایج، بیانگر عملکرد بهتر SSNMF-Q در مقایسه با سایر خوشه‌بندی‌های نیمه نظارتی مبتنی بر NMF است.



## ۱- مقدمه

در دنیای کنونی سامانه‌های پیچیده توسط شبکه‌هایی همانند شبکه‌های اجتماعی، شبکه‌های ترافیکی و شبکه‌های بیولوژیکی مدل‌سازی می‌شوند [۱-۳]. با استفاده از ابزار تئوری گراف، می‌توان شبکه‌ها را توسط مجموعه‌ای از گره‌ها و یال‌ها نمایش داد. به‌طور مثال، شبکه‌های اجتماعی یکی از سامانه‌های پیچیده در دنیای کنونی هستند که برای هر ویژگی این شبکه گراف‌هایی ترسیم می‌شود. گراف ارتباطی اشخاص از طریق دنبال کننده و دنبال شونده، گرافی جهت‌دار برای یک مجموعه انسانی را تشکیل خواهد داد که هر حساب شبکه اجتماعی به‌عنوان یک گره و ارتباط دنبال کننده و دنبال شونده میان این اعضا به‌عنوان یال‌های جهت‌دار این گراف در نظر گرفته خواهند شد. با تحلیل گراف‌ها می‌توان اطلاعات بسیار مفیدی از ساختار این شبکه‌های پیچیده به دست آورد. در این زمینه، یکی از ابزارهای متداول و مؤثر جهت بررسی ساختار شبکه‌های پیچیده، خوشه‌بندی گراف‌ها است [۲]. در مثال گراف ارتباطی شبکه اجتماعی، این تحلیل به ما نشان می‌دهد که کدام حساب‌ها به دلیل نوع ارتباط در شبکه به همدیگر نزدیک‌تر هستند و در هر خوشه قرار خواهند گرفت. از این رو، با توجه به [۳] در سال‌های گذشته، روش‌های متنوعی چون خوشه‌بندی مبتنی بر پیمانی [۴، ۵]، خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی [۶]، خوشه‌بندی مبتنی بر انتشار برچسب [۷، ۸]، خوشه‌بندی مبتنی بر پیاده‌روی تصادفی [۹، ۱۰] و خوشه‌بندی بر اساس بهینه‌سازی تکاملی [۱۱] جهت خوشه‌بندی گراف‌ها ارائه شده است.

مطابق با [۴]، معیار پیمانی مشخصاً به‌عنوان معیار صحت‌سنجی خوشه‌بندی گراف‌ها معرفی شده است. این معیار، چگالی ارتباط میان گره‌ها در داخل هر گروه و ارتباط میان گروهی را به‌عنوان شاخصی جهت بررسی صحت خوشه‌بندی‌ها در نظر می‌گیرد. از این رو پیمانی، معیاری محبوب و متداول جهت خوشه‌بندی شبکه‌های گرافی است. به‌عبارت‌دیگر، این معیار مختص خوشه‌بندی شبکه‌های گرافی است. از این رو در این سال‌ها به دلیل فراگیر شدن این معیار، روش‌های مبتنی بر معیار پیمانی متنوعی

جهت خوشه‌بندی گراف‌های پیچیده ایجاد شده است. از جمله این روش‌ها می‌توان به روش‌های چون روش حریصانه<sup>۱</sup> [۱۲]، روش حل بهینه‌سازی تکاملی [۱۳] و روش طیفی<sup>۲</sup> [۴] اشاره نمود. اگرچه خوشه‌بندی مبتنی بر معیار پیمانی روشی فراگیر جهت خوشه‌بندی است اما محدودیت‌هایی چون وابسته بودن به مجموع یال‌ها [۱۴] و عدم استفاده از دانش اولیه خوشه‌بندی را دارند. همچنین مطابق با [۴] روش حل خوشه‌بندی مبتنی بر پیمانی مسئله آن‌پی-سخت<sup>۳</sup> است.

خوشه‌بندی نیمه نظارتی، به روش‌های خوشه‌بندی اطلاق می‌شود که در آن‌ها، دانش اولیه از اعضای هر گروه و یا تعداد آن‌ها در فرآیند تشخیص گروه‌ها دخالت داده شده و از این دانش اولیه در خوشه‌بندی استفاده می‌شود. روش‌های نیمه نظارتی مبتنی بر تجزیه نامنفی ماتریس، از جمله روش‌های متداول خوشه‌بندی در راستای بهره‌گیری از دانش اولیه اعضای هر گروه هستند. در روش‌های خوشه‌بندی، ماتریس‌های مشابهت گراف (مانند ماتریس مجاورتی<sup>۴</sup>) به چندین ماتریس کاهش بعد داده‌شده تبدیل می‌شوند که با استفاده از آن‌ها می‌توان خوشه‌بندی و خوشه‌های مدنظر را به دست آورد [۱۵، ۱۶]. در این روش‌ها برای استخراج گروه‌ها، از مدل‌ها و توابع متفاوتی همچون مدل تجزیه نامنفی ماتریسی متقارن [۱۶]، مدل مقاوم تجزیه نامنفی ماتریس [۱۷، ۱۸]، مدل تجزیه نامنفی ماتریسی سه عاملی [۱۹]، مدل تجزیه نامنفی ماتریسی عمیق [۲۰-۲۲] و مدل گرافی تجزیه نامنفی ماتریسی منظم شده [۲۳، ۲۴] استفاده شده است. مطابق با [۱۷، ۱۸] مدل مقاوم تجزیه نامنفی ماتریسی توانسته است با استفاده از نرم  $l_2$  عدم قطعیت‌هایی همچون وارد کردن دانش‌های اولیه اشتباهی و خطاهای ترکیب انواع نوع داده را کم‌رنگ‌تر کند. در مدل تجزیه نامنفی ماتریسی سه عاملی [۱۹]، معیار پیمانی با مدل سه عاملی، ترکیبی جدید تشکیل داده است که منتج به بهبود عملکرد در خوشه‌بندی نهایی شده است. در مدل‌های تجزیه نامنفی ماتریسی عمیق [۲۰-۲۲]، خطای حل به روش تکراری به شدت کاهش پیدا کرده است، همچنین به کمک معیار پیمانی، بهبود قابل‌توجهی در مدل‌سازی‌ها قابل مشاهده است.

<sup>3</sup> NP hard

<sup>4</sup> Adjacency matrix

<sup>1</sup> Greedy algorithm

<sup>2</sup> Spectral algorithm



نامنفی ماتریسی است. در نتیجه، معیار پیمانی می‌تواند جایگزین ماتریس مجاورتی در خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی شود. این مزیت و نوآوری به ما کمک خواهد کرد تا از مزایای روش خوشه‌بندی بر پایه تجزیه نامنفی ماتریسی همچون روش حل تکراری و تعمیم به خوشه‌بندی نیمه نظارتی را برای ساختار جدید تعریف‌شده<sup>۱</sup> مبتنی بر معیار پیمانی استفاده شود. از این رو در این مقاله مدل جدیدی برای خوشه‌بندی نیمه نظارتی به نام  $SSNMF-Q$  پیشنهاد شده است تا بتواند خوشه‌بندی مبتنی بر معیار پیمانی را بهبود بخشد. لذا، نوآوری‌ها و مراحل رسیدن به این هدف به شرح زیر است.

۱. اثبات خواهد شد که خوشه‌بندی مبتنی بر پیمانی دارای ساختاری مشابه با خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی است. این مشابه‌سازی کمک خواهد کرد تا ترکیبی از دانش اولیه<sup>۱</sup> خوشه‌بندی گراف‌ها در این خوشه‌بندی نوین ارائه شود.

۲. ساختاری مشابه‌سازی شده، به‌عنوان تابع اصلی روش خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی انتخاب می‌شود. سپس بخش‌های موردنیاز از دانش اولیه به این تابع اضافه شده و مدل نهایی نیمه نظارتی  $SSNMF-Q$  تولید می‌شود.

۳. نشان داده می‌شود که مدل  $SSNMF-Q$  می‌تواند دانش اولیه را با روش خوشه‌بندی مبتنی بر معیار پیمانی ترکیب کند و کیفیت خوشه‌بندی را با افزایش دانش اولیه بهبود بخشد.

۴. برای الگوریتم  $SSNMF-Q$ ، روش حل تکراری طراحی شده و نتایج خوشه‌بندی بر روی پنج مجموعه داده واقعی نمایش داده می‌شود. در نهایت نیز بر اساس نتایج به‌دست‌آمده، تأثیر مدل ارائه‌شده بر روی کیفیت پاسخ بررسی می‌شود.

در ادامه ساختار این مقاله این‌گونه تنظیم شده است که در بخش ۲ به معرفی مختصری از خوشه‌بندی مبتنی بر معیار پیمانی و بررسی روش‌های پیشین خوشه‌بندی نیمه نظارتی مبتنی بر تجزیه نامنفی ماتریسی پرداخته شده است. در بخش ۳، روش حل تکراری الگوریتم نیمه نظارتی پیشنهادی بررسی شده است و در بخش ۴، معیارهای ارزیابی و مجموعه داده‌هایی جهت آزمایش معرفی شده‌اند. سپس پارامترهای موردنیاز استخراج شده و نتایج بر

نکته قابل‌توجه آن است که همه<sup>۱</sup> مدل‌های خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی از جنس حل مسئله<sup>۱</sup> بهینه‌سازی می‌باشند؛ لذا می‌توان دانش اولیه را در ماتریس‌های مشابهت گراف ترکیب کرد [۲۵، ۲۶] و یا آن را به‌صورت عبارت جداگانه‌ای به مسئله بهینه‌سازی اضافه نمود [۲۷-۲۹]. لازم به ذکر است که تمامی این روش‌ها عموماً برای پیدا کردن خوشه‌بندی بهتر، از ماتریس مجاورتی استفاده می‌کنند.

مزیت اصلی خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی، کاربرد آن برای انواع داده‌ها از جمله داده‌های صوتی، تصویری و متنی و گرافی است. تاکنون از این خوشه‌بندی برای مسائلی چون جداسازی منابع صوتی [۳۰]، اطلاعات چندوجهی [۳۱]، پردازش تصویر [۳۲]، استخراج کلمات کلیدی و عنوان اسناد [۳۳] استفاده شده است. هرچند که این عمومیت و گستردگی کاربرد برای انواع داده‌های از جنس متفاوت، یکی از مزایای مهم تجزیه نامنفی ماتریسی است، اما این روش قابلیت بررسی و استخراج ویژگی‌های نوع خاصی از داده را به‌صورت تخصصی ندارد. به عبارت دیگر نوع داده (داده صوتی، شبکه‌های گرافی، متنی، تصویر) در خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی اثری نخواهد داشت. در راستای بهبود این محدودیت، به‌عنوان نمونه در [۱۹، ۲۲] معیار پیمانی به روش خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی به‌صورت خطی اضافه شده است تا کیفیت خوشه‌بندی را برای شبکه‌های گرافی افزایش دهد.

با توجه به مطالب فوق و مزیت روش خوشه‌بندی مبتنی بر معیار پیمانی یعنی مختص خوشه‌بندی شبکه‌های گرافی بودن و محدودیت‌های این روش چون عدم استفاده از دانش اولیه و ان‌پی-سخت بودن مسئله، و همچنین مزایای خوشه‌بندی بر پایه تجزیه نامنفی ماتریسی همچون توانایی استفاده از دانش اولیه و حل تکراری و محدودیت عمومیت برای تمامی داده‌ها و عدم تأثیر نوع داده، ما را بر این داشت در این مقاله راه‌حل نوینی با بهره‌گیری از مزایای تجزیه نامنفی ماتریسی برای خوشه‌بندی مبتنی بر معیار پیمانی ارائه می‌شود. بدین منظور روش خوشه‌بندی مبتنی بر معیار پیمانی دارای ساختاری مشابه با خوشه‌بندی مبتنی بر تجزیه

<sup>1</sup> Semi-Supervised Symmetric Nonnegative Matrix Factorization based on modularity ( $Q$ )



$$\max_H Q = \max_H \frac{1}{2m} \text{tr}(H^T B H), \quad (1)$$

$$B = A - B_1$$

$$S. t. H^T H = I$$

که  $B$  ماتریس معرفی معیار پیمانگی،  $A$  ماتریس مجاورتی گراف،  $H \in R^{n \times k}$  ماتریس اعضای گروه برای هر گروه،  $(B_1)_{ij} = \frac{k_i k_j}{2m}$  و  $k$  تعداد خوشه‌ها در شبکه است. همچنین  $k_i$  و  $k_j$  درجه مربوط به گره  $i$  ام و  $j$  ام برای محاسبه ماتریس  $B_1$  می‌باشند.

## ۲-۲ خوشه‌بندی مبتنی بر پایه تجزیه نامنفی ماتریسی

در روش عمومی خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی، ماتریس مشابهت گراف به دو ماتریس کاهش بعد یافته شده تبدیل می‌شود و سپس از روی ماتریس‌های تجزیه‌شده، خوشه‌بندی نهایی به دست می‌آید. در به دست آوردن این ماتریس‌های تجزیه‌شده سعی می‌شود که ویژگی‌های پنهان ماتریس اولیه در ماتریس‌های کاهش بعد یافته حفظ شود. به عبارت دیگر با فرض تعداد خوشه‌بندی  $k$ ، ماتریس مشابهت  $X \in R^{n \times n}$  به دو ماتریس جدید  $H \in R^{n \times k}$  و  $U \in R^{n \times k}$  تجزیه می‌شود ( $X \approx HU^T$ ) که ماتریس  $H$  بیانگر ماتریس ارتباطات جامعه<sup>۱</sup> و  $U$  ماتریس اعضای هر گروه<sup>۲</sup> است. برای به دست آوردن ماتریس‌های کاهش بعد داده‌شده، مسئله بهینه‌سازی با تابع هزینه<sup>۳</sup> ذیل معرفی می‌شود:

$$\min_{H,U} J_{nmf}(H,U) = \|X - HU^T\|_F^2 \quad (2)$$

که در آن  $\| \cdot \|_F$  نرم فروبینوس است. این تابع بهینه‌سازی از مدل تابع بهینه‌سازی غیرمحدب است که در پژوهش‌های متعددی همچون [۱۵]، برای به دست آوردن دو ماتریس  $H$  و  $U$ ، روش‌های حل تکراری مبتنی بر توابع لاگرانژین پیشنهاد شده است.

یکی دیگر از مدل‌ها، تابع تجزیه نامنفی ماتریسی متقارن<sup>۴</sup> است که با کاهش متغیرها توانسته است خطای مدل‌سازی را کاهش دهد. با توجه به [۱۶]، تابع بهینه‌سازی به صورت رابطه (۳) است.

$$\min_H J_{SNMF}(H) = \|X - HH^T\|_F^2 \quad (3)$$

که ماتریس  $H \in R^{n \times k}$ ، بیانگر اعضای هر گروه خواهد بود. با مقایسه روابط (۲) و (۳) مشاهده می‌شود که مزیت این روش، گنجاندن اطلاعات ارتباطات جامعه در ماتریس  $H$  است.

روی مجموعه داده‌های واقعی نمایش داده شده است و در انتها در بخش ۵، نتیجه‌گیری کلی و خلاصه‌ای از کار بیان شده است.

## ۲- خوشه‌بندی گراف‌ها

در این بخش، ابتدا تعریف مختصری از معیار پیمانگی به‌عنوان معیار صحت‌سنجی در خوشه‌بندی ارائه می‌گردد. سپس به بررسی روش‌های خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی پرداخته می‌شود و در نهایت خوشه‌بندی گراف‌ها از دیدگاه‌های مختلف نیمه نظارتی مورد بررسی قرار می‌گیرد.

### ۲-۱ خوشه‌بندی مبتنی بر معیار پیمانگی

گراف  $G = (V, E)$  در نظر گرفته شده است که در آن  $V$  مجموعه گره‌ها با مجموع تعداد  $n$  گره و  $E$  مجموعه یال‌ها میان دو گره با مجموع تعداد  $m$  یال مجزا است. در سال‌های گذشته، برای تبدیل گراف‌ها به روابط ریاضی، ماتریس‌های مشابهت گراف همانند ماتریس لاپلاسیان<sup>۱</sup> ( $L$ ) و ماتریس مجاورتی ( $A$ ) معرفی شده است. همچنین برای بررسی خوشه‌بندی، معیار پیمانگی ( $Q$ ) بر اساس چگالی یال‌های موجود در گروه‌ها و ارتباطات میان گروهی، صحت‌سنجی گروه‌بندی‌ها را مشخص می‌کند. به‌طور مثال این معیار به صورت  $Q = \frac{1}{2m} \sum_{ij} \{A_{ij} - \frac{k_i k_j}{2m}\} \delta(x_i, x_j)$  برای صحت‌سنجی خوشه‌بندی تعریف می‌شود. که  $k_i$  و  $k_j$  درجه مربوطه گره  $i$  ام و  $j$  ام و  $\delta(x_i, x_j)$  معیار هم‌گروه بودن یا نبودن دو گره  $i$  ام و  $j$  ام است. همچنین بدیهی است که هر چه این مقدار بزرگ‌تر باشد کیفیت خوشه‌بندی بهتر خواهد بود. همچنین این معیار به صورت مستقل جهت تعیین خوشه‌های شبکه‌های پیچیده نیز استفاده می‌شود [۴، ۱۹]. از این رو به‌طور کلی خوشه‌بندی مبتنی بر معیار پیمانگی به شکل مسئله بهینه‌سازی بر روی معیار پیمانگی و شرایط ذیل قابل بازنویسی است [۱۹]:

<sup>3</sup> Community membership matrix

<sup>4</sup> Symmetric NMF (SNMF)

<sup>1</sup> Laplacian matrix

<sup>2</sup> Community relation matrix

## ۳-۲ خوشه‌بندی نیمه نظارتی مبتنی بر تجزیه نامنفی

### ماتریسی

روش‌های خوشه‌بندی نیمه نظارتی، به روش‌هایی اطلاق می‌گردد که اعمال دانش اولیه از اعضای برخی از گروه‌ها باعث بهبود خوشه‌بندی شود. تاکنون دانش اولیه از اعضای خوشه، به دو صورت عضو بودن در گروه و عضو نبودن در گروه‌های دیگر مشخص شده است. همچنین میزان دانش اولیه از اعضای هر گروه به صورت درصدی از دانش اعضای مشخص می‌شود. یکی از متداول‌ترین و فراگیرترین روش‌های خوشه‌بندی نیمه نظارتی، روش‌های خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی است. با توجه به سیر زمانی تکامل روش‌های نیمه نظارتی، این روش‌ها معمولاً دو رویکرد کلی را دنبال می‌کنند. رویکرد اول تغییر در ماتریس مشابهت خوشه‌بندی و رویکرد دوم تغییر در تابع بهینه‌سازی تجزیه نامنفی ماتریسی است. با توجه به امکان تغییرات در ترکیب خطی مدل خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی و مدل دانش اولیه، رویکرد دوم رویکرد به‌روزتری در سال‌های اخیر بوده است. در ادامه سعی شده است که برای این دو رویکرد توضیحات تکمیلی ارائه گردد.

رویکرد اول، تغییر در ماتریس مشابهت خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی است. ژانگ در [۲۵] دانش اولیه اطلاعات برچسب گروه‌ها را با ساختار اولیه گراف ترکیب نموده است. به‌صورتی که ماتریس مجاورتی تعمیم‌یافته  $\bar{A}$  را جایگزین ماتریس مجاورتی اصلی کرده است. این ماتریس تعمیم‌یافته به صورت ذیل تعریف شده است:

$$\bar{A} = \begin{cases} \alpha & \text{اگر } x_i, x_j \text{ دارای برچسب یکسان باشند} \\ 0 & \text{اگر } x_i, x_j \text{ دارای برچسب متفاوت باشند} \\ A_{ij} + I_{ij} & \text{حالات دیگر} \end{cases} \quad (4)$$

که  $x_i$  یا  $x_j$  بیانگر برچسب اعضای خوشه برای گره  $i$  ام و  $j$  ام است. و همچنین در رابطه (۴)،  $\alpha$  مقدار مثبت برابر با ۲ و ۱ ماتریس یکه است. همان‌طور که از رابطه (۴) مشخص است این روش با در نظر گرفتن دانش از اعضای هر گروه، ماتریس مشابهت

جدیدی برای گراف طراحی می‌کند تا بتواند خوشه‌بندی بهتری با در نظر دانش اولیه اعضای هر گروه ایجاد کند.

در تحقیق دیگر [۲۶]، ماتریس مجاورتی گراف بادانش اولیه جفت شدن گروه‌ها در هر گروه و تفاوت گروه‌های میان گره‌های ماتریس مجاورتی بازنویسی شده است که الگوریتم  $SNMF-SS^1$  پدید آمده است. ماتریس ورودی خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی به صورت ذیل تعریف می‌شود:

$$\bar{X} = X - \alpha W_{ML} + \beta W_{CL} \quad (5)$$

که در آن  $W_{ML}$  ماتریسی تشکیل‌شده از رابطه میان اعضای هر گروه و  $W_{CL}$  ماتریس تشکیل‌شده از رابطه میان اعضای گروه‌های متفاوت را نشان می‌دهد.

در رویکرد دوم برای استفاده از دانش اولیه، هدف اضافه کردن عبارت جدید به تابع بهینه‌سازی خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی است؛ بنابراین تابع هزینه خوشه‌بندی از مجموع بخش خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی و بخش مربوط به دانش اولیه خوشه‌بندی تشکیل شده است. به‌طور معادل رابطه (۶) تابع هزینه کلی برای این رویکرد را نشان می‌دهد.

$$J_{SNMF} = J_{clustering} + J_{pairs} \quad (6)$$

که در آن تابع  $J_{pairs}$  تابع بهینه‌سازی مربوط به دانش اولیه و  $J_{clustering}$  تابع مربوط به خوشه‌بندی است. در ادامه سعی شده است که پژوهشی‌های مرتبط با این رویکرد به‌طور مختصر مورد بررسی قرار گیرد.

به‌طور مثال، یانگ در [۲۷]، روش  $PCNMF^2$  را به‌نحوی که دانش اولیه به صورت رابطه‌ای جدید وارد بهینه‌سازی تجزیه نامنفی ماتریس شود، مطرح کرد. در این رابطه جدید، تابع بهینه‌سازی به صورت ذیل بازنویسی می‌شود:

$$J_{PCNMF} = \|X - HU^T\|_F^2 + \lambda tr(U^T LU) \quad (7)$$

که ماتریس  $X$  همان ماتریس مجاورتی است. همچنین  $L = D - A_1$  است که در آن ماتریس قطری  $D$  با اعضای  $D_{ii} = \sum_{j=1}^n (A_1)_{ij}$  است و ماتریس  $A_1$  با اطلاعات از دانش اولیه به صورت ذیل بازنویسی می‌شود:

<sup>2</sup> Pairwise Constraints-guided Nonnegative Matrix Factorization

<sup>1</sup> Semi-Supervised Nonnegative Matrix -based semi-supervised clustering

قطعی‌ها، تابع بهینه‌سازی مطابق با رابطه (۱۱) در مقابل خطا در دانش اولیه از اعضای هر گروه مقاوم شده است و خوشه‌بندی بهتری نسبت به حضور این عدم قطعیت‌ها به وجود آورده است.

### ۳- الگوریتم خوشه‌بندی پیشنهادی

با توجه به بررسی روش‌های مختلف و نیازمندی‌های مطرح‌شده اعم از ارائه روش حل تکراری برای خوشه‌بندی مبتنی بر معیار پیمانی و استفاده از دانش اولیه اعضای گروه در این روش خوشه‌بندی، در این بخش ابتدا در قضیه ۱ مشابهت روش‌های خوشه‌بندی مبتنی بر معیار پیمانی و خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی اثبات می‌شود و سپس، الگوریتم خوشه‌بندی پیشنهادی معرفی می‌گردد.

#### ۳-۱ ارتباط خوشه‌بندی مبتنی بر پیمانی و خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی

در این بخش، جهت تلفیق معیار پیمانی با خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی، مشابهت میان دو خوشه‌بندی مبتنی بر معیار پیمانی و تجزیه نامنفی ماتریسی در قالب قضیه ۱ و با توجه به مطالب بخش‌های ۲-۲ و ۲-۱ بیان و اثبات می‌شود.

قضیه ۱: خوشه‌بندی مبتنی بر پیمانی در شبکه‌های پیچیده دارای ساختاری مشابه با خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی است. یا به عبارتی خواهیم داشت.

$$\max_H Q = \min_H \|B - HH^T\|_F^2 \quad (12)$$

اثبات: مطابق با اثبات مشابهت میان دو خوشه‌بندی مبتنی بر معیار چگالی پیمانی<sup>۳</sup> و تجزیه نامنفی ماتریسی در [۳۴]، تابع بهینه‌سازی معیار پیمانی با توجه به رابطه (۱) به صورت ذیل بازنویسی می‌شود:

$$\max_H \frac{1}{2m} \text{tr}(H^T B H) \propto -\frac{1}{2m} \min_H \text{tr}(H^T B H) \quad (13)$$

اگر شرط  $H^T H = I$  و  $H$  ماتریس ثابت در نظر گرفته شود، خواهیم داشت:

$$(A_1)_{ij} = \begin{cases} \alpha & \text{اگر } x_i, x_j \text{ دارای برجسب یکسان باشند} \\ -(1-\alpha) & \text{اگر } x_i, x_j \text{ دارای برجسب متفاوت باشند} \\ 0 & \text{حالات دیگر} \end{cases} \quad (8)$$

در رابطه فوق  $\alpha$  مقدار مثبت و برابر با ۲ است. در پژوهشی دیگر [۲۸]، روش PCSNMF<sup>۱</sup> پیشنهادشده است در این روش دانش اولیه اعضا و ارتباط‌های هر گره در هر گروه و در گروه‌های مختلف با روش تجزیه نامنفی ماتریسی ترکیب شده است. تابع بهینه‌سازی مطرح‌شده برای روش PCSNMF به صورت ذیل است:

$$J_{PCSNMF} = \|X - HH^T\|_F^2 + \alpha(\text{tr}(H^T M H B_2) + \text{tr}(H^T C H)) \quad (9)$$

که در آن  $B_2 = \begin{bmatrix} 0 & 1 & \dots & 1 \\ 1 & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 1 \\ 1 & \dots & 1 & 0 \end{bmatrix}$  همان ماتریس مجاورتی تعریف می‌شود و دانش اولیه در ماتریس‌های ذیل خلاصه‌شده است:

$$M_{ij} = \begin{cases} 1 & \text{اگر } x_i, x_j \text{ دارای برجسب یکسان باشند} \\ 0 & \text{حالات دیگر} \end{cases} \quad (10)$$

$$C_{ij} = \begin{cases} 1 & \text{اگر } x_i, x_j \text{ دارای برجسب متفاوت باشند} \\ 0 & \text{حالات دیگر} \end{cases}$$

علاوه بر آن مطابق با [۱۷] با توجه به خطاهای انسانی در انتخاب اولیه خوشه‌بندی گره‌ها، روش RSSNMF<sup>۲</sup> به عنوان روشی مقاوم جهت بهبود کارایی در خطاهای انسانی مطرح‌شده است که این روش توانسته است روش PCSNMF را بهبود ببخشد. تابع هزینه بهینه‌سازی بهبود داده‌شده آن به صورت ذیل بازنویسی شده است:

$$J_{RSSNMF} = \|X - HU^T\|_2 + \alpha \text{tr}(U^T M U B_2) + \beta \text{tr}(U^T C U) \quad (11)$$

که در آن  $\| \cdot \|_2$  نرم مرتبه دوم و  $X$  همان ماتریس مجاورتی است. به دلیل توانایی‌های ذاتی نرم مرتبه دوم در حذف نویز و عدم

<sup>2</sup> Robust Semi-Supervised Nonnegative Matrix Factorization

<sup>3</sup> Modularity density

<sup>1</sup> Pairwisely Constrained Symmetric Nonnegative Matrix Factorization

ترکیبی از مجموع اعضای اولیه هر گروه و اعضای دیگر گروه‌ها خواهد بود و بخش خوشه‌بندی از روش خوشه‌بندی قضیه ۱ استفاده خواهد کرد. در نتیجه خوشه‌بندی نیمه نظارتی پیشنهادی به صورت ذیل است:

$$\min_H J_{SSNMF-Q} = \|B - HH^T\|_F^2 + \alpha \text{tr}(HH^T B_2) + \beta \text{tr}(HCH^T) \quad (17)$$

$$s. t. H > 0, \sum_{r=1}^k H_{ir} = 1$$

که در آن ماتریس‌های  $B_2$ ،  $M$  و  $C$  با توجه به رابطه‌های (۹) و (۱۰) تعریف خواهند شد.

بدین ترتیب، با در نظر گرفتن پارامترهای داخلی متغیر، بهبود کارایی در بهینه‌سازی الگوریتم SSNMF-Q نسبت به الگوریتم PCSNMF حاصل شده است و با در نظر گرفتن قضیه ۱، ماتریس اختصاصی خوشه‌بندی همچون معیار پیمانی به عنوان ماتریس مشابهت الگوریتم خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی معرفی شده است. همچنین با توجه به رابطه (۱۷) و قضیه ۱، الگوریتم ارائه شده برای خوشه‌بندی مبتنی بر معیار پیمانی یک الگوریتم نیمه نظارتی است.

### ۳-۳ روش حل تکراری تابع بهینه‌سازی SSNMF-Q

از آنجاکه در رابطه (۱۷) ماتریس  $H$  ماتریس نامنفی است در نتیجه خوشه‌بندی نهایی از طریق روش حل مبتنی بر لاگرانژین قابل استخراج است. بر این اساس رابطه (۱۷) به صورت فرم ترانهاده آن یعنی رابطه (۱۸) بازنویسی می‌شود.

$$\min_H J_{SSNMF-Q} = (\text{tr}(BB^T) - 2\text{tr}(BHH^T) + \text{tr}(HH^T HH^T) + \alpha \text{tr}(HH^T B_2) + \beta \text{tr}(HCH^T)) \quad (18)$$

$$s. t. H > 0, \sum_{r=1}^k H_{ir} = 1$$

با فرض در نظر گرفتن ماتریس  $\Phi$  به عنوان ماتریس لاگرانژین، تابع لاگرانژین به صورت ذیل تعریف می‌شود:

$$L(H) = J_{SSNMF-Q} + \text{tr}(\Phi H) \quad (19)$$

مشقات جزئی تابع لاگرانژین نسبت به متغیرهای مسئله به صورت زیر خواهد بود:

$$\max_H \frac{1}{2m} \text{tr}(H^T B H) \propto \min_H (\text{tr}(H^T H H^T H) - 2\text{tr}(H^T B H) + \text{tr}(H H^T)) \quad (14)$$

با توجه به ویژگی‌های ترانهاده‌های ماتریس چون  $\text{tr}(H^T B H) = \text{tr}(B H H^T)$  و  $\text{tr}(H^T H H^T H) = \text{tr}(H H^T H H^T)$  رابطه (۱۴) به صورت ذیل بازنویسی می‌شود:

$$\max_H \frac{1}{2m} \text{tr}(H^T B H) \propto \min_H \text{tr}(H H^T H H^T - 2B H H^T + B B^T) \propto \min_H \|B - H H^T\|_F^2 \quad (15)$$

با توجه به رابطه (۱)، رابطه (۳) و رابطه اثبات شده (۱۵)، تابع هزینه بهینه‌سازی ماتریس پیمانی دارای ساختاری مشابه با تابع هزینه بهینه‌سازی تجزیه نامنفی ماتریسی است. به عبارت دیگر خواهیم داشت:

$$\max_H Q = \max_H \frac{1}{2m} \text{tr}(H^T B H) \propto \min_H \|B - H H^T\|_F^2 \quad (16)$$

از این رو می‌توان اثبات کرد که خوشه‌بندی مبتنی بر معیار پیمانی دارای ساختاری مشابه با خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی است.

### ۳-۲ الگوریتم SSNMF-Q

با توجه به قضیه ۱، نشان داده شد که خوشه‌بندی مبتنی بر معیار پیمانی از جنس خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی با فرض ماتریس مشابهت پیمانی است. مزیت این مشابه‌سازی بهره‌گیری از معیار پیمانی در روش خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی و استفاده از ساختاری جدید برای خوشه‌بندی مبتنی بر معیار پیمانی است. در نتیجه با در نظر گرفتن این مشابه‌سازی و با توجه به حل ان‌پی-سخت بودن و عدم ارائه راه‌حل نیمه نظارتی روش حل خوشه‌بندی مبتنی بر پیمانی، در این مقاله روش نوین خوشه‌بندی نیمه نظارتی مبتنی بر ماتریس مشابهت پیمانی پیشنهاد می‌گردد. این ترکیب علاوه بر استفاده از دانش خوشه‌بندی گراف (پیمانی)، از دانش‌های اولیه برای بهبود خوشه‌بندی مبتنی بر پیمانی و تجزیه نامنفی ماتریسی استفاده خواهد کرد تا راه‌حل تکراری و روش خوشه‌بندی نیمه نظارتی مبتنی بر پیمانی ارائه گردد. بدین منظور با ترکیب بخش خوشه‌بندی و بخش دانش اولیه از برخی از اعضای هر گروه، می‌توان تابع بهینه‌سازی (۱۷) را پیشنهاد داد که این دانش به صورت

#### ۴- نتایج و شبیه‌سازی

جهت بررسی و ارزیابی روش‌ها خوشه‌بندی نیمه نظارتی SSNMF-Q با دیگر روش‌ها، ابتدا علاوه بر معیار پیمانی معیار اطلاعات متقابل ( $NMI^*$ ) معرفی می‌شوند. سپس مجموعه شبکه متداول در این راستا معرفی می‌گردد و در انتها با توجه به شرایط مختلف دانش اولیه، ارزیابی و صحت‌سنجی الگوریتم‌های خوشه‌بندی بر روی این دادگان انجام می‌گیرد.

#### ۴-۱ معیار اعتبار سنجی

در این مقاله، معیارهای پیمانی ( $Q$ ) و اطلاعات متقابل نرمالیزه به‌عنوان شاخص، جهت بررسی و ارزیابی روش‌های خوشه‌بندی در نظر گرفته شده‌اند. بر این اساس علاوه بر معرفی معیار پیمانی در رابطه (۱)، در معیار اطلاعات متقابل از دانش برچسب‌های خوشه‌بندی واقعی<sup>۳</sup> هر گره، برای ارزیابی روش‌ها استفاده می‌شود. همچنین این معیار به دلیل استفاده از دانش برچسب‌های خوشه‌بندی واقعی، معیاری فراگیر جهت ارزیابی روش‌های خوشه‌بندی گراف است. در نتیجه معیار اطلاعات متقابل می‌تواند به‌صورت ذیل بازنویسی شود:

$$NMI(C, C') = \frac{-2 \sum_{i=1}^{|C|} \sum_{j=1}^{|C'|} n_{C_i \cap C'_j} \times \log \left( \frac{n \times n_{C_i \cap C'_j}}{n_{C_i} \times n_{C'_j}} \right)}{\sum_{i=1}^{|C|} n_{C_i} \times \log \left( \frac{n_{C_i}}{n} \right) + \sum_{j=1}^{|C'|} n_{C'_j} \times \log \left( \frac{n_{C'_j}}{n} \right)} \quad (24)$$

که  $C$  و  $C'$  بیانگر دو خوشه‌بندی مختلف از اعضای شبکه است که در آن  $n_{C_i}$  تعداد اعضا در گروه  $C_i$  و  $|C|$  تعداد کل گروه‌ها است. همچنین این معیار زمانی برابر با بیشترین مقدار خود یعنی یک خواهد شد که خوشه‌بندی  $C$  به‌طور کامل مشابه با خوشه‌بندی  $C'$  باشد و همچنین اگر  $C$  و  $C'$  به‌طور کامل هیچ شباهتی در اعضای هر خوشه با همدیگر نداشته باشند این معیار برابر با کمترین مقدار خود یعنی صفر خواهد بود.

$$\frac{\partial L_{SSNMF-Q}}{\partial H} = \varphi - 4BH + 4HH^T H + 2\alpha B_2 HM + 2\beta C \quad (20)$$

از آنجا که  $B = A - B_1$  است و شرایط KKT ( $\varphi_{ij} H_{ij} = 0$ ) برقرار است. رابطه (۲۰) به‌صورت رابطه (۲۲) بازنویسی می‌شود:

$$-4(AH)_{ij} H_{ij} + 4(B_1 H)_{ij} H_{ij} + 4(HH^T H)_{ij} H_{ij} + 2\alpha(B_2 HM)_{ij} H_{ij} + 2\beta(CH)_{ij} H_{ij} = 0 \quad (21)$$

در نتیجه قانون به‌روزرسانی تکراری به‌صورت ذیل طراحی می‌شود.

$$H_{ij} = H_{ij} \frac{4(AH)_{ij}}{4(B_1 H)_{ij} + 4(HH^T H)_{ij} + 2\alpha(B_2 HM)_{ij} + 2\beta(CH)_{ij}} \quad (22)$$

با توجه به [۱۹]، شرط  $\sum_{r=1}^k H_{ir} = 1$  به شکل زیر قابل اعمال است.

$$H_{ir} := \frac{H_{ir}}{\sum_{r=1}^c H_{ir}} \quad (23)$$

در نهایت الگوریتم SSNMF-Q در ادامه تعریف می‌شود:

#### الگوریتم خوشه‌بندی: مدل SSNMF-Q

ورودی:

ماتریس مجاورتی  $A$

تعداد خوشه‌بندی  $k$

تعداد انجام حلقه  $I_t$

دانش اولیه از خوشه‌بندی

خروجی:

برچسب گروه برای هر گره

مراحل الگوریتم:

۱: مقداردهی به ماتریس  $H$

۲: محاسبه ماتریس‌های  $B_1, M, CB$ ، با توجه به رابطه (۱۷) و (۱)

۳: شروع حلقه با فرض  $I_t$  مرتبه

$$H_{ij} = H_{ij} \frac{4(AH)_{ij}}{4(B_1 H)_{ij} + 4(HH^T H)_{ij} + 2\alpha(B_2 HM)_{ij} + 2\beta(CH)_{ij}} \quad 4$$

$$H_{ir} := \frac{H_{ir}}{\sum_{r=1}^c H_{ir}} \quad 5$$

۶: انتهای حلقه

۷: برچسب گروه‌ها با رابطه  $(v_i, I_i) = \operatorname{argmax}_{r \leq k} H_{ir}^*$

<sup>3</sup> Ground truth partition labels

<sup>1</sup> Karush-Kuhn-Tucker

<sup>2</sup> Normalized Mutual Information



جدول ۱: اطلاعات مربوط به مجموعه شبکه‌های واقعی

توضیحات	$m$ (تعداد یال)	$n$ (تعداد گره)	شبکه‌ها
شبکه‌ای از مجموعه باشگاه‌های کاراته [۱۷، ۳۵]	۷۸	۳۴	کاراته
شبکه‌ای از مجموعه فوتبالی دانشگاه آمریکا [۱۷، ۳۶]	۶۱۳	۱۱۵	فوتبال
شبکه‌ای میان دلفین‌ها [۳۷، ۱۹]	۱۵۹	۶۲	دلفین
شبکه‌ای از مجموعه کتاب‌های سیاسی در آمریکا [۱۹، ۳۸]	۴۴۱	۱۰۵	کتاب‌های سیاسی
شبکه‌ای از مجموعه بلاگ‌های سیاسی در آمریکا [۲۴، ۳۹]	۳۳۴۳۰	۱۲۲۴	بلاگ‌های سیاسی

#### ۴-۲-۲ تنظیم پارامترهای داخلی

در روش‌های خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی که وابسته به انتخاب پارامترهای داخلی مدلسازان هستند، تنظیم و انتخاب پارامتر یکی از بخش‌های اساسی این الگوریتم‌ها است [۱۷، ۱۹، ۲۸]. از متداول‌ترین روش‌ها جهت تنظیم پارامترها، استفاده از روش سعی و خطا و انتخاب بهترین خوشه‌بندی با استفاده از معیارهای مانند چگالی پیمانگی [۱۷]، پیمانگی [۲۸] و اطلاعات متقابل نرمالیزه [۱۹] است. در مدل پیشنهادی SSNMF-Q،  $\alpha$  و  $\beta$  پارامترهای وزن دهی جهت کنترل و بهبود اثر دانش اولیه هم‌گروهی بودن و یا هم‌گروه نبودن گره‌ها است. همچنین این پارامترها می‌توانند در بهبود تابع بهینه‌سازی و رسیدن به بهترین پاسخ تأثیر بسزایی داشته باشند. با توجه به این موضوع، در صورتی که  $\alpha$  و  $\beta$  مقادیر بسیار بزرگی داشته باشند، در مدل پیشنهادی SSNMF-Q بخش خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریس بی‌اثر خواهد شد. از این‌رو بهترین خوشه‌بندی حاصل نخواهد شد. از طرفی در صورتی که  $\alpha$  و  $\beta$  مقادیر بسیار کوچک

#### ۴-۲ بررسی و ارزیابی الگوریتم‌های خوشه‌بندی نیمه نظارتی

در این بخش، ابتدا گراف‌های پیچیده و ساده برای ارزیابی معرفی خواهند شد. سپس الگوریتم SSNMF-Q با سایر الگوریتم‌های خوشه‌بندی نیمه نظارتی مبتنی بر تجزیه نامنفی ماتریسی با در نظر گرفتن رویکرد اضافه شدن عبارت دانش اولیه (رویکرد دوم در بخش ۲-۳) مقایسه خواهد شد.

یکی از بخش‌های اساسی و مهم برای بهره‌گیری از دانش اولیه خوشه‌بندی برخی از گره‌ها است. بدین منظور، ابتدا می‌توان جهت پیاده‌سازی و استفاده از دانش اولیه برای هر گراف غیرمستقیم،  $\frac{n(n-1)}{2}$  جفت گره را خوشه‌بندی کرد. این گره‌ها، به‌طور تصادفی با توجه به درصد‌های تعیین‌شده برای وارد کردن دانش اولیه (۱٪، ۵٪، ۱۰٪، ۲۰٪ و ۳۰٪) انتخاب می‌شوند. سپس، با توجه به دانش برجسب‌های خوشه‌بندی واقعی، جفت گره‌های انتخاب‌شده برحسب عضو بودن در هر گروه یا حضور در گروه‌های متفاوت خوشه‌بندی می‌شوند و ماتریس‌های تولیدشده بر اساس دانش اولیه گراف‌ها برای هر الگوریتم تشکیل می‌دهند.

#### ۴-۲-۱ مجموعه شبکه‌ها

در این مقاله، پنج مجموعه داده برای ارزیابی و صحت‌سنجی الگوریتم‌های خوشه‌بندی در نظر گرفته شده است. مجموعه داده انتخاب‌شده به‌گونه‌ای است که هر دو نوع داده ساده (شبکه کاراته) و پیچیده (سایر شبکه‌ها) مدنظر قرار گرفته شوند تا بتوان کارایی الگوریتم جدید پیشنهادی را برای هر دو نوع داده بررسی نمود. اطلاعات تکمیلی این پنج مجموعه داده در جدول ۱ نمایش و توضیح داده شده است.

در این جدول پارامترهای  $m$  و  $n$  به ترتیب، تعداد یال‌ها و گره‌های هر شبکه هستند که دانسته فرض می‌شوند. همان‌گونه که دیده می‌شود، شبکه‌های پیچیده‌تر دارای گره‌ها و یال‌های بیشتری هستند. لذا، با انتخاب این شبکه‌ها، می‌توان چگونگی عملکرد الگوریتم پیشنهادی را برای شبکه‌هایی با ابعاد مختلف نیز بررسی نمود.

به‌طور مثال مطابق شکل (۱)، نمونه‌ای از مقادیر مختلف معیارهای اندازه‌گیری برای مجموعه داده کتاب‌های سیاسی با فرض ۲۰٪ دانش اولیه در نظر گرفته شده است. همان‌طور که در شکل (۱) مشاهده می‌شود مقادیر متفاوت از پارامتر  $\alpha$  و  $\beta$  در کیفیت پاسخ تأثیر خواهند داشت. از این‌رو بهترین مقادیر پارامترهای  $\alpha$  و  $\beta$  برای مجموعه داده جدول ۱ با درصد دانش‌های اولیه متفاوت جهت استفاده در انجام آزمایش‌های نهایی در جدول ۲ ثبت شده است.

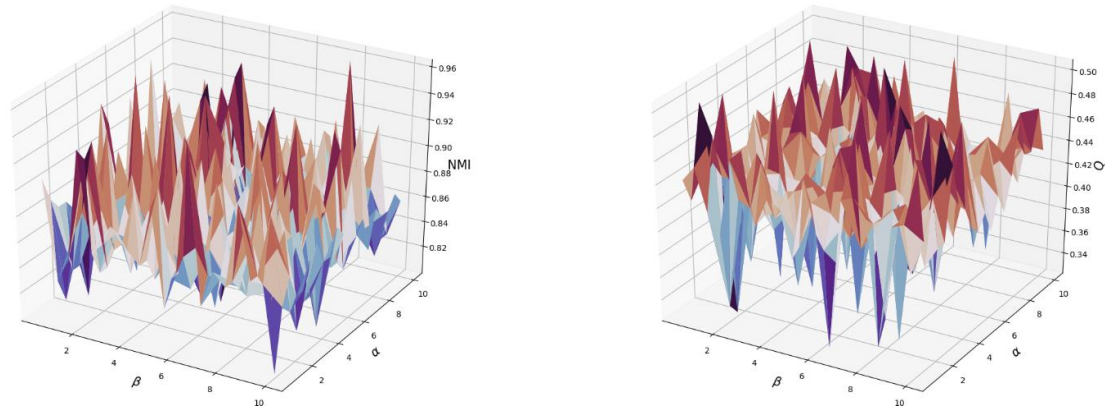
#### ۴-۲-۳ بررسی و ارزیابی الگوریتم‌های خوشه‌بندی

برای بررسی و ارزیابی الگوریتم‌های مختلف خوشه‌بندی، از الگوریتم‌های خوشه‌بندی PCNMF [۲۷]، PCSNMF [۲۸] و RSSNMF [۱۷]، به دلیل ترکیب خطی بخش دانش اولیه و بخش خوشه‌بندی تجزیه نامنفی ماتریسی استفاده شده است. روش PCNMF، روشی مینا برای ترکیب خطی دو ترم دانش اولیه و تجزیه نامنفی ماتریسی، روش PCSNMF روشی پایه جهت جداسازی انواع دانش اولیه (اعم از دانش هم‌گروه بودن یا هم‌گروه نبودن) و روش RSSNMF روشی جدید و مقاوم برای انواع دانش اولیه است. در میان الگوریتم‌های مختلف نیمه نظارتی، این روش‌ها جزو کارآمدترین آن‌ها هستند. از این‌رو مقایسه این روش‌ها با روش نوین SSNMF-Q، مقایسه با بهترین روش‌های نیمه نظارتی دوره زمانی‌های مختلف خواهد بود. برای پیاده‌سازی این الگوریتم‌ها، چنانچه الگوریتم برای شبکه‌های انتخاب شده در مرجع موردنظر شبیه‌سازی شده باشد، شبیه‌سازی با همان پارامترها تکرار شده است. در غیر این صورت، پارامترهای نامعلوم هر الگوریتم متناسب با مقادیر معین شده و با استفاده از روش‌های مطرح شده در مرجع موردنظر انتخاب شده است. از طرفی برای الگوریتم SSNMF-Q برای تنظیم پارامترها، متناسب با هر داده، پارامترهای موردنیاز مطابق با جدول ۲، تنظیم شده‌اند.

نزدیک به صفر را داشته باشند عملاً اثر دانش اولیه از گروه برخی از گره‌ها وارد الگوریتم خوشه‌بندی نخواهد شد. در نتیجه بررسی و ارزیابی بهترین مقدار برای  $\alpha$  و  $\beta$  می‌تواند منتج به بهترین خوشه‌بندی گره‌های شبکه‌های پیچیده شود. در نتیجه برای تنظیم پارامترهای مدل پیشنهادی SSNMF-Q، مطابق با روش سعی و خطا در [۱۷، ۱۹، ۲۸] مقادیر مختلف از  $\alpha$  و  $\beta$  را در بازه‌های مختلف قرار می‌دهیم سپس با توجه به معیار صحت‌سنجی، بهترین خوشه‌بندی را به‌عنوان خوشه‌های نهایی در نظر می‌گیریم. از این‌رو، برای تمامی مجموعه داده‌های جدول ۱، ابتدا آزمایش‌های متعددی با توجه به مقادیر مختلف  $\alpha$  و  $\beta$  در بازه ۰.۵ تا ۱۰ با فرض افزایش گام به مقدار ۰.۵ انجام شده و در هر مرحله مقادیر معیار اطلاعات متقابل و معیار پیمانگی، محاسبه و اندازه‌گیری شده است. سپس، یا انتخاب معیار اطلاعات متقابل به‌عنوان معیار اصلی جهت انتخاب بهترین مقدار پارامترهای  $\alpha$  و  $\beta$ ، این پارامترها مطابق با جدول ۲ برای انواع شبکه و درصد‌های دانش اولیه به دست می‌آیند.

جدول ۲: اطلاعات مربوط به مجموعه شبکه‌های واقعی

شبکه‌ها	پارامتر	درصد دانش‌های اولیه متفاوت				
		٪۱	٪۵	٪۱۰	٪۲۰	٪۳۰
کاراته	$\alpha$	۹.۵	۴.۵	۱	۳	۰.۵
	$\beta$	۷	۶	۷.۵	۳.۵	۴.۵
فوتبال	$\alpha$	۲	۳.۵	۵.۵	۱.۵	۵
	$\beta$	۳	۹	۱.۵	۹	۱۰
دلفین	$\alpha$	۶.۵	۲.۵	۹.۵	۲.۵	۱۰
	$\beta$	۷	۹	۱.۵	۰.۵	۷.۵
کتاب‌های سیاسی	$\alpha$	۶	۳.۵	۷.۵	۸	۲
	$\beta$	۹.۵	۱.۵	۲.۵	۸.۵	۸.۵
بلاگ‌های سیاسی	$\alpha$	۸.۵	۱.۵	۴	۵.۵	۷.۵
	$\beta$	۳.۵	۲.۵	۹.۵	۳.۵	۶.۵



شکل (۱): نمایشی از تغییرات  $\alpha$  و  $\beta$  با توجه به معیار پیمانگی (سمت چپ) و معیار اطلاعات متقابل (سمت راست) روی داده‌های کتاب‌های سیاسی

هیچ دانش اولیه‌ای همانند روش‌های دیگر بهترین عملکرد را از خود ارائه دهد.

- از مقایسه نتایج می‌توان دریافت که برای داده‌های پیچیده شامل جدول‌های ۴ تا ۷، الگوریتم پیشنهادی در این مقاله کارایی بهتر از سه الگوریتم دیگر دارد.
- همان‌طور که مشاهده می‌شود، الگوریتم SSNMF-Q الگوریتمی است که با بهره‌ای کمتر از دانش اولیه توانسته است به حداکثر مقدار معیار اطلاعات متقابل برسد. در واقع می‌توان ادعا کرد که به‌صورت میانگین الگوریتم ما با بهره کمتر از دانش اولیه توانسته است خوشه‌بندی بهتری نسبت به دیگر الگوریتم‌های به‌روز دنیا ارائه دهد. همچنین الگوریتم SSNMF-Q نشان داده است که افزایش دانش اولیه باعث بهبود خوشه‌بندی بر پایه معیار پیمانگی می‌شود.
- انتخاب ماتریس پیمانگی به‌عنوان ماتریس مشابهت گراف و انتخاب بهترین پارامتر  $\alpha$  و  $\beta$  باعث بهبود پاسخ خوشه‌بندی نسبت به الگوریتم‌های چون RSSNMF و PCSNMF شده است.

لذا، با در نظر گرفتن مقادیر مشخص‌شده برای این پارامترها و با توجه به جدول‌های ۳ تا ۷، نتایج زیر به دست می‌آید.

- در تمامی الگوریتم‌ها در جدول‌های ۴ تا ۷ با افزودن دانش اولیه با درصد‌های متفاوت، خوشه‌بندی مناسب‌تری از لحاظ دستیابی به خوشه‌بندی واقعی<sup>۱</sup> یا همان بهتر شدن معیار اطلاعات متقابل اتفاق افتاده است؛ اما در برخی موارد مانند جدول ۵، ۶ و ۷، خوشه‌بندی نسبت به معیار پیمانگی افت داشته است. علت این امر آن است که خوشه‌بندی‌های نیمه نظارتی مبتنی بر تجزیه نامنفی ماتریسی با وارد کردن دانش اولیه هر گره سعی در بهبود و رسیدن به خوشه‌بندی واقعی خواهند داشت اما، نوع و ویژگی‌های گراف را در نظر نمی‌گیرند. پس بدیهی است که خوشه‌بندی که بر اساس معیاری مبتنی بر خوشه‌بندی واقعی هستند (مانند معیار اطلاعات متقابل) بهبود کیفیت داشته باشد و هم‌چنین ممکن است که خوشه‌بندی که بر اساس معیارهای ویژگی گراف می‌باشند (مانند معیار پیمانگی) دچار افت کیفیت شوند.
- الگوریتم ارائه‌شده در این مقاله برای مجموعه داده ساده مانند کاراته در جدول ۳ توانسته است بدون استفاده از

<sup>1</sup> Ground truth partition labels

جدول ۳: نتایج الگوریتم‌ها بر روی مجموعه کاراته

معیار اطلاعات متقابل				معیار پیمانی				درصد
SSNMF-Q	RSSNMF [۱۷]	PCSNMF [۲۸]	PCNMF [۲۷]	SSNMF-Q	RSSNMF [۱۷]	PCSNMF [۲۸]	PCNMF [۲۷]	
۱	۱	۱	۱	۰.۳۷۱	۰.۳۷۱	۰.۳۷۱	۰.۳۷۱	%۱
۱	۱	۱	۱	۰.۳۷۱	۰.۳۷۱	۰.۳۷۱	۰.۳۷۱	%۵
۱	۱	۱	۱	۰.۳۷۱	۰.۳۷۱	۰.۳۷۱	۰.۳۷۱	%۱۰
۱	۱	۱	۱	۰.۳۷۱	۰.۳۷۱	۰.۳۷۱	۰.۳۷۱	%۲۰
۱	۱	۱	۱	۰.۳۷۱	۰.۳۷۱	۰.۳۷۱	۰.۳۷۱	%۳۰

جدول ۴: نتایج الگوریتم‌ها بر روی مجموعه فوتبالیست

معیار اطلاعات متقابل				معیار پیمانی				درصد
SSNMF-Q	RSSNMF [۱۷]	PCSNMF [۲۸]	PCNMF [۲۷]	SSNMF-Q	RSSNMF [۱۷]	PCSNMF [۲۸]	PCNMF [۲۷]	
۰.۹۲۰	۰.۹۰۴	۰.۸۳۰	۰.۸۷۸	۰.۵۴۰	۰.۵۳۶	۰.۵۰۷	۰.۵۲۸	%۱
۰.۹۳۵	۰.۹۳۰	۰.۸۷۸	۰.۹۲۰	۰.۵۸۵	۰.۵۸۰	۰.۵۲۸	۰.۵۴۰	%۵
۰.۹۶۵	۰.۹۵۶	۰.۹۰۴	۰.۹۳۰	۰.۵۹۴	۰.۵۸۷	۰.۵۴۷	۰.۵۸۰	%۱۰
۰.۹۷۳	۰.۹۷۳	۰.۹۳۰	۰.۹۵۶	۰.۶۰۱	۰.۶۰۱	۰.۵۸۰	۰.۵۸۷	%۲۰
۰.۹۷۳	۰.۹۷۳	۰.۹۶۵	۰.۹۷۳	۰.۶۰۱	۰.۶۰۱	۰.۵۹۴	۰.۶۰۱	%۳۰

جدول ۵: نتایج الگوریتم‌ها بر روی مجموعه دلفین

معیار اطلاعات متقابل				معیار پیمانی				درصد
SSNMF-Q	RSSNMF [۱۷]	PCSNMF [۲۸]	PCNMF [۲۷]	SSNMF-Q	RSSNMF [۱۷]	PCSNMF [۲۸]	PCNMF [۲۷]	
۰.۹۸۳	۰.۹۸۳	۰.۹۵۱	۰.۹۶۷	۰.۳۷۸	۰.۳۷۸	۰.۳۸۹	۰.۳۸۴	%۱
۱	۱	۰.۹۶۷	۰.۹۸۳	۰.۳۷۳	۰.۳۷۳	۰.۳۸۴	۰.۳۷۸	%۵
۱	۱	۰.۹۸۳	۱	۰.۳۷۳	۰.۳۷۳	۰.۳۷۸	۰.۳۷۳	%۱۰
۱	۱	۱	۱	۰.۳۷۳	۰.۳۷۳	۰.۳۷۳	۰.۳۷۳	%۲۰
۱	۱	۱	۱	۰.۳۷۳	۰.۳۷۳	۰.۳۷۳	۰.۳۷۳	%۳۰

جدول ۶: نتایج الگوریتم‌ها بر روی مجموعه کتاب‌های سیاسی

معیار اطلاعات متقابل				معیار پیمانی				درصد
SSNMF-Q	RSSNMF [۱۷]	PCSNMF [۲۸]	PCNMF [۲۷]	SSNMF-Q	RSSNMF [۱۷]	PCSNMF [۲۸]	PCNMF [۲۷]	
۰.۴۹۲	۰.۴۴۸	۰.۴۰۴	۰.۴۳۰	۰.۴۰۳	۰.۳۹۵	۰.۵	۰.۴۳۸	٪۱
۰.۴۹۶	۰.۴۷۰	۰.۴۱۷	۰.۴۴۰	۰.۳۳۵۷	۰.۴۰۵	۰.۴۸۸	۰.۳۵۹	٪۵
۰.۵۰۲	۰.۴۹۱	۰.۴۲۵	۰.۴۴۸	۰.۳۳۸	۰.۳۹۷	۰.۴۵۴	۰.۳۴۴	٪۱۰
۰.۵۱۲	۰.۵۰۲	۰.۴۳۰	۰.۴۶۵	۰.۳۳۸	۰.۳۳۱	۰.۴۰۳	۰.۳۵۷	٪۲۰
۰.۵۱۲	۰.۵۰۲	۰.۴۴۸	۰.۴۷۶	۰.۳۲۵	۰.۳۲۵	۰.۳۹۷	۰.۳۳۱	٪۳۰

جدول ۷: نتایج الگوریتم‌ها بر روی مجموعه بلاگ‌های سیاسی

معیار اطلاعات متقابل				معیار پیمانی				درصد
SSNMF-Q	RSSNMF [۱۷]	PCSNMF [۲۸]	PCNMF [۲۷]	SSNMF-Q	RSSNMF [۱۷]	PCSNMF [۲۸]	PCNMF [۲۷]	
۰.۹۷۳	۰.۹۶۸	۰.۹۳۶	۰.۹۵۹	۰.۴۱۹	۰.۴۲۳	۰.۴۲۶	۰.۴۲۲	٪۱
۰.۹۸۴	۰.۹۷۳	۰.۹۴۳	۰.۹۶۴	۰.۴۱۷	۰.۴۱۹	۰.۴۲۴	۰.۴۲۳	٪۵
۰.۹۸۴	۰.۹۸۴	۰.۹۵۹	۰.۹۷۲	۰.۴۱۷	۰.۴۱۷	۰.۴۲۲	۰.۴۲۳	٪۱۰
۰.۹۸۴	۰.۹۸۴	۰.۹۷۱	۰.۹۸۴	۰.۴۱۷	۰.۴۱۷	۰.۴۲۳	۰.۴۱۷	٪۲۰
۰.۹۸۴	۰.۹۸۴	۰.۹۸۴	۰.۹۸۴	۰.۴۱۷	۰.۴۱۷	۰.۴۱۷	۰.۴۱۷	٪۳۰

نمی‌گیرد. از این رو ما در این مقاله اثبات کرده‌ایم که خوشه‌بندی مبتنی بر معیار پیمانی دارای ساختاری مشابه با خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی است. بر این اساس ماتریس پیمانی به‌عنوان ماتریس مشابهت گراف در روش خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی در نظر گرفته شده است و به‌عنوان روش نیمه نظارتی برای خوشه‌بندی مبتنی بر معیار پیمانی ارائه شده است. سپس با توسعه الگوریتم SSNMF-Q امکان ترکیب دانش اولیه از اعضای گروه‌ها با روش مشابه‌سازی شده ایجاد گشته و برای حل بهینه‌سازی، روش حل تکراری پیشنهاد و محاسبه شده است. در نهایت نتایج خروجی الگوریتم‌ها بر روی پنج مجموعه

## ۵- نتیجه‌گیری

خوشه‌بندی مبتنی بر معیار پیمانی یکی از روش‌های متداول برای خوشه‌بندی گراف‌ها است. از این رو در سال‌های گذشته، روش‌های مختلفی برای استخراج گروه‌های گراف مطرح و بررسی شده است. محدودیت‌های روش خوشه‌بندی مبتنی بر پیمانی، انپی-سخت بودن حل مسئله و عدم استفاده از دانش اولیه از اعضای گروه‌های برخی از گره‌ها است. از طرفی، یکی از روش‌های رایج جهت خوشه‌بندی نیمه نظارتی، روش خوشه‌بندی مبتنی بر تجزیه نامنفی ماتریسی است. اما این روش، ویژگی‌های خاص شبکه‌ها را در نظر



- [15] R. Ghadirian and N. Bigdeli, "Hybrid Adaptive Modularized Tri-Factor Non-Negative Matrix Factorization for Community Detection in Complex Networks", *Scientia Iranica*, Vol. 14, 2022.
- [16] Z. Liu, G. Yuan and X. Luo, "Symmetry and Nonnegativity-Constrained Matrix Factorization for Community Detection," *IEEE/CAA Journal of Automatica Sinica*, vol. 9, no. 9, pp. 1691-1693, 2022.
- [17] C. He, Q. Zheng, Y. Tang, S. Liu, J. Zheng, "Community detection method based on robust semi-supervised nonnegative matrix factorization", *Physica A: Statistical Mechanics and its Applications*, Vol. 523, no. 1, pp. 279 – 291, 2019.
- [18] C. He, Y. Tang, K. Liu, H. Li and S. Liu, "A robust multi-view clustering method for community detection combining link and content information", *Physica A: Statistical Mechanics and its Applications* Vol. 514, pp. 396-411, 2018.
- [19] C. Yan and Z. Chang, "Modularized tri-factor nonnegative matrix factorization for community detection enhancement", *Physica A: Statistical Mechanics and its Applications*, Vol. 533, no. 122050, 2019.
- [20] P. M. Zheng and Z. Zhou, "Structural Deep Nonnegative Matrix Factorization for community detection", *Applied Soft Computing*, Vol. 97, no: B, Issue: 106846, 2020.
- [21] D. Handshutter, N. Gillis and X. Seibert, "A survey on deep matrix factorizations", *Computer Science Review*, Vol. 42, Issue: 100423, 2021.
- [22] J. Huang, T. Zhang, W. Yu, J. Zhu and E. Cai, "Community Detection Based on Modularized Deep Nonnegative Matrix Factorization", *International Journal of Pattern Recognition and Artificial Intelligence*, Vol. 35, no. 35, Issue: 2159006, 2021.
- [23] H. jin and S. Li, "Graph regularized nonnegative matrix tri-factorization for overlapping community detection", *Physica A: Statistical Mechanics and its Applications*, Vol. 515, pp. 376-387, 2019.
- [24] C. Chen, W. Zho and B. Peng, "Differentiated graph regularized non-negative matrix factorization for semi-supervised community detection", *Physica A: Statistical Mechanics and its Applications*, Vol. 604, no. 127692, 2022.
- [25] Z.-Y. Zhang, "Community structure detection in complex networks with partial background information", *Europhysics Letters*, Vol. 101, no. 4, pp. 48005, 2013.
- [26] X. Ma, L. Gao, X. Yong, and L. Fu, "Semi-supervised clustering algorithm for community structure detection in complex networks", *Physica A: Statistical Mechanics and its Applications*, Vol. 389, no. 1, pp. 187 – 197, 2010.
- [27] Y. Yang and B. Hu, "Pairwise constraints-guided non-negative matrix factorization for document clustering", in *Web Intelligence, IEEE/WIC/ACM International Conference on*, IEEE, pp. 250– 256, 2007.
- [28] X.H. Shi, H.T. Lu, Y.C. He and S. He, "Community detection in social network with pairwise constrained symmetric non-negative matrix factorization", *Proceedings of the 7th IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, pp. 541–546, 2015.

دادهٔ مختلف استخراج و بر اساس نتایج آن، کارایی و تأثیرگذاری الگوریتم پیشنهادی بررسی شده است.

## References

- [1] S. Kumar and R. Hanot, "Community Detection Algorithms in Complex Networks: A Survey", *Advances in Signal Processing and Intelligent Recognition Systems*, Vol. 1365, no. 202, 215, 2021.
- [2] K. D. Asim, T. Yahui and G. R. Yulia, "Community detection in complex networks: From statistical foundations to data science applications", *WIREs Computational Statistics*, Vol. 14, no. 2, 2021.
- [3] M.E.J. Newman, "Networks", Oxford university press, 2018.
- [4] M.E.J. Newman and M. Girvan, "Finding and evaluating community structure in networks", *Physical Review E*, Vol. 69, no. 2, Issue: 026113, 2004.
- [5] N. Fariahag, M. Mordi, Z. J. Wang, "Community structure detection from networks with weighted modularity", *Pattern Recognition Letters*, Vol. 122, pp. 14-22, 2019.
- [6] T. Li, X. Wang, S.H. Zhu, S.H. Zhu and C. Ding, "Community discovery using nonnegative matrix factorization", *Data Min. Knowl. Discovery*, Vol. 22 no. 3, pp. 493–521, 2011.
- [7] C. Li, H. Chen and T. Li, "A stable community detection approach for complex network based on density peak clustering and label propagation", *Applied Intelligence*, Vol. 52pp. 1188-1208, 2022.
- [8] T. Wang, S. Chen, X. Wang and J. Wang, "Label propagation algorithm based on node importance", *Physica A: Statistical Mechanics and its Applications*, Vol. 551, no. 124137, 2020.
- [9] M. Rosvall and C.T. Bergstrom, "Maps of random walks on complex networks reveal community structure", *Proceedings of the National Academy of Sciences*, Vol. 105, no.4, pp. 1118–1123, 2008.
- [10] J. Zhou, L. Li, A. Zeng, Y. Fan and Z. Di, "Random walk on signed networks", *Physica A: Statistical Mechanics and its Applications*, Vol. 508, pp. 558-556, 2018.
- [11] R. Shang, K. Zhao, W. Zhang, J. Feng, Y. Li and L. Jiao, "Evolutionary multiobjective overlapping community detection based on similarity matrix and node correction", *Applied Soft Computing*, Vol. 127 no. 109397, 2022.
- [12] J. Sanchez and A. Duarte, "Iterated Greedy algorithm for performing community detection in social networks", *Future Generation Computer Systems*, Vol. 88, pp. 785-791, 2018.
- [13] M. Guerrero, F. G. Montoya, R. Baños, A. Alcayde and C. Gil, "Adaptive community detection in complex networks using genetic algorithms", *Neurocomputing*, Vol. 266, pp. 101-113, 2017.
- [14] S. Fortunato and M. Barthlemy, "Resolution limit in community detection", *Proceedings of the National Academy of Sciences*, Vol. 104, no. 1, pp. 36–41, 2007.



- [29] X. Liu, W. Wang, D. He, P. Jiao, D. Jin and C. Vittorio, "Semi-supervised community detection based on non-negative matrix factorization with node popularity", *Information Sciences*, Vol. 381, no. 12, pp. 304–321, 2017.
- [30] C. Févotte, E. Vincent and A. Ozerov. " Single-channel audio source separation with NMF: divergences, constraints and algorithms ", *Audio Source Separation*, Springer, hal-01631185f, pp. 1-24, 2018.
- [31] G. A. Khan, J. Hu, T. li, B. Diallo and H. Wang, " Multi-view data clustering via non-negative matrix factorization with manifold regularization", *International Journal of Machine Learning and Cybernetics*, Vol. 13, pp. 677-689, 2022.
- [32] S. Peng, W. Ser, B. Chen and Z. Lin, " Robust semi-supervised nonnegative matrix factorization for image clustering", *Pattern Recognition*, Vol. 111, no. 107683, 2021.
- [33] K. Huang, X. Fu and N.D. Sidiropoulos, "Anchor-free correlated topic modeling: Identifiability and algorithm", *Advances in Neural Information Processing Systems*, pp. 1794-1802, 2016.
- [34] X. Ma, D. Dong, Q. Wang, "Community detection in multi-layer networks using joint nonnegative matrix factorization", *IEEE Transaction on Knowledge and Data Engineering*. 31 (2): 273-286, 2019.
- [35] W.W. Zachary, "An information flow model for conflict and fission in small groups", *Journal of anthropological research*, Vol. 33, no. 4, pp. 452–473, 1977.
- [36] A. Lancichinetti, S. Fortunato and F. Radicchi, "Benchmark graphs for testing community detection algorithms", *Physical Review E*, Vol. 78, no. 4, Issue: 046110, 2008.
- [37] D. Lusseau, K. Schneider, O.J. Boisseau, P. Haase, E. Slooten and S.M. Dawson, "The bottlenose dolphin community of doubtful sound features a large proportion of long-lasting associations", *Behavioral ecology and sociobiology*, Vol. 54, no. 4, pp. 396–405, 2003.
- [38] J. Kunegis. KONECT: "The Koblenz Network Collection", *Proceedings of the 22nd international conference on World Wide Web companion*, pp. 1343–1350, 2013.
- [39] L. A. Adamic, N. Glance, "The political blogosphere and the 2004 US election: divided the blog", *Proceedings of the 3rd workshop on Link discovery*, ACM, pp. 36–43, 2005.

A novel semi-supervised clustering method for complex network based on modularity

Mohammad Ghadirian<sup>1</sup>, Nooshin Bigdeli<sup>1\*</sup>

<sup>1</sup> Control Engineering Department, Imam Khomeini International University, Qazvin, Iran

Article Information

Original Research Paper

Received:

2022 November 15

Accepted:

2022 January 9

Keywords:

Non-negative matrix factorization, community detection, semi-supervised clustering, modularity criterion

Corresponding Author\*:

n.bigdeli@eng.ikiu.ac.ir

Abstract

Clustering is a powerful tool for analyzing complex networks which is widely used for modeling complex systems. Modularity is a comprehensive criterion for evaluating the quality of clusters. However, it has some limitations and challenges such as being a NP-hard problem and not using prior information. So, Modularity-based community detection cannot be extended as a semi-supervised community detection method. On the other hand, one of the most common semi-supervised methods which can use prior knowledge for clustering is community detection based on non-negative matrix factorization (NMF). But, this method is not able to consider the features of the networks. Therefore, in this paper to overcome the mentioned limitations and challenges and by presenting a new proof, a structure similar to community detection based on NMF is presented for modularity-based community detection which can employ prior knowledge and iterative solution. Therefore, a novel semi-supervised community detection based on modularity (SSNMF-Q) criteria is developed by utilizing prior information and iterative solution instead of solving a NP-hard problem. To evaluate SSNMF-Q, five real world networks are used and it is shown that the SSNMF-Q had better performance compared to other semi-supervised community detection methods based on NMF.

 : 10.22034/ABMIR.2023.19231.1018