

# ترکیب یک روش خوشه‌بندی تجمعی و یک

## معیار شباهت جدید برای مدل‌سازی

### رفتار وراثتی بیماری‌ها

موسی مجرد\*<sup>۱،۲</sup>، حمید پروین<sup>۳</sup>، صمد نجاتیان<sup>۴</sup> و کرم‌الله باقری فرد<sup>۵،۶</sup>

<sup>۱</sup>گروه کامپیوتر، واحد فیروزآباد، دانشگاه آزاد اسلامی، فیروزآباد، ایران

<sup>۲</sup>باشگاه پژوهش‌گران جوان، واحد فیروزآباد، دانشگاه آزاد اسلامی، فیروزآباد، ایران

<sup>۳</sup>گروه کامپیوتر، واحد نورآباد ممسنی، دانشگاه آزاد اسلامی، نورآباد ممسنی، ایران

<sup>۴</sup>باشگاه پژوهش‌گران جوان، واحد نورآباد ممسنی، دانشگاه آزاد اسلامی، نورآباد ممسنی، ایران

<sup>۵</sup>گروه برق، واحد یاسوج، دانشگاه آزاد اسلامی، یاسوج، ایران

<sup>۶</sup>باشگاه پژوهش‌گران جوان، واحد یاسوج، دانشگاه آزاد اسلامی، یاسوج، ایران

<sup>۷</sup>گروه کامپیوتر، واحد یاسوج، دانشگاه آزاد اسلامی، یاسوج، ایران

#### چکیده

امروزه تئوری‌های بسیاری در مورد علل بروز بیماری‌های وراثتی وجود دارد، اما پزشکان معتقدند که دو فاکتور ژنتیک و محیط زیست هم‌زمان باهم نقش مهمی در بروز و پیشرفت این بیماری‌ها ایفا می‌کنند، هرچند که چگونگی این اثرگذاری هنوز به‌طور دقیق مشخص نیست. برای اینکه بتوان ژن‌های مؤثر در بروز بیماری‌ها را تشخیص داد، باید ارتباط بین سلول‌ها/بافت‌ها را به‌دست آورد. تعامل بین سلول‌ها یا بافت‌های مختلف را می‌توان، با بیان ژن بین آنها نشان داد. با نمونه‌برداری از کروموزوم‌ها، اطلاعات مفیدی در مورد نوع بیماری و چگونگی انتقال آن استخراج می‌شود. با بررسی این اطلاعات می‌توان اختلالاتی را که منجر به تغییرات به‌شدت پرتکراری شده‌اند شناسایی کرد. در این مقاله تشخیص ارتباط‌های بین سلولی و بین بافتی در بیماری‌های مختلف با توجه به مشخصات ساختار توپولوژیکی گراف و یک روش خوشه‌بندی تجمعی بهبودیافته انجام شده است. روش پیشنهادی دو مرحله دارد؛ در مرحله نخست چندین مدل خوشه‌بندی به‌منظور تشخیص ارتباط‌های اولیه بین سلول‌ها یا بافت‌ها در جهت تولید نتایج بهتر نسبت به الگوریتم‌های انفرادی، ترکیب می‌شوند. در مرحله دوم تشابه بین سلول‌ها یا بافت‌ها در هر خوشه با استفاده از یک معیار شباهت مبتنی بر ساختار توپولوژیکی گراف محاسبه و درنهایت از بیشینه شباهت‌های بین سلول‌ها یا بافت‌ها در هر خوشه برای کشف ارتباطات بین بیماری‌ها استفاده می‌شود. به‌منظور ارزیابی عملکرد روش پیشنهادی از چندین مجموعه داده UCI و همچنین مجموعه داده فانتوم پنج استفاده شده است. نتایج روش پیشنهادی روی مجموعه داده فانتوم پنج، ضریب سیلوئت ۰/۹۰۱ را در ۱۸ خوشه برای سلول‌ها و ۰/۷۶۲ در ۱۳ خوشه برای بافت‌ها را گزارش می‌کند.

واژگان کلیدی: ارتباط بین سلولی، خوشه‌بندی بهبودیافته، ساختار توپولوژیکی گراف، مجموعه داده فانتوم ۵.

## Combining an Ensemble Clustering Method and a New Similarity Criterion for Modeling the Hereditary Behavior of Diseases

Musa Mojarad<sup>\*1,2</sup>, Hamid Parvin<sup>3,4</sup>, Samad Nejatian<sup>5,6</sup> & Karamollah Bagheri Fard<sup>6,7</sup>

<sup>1</sup>Department of Computer Engineering, Firoozabad Branch, Islamic Azad University, Firoozabad, Iran

<sup>2</sup>Young Researchers Club, Firoozabad Branch, Islamic Azad University, Firoozabad, Iran

<sup>3</sup>Department of Computer Engineering, Nourabad Mamasani Branch, Islamic Azad University, Nourabad Mamasani, Iran.

\* Corresponding author

\* نویسنده عهده‌دار مکاتبات

سال ۱۴۰۰ شماره ۲ پیاپی ۴۸

تاریخ ارسال مقاله: ۱۳۹۹/۱۲/۷ • تاریخ پذیرش: ۱۳۹۹/۵/۲۸ • تاریخ انتشار: ۱۴۰۰/۰۷/۱۷ • نوع مطالعه: کاربردی



## Abstract

**Background:** There are many theories about the causes of hereditary diseases, but physician believe that both the genetic and environmental factors simultaneously play an important role in the development and progression of these diseases, although the extent to which this effect is not yet clear. In order to detect effective genes in the development of diseases, it is necessary to achieve the relationship between cells/tissues.

**Objective:** In fact, inter-cell or inter-tissue communications indicate the hereditary relationships between patients. Detecting these communications help to identify common parts of the body that are influenced by various diseases. The interaction between different cells/tissues can be demonstrated by expressing the gene between them. By sampling chromosomes, useful information is obtained about the type of disease and how it is transmitted. By examining this information, you can identify disorders that have led to highly altered changes. In previous research, various clustering methods have been used to discover the links between diseases based on gene expression data. However, ensembl clustering approaches have not yet been used for this purpose.

**Method:** In this paper, the recognition of intercellular and inter-tissue interactions in various diseases have been done according to the characteristics of the topological structure of the graph and an improved ensembl clustering method. The proposed clustering algorithm uses an agreed similarity function to measure the similarity between objects. The proposed method has two stages; in the first step, several clustering models are combined to identify the initial relationships between cells or tissues in order to produce better results than individual algorithms. In the second stage, the similarity between cells or tissues in each cluster is calculated by using a similarity criterion based on the topological structure of the graph. Eventually, the maximum similarity between cells or tissues in each cluster is used to discover the relationship between diseases. In addition, an algorithm for improving the uncertainty of objects is evaluated by allocating them to other clusters in order to enhance the quality of the final clusters.

**Results:** To evaluate the performance of the proposed method, several UCI datasets and the FANTOM5 dataset have been used. The results of the proposed method on the phantom data set 5 report a silhouette of 0.901 in 18 clusters for cells and 0.762 in 13 clusters for tissues.

**Conclusion:** The conducted evaluations have confirmed the power of the proposed clustering algorithm in terms of accuracy. Clustering of cells or tissues has increased the accuracy and concentration of the topological similarity criterion of the graph in the range of similarity of cells or tissues.

**Keywords:** Intercellular communication, Improved clustering, Graph topological structure, FANTOM5 dataset.

اطلاق می‌شود. ژنوم انسان به مجموعه ژنتیکی و ژن‌های داخل هسته سلول‌های انسان گفته می‌شود [2]. ژن‌های موجود در یک کروموزوم حدود چند میلیون است. بیش‌تر این ژن‌ها در تمام انسان‌ها یکسان و ثابت هستند (Fix Gen). به این نوع ژن‌ها دخالتی در تغییرات ژنتیکی و سلولی/بافتی ندارند [3]. تغییرات این نوع ژن‌ها باعث ایجاد ناپایداری‌هایی در انسان خواهد شد (به‌عنوان مثال تولید انسان‌هایی با چهار چشم)؛ بنابراین این نوع ژن‌ها نباید در بدن انسان تغییر کنند؛ اما بخش کوچکی از ژن‌ها ثابت نیستند و مقدار آن‌ها در انسان‌های مختلف متغیر هستند (Variables Gen). متغیربودن این ژن‌ها باعث می‌شود انسان‌های مختلفی وجود داشته باشد. هر ژن در روی کروموزوم مکان معینی را اشغال می‌کند که در اصطلاح به آن جایگاه ژن (Locus) گفته می‌شود [4].

## ۱- مقدمه

بدن انسان از انواع مختلفی سلول و بافت تشکیل شده است. تعداد سلول‌های بدن انسان حدود  $37/2$  تریلیون تخمین زده شده است. سلول‌ها و بافت‌های بدن به‌طور مستقیم یا غیرمستقیم به یکدیگر وابسته هستند و توانایی ارتباط و تأثیر بر روی یکدیگر را دارند؛ بنابراین سازوکاری مؤثر و کارآمد به‌منظور یافتن ارتباط بین این تعداد نجومی از سلول‌ها و بافت‌ها مورد نیاز است. یافتن این ارتباطات به شناسایی بیماری‌ها مختلف و عوامل آن کمک خواهد کرد. سلول‌ها و بافت‌های بیماری می‌تواند بر اثر جهش ایجاد شوند [1]. هسته هر سلول دستورالعمل‌های کدگذاری شده لازم را برای هدایت فعالیت‌های سلول و ساخت پروتئین‌های دارد. به یک گروه کامل از این دستورالعمل‌های خام ژنوم

در بخش ۵ به شرح آزمایش‌های انجام‌شده روی مجموعه داده‌های UCI و همچنین پایگاه داده فانتوم می‌پردازیم. در نهایت نتیجه‌گیری و پیشنهادهای آینده در بخش ۶ ارائه می‌شود.

## ۲- خوشه‌بندی تجمعی و تئوری گراف

خوشه‌بندی یکی از شاخه‌های یادگیری بدون نظارت است که هدف آن یافتن خوشه‌های مشابه از اشیا در بین نمونه‌های ورودی است. می‌توان نشان داد که هیچ معیار مطلق برای بهترین خوشه‌بندی وجود ندارد، بلکه این بستگی به مسأله و نظر کاربر دارد؛ باین‌حال معیارهای مختلفی برای خوب‌بودن یک خوشه‌بندی ارائه شده است که می‌تواند کاربر را برای رسیدن به یک خوشه‌بندی مناسب راهنمایی کند. در این مقاله از یک خوشه‌بندی تجمعی بهبودیافته به‌منظور قراردادن سلول‌ها یا بافت‌های مشابه بیماری‌های مختلف در خوشه‌های یکسانی استفاده می‌شود. خوشه‌بندی تجمعی با یک خوشه به‌ازای هر مشاهده شروع شده و در هر مرحله، دو خوشه با بیشترین شباهت ادغام شده و یک خوشه بزرگ‌تر را تشکیل می‌دهند. این کار باعث می‌شود تا در سطوح بالاتر خوشه‌های کمتری وجود داشته باشد [9]. در اینجا از چند مدل مختلف خوشه‌بندی به‌منظور تولید نتایج بهتر نسبت به الگوریتم‌های انفرادی استفاده می‌شود. برای یک مجموعه داده با  $n$  نمونه به‌صورت  $X = \{x_1, x_2, \dots, x_n\}$ ،  $P_q = \{c_1^q, c_2^q, \dots, c_k^q\}$  نتیجه  $q$ -مین روش خوشه‌بندی با  $k$  خوشه است. خوشه‌بندی  $\Phi$  می‌تواند به‌وسیله  $m$  روش به‌صورت  $\Gamma = \{P_1, P_2, P_3, \dots, P_m\}$  و تابع توافقی  $F$  ساخته شود. رابطه (۱) فرم کلی تابع هدف برای خوشه‌بندی  $\Phi$  را نشان می‌دهد:

$$\Phi(F, \Gamma)P_1, P_2, P_3, \dots, P_m = F(\Gamma) \quad (1)$$

وظیفه یک خوشه‌بندی تجمعی یافتن افراز  $P^*$  از مجموعه داده  $X$  به‌وسیله ترکیب اعضای  $P$  و ارائه یک تابع مؤثر بدون دسترسی به ویژگی‌های اصلی است [10]. بنابراین  $P^*$  به‌احتمال از نظر سازگاری و کیفیت بهتر از اعضای انفرادی در تجمع است. این پژوهش در نظر دارد با توجه به مشخصات ساختار توپولوژی گراف، یک معیار برای محاسبه شباهت بین سلول‌ها/بافت‌ها کاربران ارائه دهد. شباهت بین دو سلول و یا دو بافت احتمال وجود ارتباط بین آن‌ها را نشان می‌دهد که بر اساس خصوصیت‌های موجودیت‌ها و سایر پیوندهای مشاهده‌شده انجام می‌گیرد [11].

همان‌طور که گفته شد، تمامی کروموزوم‌های یک بخش بدن یکسان هستند. به‌عنوان مثال تمام کروموزوم‌های دستگاه دفاعی بدن در تمام انسان‌ها کروموزوم  $\Gamma$  است. در ابتدا این کروموزوم برای همه انسان‌ها به‌طور کامل یکسان بودند، ولی در اثر حوادث و بروز بیماری‌های مختلف، در گذر زمان دچار تغییراتی شده‌اند. در یک انسان ممکن است، یکی از کروموزوم‌های دستگاه دفاعی بدن در یکی از سلول‌ها یا بافت‌ها دچار بیماری‌های ویروسی شده باشد. خصوصیات این ویروس باعث تغییرات و انتقال ویژگی‌های آن به سایر ژن‌های همسایه می‌شود. در نتیجه این تغییرات باعث تغییر در DNA آن سلول/بافت شده و اختلالاتی را ایجاد می‌کند [5].

با انجام نمونه‌برداری‌هایی روی این کروموزوم‌ها باعث استخراج اطلاعات مفیدی در مورد نوع بیماری و چگونگی انتقال آن به دست می‌آید. با بررسی کروموزوم‌های نمونه‌برداری‌شده از این سلول‌ها و بافت‌ها، می‌توانیم اختلالات ناشی از عملکرد ویروس را که منجر به تغییرات به‌شدت پرتکراری می‌شوند، شناسایی کنیم؛ اما چالش پیشرو شباهت بخش زیادی از این سلول‌ها/بافت‌ها و عدم توجه به تغییرات و الگوها در بخش‌های مختلف است. هدف این پژوهش ارائه یک روش خوشه‌بندی بهبودیافته به‌منظور شناسایی ارتباطات بین سلولی و بین بافتی با توجه به بیماری‌های مختلف است. وجود بیماری‌های مختلف باعث ایجاد ارتباط‌های چندوجهی بین سلول‌ها/بافت‌ها می‌شود که در نهایت این ارتباطات یک هایپر گراف را تشکیل می‌دهد [6]. به‌منظور شناسایی تغییرات ایجادشده در سلول‌ها/بافت‌ها، ارتباطات بین هر سلول/بافت را به‌صورت یک یال با وزن خاصی در گراف ایجاد می‌شود. در اینجا با استخراج اطلاعات توپولوژیکی گراف و استفاده از یک معیار شباهت ترکیبی احتمال مرتبط‌بودن ژن‌ها با توجه به مجموعه‌ای از بیماری‌ها مشخص می‌شود. به‌منظور تحقق اهداف پژوهش از مجموعه داده فانتوم پنج [7] استفاده شده است. اطلاعات مربوط به این مجموعه داده در نشانی <http://fantom.gsc.riken.jp/5> موجود است. در این پایگاه داده اطلاعات بیان ژن بیماران تحت عنوان راه‌انداز ژن در دسترس است [8].

ادامه مقاله به‌صورت زیر سازمان‌دهی شده است. بخش ۲ مسأله خوشه‌بندی تجمعی و تئوری گراف معرفی می‌شود. بخش ۳، کارهای مرتبط و بخش ۴، جزئیات روش پیشنهادی به همراه مراحل متفاوت آن شرح داده می‌شود.



رهیافت‌های ریاضی برای بررسی رفتار سلولی بیشینه طولانی دارد [12]، [13]. برای بررسی رفتار سلولی نظریه گراف سهم گسترده‌ای به دانش بیولوژی اضافه کرده است [14]، [15]. نظریه گراف به بررسی تعامل بر روی توپولوژی‌های با مقیاس بزرگ از یک شبکه تعامل تمرکز می‌کند. ترکیب دو ره‌یافت به‌منظور بررسی تأثیر تئوری بازی روی توپولوژی سلولی بافت مسطح توسط [16] انجام شد. در این پژوهش تئوری بازی و نظریه گراف برای درمان بافت به‌عنوان یک گراف دینامیکی به‌روزشده بر مبنای تعاملات محلی سلول‌های مجزا با هم ترکیب می‌شوند. وجود سازوکارهای ارتباط سلول‌به‌سلول و بافت‌به‌بافت ممکن است، قادر به طراحی دستگاه‌های شبکه ارتباطی کنترل‌شده بین نانوماشین‌ها باشد [18]. فناوری مهندسی مولکولی اخیر ممکن است اجازه اصلاح اجزا بیولوژیکی به اجزای معماری را که برای دستگاه‌های شبکه ارتباطی مناسب می‌باشند، بدهند. توصیف یک شبکه ارتباط مولکولی مبتنی بر سلول توسط [19] ارائه شد. در این شبکه یک رده‌کلاس از دستگاه‌های ارتباط مولکولی که در آن‌ها نانوماشین‌ها با استفاده از شبکه ارتباطی سلول‌به‌سلول با یکدیگر ارتباط برقرار می‌کنند، مورد توجه است. یک مدل محتوا بر مبنای تصمیم‌گیری و یک روش استخراج گراف توسط [20] ارائه شد. مدل پیشنهادی دو مرحله دارد؛ در مرحله نخست گروه‌های مشابه در گراف تشخیص داده شده سپس ادغام می‌شوند. تولید مجموعه‌ای از یال‌های تکراری بر مبنای درجه پشتیبانی در مرحله دوم انجام می‌شود. به‌منظور درک بهتر سلول‌ها و شبکه‌های ارتباطی خود جمعی و میزان پیشرفت بیماری یک مدل ریاضی یک‌پارچه سلولی سیتوکین متنوع ایجاد شد [21]. در این پژوهش از یک مدل بهینه‌سازی‌شده برای سازگاری بهتر با مجموعه داده بهره گرفته شده است. یک روش نمونه‌گیری تصادفی یکنواخت یال‌ها برای الگوریتم استخراج گراف توسط [22]، استفاده شد. در این روش علاوه بر مقادیر وزن یال‌ها از نام یال‌ها، درجه و وابستگی نیز بهره گرفته شده است. به‌کارگیری استراتژی جستجوی محلی برای تحلیل استخراج گراف برای کشف الگوهای رفتاری بیماری‌ها توسط [24] پیشنهاد شد. برای نخستین بار نقشه‌ای از ارتباطات بین ۱۴۴ نوع سلول بدن انسان ارائه شد [25]. اطلاعات این پژوهش نشان می‌دهد که بیش‌تر سلول‌ها یک بیانی از ده‌ها تا صدها گیرنده برای ایجاد یک شبکه سیگنالی از طریق

مسیرهای گیرنده چندگانه دارند. تشخیص ارتباطات ژنی در شبکه‌های چندلایه جهت آشکارسازی سرطان با استفاده از اطلاعات ژنومی برای شناسایی دارندگان بیماری سرطان نیز بررسی شده است. شبکه‌های چندلایه در نظر گرفته‌شده در [26] ترکیب فاکتور هدف‌گیری رونویسی و هدف‌گیری همکاری MICRO RNA، تعامل پروتئین-پروتئین و شبکه‌های ژنی است. در این پژوهش جهت استخراج اطلاعات بیولوژیکی مرتبط به هم از یک الگوریتم خوشه‌بندی ترکیبی استفاده شده است. ترکیبی از لایه‌های مختلف اطلاعات اجازه استخراج الگوهای نظارتی و نقش کارکردی هر دو ژن شناخته‌شده و ژن‌های نامزد جدید را می‌دهد. مدل‌ها و راه‌حلهایی به‌منظور استخراج خواص پنهان ویژگی‌های شبکه در [27] مورد بحث قرار گرفته است. این نتایج حاصل از این پژوهش به درک بهتر مفاهیم بیولوژیکی کمک خواهد کرد. شبکه‌های نظارتی حاوی اطلاعات مربوط به کنترل بیان ژن در سلول‌ها است. شبکه‌های انتقال سیگنال اغلب از چندین گراف برای نشان دادن یک‌سری فعل‌وانفعالات بین خواص زیستی بهره می‌گیرند. الگوریتم خوشه‌بندی طیفی برای پیدا کردن خوشه‌ها در گراف‌ها مورد استفاده قرار می‌گیرد؛ به‌طوری‌که گروه‌ها با یال‌های بسیار مشابه به هم متصل و ارتباط بین آن‌ها ضعیف است. تراکم شبکه تعداد اتصالات در هر گره را مطابق رابطه (۲) نشان می‌دهد:

$$density = \frac{2|E|}{|V|(|V| - 1)} \quad (2)$$

استدلال می‌شود که شبکه‌های بیولوژیکی به‌طور کلی کم‌اتصال هستند. در یک گراف بدون جهت مرکزیت به‌صورت رابطه (۳) تعریف می‌شود.

$$G_{CLOCIL} = \frac{1}{\sum_{|tev|} dist(i, j)} \quad (3)$$

در همین‌اواخر کنسرسیوم فانوم پنج از تحلیل پوششی بیان ژن cage برای ایجاد ارتقادهنده اطلس استفاده کرده است [25]. در این پژوهش به معرفی نخستین نقشه سیگنال‌دهی سلول‌به‌سلول از طریق ارائه شبکه در مقیاس کلان پرداخته شده است. الگوریتم افزابندی ابر گراف (HGPA) و همچنین الگوریتم فراخوشه‌بندی (MCLA) توسط [29] ارائه شدند. یک ابر گراف در HGPA و MCLA ساخته می‌شود و در آن هر خوشه به‌صورت یک فرایال ارائه می‌شود. HGPA به‌طور مستقیم فراگراف را به‌وسیله برش یک تعداد کمینه

عبارتی  $y_{ij} - med(y_{*j})$ ، جایی که  $y_{ij}$  لگاریتم بیان شده هر ژن از نمونه  $i$  است [37].

یک روش خوشه‌بندی مناسب نقش کلیدی در یافتن ارتباطات بین سلولی و بین بافتی به عهده دارد. هدف ما طراحی یک روش خوشه‌بندی تجمعی بهبودیافته است که کارایی و اثربخشی بیشتری نسبت به روش‌های انفرادی داشته باشد.

در پژوهش‌های قبلی یک چگالی مبتنی بر الگوریتم خوشه‌بندی برای دسته‌بندی مجموعه داده‌های دودویی پیشنهاد شد که تعامل بین سلول‌ها و بافت‌های مختلف را می‌تواند با بیان ژن بین سلول‌ها و یا بافت‌های مختلف نشان دهد. مشکل پژوهش قبلی این بود که تنها از یک معیار شباهت برای خوشه‌بندی داده‌های بیان ژن استفاده می‌شد. این مشکل در این پژوهش با در نظر گرفتن چند روش خوشه‌بندی انفرادی حل شده است. ایده اصلی روش پیشنهادی این است که به جای محاسبه شباهت بین یک زوج از اشیا در یک روش انفرادی، شباهت بین اعضای ایجاد شده در چندین روش محاسبه و سپس این شباهت در عضویت بین خوشه‌های جدید مشتق شود. عاقلانه آن است که ما شباهت بین یک زوج از اشیا را با توجه به چندین روش خوشه‌بندی انفرادی ایجاد کنیم. فرایند تجمیع به منظور محاسبه شباهت به وسیله  $m$  روش انفرادی و یک تابع توافقی  $F$  ساخته می‌شود؛ بنابراین واضح است که اثربخشی بیشتر و امکان تأثیرگذاری بیشتر برای در نظر گرفتن شباهت بین خوشه‌های اولیه به جای شباهت اشیا را دارد؛ سپس ما می‌توانیم مفهوم اطلاعات همسایه اشتراکی را از سطح شیء به سطح خوشه ارتقا دهیم؛ بنابراین دو خوشه به خوبی به یکدیگر مرتبط هستند اگر اشیا آن‌ها در روش‌های انفرادی بیشتری یکسان باشند. اگر دو خوشه از دو روش خوشه‌بندی انفرادی مختلف دارای اعضای مشترکی باشند، درجه قطعیت اعضای آن‌ها به یک خوشه بیشتر خواهد بود. با توجه به نامشخص بودن تعداد خوشه‌های نهایی، اگر دو خوشه دارای درجه بالایی از قطعیت نیز باشند، آن‌ها باید جدا نگه‌داشته شوند. در هر حال، به جای استفاده از میزان شباهت تابع توافقی در تشکیل خوشه‌ها، از آن برای تولید خوشه‌های اولیه و ادغام آن‌ها در جهت یافتن خوشه‌بندی پایانی بهره می‌گیریم. در روش پیشنهادی تعداد خوشه‌های بهینه در مرحله ادغام و با در نظر گرفتن یک حد آستانه به صورت خودکار ایجاد می‌شود.

ممکن از فرایال به  $K$  گروه متصل از تقریب همان اندازه با استفاده از افزایشی فراگراف HMETIS [30] افزایشی می‌کند. MCLA ابتدا شباهت بین زوج خوشه‌ها را از جنبه اشتراک اشیا بین آن‌ها تعریف می‌کند. این کار از طریق توسعه شاخص Jaccard انجام، سپس گراف طوری ساخته می‌شود که گروه‌ها، خوشه‌ها و یال‌ها، شباهت بین زوج خوشه‌ها را بیان می‌کنند.

یک توسعه دیگر توسط Iam-on [31] با هدف تعریف مجدد ماتریس همکاری ارائه شده است و ارتباط بین خوشه‌های تخمین شده از یک مدل پیوند شبکه و ترجمه این ماتریس به عنوان ویژگی‌های بردارها و یا یک گراف دوبخشی را استفاده می‌کند. این روش از یک الگوریتم خوشه‌بندی معمولی درون یک ماتریس شباهت استفاده می‌کند تا نتایج خوشه نهایی را تولید کند. روش دیگر خوشه بندی تقسیم تجمعی (DICLENS) است که به وسیله [32] ارائه شده است. این روش بر اساس درخت شباهت پوشای کمینه است. در این روش هر گروه یک خوشه و هر یال یک شباهت درون خوشه‌ای بین خوشه‌ها را نشان می‌دهد. این روش، یال‌های با کمینه شباهت، برش داده می‌شوند تا از هم گسیختگی متاخوشه‌ها فراهم شود. به طور کامل واضح است که اعضای ساخته شده دارای یک تأثیر مستقیم و قوی روی کارایی اجماع دارند؛ بنابراین یک الگوریتم خوشه‌بندی برای ایجاد اعضای خوشه‌ها به صورت خودکار به کار می‌رود.

#### ۴- روش پیشنهادی

نرمال‌سازی، گام بسیار مهم در تحلیل ژن‌های متمایزکننده بیان شده است. روش‌های نرمال‌سازی متفاوتی با توجه به نوع داده‌ها، معرفی شده‌اند. همانند RPKM، میانگین متوسط ترمیم شده مقادیر (TMM)، میانگین کوانتوم، RLE و غیره [33]. TMM و RLE دو روش متفاوت برای محاسبه مقیاس‌های پوششی [34] و این که روش RLE یک روش قابل مقایسه است [35]. روش نرمال‌سازی RLE، جهت نرمال‌سازی داده‌های میکروآرایه و داده‌های مربوط به مقادیر بیان ژن سلول‌ها یا بافت‌ها با ابعاد بسیار بالا به کار برده می‌شود. این داده‌ها به‌طور معمول از نمونه‌های انسانی متعددی به دست آمده‌اند. در این پژوهش برای نرمال‌سازی از روش RLE استفاده شده است [36]. در این روش برای هر ژن  $J$  میانگین را برای هر نمونه محاسبه می‌کند؛ یعنی  $med(y_{*j})$ ، به طوری که  $y_{*j}$  مقدار زامین ستون از ماتریس  $[y_{ij}]$  است؛ سپس انحراف از میانگین محاسبه می‌شود. به



نامشخص بودن تعداد خوشه‌ها، مجموعه‌داده به تعداد  $K-1$  خوشه مختلف انجام می‌شود. به طوری که  $k=2,3,\dots,K$  و  $K=N$  است (تعداد کل اشیاء در مجموعه‌داده). تابع توافقی  $F$  یک ماتریس  $N \times N$  متقارن است که شباهت بین تمام زوج از اشیاء را مطابق شکل (۱) نشان می‌دهد.

$f(x_i, x_j)$  به منظور محاسبه شباهت بین دو شیء  $x_i$  و  $x_j$  است و به صورت رابطه (۴) تعریف می‌شود:

$$f(x_i, x_j) = \sum_{P=1}^m \sum_{k=2}^N \sum_{z=1}^k \varphi_{x_i, x_j}^z(P, k), \quad (x_i, x_j) \in X \quad (4)$$

در این رابطه  $m$  روش خوشه‌بندی انفرادی به منظور افزایش دقت، تابع توافقی روی مجموعه اشیاء  $X$  به کار گرفته می‌شود. تابع  $\varphi_{x_i, x_j}^z(P, k)$  شباهت بین دو شیء  $x_i$  و  $x_j$  را در  $Z$  امین خوشه از روش خوشه‌بندی انفرادی  $P$  با  $k$  خوشه نشان می‌دهد. قرار گرفتن دو شیء در یک خوشه در مراحل مختلف خوشه‌بندی، باعث افزایش شباهت دو شیء خواهد شد. این تابع به صورت رابطه (۵) خواهد بود:

$$\varphi_{x_i, x_j}^z(P, k) = \begin{cases} \beta^{(N-k)} & \text{if } (x_i, x_j) \in C_z^P \\ 1 & \text{otherwise} \end{cases} \quad (5)$$

در این رابطه  $\beta$  یک ثابت عددی بین  $[0, 1]$  است که میزان تأثیر شباهت را نسبت به تعداد خوشه‌ها نشان می‌دهد. اهمیت شباهت بین دو شیء در روش‌های خوشه‌بندی انفرادی با تعداد خوشه‌های بالا، بیشتر است. لذا  $\beta^{(N-k)}$  به همسایگی اشیاء در یک خوشه با تعداد خوشه‌های بیشتر، اهمیت بیشتری می‌دهد.  $C_z^P$  اعضای  $Z$ -مین خوشه از  $P$  امین روش خوشه‌بندی انفرادی را نشان می‌دهد.

به منظور استفاده بهتر از شباهت بین زوج اشیاء در تابع توافقی از نرمال شده ماتریس  $F$  به صورت رابطه (۶) استفاده می‌شود.

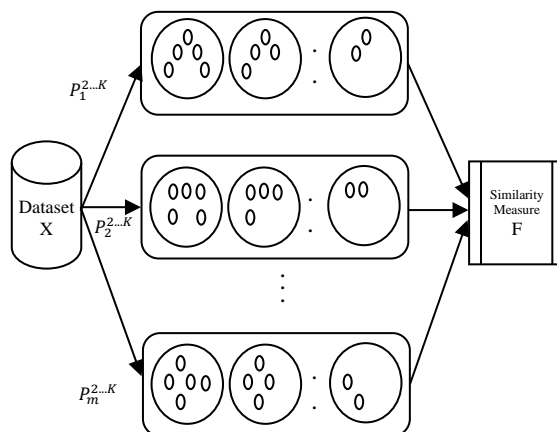
$$f(x_i, x_j) = \frac{f(x_i, x_j)}{\max(F)}, \quad (x_i, x_j) \in X \quad (6)$$

شکل (۲) شبکه‌کد الگوریتم پیشنهادی در مرحله محاسبه شباهت را نشان می‌دهد.

## ۴-۲- ایجاد خوشه‌های اولیه

به منظور ایجاد خوشه‌های اولیه از شباهت‌های بین اشیاء در تابع  $F$  و یک مقدار آستانه استفاده می‌شود. در این مرحله

یکی از مهم‌ترین ویژگی‌های گراف‌های پیچیده وجود ساختارهای اجتماعی است. به طور مشخص شناسایی این ساختارها در گراف‌های پیچیده به تحلیل ویژگی‌های ساختاری گراف کمک می‌کند. ساختار کلی روش پیشنهادی در نظر دارد با توجه به اندازه‌گیری شباهت بین سلول‌ها/بافت‌های بدن انسان، سلول‌ها/بافت‌هایی با بالاترین ارتباط را ارائه دهد. سامانه پیشنهادی شامل دو مرحله کلی است: در مرحله نخست، تمام سلول‌ها یا بافت‌ها با توجه به اندازه‌گیری شباهت در خوشه‌های یکسانی قرار می‌گیرند. در مرحله دوم، معیار شباهت توپولوژیکی به هر یک از خوشه‌ها اعمال شده و تشابه هر سلول/بافت با سایر سلول‌ها/بافت‌ها محاسبه می‌شود. در نهایت، ارتباط هر سلول/بافت با تعدادی از سلول‌ها/بافت‌ها با توجه به بالاترین نمره شباهت ارائه می‌شود. کارایی روش مطرح شده با استفاده از مجموعه‌داده‌های واقعی و همچنین مجموعه‌داده فانتوم پنج مورد ارزیابی قرار می‌گیرد. الگوریتم ارائه شده صرفاً برای مجموعه‌داده‌های واقعی طراحی شده است و حضور نوفه را به حساب نمی‌آورد؛ بنابراین در مورد وجود نوفه و تداخل آن در تعیین تعداد خوشه‌ها نگرانی وجود ندارد، چون تعداد خوشه‌ها به طور خودکار طی فرایند خوشه‌بندی به وسیله ادغام چندین خوشه کوچک استنتاج می‌شود. جزئیات الگوریتم مطرح شده به شرح زیر است:



(شکل-۱): محاسبه شباهت بین هر زوج از اشیاء.  
(Figure-1): Calculates the similarity between each pair of objects.

## ۴-۱- محاسبه شباهت

ارائه یک تابع توافقی مناسب در تعیین میزان شباهت بین هر زوج از اشیاء و تشکیل خوشه‌های نهایی بسیار مؤثر است. شباهت بین دو شیء بر مبنای حضور بیشتر در خوشه‌بندی‌های مختلف محاسبه می‌شود. با توجه به

و  $C_j$  دو نمونه از خوشه اولیه با مقادیر  $C_i = \{a_1, a_2, \dots, a_{z1}\}$  و  $C_j = \{b_1, b_2, \dots, b_{z2}\}$  هستند. فرایند ترکیب دو خوشه بر اساس رابطه (۸) انجام می‌شود:

$$if \sigma(C_i, C_j) \geq th_2 \Rightarrow \begin{cases} \text{hence merged} & \text{true} \\ \text{not merged} & \text{false} \end{cases} \quad (\lambda)$$

پارامتر  $th_2$  حد آستانه برای ترکیب خوشه‌ها است و تأثیر و حساسیت بالایی روی کیفیت نتیجه خوشه‌های پایانی و همچنین تعداد خوشه‌ها دارد؛ بنابراین مقدار آن به‌وسیله یک الگوریتم بهینه‌سازی تپه‌نوردی مورد مطالعه قرار گرفته است. برای  $\sigma$ ، در هر تکرار تنها دو خوشه با بالاترین شباهت که معیار داده‌شده در رابطه (۸) را ارضا می‌کند، باید به‌وسیله جایگزین کردن آن‌ها با یک خوشه جدید  $C_k$  ترکیب شوند. میانگین شباهت درون خوشه‌های  $C_i$  و  $C_j$  به‌صورت رابطه (۹) تعریف می‌شود:

$$\sigma(C_i, C_j) = \frac{\sum_{i=1}^{z1} \sum_{j=1}^{z2} F(a_i, b_j)}{z_1 \times z_2} \quad (9)$$

به‌طوری‌که  $F(a_i, b_j)$  شباهت دو شیء  $a_i$  از خوشه  $C_i$  و  $b_j$  از خوشه  $C_j$  را نشان می‌دهد. برای خوشه نخست، شباهت میان خوشه با سایر خوشه‌ها به‌صورت  $\sigma(C_i) = [\sigma(C_i, C_1), \sigma(C_i, C_2), \dots, \sigma(C_i, C_Z)]$  تعیین می‌شود. بالاترین مقدار از نمونه‌های  $\sigma(C_i)$  نشان‌دهنده بیشترین شباهت بین خوشه  $C_1$  و  $C_z$  برای همه  $z \in (1, 2, \dots, Z)$  است؛ بنابراین با فرض بالاترین شباهت نمونه خوشه‌های  $C_i$  و  $C_j$  به‌وسیله جایگزین شدن با خوشه جدید  $C_k = \{a_1, a_2, \dots, a_{z1}, b_1, b_2, \dots, b_{z2}\}$  ترکیب می‌شوند. ترکیب خوشه‌ها که در نهایت به‌صورت یک خوشه بزرگ نشان داده می‌شوند، راه‌حل نهایی مسأله است. فرایند ترکیب خوشه‌ها تا زمانی که هیچ زوج خوشه‌ای شبیه به هم وجود نداشته باشد، ادامه می‌یابد.

#### ۴-۴- بهبود خوشه‌های نهایی

آنچه در یک روش خوشه‌بندی حائز اهمیت است، کمینه داده‌هایی است که به‌اشتباه به خوشه‌ای تخصیص داده شده است. هدف اصلی از این مرحله برطرف کردن خوشه‌اشیایی است که در مرحله ترکیب خوشه‌ها به‌اشتباه به یک خوشه تخصیص داده شده‌اند. برای شیء  $x_i$  از خوشه  $z$ ، اگر میانگین مقدار شباهت  $x_i^z$  نسبت به سایر اعضای خوشه  $z$  از مقدار آستانه  $th_3$  کمتر باشد  $f(x_i, x_j) \leq th_3$  ( $\forall j = 1, 2, \dots, z_1$ ) شیء از خوشه حذف و یک خوشه جدید به آن تخصیص داده می‌شود.  $\rho_z$  اندیس

اشیا با بیشینه شباهت در خوشه‌های یکسانی قرار داده می‌شوند. برای محاسبه بیشینه شباهت اشیا، هر سطر در ماتریس  $F$  مورد توجه است. رابطه (۷) به‌منظور محاسبه حداکثر شباهت‌ها ارائه شده است:

$$F_{max}^{th_1} = \begin{bmatrix} f_{max}^1 \\ f_{max}^2 \\ \dots \\ f_{max}^N \end{bmatrix} \quad (7)$$

$$= \begin{bmatrix} \max(f_{1,1}, f_{1,2}, \dots, f_{1,N}) \\ \max(f_{2,1}, f_{2,2}, \dots, f_{2,N}) \\ \vdots \\ \max(f_{N,1}, f_{N,2}, \dots, f_{N,N}) \end{bmatrix}_{N \times N}$$

تابع  $F_{max}^{th_1}$  شاخصی برای تعریف بیشینه مقادیر شباهت با توجه به مقدار آستانه  $th$  است؛ بنابراین  $f_{max}^i$  فهرستی از اشیای سطر  $i$ -ام را نشان می‌دهد که  $f_i^i \geq th_1$  جایی که  $\Gamma = 1, 2, \dots, N$  است. محتوای نخستین سطر شامل نخستین شیء است؛ بنابراین اشیایی که به  $f_{max}^1$  رسیده‌اند، به خوشه نخست  $C_1$  تعلق دارند. برای سایر اشیا ( $x_i$ ) از  $2$  تا  $N$  به‌صورت زیر عمل می‌شود:

(۱) تعیین می‌کنید که آیا شیء  $x_i$  در خوشه قبلی دیده شده است یا خیر؟ اگر چنین است پس شیء در این خوشه قرار نمی‌گیرد.

(۲) در غیر این صورت شیء  $x_i$  در خوشه بعدی  $C_z$  قرار داده می‌شود؛ بنابراین اشیایی که به بیشینه شاخص شباهت برای شیء  $x_i$  رسیده‌اند، به خوشه  $C_z$  تعلق دارند؛ بنابراین خوشه‌های اولیه  $C_z$ ، برای همه  $z \in (1, 2, \dots, Z)$  به‌طوری‌که  $Z \leq N$  ایجاد می‌شوند.

Similarity Calculation Algorithm :  
Input : X dataset with N samples and M features; m clustering method and  $\beta$ .  
Output : Similarity matrix  $F_{N \times N}$

```

For P = 1 to m Do
  For k = 2 to N Do
    For z = 1 to k Do
      If  $(x_i, x_j) \in Cluster_z$  from p method Then
         $F(x_i, x_j) = \beta^{(N-k)}$ 
      Else
         $F(x_i, x_j) = 1$ 
      End If
    End For
  End For
End For
F ← Normalized F;
    
```

(شکل-۲): الگوریتم محاسبه شباهت.

(Figure-2): Similarity Calculation Algorithm

#### ۴-۳- ترکیب خوشه‌های اولیه

این مرحله با  $Z$  خوشه اولیه شروع و ترکیب خوشه‌ها بر اساس مشابهت آن‌ها به‌صورت سلسله‌مراتبی تا رسیدن به یک حد آستانه تکرار می‌شود. در هر تکرار دو خوشه با بالاترین مشابهت بین اعضا ترکیب می‌شوند. فرض کنید  $C_i$



با اندیس  $\psi_z$  در صورت کمتر بودن نسبت به پارامتر حد آستانه  $th_3$  از خوشه  $Z$  حذف می‌شود.

برای تعیین خوشه‌های مناسب برای اشیای حذف‌شده مرحله ترکیب خوشه‌ها دوباره تکرار می‌شود. دو مرحله بهبود خوشه‌های نهایی و ترکیب خوشه‌ها، به صورت سلسله‌مراتبی تکرار شده تا زمانی که هیچ تغییری در خوشه‌های نهایی حاصل نشود. شکل (۳) شبه‌کد الگوریتم خوشه‌بندی پیشنهادی را نشان می‌دهد.

شیء از خوشه  $Z$  را با کمترین مشابهت نسبت به سایر اعضا نشان می‌دهد و به صورت رابطه (۱۰) تعریف می‌شود:

$$\psi_z = \arg \min_{1 \leq j \leq z_1} \left\{ \frac{1}{z_1} \sum_{i=1}^{z_1} f(x_i, x_j) \right\} \quad (\forall z = 1, \dots, Z) \quad (10)$$

به طوری که  $z_1$  تعداد اعضای خوشه  $Z$  و  $\arg \min\{*\}$  اندیس شیء با کمترین مقدار شباهت را نشان می‌دهند. شیء

```

Clustering Algorithm :
Input :  $F_{N \times N}$ ;  $th_1$ ;  $th_2$ ;  $th_3$ 
Output : Cluster members,  $C_z (\forall j = 1, \dots, Z)$ 

//Creation of initial clusters :
 $F_{max}^{th_1}$  = Maximum similarity in the F for each  $x_i$  sample, ( $\forall i = 1, \dots, N$ ).
 $z = 1, C_z = \{1^{st} \text{ sample and } f_{max}^1\}$ ;
For  $i = 2$  to  $N$  Do
    If sample  $x_i$  has not included in any previous clusters Then
         $z = z + 1$ ;
         $f_{max}^i$  = samples of  $f_{max}^i$  has not included in any previous clusters.
         $C_z = \{i^{st} \text{ sample and } f_{max}^i\}$ ;
    End If
End For
Repeat until no change in the final cluster
    //Merging of the initial clusters :
     $Z \leftarrow$  Total number of clusters;
     $\sigma^{max} = \infty$ ;
    While  $\sigma^{max} \geq th^2$  Do
        For  $z = 1$  to  $Z$  Do
             $\sigma(C_z, C_j) = eq. (9), (\forall j = 1, \dots, Z)$ 
        End For
         $[C_i, C_j, \sigma^{max}] = \max_{1 \leq i \leq Z} \sigma(C_i, C_j), (\forall j = 1, \dots, Z)$ 
         $C_{new} = \{C_i, C_j\}$ ;
        Remove  $C_i$  and  $C_j$  clusters.
    End While
    //Improve the final cluster :
    For  $z = 1$  to  $Z$  Do
         $\psi = eq. (10)$ ;
        If  $\psi \leq th_3$  Then
            Remove  $\psi$  sample from  $z$  cluster.
             $C_{new} = \{C_\psi\}$ ;
        End If
    End For
End Repeat
    
```

(شکل-۳): الگوریتم خوشه‌بندی پیشنهادی  
(Figure-3): Proposed Clustering Algorithm

در گراف  $G$  را نشان می‌دهند. میزان شباهت در گراف وزن یال‌ها را مشخص می‌کند؛ بنابراین  $W = F$ . ماتریس مجاورت  $A$  از گراف  $G$  یک ماتریس  $N \times N$  متقارن با برچسب‌های رئوس گراف و با مقادیر صفر و یک است. مقدار صفر و یک در موقعیت  $e(V_i, V_j)$  با توجه به این که آیا  $V_i$  و  $V_j$  مرتبط هستند یا نه تعیین می‌شود. ارتباط بین رئوس (سلول‌ها/بافت‌ها) با توجه به معیار تعریف‌شده در رابطه (۱۱) تعیین می‌شود:

$$A_{i,j} \Rightarrow \begin{cases} 1 & W_{i,j} \geq th_4 \\ 0 & otherwise \end{cases} \quad (11)$$

## ۵-۴- تئوری گراف در تشخیص ارتباطات بین سلولی/بافتی

یک گراف  $G = (V, E, W)$  مجموعه‌ای از رئوس  $V$  و مجموعه‌ای از یال‌های  $E$  است؛ به طوری که یال به یک جفت از رئوس می‌پیوندد.  $W(e)$  برای هر یال  $e \in E$  متناظر با وزن است. در اینجا، به طور کلی  $G$  همیشه یک گراف بدون جهت وزن‌دار خواهد بود؛ بنابراین  $e(V_i, V_j)$  و  $e(V_j, V_i)$  نشان‌دهنده یال مشابه در گراف هستند. رئوس گراف نماینده سلول‌ها/بافت‌ها و یال‌ها ارتباط بین سلول‌ها/بافت‌ها

نمی‌کند و به جای آن از مسیرهایی با بیشینه طول  $l \in [2, L]$  بین رئوس استفاده می‌کند. برای محاسبه شباهت بین سلول‌ها یا بافت‌ها در داخل هر خوشه، ابتدا ماتریس مجاورت برای هر خوشه ایجاد می‌شود. رابطه (۱۲) ماتریس مجاورت خوشه  $Z$  را نشان می‌دهد:

$$A_{i,j}(z) \Rightarrow \begin{cases} A_{i,j} & V_i \in C_z \text{ and } V_j \in C_z \\ 0 & \text{otherwise} \end{cases} \quad (12)$$

به طوری که  $C_z$  فهرست رئوس گراف در خوشه  $Z$  را مشخص می‌کند.

ماتریس مجاورت گراف  $G$  هنگامی که به توان  $n$  می‌رسد، تعدادی از مسیرهای با طول  $n$  بین دو رأس را نشان می‌دهد. در این مقاله، با توجه به در نظر گرفتن  $L = 4$ ، برای هر خوشه، ماتریس مجاورت به توان‌های ۲، ۳ و ۴ می‌رسد. برای مثال  $A_{i,j}^2(z)$  تعداد مسیرهای غیر مجاور بین دو رأس  $V_i$  و  $V_j$  با طول مسیر ۲ را در خوشه  $Z$  مشخص می‌کند.

در مرحله بعد برای تعیین شباهت بین سلول‌ها/بافت‌ها، یک معیار ترکیبی جدید به نام  $KFL$  پیشنهاد می‌دهیم که ترکیبی از دو معیار شباهت  $Katz$  و  $FriendLink$  است. با در نظر گرفتن این دو معیار و نیز گره‌های مشترک رأس میانی با دو رأس اصلی محاسبه می‌شود. معیار شباهت پیشنهادی بین دو گره  $v_x$  و  $v_y$  به صورت رابطه (۱۳) محاسبه می‌شود.

$$KFL(v_x, v_y) = \sum_{i=2}^l \frac{1 - (|paths_{v_x, v_y}^i| \times \alpha^i)}{i-1} \cdot \frac{|paths_{v_x, v_y}^i|}{\prod_{j=2}^i (n-j)} \quad (13)$$

در این رابطه  $l$  بیشینه طول مسیر مورد بررسی است.  $\alpha$  پارامتر ثابتی بین  $[0,1]$  است که با توجه به اینکه به توان  $i$  می‌رسد، مقدار آن به نسبت طول ارتباطات متغیر است. هرچه طول ارتباطها بزرگ‌تر شود، ارزش این ترم به دلیل کم‌اهمیت شدن و کاهش احتمال ارتباط، کم خواهد شد.  $N$  تعداد کل گره‌ها و  $|paths_{v_x, v_y}^i|$  تعداد مسیرهای بین  $v_x$  و  $v_y$  با طول  $i$  است. در ارزیابی‌های انجام شده معیار  $KFL$  با مقدار پارامترهای  $l = 3$  و  $\alpha = 0.05$  عملکرد بهتری دارد؛ در نهایت، با استفاده از ماتریس شباهت  $sim$  مبتنی بر الگوریتم ابتکاری  $KFL$ ،  $k$  تعداد سلول‌ها/بافت با بالاترین میزان شباهت به عنوان سلول‌ها/بافت‌های مرتبط به سلول/سلول هدف توصیه می‌شوند. رابطه (۱۴) سلول‌ها/بافت‌های مرتبط با سلول/بافت  $x_i$  را نشان می‌دهد. به طوری که  $rank[*]$  رتبه شباهت را با توجه به سلول/بافت  $x_i$  و تمام سلول‌ها/بافت‌های غیر مجاور آن را در خوشه

به طوری که  $W_{i,j}$  وزن بین رئوس  $V_i$  و  $V_j$  را نشان می‌دهد و  $th_4$  حد آستانه وزن برای تشکیل ماتریس مجاورت گراف  $G$  است. هدف این مرحله یافتن فهرستی از سلول‌ها یا بافت‌های مشابه برای یک سلول یا بافت  $V_c$  ( $\forall c = 1, 2, \dots, N$ ) است.

سامانه‌های توصیه‌گر دوست (FRS) برنامه‌های بسیار محبوب در خدمات شبکه‌های اجتماعی برخط هستند که به کاربران کمک می‌کند تا دوستان جدید با علایق مشابه بیابند. ساختار کلی سامانه در نظر دارد با توجه به اندازه‌گیری مشابهت برای محاسبه شباهت بین کاربران و هر کاربر، دوستان با بالاترین شباهت را ارائه دهد. در این مقاله، یک FRS با استفاده از ترکیب «شباهت توپولوژیکی ساختار گراف» و «اطلاعات خوشه‌بندی سلول‌ها/بافت‌ها» به منظور معرفی سلول‌ها/بافت‌های مرتبط با ویژگی‌ها و کاربردهای مشابه معرفی شده است. در واقع از این روش برای شناسایی جوامع (سلول‌ها یا بافت‌های مشابه) بهره می‌گیریم.

اگرچه تعریف کامل و مستقلی برای یک جامعه وجود ندارد، اما بیشتر شبکه‌های دنیای واقعی ساختارهای جامعه را نشان می‌دهند. یک جامعه (که همچنین خوشه، ماژول یا گروه نامیده می‌شود) یک زیرمجموعه از گره‌ها است که وجه تشابه زیادی با یکدیگر دارند [38]. جوامع در شبکه‌ها گروهی از گره‌ها هستند که بین آن‌ها ارتباطات کم است؛ در حالی که درون آن اتصالات چگال است [39]. تشخیص ساختار جامعه در شبکه‌های پیچیده توجه زیادی از پژوهش‌گران را به خود جلب کرده است؛ زیرا ساختار جامعه در شبکه‌های دنیای واقعی مثل ارتباط بین بیان ژن‌ها، پروتئین‌ها [40] و راه‌انداز ژن برای تشخیص انواع بیماری، شبکه‌های اجتماعی [39]؛ تجزیه و تحلیل سند، که در آن رأس‌ها نشان‌دهنده نویسندگان یا اسناد هستند [41] و غیره به کار برده می‌شوند.

سامانه پیشنهادی این پژوهش با استفاده از یک معیار شباهت جدید که ترکیبی از الگوریتم  $FriendLink$  [42] و  $Katz$  [43] به هر یک از خوشه‌ها اعمال شده و تشابه هر سلول/بافت با تمام سلول‌ها/بافت‌های غیر مجاور محاسبه و در نهایت، به هر سلول/بافت با توجه به بالاترین نمره شباهت، تعدادی سلول/بافت به عنوان سلول‌ها/بافت‌های مرتبط ارائه می‌شود. الگوریتم  $FriendLink$  و  $Katz$  یک روش برای پیش‌بینی پیوندهای جدید در گراف است که با پیمودن تمام مسیرها با طول محدود عمل می‌کند، به طوری که از همه جهت‌ها با طول‌های مختلف در شبکه استفاده



جاری نشان می‌دهد. شکل (۴) شبه‌کد الگوریتم پیشنهادی تئوری گراف در تشخیص ارتباط‌های بین سلولی را نشان می‌دهد. تشخیص ارتباط‌های بین بافتی نیز به صورت مشابه انجام می‌شود:

$$com(x_i, x_j) = \begin{cases} 1 & \text{rank}[sim(x_i, x_j)] \geq k, (\forall j = 1, 2, \dots, N) \\ 0 & \text{otherwise} \end{cases} \quad (14)$$

## ۵- نتایج و آزمایش‌ها

در این بخش عملکرد الگوریتم پیشنهادی به منظور کشف روابط بین سلول‌ها/بافت‌ها با استفاده از پنج مجموعه داده واقعی از مخزن یادگیری ماشین UCI مورد ارزیابی قرار گرفته است. Iris, Wine, Vowel, Glass و Seed مجموعه داده‌هایی هستند که ما استفاده کرده‌ایم. جدول (۱) جزئیات این مجموعه داده‌ها را نشان می‌دهد. برخی از مجموعه داده‌ها دارای مقادیر گم‌شده در برخی از اشیاء هستند که حذف شده است.

Cell-to-Cell Communication with Graph topological Algorithm Input : Cluster members, $C_z \forall j = 1, \dots, Z$ and Adjacency matrix $A$ Output : $com(x_i, x_j) \forall x_i, x_j \in X$
$A(z)$ : Adjacency matrix calculate for the $z$ cluster $\forall z = 1, \dots, Z$ . <b>For</b> $z = 2$ to $Z$ <b>Do</b> <b>For</b> $i = 2$ to $N$ <b>Do</b> $sim(v_i, v_j) = eq.(13), (\forall j = 1, 2, \dots, N)$ //KFL algorithm is applied to the $A(z)$ . $com(x_i, x_j) = eq.(14)$ // For each cell $k$ number of cells with the highest similarity score obtained are offered as a cell-to-cell communication. <b>End For</b> <b>End For</b>

(شکل-۴): الگوریتم پیشنهادی تئوری گراف در تشخیص ارتباط‌های بین سلولی و بین بافتی.

(Figure-4): The proposed algorithm for graph theory in detecting intercellular and interconnected relationships.

(جدول-۱): مجموعه داده‌های استفاده شده در آزمایش‌ها.

(Table-1): The data set used in experiments

مجموعه داده	تعداد نمونه‌ها	تعداد ویژگی‌ها	تعداد خوشه‌ها	تعداد نمونه در هر خوشه	نوع ویژگی‌ها
Iris	۱۵۰	۴	۳	۵۰-۵۰-۵۰	(۴ R)
Wine	۱۷۸	۱۳	۳	۴۸-۷۱-۵۹	(۱۱ R), (۲۱)
Thyroid	۲۱۵	۵	۳	۳۰-۳۵-۱۵۰	(۵ R)
Glass	۲۱۴	۹	۶	۱۳-۱۷-۷۸-۷۰-۲۹-۹	(۹ R)
Bcw	۶۸۳	۹	۲	۲۳۹-۴۴۴	(۹۱)

خوشه‌بندی تجمعی، روشی جدید در خوشه‌بندی است که از ترکیب نتایج روش‌های خوشه‌بندی مختلف به دست می‌آید. صحت، درستی و پایداری از مشخصه‌های مهم یک سامانه خوشه‌بندی تجمعی در مقایسه با روش‌های کلاسیک خوشه‌بندی است. در این مقاله از سه الگوریتم K-Means, FCM و K-Medoids برای خوشه‌بندی اسناد استفاده شده است. الگوریتم پیشنهادی با ترکیب و استفاده از نتایج این سه روش سعی در بهبود نتایج بر روی مجموعه داده‌ها دارد.

با توجه به استفاده از پارامترهای حد آستانه در الگوریتم پیشنهادی و تأثیر آن‌ها در تخمین تعداد صحیح خوشه‌ها به صورت خودکار، انتخاب بهترین مقادیر برای این پارامترها یک چالش بزرگ برای ما است. راه‌حلی که برای این مشکل ارائه شده، بهینه‌سازی مقدار پارامترها با استفاده از الگوریتم تپه‌نوردی است. تپه‌نوردی یک فن بهینه‌سازی متعلق به خانواده الگوریتم‌های جستجوی محلی است؛ یک فن تکرارشونده که با یک راه‌حل دلخواه شروع به کار کرده و سپس تلاش می‌کند تا با تغییر بر روی یک عنصر از راه حل، به پاسخ بهتری دست پیدا کند. اگر این تغییر منجر به ایجاد یک راه حل بهتر شود، تغییر دیگری بر روی این راه حل جدید انجام خواهد گرفت. این روال تا زمانی که بهبود بیشتری در راه حل میسر نباشد ادامه می‌یابد.

با توجه به استفاده از چندین مجموعه داده واقعی در آزمایش‌ها،  $fitness_{HC}$  عملکرد تخمین مقادیر آستانه  $th_1, th_2, th_3$  و  $th_4$  به وسیله الگوریتم تپه‌نوردی را نشان می‌دهد و با رابطه (۱۵) تعریف می‌شود. این رابطه میانگین دقت خوشه‌بندی انجام شده را به وسیله روش پیشنهادی روی تمام مجموعه داده‌های استفاده شده محاسبه می‌کند:

$$fitness_{HC} = \frac{1}{H} \sum_{h=1}^H FC(th_i, ds_h), (\forall i = 1, 2, 3, 4) \quad (15)$$

به طوری که  $H$  تعداد مجموعه داده‌های استفاده شده و  $ds_h$   $n$  آمین مجموعه داده را نشان می‌دهد. تابع  $FC$  نقش الگوریتم پیشنهادی را دارد و دقت خوشه‌بندی را با ورودی‌های تعیین شده محاسبه می‌کند.

در آزمایش‌ها مقادیر زیر برای پارامترها در نظر گرفته شده و سپس چارچوب خوشه‌بندی تجمعی پیشنهادی دنبال می‌شود:

$$th_1 = 0.75, th_2 = 0.0034, th_3 = 0.13, th_4 = 0.5, L = 4, \beta = 0.95$$

این الگوریتم‌ها شامل روش همکاری (CO) با استفاده از پیوند میانگین [28] و ONCE که آن‌هم از پیوند میانگین استفاده می‌کند [23] و MCLA [29] است. روش همکاری پیوندی (CO) از مشکل تطبیق برچسب به‌وسیله نگاشت اعضای اجماع به داخل یک نمایش جدید که در آن ماتریس مشابهت بین یک زوج از اشیا با توجه به چگونگی تعداد اعضای یک زوج مخصوص با یکدیگر خوشه می‌شوند، جلوگیری می‌شود. در این کار، افزاز نهایی به‌وسیله به‌کارگیری ادغام الگوریتم‌های خوشه‌بندی سلسله‌مراتبی پیوند منفرد یا میانگین برای ماتریس تطبیق به‌دست می‌آید. در روش ONCE یک مسأله توافق غیرقطعی بین اعضا مطرح شده است. روش جدید خوشه‌بندی تجمعی اعضای همسایه یا خوشه‌بندی تجمعی نامیده شده است. هدف از این روش استفاده از رابطه همسایگی بین زوج‌اشیا است که این ارتباط بین زوج‌ها به‌خوبی برای محاسبه ماتریس مشابهت به کار می‌رود. در MCLA ابتدا شباهت بین زوج خوشه‌ها از جنبه اشتراک اشیا بین آن‌ها را تعریف می‌کند. این کار از طریق توسعه شاخص Jaccard انجام، سپس گراف طوری ساخته می‌شود که گره‌ها، خوشه‌ها و یال‌ها، شباهت بین زوج خوشه‌ها را بیان می‌کنند. ما الگوریتم را ده مرتبه اجرا و هر مرتبه کارایی را به‌وسیله MMI, ARI اندازه‌گیری کردیم. هنگامی که NMI, ARI برای ارزیابی نتایج خوشه‌بندی به کار می‌روند، یکی از افزازهای خوشه‌بندی باید به‌عنوان افزاز درست داده‌ها در نظر گرفته شود که با  $p^t$  نمایش داده می‌شود. درعمل به‌طور معمول برچسب‌های رده در نظر گرفته می‌شود که برای بقیه پاسخ‌هایی که می‌تواند برای تأیید کیفیت (دقت) نتیجه خوشه‌بندی استفاده شود، وجود ندارد. افزاز دیگر نتیجه خوشه‌بندی تجمعی است که نیاز به محاسبه  $p^*$  دارد. ARI نسخه وابسته به شاخص (RI) است که به‌صورت رابطه (۱۶) تعریف می‌شود:

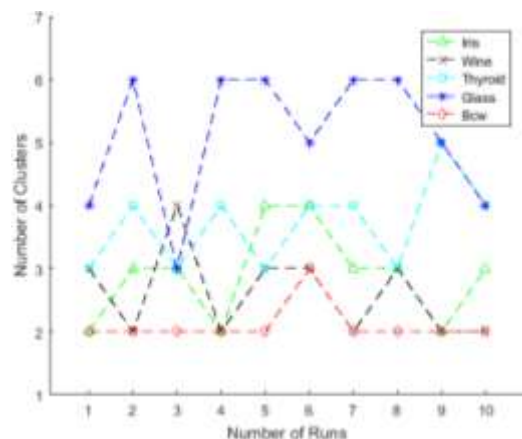
$$ARI(P^*, P^t) = \frac{RI(P^*, P^t) - Expected[RI]}{1 - Expected[RI]} \quad (16)$$

که در آن RI از رابطه (۱۷) محاسبه می‌شود:

$$RI(P^*, P^t) = \frac{n_{11} + n_{00}}{n_{00} + n_{11} + n_{10} + n_{01}} \quad (17)$$

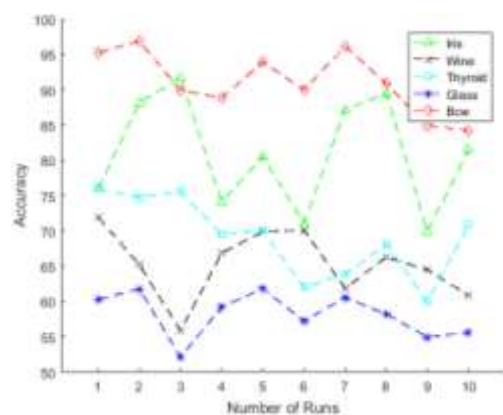
در این رابطه  $n_{11}$  تعداد زوج شیء اختصاص داده به خوشه‌های مشابه در هر دو  $p^*$  و  $p^t$  است.  $n_{00}$  تعداد زوج شیء اختصاص داده‌شده به خوشه‌های متفاوت در  $p^*$  و  $p^t$  است.  $n_{10}$  تعداد زوج‌اشیای اختصاص داده‌شده به همان

در تمام آزمایش‌های انجام‌شده از همه نمونه‌ها در مجموعه داده‌ها استفاده شده است. هدف اصلی این آزمایش‌ها، آزمایش کردن کارایی الگوریتم خوشه‌بندی پیشنهادی و همچنین مشاهده چگونگی تأثیر الگوریتم ما و مقایسه آن با سایر الگوریتم‌های رقیب است. شکل (۵) تعداد خوشه‌های کشف‌شده به‌وسیله الگوریتم پیشنهادی در همه مجموعه داده‌های آزمایش‌شده در ده اجرا را نشان می‌دهد. نتیجه به‌دست‌آمده حاکی از آن است که خوشه‌ها در بیشتر مجموعه داده‌ها درست و یا نزدیک به مقدار واقعی تخمین زده شده اما همچنان ناپایدار و قابل تغییر در ده اجرا است.



(شکل-۵): تعداد خوشه‌های تولیدشده در ده اجرای مجزا. (Figure-5): The number of clusters in 10 separate runs

دقت نتایج روش پیشنهادی مطابق با شکل (۵) از ده اجرای مجزا برای همه مجموعه داده‌های آزمایش‌شده در شکل (۶) نشان داده شده است. نتیجه به‌دست‌آمده حاکی از عملکرد خوب روش پیشنهادی در تعیین خوشه‌های درست در بیشتر مجموعه داده‌ها است. به‌منظور بررسی دقیق‌تر روش پیشنهادی ما نتایج را با چندین روش مشابه مقایسه و عملکرد آن را تجزیه و تحلیل می‌کنیم.



(شکل-۶): دقت خوشه‌بندی تولیدشده در ده اجرای مجزا. (Figure-6): The accuracy of the cluster in 10 separate runs.

خوشه در  $p^*$  و خوشه‌های متفاوت در  $p^t$  است.  $n_{01}$  تعداد زوج‌اشیای اختصاص داده‌شده به خوشه‌های متفاوت  $p^*$  و همان خوشه در  $p^t$  است. بیشینه مقدار ARI برابر با یک است. این بدان معنا است که  $p^*$  با  $p^t$  یکسان (همگون) است و یک مقدار مورد انتظار صفر برای خوشه‌بندی مستقل دارد. NMI بر اساس میانگین اطلاعات دوبه‌دو بین هر زوج از خوشه‌ها و برچسب‌های رده، محاسبه می‌شود. در نظر بگیرید که  $p^*$  نتیجه آخر خوشه‌بندی اجماع و  $p^t$  خوشه‌بندی ground-truth برای مجموعه داده X باشد. NMI دو افراز به صورت رابطه (۱۸) محاسبه می‌شود:

$$NMI(P^*, P^t) = \frac{\sum_{i=1}^{k^*} \sum_{j=1}^{k^t} n_{ij} \log \left( \frac{n_{ij}}{n_i \cdot n_j} \right)}{\left( \sum_{i=1}^{k^*} n_i \log \left( \frac{n_i}{n} \right) \right) \left( \sum_{i=1}^{k^t} n_i \log \left( \frac{n_i}{n} \right) \right)} \quad (18)$$

(جدول-۲): میانگین کارایی و انحراف استاندارد معیار ARI

برای ۱۰ اجرا.

(Table-2): Average performance and standard deviation of ARI standard for 10 executions.

مجموعه داده	Co-Average	ONCE-Average	MCLM	روش پیشنهادی
Iris	۰/۶۶۹±۰/۰۶۵	۰/۶۷۴±۰/۰۵۷	۰/۷۲۲±۰/۰۴۳	۰/۶۸۱±۰/۰۴۲
Wine	۰/۳۲۴±۰/۰۴۵	۰/۳۴۴±۰/۰۶۰	۰/۳۹۳±۰/۰۰۸	۰/۴۲۲±۰/۰۱۹
Thyroid	۰/۲۲۵±۰/۱۷۵	۰/۱۸۹±۰/۱۲۱	۰/۴۴۸±۰/۱۱۹	۰/۳۸۳±۰/۰۴۵
Glass	۰/۲۶۵±۰/۰۰۶	۰/۲۵۶±۰/۰۰۸	۰/۱۵۲±۰/۰۲۲	۰/۲۷۱±۰/۰۰۹
Bcw	۰/۸۶۶±۰/۰۱۸	۰/۸۶۰±۰/۰۱۶	۰/۸۶۴±۰/۰۱۴	۰/۸۶۸±۰/۰۱۲
Ave-P	۰/۴۱۴	۰/۴۰۳	۰/۴۳۳	۰/۴۱۷
Ave-C	۰/۰۴۵	۰/۰۳۷	۰/۰۳۵	۰/۰۴۲

بیشینه مقدار NMI برابر یک است که به این معنی است که  $p^*$  همگون با  $p^t$  و کمینه مقدار آن برابر صفر است، که مربوط به زمانی است که  $p^*$  به‌طور کامل با  $p^t$  متفاوت است. جدول (۲) میانگین کارایی اندازه‌گیری شده به‌وسیله شاخص ARI از طریق انحراف استاندارد هر مجموعه داده و میانگین کارایی اعضای تولیدشده را نشان می‌دهد. این اندازه‌گیری شامل میانگین کارایی (Ave-P) هر یک از روش‌های اجماع روی پنج مجموعه داده و همچنین میانگین سازگاری (Ave-C) است. نتایج نشان می‌دهد که الگوریتم پیشنهادی اغلب بهتر از الگوریتم‌های خوشه‌بندی تجمعی بررسی شده کار

می‌کند. این موضوع به‌خصوص در مورد مجموعه داده‌های Wine و Bcw صحیح است. این در حالی است که برای Iris و Thyroid نتایجی نزدیک به بالاترین کارایی در این مجموعه داده‌ها ارائه می‌دهد. هرچند نتیجه روی Bcw نشان می‌دهد که روش پیشنهادی برای مجموعه داده‌های به‌نسبه بزرگ مناسب است. درخصوص سازگاری روش پیشنهادی، در مورد دو مجموعه داده Bcw, Glass بسیار سازگارتر است، درحالی که روی Wine در مقایسه با سایر الگوریتم‌ها از نظر سازگاری در مکان دوم قرار دارد. به‌طور میانگین، دو الگوریتم نتایج بسیار نزدیکی از نظر سازگاری داشتند که عبارت‌اند از MCLA, ONCE. نتایج این دو الگوریتم به ترتیب عبارت‌اند از ۰/۰۳۵ و ۰/۰۳۷.

نتایج آزمایش‌های مشابه با استفاده از شاخص NMI

در جدول (۳) نشان داده شده است. روش پیشنهادی دارای بالاترین کارایی روی دو مجموعه داده Iris و Bcw است. هرچند در Wine و Glass این الگوریتم نتایجی بسیار نزدیک به بالاترین کارایی دارد. در اینجا، چندین نتیجه جالب وجود دارد. نخست این که کارایی روش پیشنهادی بهتر از ONCE و CO در پنج مجموعه داده است، درحالی که عملکرد آن روی سایر مجموعه داده‌ها بسیار نزدیک به آن‌ها است. دوم این که روش پیشنهادی همان کارایی الگوریتم‌های CO و MCLA را در مجموعه داده Bcw دارد و این که بالاترین دقت مربوط به همین مجموعه داده است. درخصوص سازگاری اندازه‌گیری شده به‌وسیله انحراف استاندارد، روش ما در مقایسه با سایر الگوریتم‌ها روی مجموعه داده thyroid دارای بیشترین سازگاری است و نتایج آن بسیار نزدیک به سازگارترین الگوریتم در مجموعه داده‌های دارای بیشترین استفاده نظیر Wine و Bcw است. با ملاحظه میانگین کارایی اعضای ایجادشده، خواهیم یافت که عملکرد همه روش‌های تجمعی به‌اندازه میانگین اعضای روی همه مجموعه داده‌ها است.

(جدول-۳): میانگین کارایی و انحراف استاندارد معیار NMI

برای ۱۰ اجرا.

(Table-3): Average performance and standard deviation of NMI criteria for 10 executions.

مجموعه داده	Co-Average	ONCE-Average	MCLM	روش پیشنهادی
Iris	۰/۷۵۳±۰/۰۱۷	۰/۷۴۹±۰/۰۰۲	۰/۷۵۵±۰/۰۰۳	۰/۷۶۱±۰/۰۲۰
Wine	۰/۴۰۶±۰/۰۱۰	۰/۱۴۵±۰/۰۰۲	۰/۴۱۵±۰/۰۰۰	۰/۴۱۹±۰/۰۰۱

(فیلتر) به ۸۶۴۲۷ کاهش می‌یابد که در ۷۰۲ نمونه برای سلول‌ها و ۱۲۵ نمونه برای بافت‌ها قرار دارند. از این تعداد نمونه، ۱۰۸ سلول منحصربه‌فرد و چهل بافت منحصربه‌فرد وجود دارد که از مجموعه داده انتخاب می‌شوند؛ بنابراین مجموعه داده فانتوم مورد آزمایش دارای ۷۰۲ نمونه برای سلول‌ها، ۱۲۵ نمونه برای بافت‌ها با ۸۶۴۲۷ راه‌انداز است که برای انجام آزمایش‌ها مشخص شده‌اند.

راه‌اندازهایی از مجموعه داده که تمامی مقادیر Expression آن‌ها از عدد آستانه کوچک‌تر است، برای یافتن ارتباط بین سلول‌ها/بافت‌ها مفید نیست. دلیل این امر نزدیک بودن میانگین مقادیر این راه‌اندازها در تمام نمونه‌ها با بیشینه مقدار راه‌انداز در همان نمونه است (انحراف معیار). این امر نشان می‌دهد که میزان Expression در این بخش از کروموزوم در تمام سلول‌ها/بافت‌ها به‌طور تقریبی یکسان است. ارتباط دو سلول یا دو بافت زمانی به وجود می‌آید که در تعداد از سلول‌ها یا بافت‌ها مقدار راه‌اندازها به میزان قابل توجهی Expression شده باشند. تجزیه و تحلیل نشان می‌دهد که مقادیر Expression کمتر از هزار با توجه به بیشینه مقدار Expression در کل مجموعه داده که گاهی به حدود پنج‌هزار نیز می‌رسد برای یافتن ارتباط بین سلول‌ها یا بافت‌ها مفید نیست. جدول (۴) ارتباطات بین سلولی را با توجه به بخش ۴-۵ نشان می‌دهد. ارتباط بین سلول‌ها با توجه به ده مورد از قوی‌ترین ارتباط‌ها بیان شده است. نماد  $G_i$  معرف سلول  $i$ -ام است که در پیوست یک آورده شده است. میزان شباهت بین دو سلول با رابطه (۶) محاسبه شده است. تعداد راه‌اندازهای بیان‌شده در دو سلول با توجه به تعداد شباهت‌ها بین دو نمونه (سلول) اندازه‌گیری شده است. ژنی با بیشترین راه‌اندازهای بیان‌شده نشان‌دهنده منطقه (مکان) ارتباط بین دو سلول است. به‌منظور بررسی حد آستانه مناسب جهت فیلتر کردن برخی از ارتباطات راه‌اندازها و استخراج ارتباط‌های بین سلولی مناسب ما با چهار حد آستانه مختلف ۵۰۰، ۱۰۰۰، ۱۵۰۰ و ۲۰۰۰ نتایج را بررسی می‌کنیم.

(جدول-۴): ارتباط‌های بین سلولی استخراج‌شده از مجموعه داده فانتوم با استفاده از الگوریتم پیشنهادی.

(Table-4): Intercellular connections extracted from the phantom dataset using the proposed algorithm.

نام سلول	میزان شباهت	تعداد راه‌اندازهای بیان‌شده	ژن با بیشترین راه‌انداز بیان‌شده	نام سلول	میزان شباهت	تعداد راه‌اندازهای بیان‌شده	ژن با بیشترین راه‌انداز بیان‌شده
حد آستانه ۱۰۰۰				حد آستانه ۵۰۰			
C <sub>1</sub>	۰/۹۸۱	۶۴۵۸۰	ABLIM1	C <sub>7</sub>	۰/۹۱۱	۶۴۵۸۰	ABLIM1
C <sub>1</sub>	۰/۹۷۸	۶۴۳۲۵	ABLIM1	C <sub>1</sub>	۰/۹۰۹	۶۴۱۶۶	ABLIM1
C <sub>1</sub>	۰/۹۶۳	۶۴۱۶۶	ABLIM1	C <sub>5</sub>	۰/۹۰۹	۵۷۸۵۲	TACC2
C <sub>2</sub>	۰/۹۴۳	۶۱۸۲۳	ABLIM1	C <sub>2</sub>	۰/۸۹۶	۵۸۱۴۳	KIAA1217

Thyroid	۰/۲۹۳±۰/۰۷۷	۰/۲۵۰±۰/۰۰۶	۰/۳۵۶±۰/۰۰۴	۰/۳۱۶±۰/۰۰۴
Glass	۰/۴۴۱±۰/۰۰۰	۰/۴۴۹±۰/۰۰۱	۰/۳۰۷±۰/۰۰۳	۰/۴۵۰±۰/۰۰۱
Bcw	۰/۷۷۳±۰/۰۰۲	۰/۷۶۵±۰/۰۰۲	۰/۷۷۰±۰/۰۰۱	۰/۷۷۵±۰/۰۰۲
Ave-P	۰/۴۸۸	۰/۴۸۱	۰/۴۷۷	۰/۵۰۱
Ave-C	۰/۰۲۳	۰/۰۲۱	۰/۰۲۲	۰/۰۳۷

در حال، الگوریتم ما بهترین الگوریتم از نظر مقایسه با سایر الگوریتم‌ها است. با توجه به نتایج شاخص NMI، این نتایج بسیار نزدیک به نتایج ارائه‌شده به‌وسیله ARI است. نتایج حاصل از ارزیابی الگوریتم خوشه‌بندی پیشنهادی، کارایی این الگوریتم را نشان می‌دهد. در ادامه، از این الگوریتم در جهت خوشه‌بندی مجموعه داده فانتوم و استخراج ارتباطات بین سلولی یا بین بافتی استفاده می‌شود. در این مقاله از داده‌های فانتوم ۵ استفاده شده است. در داده‌های فانتوم ۱۸۲۹ نمونه با ۲۰۱۸۰۲ راه‌انداز وجود دارد. ویژگی‌ها در فانتوم راه‌اندازها هستند که در واقع حاوی اطلاعاتی در مورد ژن‌ها است که منجر به تولید آن شده است. هدف این پژوهش این است که مشخص شود چه ژن‌هایی برای چه بیماری‌هایی Expression شده‌اند. در واقع هدف شناسایی سلول‌ها یا بافت‌هایی است که در یک یا چند بیماری Gene Expression مشابهی دارند؛ بنابراین ویژگی‌های مجموعه داده راه‌اندازها هستند که از قسمت خاصی از یک کروموزوم گرفته شده است و میزان Gene Expression را نشان می‌دهد. بررسی‌های انجام‌شده نشان می‌دهد که مقادیر برخی از راه‌اندازها (در برخی از ژن‌ها) در تمام سلول‌ها/بافت‌ها به میزان بسیار کمی Gene Expression شده‌اند؛ بنابراین با در نظر گرفتن یک حد آستانه می‌توان تنها راه‌اندازهایی از مجموعه داده را انتخاب کرد که در آن‌ها Gene Expression به میزان مناسبی باشد. یک مقدار حد آستانه در نظر گرفته شده است و تنها راه‌اندازهایی که دست‌کم در یکی از سلول‌ها یا بافت‌ها دارای مقدار Gene Expression بالاتر از حد آستانه باشد، انتخاب می‌شوند. تعداد کل راه‌اندازها با در نظر گرفتن آستانه هزار

C <sub>3</sub>	ABLIM1	۶۱۶۱۲	۰/۹۳۰	C <sub>6</sub>	C <sub>2</sub>	ABLIM1	۶۱۸۲۳	۰/۸۹۰	C <sub>6</sub>
C <sub>4</sub>	ABLIM1	۶۱۲۹۱	۰/۹۰۸	C <sub>6</sub>	C <sub>6</sub>	ABLIM1	۶۰۸۸۳	۰/۸۶۵	C <sub>11</sub>
C <sub>5</sub>	TACC2	۶۱۰۹۶	۰/۸۹۶	C <sub>2</sub>	C <sub>9</sub>	ABLIM1	۴۶۵۹۲	۰/۸۶۱	C <sub>11</sub>
C <sub>6</sub>	ABLIM1	۶۰۹۸۲	۰/۸۹۱	C <sub>10</sub>	C <sub>18</sub>	ABLIM1	۵۷۲۵۰	۰/۸۵۸	C <sub>10</sub>
C <sub>6</sub>	ABLIM1	۶۰۸۸۳	۰/۸۸۷	C <sub>11</sub>	C <sub>15</sub>	KIAA1217	۵۷۵۰۴	۰/۸۵۷	C <sub>11</sub>
C <sub>2</sub>	TACC2	۶۰۸۲۳	۰/۸۸۰	C <sub>12</sub>	C <sub>15</sub>	KIAA1217	۵۲۲۰۹	۰/۸۵۳	C <sub>16</sub>
حد آستانه ۲۰۰۰					حد آستانه ۱۵۰۰				
C <sub>1</sub>	ABLIM1	۶۴۵۸۰	۰/۹۷۰	C <sub>7</sub>	C <sub>1</sub>	ABLIM1	۶۴۵۸۰	۰/۹۷۹	C <sub>7</sub>
C <sub>1</sub>	ABLIM1	۶۴۱۶۶	۰/۹۶۱	C <sub>9</sub>	C <sub>1</sub>	ABLIM1	۶۴۳۲۵	۰/۹۷۴	C <sub>8</sub>
C <sub>3</sub>	ABLIM1	۶۱۶۱۲	۰/۹۴۳	C <sub>6</sub>	C <sub>1</sub>	ABLIM1	۶۴۱۶۶	۰/۹۶۸	C <sub>9</sub>
C <sub>5</sub>	TACC2	۶۱۰۹۶	۰/۹۲۱	C <sub>2</sub>	C <sub>2</sub>	ABLIM1	۶۱۸۲۳	۰/۹۴۰	C <sub>6</sub>
C <sub>2</sub>	ABLIM1	۶۰۸۲۳	۰/۹۰۹	C <sub>12</sub>	C <sub>3</sub>	ABLIM1	۶۱۶۱۲	۰/۹۲۲	C <sub>6</sub>
C <sub>9</sub>	ABLIM1	۴۶۵۹۲	۰/۸۸۱	C <sub>11</sub>	C <sub>4</sub>	ABLIM1	۶۱۲۹۱	۰/۸۹۵	C <sub>6</sub>
C <sub>12</sub>	TACC2	۵۴۷۶۵	۰/۸۸۰	C <sub>14</sub>	C <sub>5</sub>	TACC2	۶۱۰۹۶	۰/۸۸۶	C <sub>2</sub>
C <sub>15</sub>	KIAA1217	۵۲۲۰۹	۰/۸۳۴	C <sub>16</sub>	C <sub>11</sub>	MUC6	۵۷۱۰۸	۰/۸۵۸	C <sub>13</sub>
C <sub>15</sub>	KIAA1217	۵۷۵۰۴	۰/۸۱۳	C <sub>11</sub>	C <sub>15</sub>	KIAA1217	۵۷۵۰۴	۰/۸۵۷	C <sub>11</sub>
C <sub>15</sub>	KIAA1217	۵۴۱۸۷	۰/۷۹۸	C <sub>8</sub>	C <sub>15</sub>	KIAA1217	۵۲۲۰۹	۰/۸۵۳	C <sub>16</sub>

جدول ۵-): ارتباط‌های بین بافتی استخراج شده از مجموعه داده فانتوم با استفاده از الگوریتم پیشنهادی.

(Table-5): Interconnecting extracted from the phantom dataset using the proposed algorithm.

نام بافت	میزان شباهت	تعداد راه‌اندازهای بیان شده	ژن با بیشترین راه‌انداز بیان شده	نام بافت	میزان شباهت	تعداد راه‌اندازهای بیان شده	ژن با بیشترین راه‌انداز بیان شده
حد آستانه ۱۰۰۰				حد آستانه ۵۰۰			
T <sub>7</sub>	ACSL	۷۵۳۵۶	۰/۹۷۰	T <sub>11</sub>	T <sub>7</sub>	ACSL	۷۵۳۵۶
T <sub>5</sub>	CELF	۷۴۵۷۳	۰/۹۶۱	T <sub>11</sub>	T <sub>5</sub>	CELF	۷۴۵۷۳
T <sub>4</sub>	TCF7L	۷۴۰۳۰	۰/۹۱۵	T <sub>7</sub>	T <sub>7</sub>	ACSL	۷۴۳۵۸
T <sub>7</sub>	NEURL	۷۳۵۷۲	۰/۸۸۱	T <sub>13</sub>	T <sub>6</sub>	TCF7L	۷۴۰۸۴
T <sub>4</sub>	ACSL	۷۳۵۲۶	۰/۸۵۸	T <sub>5</sub>	T <sub>4</sub>	TCF7L	۷۴۰۳۰
T <sub>10</sub>	TCF7L	۷۳۴۸۳	۰/۸۴۱	T <sub>11</sub>	T <sub>7</sub>	NEURL	۷۳۸۷۲
T <sub>1</sub>	ADD	۷۳۴۹۵	۰/۸۰۸	T <sub>11</sub>	T <sub>5</sub>	TCF7L	۷۳۸۴۰
T <sub>2</sub>	SCD	۷۳۴۵۸	۰/۸۰۰	T <sub>11</sub>	T <sub>4</sub>	ACSL	۷۳۵۲۶
T <sub>6</sub>	NEURL	۷۳۴۰۲	۰/۷۸۰	T <sub>14</sub>	T <sub>3</sub>	ABLIM1	۷۳۵۱۷
T <sub>7</sub>	TACC2	۷۳۲۷۴	۰/۷۷۱	T <sub>9</sub>	T <sub>10</sub>	TCF7L	۷۳۴۸۳
حد آستانه ۲۰۰۰				حد آستانه ۱۵۰۰			
T <sub>7</sub>	ACSL	۷۵۳۵۶	۰/۹۷۰	T <sub>11</sub>	T <sub>7</sub>	ACSL	۷۵۳۵۶
T <sub>5</sub>	CELF	۷۴۵۷۳	۰/۹۶۱	T <sub>11</sub>	T <sub>5</sub>	CELF	۷۴۵۷۳
T <sub>7</sub>	NEURL	۷۳۸۷۲	۰/۹۴۴	T <sub>13</sub>	T <sub>4</sub>	TCF7L	۷۴۰۳۰
T <sub>5</sub>	ACSL	۷۳۸۴۰	۰/۹۱۸	T <sub>14</sub>	T <sub>4</sub>	ACSL	۷۳۵۲۶
T <sub>10</sub>	TCF7L	۷۳۴۸۳	۰/۸۹۱	T <sub>11</sub>	T <sub>10</sub>	TCF7L	۷۳۴۸۳
T <sub>1</sub>	ADD	۷۳۴۹۵	۰/۸۷۳	T <sub>11</sub>	T <sub>2</sub>	SCD	۷۳۴۵۸
T <sub>7</sub>	TACC2	۷۳۲۷۴	۰/۸۶۲	T <sub>9</sub>	T <sub>7</sub>	TACC2	۷۳۲۷۴
T <sub>7</sub>	ABLIM1	۷۳۲۴۵	۰/۷۸۹	T <sub>12</sub>	T <sub>7</sub>	SCD	۷۳۲۴۵
T <sub>5</sub>	SCD	۷۳۰۵۵	۰/۷۶۵	T <sub>13</sub>	T <sub>8</sub>	NEURL	۷۳۱۶۲
T <sub>4</sub>	ACSL	۷۳۰۳۲	۰/۷۸۵	T <sub>6</sub>	T <sub>5</sub>	ABLIM1	۷۳۰۵۵

نسبت به موارد دیگر است. نتایج با این آستانه نشان‌دهنده بیشترین ارتباط بین دو سلول C<sub>1</sub> و C<sub>7</sub> در بیماری‌های

با توجه به نتایج به‌دست‌آمده میزان شباهت بین سلول‌ها با در نظر گرفتن آستانه هزار در بهینه‌ترین حالت

در مجموعه داده فانتوم ممکن است از یک سلول یا بافت چند نمونه و از چند شخص بیماری گرفته شده باشد. همچنین ممکن است از یک سلول یا بافت در بیماری‌های مختلفی نمونه گرفته شده باشد؛ بنابراین هر سلول/بافت می‌تواند با چند یال به سایر سلول‌ها/بافت‌ها مرتبط باشد. این ارتباطها در نهایت تشکیل یک ابر گراف را می‌دهد. با توجه به نبود رده هدف در داده‌های بیان ژن (فانتوم ۵)، نیاز به شاخص‌های اعتبارسنجی برای سنجش میزان صحت نتایج خوشه‌بندی داریم. در این پژوهش از شاخص درونی Silhouette [17] برای این منظور استفاده شده است. این معیار عمل ارزیابی خوشه‌ها را با استفاده از مقادیری که از خوشه‌ها و نمای آن‌ها محاسبه می‌شود، انجام می‌دهد. ارتباط دو سلول یا دو بافت زمانی به وجود می‌آید که در تعدادی از سلول‌ها/بافت‌ها (نمونه‌ها)، مقدار راه‌اندازها به میزان قابل توجهی بیان شده باشند. تجزیه و تحلیل نشان می‌دهد که مقادیر بیان ژن با حد آستانه‌های مختلف مقادیر متفاوتی در ضریب Silhouette ایجاد می‌کند. برای بررسی این موضوع، نتایج خوشه‌بندی در جدول (۶) برای بررسی حد آستانه‌های مختلف گزارش شده و برای ارزیابی نتایج با پژوهش‌های قبلی و خوشه‌بندی دودویی مقایسه شده است. حد آستانه‌ها با توجه به سایر مقالات مشخص شده است و ضریبی از صد است.

مختلف است. تعداد راه‌اندازهای Expression شده در این دو سلول ۶۴۵۸۰ است. این تعداد رفتار مشابه این دو سلول را در مواجهه با بیماری‌های مختلف نشان می‌دهد. همچنین سلول  $C_1$  با سلول  $C_8$  نیز دارای ارتباط بالایی هستند. این اطلاعات می‌تواند در استخراج الگوهای رفتاری از یک ویروس خاص کمک کند. به طور کلی بیشتر ارتباطات در ژن ABLIM1 بیان شده‌اند. جدول (۵) ارتباطات بین بافتی را با توجه به ۱۰ مورد از قوی‌ترین ارتباطات بیان شده نشان می‌دهد. نماد  $T_i$  معرف بافت  $i$ -ام است که در پیوست ۲ آورده شده است. ژنی با بیشترین بیان ژن بیان شده نشان‌دهنده منطقه (مکان) ارتباط بین دو بافت است. به منظور بررسی حد آستانه مناسب جهت فیلتر کردن برخی از ارتباطات راه‌اندازها و استخراج ارتباط‌های بین بافتی مناسب ما با چهار حد آستانه مختلف ۵۰۰، ۱۰۰۰، ۱۵۰۰ و ۲۰۰۰ نتایج را بررسی می‌کنیم. نتایج به دست آمده میزان شباهت بین بافت‌ها با در نظر گرفتن آستانه پانصد در بهینه‌ترین حالت نسبت به موارد دیگر را نشان می‌دهد که بیشترین ارتباط مربوط به دو بافت  $T_7$  و  $T_{11}$  در بیماری‌های مختلف است. تعداد راه‌اندازهای Expression شده در این دو بافت ۷۵۳۵۶ است. به طور کلی بیشتر ارتباطات در ژن ACSL بیان شده‌اند.

(جدول-۶): مقایسه نتایج خوشه‌بندی با حد آستانه‌های مختلف.

(Table-6): Compare clustering results with different thresholds.

	روش‌ها	حد آستانه	۰/۵/۰	۱/۰	۱۵/۰	۲/۰	۲۵/۰	۳/۰	۳۵/۰	۰/۴
سلول‌ها	خوشه‌بندی	تعداد خوشه‌ها	۱۶	۱۹	۱۹	۱۲	۱۰	۱۲	۸	۱۱
		ضریب Silhouette	۰/۶۴۳	۰/۷۵۸	۰/۶۰۱	۰/۷۵۴	۰/۶۹۹	۰/۴۳۵	۰/۳۸۱	-۰/۵۳۱
	بایتری	تعداد خوشه‌ها	۱۵	۱۸	۱۵	۱۴	۱۱	۱۱	۱۶	۱۴
		ضریب Silhouette	۰/۷۱۱	۰/۹۰۱	۰/۵۹۸	۰/۶۷۶	۰/۶۳۲	۰/۶۴۳	۰/۵۴۳	۰/۳۴۲
بافت‌ها	خوشه‌بندی	تعداد خوشه‌ها	۹	۸	۱۱	۱۰	۱۱	۹	۱۳	۱۳
		ضریب Silhouette	۰/۴۹۸	۰/۵۱۶	۰/۴۷۸	۰/۶۰۰	۰/۶۵۴	۰/۶۷۵	۰/۶۶۱	۰/۶۷۳
	بایتری	تعداد خوشه‌ها	۱۱	۱۱	۹	۹	۱۳	۱۳	۱۴	۱۵
		ضریب Silhouette	۰/۵۹۸	۰/۵۲۰	۰/۶۱۹	۰/۶۰۴	۰/۷۲۳	۰/۷۶۲	۰/۷۰۰	۰/۷۳۷

Silhouette ۰/۷۶۲ در سیزده خوشه نسبت به روش دودویی با ضریب Silhouette ۰/۶۷۵ در نه خوشه با حد آستانه ۲/۵ در بافت‌ها به وضوح قابل مشاهده است.

## ۶- نتیجه‌گیری

ارتباطات بین سلول‌ها و بافت‌ها به شناسایی بیماری‌ها مختلف و عوامل آن‌ها کمک خواهد کرد. در واقع ارتباط بین

با توجه به نتایج آزمایش، حد آستانه ۱/۰ دارای ضریب Silhouette با بهترین مقدار ۰/۸۵۸ در ۱۹ خوشه برای روش خوشه‌بندی دودویی در سلول‌ها است. این در حالی است که روش خوشه‌بندی تجمعی پیشنهاد شده در این پژوهش با این حد آستانه به ضریب Silhouette ۰/۹۰۱ در ۱۸ خوشه رسیده است و عملکرد بهتری دارد. همچنین برتری روش پیشنهادی در بهترین حالت با ضریب

- [3] W. Cook, and J. Palsberg, "A denotational semantics of inheritance and its correctness", Vol. 24. No. 10. ACM, 1989.
- [4] H. C. Lukaski, "Methods for the assessment of human body composition: traditional and new", *The American journal of clinical nutrition*, vol. 46 (4), pp. 537-556, 1987.
- [5] R. A. Lewis, B. Otterud, D. Stauffer, J. M. Lalouel, & M. Leppert, "Mapping recessive ophthalmic diseases: linkage of the locus for Usher syndrome type II to a DNA marker on chromosome 1q", *Genomics*, vol. 7(2), pp.250-256, 1990
- [6] A. Bretto, H. Cherifi, D. Aboutajdine, "Hypergraph imaging: an overview", *Pattern Recognition*, vol. 35(3), pp. 651-658, 2002.
- [7] A. R. Forrest, H. Kawaji, M. Rehli, J. K. Baillie, M. J. De Hoon, V. Haberle, R. Andersson, "A level mammalian expression atlas", *Nature*, vol. 507(7493), pp. 462, 2012.
- [8] M. E. Hegi, A. C. Diserens, S. Godard, , P. Y. Dietrich, , L. Regli, , S. Ostermann, R. Stupp, "Clinical trial substantiates the predictive value of O-6-methylguanine-DNA methyltransferase promoter methylation in glioblastoma patients treated with temozolomide", *Clinical cancer research*, vol.10(6), pp. 1871-1874, 2008.
- [9] Karimi A, Hoseini L S. An Optimal Algorithm for Dividing Microscopic Images of Blood for the Diagnosis of Acute Pulmonary Lymphoblastic Cell Using the FCM Algorithm and Genetic Optimization. *JSDP*. 2018; 15 (2) :45-54
- [10] Vahidi Ferdosi S, Amirkhani H. Weighted Ensemble Clustering for Increasing the Accuracy of the Final Clustering. *JSDP*. 2020; 17 (2) :100-85.
- [11] M. A. Ahmad, Z. Borbora, J. Srivastava, & N. Contractor, "Link prediction across multiple social networks. In Data Mining Workshops (ICDMW)", 2010 *IEEE International Conference on*, pp. 911-918, 2010.
- [12] H. M. Byrne, "Dissecting cancer through mathematics: from the cell to the animal

سلول‌ها/بافت‌ها روابط وراثتی بین بیماران را نشان می‌دهد. این روابط به شناسایی نقاط مشترک بدن که تحت تأثیر بیماری‌های مختلفی قرار می‌گیرد، کمک می‌کند. در این مقاله تشخیص ارتباط‌های بین سلولی و بین بافتی در بیماری‌های مختلف با توجه به یک معیار شباهت توپولوژیکی گراف و یک روش خوشه‌بندی تجمعی بهبودیافته ارائه شده است. ارزیابی عملکرد الگوریتم خوشه‌بندی پیشنهادی دقت بالای آن را به اثبات رسانده است. خوشه‌بندی سلول‌ها/بافت‌ها باعث افزایش دقت و تمرکز معیار شباهت توپولوژیکی گراف در محدوده‌های مشابهت سلول‌ها/بافت‌ها شده است. الگوریتم خوشه‌بندی ارائه‌شده از یک تابع توافقی شباهت به‌منظور محاسبه است مشابهت بین اشیا استفاده می‌کند. ایجاد خوشه‌های اولیه با دقت بالایی با استفاده از مشابهت بین اشیا تولید و سپس مشابهت بین خوشه‌های اولیه محاسبه‌شده و با خوشه‌هایی با بالاترین ارتباط با یکدیگر ادغام می‌شوند. به‌منظور اصلاح و افزایش کیفیت خوشه‌های نهایی، یک الگوریتم بهبود عدم قطعیت اشیا را به‌وسیله تخصیص آن‌ها به سایر خوشه‌ها بررسی می‌کند. ما روش خود را روی پنج مجموعه‌داده دنیای واقعی آزمایش کردیم. نتایج نشان می‌دهد که به‌طور میانگین الگوریتم خوشه‌بندی پیشنهادی ما بهتر از سایر الگوریتم‌های مشابه کار می‌کند. از مزیت‌های روش پیشنهادی تعیین تعداد بهینه خوشه‌ها به‌صورت خودکار است. به‌جای محاسبه مشابهت بین اشیا به‌وسیله یک روش خوشه‌بندی انفرادی، از چندین روش بهره می‌گیرد. همچنین استفاده از الگوریتم تپه‌نوردی در تعیین مقادیر حد آستانه باعث کاهش پارامترهای ورودی مسأله و افزایش دقت در تعیین تعداد خوشه‌های بهینه‌شده است. ارزیابی کیفیت نتایج به‌وسیله یک استراتژی تطبیقی برای مقادیر آستانه می‌تواند در پژوهش‌های آینده مورد بررسی قرار بگیرد و میزان تأثیر و حساسیت روی کیفیت نتیجه خوشه‌های پایانی با روش حاضر مقایسه شود.

## 7- References

## ۷- مراجع

- [1] I. H. Park, N. Arora, H. Huo, N. Maherali, T. Ahfeldt, A. Shimamura, G. Q. Daley, "Disease-specific induced pluripotent stem cells", *cell*, no. 134(5), 877-886, 2008.
- [2] D. R. Bentley, S. Balasubramanian, H. Swerdlow, G. P. Smith, J. Milton, C. G. Brown, J. M. Boutell, "Accurate whole human genome sequencing using reversible terminator chemistry", *nature*, vol. 456(7218), pp.53, 2008.

- [24] W. Cui, Y. Xiao, H. Wang, W. Wang, "Local search of communities in large graphs", *In Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, 2014, pp. 991-1002.
- [25] J. A. Ramilowski, T. Goldberg, J. Harshbarger, E. Kloppman, M. Lizio, V. P. Satagopam, A. Forrest, "A draft network of ligand-receptor-mediated multicellular signalling in human", *Nature communications*, 2015.
- [26] L. Cantini, E. Medico, S. Fortunato, M. Caselle, "Detection of gene communities in multi-networks reveals cancer drivers", *Scientific reports*, vol. 5, srep17386, 2015.
- [27] G. A. Pavlopoulos, M. Secrier, C. Moschopoulos, N. Soldatos, S. Kossida, J. Aerts, P. G. Bagos, "Using graph theory to analyze biological networks", *BioData mining*, vol. 4(1), no.10, 2010.
- [28] A. L. Fred, A. K. Jain, "Combining multiple clusterings using evidence accumulation", *IEEE transactions on pattern analysis and machine intelligence*, vol. 27(6), pp. 835-850, 2005.
- [29] A. Strehl, J. Ghosh, "Relationship-based clustering and visualization for high-dimensional data mining", *INFORMS Journal on Computing*, vol. 15(2), pp. 208-230, 2003.
- [30] G. Karypis, V. Kumar, A hypergraph partitioning package, 1998.
- [31] N. Iam-On, T. Boongoen, S. Garrett, C. Price, "A link-based approach to the cluster ensemble problem", *IEEE transactions on pattern analysis and machine intelligence*, vol.33(12), pp. 2396-2409, 2011.
- [32] S. Mimaroglu, E. Aksehirli, "Dicens: Divisive clustering ensemble with automatic cluster number", *IEEE/ACM transactions on computational biology and bioinformatics*, vol. 9(2), pp. 408-420, 2012.
- [33] M. A. Dillies, A. Rau, J. Aubert, C. Hennequet-Antier, M. Jeanmougin, N. Servant, G. Guernec, "A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis", *Briefings in bioinformatics*, vol. 14(6), pp. 671-683, 2013.
- [34] R. Reddy, "A comparison of methods: normalizing high-throughput RNA sequencing data", *bioRxiv*, 026062, 2015.
- [35] S. Hawinkel, "Evaluation of normalization and analysis methods for microbiome data", 2015.
- [36] S. Noguchi, T. Arakawa, S. Fukuda, M. Furuno, A. Hasegawa, F. Hori, T. Kawashima, "FANTOM5 CAGE profiles of human and mouse samples", *Scientific data*, vol. 4, pp. 112-170, 2017.
- model", *Nature reviews. Cancer*, vol.10 (3), pp. 221, 2010.
- [13] A. Csikász-Nagy, "Computational systems biology of the cell cycle", *Briefings in bioinformatics*, vol. 10(4), pp. 424-434, 2004.
- [14] J. Hofbauer, K. Sigmund, "Evolutionary game dynamics", *Bulletin of the American Mathematical Society*, vol. 40(4), pp.479-519, 2003.
- [15] R. P. Araujo, D. S. McElwain, "A history of the study of solid tumour growth: the contribution of mathematical modelling", *Bulletin of mathematical biology*, vol. 66(5), pp. 1039, 2004.
- [16] A. Csikász-Nagy, M. Cavaliere, S. Sedwards, "Combining game theory and graph theory to model interactions between cells in the tumor microenvironment", *In New Challenges for Cancer Systems Biomedicine* Springer Milan, pp. 3-18, 2012.
- [17] S. Aranganayagi, K. Thangavel, "Clustering categorical data using Silhouette coefficient as a relocating measure", *In Conference on Computational Intelligence and Multimedia Applications, 2007*, International Conference on vol. 2, pp. 13-17, 2007.
- [18] T. Nakano, T. Suda, M. Moore, R. Egashira, A. Enomoto, K. Arima, "Molecular communication for nanomachines using intercellular calcium signaling. In Nanotechnology," 5th IEEE Conference on, pp. 478-481, 2005.
- [19] T. Nakano, M. Moore, A. Enomoto, T. Suda, T. Koujin, T. Haraguchi, Y. Hiraoka, A cell-based molecular communication network, *In Proceedings of the 1st international conference on Bio inspired models of network, information and computing systems*, pp. 23, 2006.
- [20] H. Jiang, J. Liu, Z. Zhao, Graph mining based knowledge discovery in designing decision-making context models, *In Systems and Informatics (ICSAI), 2014 2nd International Conference on*, pp. 948-953, 2014.
- [21] H. Shao, Y. He, K. C. Li, X. Zhou, "A system mathematical model of a cell-cell communication network in amyotrophic lateral sclerosis", *Molecular BioSystems*, vol. 9(3), pp.398-406, 2013.
- [22] R. Gao, H. Xu, P. Hu, W. C. Lau, "Accelerating graph mining algorithms via uniform random edge sampling", *In Communications (ICC), 2016 IEEE International Conference on*, pp. 1-6, 2016.
- [23] T. Alqurashi, W. Wang, "Object-neighbourhood clustering ensemble method", *In International Conference on Intelligent Data Engineering and Automated Learning*, pp. 142-149, 2014.





**موسی مجرد**، در سال ۱۳۹۷ دکترای خود را در رشته رایانه گرایش سیستم‌های نرم‌افزاری از دانشگاه آزاد اسلامی واحد یاسوج دریافت کرد. ایشان هم‌اکنون عضو هیأت علمی دانشگاه آزاد اسلامی واحد فیروزآباد و زمینه‌های پژوهشی مورد علاقه ایشان داده‌کاوی، الگوریتم‌های تکاملی، گراف کاوی، شبکه‌های اجتماعی و بهینه‌سازی است. نشانی رایانامه‌های ایشان عبارتند از:

**m.mojarad@iauf.ac.ir,**  
**mosa.mojarad@gmail.com**



**حمید پروین**، در سال ۱۳۹۰ دکترای خود را در رشته رایانه گرایش هوش مصنوعی از دانشگاه علم و صنعت دریافت کرد و هم‌اکنون عضو هیأت علمی دانشگاه آزاد اسلامی واحد نورآباد ممسنی است. زمینه‌های پژوهشی مورد علاقه ایشان یادگیری ماشین، داده‌کاوی و الگوریتم‌های تکاملی است. نشانی رایانامه ایشان عبارت است از:

**parvinhamid@gmail.com**



**صمد نجاتیان**، در سال ۱۳۹۳ دکترای خود را در رشته برق (مخابرات) از دانشگاه یو تی ام مالزی دریافت کرد و هم‌اکنون عضو هیأت علمی دانشگاه آزاد اسلامی واحد یاسوج است. زمینه‌های پژوهشی مورد علاقه ایشان سامانه‌های هوشمند، تشخیص الگو، ارتباطات بی‌سیم، نرم‌افزار تعریف رادیو، رادیو شناختی و آماری است.

نشانی رایانامه ایشان عبارت است از:

**nejatian@iauyasooj.ac.ir**



**کرم‌الله باقری‌فرد**، مدرک دکترای خود را در سال ۱۳۹۶ در گرایش نرم‌افزار در دانشگاه اراک دریافت کرد و هم‌اکنون عضو هیأت علمی دانشگاه آزاد اسلامی واحد یاسوج است. زمینه‌های پژوهشی مورد علاقه ایشان یادگیری ماشین، داده‌کاوی و سامانه‌های پیشنهاددهنده است. نشانی رایانامه ایشان عبارت است از:

**k.bagheri@iauyasooj.ac.ir**

- [37] L. C. Gandolfo, T. P. Speed, RLE Plots: Visualising Unwanted Variation in High Dimensional Data. arXiv preprint arXiv:1704.03590, 2017.
- [38] M. Arab, M. Hasheminezhad, Limitations of quality metrics for community detection and evaluation, In Web Research (ICWR), 2017 3th International Conference, pp. 7-14, 2017.
- [39] Y. Qiao, H. Wang, D. Wang, "Parallelizing and optimizing overlapping community detection with speaker-listener Label Propagation Algorithm on multi-core architecture", In *Cloud Computing and Big Data Analysis (ICCCBDA)*, 2017 IEEE 2nd International Conference, 2017, pp. 439-443.
- [40] S. Mohammadi, J. Davila-Velderrain, M. Kellis, A. Grama, DECODE-ing sparsity patterns in single-cell RNA-seq. bioRxiv, 241646, 2018
- [41] D. J. D. S. Price, Networks of scientific papers. Science, pp. 510-515, 1965.
- [42] A. Papadimitriou, P. Symeonidis, Y. Manolopoulos, "Fast and accurate link prediction in social networking systems", *Journal of Systems and Software*, vol. 85(9), pp. 2119-2132, 2012.
- [43] L. Katz, "A new status index derived from sociometric analysis", *Psychometrika*, vol. 18 (1), pp. 39-43, 1953.

## پیوست ۱

Full names of cells ( $C_i$ ) in Table(4-19).

- $C_1$ : (hes.gfp.embryonic.stem.cells),  $C_2$ : (basophils),  $C_3$ : (cd14..monocytes...treated.with.lipopolysaccharide),  $C_4$ : (h9.embryonic.stem.cells),  $C_5$ : (lens.epithelial.cells),  $C_6$ : (ciliary.epithelial.cells),  $C_7$ : (cd14.cd16..monocytes.2),  $C_8$ : (fibroblast...aortic.advantitial.donor2..cytoplasmic.f raction),  $C_9$ : (cd14.cd16..monocytes.1),  $C_{10}$ : (b.lymphoblastoid.cell.line..gm12878.encode),  $C_{11}$ : (cd14..monocytes),  $C_{12}$ : (astrocyte...cerebral.cortex),  $C_{13}$ : (mast.cell),  $C_{14}$ : (melanocyte),  $C_{15}$ : (mesothelial.cells),  $C_{16}$ : (fibroblast...villous.mesenchymal),  $C_{17}$ : (melanocyte...dark),  $C_{18}$ : (melanocyte...dark).

## پیوست ۲

Full names of tissues ( $T_i$ ) in Table(4-20).

- $T_1$ : (amygdala),  $T_2$ : (caudate.nucleus),  $T_3$ : (hippocampus),  $T_4$ : (kidney),  $T_5$ : (locus.coeruleus),  $T_6$ : (medial.frontal.gyrus),  $T_7$ : (medulla.oblongata),  $T_8$ : (occipital.cortex),  $T_9$ : (small.intestine),  $T_{10}$ : (temporal.lobe),  $T_{11}$ : (testis),  $T_{12}$ : (throat),  $T_{13}$ : (thymus),  $T_{14}$ : (trachea).