

# پیش بینی پیوند در شبکه‌های علمی با استفاده از یادگیری ماشین و گراف‌های وزن‌دار

سیدمهدی وحیدی پور<sup>\*</sup>، استادیار، علیرضا محمدی<sup>۱</sup>

۱ دانشکده مهندسی برق و کامپیوتر، دانشگاه کاشان، کاشان، ایران، vahidipour@kashanu.ac.ir

چکیده: با سرعت گرفتن رشد علم و انتشار مقالات و افزایش زمینه‌های علمی، یافتن همکار پژوهشی مناسب، یافتن منابع تحقیق و زمینه تحقیق برای محققان و نهادهای مربوطه، روز به روز سخت‌تر می‌شود. با انتخاب درست این موارد، می‌توان بیشترین بازدهی را از هزینه و زمان صرف شده برای پژوهش کسب کرد. برای حل این مسئله می‌توان با ایجاد شبکه‌ای شامل مقالات، دانشمندان و سایر موجودیت‌های علمی و ارتباطات بین آن‌ها، یک شبکه علمی ایجاد کرد و با استفاده از پیش‌بینی پیوند ارتباطاتی که در آینده شکل می‌گیرد را پیش‌بینی کرد. در این مقاله چارچوبی مبتنی بر یادگیری ماشین برای پیش‌بینی پیوند در شبکه‌های علمی ارائه شده است. در این چارچوب با وزندهی شبکه بر اساس زمان و محتوا، محاسبه ویژگی‌های ساختاری و متنی جاسازی شده و انتخاب و استخراج ویژگی انجام می‌شود. در نهایت نمونه‌گیری منفی با استفاده از خوشه‌بندی تولید می‌شود تا یک مدل یادگیری ماشین برای پیش‌بینی پیوند آموزش داده شود. هر یک از مراحل این چارچوب به صورت جدا و همه با هم آزمایش شدند و نتایج نشان داد روش وزندهی پیشنهاد شده برای شبکه ارجاعات و همکاری نویسندگان باعث افزایش دقت معیارهای شباهت وزن‌دار و در نتیجه افزایش دقت کل الگوریتم می‌شود. همچنین نمونه‌گیری منفی با استفاده از خوشه‌بندی باعث بهتر آموزش داده شدن الگوریتم یادگیری ماشین می‌شود. ویژگی‌های متنی داده‌های علمی مانند عنوان و چکیده مقالات نیز نقش مؤثری در پیش‌بینی پیوندهای آینده دارند.

واژه‌های کلیدی: پیش‌بینی پیوند، شبکه ارجاعات، شبکه همکاری نویسندگان، یادگیری ماشین

\* سیدمهدی وحیدی پور، vahidipour@kashanu.ac.ir

# Link prediction in scientific networks using machine learning and weighted graphs

---

S. Mehdi Vahidipour<sup>1\*</sup>, Assistant Professor, Alireza Mohamadi

1 Electrical and Computer Engineering Faculty, University of Kashan, Kashan, Iran,  
vahidipour@kashanu.ac.ir.

**Abstract:** With the acceleration of the development of science and the publication of articles and the increase of scientific fields, finding suitable research partners, finding research sources and research fields for researchers and relevant institutions is becoming more and more difficult. By choosing these things correctly, you can get the most efficiency from the cost and time spent on research. To solve this problem, a scientific network can be created by creating a network including articles, scientists, and other scientific entities and the connections between them, and predicting the connections that will be formed in the future using link prediction. In this paper, a framework based on machine learning is presented for link prediction in scientific networks. In this framework, by weighting the network based on time and content, calculating embedded structural and textual features, feature selection and extraction, and finally negative sampling using clustering, a machine learning model is trained for link prediction. Each of the steps of this framework was tested separately and all together, and the results showed that the proposed weighting method for the network of references and authors' collaboration increases the accuracy of the weighted similarity criteria and, as a result, increases the accuracy of the entire algorithm. Also, negative sampling using clustering makes the machine learning algorithm better trained. The textual features of scientific data such as the title and abstract of articles also play an effective role in predicting future links.

**Keywords:** *Link Prediction; Citation Networks; Author Collaboration Networks; Machine Learning.*

\* S. Mehdi Vahidipour, [vahidipour@kashanu.ac.ir](mailto:vahidipour@kashanu.ac.ir)

## ۱- مقدمه

نشریات را با مفاهیمی مانند ارجاع<sup>۲</sup>، همکاری<sup>۳</sup> و نویسندگی بازنمایی می کنند [۷].

با افزایش سرعت رشد علم و ایجاد زمینه های علمی جدید در سال های اخیر و همکاری دانشمندان در حوزه های مختلف علوم با یکدیگر، این نگرانی برای دانشمندان و سازمان های دانش بنیان به وجود آمده است که در سال های آینده، زمان و بودجه پژوهشی را به چه زمینه هایی اختصاص دهند. همچنین یافتن مقالات مرتبط با یک زمینه علمی خاص [۸] و یافتن دانشمندان متخصص برای همکاری های علمی [۹] از جمله مسائلی هستند که در شبکه های علمی و با استفاده از پیش بینی پیوند راه حل هایی برایشان ارائه شده است.

در پژوهش های انجام شده در زمینه شبکه های علمی، بیشتر تمرکز روی ساختار شبکه و چگونگی پیوند عناصر با یکدیگر بوده است و محتوای عناصر کمتر مورد توجه قرار گرفته است. در برخی پژوهش ها که محتوای متنی برای پیش بینی پیوند استفاده شده نیز از روش های ساده مانند شباهت کسینوسی بردار tf-idf استفاده شده است [۶]. همچنین در روش های یادگیری ماشین، با توجه به حجم زیاد شبکه های علمی، انتخاب داده های آموزشی به صورت تصادفی انجام شدند [۱۰]. در این مقاله برای بهبود این نقاط ضعف، روش های جاسازی کلمات<sup>۴</sup> برای تحلیل محتوای متنی و استفاده از خوشه بندی برای انتخاب داده های آموزشی برای پیش بینی وزن پیشنهاد شده است. در نهایت با ترکیب روش های موجود و روش های پیشنهادی که نتایج را بهبود دادند، یک چارچوب کلی برای پیش بینی پیوند در شبکه های علمی پیشنهاد می کند.

شبکه ها در زندگی امروزه ما نقش مهمی دارند و تحلیل شبکه ها اطلاعات مفیدی در اختیار ما قرار می دهد [1-2]. این اطلاعات این امکان را فراهم می کند که رفتار شبکه ها را بهتر درک کنیم و چگونگی تکامل آن ها را در آینده شبیه سازی کنیم. پیش بینی پیوند یکی از تحلیل هایی است که روی شبکه ها انجام می شود و امروزه اهمیت زیادی پیدا کرده است. پیش بینی پیوند به معنای یافتن پیوندهای از دست رفته و یا پیش بینی پیوندهای آینده است [۳-۴]. این پیش بینی می تواند بر اساس پیوندهای موجود در حال حاضر شبکه، ویژگی های رأس های شبکه و یا هر دو انجام شود [۵].

افراد جدیدی که یک شبکه اجتماعی پیشنهاد دنبال کردن آن ها را به کاربرانش می دهد مثالی از کاربرد پیش بینی پیوند در شبکه های اجتماعی است. اطلاعات گذشته کاربر در شبکه اجتماعی می تواند مبنای این پیشنهاد باشد؛ مثلاً افرادی را که دنبال می کند (ویژگی های ساختاری) و یا ویژگی هایی مانند شغل و محل سکونت (ویژگی های محتوایی) [۶].

پیش بینی پیوند تا کنون در زمینه های بسیاری مانند شبکه های اجتماعی، مسائل بیولوژیک و سیستم های توصیه گر به کار گرفته شده و منجر به نتایج قابل قبولی شده است [۵]. یکی از کاربردهای پیش بینی پیوند که به تازگی اهمیت پیدا کرده ولی کمتر مورد توجه قرار گرفته است، پیش بینی پیوند در شبکه های علمی است؛ شبکه های علمی نوع خاصی از شبکه های پیچیده<sup>۱</sup> هستند که روابط بین موجودیت های علمی مانند دانشمندان، مقالات و

<sup>2</sup> Citation

<sup>3</sup> Co-authorship

<sup>4</sup> Embedding

<sup>1</sup> Complex Networks

هستند [۵]. شبکه‌های دنیای واقعی که دارای ویژگی‌های ساختاری غیر بدیهی<sup>۲</sup> هستند و تفاوت‌های زیادی با گراف‌های تصادفی و گراف‌های مشبک<sup>۳</sup> دارند، شبکه‌های پیچیده هستند [۱۱]. شبکه‌هایی که پیش‌تر مثال زده شد، همگی شبکه‌های پیچیده هستند. با بررسی رفتار شبکه‌ها و پردازش آن‌ها، می‌توان به اطلاعات مفیدی دست یافت. تشخیص جامعه<sup>۴</sup>، بصری‌سازی شبکه<sup>۵</sup>، بررسی ساختار شبکه<sup>۶</sup> و پیش‌بینی پیوند<sup>۷</sup> از جمله پردازش‌هایی است که می‌تواند روی یک شبکه صورت گیرد [۴].

شبکه‌های پیچیده را می‌توان با داده‌ساختار گراف نشان داد. مدل کردن شبکه‌های پیچیده با استفاده از گراف این مزیت را فراهم می‌کند که از الگوریتم‌های پردازش گراف برای تحلیل و بررسی این شبکه‌ها استفاده کرد.

## ۲-۲ شبکه‌های علمی

در این مقاله، به شبکه‌هایی که متشکل از عناصر علمی مانند مقالات، نویسندگان، کلمات کلیدی، انتشارات و دانشگاه‌ها به همراه روابط بین آن‌ها هستند شبکه‌های علمی گفته می‌شود. شبکه‌های علمی چندان عبارت پر کاربردی در ادبیات پژوهش نیست و معمولاً در هر پژوهش از عناوین خاص منظوره برای نامیدن شبکه مورد بحث مانند شبکه همکاری نویسندگان<sup>۸</sup> [۱۲] و شبکه‌ی ارجاعات (استنادات)<sup>۹</sup> [۸، ۱۰] استفاده شده است. ولی با توجه به تعدد این شبکه‌ها که وجه مشترک همه‌ی آن‌ها «علم محور» بودن

ساختار مقاله در ادامه بدین صورت است. در بخش دوم تعریف‌ها، اصطلاحات و مفاهیم پایه مربوط به شبکه‌های پیچیده، شبکه‌های علمی و پیش‌بینی پیوند شرح داده می‌شود و سپس چند الگوریتم و روش حل مسئله پیش‌بینی پیوند که مبنای روش‌های پیشنهادی و آزمایشات هستند بررسی خواهد شد. در بخش سوم کارهای پیشین در زمینه پیش‌بینی پیوند در شبکه‌های مبتنی بر مفاهیم علمی مورد بررسی قرار می‌گیرند. در بخش چهارم روش‌های پیشنهادی ارائه می‌شود و در فصل پنجم آزمایش‌هایی برای مقایسه روش‌های رایج موجود و روش‌های پیشنهادی طراحی می‌شوند. در بخش ششم، نتیجه‌گیری مقاله ارائه خواهد شد.

## ۲- مفاهیم پایه

ابتدا مفهوم شبکه‌های پیچیده بررسی شده و در ادامه، نوع خاصی از شبکه‌های پیچیده که به عنوان شبکه‌های علمی در این مقاله از آن یاد می‌شود، شرح داده می‌شود. سپس تعریف پیش‌بینی پیوند به طور کلی و کاربردهای آن در شبکه‌های علمی مرور می‌شوند. پس از آن رویکردها و روش‌های مورد استفاده در پیش‌بینی پیوند شرح داده شده است. در انتهای این بخش دو الگوریتم جاسازی کلمات برای تحلیل محتوای متنی و جاسازی رأس برای تحلیل ساختاری گراف تشریح شده است.

## ۲-۱ شبکه‌های پیچیده

امروزه شبکه‌ها به یکی از مدل‌های مهم برای بازنمایی<sup>۱</sup> بسیاری از مسائل تبدیل شده‌اند. سیستم‌های اجتماعی، زیستی و اطلاعاتی بسیاری را می‌توان با شبکه‌ها توصیف کرد. رأس‌های این شبکه‌ها کاربران، کامپیوترها، افراد، عناصر زیستی (پروتئین، ژن و غیره) و مانند این‌ها هستند. یال‌ها نشانگر ارتباطات یا تعاملات بین رأس‌ها

<sup>2</sup> Non-trivial

<sup>3</sup> lattice graph

<sup>4</sup> Community detection

<sup>5</sup> Network visualization

<sup>6</sup> Network structure analysis

<sup>7</sup> Link Prediction

<sup>8</sup> Co-authorship network

<sup>9</sup> Citation network

<sup>1</sup> representation

است از عبارت «شبکه‌های علمی» برای نامیدن این شبکه‌ها استفاده می‌کنیم.

در شکل ۱ یک شبکه علمی وجود دارد که متشکل از مجلات<sup>۱</sup> علمی و روابط بین آن‌هاست. در صورتی که مقاله‌ای از یک مجله به مقاله دیگری از مجله دیگر استناد کرده باشد، بین این دو مجله یال برقرار خواهد شد؛ برای مشخص تر شدن شکل تمام اطلاعات نشان داده نشده است. همچنین رأس‌های هم‌رنگ زمینه علمی یکسانی دارند و رأس‌های بزرگ‌تر مجلاتی هستند که تعداد مقاله بیشتری منتشر کرده‌اند.

نشان داده نشده در مجله محاسبات نرم

---

<sup>1</sup> Journal



معیارهای مشابهت محلی، تنها دو رأس هدف و دو لایه از همسایه های آن‌ها را در نظر می‌گیرند و بقیه‌ی قسمت‌های گراف تأثیری در نتیجه ندارند. به همین دلیل از روش‌های سراسری و شبه محلی پیچیدگی زمانی و مکانی کمتری دارند. با این وجود دقت قابل قبولی دارند و پژوهش‌ها نشان می‌دهد نسبت به معیارهای شبه محلی و سراسری کارآمدتر هستند [۴].

با توجه به حجم داده زیاد شبکه‌های علمی (بیش از صد هزار رأس و یک میلیون یال) مشابهت‌های شبه محلی و سراسری از نظر زمانی چندان کارآمد و مقیاس پذیر نیستند. به همین دلیل، روش‌های پیشنهادی و آزمایشات مبتنی بر مشابهت در بخش‌های بعدی از مشابهت‌های محلی استفاده می‌کنند و منابع مطالعه شده نیز از این دسته مشابهت‌ها هستند.

در ادامه  $S(x, y)$  به معنای مشابهت بین رأس  $x$  و  $y$  است، نماد  $\Gamma_x$  نمایانگر همسایه‌های رأس  $x$  می‌باشد و  $w(x, y)$  به معنای وزن یالی است که بین رؤس  $x$  و  $y$  قرار دارد.

#### همسایه مشترک<sup>۱</sup>

ساده‌ترین معیار مشابهت در پیش‌بینی پیوند همسایه مشترک است، همسایه مشترک به صورت تعداد همسایه‌های مشترک دو رأس تعریف می‌شود. این معیار با وجود سادگی، در کاربردهای دنیای واقعی نتیجه مطلوبی کسب می‌کند و پایه معیارهای پیچیده تر است [۴]. در این معیار فرض می‌شود هر چقدر دو رأس همسایه های مشترک بیشتری داشته باشند، احتمال این که خود نیز با یکدیگر همسایه شوند بیشتر است (معادله ۱) [۱۶]. در نسخه وزن دار که در معادله ۲ نشان داده شده است، به جای در نظر گرفتن

از دو جنبه می‌توان گراف مورد مطالعه برای پیش‌بینی پیوند را بررسی کرد: ساختار گراف و محتوای گراف. رویکرد ساختاری گراف، فارغ از این که گراف چه پدیده‌ای از دنیای واقعی یا تئوری را بازنمایی می‌کند، صرفاً به ارتباط رأس‌های گراف می‌پردازد. در بسیاری از گراف‌های دنیای واقعی رأس‌ها و در بعضی موارد یال‌ها، داده‌هایی در خود دارند که نوع، ویژگی‌ها و مقادیر رأس یا یال را نشان می‌دهند. در رویکرد محتوایی گراف از این ویژگی‌ها برای پیش‌بینی پیوند استفاده می‌شود.

وزن یال را می‌توان جزئی از ساختار گراف دانست. همچنین می‌توان آن را یک ویژگی دانست که برای هر یال مقدار مشخصی دارد؛ در این صورت جزئی از محتوای گراف محسوب می‌شود. در این مقاله، فرض شده وزن جزئی از ساختار گراف است و معیارهای مشابهت وزن‌دار به عنوان راه حل پیش‌بینی پیوند در گراف‌های وزن‌دار بررسی شده است.

در ادامه هر یک از رویکردهای مبتنی بر ساختار و محتوا شرح داده شده است و در انتها روش‌هایی که رویکردهای ساختاری و محتوایی را ترکیب می‌کنند و از مزایای هر دو بهره می‌برند به طور خلاصه بیان شده است.

#### ۲-۴-۱ رویکرد مبتنی بر ساختار گراف

در رویکردهای مبتنی بر ساختار گراف عموماً از معیارهای مشابهت برای پیش‌بینی پیوند استفاده می‌شود. معیار مشابهت بین دو رأس گراف محاسبه می‌شود و هر چقدر حاصل عددی بزرگتر باشد مشابهت بین دو رأس و در نتیجه احتمال ایجاد پیوند بالاتر است. معیارهای مشابهت به سه دسته محلی، شبه محلی و سراسری تقسیم می‌شوند.

<sup>1</sup> Common neighbors

همسایه‌های کمتری دارد احتمال بیشتری برای برقراری پیوند دارد. تعریف استاندارد معیار ژاکارد در معادله ۵ و تعریف وزن دار آن در معادله ۶ آمده است [۱۷].

$$s_{JA}(x, y) = \frac{|\Gamma_x \cap \Gamma_y|}{|\Gamma_x \cup \Gamma_y|} \quad 5$$

$$s_{JA}(x, y) = \frac{\sum_{z \in \Gamma_x \cap \Gamma_y} w(x, z) + w(z, y)}{\sum_{a \in \Gamma_x} w(a, x) + \sum_{b \in \Gamma_y} w(b, y)} \quad 6$$

### تخصیص منابع

این شاخص از فرایند تخصیص منابع که در شبکه‌های پیچیده صورت می‌گیرد الهام گرفته شده است. شاخص تخصیص منابع، انتقال واحدهای منابع بین دو رأس غیر متصل  $x$  و  $y$  از طریق همسایگان آن‌ها را مدل می‌کند، هر رأس همسایه یک واحد از منابع را از رأس  $x$  دریافت می‌کند و آن را به صورت مساوی بین همسایگانش توزیع می‌کند. میزان منابع دریافت شده توسط  $y$  به عنوان مشابهت بین این دو رأس در نظر گرفته می‌شود [۱۹]. در نسخه وزن دار این شاخص، توزیع منابع از رأس  $x$  به همسایگان به نسبت وزن یال هر همسایه صورت می‌گیرد [۱۷].

$$s_{RA}(x, y) = \sum_{z \in \Gamma_x \cap \Gamma_y} \frac{1}{|\Gamma_z|} \quad 7$$

$$s_{WRA}(x, y) = \sum_{z \in \Gamma_x \cap \Gamma_y} \frac{w(x, z) + w(z, y)}{\sum_{z' \in \Gamma_z} w(z', z)} \quad 8$$

### اتصال ترجیحی

شاخص اتصال ترجیحی بر این مبنا کار می‌کند که رأس‌هایی که پیوندهای زیادی تشکیل داده‌اند با احتمال بیشتری با یکدیگر پیوند تشکیل می‌دهند. مقدار این شاخص برای دو رأس برابر با حاصلضرب تعداد همسایه‌های دو رأس است [۲۰]. این شاخص بر خلاف شاخص‌های پیشین بر مبنای همسایه‌های مشترک کار

تعداد همسایه‌های مشترک، مجموع وزن همسایه‌های مشترک تا دو رأس  $x$  و  $y$  محاسبه می‌شود [۱۷].

$$s_{CN}(x, y) = |\Gamma_x \cap \Gamma_y| \quad 1$$

$$s_{WCN}(x, y) = \sum_{z \in \Gamma_x \cap \Gamma_y} w(x, z) + w(z, y) \quad 2$$

### آدامیک آدار<sup>۱</sup>

این معیار، بهبود یافته معیار همسایه مشترک است، در معیار آدامیک آدار (معادله ۳)، هر همسایه مشترک با لگاریتم درجه‌اش جریمه می‌شود. یعنی همسایه‌های مشترکی که با رأس‌های زیادی یال برقرار کردند ارزش کمتری نسبت به همسایه‌های مشترکی دارند که با رأس‌های کمتری پیوند دارند [۱۸]. در نسخه وزن دار آدامیک آدار کفایت به جای جریمه کردن با لگاریتم درجه همسایه مشترک، لگاریتم مجموع وزن‌های یال‌های هر همسایه مشترک را جریمه کنیم [۶] (به معادله ۴ رجوع کنید).

$$s_{AA}(x, y) = \sum_{z \in \Gamma_x \cap \Gamma_y} \frac{1}{\log |\Gamma_z|} \quad 3$$

$$s_{WAA}(x, y) = \sum_{z \in \Gamma_x \cap \Gamma_y} \frac{w(x, z) + w(z, y)}{\log \sum_{z' \in \Gamma_z} w(z', z)} \quad 4$$

### ژاکارد

معیار ژاکارد یک معیار قدیمی و بسیار پرکاربرد در زمینه‌های مختلف است که اولین بار برای بررسی مشابهت مجموعه‌ها مورد استفاده قرار گرفت. این معیار نسبت همسایه‌های مشترک به تمام همسایه‌های دو رأس را حساب می‌کند. این معیار نیز بهبود یافته همسایه مشترک است که تعداد همسایه‌های مشترک با استفاده از تعداد کل همسایه‌ها جریمه می‌شوند. در واقع اگر دو جفت رأس تعداد همسایه‌های مشترکی داشته باشند، جفت رأسی که تعداد کل

<sup>1</sup> Adamic-Adar



## ۲-۴-۳ رویکردهای ترکیبی

ساختار و محتوای گراف هر دو داده‌های ارزشمندی برای پیش‌بینی پیوند محسوب می‌شوند. روش‌هایی مانند مشابهت‌ها که فقط ساختار گراف را در نظر می‌گیرند یا روش‌های مبتنی بر محتوا که فقط بر محتوا تمرکز دارند و ساختار را نادیده می‌گیرند در مسائلی که هم محتوا و هم ساختار اهمیت زیادی دارند منجر به نتیجه قابل قبولی نمی‌شوند. به همین خاطر روش‌هایی ارائه شده‌اند که ساختار و محتوا را هم‌زمان در پیش‌بینی پیوند لحاظ می‌کنند.

دسته‌ای از این روش‌ها که به عنوان روش‌های مبتنی بر وزن‌دار کردن شناخته می‌شوند، محتوای گراف را در قالب وزن در ساختار گراف تعبیه می‌کنند و سپس با روش‌های مبتنی بر ساختار مانند معیارهای مشابهت وزن‌دار به حل مسئله پیش‌بینی پیوند در گراف‌های دارای محتوا می‌پردازند [۶, ۱۲].

دسته‌ای از روش‌های دیگر نیز برای هر یال بردارهایی ایجاد می‌کنند که درایه‌های آن مقادیر ویژگی‌های ساختاری و محتوایی گراف است. این بردارها را با استفاده از معیارهای مشابهت برداری یا الگوریتم‌های یادگیری ماشین پردازش می‌کنند و خروجی آن را برای حل مسئله پیش‌بینی پیوند استفاده می‌کنند [۱۰, ۱۴, ۲۱].

## ۲-۵-۵ جاسازی کلمات<sup>۲</sup>

در پردازش زبان طبیعی<sup>۳</sup> دسته‌ای از روش‌ها وجود دارند که هر کلمه را به یک بردار با طول ثابت نگاشت می‌کنند. الگوریتم‌های زیادی هستند که وظیفه نگاشت هر کلمه به یک بردار را انجام می‌دهند که از نظر روش انجام کار با یکدیگر تفاوت دارند ولی وجه مشترک همه این الگوریتم‌ها تبدیل کلمه به یک بردار با طول ثابت

نمی‌کند و معمولاً در پژوهش‌ها نتایج دقیقی از آن حاصل نمی‌شود ولی با توجه به پیچیدگی زمانی کمتر از شاخص‌های مبتنی بر همسایه مشترک، می‌تواند در کنار آن به عنوان مکمل استفاده شود [۴]. نسخه وزن‌دار این شاخص، حاصلضرب مجموع یال‌های دو رأس را به عنوان مقدار شاخص محاسبه می‌کند [۱۷].

$$S_{PA}(x, y) = |E_x| \times |E_y| \quad 9$$

$$S_{WPA}(x, y) = \sum_{x' \in E_x} (x, x') \times \sum_{y' \in E_y} (y, y') \quad 10$$

## ۲-۴-۲ رویکرد مبتنی بر محتوا و زمان

معیارهای مشابهت گفته شده صرفاً ساختار شبکه را در نظر می‌گیرند. منظور، ساختار<sup>۱</sup> دو رأسی است که شباهتشان حساب می‌شود و همسایه‌های آن‌هاست. در صورتی که علاوه بر ساختار گراف، موارد دیگری نیز وجود دارند که می‌توان از آن‌ها برای پیش‌بینی پیوند کمک گرفت. به عنوان مثال در یک شبکه اجتماعی ممکن است دو نفر هیچ همسایه مشترک و در نتیجه مشابهت ساختاری نداشته باشند ولی به واسطه شغلشان تمایل به برقراری پیوند داشته باشند و در آینده پیوند برقرار کنند، در اینجا شغل، محتوای رئوس در نظر گرفته می‌شود و مستقل از ساختار است.

همچنین زمان برقراری پیوند در گراف می‌تواند حائز اهمیت باشد. مثلاً در شبکه همکاری نویسندگان می‌توانیم این فرضیه را مطرح کنیم که هر چقدر زمان به وجود آمدن یک پیوند بیشتر باشد (سابقه همکاری بیشتر باشد)، آن پیوند اهمیت بیشتری دارد. به طور کلی، با ترکیب اطلاعات زمانی و محتوایی با ساختار شبکه، نتیجه پیش‌بینی پیوند بهبود چشمگیری می‌یابد [۶].

<sup>2</sup> Word embedding

<sup>3</sup> Natural Language Processing (NLP)

<sup>1</sup> Topology, Structure

پیش از این پژوهش‌هایی در زمینه پیش‌بینی پیوند در شبکه‌های علمی انجام شده که در هر یک از این پژوهش‌ها یک مسئله روی یک نوع شبکه خاص (داده مورد مطالعه) برای یک کاربرد مشخص تعریف شده و با یک روش پیشنهادی پاسخی به آن مسئله داده شده است. بنابراین، مرور کارهای گذشته از دو جنبه کاربرد و شبکه مورد مطالعه و بررسی می‌شوند.

### ۱-۳ شبکه‌ها و کاربردهای مورد مطالعه

یکی از وجوه تمایز مهم پژوهش‌های انجام شده در زمینه پیش‌بینی پیوند شبکه‌های علمی، شبکه‌ای (گرافی) است که مورد پردازش قرار گرفته است. با مقایسه این پژوهش‌ها یک دسته‌بندی کلی از انواع شبکه‌های مورد پردازش و کاربردهایی از دنیای واقعی که برای پیش‌بینی پیوند روی آن شبکه تعریف شده ارائه می‌شود.

است. این روش‌ها از این جهت سودمند هستند که به راحتی می‌توان عملیات محاسباتی<sup>۱</sup> را روی کلمات، جملات و پرونده‌ها انجام داد. با تبدیل متن به بردار، این امکان فراهم می‌شود که عملیاتی مانند جمع، تفریق، محاسبه فاصله و محاسبه مشابهت را از طریق عملیات محاسباتی برداری انجام دهیم [۲۲].

الگوریتم‌های word2vec<sup>۲</sup> و GloVe<sup>۳</sup> از جمله الگوریتم‌های مطرح در این زمینه هستند که منجر به نتایج خوبی شده‌اند [۲۳]. از آنجایی که بخش بزرگی از محتوای شبکه‌های علمی مانند عنوان مقاله و چکیده در قالب متن می‌باشد، در روش‌های پیشنهادی از الگوریتم‌های جاسازی کلمات برای پردازش محتوای متنی شبکه استفاده خواهد شد.

### ۲-۶ جاسازی رأس‌ها<sup>۳</sup>

پردازش ساختار داده گراف با استفاده از الگوریتم‌های استاندارد یادگیری ماشین و داده‌کاوی کار مشکلی است، به این خاطر که ورودی این الگوریتم‌ها داده‌های عددی ساختار یافته هستند ولی گراف، متشکل از تعدادی رأس و ارتباطات میان این رأس‌هاست [۲۴، ۲۵]. برای ایجاد نگاهت بین داده‌های رابطه‌ای مانند گراف و داده‌های ساختار یافته که قابلیت انجام عملیات محاسباتی روی آن‌ها فراهم باشد الگوریتم‌هایی ارائه شده‌است. یکی از این الگوریتم‌ها که در بسیاری از مسائل گراف مانند پیش‌بینی پیوند، تشخیص جامعه و دسته‌بندی رئوس نتایج خوبی به دست آورده‌است الگوریتم node2vec است [۲۶]. برای مطالعه جزئیات این روش می‌توانید به پیوست الف مراجعه کنید.

### ۳- مرور کارهای گذشته

<sup>1</sup> Arithmetic

<sup>2</sup> Global Vectors for word representation

<sup>3</sup> Node embedding

## شبکه‌ی ارجاعات



شکل ۲ شبکه‌های پر کاربرد در پیش‌بینی پیوند شبکه‌های علمی

جدول ۱: مجموعه داده‌های مورد استفاده در کارهای پیشین

نام مجموعه داده	تعداد رأس‌ها	تعداد پیوندها	نوع شبکه
arXiv API <sup>1</sup>	به روز میشود		شبکه ارجاعات
DBLP <sup>2</sup>	به روز میشود		شبکه ارجاعات
HEP-TH <sup>3</sup>	27,770	352,807	شبکه ارجاعات
SCI- Expanded (Web of Science) <sup>4</sup>	۵۳ میلیون	۱ میلیارد	شبکه ارجاعات
CiteSeer <sup>5</sup>	3,312	4,732	شبکه ارجاعات
ogbl-citation2 <sup>6</sup>	2,927,963	30,561,187	شبکه ارجاعات
ogbl-collab	235,868	1,285,465	شبکه همکاری نویسندگان

<sup>1</sup> <https://info.arxiv.org/help/api/>

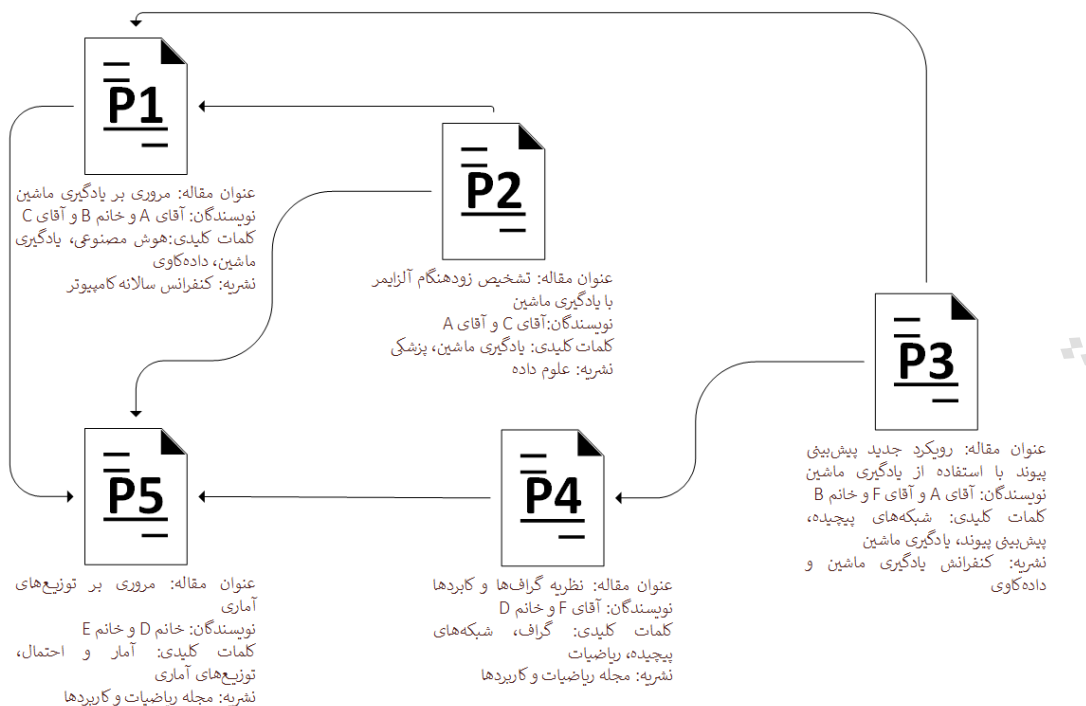
<sup>2</sup> <https://dblp.org/>

<sup>3</sup> <https://snap.stanford.edu/data/cit-HepTh.html>

<sup>4</sup> <https://clarivate.com/webofsciencegroup/solutions/webofscience-scie/>

<sup>5</sup> <https://relational.fit.cvut.cz/dataset/CiteSeer>

<sup>6</sup> <https://ogb.stanford.edu/docs/linkprop/>



شکل ۳ نمونه‌ای از شبکه ارجاعات

دید جامع نسبت به شبکه‌هایی که در این مقاله مورد مطالعه قرار خواهند گرفت کمک خواهد کرد.

با توجه به ویژگیهای مشترک شبکه‌ها، این دسته‌بندی برای یافتن الگوریتم مناسب پیش‌بینی پیوند مفید است. علاوه بر این، به ایجاد دید جامع نسبت به شبکه‌هایی که در این مقاله مورد استفاده خواهند بود کمک خواهد کرد. در جدول ۱ مجموعه شبکه علمی در دنیای واقعی که در کارهای پیشین مورد استفاده قرار گرفتند به همراه تعداد رأس‌ها و پیوندها نام برده شده است. مجموعه داده ogbl-collab شبکه همکاری نویسندگان و بقیه شبکه ارجاعات هستند.

### ۳-۱-۱ شبکه ارجاعات

شبکه ارجاعات، جامع ترین شبکه علمی است که در آن مقاله‌ها رأس‌ها و ارجاعات یال‌های شبکه هستند. اگر مقاله‌ای به مقاله دیگری ارجاع داده باشد، میان دو رأس متناظر یال برقرار می‌شود.

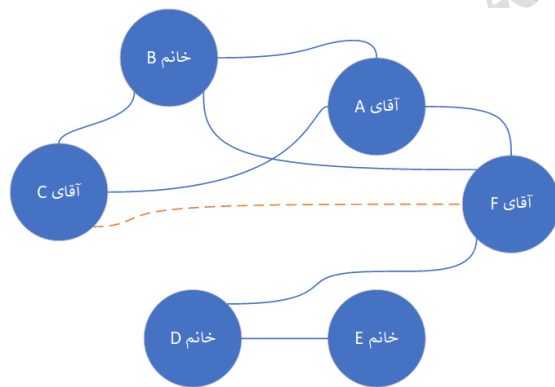
جامع ترین شبکه علمی موجود در دسترس، شبکه ارجاعات است که در صورت کامل بودن فرا داده‌های آن (نام نویسندگان، کلمات کلیدی و غیره) می‌توان برخی شبکه‌های علمی دیگر مانند شبکه کلمات کلیدی و شبکه همکاری نویسندگان را بر اساس آن ایجاد کرد (شکل ۲). مواردی که در شکل با رنگ آبی مشخص شدند در مقالات مطالعه شده صریحاً مورد استفاده قرار گرفته‌اند و مواردی که با رنگ خاکستری نشان داده شده‌اند، مواردی هستند که در پژوهش‌های پیشین مورد مطالعه ما به آن‌ها اشاره نشده ولی از طریق مجموعه داده‌های موجود قابل به دست آمدن هستند و می‌توانند بالقوه مفید باشند. در بخش‌های بعدی هر یک از شبکه‌هایی که در پژوهش‌های پیشین از آن استفاده شده بررسی می‌شود.

با توجه به ویژگیهای مشترک شبکه‌ها، این دسته‌بندی برای یافتن الگوریتم مناسب پیش‌بینی پیوند مفید است. علاوه بر این، به ایجاد

یال برقرار می‌شود. در نهایت شبکه نویسندگان به صورت شکل ۴ ساخته می‌شود.

در مقاله‌های [۳، ۹] پیش‌بینی همکاری نویسندگان با استفاده از یک روش ترکیبی مبتنی بر وزن‌دار کردن شبکه همکاری با مشابهت محتوایی و سپس پیش‌بینی پیوند با معیارهای شباهت وزن‌دار با استفاده از شبکه‌های همکاری نویسندگان استخراج شده از HEP-TH و arXiv انجام شده است. در مقاله [۲۷] نیز با افزودن کلمات کلیدی به شبکه همکاری نویسندگان مستخرج از شبکه SCI-Expanded و سپس اعمال الگوریتم قدم‌زدن تصادفی، پیش‌بینی انجام شده است.

پیش‌بینی پیوند در شبکه همکاری نویسندگان، افرادی را به هم معرفی می‌کند که می‌توانند در پژوهش‌های آتی همکاری مؤثر داشته باشند. برای مثال پیش‌بینی همکاری آتی بین «آقای F» و «آقای C» در شکل ۴ با نقطه چین قرمز مشخص شده است.



شکل ۴ شبکه نویسندگان مبتنی بر شکل ۳

### ۳-۱-۳ شبکه کلمات کلیدی

رأس‌های این شبکه، کلمات کلیدی مقالات هستند و در صورتی که دو کلمه کلیدی در یک مقاله استفاده شده باشند، بین آن‌ها یال ایجاد می‌شود. همان‌طور که پیش‌تر اشاره شد در صورت وجود

مجموعه داده‌هایی مانند DBLP در دسته شبکه ارجاعات قرار می‌گیرند. در شکل ۳ نمونه‌ای از شبکه ارجاعات مشاهده می‌شود. در این شبکه، مقاله P2 به مقاله P1 ارجاع داده است، بنابراین یک یال جهت‌دار از مقاله P2 به P1 در شبکه تشکیل می‌شود. به همین ترتیب سایر یال‌های شبکه بر اساس ارجاعات بین مقالات ساخته می‌شوند.

در مقاله [۱۲]، پیش‌بینی ارجاع متقابل با استفاده از یادگیری ماشین روی ویژگی‌های ساختاری و متنی شبکه ارجاعات Citeseer انجام شده است. پیش‌بینی تعداد ارجاعات به یک مقاله نیز بر اساس ویژگی‌های ساختاری در مجموعه داده HEP-TH با استفاده از یادگیری ماشین در مقاله [۶] انجام شده است. در مقاله [۵] روی شبکه ارجاعات TEP-TH، برای مسئله یافتن مقالات مرتبط با یک مقاله، یک راه حل با استفاده از معیارهای شباهت محلی ارائه شده است. مقاله [۷] نیز با استفاده از tf-idf چکیده مقالات و ساختار شبکه، روی پنج پایگاه داده ارجاعات SCI-Expanded پیش‌بینی پیوند انجام داده است.

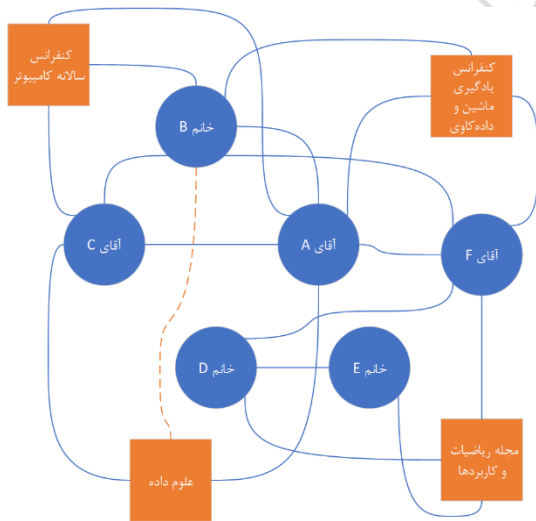
### ۳-۱-۳ شبکه همکاری نویسندگان

در این شبکه، نویسندگان رأس‌های شبکه را تشکیل می‌دهند و در صورت انجام همکاری بین دو نویسنده، بین آن‌ها یال تشکیل می‌شود. شبکه همکاری نویسندگان می‌تواند وزن‌دار باشد؛ وزن هر یال نشان دهنده تعداد همکاری دو نویسنده است.

با توجه به شبکه ارجاعات شکل ۳، «آقای A»، «خانم B» و «آقای C» نویسندگان همکار در مقاله P1 هستند. به همین خاطر بین آن‌ها یال برقرار می‌شود، به همین ترتیب برای سایر مقالات، نویسنده‌ها به گراف اضافه می‌شوند و بین نویسندگان مشترک یک مقاله

دسته دیگری از شبکه‌های علمی شبکه‌های ناهمگن<sup>۱</sup> که رأس‌های آن موجودیت‌هایی از انواع مختلف هستند. برای مثال موجودیت نشریات، نویسندگان، کلیمات کلیدی و ارتباط بین آن‌ها نمونه‌ای از شبکه علمی ناهمگن است (شکل ۶). در این شبکه دو موجودیت نویسنده و نشریه وجود دارد که نویسندگان با یکدیگر و نشریات ارتباط دارند.

پیش‌بینی پیوند در شبکه‌های ناهمگن می‌تواند بین رأس‌های هم نوع یا متفاوت صورت گیرد. در مقاله [۲۷] با افزودن کلمات کلیدی به شبکه همکاری نویسندگان، یک شبکه ناهمگن ایجاد می‌شود که باعث بهبود نتیجه پیش‌بینی پیوند بین نویسندگان شده است. در شکل ۶ نقطه چین قرمز نشان دهنده یک پیش‌بینی پیوند بین نویسنده و نشریه است؛ نویسنده مورد نظر در صورت چاپ مقاله خود در نشریه پیش‌بینی شده می‌تواند موفقیت بیشتری کسب

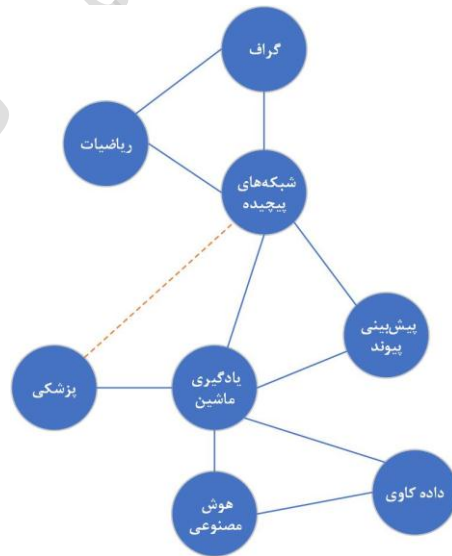


شکل ۶ یک گراف ناهمگن متشکل از نویسنده‌ها و نشریات و ارتباط بینشان

### ۲-۳ روش‌های حل مسئله پیش‌بینی پیوند

فراداده‌های مربوط به کلمات کلیدی مقاله در شبکه ارجاعات، می‌توان شبکه کلمات کلیدی از روی شبکه ارجاعات ایجاد کرد.

در مقاله [۱۴] برای پیش‌بینی روند تولید علم، با یک روش یادگیری ماشین و انتخاب ویژگی، روی شبکه کلمات کلیدی پیش‌بینی پیوند انجام شده است. به عنوان مثال در شکل ۵ در صورتی که الگوریتم پیش‌بینی پیوند، ایجاد پیوند بین دو کلمه کلیدی «شبکه‌های پیچیده» و «پزشکی» را پیش‌بینی کند می‌توان انتظار داشت در آینده مقالاتی که در زمینه ترکیب شبکه‌های پیچیده و پزشکی منتشر شوند مورد توجه قرار خواهند گرفت.



شکل ۵ شبکه‌ی کلمات کلیدی

### ۳-۱-۴ سایر شبکه‌های علمی

با توجه به داده‌هایی که در شبکه ارجاعات وجود دارد، می‌توان شبکه‌های دیگری ایجاد کرد که در منابع مطالعه شده به آن‌ها پرداخته نشده. شبکه نشریات و شبکه دانشگاه‌ها مثال‌هایی هستند که در شکل ۲ با رنگ خاکستری مشخص شده‌اند. ایجاد این شبکه‌ها و انجام پیش‌بینی پیوند در آن‌ها می‌تواند ارزشمند باشد.

<sup>1</sup> non-homogeneous

موارد آبی رنگ صریحاً در منابع مطالعه شده در زمینه شبکه‌های علمی استفاده شده‌اند، موارد سفید رنگ در مقالات مروری مربوط به پیش‌بینی پیوند به کار رفته‌اند [۴، ۵] و برای کامل بودن نمودار آورده شده‌اند.

در این بخش مروری بر روش‌های پرکاربرد پیش‌بینی پیوند در شبکه‌های علمی خواهد شد و مزایا و معایب هر روش بیان می‌شود. در شکل ۷ یک دسته‌بندی از راه‌حل‌های پیش‌بینی پیوند مورد استفاده در مقالات مروری و منابع مربوط به پیش‌بینی پیوند در شبکه‌های علمی مورد مطالعه آورده شده‌است. در این دسته‌بندی



شکل ۷ راه‌حل‌های مسئله پیش‌بینی پیوند در شبکه‌های علمی

### ۱-۲-۳ مبتنی بر مشابهت

روش‌های مبتنی بر مشابهت، یکی از رایج‌ترین روش‌های پیش‌بینی پیوند هستند. جفت رأس‌هایی که امکان وجود پیوند برای آن‌ها قرار است بررسی شود، از بین تمام جفت رأس‌های موجود که پیوندی بین آن‌ها وجود ندارد مقدار معیار مشابهت حساب می‌شود. جفت رأسی که مقدار مشابهت‌شان بیشتر از حد آستانه باشد به عنوان پیوند پیش‌بینی شده در نظر گرفته می‌شوند.

در مقاله [۸] مقدار آستانه به صورت پویا محاسبه شده است، به این ترتیب که در شبکه اولیه مشابهت تمام جفت رئوسی که پیوند بین آن‌ها برقرار است محاسبه می‌شود و میانگین این مشابهت به عنوان حد آستانه در نظر گرفته می‌شود. در پژوهش [۴] مشابهت‌ها بر اساس زیاد به کم مرتب شده و  $K$  یال که بالاترین مشابهت را دارد به عنوان پیوند پیش‌بینی شده انتخاب شده‌اند. همچنین در

مقالات [۱۰، ۱۴] از معیارهای مشابهت محلی به همراه داده‌های ساختاری و محتوایی دیگر برای ورودی الگوریتم یادگیری ماشین استفاده شده ولی مستقیماً مورد استفاده قرار نگرفته است.

روش مشابهت محلی ساده‌ترین روش پیش‌بینی پیوند است. از جمله معایب این روش می‌توان به عدم امکان استفاده همزمان از چند معیار مشابهت اشاره کرد. از طرف دیگر معیارهای مشابهت صرفاً روی ساختار گراف تعریف شده‌اند و در صورت نیاز به ترکیب زمان و محتوا با ساختار، این روش غیر قابل استفاده است. این معایب تا حدودی در روش وزن‌دار کردن و به طور کامل در روش مبتنی بر یادگیری ماشین رفع شده‌است. در ادامه کارهای پیشین مبتنی بر روش‌های وزن‌دار کردن و یادگیری ماشین بررسی شده‌اند.

## ۳-۲-۲ مبتنی بر وزن دار کردن

شبکه‌های وزن دار در همه کاربردهای دنیای واقعی از جمله شبکه‌های اجتماعی می‌تواند وجود داشته باشد و منحصر به شبکه‌های علمی نیست ولی از آنجایی که برخی از پژوهش‌ها در زمینه شبکه‌های علمی منحصراً روی شبکه‌های وزن دار پژوهش است، در این بخش جداگانه به آن‌ها پرداخته می‌شود. در پژوهش‌های پیشین مربوط به شبکه‌های علمی وزن دار، برخی از شبکه‌ها ذاتاً وزن دار بوده‌اند و در برخی دیگر با استفاده از روشی یک وزن به یال‌ها نگاشت شده است. مقاله [۲۸] نشان می‌دهد با حذف وزن یال‌ها از گراف‌های وزن دار، نتیجه به صورت چشم‌گیری تغییر می‌کند و دقت کاهش می‌یابد.

در روش‌های مبتنی بر وزن دار کردن، با در نظر گرفتن ساختار، محتوا و برخی موارد نیز زمان، برای هر رأس یک وزن محاسبه می‌کنند و سپس با استفاده از معیارهای مشابهت وزن دار مانند آنچه پیش‌تر گفته شد، وجود یا عدم وجود پیوند بین رئوس را پیش‌بینی می‌کنند. در این روش، مشابهت ساختاری، مشابهت محتوایی و زمان در یک عدد که وزن یال است خلاصه می‌شود.

یک نمونه از رابطه‌های وزن دار کردن گراف در معادله ۱۱ قابل مشاهده است، توسط مقاله [۶] روی شبکه همکاری نویسندگان ارائه شده است، این رابطه از ضرب سه مقدار در یکدیگر محاسبه می‌شود.  $\alpha$  و  $\beta$  دو ضریب هستند به طوری که  $\alpha + \beta = 1$ ؛ با افزایش (کاهش) تأثیر محتوا، تأثیر زمان کم (زیاد) می‌شود.  $|E(u, v)|$  وزنی است که یال گراف از ابتدا داشته، در گراف‌های بدون وزن این مقدار ۱ است. در قسمت دوم

$$\beta \frac{CTime - \max(t_{(u,v)})}{CTime - \min(t)}$$

زمان دیده شده در شبکه (زمان به وجود آمدن شبکه) و

$\max(t_{(u,v)})$  از بین زمان به وجود آمدن رأس  $u$  و  $v$  ماکسیمم را انتخاب می‌کند. دلیل این که  $CTime - \max(t_{(u,v)})$  بر  $CTime - \min(t)$  تقسیم شده، این است که مقدار کل نرمال شده و عددی بین ۰ و ۱ شود. به این ترتیب با ضرب این سه مقدار عددی به دست می‌آید که از ساختار، زمان و محتوا تأثیر گرفته است و به عنوان وزن یال  $u, v$  در نظر گرفته می‌شود. در ادامه با استفاده از معیارهای مشابهت وزن دار، مشابهت بین  $u, v$  محاسبه می‌شود و اگر بالاتر از حد آستانه باشد، بین آن‌ها پیوند برقرار می‌شود. نتایج پژوهش نشان می‌دهد با این روش، نتایج نسبت به هنگامی که فقط ساختار یا فقط محتوا در نظر گرفته شده بهبود می‌یابد.

$$w^{CTT}(u, v) = |E(u, v)| * \beta \frac{CTime - \max(t_{(u,v)})}{CTime - \min(t)} * \alpha^{1 - \cos(u,v)} \quad 11$$

در مقاله [۱۲] نیز با یک معیار مشابهت جدید به نام LDACosine مشابهت کسینوسی دو رأس با در نظر گرفتن ساختار و محتوا محاسبه می‌شود و نتایج بهتری نسبت به روش‌های بدون محتوا به دست آمده است.

## ۳-۲-۳ مبتنی بر یادگیری ماشین

در پژوهش‌هایی که مسئله پیش‌بینی پیوند را با الگوریتم‌های یادگیری ماشین حل کردند یک رویکرد کلی مشاهده می‌شود. در همه این پژوهش‌ها برای هر یال تعدادی ویژگی محاسبه شده و هر یال به یک بردار چند بعدی تبدیل شده، هر یال با توجه به این که در گراف آموزشی وجود دارد یا خیر برچسب وجود یا عدم وجود پیوند می‌گیرد و به این ترتیب مسئله پیش‌بینی پیوند به یک مسئله دسته‌بندی تبدیل می‌شود.

در مقاله [۱۰] روی شبکه ارجاعات، برای یال‌های ممکن یک بردار متشکل از ۱۱ ویژگی محاسبه شده، برخی از این ویژگی‌ها مانند



الگوریتم آموزش داده شده پیوندهای مثبت و منفی پیش‌بینی می‌شوند و بر اساس معیار مد نظر ارزیابی می‌شوند.

جدول ۲ ویژگی‌های پیشنهادی شبکه‌های ارجاعات و همکاری نویسنندگان

نوع	شبکه قابل اعمال		نوع	نام ویژگی	نوع
	همکاری	ارجاعات			
int	*	*	CN	نزدیک‌ترین همسایه	ساختاری
float	*	*	AA	آدامیک آدار	
float	*	*	RA	تخصیص منابع	
float	*	*	JC	ژاکارد	
int	*	*	PA	اتصال ترجیحی	
int	*	*	W_CN	نزدیک‌ترین همسایه وزن‌دار	
float	*	*	W_AA	آدامیک آدار وزن‌دار	
float	*	*	W_RA	تخصیص منابع وزن‌دار	
float	*	*	W_JC	ژاکارد وزن‌دار	
int	*	*	W_PA	اتصال ترجیحی وزن‌دار	
float	*	*	N2V_COS_SIM	شباهت کسینوسی بردار Node2Vec جفت رأس	محتوایی
float	*	*	TITLE_COS_SIM	شباهت کسینوسی عنوان مقالات دو نویسنده	
float		*	ABS_COS_SIM	شباهت کسینوسی چکیده دو مقاله	زمانی
int		*	PUB_YEAR_DIFF	اختلاف سال انتشار دو مقاله	

همسایه مشترک و ژاکارد معیارهای ساختاری و برخی مانند اختلاف سال انتشار ۲ مقاله و تعداد نویسندگان مشترک ویژگی‌هایی محتوایی بودند. در انتها بردارها به الگوریتم ماشین بردار پشتیبان داده شده و نتیجه نشان دهنده بهبود نسبت به روش‌های پیشین روی مجموعه داده‌های بررسی شده است.

همچنین در پژوهشی که روی شبکه کلمات کلیدی انجام گرفته [۱۴] و از میان ۱۲ ویژگی محاسبه شده برای یال‌ها، ویژگی‌های ژاکارد، تخصیص منابع<sup>۱</sup>، میانگین ضریب خوشگی<sup>۲</sup> و میانگین مرکزیت<sup>۳</sup> دو رأس انتخاب شدند؛ بهترین نتیجه با الگوریتم جنگل تصادفی به دست آمده است.

#### ۴- روش‌های پیشنهادی

بر اساس شکل ۸، در روش‌های پیشنهادی گراف مورد پردازش ابتدا وزن‌دهی می‌شود، ابتدا ویژگی‌های ساختاری و محتوایی جدول ۲ برای جفت رأس‌های آن‌ها حساب می‌شود. سپس با استفاده از یک الگوریتم انتخاب یا استخراج ویژگی مانند PCA داده‌ها پردازش می‌شوند و ابعاد زائد حذف می‌شوند. سپس به دو قسمت آموزشی و آزمایشی تقسیم می‌شود. گراف آموزشی برای آموزش الگوریتم یادگیری ماشین استفاده می‌شود. در ادامه برای داده‌های منفی، جفت رأسهایی که در داده‌های آموزشی پیوندی بین آنها نیست، (و یا مثبت و منفی) خوشه بندی صورت می‌گیرد و مراکز خوشه‌ها به عنوان داده آموزشی برای آموزش الگوریتم یادگیری ماشین مورد استفاده قرار می‌گیرند. از سوی دیگر ویژگی‌ها برای داده‌های آزمایشی محاسبه شده و بر اساس

<sup>1</sup> Resource Allocation

<sup>2</sup> Cluster Coefficient

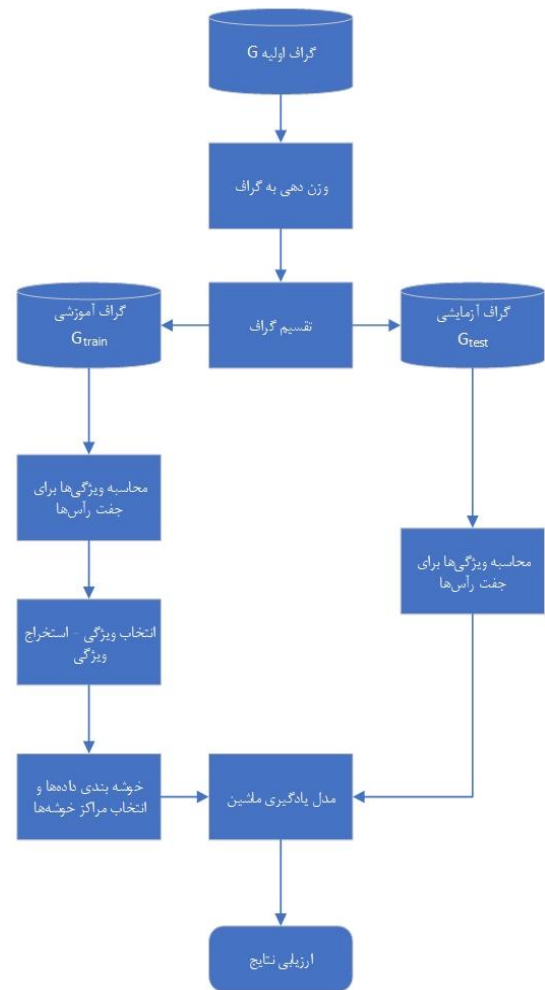
<sup>3</sup> Centrality

پیش بینی پیوند با استفاده از یادگیری ماشین و دسته بندی شرح داده شده است. در انتها برای حل مشکل نمونه گیری منفی دو روش خوشه بندی و رگرسیون پیشنهاد شده است.

#### ۴-۱ وزن دهی شبکه

با وجودی که در روش یادگیری ماشین هر ویژگی محتوایی را به عنوان یک درایه از بردار ویژگی در نظر می گیریم ولی با این حال در روش پیشنهادی از ویژگی های محتوایی برای وزن دهی به گراف نیز استفاده می کنیم. این کار به دقیق تر شدن نتایج شباهت های محلی وزن دار و در نتیجه افزایش دقت یادگیری ماشین کمک می کند. همچنین در صورتی که از روش یادگیری ماشین استفاده نشود، ویژگی های محتوایی در وزن یال ها در نظر گرفته شده و نادیده گرفته نخواهد شد.

با این فرضیه که در شبکه های علمی پیوندهای جدیدتر اهمیت بیشتری نسبت به پیوندهای قدیمی تر دارند و پیوندهایی با شباهت محتوایی بیشتر باید بیش از پیوندهایی با شباهت محتوایی کمتر مورد توجه قرار بگیرند، دو معیار وزن دهی برای شبکه همکاری نویسندگان (معادله ۱۲) و شبکه ارجاعات (معادله ۱۳) پیشنهاد شده است. هر دو معیار وزن دهی پیشنهاد شده بر اساس ترکیب خطی تأثیر زمان ایجاد پیوند و شباهت محتوایی پیوند است.



شکل ۸ چارچوب پیشنهادی پیش بینی پیوند مبتنی بر یادگیری ماشین

بر اساس روش پیشنهادی، در ادامه روش پیشنهادی برای وزن دهی شبکه، سپس ویژگی هایی برای جفت رأس های شبکه همکاری نویسندگان و ارجاعات پیشنهاد شده است. سپس روش

$$weight_{collab}(x, y) = \alpha \left( cosine\_sim(\overrightarrow{w2v(x)}, \overrightarrow{w2v(y)}) \right) + \beta \left( \sum_{i \in collab(x, y)} \frac{1}{\log(year_{current} - year_{collab}(i))} \right) \quad 12$$

$$weight_{citation}(x, y) = \alpha \left( cosine\_sim(\overrightarrow{w2v(x)}, \overrightarrow{w2v(y)}) \right) + \beta \left( \frac{1}{\log|year_x - year_y|} \right) \quad 13$$

$$\alpha, \beta > 0 \text{ and } \alpha + \beta = 1$$

می‌دهد؛ ویژگیهایی که بازه مقداری برای آنها مشخص نشده است می‌توانند هر مقداری بزرگتر از صفر داشته باشند. ویژگیهایی که در جدول ۲ آمده است، ویژگیهایی هستند که محاسبه آنها از نظر زمان اجرا مقیاس‌پذیر است. ویژگیهای دیگری مانند معیارهای مرکزیت وجود دارند اما با افزایش اندازه گراف زمان محاسبه آنها به صورت نمایی افزایش می‌یابد. به همین دلیل چنین ویژگی‌هایی در اینجا استفاده نشده است.

#### ۳-۴ پیش‌بینی پیوند با استفاده از یادگیری ماشین

گراف  $G = (V, E)$  به طوری که  $V$  مجموعه رأس‌های گراف؛ در اینجا نویسندگان یا مقالات و  $E$  پیوندهای موجود میان رأس‌های شبکه می‌باشد. مجموعه  $E'$ ، مجموعه پیوندهایی است که در گراف وجود ندارد ولی به صورت بالقوه امکان ایجاد در آینده را دارد. برای پیش‌بینی پیوند با استفاده از یادگیری ماشین مدل  $M$  آموزش داده می‌شود به طوری که هر پیوند عضو  $E$  را به یک احتمال بین صفر و یک نگاشت می‌کند. مدل  $M$  در این روش پیشنهادی می‌تواند یک الگوریتم دسته‌بندی یا رگرسیون باشد.

برای حل مسئله‌ی پیش‌بینی پیوند با استفاده از یادگیری ماشین، گراف باید در قالب مجموعه‌ای از بردارها بازنمایی شود. این بردارها برای هر جفت رأس محاسبه می‌شوند درایه‌های این بردار می‌تواند ویژگی‌های ساختاری یا محتوایی مربوط به جفت رأس باشد. مجموعه‌ای از ویژگی‌های قابل محاسبه برای بردار ویژگی‌ها در جدول ۲ آورده شده است.

بعد از آنکه برای هر جفت رأس، بردار ویژگی حساب شد، یک برچسب وجود یا عدم وجود پیوند برای آن مشخص می‌شود. مدل یادگیری ماشین این برچسب را برای جفت رأس‌های جدید پیش‌بینی می‌کند. به طور خلاصه مراحل در زیر بیان شده است:

در هر دو معادله  $\alpha$  و  $\beta$  هر دو بزرگتر از صفر و مجموع آنها برابر یک است. با افزایش  $\alpha$  و در نتیجه کاهش  $\beta$  تأثیر محتوا نسبت به زمان افزایش پیدا می‌کند. همچنین  $\overrightarrow{w2v(x)}$  یک بردار ۱۲۸ بعدی از نهفته‌سازی محتوای متنی هر رأس است. همچنین  $collab(x, y)$  مجموعه همکاری‌های دو نویسنده  $x$  و  $y$  است و  $year_{collab}(i)$  سال انجام همکاری  $i$  بین این دو نویسنده است. علت استفاده از تابع لگاریتم در جریمه اختلاف سال، کم کردن تأثیر اختلاف سال زیاد ارجاع به مقاله یا همکاری است. مثلاً نسبت اختلاف ۱۰ سال و ۲۰ سال بدون استفاده از لگاریتم ۲ برابر است ولی در دنیای واقعی چندان تفاوتی بین این دو وجود ندارد، با استفاده از لگاریتم این نسبت به ۱,۳ کاهش می‌یابد.

#### ۲-۴ تعریف ویژگی‌ها

الگوریتم‌های یادگیری ماشین برای پیش‌بینی یک برچسب یا امتیاز نیاز به داده‌های آموزشی دارند که در قالب یک بردار  $n$  بعدی به آنها داده می‌شود. این ویژگی‌ها می‌توانند مبتنی بر ساختار شبکه، مبتنی بر محتوای رأس‌ها و یال‌ها و/یا مبتنی بر زمان ایجاد یک رأس یا یال باشد. با توجه به ساختار و محتوای شبکه همکاری نویسندگان و شبکه ارجاعات، مجموعه‌ای از ویژگی‌هایی که در وجود یا عدم وجود یک پیوند نقش دارد، در جدول ۲ ارائه شده است. نوع ویژگی (ستون اول) مشخص می‌کند که ویژگی عنوان شده ساختاری، محتوایی و یا زمانی است. ستون دوم جدول نام ویژگی است و ستون سوم علامت اختصاری است که در نتایج آزمایشها از آن استفاده خواهد شد. از آنجا که در شبکه همکاری نویسندگان متن مقاله وجود ندارد، برخی ویژگی‌ها در این شبکه قابل استفاده نیست؛ مانند شباهت محتوایی دو مقاله. بنابراین در شبکه همکاری نویسندگان تنها از شباهت میان عنوان مقاله استفاده شده است. آخرین ستون جدول ۲، نوع هر ویژگی را نشان

گراف‌های دنیای واقعی معمولاً پراکنده<sup>۱</sup> هستند؛ یعنی تعداد پیوندهایی که وجود دارند بسیار بیشتر از تعداد پیوندهایی هستند که وجود ندارند. مثلاً در گراف همکاری نویسندگان ogbl-collab که روش پیشنهادی روی آن آزمایش خواهد شد، میانگین درجه‌ی رئوس ۸ و تعداد رئوس ۲۳۵,۸۶۸ است. یعنی هر نویسنده به طور میانگین با ۸ نویسنده همکاری کرده و با ۲۳۵,۸۶۰ نویسنده همکاری نکرده است.

از طرفی برای آموزش مدل یادگیری ماشین به تعدادی نمونه‌ی منفی و تعدادی نمونه‌ی مثبت نیاز است. با توجه به پراکنده بودن گراف مورد پردازش، می‌توان تمام نمونه‌های مثبت را به عنوان داده‌ی آموزشی به مدل یادگیری ماشین در نظر گرفت ولی تعداد نمونه‌های منفی در این گراف بسیار زیاد است. به عنوان مثال، در مجموعه داده ogbl-collab تعداد ۲۳۵,۸۶۸ رأس داریم که بر اساس معادله ۱۴ تعداد رأس‌های گراف کامل، تقریباً ۲۸ میلیارد یال ممکن وجود دارد که از این بین ۱,۲۸۵,۴۶۵ یال مثبت و سایر یال‌ها منفی هستند.

$$\begin{aligned} \text{complete\_graph\_edges} &= \frac{n \times (n - 1)}{2} \\ &= \frac{235,868 \times 235,867}{2} \\ &\cong 28 \times 10^9 \end{aligned} \quad 14$$

محاسبه ویژگی‌ها برای این تعداد داده‌ی آموزشی منفی بسیار زمان‌بر است و نیاز به منابع سخت‌افزاری زیادی دارد، همچنین الگوریتم‌های یادگیری ماشین معمول (غیر از آن‌هایی که برای داده‌های حجیم پیاده‌سازی شدند) برای پردازش این تعداد داده نیاز به زمان و منابع سخت‌افزاری زیادی دارند. از طرف دیگر، بسیاری از این داده‌ها تکراری هستند و اطلاعات تازه‌ای در اختیار

۱ ویژگی‌هایی برای بردارها تعریف می‌شود (مانند معیارهای شباهت، میزان مشابهت محتوایی و زمان ایجاد پیوند).

۲ به ازای تمام جفت رأس‌های ممکن در گراف آموزشی بردار ایجاد و مقادیر ویژگی‌های هر جفت رأس محاسبه می‌شود.

۳ با توجه به این که جفت رأس تشکیل دهنده بردار با یکدیگر پیوند دارند یا خیر، برچسب «وجود پیوند» یا «عدم وجود پیوند» انتصاب داده می‌شود.

۴ الگوریتم یادگیری ماشین با استفاده از بردارهای به دست آمده آموزش داده می‌شود.

۵ مدل نهایی برای هر جفت رأس جدید در گراف، پیش‌بینی پیوند می‌کند.

به عنوان مثال در جدول ۳ برای هر جفت رأس در شبکه یک بردار ویژگی (یک سطر جدول) بدست آمده است. الگوریتم پس از یادگیری این داده‌ها می‌تواند برچسب یک داده جدید را پیش‌بینی کند.

جدول ۳ مثالی از بردارهای ورودی الگوریتم پیش‌بینی پیوند (مجموعه

آموزشی)  $V$ : وجود پیوند و  $X$ : عدم وجود پیوند

جفت رأس	JC	N2V_COS_S IM	TITLE_COS_ SIM	...	برچسب
A - B	۰,۶۷	۰,۴۰	۰,۵۱	...	$V$
A - C	۰,۲۳	۰,۵۹	۰,۳۱	...	$X$
B - C	۰,۶۳	۰,۲۲	۰,۳۹	...	$X$
A - D	۰,۹۵	۰,۸۵	۰,۷۱	...	$V$
...	...	...	...	...	...

۴-۴ حل مشکل نمونه‌گیری منفی با خوشه‌بندی

<sup>1</sup> Sparse

ارزیابی مورد بررسی قرار می‌گیرند و در ادامه برای هر آزمایش، شرح آزمایش، نتایج و تحلیل ارائه می‌شود.

## ۵-۱ ابزارهای مورد استفاده

سخت‌افزار مورد استفاده برای آزمایشات این فصل یک رایانه‌ی خانگی مشخصات Intel Pentium Silver N5000 2.7GHZ 8GB Ram بوده است. البته برخی آزمایشات مربوط به روش‌های جاسازی با الگوریتم Node2Vec به دلیل زمان زیاد اجرا روی این سیستم قابل اجرا نبودند. این آزمایشات روی سرور محاسبات ابری دانشگاه فردوسی مشهد<sup>۱</sup> با پردازنده ۳۲ هسته‌ای و ۶۴ گیگابایت رم اجرا شد.

آزمایشات با زبان پایتون و در محیط Jupyter-Lab پیاده‌سازی شدند. برای پردازش‌های مربوط به گراف از NetworkX [۲۹] و برای پیاده‌سازی الگوریتم‌های مربوط به یادگیری ماشین و انتخاب و استخراج ویژگی از scikit-learn استفاده شد. برای بارگزاری داده‌ها، تقسیم گراف به مجموعه‌ی آموزشی، اعتبارسنجی و آزمایشی و در نهایت ارزیابی و مقایسه نتایج از چارچوب OGB<sup>۲</sup> استفاده کردیم. توضیحات این چارچوب در پیوست ب ارایه شده است.

## ۵-۲ مجموعه داده‌ها

در این بخش از دو مجموعه داده OGB استفاده می‌شود که در زمینه پیش‌بینی پیوند و شبکه‌های علمی هستند. مجموعه داده ogbl-collab یک شبکه از نویسندگان و همکاری‌های بین آن‌هاست. عنوان مقالات هر نویسنده در رأس‌ها و سال همکاری

الگوریتم برای تشخیص پیوندهای منفی قرار نمی‌دهند و حتی ممکن است باعث بیش برآزش مدل یادگیری ماشین نیز بشود.

برای کم کردن تعداد نمونه‌های منفی می‌توان از میان نمونه‌های منفی، تعدادی را به صورت تصادفی انتخاب کرد. با انتخاب تصادفی ممکن است برخی از نمونه‌هایی که در آموزش بهتر مدل می‌توانند تأثیر گزار باشند انتخاب نشوند و از طرفی داده‌هایی که شبیه به یکدیگر هستند چندین بار انتخاب شوند.

برای حل این مشکل دو روش مبتنی بر خوشه‌بندی و رگرسیون پیشنهاد شده است. در روش مبتنی بر خوشه‌بندی ابتدا نمونه‌های منفی به K خوشه تقسیم و از هر خوشه یک نماینده انتخاب می‌شود. به این ترتیب از تمام فضای نمونه‌های منفی به طور یکنواخت نمونه‌گیری شده و نمونه‌های منفی با یکدیگر متفاوت خواهند بود. با این روش تعداد نمونه‌های منفی کاهش و کیفیت آن‌ها افزایش پیدا خواهد کرد. به این ترتیب می‌توان با منابع سخت‌افزاری در دسترس و الگوریتم‌های یادگیری ماشین رایج مسئله دسته‌بندی را حل کرد. پیش‌بینی می‌شود استفاده از همین روش برای داده‌های مثبت نیز تأثیر مثبتی داشته باشد. در هر گراف مورد پردازش ممکن است مقدار بهینه برای K متفاوت باشد. پیشنهاد ما برای یافتن K بهینه استفاده از آزمون و خطا و یا معیارهای ارزیابی خوشه‌بندی است.

## ۵-۳ آزمایش‌ها

در این بخش پنج آزمایش برای مقایسه روشهای پیشنهادی و سایر روشهای موجود بر اساس معیارهای ارزیابی عملکرد HITS و MRR بر روی دو مجموعه داده ارجاعات و همکاری نویسندگان طراحی شده است. ابتدا ابزارهای نرم‌افزاری و سخت‌افزاری شرح داده می‌شوند. سپس مجموعه داده‌های مورد استفاده و معیارهای

<sup>1</sup> <https://ferdowsi.cloud/>

<sup>2</sup> Open Graph Benchmark

عدد کوچکتری باشد معیار سختگیرانه تر می شود. این معیار در معادله ۱۵ نشان داده شده است. در این معادله  $t$  داده مورد آزمایش،  $S_{test}$  مجموعه داده‌های آزمایشی و  $rank(t)$  رتبه پیوند مثبت در مجموعه آزمایشی است. هر چقدر الگوریتم پیش‌بینی پیوند بتواند امتیازی بالاتری به پیوند مثبت نسبت به پیوندهای منفی انتصاب دهد، پیوند مثبت در رتبه‌ی بالاتری قرار می‌گیرد و مقدار معیار Hits افزایش می‌یابد (برای اطلاعات بیشتر به پیوست ب مراجعه فرمایید).

$$Hits@k = \frac{|S_{test} \cap \{rank(t) \leq k\}|}{|S_{test}|} \quad 15$$

### ۲-۳-۵ معیار MRR

میانگین رتبه متقابل یا MRR یک امتیاز نسبی است که میانگین یا میانگین معکوس رتبه‌هایی را که در آن اولین سند مرتبط برای مجموعه‌ای از جستارها بازیابی شده است. این معیار در اصل مربوط به بازیابی اطلاعات و جستجو است ولی در پیش‌بینی پیوند به معنای میانگین رتبه یک پیوند مثبت در میان مجموعه‌ای پیوندهای منفی است. این معیار با معادله ۱۶ نشان داده می‌شود. در این معادله  $Q$  مجموعه پیوندهای منفی به علاوه پیوند مثبت مد نظر و  $rank(t)$  به معنای رتبه پیوند مثبت در میان پیوندهای منفی پس از مرتب‌سازی نزولی است. هر چقدر مدل پیش‌بینی کننده به پیوند مثبت امتیاز بیشتری نسبت به پیوندهای منفی بدهد و رتبه بالاتری کسب کند مقدار MRR بیشتر می‌شود.

$$MRR = \frac{1}{|Q|} \sum_{t=1}^{|Q|} \frac{1}{rank(t)} \quad 16$$

### ۴-۵ آزمایش اول: معیارهای شباهت محلی

در یال‌ها به عنوان فرا داده در مجموعه داده موجود است. مجموعه داده ogbl-citation2 نیز یک شبکه ارجاعات است که متشکل از مقالات و ارجاعات بین آن‌ها است. در رأس‌های این شبکه سال انتشار مقاله و عنوان چکیده مقاله به عنوان ویژگی وجود دارد. هر دو این مجموعه داده‌ها بر اساس زمان به دو مجموعه آموزشی و آزمایشی تقسیم می‌شود.

جدول ۴ مشخصات مجموعه داده‌های مربوط به شبکه‌های علمی در OGB

ogbl-citation2	ogbl-collab	
بزرگ	متوسط	مقیاس
۲,۹۲۷,۹۶۳	۲۳۵,۸۶۸	# رأس‌ها
۳۰,۵۶۱,۱۸۷	۱,۲۸۵,۴۶۵	# یال‌ها
زمان	زمان	تقسیم بر اساس
عنوان و چکیده مقاله، سال انتشار	عنوان مقالات نویسنده، سال همکاری	داده محتوایی
MRR	Hits@50	معیار ارزیابی

### ۳-۵ معیارهای ارزیابی

برای پیش‌بینی پیوند معیارهای ارزیابی زیادی وجود دارد و همچنین می‌توان از معیارهای ارزیابی دسته‌بندی برای پیش‌بینی پیوند استفاده کرد. ولی با توجه به این که در آزمایشات از چارچوب OGB استفاده می‌شود، در انتخاب معیار ارزیابی قدرت انتخاب وجود ندارد. OGB برای مجموعه داده ogbl-collab از معیار Hits@50 و برای ogbl-citation2 از معیار MRR<sup>1</sup> استفاده می‌کند. در ادامه هر دو معیار شرح داده می‌شوند.

### ۱-۳-۵ معیار Hits@k

معیار Hits@k نشان می‌دهد چه نسبتی از پیوندهای آزمایشی مثبت در میان مجموعه‌ای از پیوندهای منفی پس از مرتب‌سازی نزولی میان  $k$  پیوند اول قرار می‌گیرد. در این معیار هر چقدر  $k$

<sup>1</sup> Mean Reciprocal Rank

عمل کرده است. شباهت کسینوسی بردار نهفته‌سازی شده دو رأس ( $W2V\_SIM$ ) بدون در نظر گرفتن ساختار گراف، در هر دو مجموعه داده از پیش‌بینی کننده تصادفی بهتر عمل کرده است. این موضوع نشان می‌دهد محتوای متنی رأس‌های شبکه علمی می‌توانند در پیش‌بینی پیوند تأثیر گذار باشند. شباهت محتوا در شبکه‌ی همکاری نویسندگان بهتر از شبکه ارجاعات عمل کرده است، از این موضوع نتیجه می‌گیریم همکاری که بین دو نویسنده شکل می‌گیرد در مقایسه با ارجاع مقالات، وابستگی بیشتر به موضوع و محتوا دارد.

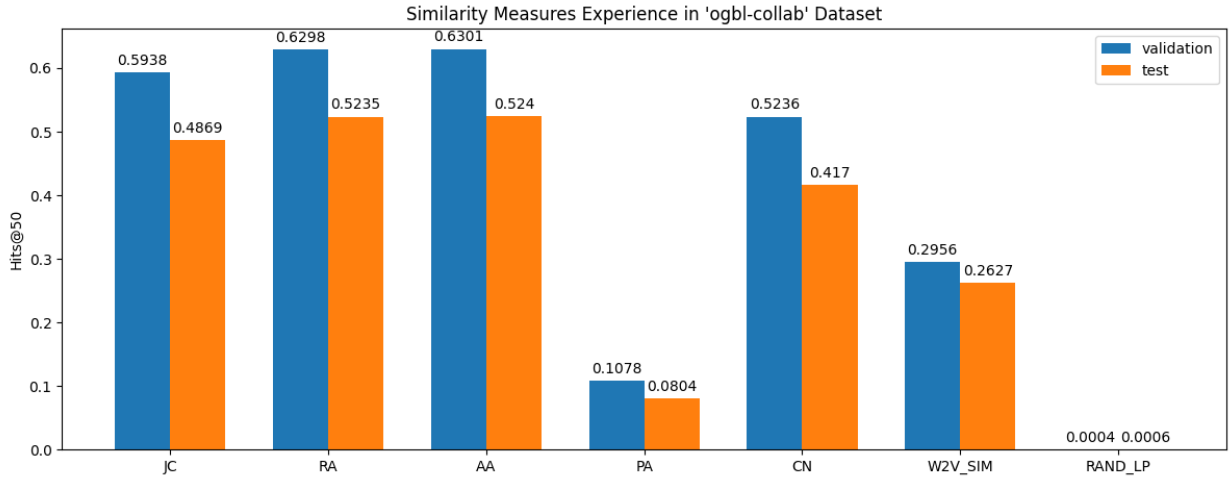
نتایج مدل تصادفی را با استفاده از احتمال می‌توان توجیه کرد. در  $ogbl-collab$  تعداد پیوندهای منفی ساختگی ۱۰۰,۰۰۰ و پیوند مثبت واقعی در  $Hits@50$  باید در بین ۵۰ داده اول رتبه‌بندی شود. بر اساس معادله ۱۷، نتیجه به دست آمده از طریق احتمال با نتیجه حاصل از آزمایش تفاوت زیادی ندارد. همچنین در  $ogbl-citation2$  میانگین معکوس رتبه پیوند هدف در میان ۱۰۰۰ رأس منفی به عنوان  $MRR$  محاسبه می‌شود. نتیجه پیش‌بینی تصادفی در معادله ۱۸ با نتایج آزمایش مطابقت دارد. با افزایش تعداد آزمایش‌ها و میانگین گرفتن از نتایج، عدد میانگین به عدد به دست آمده از طریق احتمال همگرا می‌شود.

$$Hits@50_{rand} = \frac{50}{100,000} = 0.0005 \quad 17$$

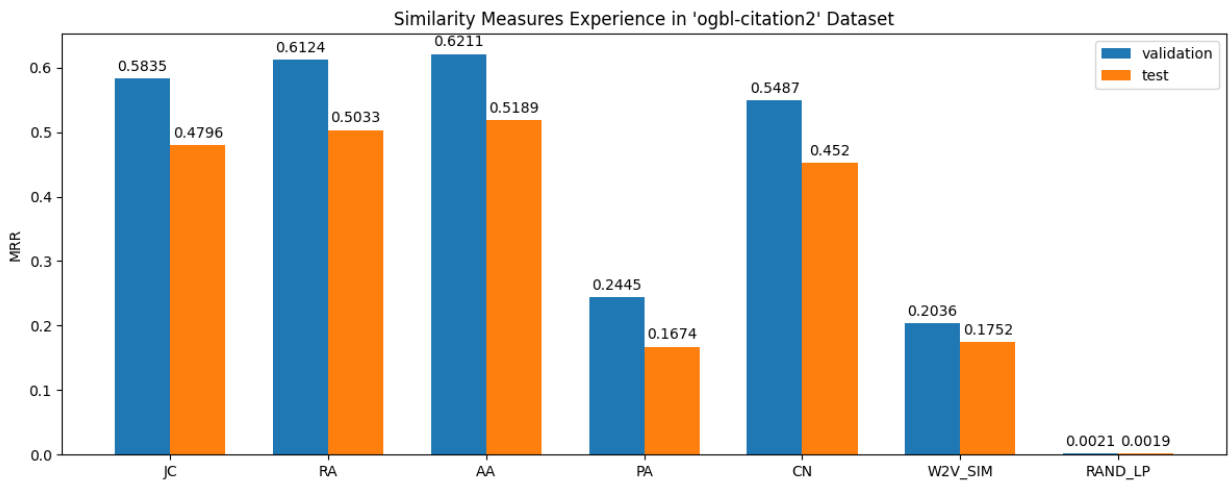
$$MRR_{rand} = \frac{1}{1000} \sum_{i=1}^{1000} \frac{1}{500} = 0.002 \quad 18$$

در این بخش، یک آزمایش برای مقایسه روشهای مبتنی بر شباهت محلی گراف طراحی شده است. در این آزمایش و آزمایش‌های بعدی داده مورد آزمایش  $ogbl-collab$  و  $ogbl-citation2$  خواهد بود که به ترتیب بر اساس معیارهای  $MRR$  و  $Hits@50$  ارزیابی می‌شوند. همچنین این آزمایش با یک پیش‌بینی کننده تصادفی با نام  $random\_lp$  تکرار شده که یک عدد تصادفی بین ۰ و ۱ به هر جفت رأس انتصاب می‌دهد. معیار  $random\_lp$  را برای این به آزمایش اضافه کردیم که ببینیم معیارهای شباهت محلی نسبت به یک روش تصادفی چقدر بهتر عمل می‌کنند. همچنین از شباهت محتوایی  $W2V\_SIM$  نیز استفاده کردیم که فارغ از ساختار شبکه، صرفاً شباهت کسینوسی بردار نهفته‌سازی دو رأس را محاسبه می‌کند.

همانطور در شکل ۹ و شکل ۱۰ مشخص است معیار آدامیک آدار ( $AA$ ) در هر دو مجموعه داده و پس از آن تخصیص منبع ( $RA$ ) بهترین نتایج را کسب کرده است. به نظر می‌رسد این دو معیار به دلیل جریمه کردن درجه همسایگان مشترک توانستند نسبت به بقیه نتیجه بهتری کسب کنند. جریمه کردن درجه همسایگان مشترک در کاربردهای دیگر مانند شبکه‌های اجتماعی تأثیر مثبتی داشته و می‌بینیم در شبکه‌های علمی نیز این قاعده برقرار است. یعنی مقالاتی که ارجاعات زیادی دارند و نویسندگانی که همکاری‌های زیادی انجام دادند نسبت به مقالات کم ارجاع و نویسندگان با تعداد کم همکاری، همسایه‌های کم اهمیت تری هستند. همچنین معیار اتصال ترجیحی ( $PA$ ) بدترین نتایج را کسب کرده است ولی با این حال با اختلاف زیادی از پیش‌بینی کننده‌ی تصادفی بهتر

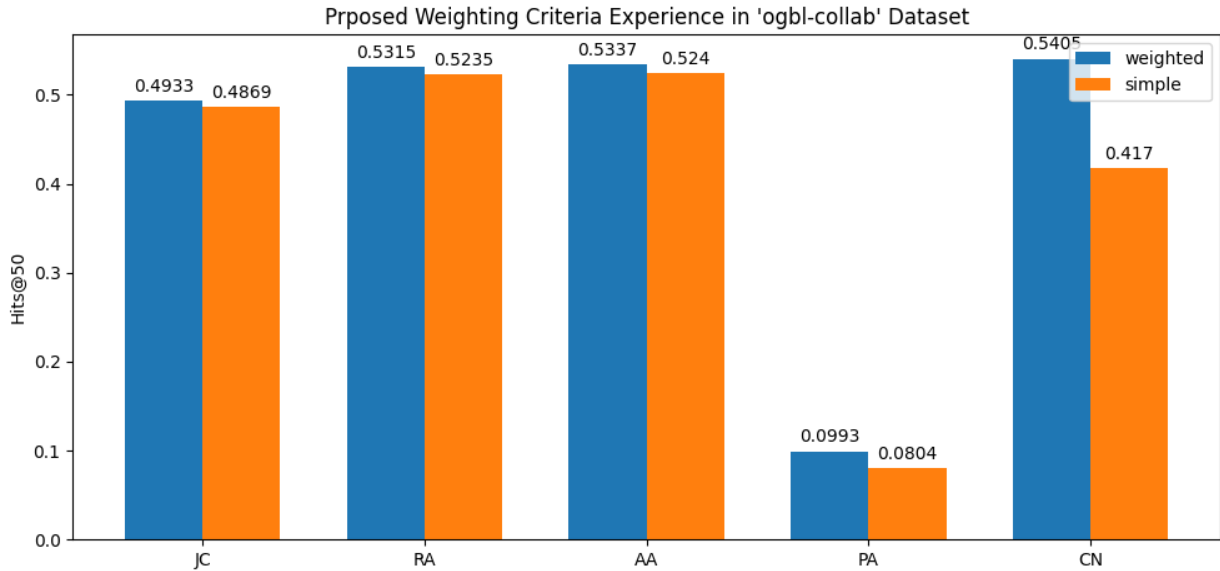


شکل ۹ مقادیر Hits@50 بر اساس استفاده از معیارهای شباهت محلی در مجموعه داده ogbl-collab

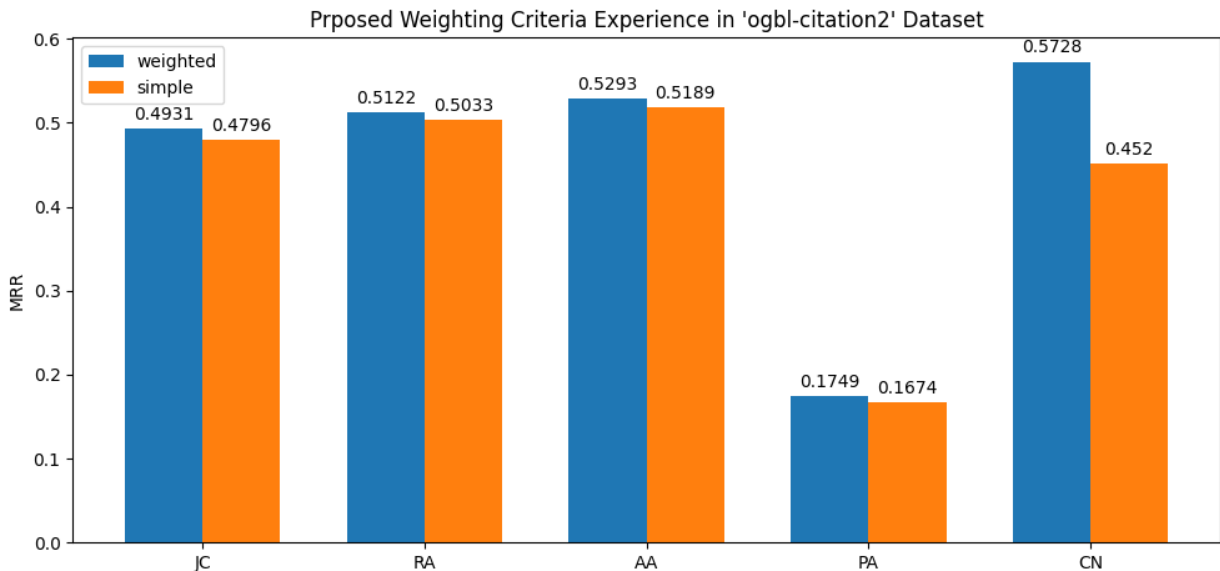


شکل ۱۰ مقادیر MRR بر اساس استفاده از معیارهای شباهت محلی در مجموعه داده ogbl-citation2





شکل ۱۱ مقادیر Hits@50 بر اساس استفاده از روش پیشنهادی وزن دهی و مقایسه با روش بدون وزن در داده ogbl-collab



شکل ۱۲ مقادیر MRR بر اساس استفاده از روش پیشنهادی وزن دهی و مقایسه با روش بدون وزن در داده ogbl-citation2

شبکه‌های همکاری نویسندگان و ارجاعات بر اساس روش پیشنهادی وزن دار شدند و با استفاده از معیارهای شباهت وزن دار که در مفاهیم پایه شرح داده شدند مورد پردازش قرار گرفتند. سپس نتایج هر معیار با نسخه بدون وزن (آزمایش اول) مقایسه

### ۵-۵ آزمایش دوم: روش وزن دهی پیشنهادی

این آزمایش برای بررسی روش‌های وزن دهی پیشنهاد شده در معادله ۱۲ و معادله ۱۳ طراحی شده است. در این آزمایش

همچنین جدول ۵ نشان می‌دهد در شبکه همکاری نویسندگان تأثیر زمان و محتوا به یک اندازه است. ولی در شبکه ارجاعات تأثیر زمان بیشتر از محتوا است؛ به عبارت دیگر تمایل ایجاد پیوند در شبکه ارجاعات بیشتر با رأس‌های جدیدتر است تا رأس‌های با شباهت محتوایی بیشتر.

## ۵-۶ آزمایش سوم: پیش‌بینی پیوند با SVM و نمونه‌گیری

### منفی تصادفی

این آزمایش برای بررسی روش‌های پیشین مبتنی بر تک ویژگی (شباهت محلی وزن‌دار و بدون وزن) با روش یادگیری ماشین طراحی شده است. در این آزمایش، به تعداد یال‌های آموزشی با برچسب مثبت، یال‌هایی با برچسب منفی به طور تصادفی انتخاب کردیم، و ویژگی‌های ذکر شده در جدول ۲ را برای آن‌ها محاسبه کردیم. به این ترتیب هر پیوند مثبت یا منفی را به یک بردار ۱۳ بعدی تبدیل کردیم. با توجه به این که مقادیر ویژگی‌های محاسبه شده در بازه‌های متفاوتی قرار دارند و بدون نرمال سازی باعث می‌شود برخی الگوریتم‌های یادگیری ماشین به سمت ویژگی‌هایی با مقدار بزرگتر سوگیری کنند. به همین خاطر مقادیر را با تابع StandardScaler از بسته scikit-learn نرمال کردیم.

در این آزمایش از SVM در بسته scikit-learn با هسته linear به عنوان الگوریتم یادگیری ماشین استفاده شده است. با توجه به این که معیارهای MRR و Hits@k هر دو بر اساس مرتب سازی بر حسب امتیاز عمل می‌کنند و بر خلاف مسئله دسته‌بندی، داده‌های آزمایشی بر حسب برچسب مثبت و منفی مطلق مورد ارزیابی قرار نمی‌گیرد و احتمال وجود پیوند در ارزیابی استفاده می‌شود، پارامتر probability را برابر true قرار دادیم که به

می‌شود. در این آزمایش و آزمایش‌های بعدی فقط داده‌های آزمایشی در نتایج آورده می‌شوند و داده‌های اعتبار سنجی را برای جلوگیری از نمودارهای اضافی و آسان‌تر شدن مقایسه نتایج گزارش نمی‌کنیم. مقدار بهینه متغیرهای  $\alpha$  و  $\beta$  برای این آزمایش که با آزمون و خطا به دست آمدند در جدول ۵ آمده است. نتایج این آزمایش در مقایسه با آزمایش قبل نیز در شکل ۱۱ و شکل ۱۲ قابل مشاهده است.

جدول ۵ مقدار  $\alpha$  و  $\beta$  در شبکه همکاری نویسندگان و ارجاعات

مقدار	متغیر	شبکه
0.5	$\alpha$ (تأثیر محتوا)	همکاری نویسندگان (معادله ۱۲)
0.5	$\beta$ (تأثیر زمان)	
0.4	$\alpha$ (تأثیر محتوا)	ارجاعات (معادله ۱۳)
0.6	$\beta$ (تأثیر زمان)	

نتایج نشان می‌دهد تمام معیارهای برای هر مجموعه داده تقریباً به میزان ثابتی بهبود پیدا کردند. البته معیار همسایه مشترک (CN) در هر دو مجموعه داده حدود ۱۴٪ افزایش داشته. همسایه مشترک پایه معیارهای دیگر مانند ژاکارد و آدامیک آدار و تخصیص منبع است منتها نسبت به آن‌ها رفتار خیلی ساده تری دارد. احتمالاً به همین دلیل رفتار ساده (تعداد همسایه‌های مشترک) در مقایسه با آدامیک آدار و تخصیص منبع که همسایگان مشترک را بر اساس درجه جریمه می‌کنند یا ژاکارد که نسبت همسایه‌های مشترک به تمام همسایه‌ها را به عنوان شباهت محاسبه می‌کند عملکرد ضعیف تری دارد. با معیار وزن دهی پیشنهاد شده، ارزش هر همسایه مشترک تفاوت خواهد کرد. همسایه مشترک وزن‌دار مجموع ارزش همسایه‌ها را عنوان شباهت محاسبه می‌کند و به همین خاطر در مقایسه با همسایه مشترک ساده بسیار هوشمندانه تر عمل می‌کند و توانسته نتایجی مشابه سایر معیارها (به غیر از اتصال ترجیحی) کسب کند.

نتایج آزمایش را تکرار کردیم و نتایج هر تکرار را در جدول ۷ گزارش کردیم.

جدول ۷ نتایج پیش‌بینی پیوند با PCA و SVM

SVM	تکرار ۱	تکرار ۲	تکرار ۳	تکرار ۴	تکرار ۵	میانگین
Hits@50	0.540	0.536	0.531	0.545	0.537	$0.538 \pm 0.004$
MRR	0.525	0.528	0.529	0.526	0.526	$0.527 \pm 0.001$

همان طور که انتظار می‌رفت، مشاهده می‌شود نتایج نسبت به آزمایش پیش که در شرایط یکسان و بدون PCA انجام شد بهبود پیدا کرده و مانند بسیاری از مسائل یادگیری ماشین PCA باعث افزایش دقت مدل یادگیری ماشین شده. با توجه به این تأثیر مثبت، در آزمایشات بعدی که از خوشه‌بندی و رگرسیون استفاده خواهد شد، از ویژگی‌های استخراج شده توسط PCA استفاده می‌شود.

#### ۵-۸ آزمایش پنجم: نمونه‌گیری با خوشه‌بندی

این آزمایش برای بررسی نتایج روش پیشنهادی برای نمونه‌گیری منفی با استفاده از خوشه‌بندی طراحی شده است. به دلیل تعداد زیاد داده الگوریتم‌های خوشه‌بندی معمولی زمان زیادی برای اجرا نیاز دارند به همین خاطر در این آزمایش از الگوریتم MiniBatchKMeans با  $batch\_size=1024$  استفاده شده است. داده‌های مثبت و منفی پس از استخراج ویژگی با PCA خوشه‌بندی شدند و الگوریتم یادگیری ماشین بر اساس مرکز هر خوشه آموزش داده شد. برای این کار آزمایش را برای  $k=100$  تا  $k=1000$  تکرار کردیم و نتایج را در شکل ۱۳ نشان دادیم. همچنین برای درک بهتر روش پیشنهادی یک بار ۵۰۰ داده مثبت و منفی را از مجموعه داده ogbl-collab یک بار به صورت تصادفی و بار دیگر با استفاده از خوشه‌بندی انتخاب کردیم و داده‌ها را به صورت دو بعدی در شکل ۱۴ نشان دادیم.

جای یک برچسب، احتمال تعلق به آن برچسب (وجود یا عدم وجود پیوند) را به عنوان خروجی محاسبه کند.

در این آزمایش با توجه به انتخاب تصادفی داده‌های منفی، ممکن است نتیجه با هر بار اجرا تغییر کند، به همین دلیل این آزمایش را پنج بار با داده‌های تصادفی متفاوت تکرار کردیم و نتایج را در جدول ۶ نشان دادیم.

جدول ۶ نتایج پیش‌بینی پیوند با SVM و نمونه‌گیری منفی تصادفی

	تکرار ۱	تکرار ۲	تکرار ۳	تکرار ۴	تکرار ۵	میانگین
Hits@50	0.513	0.503	0.528	0.508	0.503	$0.511 \pm 0.009$
MRR	0.490	0.506	0.493	0.509	0.503	$0.500 \pm 0.007$

با توجه به اختلاف نتایج در هر بار تکرار، به نظر می‌رسد داده‌های آموزشی در کیفیت پیش‌بینی پیوند تأثیر داشته باشد و به دلیل تصادفی انتخاب شدن در هر تکرار، نتایج با تکرارهای قبل تفاوت دارد. به طور کلی نتایج این آزمایش نسبت به بهترین شباهت‌های وزن دار آزمایش قبل بهبود خاصی نداشته. دلیل این موضوع می‌تواند ناکافی بودن داده‌های آموزشی و انتخاب تصادفی این داده‌ها باشد. همچنین ممکن است وجود ویژگی‌هایی که با سایر ویژگی‌ها همبستگی دارند نیز یکی از علت‌های نتیجه نامطلوب باشد. این موارد در آزمایش‌های بعدی بررسی خواهند شد.

#### ۵-۷ آزمایش چهارم: استخراج ویژگی با PCA

از آنجایی که برخی ویژگی‌های جدول ۲ ممکن است با یکدیگر همبستگی داشته باشند، این آزمایش را طراحی کردیم که با استفاده از الگوریتم استخراج ویژگی PCA ابعاد داده‌ها را کاهش دهیم در این آزمایش بهترین نتایج با  $n\_component=5$  حاصل شد، این مقدار با آزمون و خطا به دست آمد. مانند آموزش قبل با توجه به تصادفی انتخاب شدن داده‌های آموزشی، به دلیل امکان اختلاف

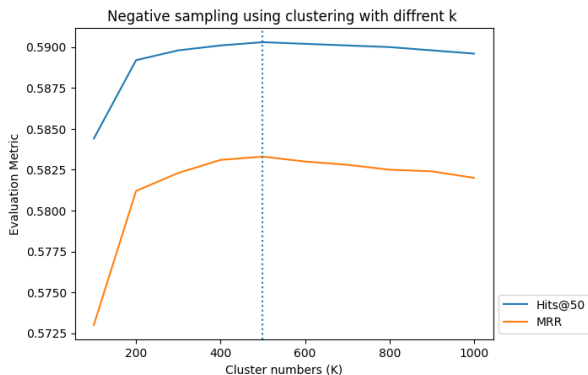
اصلی وجود دارد آموزش داده نخواهد شد. در صورتی که نمونه گیری با خوشه بندی حالات بیشتری از هر برچسب مثبت یا منفی را پوشش می دهد.

### ۶- نتیجه گیری

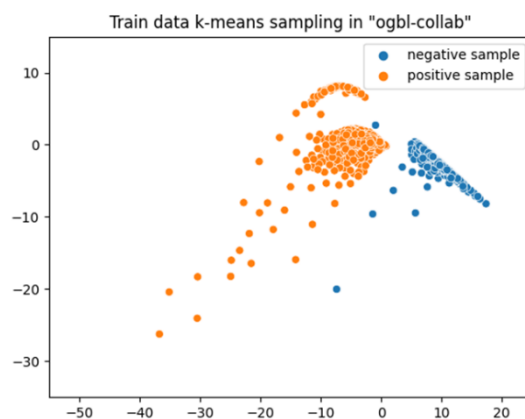
شبکه های علمی مانند شبکه همکاری نویسندگان و شبکه رجاعات شبکه هایی هستند که نحوه تکامل آن ها در آینده قابل پیش بینی است. با پیش بینی این تکامل، می توان برای یافتن متخصص و منابع مناسب برای پژوهش های آتی کمک گرفت. در این مقاله یک چارچوب کلی برای مسائل پیش بینی پیوند در شبکه های علمی پیشنهاد شد که اساس آن ترکیب ویژگی های محتوایی و ساختاری در قالب یک بردار ویژگی و استفاده از الگوریتم های یادگیری ماشین است؛ برای کاهش ابعاد از PCA استفاده شد. با توجه به حجم زیاد داده ها در شبکه های علمی دنیای واقعی و پراکنده بودن پیوندهای شبکه، نمونه های منفی بسیار بیشتر از نمونه های مثبت هستند که با روش پیشنهادی در این مقاله، استفاده از خوشه بندی داده های منفی و انتخاب مراکز خوشه ها، نمونه های محدودی به عنوان نماینده انتخاب شد که با این روش نیز نتایج بار دیگر بهبود یافت.

پایه سازی روش های یادگیری ماشین روی شبکه های علمی با توجه به برخی ویژگی های آن مفید باشد. به عنوان مثال توجه به حوزه های تحقیقاتی که می تواند در گراف های علمی به صورت زیرگراف (یا جامعه) در نظر گرفته شود. همچنین آزمایش بیشتر روی یافتن ترکیب ماسبتری از محتوا و ساختار و شاید مشارکت سایر ویژگی های شبکه های علمی می تواند یکی از کارهای آینده در این حوزه باشد.

### مراجع



شکل ۱۳ ارزیابی نتایج پیش بینی پیوند با نمونه گیری منفی با خوشه بندی به ازای تعداد خوشه های متفاوت



شکل ۱۴ مقایسه نمونه گیری تصادفی و استفاده از خوشه بندی در داده های آموزشی

در شکل ۱۴ سمت چپ یک بار از داده های آموزشی به صورت تصادفی ۵۰۰ داده مثبت و منفی انتخاب کردیم، و یکبار همین داده ها را به ۵۰۰ خوشه تقسیم کردیم و مراکز خوشه ها را در سمت راست نشان دادیم. با وجود این که در هر دو سمت شکل ۵۰۰ داده مثبت و منفی انتخاب شده ولی داده های انتخابی در نمونه گیری تصادفی بسیار نزدیک به هم تر و متراکم تر نسبت به روش خوشه بندی هستند که نشان می دهد بسیاری از داده ها تقریباً تکراری هستند. هنگامی که الگوریتم با نمونه گیری تصادفی آموزش داده شود بسیاری از حالات مثبت یا منفی که در داده های

- [10]. N. Shibata, Y. Kajikawa, and I. Sakata, "Link prediction in citation networks," *Journal of the American society for information science and technology*, vol. 63, no. 1, pp. 78-85, 2012, <https://doi.org/10.1002/asi.21664>.
- [11]. V. Latora, V. Nicosia, and G. Russo, *Complex networks: principles, methods and applications*: Cambridge University Press, 2017.
- [12]. P. M. Chuan, M. Ali, T. D. Khang, and N. Dey, "Link prediction in co-authorship networks based on hybrid content similarity metric," *Applied Intelligence*, vol. 48, no. 8, pp. 2470-2486, 2018.
- [13]. E. Bütün, M. Kaya, and R. Alhajj, "Extension of neighbor-based link prediction methods for directed, weighted and temporal social networks," *Information Sciences*, vol. 463, pp. 152-165, 2018, <https://doi.org/10.1016/j.ins.2018.06.051>.
- [14]. S. Behrouzi, Z. S. Sarmoor, K. Hajsadeghi, and K. Kavousi, "Predicting scientific research trends based on link prediction in keyword networks," *Journal of Informetrics*, vol. 14, no. 4, pp. 101079, 2020, <https://doi.org/10.1016/j.joi.2020.101079>.
- [15]. A. Daud, W. Ahmed, T. Amjad, J. A. Nasir, N. R. Aljohani, R. A. Abbasi, and I. Ahmad, "Who will cite you back? Reciprocal link prediction in citation networks," *Library Hi Tech*, vol. 35, no. 4, pp. 509-520, 2017, <https://doi.org/10.1108/LHT-02-2017-0044>.
- [16]. D. Liben-Nowell, and J. Kleinberg, "The link-prediction problem for social networks. journal of the Association for Information Science and Technology (2007)," *Google Scholar Google Scholar Digital Library Digital Library*, 2007,
- [17]. S. Martinčić-Ipšić, E. Močibob, and M. Perc, "Link prediction on Twitter," *PloS one*, vol. 12, no. 7, pp. e0181079, 2017, <https://doi.org/10.1371/journal.pone.0181079>.
- [18]. L. A. Adamic, and E. Adar, "Friends and neighbors on the web," *Social networks*, vol. 25, no. 3, pp. 211-230, 2003, [https://doi.org/10.1016/S0378-8733\(03\)00009-1](https://doi.org/10.1016/S0378-8733(03)00009-1).
- [1]. R. Taimourei-Yansary, M. Mirzarezaee, M. Sadeghi, and B. N. Araabi, "Predicting invasive disease-free survival time in breast cancer patients using semi-supervised graph-based machine learning techniques", *Soft Computing Journal*, vol. 10, no. 1, pp. 48-69, 2022.
- [2]. E. Mahfooz and G. Fath-Tabar, "Sum of distance between vertices of graphs", *Soft Computing Journal*, vol. 5, no. 2, pp. 28-33, 2021.
- [3]. A. Keypour, "Link Prediction in Social Networks through classifiers combination", *Soft Computing Journal*, vol. 4, no. 2, pp. 2-17, 2021.
- [4]. V. Martínez, F. Berzal, and J.-C. Cubero, "A survey of link prediction in complex networks", *ACM computing surveys (CSUR)*, vol. 49, no. 4, pp. 1-33, 2016, DOI: <https://doi.org/10.1145/3012704>.
- [5]. L. Lü, and T. Zhou, "Link prediction in complex networks: A survey," *Physica A: statistical mechanics and its applications*, vol. 390, no. 6, pp. 1150-1170, 2011, <https://doi.org/10.1016/j.physa.2010.11.027>.
- [6]. C. P. Muniz, R. Goldschmidt, and R. Choren, "Combining contextual, temporal and topological information for unsupervised link prediction in social networks," *Knowledge-Based Systems*, vol. 156, pp. 129-137, 2018, <https://doi.org/10.1016/j.knosys.2018.05.027>.
- [7]. M. Nikkar, R. Alijani, and K. M. H. GHAZIZADEH, "Investigation of the presence of surgery researchers in research gate scientific network: An altmetrics study," *Iranian Journal of Surgery*, vol. 25, no. 2, pp. 76-82, 2017.
- [8]. H. Liu, H. Kou, C. Yan, and L. Qi, "Link prediction in paper citation network to construct paper correlation graph," *EURASIP Journal on Wireless Communications and Networking*, vol. 2019, no. 1, pp. 1-12, 2019.
- [9]. E. Bütün, and M. Kaya, "Predicting citation count of scientists as a link prediction problem," *IEEE transactions on cybernetics*, vol. 50, no. 10, pp. 4518 - 4529, 2019, <https://doi.org/10.1109/TCYB.2019.2900495>.

*function using NetworkX*, Los Alamos National Lab.(LANL), Los Alamos, NM (United States), 2008.

### پیوست الف: الگوریتم **node2vec**

الگوریتم **node2vec** بر پایه الگوریتم **word2vec** که در پردازش متن برای تبدیل کلمات به فضای برداری استفاده می شود ارائه شده است. در الگوریتم **node2vec** رأس ها معادل کلمات در الگوریتم **word2vec** هستند. الگوریتم **word2vec** دنباله ای از کلمات (جملات) را به عنوان ورودی می پذیرد. با انجام تعداد زیادی قدم زدن تصادفی و ثبت دنباله رأس های طی شده، این دنباله برای رأس ها ساخته می شود (شکل الف-۱، قسمت ب).

در این الگوریتم نسخه خاصی از قدم زدن تصادفی به نام قدم زدن تصادفی درجه دوم<sup>۱</sup> استفاده می شود. در قدم زدن تصادفی درجه اول، تاریخچه قدم زدن نگهداری نمی شود. یعنی هنگامی که از یک رأس به رأس مبدأ انتقال صورت می گیرد، در مبدأ الگوریتم نمی داند از چه رأسی به حالت فعلی رسیده ولی در قدم زدن تصادفی درجه دوم، رأسی که از آن به مبدأ انتقال صورت گرفته اهمیت پیدا می کند.

در قدم زدن تصادفی درجه دوم، احتمال انتقال از یک رأس به رأس های مجاور از طریق تابع  $\alpha$  که در معادله الف-۱ نشان داده شده است صورت می گیرد. در این رابطه  $t$  رأسی است که انتقال از آن صورت گرفته،  $v$  رأسی است که در زمان فعلی الگوریتم در آن قرار دارد و  $x$  رأس هایی که هستند که امکان انتقال به آن ها از رأس  $v$  وجود دارد. احتمال بازگشت به رأس قبلی  $p$  و احتمال

- [19]. T. Zhou, L. Lü, and Y.-C. Zhang, "Predicting missing links via local information," *The European Physical Journal B*, vol. 71, no. 4, pp. 623-630, 2009, <https://doi.org/10.1140/epjb/e2009-00335-8>.
- [20]. M. E. Newman, "Clustering and preferential attachment in growing networks," *Physical review E*, vol. 64, no. 2, pp. 025102, 2001, <https://doi.org/10.1103/PhysRevE.64.025102>
- [21]. N. Benchettara, R. Kanawati, and C. Rouveirol, "A supervised machine learning link prediction approach for academic collaboration recommendation." pp. 253-256, <https://doi.org/10.1145/1864708.1864760>.
- [22]. F. Almeida, and G. Xexéo, "Word embeddings: A survey," *arXiv preprint arXiv:1901.09069*, 2019.
- [23]. J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation." pp. 1532-1543.
- [24]. J. Zhou, L. Liu, W. Wei, and J. Fan, "Network representation learning: from preprocessing, feature extraction to node embedding," *ACM Computing Surveys (CSUR)*, vol. 55, no. 2, pp. 1-35, 2022.
- [25]. M. Grohe, "word2vec, node2vec, graph2vec, x2vec: Towards a theory of vector embeddings of structured data." pp. 1-16.
- [26]. A. Grover, and J. Leskovec, "node2vec: Scalable feature learning for networks." pp. 855-864, <https://doi.org/10.1145/2939672.2939754>.
- [27]. D. Lande, M. Fu, W. Guo, I. Balagura, I. Gorbov, and H. Yang, "Link prediction of scientific collaboration networks based on information retrieval," *World Wide Web*, pp. 1-19, 2020, <https://doi.org/10.1007/s11280-019-00768-9>.
- [28]. B. Liu, S. Xu, T. Li, J. Xiao, and X.-K. Xu, "Quantifying the Effects of Topology and Weight for Link Prediction in Weighted Complex Networks," *Entropy*, vol. 20, no. 5, pp. 363, 2018, <https://doi.org/10.3390/e20050363>.
- [29]. A. Hagberg, P. Swart, and D. S Chult, *Exploring network structure, dynamics, and*

<sup>1</sup> Second order random-walk

می‌شود. مراحل انجام یک پژوهش در چارچوب OGB در شکل ب-۱ نشان داده شده است. ابتدا یک مجموعه داده برای پژوهش مورد نظر انتخاب می‌شود. سپس OGB این مجموعه داده را به سه قسمت آموزشی، اعتبارسنجی و آزمایشی تقسیم و بارگزاری می‌کند. سپس با استفاده از یک الگوریتم داده‌های اعتبارسنجی و آزمایشی پردازش می‌شوند و در نهایت نتایج توسط ارزیاب مختص به مجموعه داده انتخاب شده ارزیابی می‌شود و در صورت انتشار مقاله در تابلوی امتیازات درج می‌شود.

مسائلی که در OGB برای آن‌ها مجموعه داده و ارزیاب فراهم شده اعم از پیش‌بینی پیوند، دسته‌بندی رأس، تکمیل گراف دانش و مسائل دیگر هستند. این مجموعه‌داده‌ها در زمینه‌هایی مانند شبکه‌های علمی، گراف‌های بیولوژیکی و شبکه‌های دانش هستند. در این بخش فقط به مجموعه‌داده‌های مرتبط با شبکه‌های علمی و مسئله پیش‌بینی پیوند می‌پردازیم.

استفاده از چنین چارچوب‌هایی به پژوهشگران کمک می‌کند نتایج پژوهش‌های خود را با سایر الگوریتم‌های پیشین مقایسه کنند. در این مقاله آزمایشات و ارزیابی نتایج مبتنی بر این چارچوب انجام می‌شود. به عنوان نمونه تابلوی نمایش و رتبه‌بندی ارزیابی تعدادی از مقالات، بر اساس معیار Hits@50 با مجموعه داده ogbl-collab در شکل ب-۲ نشان داده شده است.

انتقال به قسمتی از گراف که تا به حال دیده نشده است  $q$  می‌باشد. همچنین کوتاه‌ترین فاصله بین  $t$  و  $x$  است.

$$\alpha_{pq}(t, x) = \begin{cases} \frac{1}{p} & \text{if } d_{tx} = 0 \\ 1 & \text{if } d_{tx} = 1 \quad (p, q < 1) \\ \frac{1}{q} & \text{if } d_{tx} = 2 \end{cases} \quad \text{معادله الف-۱}$$

با تغییر  $q$  و  $p$  می‌توان سرعت کاوش در گراف را تنظیم کرد. به عبارت دیگر، با کاهش  $q$  قدم زدن تصادفی به سمت پیمایش عمقی<sup>۱</sup> گراف متمایل می‌شود و با کاهش  $p$  الگوریتم به سمت پیمایش سطحی<sup>۲</sup> گراف متمایل می‌شود.

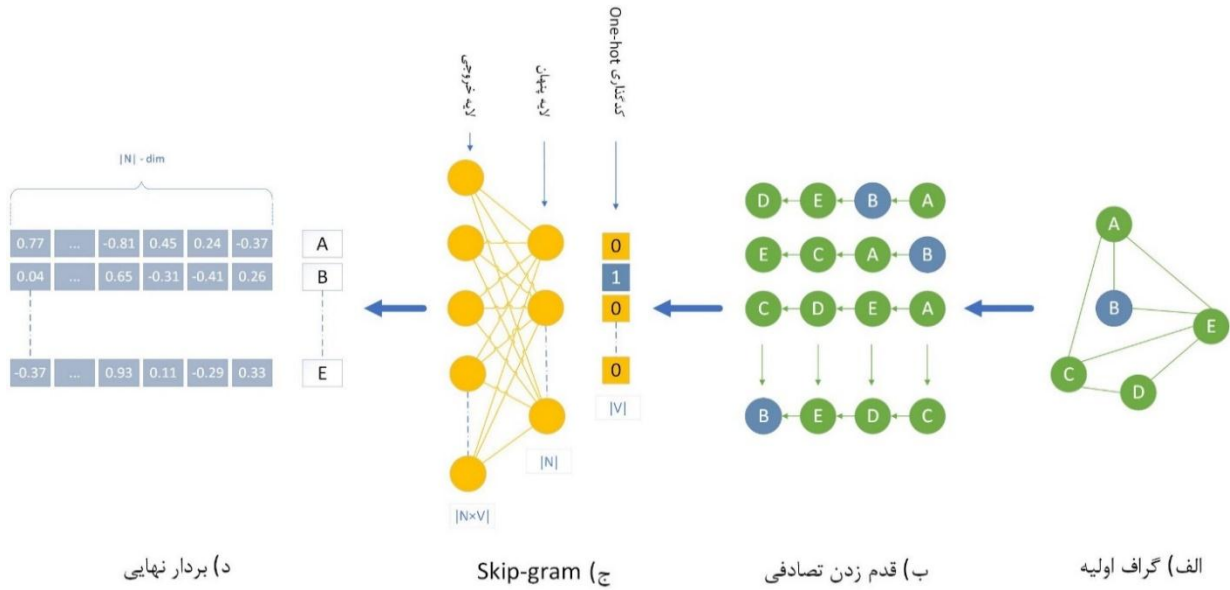
در ادامه هر یک از این دنباله‌های ساخته شده توسط قدم زدن تصادفی به عنوان ورودی به شبکه عصبی skip-gram داده می‌شود. مدل skip-gram یک شبکه عصبی با یک لایه پنهان است. این مدل آموزش داده می‌شود که احتمال وجود یک رأس (کلمه) را در صورت وجود یک رأس دیگر در دنباله محاسبه کند. خروجی مدل skip-gram یک نهفته سازی برای هر رأس است (شکل الف-۱، قسمت ج و د).

### پیوست ب: چارچوب OGB

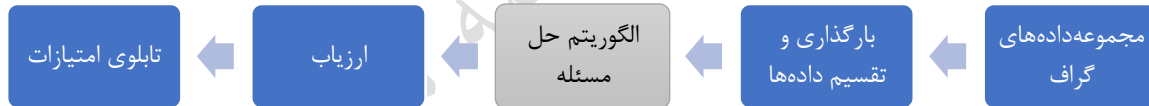
OGB یک چارچوب نرم‌افزاری متن باز برای انجام آزمایشات و ارزیابی نتایج برای مسائل مربوط به پردازش گراف است. این چارچوب با زبان پایتون و توسط دانشگاه استنفورد توسعه داده شده است [۲۷]. به طور کلی این چارچوب ارائه مجموعه داده‌ها و ارزیابی نتایج را در یک قالب استاندارد برای پژوهشگران فراهم می‌کند. در نهایت بهترین نتایج ثبت شده که موفق به ثبت مقاله شده باشند در جدول برترین نتایج در سایت OGB نمایش داده

<sup>1</sup> DFS

<sup>2</sup> BFS



شکل الف-۱) الگوریتم node2vec



شکل ب-۱) مراحل انجام یک مسئله گراف در OGB



Rank	Method	Ext. data	Test Hits@50	Validation Hits@50	Contact	References	#Params	Hardware	Date
1	GIDN@YITU	No	0.7096 ± 0.0055	0.9620 ± 0.0040	Zixiao Wang, Yu Zhang(ZhejiangLab, HUST)	Paper, Code	60,449,025	DepGraph@SCTS/CGCL	Oct 10, 2022
2	PLNLP + SIGN	No	0.7087 ± 0.0033	1.0000 ± 0.0000	Liang Yao (Tencent)	Paper, Code	34,980,864	Tesla-P40 (24G GPU)	Apr 7, 2022
3	PLNLP (random walk aug.)	No	0.7059 ± 0.0029	1.0000 ± 0.0000	Zhitao Wang (WeChat@Tencent)	Paper, Code	34,980,864	Tesla-P40 (24G GPU)	Dec 21, 2021
4	HOP-REC	No	0.7012 ± 0.0016	1.0000 ± 0.0000	Bo-Yu Lin (CFDA & CLIP Labs)	Paper, Code	30,191,104	CPU	Oct 21, 2021
5	PLNLP+ LRGA	No	0.6909 ± 0.0055	1.0000 ± 0.0000	Hao Xu	Paper, Code	35,200,656	NVIDIA Tesla V100 (32GB GPU)	Jun 26, 2022
17	DeepWalk	No	0.5037 ± 0.0034	Please tell us	Hao Xiong (DGL)	Paper, Code	61,390,187	g4dn.2xlarge, T4 (15GB GPU)	Jun 30, 2020
18	Node2vec	No	0.4888 ± 0.0054	0.5703 ± 0.0052	Matthias Fey – OGB team	Paper, Code	30,322,945	GeForce RTX 2080 (11GB GPU)	Jun 22, 2020
19	GraphSAGE	No	0.4810 ± 0.0081	0.5688 ± 0.0077	Matthias Fey – OGB team	Paper, Code	460,289	GeForce RTX 2080 (11GB GPU)	Jun 24, 2020
20	GCN (val as input)	No	0.4714 ± 0.0145	0.5263 ± 0.0115	Matthias Fey – OGB team	Paper, Code	296,449	GeForce RTX 2080 (11GB GPU)	Oct 19, 2020
21	GCN	No	0.4475 ± 0.0107	0.5263 ± 0.0115	Matthias Fey – OGB team	Paper, Code	296,449	GeForce RTX 2080 (11GB GPU)	Jun 24, 2020
22	Matrix Factorization	No	0.3886 ± 0.0029	0.4896 ± 0.0029	Matthias Fey – OGB team	Paper, Code	60,514,049	GeForce RTX 2080 (11GB GPU)	Jun 22, 2020

شکل ب-۲: تابلوی امتیازات برای مجموعه داده ogbl-collab