



**Journal of National Studies on
Librarianship and Information
Organization
(NASTINFO)**



Research Article

Analysis of Citation-based Indicators to Determine the Relevance of Articles

M. Goltaji¹

J. Abbaspour²

A. Jowkar³

S. M. Fakhrahmad⁴

A. Nikseresht⁵

- ¹. PhD Candidate, Knowledge and Information Science Department, International Division, Shiraz University, Shiraz, Iran, marzieh.goltaji@gmail.com
- ². PhD in Knowledge & Information Science, Assistant Professor, School of Education & Psychology, Shiraz University, Shiraz, Iran, (Corresponding Author), javad.abbaspour@gmail.com
- ³. PhD in Knowledge & Information Science, Professor, School of Education & Psychology, Shiraz University, Shiraz, Iran, ajowkar2003@yahoo.com
- ⁴. PhD in Computer Engineering, Associate Professor, Department of Computer Science and Engineering & IT, Shiraz University, Shiraz, Iran, fakhrahmad@shirazu.ac.ir
- ⁵. PhD in Computer Engineering, Assistant Professor, School of Education & Psychology, Shiraz University, Shiraz, Iran, nikseresht@gmail.com

Abstract

Purpose: The present study aimed to investigate the potential of citation-based indicators (Co-Citation, Bibliographic Coupling, Amsler, PageRank, HITS) to determine the relevance of articles.

Method: This is applied research with correlational approach. The population consisted of 26,262 articles in the PubMed Central open access subset of the CITREC, which had citation relationship with other articles based on all three traditional citation-based indicators (Co-Citation, Bibliographic coupling, Amsler). From among the citations in the research population, 30 were selected as basic ones, and the full-text of them were retrieved based on the mesh similarity. Then the similarities among the retrieved documents were extracted based on citation-based indicators. Each of the citation-based metrics was considered as independent variable and the mesh similarity as dependent variable. A MySQL database was created using WampServer simulation software and PHP My Admin. Then, using online demo of the CITREC test collection, an output was prepared. By entering the output into the MySQL database which contains the research data set, the main structure of its tables was created. Finally, by studying all the required codes from the CITREC source code package, we attempted to enter the required codes by applying necessary changes. The results were entered in the created MySQL database. By writing a query in SQL language, the set citation network was completely extracted and stored in a Comma-separated values (CSV) file. Then, a program was written in Python that could open and process this large file and calculate PageRank and HITS numbers (authority and Hub).

Findings: The results showed that all six measures studied had a significant and positive correlation with the relevance of articles. In other words, with increasing the values of each measure, the degree of relevance of the articles also increased. The highest correlation with the relevance of the articles belonged to the Amsler measure, followed by the Bibliographic Coupling. After Amsler and Bibliographic Coupling, the highest correlation was observed in the HITS(Authority) variable, and the PageRank variable was in the fourth place; Finally, the lowest correlation with the relevance of the articles was related to the Co-Citation and the HITS (Hub). Therefore, among the known Citation- based measure studied here, Amsler, Bibliographic Coupling, HITS(Authority) and PageRank metrics, respectively, had more potential to determine the relevance of articles rather than others.

Conclusion: Based on the findings, it can be concluded that the citation-based metrics studied are able to estimate the degree of relevance of articles. Therefore, they can be used in various information retrieval platforms, including search engines, citation- based databases, recommender systems, and even digital libraries to access articles, suggest similar articles, and rank retrieved results; Also, the Amsler measure as the less used in information retrieval systems than the two traditional Measure (Co- Citations and Bibliographic Coupling) needs to be considered more than ever. On the other hand, despite the fact that Co- Citations measure is used in some international information retrieval databases (such as Science Direct and CiteSeer) to retrieve relevant documents and suggest similar documents, it is less efficient than other metrics.

Keywords: Citation- Based Metrics, Relevance of Articles, Co-citation, Bibliographic Coupling, Amsler, PageRank, HITS

Follow this and additional works at: <http://nastinfo.nlai.ir/>



This work is licensed under the terms of the Creative Commons Attribution 4.0 International License: <http://creativecommons.org/licenses/by/4.0>

Recommended Citation

Goltaji, M., Abbaspour, J., Jowkar, A., Fakhrahmad, S., nikseresht, A. (2021). Analysis of Citation-based Indicators to Determine the Relevance of Articles. *Journal of National Studies on Librarianship and Information Organization (NASTINFO)*, 32 (3): 56-76

This Review Article is brought to you for free and open access by Journal of National Studies on Librarianship and Information Organization (NASTINFO).

Received:06, Aug. 2021; accepted:18, Sep. 2021



فصلنامه علمی - پژوهشی
مطالعات ملی کتابداری و سازماندهی اطلاعات



مقاله پژوهشی

تحلیل سنجه‌های استنادمحور برای تعیین میزان ربط مقاله‌ها

مرضیه گل‌تاجی^۱

جواد عباس‌پور^۲

عبدالرسول جوکار^۳

سیدمصطفی فخراحمد^۴

علیرضا نیک‌سرشت^۵

- ^۱ دانشجوی دکتری علم اطلاعات و دانش‌شناسی، واحد بین‌الملل دانشگاه شیراز، شیراز، ایران <mailto:University of Cambridge>
- ^۲ استادیار گروه علم اطلاعات و دانش‌شناسی، دانشکده علوم تربیتی و روان‌شناسی، دانشگاه شیراز، شیراز، ایران (نویسنده مسئول)، javad.abbaspour@gmail.com
- ^۳ استاد گروه علم اطلاعات و دانش‌شناسی، دانشکده علوم تربیتی و روان‌شناسی، دانشگاه شیراز، شیراز، ایران jowkar2003@yahoo.com
- ^۴ دانشیار گروه مهندسی و علوم کامپیوتر و فناوری اطلاعات، دانشکده مهندسی برق و کامپیوتر، دانشگاه شیراز، شیراز، ایران mfakhrmahmad@gmail.com
- ^۵ استادیار گروه علم اطلاعات و دانش‌شناسی، دانشکده علوم تربیتی و روان‌شناسی، دانشگاه شیراز، شیراز، ایران nikseresht@gmail.com

چکیده

هدف: شناخت توانایی سنجه‌های استنادمحور (هم‌استنادی، زوج کتاب‌شناختی، امسلر، پیچ‌رنک و هیتس (اعتبار و کانون)) برای تعیین میزان ربط مقاله‌ها با یکدیگر.

روش: پژوهش حاضر از نظر هدف، کاربردی و از لحاظ شیوه گردآوری داده‌ها، پژوهشی توصیفی از نوع همبستگی است. جامعه آماری، مجموعه مقالات موجود در زیرمجموعه دسترسی آزاد پاب‌مد سنترال مجموعه آزمون سایترک بود که بر اساس سه سنجه هم‌استنادی، زوج کتاب‌شناختی و امسلر با سایر مقالات رابطه استنادی داشتند. از میان ۲۶۲۶۲ مقاله، ۳۰ مقاله به‌عنوان مقالات پایه انتخاب شد و مقالات مرتبط با هر یک از آن‌ها بر اساس سنجه ربط مش‌بازیابی گردید؛ هر یک از سنجه‌های استنادمحور متغیر مستقل و سنجه ربط مش متغیر وابسته بود. با استفاده از نرم‌افزار شبیه‌ساز ومپ‌سرور و پی‌اچ‌پی.مای‌ادمین یک پایگاه مای.اس.کیوال ایجاد شد؛ سپس، با مطالعه کلیه کدهای مورد نیاز از بسته کد منبع سایترک، کدهای لازم با اعمال

تغییرات ضروری، اجرا و نتایج حاصل در پایگاه مای.اس.کیوال وارد شد. با نوشتن پرس و جو به زبان اس.کیوال، شبکه استنادی مجموعه به صورت کامل استخراج شد سپس با کدنویسی به زبان پایتون اعداد مربوط به پیچ‌رنک و هیتس (اعتبار و کانون) به صورت جداگانه محاسبه گردید.

یافته‌ها: نتایج نشان داد تمامی شش سنجه در سطح یک صدم همبستگی معنادار و مثبت با میزان ربط مقاله‌ها داشت؛ به عبارت دیگر، با افزایش مقادیر هر یک از سنجه‌ها، درجه ربط مقاله‌ها نیز افزایش یافت. بیشترین میزان همبستگی مربوط به سنجه امسلر و پس از آن، زوج کتاب‌شناختی بود. پس از سنجه‌های امسلر و زوج کتاب‌شناختی، بیشترین همبستگی میان متغیر هیتس (اعتبار) با ربط مقاله‌ها بود. متغیر پیچ‌رنک در مرتبه چهارم قرار داشت؛ در نهایت، کم‌ترین میزان همبستگی با ربط مقاله‌ها، مربوط به سنجه‌های هم‌استنادی و هیتس (کانون) بود؛ بنابراین، از میان سنجه‌های استنادی بررسی شده در این پژوهش، سنجه‌های امسلر، زوج کتاب‌شناختی، هیتس (اعتبار) و پیچ‌رنک بیش از سایر سنجه‌ها از پتانسیل لازم برای تعیین میزان ربط مقاله‌ها برخوردار بودند.

نتیجه‌گیری: بر اساس یافته‌های پژوهش می‌توان گفت سنجه‌های استنادمحور مطالعه شده قادرند درجه ربط مقاله‌ها را برآورد کنند و در بافتارهای مختلف بازیابی اطلاعات شامل موتورهای جست‌وجو، پایگاه‌های اطلاعاتی و استنادی، سامانه‌های پیشنهاددهنده و حتی کتابخانه‌های دیجیتالی برای دسترسی به مقالات مرتبط، پیشنهاد مقالات مشابه و رتبه‌بندی نتایج بازیابی کاربرد داشته باشند؛ همچنین، لازم است به سنجه امسلر که نسبت به دو سنجه سنتی هم‌استنادی و زوج کتاب‌شناختی، در سامانه‌های اطلاعاتی کمتر استفاده شده است، بیش از پیش توجه شود؛ از طرفی، علیرغم اینکه سنجه هم‌استنادی در برخی از پایگاه‌ها و سامانه‌های بازیابی اطلاعات بین‌المللی (مانند ساینس دایرکت و سایت سیر) برای بازیابی مدارک مرتبط و پیشنهاد مدارک مشابه استفاده می‌شود در مقایسه با سایر سنجه‌ها از کارایی کمتری برخوردار است.

کلیدواژه‌ها: ربط مقاله‌ها، هم‌استنادی، زوج کتاب‌شناختی، امسلر، پیچ‌رنک، هیتس، سنجه‌های استنادمحور

استناد به این مقاله:

گل تاجی، مرضیه، عباس پور، جواد، جوکار، عبدالرسول، فخر احمد، سید مصطفی، نیک سرشت، علیرضا. (۱۴۰۰). تحلیل سنجه‌های استنادمحور برای تعیین میزان ربط مقاله‌ها (۱۴۰۰). فصلنامه مطالعات ملی کتابداری و سازماندهی اطلاعات، ۳۲

(۳): ۵۶-۷۶

10.30484/NASTINFO.2021.2942.2069

دریافت: ۱۴۰۰/۰۵/۱۵؛ پذیرش: ۱۴۰۰/۰۶/۲۷

مقدمه

(Hariri, 2011; Levene, 2007)؛ از طرفی، نتایج

پژوهش‌های متعدد نشان داده است که کاربران تنها به چند نتیجه نخست توجه می‌کنند و سایر نتایج بازیابی شده نادیده می‌گیرند (Bar-Ilan, Levene & Mat-Hassan, 2006; Lewandowski, 2008; Agrahri, Manickam, & Riedl, 2008; Bar-Ilan, Keenoy, Levene, & Yaari, 2009; Lewandowski, 2017). مهم‌ترین عواملی که می‌تواند بر این امر تأثیر مستقیم داشته باشد، کیفیت سنجه‌های ربط و شباهت است؛ به عبارت دیگر، هرچه سنجه‌های بهتری به کار گرفته شود، می‌توان انتظار نتایج بازیابی بهتری را داشت؛ بنابراین، سنجه‌های به کاررفته در سامانه‌ها برای شباهت‌یابی و ربط مدارک بایستی به حدی اثربخش باشد که نه تنها کلیه مدارک مرتبط با پرس‌وجو کاربر را بازیابی کند؛ بلکه مدارک را به نحوی رتبه‌بندی کند که با حداقل کوشش و زمان، نیاز اطلاعاتی کاربر پاسخ داده شود.

(McGill, Koll, & Noreault, 1979) ۶۷ سنجه ربط و شباهت مختلف به کاررفته در بازیابی اطلاعات را بررسی و مقایسه کردند (Quated in Zhang and Korfage, 1999). بر اساس جست‌وجوهای انجام‌شده، در حال حاضر آمار دقیقی از تعداد سنجه‌های شباهت و ربط ابداع‌شده توسط متخصصان در دسترس نیست؛ ولی اکنون که در حدود چهار دهه از انتشار اثر آن‌ها می‌گذرد، به احتمال فراوان تعداد آن‌ها بسیار بیشتر از گذشته خواهد بود. سنجه‌های شباهت را بر اساس آنچه مبنای محاسبه شباهت قرار می‌گیرد، به سه گروه اصلی دسته‌بندی کرده‌اند: سنجه‌های متن‌محور، سنجه‌های استنادمحور و سنجه‌های ترکیبی (Boyack and Klavans, 2010). هرچند امروزه به دلیل گسترش و بهبود نظام‌های پردازش متن، سنجه‌های متن‌محور تحولات چشم‌گیری داشته‌اند؛ ولی سنجه‌های استنادمحور

سامانه‌های بازیابی اطلاعات شامل موتورهای کاوش، کتابخانه‌های دیجیتال و پایگاه‌های اطلاعاتی، از هر نوع و در هر سطح از پیچیدگی، هدف بنیادی مشترکی را دنبال می‌کنند و آن نمایش مدارک یا اطلاعات مرتبط با درخواست یا نیاز اطلاعاتی کاربران است (نشاط، ۱۳۸۲). در این سامانه‌ها برای تعیین میزان ربط پرس‌وجو^۱ با مدرک یا مدارک از معیارهایی با عنوان "سنجه‌های شباهت" استفاده می‌شود و چنانچه مقدار شباهت از یک حد آستانه بیشتر باشد، مدارک به عنوان مدارک مرتبط بازیابی می‌شوند؛ همچنین، این سنجه‌ها شاخص‌های مناسبی برای رتبه‌بندی مدارک بازیابی‌شده‌اند و به سامانه‌ها اجازه می‌دهند تا با مرتب ساختن مدارک بر اساس میزان ربط، مدارک مرتبط‌تری را در صدر نتایج قرار دهند (Zhang and Korfage, 1999)؛ همچنین سنجه‌های شباهت در بسترهای دیگری مانند طبقه‌بندی خودکار مدارک؛ تشکیل خوشه‌های مدارک مشابه؛ ایجاد سامانه‌های بازیابی گرافیکی و ترسیم نقشه‌های علوم نیز به کار می‌روند (Torres, Basnet, Sung, Mukkamala, & Ribeiro, 2009; Jouili, Tabbone, & Valveny, 2010; Lin, Jiang, & Lee, 2013; Amer and Abdalla, 2020; Eminagaoglu, 2020).

علی‌رغم کاربرد وسیع و اهمیت سنجه‌های شباهت در بازیابی اطلاعات به طور گسترده و سامانه‌های بازیابی اطلاعات به صورت محدود، چالش اساسی این موضوع، مسئله ربط و رضایت کاربر است؛ نتایج پژوهش‌ها نشان می‌دهد کاربران از نتایج بازیابی‌شده توسط سامانه‌های بازیابی اطلاعات رضایت ندارند و بین نتایج بازیابی‌شده از طریق این سامانه‌ها با نظر کاربران قرابت چندانی وجود ندارد (سعدین‌خرم و عباس‌پور، ۱۳۹۸؛ حاجیان و چشمه‌سهرابی، ۱۳۹۹؛ Bar-Ilan, Keenoy, Yaari,

1. Query

یکدیگر در مدارک بعدی استناد می‌شود (Small, 1973). در زوج کتاب‌شناختی، به مآخذ مشترک توجه می‌شود؛ به عبارتی، وجود یک مآخذ در دو مقاله، واحد اندازه‌گیری حد اشتراک آن دو مقاله محسوب می‌شود و هرچه دو مقاله در تعداد بیشتری از مآخذ خود مشترک باشند از لحاظ محتوایی به یکدیگر نزدیک‌ترند (Kessler, 1963). برای محاسبه شباهت مدارک با سنجه امسler که شباهت پیوندی^۴ نیز نامیده می‌شود، هم‌زمان هر دو سنجه زوج کتاب‌شناختی و هم‌استنادی لحاظ می‌شود (Bichteler and Eaton, 1980).

در کنار سه سنجه استنادمحور پیش‌گفته، که از قدمت بیشتری برخوردارند، دو سنجه هیتس^۵ و پیچرنک^۶ نیز وجود دارد که در اصل برای محیط وب و به‌ویژه رتبه‌بندی صفحات وب ایجاد شدند؛ اما پژوهشگران قابلیت به‌کارگیری آنان برای شناسایی مدارک علمی مهم و معتبر و رتبه‌بندی آن‌ها را بررسی و عملکرد آن‌ها را برای سنجش شباهت مدارک تأیید کردند (Zhuge and Zhang, 2010; Yin, Huang, & Li, 2011; Liu, Zhang, & Guo, 2012).

هر دو سنجه پیچرنک و هیتس، سنجه‌های بازگشتی و تکرارشونده مبتنی بر پیوند مدارک در محیط وب هستند؛ با این تفاوت که گرایش سنجه پیچرنک به سمت بازتاب مقالات قدیمی و سنجه هیتس به مقالات جدید است (Jiang, Sun, Yang, Zhuge, and Yao, 2014; Devi, Gupta and Dixit, 2014). در رتبه‌بندی صفحات وب با هیتس، ابتدا جست‌وجوی موضوعی با یک یا چند عبارت پرسشی انجام می‌گیرد، و تعدادی صفحه دارای بالاترین رتبه استخراج می‌شوند که

نیز دارای برخی از ویژگی‌ها و قابلیت‌های منحصر به فردند که آن‌ها را از برخی جهات از رویکردهای متن‌محور متمایز می‌سازد؛ از جمله اینکه، سنجه‌های استنادمحور که در قالب پرس‌وجو توسط کاربران به نظام داده می‌شوند تحت تأثیر عبارت‌ها و کلمات گوناگونی قرار نمی‌گیرند؛ دسترسی به منابع و متون بین‌رشته‌ای را فراهم می‌سازند؛ و منابعی را که با یک موضوع خاص ارتباط دارند ولی با سامانه‌های بازیابی متنی ردگیری نمی‌شوند، نشان می‌دهند (Smith, 1981). علاوه بر این، داده‌های استنادی به‌راحتی قابل درک هستند و به‌سرعت هم‌گردآوری می‌شوند (Su and et al., 2011).

در برخی از سامانه‌های بازیابی اطلاعات، به‌ویژه پایگاه‌های استنادی و برخی موتورهای جست‌وجو، به‌منظور محاسبه میزان ربط مدارک با پرس‌وجو یا سایر مدارک، از پیوندهای استنادی میان مدارک و ساختار گراف آن‌ها استفاده می‌شود (Reyhani Hamedani, Lee, and Kim, 2013; Reyhani Hamedani and Kim, 2021). برخی از متخصصان معتقدند که استنادها اغلب برای اندازه‌گیری شباهت ضمنی میان مدارک علمی استفاده می‌شوند و بازیابی مدارک با استفاده از پیوندهای استنادی قادر است مدارک شناسایی‌نشده با پردازش متنی را بیابد (Eto, 2019)؛ در مقابل، عده‌ای هم عقیده دارند که این‌دسته از سنجه‌ها، تنها روابط استنادی میان مقالات علمی را در نظر می‌گیرند و از محتوای مقالات غافل می‌شوند (Reyhani Hamedani and et al., 2013).

از مهم‌ترین سنجه‌های سنتی مبتنی بر استناد می‌توان به سنجه‌های هم‌استنادی^۱، زوج کتاب‌شناختی^۲، و امسler^۳ اشاره کرد. هم‌استنادی پیوندی است که توسط نویسندگان جدید، میان مقالات پیشین برقرار می‌شود؛ بنابراین، هم‌استنادی، بسامد تعداد دفعاتی است که به دو مدرک همراه

4. Linkage Similarity
5. HITS
6. PageRank

1. Cocitation
2. Bibliographic Coupling
3. Amster

کتاب‌شناختی و استناد مستقیم^۵ را از نظر قابلیت شناسایی بهینه پیشگامان پژوهش و مدارک هسته با یکدیگر مقایسه کردند. یافته‌ها نشان داد استناد مستقیم در مجموع بهترین و هم‌استنادی بدترین عملکرد را نسبت به سایرین داشتند.

پس از ظهور شبکه جهانی وب- که اساس آن را فرایوندها^۶ تشکیل می‌دادند- سوالی که مطرح شد این بود که آیا از سنج‌های پیشین، به‌ویژه سنج‌های استنادمحور که قرابت نزدیکی با فرایوندها داشتند، می‌توان در محیط وب هم استفاده کرد یا خیر؟ و در صورت مثبت بودن پاسخ این سوال، هر یک از این سنج‌ها در چه بستری کارایی بهتری دارند؟ با توجه به اینکه استناد و سنج‌های استنادمحور، اعتباری که مدارک به یکدیگر می‌دهند را در نظر نمی‌گیرند، توانایی سنج‌های مبتنی بر گراف مانند الگوریتم‌های پیچ‌رنک و هیتس در شباهت‌یابی مدارک مدنظر پژوهش‌ها قرار گرفت و قابلیت این سنج‌های جدید استنادمحور در قیاس باهم، و در مقایسه با سنج‌های پیشین آزموده شد.

[Ma, Guan and Zhao \(2008\)](#) با مقایسه دو سنج

پیچ‌رنک و تعداد استنادات، تلاش کردند تا روش جدیدی برای اندازه‌گیری اهمیت مقالات علمی مبتنی بر پیچ‌رنک گوگل عرضه کنند. آن‌ها نتایج حاصل از رتبه‌بندی مقالات علمی با پیچ‌رنک را با تعداد استنادات به آن‌ها مقایسه و همبستگی زیادی بین استناد و پیچ‌رنک مشاهده کردند.

مقایسه این دو سنج در پژوهش [Herskovic and](#)

[Bernstam \(2005\)](#) تکرار شد و نتایج نشان داد

عملکرد پیچ‌رنک نسبت به تعداد استناد بهتر بود. در

پژوهشی دیگر [Lin \(2008\)](#) دو سنج جدید

استنادمحور پیچ‌رنک و هیتس را مقایسه کرد. وی قابلیت

"مجموعه ریشه"^۱ را شکل می‌دهند. این صفحات مرتبط‌ترین صفحات به پرس‌وجو مطرح‌شده فرض می‌شوند. سپس، صفحاتی که به وسیله پیوند، به مجموعه ریشه پیوند داده‌اند و یا مجموعه ریشه با آن‌ها پیوند برقرار کرده است به مجموعه افزوده می‌شوند و مجموعه‌ای بزرگ‌تر با نام "مجموعه پایه"^۲ به وجود می‌آید که مبنای محاسبه الگوریتم هیتس است. این الگوریتم به هر صفحه در این مجموعه یک امتیاز اعتبار^۳ و یک امتیاز کانون^۴ اختصاص می‌دهد ([Kleinberg, 1999](#)).

هیتس دو مقدار اعتبار و کانون را به ازای هر صفحه که به صورت بازگشتی متقابل روی هم اثر دارند محاسبه می‌کند؛ به عبارت دیگر، یک صفحه با امتیاز اعتبار بالا به وسیله تعدادی از صفحات با امتیاز کانون بالا (تعداد زیادی از صفحات با امتیاز کانون بالا به آن پیوند داده باشند) و یک صفحه با امتیاز کانون بالا به تعدادی صفحه با امتیاز اعتبار بالا اشاره می‌کند (صفحه‌ای است که به صفحات با امتیاز اعتبار بالا پیوند داده باشد) ([Kleinberg, 1999; Kirsch, Gnasa, Won, and Cremers, 2008](#)).

عملکرد سنج‌های استنادمحور در بازیابی از گذشته تا به امروز همواره مدنظر پژوهشگران بوده است. در پژوهش‌های اولیه، عملکرد سنج‌های تعداد استنادات، زوج کتاب‌شناختی، هم‌استنادی، و امسلر در موضوع‌های مختلف بازیابی اطلاعات و همچنین در قیاس با یکدیگر آزموده شد. یافته‌های پژوهش [Bichteler and Eaton \(1980\)](#) بیانگر برتری مقیاس شباهت پیوندی امسلر (ترکیب زوج کتاب‌شناختی و هم‌استنادی) در مقایسه با سنج زوج کتاب‌شناختی به‌تنهایی بود. در پژوهش دیگری، [Shibata, Kajikawa, Takeda, and Matsushima \(2009\)](#) سه سنج هم‌استنادی، زوج

4. Hub

5. Direct citation

6. Hyperlinks

1. Root Set

2. Base Set

3. Authority

زیست‌فناوری و همچنین ناتوانی سنجه‌های متن‌محور در بازیابی همه مدارک مرتبط، الگوریتم گراف چندتایی محتوا^۹ را پیشنهاد دادند. نتایج آزمون این الگوریتم روی داده‌های موجود در پاب‌مد نشان داد که این الگوریتم به‌لحاظ بازیابی و رتبه‌بندی مقالات به‌طور معناداری از الگوریتم سنتی متن‌محور پاب‌مد بهتر عمل می‌کند.

Janssens, Gwinn, Brockman, Powell, & Goodman (2020) روش جست‌وجوی استنادمحور را

گسترش دادند. این روش که برای کارایی بیشتر نسبت به کلیدواژه سنتی طراحی شده بود هم‌استنادها^{۱۰} نام داشت. نتایج نشان داد این روش، روندی مؤثر و درست برای یافتن مقالات مرتبط است.

با مرور پژوهش‌های انجام‌شده پیشین، می‌توان گفت همان‌گونه که Thompson, Panchev & Oakes,

(2015) اشاره می‌کنند گرچه به سنجه‌های ربط و

شباهت‌یابی در بسیاری از نظام‌های بازیابی و فرایندهای پردازش زبان طبیعی برای یافتن مرتبط‌ترین مدارک و

مدارک مشابه با مدرک اصلی نیاز است و موفقیت این سامانه‌ها نیز تا حد زیادی به عملکرد صحیح این سنجه‌ها

وابسته است؛ اما عملکرد این سنجه‌ها از سامانه‌ای به سامانه دیگر می‌تواند متفاوت باشد؛ سنجه‌ای که در یک

سامانه اثربخش بود، در سامانه دیگر مؤثر نبود. همچنین در عمل نیز هر یک از این سنجه‌ها به‌صورت مجزا در

پایگاه‌های اطلاعاتی، کتابخانه‌های دیجیتال و موتورهای جست‌وجو به کار رفته‌اند؛ برای مثال، ساینس‌دایرکت^{۱۱}،

از سنجه هم‌استنادی برای پیشنهاد مدارک مشابه استفاده می‌کند (Eto, 2013; Eto, 2019)؛ در پایگاه‌های وب

آو ساینس^{۱۲} و اسکوپوس^{۱۳} از سنجه زوج کتاب‌شناختی Burnham, 2006; Ahlgren and Jarneving,)

پیچ‌رنک و هیتس را در شبکه مقالات مرتبط پاب‌مد^۱ برای بازیابی متون زیست‌پزشکی بررسی کرد. نتایج این پژوهش نشان داد که پیچ‌رنک نسبت به هیتس برای تحلیل ساختار پیوند شبکه‌های مدارک مرتبط مؤثرتر است.

Yin and et al.,(2011) سه الگوریتم تحلیل پیوندی

را که عبارت بودند از توزیع درجه^۲، پیچ‌رنک و هیتس از نظر کارآمدی در بازیابی متون زیست‌پزشکی

بررسی کردند. تحلیل داده‌ها نشان داد که هرچند هر سه الگوریتم مذکور منجر به بهبود بازیابی متون زیست‌پزشکی

می‌شوند؛ اما الگوریتم توزیع درجه نسبت به سایر موارد عملکرد بهتری دارد و پیچ‌رنک و هیتس به‌ترتیب در

رتبه‌های بعدی قرار می‌گیرند؛ آن‌ها همچنین یک چارچوب ترکیبی احتمالی^۳ را پیشنهاد کردند که اطلاعات

استنادی را با مدل وزن‌دهی احتمالی محتوامحور^۴ ترکیب می‌کند. بر اساس بررسی‌های آنان این چارچوب باعث

بهبود بازیابی متون زیست‌پزشکی می‌شود. Yoon, Kim & Park(2016) سنجه‌های رتبه شباهت^۵، شباهت

معکوس (آر.وی.اس رتبه شباهت)^۶، رتبه نفود (رتبه پی)^۷ و رتبه رابط (رتبه سی)^۸ را معرفی کردند که به‌ترتیب

نسخه بازگشتی هم‌استنادی، زوج کتاب‌شناختی و امسلرند. رتبه رابط نمره شباهت را بر اساس تعداد رابط‌ها در گراف

غیرجهت‌دار محاسبه می‌کرد، این سنجه‌ها در مقایسه با نسخه‌های غیربازگشتی خود عملکرد بهتری داشتند و از

میان همه آن‌ها رتبه رابط بهترین عملکرد را داشت. Colavizza, Boyack, Van Eck, & Waltman

(2018) در چهار نشریه از رشته‌های مختلف شباهت جفت مقالات هم‌استناد در سطوح مختلف را مطالعه کردند.

نتایج نشان داد که با کاهش سطح هم‌استنادی شباهت مقالات افزایش می‌یابد. (Jiang and et al. (2019) با توجه به نیاز پژوهشگران به دسترسی به مقالات

7. Penetrating (P-rank)
8. Connectors rank (C-rank)
9. Content Tuple Graph Algorithm (CTGA)
10. CoCites
11. ScienceDirect
12. Web of Scienc(WoS)
13. Scopus

1. PubMed Related Articles
2. Degree Distribution
3. probabilistic combination framework
4. Content-based probabilistic weighting model
5. Simrank
6. Reverse Simrank(rvs-Simrank)

۲۵۵۳۳۹ عنوان مقاله است. سایترک کد دسترسی آزاد جاوا برای محاسبه ۳۵ سنجه شباهت مبتنی بر استناد و متن را عرضه می‌کند که زوج کتاب‌شناختی، هم‌استنادی و امسلر از جمله سنجه‌های شباهت مبتنی بر استناد در این مجموعه هستند. در این پژوهش، برای فراهم کردن امکان مقایسه شاخص‌ها با یکدیگر مقالاتی انتخاب شد که شاخص‌های استنادمحور بررسی شده را دارا باشند؛ به عبارت دیگر، بر اساس سه سنجه استنادمحور هم‌استنادی، زوج کتاب‌شناختی و امسلر با سایر مقالات رابطه استنادی داشته باشند (پیوست ۱)؛ در مرحله بعد، با استفاده از نرم‌افزار شبیه‌ساز و مپ‌سرور^۶ و پی‌اچ‌پی.مای‌ادمین^۷ یک پایگاه مای.اس.کیوال^۸ ایجاد شد؛ سپس، با استفاده از نسخه نمایشی^۹ برخط مجموعه آزمون سایترک، یک خروجی تهیه شد و با وارد کردن آن به پایگاه مای.اس.کیوال حاوی مجموعه داده‌های پژوهش، ساختار اصلی جدول‌های آن ایجاد شد؛ در نهایت، با مطالعه کدهای لازم از بسته کد منبع^{۱۰} سایترک، تلاش شد تا این کدها با اعمال تغییرات ضروری، اجرا و نتایج حاصل در پایگاه مای.اس.کیوال وارد شود.

در زیرمجموعه‌ی دسترسی آزاد پاب‌مد سنترال مجموعه آزمایشی سایترک، برای سنجش میزان شباهت مدارک با پرس‌وجوها، از سرعنوان‌های موضوعی پزشکی (مش) تخصیص یافته به هریک از مدارک استفاده می‌شود. سرعنوان‌های موضوعی پزشکی (مش)، اصطلاحنامه چند سلسله‌مراتبی از توصیفگرهای موضوعی است. متخصصان

2008; Nicolaisen and Frandsen, 2012; Char and Ajiferuke, 2013)، در کتابخانه دیجیتالی ای.سی.ام^۱ و در پایگاه استنادی سایت سیر^۲ از هم‌استنادی برای پیشنهاد مدارک مشابه و مرتبط به کاربر استفاده می‌شود (Wanjantuk and Keane, 2004; Lu, Janssen, Milios, Japkowicz, and Zhang, 2007; Gipp and Beel, 2009; Eto, 2013; Eto, 2019). موتور جست‌وجوی گوگل از پیچ‌رنک و موتور جست‌وجوی اسک^۳ از هیتس برای بازیابی و رتبه‌بندی اسناد مرتبط با پرس‌وجو کاربر استفاده می‌نمایند (Thelwall, 2003; Goswami, Mantri, & Bhattacharya, 2017).

در پژوهش‌های اشاره‌شده توانایی هر یک از این سنجه‌ها برای تعیین میزان ربط مدارک با پرس‌وجو و سایر مدارک اغلب به‌تنهایی آزمون شده است و اطلاعات چندانی از میزان اثربخشی آن‌ها بر روی یک مجموعه مدارک واحد وجود ندارد. با بررسی این سنجه‌ها بر روی یک مجموعه یکسان می‌توان درک دقیق‌تری از توانایی هر یک از آن‌ها به دست آورد. براین اساس پژوهش حاضر قصد دارد تا توانایی این سنجه‌ها را برای تعیین میزان ربط مقاله‌ها بررسی کند. با توجه به مسئله گفته‌شده این پژوهش در پی پاسخ‌گویی به این سوال است که میزان همبستگی سنجه‌های هم‌استنادی، زوج کتاب‌شناختی، امسلر، پیچ‌رنک و هیتس با ربط مقاله‌ها چقدر است؟

روش پژوهش

جامعه آماری پژوهش، مجموعه مقالاتی بود که با نام "زیرمجموعه دسترسی آزاد پاب‌مد سنترال"^۱ در مجموعه آزمون سایترک^۲ آمده است. این مجموعه مشتمل بر

6. WampServer
7. phpMyAdmin
8. MySQL
9. Demo
10. Source Code

1. ACM Digital Library
2. CiteSeer
3. ASK
4. PubMed Central Open Access Subset (PMC OAS)
5. CITREC (<https://dke.uni-wuppertal.de/en/projects/citrec.html>)

یکدیگر به دست آمد که به صورت نزولی مرتب و در پایگاه مای.اس.کیوال با نام سیم مش اینترسکشن^۴ ذخیره شد؛ در مرحله بعد، با نوشتن پرس‌وجو به زبان اس.کیوال^۵، شبکه استنادی مجموعه به صورت کامل استخراج شد و در قالب یک فایل مقادیر جداشده با کاما^۶ ذخیره شد؛ به دلیل بزرگ بودن فایل، استفاده از نرم‌افزارهای موجود برای محاسبه پیچ‌رنک و هیتس امکان‌پذیر نبود، بنابراین یک برنامه به زبان پایتون نوشته شد که قادر بود این فایل بزرگ را باز و پردازش کند و نیز سنجه‌ها را اجرا کند. بدین ترتیب، اعداد مربوط به پیچ‌رنک و هیتس (اعتبار و کانون) به صورت جداگانه محاسبه و ذخیره گردید.

برای انتخاب نمونه، ابتدا به صورت تصادفی، ۳۰ مقاله از میان کلیه مقالات جامعه پژوهش، به عنوان مدارک پایه انتخاب شد و هر یک به منزله یک پرس‌وجو در نظر گرفته شد. با جست‌وجوی هر مقاله در جدول‌های مای.اس.کیوال ساخته شده در مرحله قبل، مدارک مرتبط با هر یک از ۳۰ مدرک پایه بر اساس سنجه شباهت مش شناسایی و بازیابی گردید؛ سپس، برای آن مقالات میزان ربط مقالات، سنجه‌های استنادمحور هم‌استنادی، زوج کتاب‌شناختی، امسلر، پیچ‌رنک، هیتس (اعتبار و کانون) از جدول‌های مربوط استخراج گردید. برای سنجش توانایی شش سنجه زوج کتاب‌شناختی، هم‌استنادی، امسلر، پیچ‌رنک، هیتس (اعتبار) و هیتس (کانون) برای تعیین میزان ربط مقاله‌ها با یکدیگر، این شش سنجه به عنوان متغیرهای مستقل و ربط مقاله‌ها بر اساس شاخص مش به عنوان متغیر وابسته وارد نرم افزار اس.پی.اس. اس نسخه ۲۳ شد. قبل از انجام آزمون همبستگی، فرض نرمال بودن

کتابخانه ملی پزشکی (ان.ال.ام) آمریکا^۱ این اصطلاحنامه را ایجاد و به صورت دستی، مناسب‌ترین توصیفگرها را به مدارک موجود در مجموعه دیجیتال ان.ال.ام مدلین^۲ اختصاص می‌دهند. سایترک، مش را به عنوان استاندارد طلایی و قضاوت موضوعی متخصصان در نظر می‌گیرد. این قضاوت‌ها می‌تواند برای ایجاد یک استاندارد طلایی برای ربط موضوعی مناسب باشند. به این صورت که هر چه تعداد اشتراک توصیفگرهای اختصاص یافته به دو مدرک بیشتر باشد آن دو مقاله با یکدیگر شباهت بیشتری دارند؛ به عبارتی، ربط آن دو مقاله با هم بیشتر است. این استاندارد طلایی، پژوهشگران را قادر می‌سازد تا میزان انعکاس ربط موضوعی سنجه‌های شباهت مبتنی بر استناد را بسنجند؛ بنابراین، این سنجه به عنوان شاخص استاندارد پژوهش حاضر در نظر گرفته شد و متغیر ملاک پژوهش بود. از این تکنیک در پژوهش‌های متعدد دیگری نیز برای سنجش شباهت مدارک استفاده شده است (Batet et al., 2010, Eto, 2012, Lin and Wilbur, 2007; Zhu et al., 2009 quoted in Gipp, Meuschke & Lipinski, 2015). در این مجموعه آزمون، سنجه شباهت مش (ربط متخصص) با اندازه‌گیری میزان شباهت میان هر زوج پرسش - مدرک، با استفاده از سنجه ژاکارد و با تقسیم اشتراک اصطلاحات مش بر اجتماع آنان محاسبه شده است.

برای تعیین ربط مدارک با یکدیگر، باید کدنویسی انجام و کد مش اینترسکشن^۳ اجرا می‌شد. بدین منظور، یک جدول جدید در پایگاه مای.اس.کیوال ایجاد و هر مدرک با کل مدارک موجود در مجموعه مقایسه شد؛ سپس میزان ربط مدارک بر اساس میزان اشتراک توصیفگرهای مش آنان با

4. Sim-Mesh-Intersections
5. SQL
6. Comma-separated values (CSV)

1. U.S. National Library of Medicine (NLM)
2. NLM's digital collection MEDLINE
3. Mesh-Intersections

متغیرها از طریق آزمون کولموگروف اسمیرنوف بررسی

شد (جدول ۱).

جدول ۱- نتایج آزمون کولموگروف اسمیرنوف سطح نرمال بودن متغیرها

متغیرها	شاخص آماری	هم استنادی	امسلر	زوج کتابشناختی	پیچ رنگ	هیتس (اعتبار)	هیتس (کانون)	ربط
تعداد نمونه	۱۵۰۰۰	۱۵۰۰۰	۱۵۰۰۰	۱۵۰۰۰	۱۵۰۰۰	۱۵۰۰۰	۱۵۰۰۰	۱۵۰۰۰
میانگین	۰/۱۲۵	۰/۷۹۴	۰/۶۵۷	۰/۰۰۰۰۰۰۰۳	۰/۰۰۰۰۰۰۰۳	۰/۰۰۰۰۱	۰/۰۰۰۰۴	۰/۲۳
انحراف استاندارد	۰/۷۷۶	۲/۲۳۸	۱/۷۵۹	۰/۰۰۰۰۰۰۰۳	۰/۰۰۰۰۰۰۰۳	۰/۰۰۰۰۴	۰/۰۰۰۰۸	۰/۰۷۸
سطح معناداری	۰/۰۰۱	۰/۰۰۱	۰/۰۰۱	۰/۰۰۱	۰/۰۰۱	۰/۰۰۱	۰/۰۰۱	۰/۰۰۱

مربوط به سنجه امسلر است. میزان این همبستگی ۰/۱۷۴ است و با اطمینان ۹۹ درصد می توان گفت که بین سنجه امسلر و ربط مقاله ها رابطه معناداری وجود دارد. با توجه به مثبت بودن جهت رابطه می توان گفت که با افزایش مقدار سنجه امسلر، میزان ربط مقاله ها نیز افزایش پیدا می کند ($r=0/174$ ، $p \leq 0/01$).

با توجه به این که برای تمامی متغیرها سطح معناداری از پنج صدم کمتر است، می توان گفت داده های پژوهش نرمال نیستند و بر این اساس از آزمون همبستگی اسپیرمن استفاده شد.

یافته های پژوهش

برای سنجش رابطه بین سنجه های استناد محور با میزان ربط مقاله ها، با توجه به نرمال نبودن توزیع داده ها از آزمون ضریب همبستگی اسپیرمن استفاده شد. نتایج این آزمون نشان داد که تمامی شش سنجه در سطح یک صدم همبستگی معناداری با میزان ربط مقاله ها دارد. همان گونه که جدول ۲ نشان می دهد بیشترین میزان همبستگی

جدول ۲- نتایج آزمون همبستگی اسپیرمن برای سنجش رابطه بین سنجه‌های استاندارد محور با میزان ربط مقاله‌ها

متغیر	۱	۲	۳	۴	۵	۶	۷
۱ هیتس (اعتبار)	۱						
۲ هیتس (کانون)	۰/۴۵۶**	۱					
۳ پیج رنک	۰/۸۳۶**	۰/۲۴۱**	۱				
۴ هم استنادی	۰/۳۵۴**	۰/۲۴۳**	۰/۲۶۸**	۱			
۵ زوج کتاب‌شناختی	۰/۲۷۹**	۰/۵۶۴**	۰/۱۵۱**	۰/۳۰۵**	۱		
۶ امسلر	۰/۳۳۱**	۰/۵۶۶**	۰/۱۹۶**	۰/۴۵۹**	۰/۹۶۵**	۱	
۷ ربط	۰/۱۵۳**	۰/۰۸۹**	۰/۱۱۳**	۰/۰۸۹**	۰/۱۶۱**	۰/۱۷۴**	۱

** $p \leq 0/01$

با توجه به یافته‌های به دست آمده می‌توان گفت که کم‌ترین میزان همبستگی با ربط مقاله‌ها، مربوط به سنجه‌های هم‌استنادی و هیتس (کانون) است. براین اساس می‌توان گفت که سنجه امسلر که ترکیبی از سنجه‌های زوج کتاب‌شناختی و هم‌استنادی است، بهتر از سایر سنجه‌ها، می‌تواند میزان ربط مقاله‌ها را تعیین کند.

بحث و نتیجه‌گیری

در این پژوهش توانایی سنجه‌های استاندارد محور هم‌استنادی، زوج کتاب‌شناختی، امسلر، پیج‌رنک و هیتس (اعتبار و کانون) برای تعیین میزان ربط مقاله‌ها تخمین زده شد. نتایج نشان داد که در بین سنجه‌های استاندارد محور، سنجه امسلر، زوج کتاب‌شناختی به ترتیب بیشترین توان را در تعیین میزان ربط مقاله‌ها داشتند. سنجه امسلر، هر دو سنجه هم‌استنادی و زوج کتاب‌شناختی را در خود لحاظ می‌کند. در این سنجه، به جای این‌که از هریک از دو شاخص پیشین به تنهایی استفاده شود، آن دو با یکدیگر ترکیب می‌شوند. به این ترتیب از تمامی علائمی که به نحوی نشانگر ارتباط موضوعی مدارک هستند استفاده می‌شود (Bichteler and Eaton, 1980). این یافته با نتایج به دست آمده از پژوهش Bichteler and Eaton (1980) همسو است. نتایج این پژوهش نیز نشان داد استفاده از این مقیاس نسبت به زوج

پس از سنجه امسلر، بیشترین میزان همبستگی، مربوط به سنجه زوج کتاب‌شناختی است. با توجه به ضریب همبستگی به دست آمده می‌توان گفت که بین دو سنجه زوج کتاب‌شناختی و ربط مقاله‌ها با اطمینان ۹۹ درصد رابطه معناداری وجود دارد؛ همچنین، با توجه به مثبت بودن جهت رابطه می‌توان گفت که با افزایش سنجه زوج کتاب‌شناختی، ربط مقاله‌ها نیز افزایش می‌یابد ($p \leq 0/01$ ، $r = 0/161$).

پس از سنجه‌های امسلر و زوج کتاب‌شناختی، بیشترین همبستگی بین متغیر هیتس (اعتبار) با ربط مقاله‌ها است. با توجه به مقدار ضریب همبستگی می‌توان گفت که با احتمال خطای کمتر از یک صدم درصد بین دو متغیر رابطه مثبت و معناداری وجود دارد. به عبارتی با افزایش هر واحد متغیر اعتبار هیتس، به میزان ۰/۱۵۳ متغیر ربط مقاله‌ها افزایش خواهد داشت ($p \leq 0/01$ ، $r = 0/153$).

در مرتبه چهارم میزان همبستگی ربط مقاله‌ها با سنجه‌های استنادی، متغیر پیج‌رنک قرار دارد که با توجه به مقدار ضریب همبستگی می‌توان گفت که با میزان خطای کمتر از یک صدم درصد بین دو متغیر رابطه مثبت و معناداری وجود دارد. به عبارتی با افزایش هر واحد متغیر پیج‌رنک، به میزان ۰/۱۱۳ متغیر ربط مقاله‌ها افزایش خواهد داشت ($p \leq 0/01$ ، $r = 0/113$).

پرس و جو است؛ بنابراین، در مقایسه با الگوریتم پیچ‌رنک که صرفاً ساختار گراف را در نظر گرفته و مستقل از پرس و جو است، در تعیین ربط مقاله‌ها موفق‌تر بوده است. همچنین نتایج نشان داد در بین سنجه‌های استنادی سنتی، سنجه هم‌استنادی کمترین توان را در تعیین میزان ربط مقاله‌ها داشت. عملکرد ضعیف‌تر سنجه هم‌استنادی نسبت به سایر سنجه‌ها در پژوهش‌های پیشین نیز به‌دست‌آمده است و قابل انتظار بود (Shibata and et al., 2009; Ahlgren and et al., 2020). این مسئله به این دلیل است که مقاله تازه انتشار یافته، هنوز نتوانسته است استناد قابل قبولی دریافت کند و بر این اساس نمی‌تواند با سایر مقالات هم‌استناد شود. البته باید به این نکته هم توجه کرد که مجموعه آزمون سایترک پویا نیست و قابلیت افزایش تعداد مقالات هم‌استناد نیز در آن وجود ندارد.

در این پژوهش گرچه تمامی سنجه‌های استنادمحور با ربط مقاله‌ها براساس سنجه مش، همبستگی در سطح یک صدم داشتند، اما ضرایب همبستگی چندان قوی نبود. در متون بر این نکته تأکید شده است که کاربران - به‌ویژه متخصصان موضوعی و پژوهشگران - به دلایل مختلف نیازمند سامانه‌های بازیابی کارآمدی هستند که تا حد امکان فقط مرتبط‌ترین و مهم‌ترین مدارک را در حداقل زمان برایشان بازیابی نمایند؛ برای مثال، بر اساس پژوهش‌ها کاربران تنها به چند نتیجه نخست توجه می‌کنند و سایر نتایج بازیابی شده را نادیده می‌گیرند (Bar-Ilan, et al. 2006; Lewandowski, 2008; Bar-Ilan, et al., 2009; Lewandowski, 2017). از جمله مهم‌ترین عواملی که بر تحقق این امر تأثیر مستقیم دارد، کیفیت سنجه‌هایی است که ربط بین مقاله‌ها با یکدیگر یا مدارک با پرس و جو را تعیین می‌کند؛ به عبارت دیگر، هرچه سنجه‌های بهتری به کار گرفته شود، می‌توان انتظار داشت

کتاب‌شناختی، منجر به بهبود بازیابی مدارک مرتبط می‌شود. گرچه به نظر می‌رسد از این سنجه نسبت به دو سنجه هم‌استنادی و زوج کتاب‌شناختی، در سامانه‌های اطلاعاتی کمتر استفاده شده است، بنابراین می‌توان این سنجه را بیش از پیش مدنظر قرار داد.

سنجه زوج کتاب‌شناختی، همان‌گونه که در مقدمه نیز بیان شد، علاوه بر این که عمدتاً در متون به‌عنوان یک سنجه کارآمد برای تعیین مدارک مرتبط یاد می‌شود؛ پایگاه‌های اطلاعاتی همانند وب آو ساینس و اسکوپوس نیز از آن برای بازیابی مدارک مرتبط استفاده می‌کنند (Burnham, 2006; Ahlgren and Jarneving, 2008; Nicolaisen and Frandsen, 2012; Char and Ajiferuke, 2013). براین اساس، همانند نتایج پژوهش‌های پیشین (برای مثال، Boyack and Klavans, 2010; Ahlgren, Chen, Colliander, & van Eck, 2020) که نشان داده بودند که عملکرد سنجه زوج کتاب‌شناختی نسبت به هم‌استنادی برای شناسایی شباهت میان مدارک بهتر است، در این پژوهش نیز یافته‌ها تأییدکننده این مهم بود.

در کنار دو سنجه امسلر و زوج کتاب‌شناختی، نتایج نشان داد دو سنجه هیتس (اعتبار) و پیچ‌رنک نیز توانایی مناسبی برای تعیین میزان ربط مقاله‌ها دارند. این دو سنجه که از سنجه‌های جدیدتر استنادمحور هستند، با استفاده از ساختار پیوندی شبکه‌های مقالات، رویت‌پذیری و اعتبار را به‌طور هم‌زمان لحاظ می‌کنند. این یافته، یافته‌های پژوهش‌های پیشین مبنی بر قدرت الگوریتم‌های محیط وب از جمله پیچ‌رنک و هیتس را در بازیابی و رتبه‌بندی مدارک مرتبط تأیید می‌کند (Liu and Lin, 2007; Chen and et al., 2007, Lin, 2008, Yin, Huang, Hu, & Li 2009; Zhuge and Zhang, 2010; Yin and et al., 2011; Liu, Zhang, & Guo, 2012). الگوریتم هیتس، علاوه بر ساختار گراف، از روش محتواکاوای نیز استفاده می‌کند و وابسته به

سنجه‌های شباهت‌یابی مانند سنجه‌های مبتنی بر استناد برای پیشنهاد مقالات مشابه در پایگاه‌های اطلاعاتی و در بافتارهای مختلف شامل موتورهای جست‌وجو، پایگاه‌های اطلاعاتی و استنادی، سامانه‌های پیشنهاددهنده و حتی کتابخانه‌های دیجیتالی کاربردهای بسیاری دارند و می‌توان از پتانسیل این سنجه‌ها، برای پیشنهاد مقالات مشابه و رتبه‌بندی نتایج بازیابی بهره برد.

در نهایت، باید به این امر اشاره کرد که با توجه به خلاءهای پژوهشی موجود در زمینه مقایسه سنجه‌های شباهت‌یابی، در این پژوهش، توانایی تعدادی از این سنجه‌ها برای تعیین میزان ربط مقاله‌ها بررسی شد؛ لازم است نتایج ذکرشده بر روی دیگر انواع منابع اطلاعاتی و دیگر مجموعه دادگان به دست آید تا بتوان با قطعیت بیشتری در مورد نتایج به‌دست‌آمده قضاوت کرد. الگوریتم‌ها و مدل‌های پیش‌بینی حاصل از ترکیب این سنجه‌ها در کنار یکدیگر می‌بایست در پژوهش‌های آتی مدنظر قرار گیرد.

نتایج بازیابی بهتر شود. ضرایب همبستگی پایین سنجه‌های استنادمحور، می‌تواند نشان‌دهنده این باشد که استفاده از این سنجه‌ها به‌تنهایی برای یافتن مدارک مرتبط کفایت نمی‌کند و بایستی توانایی سایر سنجه‌ها مانند سنجه‌های متن‌محور و حتی ترکیب سنجه‌های متن‌محور و استنادمحور با یکدیگر مدنظر قرار گیرد. نتایج پژوهش‌های مختلف نیز نشان می‌دهد مدارک بازیابی‌شده با استفاده از سنجه‌های متن‌محور، با مدارکی که از طریق سنجه‌های استنادمحور بازیابی می‌شوند، متفاوت هستند (Mccain, 1989; Pao and Worthen, 1989,) (Ahlgren and Jarneving, 2008). توصیه به استفاده ترکیبی از دو گروه سنجه‌های متن‌محور و استنادمحور در تعدادی از پژوهش‌های پیشین تاکید شده است (Yin and et al., 2009; Boyack and Klavans, 2010; Reyhani Hamedani and et al., 2013; Reyhani Hamedani, Kim, & Kim, 2016; Jiang and et al., 2019; Ahlgren and et al., 2020; Boyack and Klavans, 2020).

References

مآخذ

- Agrahri, A. K., Manickam, D. A. T., & Riedl, J. (2008, October). Can people collaborate to improve the relevance of search results?. *In Proceedings of the 2008 ACM conference on Recommender systems* (pp. 283-286). ACM. <https://doi.org/10.1145/1454008.1454052>
- Ahlgren, P., Chen, Y., Colliander, C., & van Eck, N. J. (2020). Enhancing direct citations: A comparison of relatedness measures for community detection in a large set of PubMed publications. *Quantitative Science Studies*, 1(2), 714-729. https://doi.org/10.1162/qss_a_00027
- Ahlgren, P., & Jarneving, B. (2008). Bibliographic coupling, common abstract stems and clustering: A comparison of two document-document similarity approaches in the context of science mapping. *Scientometrics*, 76(2), 273-290. <https://doi.org/10.1007/s11192-007-1935-1>
- Amer, A. A., & Abdalla, H. I. (2020). A set theory based similarity measure for text clustering and classification. *Journal of Big Data*, 7(1), 1-43. <https://doi.org/10.1186/s40537-020-00344-3>

- Bar-Ilan, J., Keenoy, K., Levene, M., & Yaari, E. (2009). Presentation bias is significant in determining user preference for search results—A user study. *Journal of the American Society for Information Science and Technology*, 60(1), 135-149. <https://doi.org/10.1002/asi.20941>
- Bar-Ilan, J., Keenoy, K., Yaari, E., & Levene, M. (2007). User rankings of search engine results. *Journal of the American Society for Information Science and Technology*, 58(9), 1254-1266. <https://doi.org/10.1002/asi.20608>
- Bar-Ilan, J., Levene, M., & Mat-Hassan, M. (2006). Methods for evaluating dynamic changes in search engine rankings: a case study. *Journal of Documentation*, 62(6), 708-729. <https://doi.org/10.1108/00220410610714930>
- Bichteler, J., & Eaton, E. A. (1980). The combined use of bibliographic coupling and cocitation for document retrieval. *Journal of the American Society for Information Science*, 31(4), 278-282. <https://doi.org/10.1002/asi.4630310408>
- Boyack, K. W., & Klavans, R. (2010). Co-citation analysis, bibliographic coupling, and direct citation: Which citation approach represents the research front most accurately? *Journal of the American Society for Information Science and Technology*, 61(12), 2389-2404. <https://doi.org/10.1002/asi.21419>
- Boyack, K. W., & Klavans, R. (2020). A comparison of large-scale science models based on textual, direct citation and hybrid relatedness. *Quantitative Science Studies*, 1-16. https://doi.org/10.1162/qss_a_00085
- Burnham, J. F. (2006). Scopus database: A review. *Biomedical Digital Libraries*, 3(1), 1-8. <https://doi.org/10.1186/1742-5581-3-1>
- Char, D. C., & Ajiferuke, I. (2013, October). Comparison of the effectiveness of related functions in Web of Science and Scopus. In *Proceedings of the Annual Conference of CAIS/Actes du congrès annuel de l'ACSI*. <https://doi.org/10.29173/cais353>
- Chen, P., Xie, H., Maslov, S., & Redner, S. (2007). Finding scientific gems with Google's PageRank algorithm. *Journal of Informetrics*, 1(1), 8-15. <https://doi.org/10.1016/j.joi.2006.06.001>
- Colavizza, G., Boyack, K. W., Van Eck, N. J., & Waltman, L. (2018). The closer the better: Similarity of publication pairs at different cocitation levels. *Journal of the Association for Information Science and Technology*, 69(4), 600-609. <https://doi.org/10.1002/asi.23981>
- Devi, P., Gupta, A., & Dixit, A. (2014). Comparative study of hits and pagerank link based ranking algorithms. *International Journal of Advanced Research in Computer and Communication Engineering*, 3(2), 5749-5754.
- Eminagaoglu, M. (2020). A new similarity measure for vector space models in text classification and information retrieval. *Journal of Information Science*. <https://doi.org/10.1177/0165551520968055>
- Eto, M. (2013). Evaluations of context-based co-citation searching. *Scientometrics*, 94(2), 651-673. <https://doi.org/10.1007/s11192-012-0756-z>
- Eto, M. (2019). Extended co-citation search: Graph-based document retrieval on a co-citation network containing citation context information. *Information Processing & Management*, 56(6), 102046. <https://doi.org/10.1016/j.ipm.2019.05.007>
- Gipp, B., & Beel, J. (2009). Citation proximity analysis (CPA): a new approach for identifying related work based on co-citation analysis. In *ISSI'09: 12th International Conference on Scientometrics and Informetrics* (pp. 571-575). Retrieved June 20, 2020 from [https://isg.beel.org/pubs/Citation%20Proximity%20Analysis%20\(CPA\)%20-%20A%20new%20approach%20for%20identifying%20related%20work%20based%20on%20Co-Citation%20Analysis%20--%20preprint.pdf](https://isg.beel.org/pubs/Citation%20Proximity%20Analysis%20(CPA)%20-%20A%20new%20approach%20for%20identifying%20related%20work%20based%20on%20Co-Citation%20Analysis%20--%20preprint.pdf)
- Gipp, B., Meuschke, N., & Lipinski, M. (2015). CITREC: An Evaluation Framework for Citation-Based Similarity Measures based on TREC Genomics and PubMed Central. In *iConference 2015*, Newport Beach, California. <https://doi.org/10.5281/zenodo.3547372>

- Goswami, P., Mantri, M., & Bhattacharya, M. (2017). Web Page Ranking Based On User Query. *International Journal of HIT Transaction on ECCN*, 3(2A), 21-33.
 - Hajian, A., & CheshmehSohrabi, M. (2020). Ranking and Relevance in Noormags and RICEST Databases. *National Studies on Librarianship and Information Organization*, 31(3), 72-92. DOI: 10.30484/nastinfo.2020.2472.1934 [In Persian]
- [حاجیان، آزاده؛ چشمه‌سهرابی، مظفر (۱۳۹۹). رتبه‌بندی و ربط مقالات در پایگاه‌های اطلاعاتی نورمگز و رایست. *مطالعات ملی کتابداری و سازماندهی اطلاعات*، ۳۱(۳)، ۷۲-۹۲.]
- Hariri, N. (2011). Relevance ranking on Google: Are top ranked results really considered more relevant by the users?. *Online Information Review*, 35(4), 598-610. <https://doi.org/10.1108/14684521111161954>
 - Herskovic, J. R., & Bernstam, E. V. (2005). Using incomplete citation data for MEDLINE results ranking. In *AMIA Annual Symposium proceedings* (Vol. 2005, p. 316-320). American Medical Informatics Association.
 - Janssens, A. C. J., Gwinn, M., Brockman, J. E., Powell, K., & Goodman, M. (2020). Novel citation-based search method for scientific literature: a validation study. *BMC medical research methodology*, 20(1), 1-11. <https://doi.org/10.1186/s12874-020-0907-5>
 - Jiang, X., Sun, X., Yang, Z., Zhuge, H., & Yao, J. (2016). Exploiting Heterogeneous Scientific Literature Networks to Combat Ranking Bias: Evidence From the Computational Linguistics Area. *Journal of the American Society for Information Science and Technology*, 67(7), 1679-1702. <https://doi.org/10.1002/asi.23463>
 - Jiang, T., Zhang, Z., Zhao, T., Qin, B., Liu, T., Chawla, N. V., & Jiang, M. (2019, November). CTGA: Graph-based Biomedical Literature Search. In *2019 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)* (pp. 395-400). IEEE. DOI: [10.1109/BIBM47256.2019.8983173](https://doi.org/10.1109/BIBM47256.2019.8983173)
 - Jouili S., Tabbone S., Valveny E. (2010) Comparing Graph Similarity Measures for Graphical Recognition. In: Ogier JM., Liu W., Lladós J. (eds) *Graphics Recognition. Achievements, Challenges, and Evolution. GREC 2009. Lecture Notes in Computer Science*, vol 6020. Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-642-13728-0_4
 - Kirsch, S. M., Gnasa, M., Won, M., & Cremers, A. (2008). From PageRank to Social Rank: Authority-Based Retrieval in Social Information Spaces. In *Social Information Retrieval Systems: Emerging Technologies and Applications for Searching the Web Effectively* (pp. 134-154). IGI Global.
 - Kessler, M.M. (1963). Bibliographic coupling between scientific papers. *American Documentation*, 14(1), 10-25. <https://doi.org/10.1002/asi.5090140103>
 - Kleinberg, J. M. (1999). Authoritative sources in a hyperlinked environment. *Journal of the ACM (JACM)*, 46(5), 604-632.
 - Lewandowski, D. (2008). The retrieval effectiveness of web search engines: considering results descriptions. *Journal of documentation*, 64(6), 915- 937. Retrieved June 20, 2020 from <https://arxiv.org/ftp/arxiv/papers/1511/1511.05800.pdf>.
 - Lewandowski D. (2017) Is Google Responsible for Providing Fair and Unbiased Results? In: Taddeo M., Floridi L. (eds) *The Responsibilities of Online Service Providers. Law, Governance and Technology Series*, vol 31, pp.61-77 Springer, Cham. https://doi.org/10.1007/978-3-319-47852-4_4
 - Lin, J. (2008). PageRank without hyperlinks: Reranking with PubMed related article networks for biomedical text retrieval. *BMC Bioinformatics*, 9, 270. <https://doi.org/10.1186/1471-2105-9-270>
 - Lin, Y. S., Jiang, J. Y., & Lee, S. J. (2013). A similarity measure for text classification and clustering. *IEEE transactions on knowledge and data engineering*, 26(7), 1575-1590. doi: 10.1109/TKDE.2013.19.

- Liu, Y., & Lin, Y. (2007, October). Supervised HITS algorithm for MEDLINE citation ranking. In *Bioinformatics and Bioengineering (BIBE)*, 2007. *Proceedings of the 7th IEEE International Conference on* (pp. 1323-1327). IEEE.
- Liu, X., Zhang, J., & Guo, C. (2012, October). Full-text citation analysis: enhancing bibliometric and scientific publication ranking. In *Proceedings of the 21st ACM international conference on Information and knowledge management* (pp. 1975-1979). ACM. <https://doi.org/10.1145/2396761.2398555>
- Lu, W., Janssen, J., Milios, E., Japkowicz, N., & Zhang, Y. (2007). Node similarity in the citation graph. *Knowledge and Information Systems*, 11(1), 105-129. <https://doi.org/10.1007/s10115-006-0023-9>
- Ma, N., Guan, J., & Zhao, Y. (2008). Bringing PageRank to the citation analysis. *Information Processing and Management*, 44:800-810. <https://doi.org/10.1016/j.ipm.2007.06.006>
- McCain, W. (1989). Descriptor and Citation Retrieval in the Medical Behavioral Sciences Literature: Retrieval Overlaps and Novelty Distribution. *Journal of the American Society for Information Science*, 40(2), 110-114. [https://doi.org/10.1002/\(SICI\)1097-4571\(198903\)40:2<110::AID-ASI5>3.0.CO;2-T](https://doi.org/10.1002/(SICI)1097-4571(198903)40:2<110::AID-ASI5>3.0.CO;2-T)
- McGill, M., Koll, M., & Noreault, T. (1979). An evaluation of factors affecting document ranking by information retrieval systems, Syracuse, NY: School of Information Studies, Syracuse University.
- Nicolaisen, J., & Frandsen, T. F. (2012). Consensus formation in science modeled by aggregated bibliographic coupling. *Journal of Informetrics*, 6(2), 276-284. DOI: [10.1016/j.joi.2011.08.001](https://doi.org/10.1016/j.joi.2011.08.001)
- Neshat, N. (2003). Hermeneutic and Information Retrieval. *Informology*, 2, 31-46. [http://ensani.ir/file/download/article/20110108110840-0%20\(1\)](http://ensani.ir/file/download/article/20110108110840-0%20(1))

[نشاط، نرگس (۱۳۸۲). هرمنوتیک و بازیابی اطلاعات. اطلاع‌شناسی، ۲، ۳۱-۴۶.]

- Pao, M. L., & Worthen, B. (1989). Retrieval Effectiveness by Semantic and Citation Searching. *Journal of the American Society for Information Science*, 40(4), 226-235. [https://doi.org/10.1002/\(SICI\)1097-4571\(198907\)40:4<226::AID-ASI2>3.0.CO;2-6](https://doi.org/10.1002/(SICI)1097-4571(198907)40:4<226::AID-ASI2>3.0.CO;2-6)
- Reyhani Hamedani, M., & Kim, S. W. (2021). On Investigating Both Effectiveness and Efficiency of Embedding Methods in Task of Similarity Computation of Nodes in Graphs. *Applied Sciences*, 11(1), 162. <https://doi.org/10.3390/app11010162>
- Reyhani Hamedani, M., Kim, S. W., & Kim, D. J. (2016). SimCC: A novel method to consider both content and citations for computing similarity of scientific papers. *Information Sciences*, 334, 273-292. <https://doi.org/10.1016/j.ins.2015.12.001>
- Reyhani Hamedani, M., Lee, S. C., & Kim, S. W. (2013, October). On combining text-based and link-based similarity measures for scientific papers. In *Proceedings of the 2013 Research in Adaptive and Convergent Systems* (pp. 111-115). ACM. DOI: [10.1145/2513228.2513321](https://doi.org/10.1145/2513228.2513321)
- Sadein Khorram, S., & Abbaspour, J. (2019). Article Ranking by Recommender Systems vs. Users' Perspectives. *National Studies on Librarianship and Information Organization*, 30(3), 46-57. DOI: 10.30484/nastinfo.2019.2187.1838 http://nastinfo.nlai.ir/article_2346.html

[سعدین خرم، صبا؛ عباس پور، جواد (۱۳۹۸). سنجش رتبه‌بندی سامانه‌های پیشنهاددهنده مقاله در تقابل با رتبه‌بندی

کاربران. مطالعات ملی کتابداری و سازماندهی اطلاعات، ۳۰ (۳)، ۴۶-۵۷.]

- Shibata, N., Kajikawa, Y., Takeda, Y., & Matsushima, K. (2009). Comparative study on methods of detecting research fronts using different types of citation. *Journal of the American Society for Information Science and Technology*, 60(3), 571-580. <https://doi.org/10.1002/asi.20994>
- Small, H. (1973). Co-citation in the scientific literature: A new measure of the relationship between two documents. *Journal of the American Society for Information Science*, 24(4), 265-269. <https://doi.org/10.1002/asi.4630240406>
- Smith, L. C. (1981). Citation Analysis. *Library Trends*, 30(1), 83-106.

- Su, C., Pan, Y., Zhen, Y., Ma, Z., Yuan, J., Guo, H., ... & Wu, Y. (2011). PrestigeRank: A new evaluation method for papers and journals. *Journal of Informetrics*, 5(1), 1-13. <https://doi.org/10.1016/j.joi.2010.03.011>
- Thelwall, M. (2003). Can Google's PageRank be used to find the most important academic Web pages?. *Journal of Documentation*, 59(2), 205-217. <https://doi.org/10.1108/00220410310463491>
- Thompson, V. U., Panchev, C., & Oakes, M. (2015, November). Performance evaluation of similarity measures on similar and dissimilar text retrieval. In *2015 7th International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management (IC3K)* (Vol. 1, pp. 577-584). IEEE.
- Torres, G. J., Basnet, R. B., Sung, A. H., Mukkamala, S., & Ribeiro, B. M. (2009). A similarity measure for clustering and its applications. *Int J Electr Comput Syst Eng*, 3(3), 164-170.
- Wanjantuk, P., & Keane, J. A. (2004, October). Finding related documents via communities in the citation graph. In *Communications and Information Technology, 2004. ISCIT 2004. IEEE International Symposium on* (Vol. 1, pp. 445-450). IEEE. doi: [10.1109/ISCIT.2004.1412885](https://doi.org/10.1109/ISCIT.2004.1412885).
- Yin, X., Huang, X., Hu, Q., & Li, Z. (2009, April). Boosting biomedical information retrieval performance through citation graph: An empirical study. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining* (pp. 949-956). Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-642-01307-2_100
- Yin, X., Huang, J. X., & Li, Z. (2011). Mining and modeling linkage information from citation context for improving biomedical literature retrieval. *Information processing & management*, 47(1), 53-67. <https://doi.org/10.1016/j.ipm.2010.03.010>
- Yoon, S. H., Kim, S. W., & Park, S. (2016). C-Rank: A link-based similarity measure for scientific literature databases. *Information Sciences*, 326, 25-40. <https://doi.org/10.1016/j.ins.2015.07.036>
- Zhang, J. & Korfhage, R. (1999). A Distance and Angle Similarity Measure Method. *Journal of the American Society for Information Science*, 50(9), 772-778. [https://doi.org/10.1002/\(SICI\)1097-4571\(1999\)50:9<772::AID-ASI5>3.0.CO;2-E](https://doi.org/10.1002/(SICI)1097-4571(1999)50:9<772::AID-ASI5>3.0.CO;2-E)
- Zhuge, H., & Zhang, J. (2010). Topological centrality and its e-Science applications. *Journal of the Association for Information Science and Technology*, 61(9), 1824-1841.
- <https://doi.org/10.1002/asi.21353>

پیوست ۱

شرایط انتخاب مقالات:

- ۱- مدرک انتخاب شده دارای شباهت امسلر با مدارک دیگر باشد (مجموعه A)
- ۲- مدرک انتخاب شده دارای شباهت زوج کتابشناختی با مدارک دیگر باشد (مجموعه B)
- ۳- مدرک انتخاب شده دارای شباهت هم‌استنادی با مدارک دیگر باشد (مجموعه C)
- ۴- مدرک انتخاب شده دارای شباهت مش با مدارک دیگر باشد (مجموعه M)
- ۵- بین مجموعه‌های A، B، C و M اشتراک وجود داشته باشد.

پرسشی که با توجه به شرایط بالا برای انتخاب مدارک استفاده شد به شرح زیر است:

```
select distinct sim_amsler.document1,sim_amsler.document2
from
document, sim_amsler, sim_bibco, sim_cocit,
sim_mesh_intersections_27000
where document.pmid = sim_amsler.document1
and ((sim_amsler.document1 = sim_bibco.document1 and
sim_amsler.document2 = sim_bibco.document2) )
and ((sim_amsler.document1 = sim_cocit.document1 and sim_amsler.document2
= sim_cocit.document2) )
and (sim_amsler.document1 = sim_mesh_intersections_27000.document1 )
and sim_bibco.value>0 and sim_cocit.value>0
and document.id=@Random
```

هر دیتاست شامل ۳۰ مدرک بود و به ازای هر مدرک، بر اساس سنجه‌های شباهت یادشده، تعدادی مدرک دیگر وجود داشت که دارای میزانی از شباهت بودند. برای ایجاد شبکه استنادی برای استفاده در الگوریتم‌های پیچ رنگ و هیتس، از جدول **reference** که حاوی مدارک استنادکننده و استنادشونده است استفاده شد و از بین آن‌ها مدارکی انتخاب شد که هم خود مدرک و هم استنادهای آن دارای سنجه شباهت مش باشند. شبکه استنادی به دست آمده گرافی با بیش از دو میلیون و سیصد هزار یال (ارتباط استنادی) بود. پرسشی که این شبکه استنادی را ایجاد کرد به شرح زیر است:

```
select distinct docpmid, refpmid from
(select distinct document as docpmid,refpmid from reference where refpmid is not NULL
and refpmid in
(select distinct doc1 from (
select distinct document2 as doc1 from citrec2.sim_mesh_intersections_27000
union
select distinct document1 as doc1 from citrec2.sim_mesh_intersections_27000
)t)
union
select distinct document as docpmid,refpmid from reference where refpmid is not NULL and
document in
(select distinct doc1 from (
select distinct document2 as doc1 from citrec2.sim_mesh_intersections_27000
union
select distinct document1 as doc1 from citrec2.sim_mesh_intersections_27000
)t) )tt
```