

Feasibility Study of Ontological Development Using Semi-Automatic Method based on Lexical Frequency Analysis: A Case Study of "Glaucoma"

Somayeh Tamjid

PhD Candidate in Knowledge and Information Science;
Department of Communication and Knowledge Sciences;
Islamic Azad University; Science and Research Branch;
Tehran, Iran Email: s.tamjid@srbiau.ac.ir; tamjid.s@iums.ac.ir

Fatemeh Nooshinfard*

PhD in Knowledge and Information Science; Associate Professor;
Department of Communication and Knowledge Sciences;
Islamic Azad University; Science and Research Branch;
Tehran, Iran Email: nooshinfard@srbiau.ac.ir; f.nooshinfard@gmail.com

Molouk Sadat Hosseini Beheshti

PhD in Linguistics; Associate Professor; Iranian Research Institute
for Information Science and Technology (IranDoc); Tehran, Iran;
Email: Beheshti@irandoc.ac.ir

Nadjla Hariri

PhD in Knowledge and Information Science; Professor;
Department of Communication and Information Science;
Islamic Azad University; Science and Research Branch;
Tehran, Iran Email: N_hariri@srbiau.ac.ir

Fahimeh Babalhavaeji

PhD in Knowledge and Information Science; Associate Professor;
Department of Communication and Knowledge Sciences Islamic
Azad University; Science and Research Branch; Tehran, Iran;
Email: f.babalhavaeji@gmail.com

**Iranian Journal of
Information
Processing and
Management**

Received: 13, Mar. 2022 Accepted: 21, Aug. 2022

Abstract: Following recent trends in information management systems, conventional word-based information retrieval methods are changing to concept-based approaches by means of the broad application of ontologies. More specifically, the use of ontologies for knowledge management is significant in the medical sciences and human disease domains due to the diversity and necessity of information sharing between numerous data repositories such as medical records, health record

Iranian Research Institute
for Information Science and Technology
(IranDoc)

ISSN 2251-8223

eISSN 2251-8231

Indexed by SCOPUS, ISC, & LISTA

Vol. 39 | No. 1 | pp. 131-156

Autumn 2023

<https://doi.org/jipm.39.1>

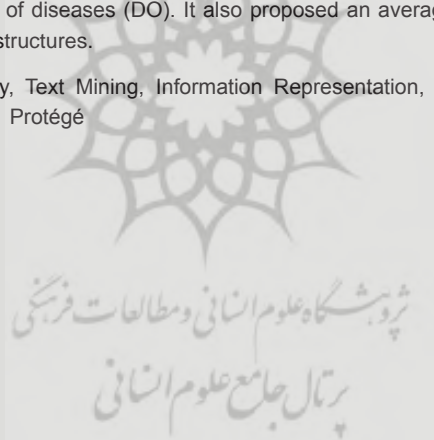


* Corresponding Author

systems, and so on. Furthermore, ontologies make natural language processing approaches more feasible by reducing semantic ambiguity and making concepts comprehensible to computer-based deductions. In this research, a semi-automated approach for ontology development is proposed, which assists in identifying structural components of an ontology and determining possible relations between them based on scientific text records. The proposed approach, in a general view, includes the gathering of a large volume of technical data in text format, processing, and extraction of results with a minimal contribution of human-based supervision. The processing stage is coded in Matlab code named TmbOnt_Alfa and applies two main techniques including word frequency and Lexico-Synactic patterns analysis, to identify concepts and relations, respectively. The role of the human supervisor is narrowed to entering target terms, eliminating unnecessary outputs, and finalizing the ontology structure. In order to evaluate the efficiency of the proposed method, a case study for ontological development in the field of glaucoma has been conducted, and results are compared with medical subject headings of MESH descriptors, the Persian medical thesaurus, ontology of diseases, and Bioassay ontology (BAO).

According to results, the developed ontology, when compared by Glaucoma entry, covered 80% of the medical titles in Mesh, 100% of the medical terms developed in the Persian Medical Thesaurus, and 100% of the Persian medical descriptors. Moreover, the resultant ontology structure is compatible with more than 90% of the same ontology represented in Bioassay and 57% of the ontology of diseases (DO). It also proposed an average of 30% more terms for existing ontological structures.

Keywords: Ontology, Text Mining, Information Representation, Glaucoma, Eye Disease, Medical Thesaurus, Protégé



امکان‌سنجی توسعه هستی‌شناسی به روش نیمه‌خودکار مبتنی بر تحلیل بسامد واژگان: مطالعه موردی بیماری «گلوکوم»

سمیه تمجید

دانشجوی دکتری علم اطلاعات و دانش‌شناسی؛
دانشگاه آزاد اسلامی؛ واحد علوم و تحقیقات؛ تهران، ایران؛
s.tamjid@srbiau.ac.ir

فاطمه نوشین فرد

دکتری علم اطلاعات و دانش‌شناسی؛ دانشیار؛
دانشگاه آزاد اسلامی؛ واحد علوم و تحقیقات؛ تهران، ایران؛
nooshinfard@sriau.ac.ir

ملوک‌السادات حسینی بهشتی

دکتری زبان‌شناسی؛ دانشیار؛
پژوهشگاه علوم و فناوری اطلاعات ایران (ایرانداک)؛
تهران، ایران beheshti@irandoc.ac.ir

نجلا حریری

دکتری علم اطلاعات و دانش‌شناسی؛ استادی؛
دانشگاه آزاد اسلامی؛ واحد علوم و تحقیقات؛ تهران، ایران؛
N_hariri@srbiau.ac.ir

فهیمه باب‌الحوایجی

دکتری علم اطلاعات و دانش‌شناسی؛ دانشیار؛
دانشگاه آزاد اسلامی؛ واحد علوم و تحقیقات؛ تهران، ایران؛
f.babalhavaeji@gmail.com



دریافت: ۱۴۰۰/۱۲/۲۲ پذیرش: ۱۴۰۱/۰۵/۳۰ مقاله برای اصلاح به مدت ۵۳ روز نزد پدیدآوران بوده است.

چکیده: تغییر رویکرد نظام‌های اطلاعاتی از پردازش واژه به پردازش مفهوم، موجب توجه به هستی‌شناسی‌ها شده است. در علوم پزشکی و بیماری‌های انسان به لحاظ وجود تنوع در اصطلاحات و لزوم اشتراک اطلاعات از طریق نرم‌افزارهای مختلف مانند پرونده‌های پزشکی، سامانه‌های ثبت سوابق بهداشتی و ... به کارگیری هستی‌شناسی‌ها ضروری به نظر می‌رسد. در پژوهش حاضر، رویکردی نیمه‌خودکار برای توسعه هستی‌شناسی پیشنهاد شده است که می‌تواند با استفاده از ابزارهای متن‌کاوی، شناسایی مؤلفه‌های ساختاری هستی‌شناسی و تعیین نسبی روابط را از متون علمی تسهیل کند. مدل پیشنهادی در قالب کد نرم‌افزاری با نام اختصاری TmbOnt_Alfa ارائه شده است. این کد با استفاده از رابط کاربری،

نشریه علمی | رتبه بین‌المللی
پژوهشگاه علوم و فناوری اطلاعات ایران
(ایرانداک)

شاپا (چاپی) ۸۲۲۳-۲۲۵۱

شاپا (الکترونیکی) ۸۲۳۱-۲۲۵۱

نمایه در SCOPUS، LISTA، ISC، و

ijpm.irandoc.ac.ir

دوره ۳۹ | شماره ۱ | صص ۱۳۱-۱۵۶

پاییز ۱۴۰۲

<https://doi.org/ijpm.39.1>



اطلاعات در یک دامنه دارند» تعریف شده است (Noy and McGuinness 2001). به‌طور خلاصه، با حرکت از تاکسونومی‌ها و اصطلاحنامه‌ها به سوی الگوی مفهومی و نظریه منطقی در برقراری ارتباط بین مفاهیم، به هستی‌شناسی‌ها می‌رسیم (حسینی بهشتی ۱۳۹۲). این تأثیر در حوزه سلامت، با توجه به تعدد و پیچیدگی اصطلاحات تخصصی و لزوم دستیابی به تعاملات بین سیستمی در بازبایی دانش می‌تواند بسیار کارآمد باشد. به همین دلیل، شاهد رشد روزافزون هستی‌شناسی‌ها در حوزه زیست‌پزشکی و سلامت به نسبت سایر حوزه‌ها هستیم (Noy et al. 2009). توسعه هستی‌شناسی در این حوزه کمک می‌کند که اطلاعات گسسته به‌دست آمده از بطن فعالیت‌های تجربی و نظری قابلیت بازبایی بهتری یافته و به دانش قابل ادراک توسط سامانه‌های نرم‌افزاری تبدیل شود. در مدت کوتاه مواجهه با بیماری «کووید ۱۹»، محققان با ساخت چندین هستی‌شناسی در سراسر دنیا تلاش کردند تا با ایجاد مفاهیم مشترک بتوانند داده‌های بیماری را به اشتراک بگذارند. هستی‌شناسی‌های «CODO»، «CIDO»، «VIDO» و «IDO-COVID-19»، همگی طی سال‌های ۲۰۲۰ و ۲۰۲۱ ایجاد شده‌اند (Babcock et al. 2021 و Dutta and DeBellis 2020). از منظر اقتصادی، اهمیت هستی‌شناسی از میزان بودجه تخصیص یافته به آن در جوامعی که هستی‌شناسی‌ها را توسعه می‌دهند، مشهود است. اختصاص بودجه ۴/۴ میلیون دلاری در سال ۲۰۰۹ برای «هستی‌شناسی ژن»، تخصیص بودجه ۱۸/۸ میلیون دلاری برای پنج سال کار در حوزه «هستی‌شناسی زیست‌پزشکی» و نیز تخصیص ۳۲/۴ میلیون دلار در سال ۲۰۰۳ و پس از آن ۶ میلیون دلار بودجه سالانه توسط «کتابخانه ملی پزشکی» برای پروژه «SNOMED-CT» می‌تواند منعکس‌کننده بخشی از حجم سرمایه‌گذاری در این حوزه باشد (Liu, Hogan and Crowley 2011). به‌طور معمول، توسعه مداوم و نگاه‌داشت مؤثر هستی‌شناسی‌ها توسط عامل انسانی انجام می‌شود که زمان‌بر بوده و مستعد اشتباهات متعددی است. امروزه، تکنیک‌های نیمه‌خودکار بازبایی دانش در حال تبدیل شدن به ابزاری برای کمک به عوامل انسانی است. اتخاذ چنین رویکردی می‌تواند فعالیت‌های تکراری و خسته‌کننده را به اقدامات نظارتی تغییر داده و منجر به کاهش زمان و بهبود کیفیت شود.

مزیت‌های حاصل از به‌کارگیری هستی‌شناسی‌ها در نرم‌افزارهای جست‌وجو و بازبایی اطلاعات، توسعه‌دهندگان هستی‌شناسی در حوزه پزشکی را بر آن داشت تا از این مفهوم برای تسهیل ادغام داده‌های زیست‌پزشکی مرتبط با بیماری‌های انسان بهره

بیرند. برای مثال، از سال ۲۰۱۲، پایگاه داده پروژۀ «هستی‌شناسی بیماری‌ها»^۱، ۱۹۲ بار به‌روزرسانی شده و طی این فرایند، حدود ۷۶۰ اصطلاح جدید، که برخی از آن‌ها در تعاریف به‌کار رفته‌اند، بدان اضافه شده است (Kibbe et al. 2015). یکی از مشکلات توسعه هستی‌شناسی از این واقعیت ناشی می‌شود که ایجاد هستی‌شناسی بر پایه اصطلاحنامه، هر چند در شناسایی اصطلاحات مرجح به‌عنوان ریشه مفید است، اما بسیار هزینه‌بر بوده و به‌طور کامل به عامل انسانی متخصص وابسته است و امکان اشتباه در مراحل مختلف دور از انتظار نیست (Kless et al. 2016). برای رفع این مشکل در این پژوهش از طراحی هستی‌شناسی به‌صورت نیمه‌خودکار و با روش متن‌کاوی و پردازش زبان طبیعی استفاده شده و مفاهیم و اصطلاحات حاصل با هستی‌شناسی‌ها و اصطلاحنامه‌های پزشکی معتبر مقایسه شده است. مطالعه موردی برای امکان‌سنجی توسعه هستی‌شناسی روی موضوع گلوکوم^۲ انجام شده است. گلوکوم بر اساس مطالعات، دومین عامل نابینایی افراد بالای ۵۰ سال اعلام شده و مهم‌ترین دلیل نابینایی قابل کنترل در افراد بالغ است (معصومی و همکاران ۱۳۹۱). جست‌وجو در پایگاه‌های اطلاعاتی نشان‌دهندۀ روند صعودی تولیدات علمی در مورد این بیماری است (شکل ۲). این در حالی است که سرعنوان‌های پزشکی صرفاً پنج زیرشاخه برای آن در نظر گرفته، ولی این پنج مورد به‌طور کامل نمی‌تواند ساختار دانش این حوزه را دربرگیرد (شکل ۵). این پژوهش در صدد توسعه هستی‌شناسی جهت بازنمون دامنه دانش در حیطه گلوکوم، از طریق متن‌کاوی جدیدترین مقالات منتشرشده در پایگاه‌های اطلاعاتی است. رجوع به متن مندرج در مقالات تازه انتشاریافته می‌تواند مفاهیم جدید مورد استفاده در حوزه دانش بر اساس آخرین مطالعات در حال انجام در مرز دانش این حوزه را شناسایی کند که در توسعه و به‌روزرسانی هستی‌شناسی بسیار مؤثر خواهد بود.

مقاله حاضر در شش بخش اصلی سازماندهی شده است. در بخش دوم، پیشینه پژوهش بررسی شده است. بخش سوم، به روش و رویکردهای خودکار موجود برای توسعه هستی‌شناسی پرداخته است. در بخش چهارم، جزئیات مطالعه موردی انجام‌شده در حوزه گلوکوم ارائه شده و در بخش پنجم، نتایج به‌دست آمده بحث شده است. سرانجام، در بخش ششم، نتیجه‌گیری و پیشنهادات برای پژوهش‌های آتی ارائه شده است.

1. disease ontology (DO)

2. glaucoma

۲. پیشینه پژوهش

توسعه و غنی‌سازی هستی‌شناسی از بطن متون انبوه، رویکردی منطقی است که طی سال‌های اخیر توسط بسیاری از محققان مورد اقبال قرار گرفته است. «میسیکوف، ولاردی و فابریانی» با استفاده از تکنیک متن‌کاوی و تحلیل تشابه متن، دستیاری برای توسعه هستی‌شناسی را در جهت کمک به متخصصان هستی‌شناسی ارائه کردند (Missikoff, Velardi and Fabriani 2003). انگیزه اصلی این تحقیق، کاهش زمان، کاهش هزینه و همسان‌سازی تفاسیر مبتنی بر رویکردهای انسانی بود. «جیانگ و تان» روشی را برای استخراج مفاهیم و روابط معنایی از مجموعه متون دامنه ارائه دادند. آن‌ها یک سیستم جدید به نام CRCTOL برای استخراج دانش معنایی اسناد متنی در دامنه موضوعی خاص طراحی کردند که به جای استخراج مفهوم و روابط به صورت سطحی، دانش معنایی غنی در قالب هستی‌شناسی را استخراج می‌کرد (Jiang and Tan 2005). «شارلت، بچیمونت و ژالت» با استفاده از ابزار پردازش زبان طبیعی به نام «SYNTAX»، روشی را توسعه داده‌اند که هستی‌شناسی پزشکی را از درون گزارش‌های متنی استخراج می‌کند (Charlet, Bachimont and Jaulent 2006). «کیان، لان و لیجان» با استفاده از پیکره متنی و روش‌های پردازش زبان طبیعی خودکار، روشی برای ساخت هستی‌شناسی کشاورزی پویا بر اساس متن کاوی متون نمایه‌شده با اصطلاحنامه AGROVOC پیشنهاد کردند (Qian, Lan and Lijun 2007). «لیو، هوگان و کرولی» روش‌های پردازش زبان طبیعی^۱ را مورد بررسی قرار داده و بر ضرورت این روش جهت پردازش حوزه‌های گسترده اطلاعات امروزی تأکید کرده‌اند. نتیجه بررسی آن‌ها برتری رویکردهای نیمه‌خودکار موجود نسبت به رویکردهای کاملاً خودکار را نشان داد (Liu, Hogan and Crowley 2011). «واچر و شرودر» یک مولد نیمه‌خودکار هستی‌شناسی را با هدف استخراج اصطلاحات، تعاریف و روابط والد-فرزند از وب، «پابمد»^۲ و یا فایل PDF ایجاد کردند. این طرح در قالب افزونه‌ای برای «پروتش»^۳ نهایی شد. عملکرد آن به این صورت است که عبارات اسمی را با توجه به طول و دفعات تکرار می‌یابد و شواهدی در خصوص شناسایی آن به‌عنوان یک اصطلاح، تعریف و یا رابطه معتبر ارائه می‌کند (Wächter & Schroeder 2010). «سیه» و همکاران طرح توسعه هستی‌شناسی نیمه‌خودکار از دانشنامه‌های یک حوزه موضوعی معین ارائه

1. natural language processing (NLP)

2. Pubmed

3. Protégé

کردند. رویکرد آن‌ها محدود به فضای دانشنامه چاپی بود و با شروع از واژه‌نامه و سپس افزودن زیرکلاس‌ها و روابط، تلاش کردند تا مفاهیم و روابط بین آن‌ها را استخراج کنند (Hsieh et al. 2011). «فایان، واچر و شرودر» روشی را برای توسعه هستی‌شناسی موجود، از طریق جست‌وجوی اینترنتی واژگان موجود هستی‌شناسی و متن‌کاوی بدون ساختار ارائه دادند (Fabian, Wächter and Schroeder 2012). «ساتسارونیس» و همکاران روشی نیمه‌خودکار برای استخراج تعاریف متداول از مفاهیم زیست‌پزشکی به‌منظور دستیابی به روابط بین مفاهیم ارائه دادند (Tsatsaronis et al. 2013).

«شیاکس» و همکاران روش جمع‌آوری سریع جزئیات تاکسونومی را با روش متن‌کاوی اجرا کردند. آن‌ها زیرساخت متن‌کاوی «آلویس» را روی مجموعه متن به‌دست‌آمده از جست‌وجو در «پابمد» و حوزه میکروبیولوژی مواد غذایی به‌کار گرفتند (Chaix et al. 2019). «شی، جین و سیه» روشی نیمه‌خودکار برای توسعه هستی‌شناسی پایه با استفاده از یک استراتژی ترکیبی شامل روش بالا به پایین (غنی‌سازی هستی‌شناسی موجود) و روش پایین به بالا (کشف مفاهیم و روابط از مجموعه اسناد خاص دامنه) پیشنهاد دادند (Chi, Jin and Hsieh 2019). به‌دلیل پیچیدگی بسیار زیاد در اصطلاحات و روابط تاکسونومی یا سلسله‌مراتبی در حوزه‌های فنی، ایده توسعه کاملاً خودکار هستی‌شناسی هنوز عملیاتی به نظر نمی‌رسد. این است که در مقاله حاضر رویکرد نیمه‌خودکار اتخاذ شده و مشابه این انتخاب در مقالات (Jiang and Tan 2005, Missikoff, Velardi, and Fabriani 2003, 1992 Hearst, Zhang et al. 2007, Hsieh et al. 2011, Chi, Jin and Hsieh 2019) مورد تأکید قرار گرفته است. بر اساس مقالات بررسی‌شده، اتخاذ رویکرد نیمه‌خودکار در راستای توسعه و بهبود مستمر هستی‌شناسی توسط بسیاری از محققان مورد توجه قرار گرفته است. چنین به نظر می‌رسد که این رویکرد همچنان برای به‌روزرسانی هستی‌شناسی بر مبنای اسناد و انتشاراتی که حجم انبوهی از متن را دربردارند، ترجیح بیشتری داشته باشد. از سوی دیگر، هر حوزه موضوعی با وجود اشتراک با دیگر حوزه‌های موضوعی، دارای مفاهیم، اصول، الگوها و نظریه‌های یکتاست. این مسئله به‌ویژه در حوزه‌های تخصصی پزشکی ملموس‌تر است. بنابراین، لازم است فعالیت‌هایی به‌منظور تحلیل حوزه صورت گیرد و ضمن شناسایی مفاهیم و روابط آن‌ها، چارچوبی برای بیان مفاهیم حوزه ارائه شود. بر این اساس، در مقاله

حاضر، برای انتخاب مفاهیم هستی‌شناسی از متن کاوی انتشارات علمی معتبر با رویکرد نیمه‌خودکار استفاده شده است.

گفته می‌شود که ایجاد هستی‌شناسی در هر حوزه نیازمند احاطه کامل موضوعی است؛ اما اغلب متخصصان علم اطلاعات و دانش‌شناسی، بر اساس چشم‌انداز و مأموریت شغلی، در مراکز اطلاعاتی تخصصی امروزی مانند پزشکی، صرفاً با احاطه نسبی به موضوع با مبحث هستی‌شناسی‌ها مواجهه دارند. بنابراین، طرح حاضر با هدف ارائه راه‌حلی سریع و کاربردی برای توسعه هستی‌شناسی بر مبنای اطلاعات متنی انبوه حاصل از جست‌وجوی هدفمند بسامد واژگان در پایگاه‌های دانش ارائه شده است تا با تسلط نسبی عامل انسانی، بیشترین پوشش موضوعی به دامنه تحت بررسی را فراهم سازد. بخش بعد به مرور رویکردهای خودکار به کاررفته برای توسعه هستی‌شناسی پرداخته و سپس، روش مقاله فعلی در طرح توسعه هستی‌شناسی نیمه‌خودکار ارائه شده است. در ادامه، به تشریح جزئیات مطالعه موردی که برای توسعه هستی‌شناسی «گلوکوم» انجام شده، پرداخته‌ایم.

۳. روش

این پژوهش از نوع بنیادی-کاربردی است که با هدف طراحی هستی‌شناسی به‌عنوان ابزار بازنمون دانش در یکی از حوزه‌های تخصصی پزشکی و بیماری‌های چشم (گلوکوم) یا آب سیاه انجام می‌گردد. بخش بنیادی آن نرم‌افزار کدنویسی‌شده‌ای است که با نام "TmbOnt_Alfa" معرفی شده است و می‌توان از آن برای پژوهش‌های دیگر نیز استفاده کرد. بخش کاربردی آن بالفعل بوده و می‌توان آن را برای سازماندهی، آموزش و ... به کار بست. بنابراین، هستی‌شناسی به‌دست آمده نه تنها قابلیت کاربردی دارد، بلکه از ابزار متن‌کاوی طراحی‌شده نیز می‌توان در گسترش هستی‌شناسی‌های سایر حوزه‌های پزشکی و غیرپزشکی هم استفاده کرد.

طبق نتایج حاصل از پیشینه‌ها، در پژوهش حاضر رویکرد نیمه‌خودکار در استخراج مفاهیم و روابط هستی‌شناسی اتخاذ شده است. بدین منظور، یک رابط کاربری نرم‌افزاری با نام "TmbOnt_Alfa" در بستر نرم‌افزار «متلب»^۱ ایجاد شده که قادر است پردازش‌های لازم به روی متن انبوه را، بدون محدودیت در حجم و فرمت متن، به‌صورت متمرکز انجام دهد.

1. Matlab

فراهم‌سازی متن انبوه می‌تواند با استفاده از موتورهای جست‌وجوی متداول مانند «پابمد» انجام شود. نرم‌افزار "TmbOnt_Alfa"، بلوک فایل متنی را به‌عنوان ورودی دریافت کرده، محتوای آن را مطابق تنظیمات تعریف‌شده توسط کاربر پردازش نموده و واژگان مؤثر را بر اساس تعدد تکرار و میزان مرتبط بودن آن‌ها با کلمه سرشاخه هستی‌شناسی، رتبه‌بندی می‌کند. بر اساس نتایج ارائه‌شده توسط این نرم‌افزار، کاربر خبره می‌تواند کلمات پیشنهادشده را بررسی نموده و در خصوص اینکه آیا این عبارت از نظر متنی مربوط به کلمه اصلی است یا خیر، تصمیم‌گیری کند. این رویکرد، افزون بر امکان ایجاد هستی‌شناسی جدید و مبنای می‌تواند به‌منظور ارتقای هستی‌شناسی‌های موجود و جهت به‌روزرسانی آن، مطابق با انتشارات علمی حوزه دانش مربوط مورد استفاده قرار گیرد. در این پژوهش، طرح پیشنهادی در قالب مطالعه موردی برای توسعه هستی‌شناسی «گلوکوم» به کار گرفته شده است. بلوک فایل متنی در این مطالعه، صرفاً از طریق جست‌وجوی هدفمند واژه "Glaucoma" در «پابمد»، فراهم شده است. در شکل ۱، طرح کلی رویکرد توسعه نیمه‌خودکار هستی‌شناسی نشان داده شده است. این طرح مشتمل بر سه ماژول اصلی شامل ماژول‌های پیش‌پردازش، ماژول پردازش اصلی و ماژول رابط کاربری است. ماژول پیش‌پردازش شامل زیرماژول توکن‌ساز (واحدسازی)^۱، برچسب‌گذار، نرمال‌ساز کلمات، حذف‌کننده علائم نگارشی، حروف اضافه و محدودیت طول یا تعداد کاراکتر کلمات است. فایل داده متنی، به‌صورت بلوک داده و از طریق رابط کاربری وارد فرایند پردازش می‌شود. زیرماژول توکن‌ساز، آرایه متنی^۲ را تبدیل به مجموعه کلمات متسلسل می‌کند و این امکان فراهم می‌آید که با استفاده از برنامه کامپیوتری، بروی تک‌تک کلمات کار پردازشی صورت پذیرد. توسط زیرماژول برچسب‌گذار، برای هر کدام از کلمات توکن‌سازی‌شده، بر اساس گرامر و موقعیت کلمه در جمله و ساختار آن، برچسب گفتاری معینی (فعل، فاعل، صفت و ...) تخصیص داده می‌شود. زیرماژول نرمال‌ساز، کلمات برچسب‌گذاری‌شده را به‌صورت ریشه‌ای آن بازمی‌گرداند تا تغییرات ناشی از اعمال زمان در افعال و صورت جمع/ منفرد کلمات را حذف و همسان نماید. زیرماژول حذف‌کننده علائم نگارشی، مواردی از قبیل نقطه، ویرگول، ممیز و ... را از ساختار جمله حذف می‌کند. در ادامه، زیرماژول حذف‌کننده حروف اضافه، اقدام به حذف حروف اضافه از قبیل "the, at, of, ..." می‌نماید. سرانجام، این امکان برای کاربر فراهم شده است که در صورت نیاز، بتواند تعداد

1. tokenized

2. text string

حداقلی کاراکترهای کلمات را معین کند تا کلماتی که طول کوتاهی دارند، از فرایند پردازش اصلی حذف گردند. کلیه این زیرماژول‌ها به شکل انتخابی برای کاربر تدارک دیده شده تا در صورت صلاحدید بتواند در پردازش‌های خاص، برخی از زیرماژول‌ها را تغییر دهد یا عملیاتی نکند. هرچند که حذف هر یک از مراحل پیش‌پردازش منجر به افزایش حجم داده و به‌دنبال آن، منجر به افزایش زمان پردازش خواهد شد. پس از اتمام عملیات پیش‌پردازش و در فرایند پردازش تحت نظارت کاربر، ساختار سلسله‌مراتبی مفاهیم و برخی روابط هستی‌شناسی معین می‌شود. طرح پیشنهادی در عین سادگی، دسترس‌پذیری، و مقرون‌به‌صرفه بودن، در ارزیابی با هستی‌شناسی‌ها و اصطلاحنامه‌های موجود پزشکی، ساختاری قابل قبول ارائه می‌کند و امکان به‌روزرسانی و توسعه ساختار هستی‌شناسی‌های موجود را بر مبنای آخرین گزارشات تخصصی فراهم می‌آورد. این طرح می‌تواند توسط کاربر غیرمتخصص با دانش دامنه مورد بررسی که اطلاعاتی کلی در خصوص هستی‌شناسی دارد نیز به کار گرفته شود.

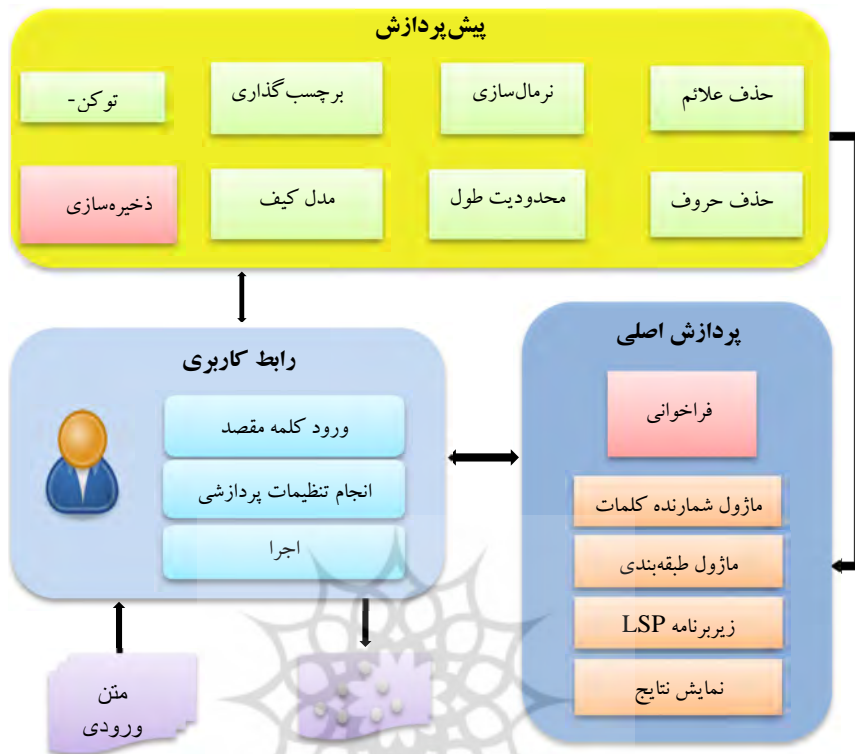
۳-۱. گردآوری داده از دامنه

بر اساس موضوع هستی‌شناسی، دانش ارزشمند باید از منابع معتبر از قبیل دایره‌المعارف‌ها، کتب مرجع، گزارش‌های علمی و یا مقالات علمی معتبر گردآوری شود. انتخاب منبع جامع برای ورودی فرایند پردازش می‌تواند ساختار هستی‌شناسی کامل و جامع را به‌دست دهد. امروزه، به لطف مجلات تحت وب و پایگاه‌های اطلاعاتی علمی، امکان استفاده گسترده از مقالات علمی برای اعتباربخشی به داده متنی امکان‌پذیر است. در پژوهش حاضر، شناسایی دامنه صرفاً با به‌کارگیری کلمه سرشاخه هستی‌شناسی^۱ و جمع‌فایل مقالات علمی حول این کلیدواژه صورت می‌پذیرد. بدین نحو که فایل متنی انبوه مورد نیاز برای ورودی فرایند متن‌کاوی، با انتخاب کلیدواژه اصلی دامنه هدف در «پابمد» و ذخیره‌سازی مرتبط‌ترین نتایج در قالب تک‌فایل متنی تهیه و استخراج می‌گردد. تأمین متن ساده و خالص (عاری از آدرس‌های اینترنتی، جداول، تصاویر و ...) منجر به کاهش زمان پردازش می‌شود. از این رو، برای سهولت پردازش، لازم است فایل متنی خام در ساده‌ترین قالب و به فرمت «اسکی»^۲ تهیه شود.

1. Target OR Seed

2. ASCII

2.



شکل ۱. طرح کلی توسعه هستی شناسی نیمه خودکار

۲-۳. پردازش داده

رویکردهای توسعه هستی شناسی خودکار بایستی بتوانند راهکار مناسبی برای دو چالش اصلی شناسایی مفاهیم و شناسایی روابط ارائه دهند (Liu, Hogan and Crowley 2011). این راهکارها را می توان به دو دسته کلی شامل روش های نمادین^۱ و یا روش های آماری تفکیک کرد. فارغ از نوع راهکار، به طور کلی استخراج مفاهیم و مترادف ها نسبت به استخراج روابط با چالش کمتری مواجه است (Liu, Hogan and Crowley 2011; Maedche, Pekar & Staab 2003).

روش نمادین اغلب مبتنی بر رویکرد تطابق الگوی واژگانی است که توسط Hearst (1992) ارائه شده و سعی دارد رابطه اصطلاحات را از طریق سبک های واژگانی که از قبل

1. symbolic

تبیین شده‌اند، شناسایی کند. شکل کلی و بنیادین الگوهای ترکیبی^۱ به صورت قوانین «اگر-آنگاه» در جدول ۱، نشان داده شده است.

از سوی دیگر، در روش‌های آماری از خوشه‌بندی و یا یادگیری ماشین برای شناسایی مفاهیم جدید و یافتن روابط آن‌ها استفاده می‌شود (Alfonseca & Manandhar 2002). روش‌های نمادین اگرچه موفقیت بیشتری در شناسایی رابطه بین مفاهیم در قیاس با روش‌های آماری نشان می‌دهند، اما روش‌های آماری به دلیل توسعه منظم تکنیک‌های یادگیری ماشین می‌توانند در راستای توسعه سیستم‌های خودکار کارایی بهتری داشته باشند (Liu, Hogan and Crowley 2011).

در تحقیق حاضر، ابزاری جهت پردازش متن با شیوه ترکیبی آماری-نمادین ارائه شده است. ماژول پردازش اصلی، داده‌هایی را که پردازش اولیه شده‌اند، به صورت مستمر فراخوانی کرده و مطابق تنظیمات تعریف شده توسط کاربر، به صورت نهایی پردازش می‌کند. ماژول اصلی مشتمل بر زیربرنامه‌های پردازش اصلی است که عملکرد جست‌وجو، خوشه‌بندی و LSP را اجرا می‌کند. این ماژول، در انتهای فرایند پردازش، قابلیت ارائه گزارش‌ها و تولید نمودار را نیز دارد.

جدول ۱. برخی از الگوهای واژگان در قالب قوانین منطقی اگر-آنگاه (Hearst 1992)

ردیف	نمونه الگوهای واژگانی
۱	<p>IF <i>NP such as {NP₁ (and/or) NP₂, ..., NP_p, ..., NP_n}</i> Then <i>NP_i ISA NP</i></p> <p>مثال: ... eye disease such as cataracts, glaucoma, macular degeneration, retinal detachment, and vision loss due to diabetes and high blood pressure is that ...^۱ ◇ Glaucoma ISA eye disease ◇ Macular degeneration ISA eye disease</p>

1. lexica-syntactic patterns (LSP)

1. <https://williamsburgeye.com/patient-education/eye-disease/>

ردیف نمونه الگوهای واژگانی

<p>IF <i>such NP as {NP₁ (and/or) NP₂, ..., NP_p, ..., NP_n}</i> Then <i>NP_i ISA NP</i></p>	۲
<p>... works by such authors as Herrick, Goldsmith, and Shakespeare... [17] ◇ Herrick ISA author ◇ Goldsmith ISA author</p>	مثال:
<p>IF <i>{NP₁ (and/or) NP₂, ..., NP_p, ..., NP_n} or/and other NP</i> Then <i>NP_i ISA NP</i></p>	۳
<p>This medication is used alone or with other medications to treat high pressure inside the eye due to glaucoma or other eye diseases (e.g., ocular hypertension)¹. ◇ Glaucoma ISA eye disease</p>	مثال:
<p>IF <i>NP including {NP₁ (and/or) NP₂, ..., NP_p, ..., NP_n}</i> Then <i>NP_i ISA NP</i></p>	۴
<p>◇ This suggests that a diabetic retinopathy screening program needs to detect and report other eye disease, including glaucoma and macular disease.² ◇ Glaucoma ISA eye disease ◇ Macular ISA eye disease</p>	مثال:
<p>IF <i>NP specially {NP₁ (and/or) NP₂, ..., NP_p, ..., NP_n}</i> Then <i>NP_i ISA NP</i></p>	۵
<p>... most European countries, specially France, England, and Spain. [17] ◇ France ISA European country ◇ England ISA European country</p>	مثال:

۳-۲-۱. شناسایی اصطلاحات

در متن‌های حجیم، شناسایی اصطلاحات و مفاهیم به‌طور معمول توسط روش‌های مبتنی بر متن‌کاوی یا روش‌های پردازش زبان طبیعی انجام می‌شود (Alfonseca & &)

1. <https://www.rxlist.com/fdb/drugs/60722/pilosol-ophthalmic-eye-drug.htm>

2. <https://www.ncbi.nlm.nih.gov/pubmed/27328169>

Manandhar 2002 و Chi, Jin, and Hsieh 2019). فرایند شناسایی اصطلاح در رویکرد حاضر، توسط یک زیربرنامه هدفمند محاسبه بسامد واژگانی انجام می‌شود. این زیربرنامه برای یافتن پیش‌واژه‌های اصطلاح مورد جست‌وجو توسعه داده شده است. زیربرنامه مذکور، به‌عنوان بخشی از برنامه اصلی، پیش‌واژه‌ها را بر اساس میزان تکرار و نزدیک بودن به واژه اصلی، اولویت‌بندی می‌کند. در واقع، بسامد واژگانی که به‌عنوان پیش‌واژه کلمه هدف قرار گرفته‌اند، محاسبه می‌شود. این است که در این پژوهش، افزون بر توجه به بسامد واژگان، کیفیت واژگان و ارتباط آن‌ها که سرانجام، چیدمان ساختار هستی‌شناسی را شکل می‌دهد نیز یکی از مهم‌ترین مراحل پردازش بوده است. طبق پژوهش «احمدی»، محل استخراج واژه در مدرک، توجه به مسائل زبانی و ترکیب واژگان، و ارتباط معنایی از جمله موارد مطرح در انتخاب واژه با کیفیت است (۱۳۹۴). مراحل مذکور به نحوی برنامه‌ریزی شده است که واژه‌ها از چکیده، که عصاره مطالب مدرک را دربردارد، استخراج شده و پس از برجسب‌گذاری^۱ بر اساس همستگی با واژه هدف و با در نظر گرفتن پیش‌واژه یا پس‌واژه ارزش‌گذاری گردند.

۳-۲-۲. شناسایی روابط

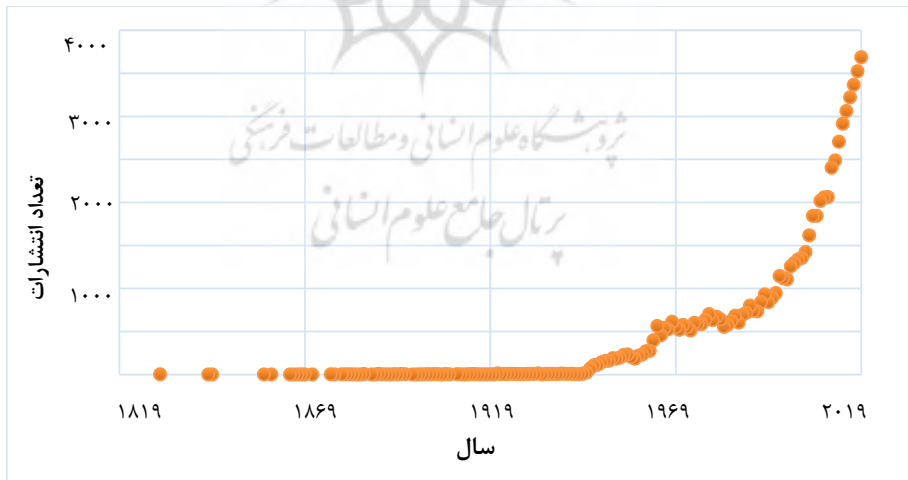
فرایند شناسایی ارتباط در این پژوهش بر مبنای دو استراتژی مجزا شامل شناسایی واژه پسین با استفاده از منطق جست‌وجوی پس‌رونده^۲ و شناسایی واژه پیشین با استفاده از منطق جست‌وجوی پیش‌رونده نسبت به واژه هدف استوار است. رویکرد شناسایی واژه پسین مبتنی بر پردازش اسم مرکب است که پیش‌تر توسط Morin (2004) مورد بحث قرار گرفته است. فرض اصلی این رویکرد بر این است که اصطلاحات حاصل از ترکیب چند واژه با یک واژه اصلی، اشاره به واژه اخص تری نسبت به شکل منفرد واژه اصلی دارند. برای مثال اصطلاح "Open-angle Glaucoma" اشاره به اصطلاح اخص‌تر در قیاس با اصطلاح "Glaucoma" دارد و با استفاده از جست‌وجوی پس‌رونده نسبت به واژه هدف "Glaucoma" در بین عباراتی که در ترکیب با این واژه در متن‌ها وجود دارد، شناسایی می‌شود. در همین حال، رویکرد پیش‌رو، متکی به روش استخراج رابطه LSP از Hearst (1992) است که در جدول ۱، ارائه شده است. این دو رویکرد، به‌صورت دو زیربرنامه مجزا برای شناسایی روابط، کدنویسی شده‌اند.

1. tagging

2. backward chaining

۴. مطالعه موردی: «گلوکوم»

طرح پیشنهادشده برای توسعه هستی‌شناسی، از جهت ارزیابی قابلیت‌های عملکردی آن در شناسایی اصطلاحات و روابط بین آن‌ها، در دامنه اطلاعاتی «گلوکوم» مورد بررسی قرار گرفت. «مرکز کنترل و پیشگیری از بیماری‌ها»^۱ در صفحه بیماری‌های چشم بیان کرده است که بیش از ۴/۲ میلیون آمریکایی با سن ۴۰ سال و بیشتر رسماً نابینا هستند. «گلوکوم» دومین عامل نابینایی بعد از آب مروارید در دنیاست. بیماری «گلوکوم» ناشی از آسیب ساختاری به عصب‌های بینایی است و دارای الگوی معینی در آسیب رأس عصب‌های بینایی است (Foster et al. 2002). بهبود هستی‌شناسی‌های موجود در این زمینه، به جامع‌تر شدن هستی‌شناسی و فراگیرتر شدن آن بر اساس آخرین یافته‌های پژوهشی کمک شایانی خواهد کرد. جست‌وجوی اولیه در «پابمد» حاکی از حجم وسیع پژوهش‌های انجام‌شده طی سال‌های اخیر در این حوزه است. در نمودار ارائه‌شده در شکل ۲، حجم انتشارات این حوزه بر اساس سال نمایش داده شده که نشانگر رشد شتابان این دسته از انتشارات طی سال‌های اخیر است. حجم انتشارات و میانگین نرخ رشد آن به وضوح رشد شتابان مرزهای علمی حوزه تخصصی «گلوکوم» را نمایش می‌دهند و از این منظر می‌تواند گزینه مناسبی برای بررسی امکان توسعه هستی‌شناسی آن بر مبنای پژوهش‌های علمی باشد.



شکل ۲. توزیع فراوانی انتشارات با محوریت «گلوکوم» بر اساس سال

1. Centers for Disease Control and Prevention (CDC)

۴-۱. ایجاد فایل متن انبوه^۱

در طرح پیشنهادی، تهیه فایل متن انبوه از بین پژوهش‌های علمی معتبر و جدید صورت گرفته است که با به‌کارگیری کلیدواژه اصلی در جست‌وجوی هدفمند این دسته از انتشارات به‌دست می‌آید. در مطالعه موردی حاضر، جست‌وجوی عبارت "Glaucoma" در «پابمد»، بیش از ۷۰،۰۰۰ نتیجه را گزارش داد که بیش از ۱۴،۰۰۰ مورد دارای خلاصه مقاله و یا متن کامل در این پایگاه اطلاعاتی بودند. از این میان، ۱۰،۰۰۰ رکورد مرتبط‌تر نخست که بالغ بر ۳،۰۰۰،۰۰۰ عبارت و اصطلاح است، برای تأمین فایل متن انبوه جمع‌گردید. این فایل با فرمت txt ذخیره شد تا در پردازش توسط نرم‌افزار "TmbOnt_Alfa" مورد استفاده قرار گیرد.

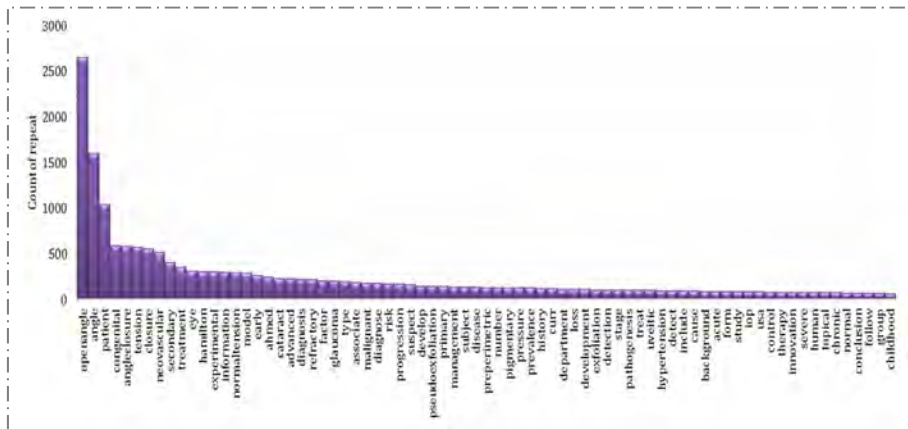
۴-۲. پردازش مقدماتی

فایل متن انبوه در رابط کاربر "TmbOnt_Alfa"، وارد فرایند پردازش گردید. ماژول پردازش مقدماتی با استفاده از تنظیمات پیش‌فرض اقدام به اعمال زیربرنامه‌های نرمال‌ساز، حذف‌کننده علائم نگارشی و حذف‌کننده حروف اضافه نمود. افزون بر این، در این مرحله، طول اصطلاحات نیز بر مبنای حداقل طول دو کاراکتر تعیین شد و برچسب‌های گفتاری برای استفاده‌های آتی اضافه گردید. خروجی مرحله پردازش مقدماتی به‌منظور استفاده مکرر در فازهای بعدی، به‌صورت فایل متغیر در فضای «متلب» ذخیره شد.

۴-۳. پردازش اصلی

سطح اول ساختار هستی‌شناسی با درج موجودیت هدف، یعنی واژه «گلوکوم» شکل گرفت. این ماژول مبتنی بر این فرض است که بیشترین شانس عضویت در لایه اول هستی‌شناسی به پیش‌واژه‌های سطح اول موجودیت هدف اختصاص دارد. بدین ترتیب، تمام پیش‌واژه‌های منسوب به واژه هدف، جمع‌شده و برای ارزیابی توسط متخصص بازنمایی شدند. در خصوص واژه «گلوکوم»، خروجی حاصل از این ماژول در شکل نشان داده شده است.

1. corpus file



شکل ۳. نتیجه جست و جوی پیش‌واژه لایه اول برای موجودیت «گلوکوم»

با توجه به توزیع فراوانی پیش‌واژه‌ها، می‌توان با تقریب مناسبی اظهار کرد که بعد از ۱۰۰ رکورد اول، میزان بروز اصطلاحات پیش‌واژه در سطح فراوانی بسیار کمتری نسبت به رکوردهای اولیه هستند. با این حال، به جهت حصول اطمینان حداکثری برای از دست ندادن هیچ موجودیت احتمالی، در اینجا ۲۵۰ کلمه برتر برای تولید سطح اول هستی‌شناسی انتخاب گردید. در این مرحله عامل انسانی می‌تواند واژگان با بسامد بالا را به‌عنوان گزینه‌های بعدی سرشاخه انتخاب کند. این رویه تا زمانی که به محدودیت تعیین شده توسط کاربر (۲۵۰ مورد در این پژوهش) و شانس ناچیز برای انتخاب شدن (بر مبنای قضاوت کاربر) برسد، ادامه می‌یابد.

۴-۳-۱. شناسایی روابط با استفاده از LSP

افزون بر رویکرد تحلیل بسامد واژگان، که در خصوص واژه «گلوکوم» شاکله کلی ساختار هستی‌شناسی را برپا کرد، کارایی شیوه LSP نیز بررسی گردید تا امکان شناسایی روابط بیشتر بین واژگان ارزیابی شود. در به‌کارگیری این رویکرد، افزون بر الگوهای LSP نمایه‌شده در جدول، الگوهای مضاعفی نیز متناسب با واژه گلوکوم که در اثنای جست‌وجو در متون تخصصی غالب به نظر می‌رسیدند، انتخاب و مطابق با جدول به فرایند جست‌وجو افزوده شدند.

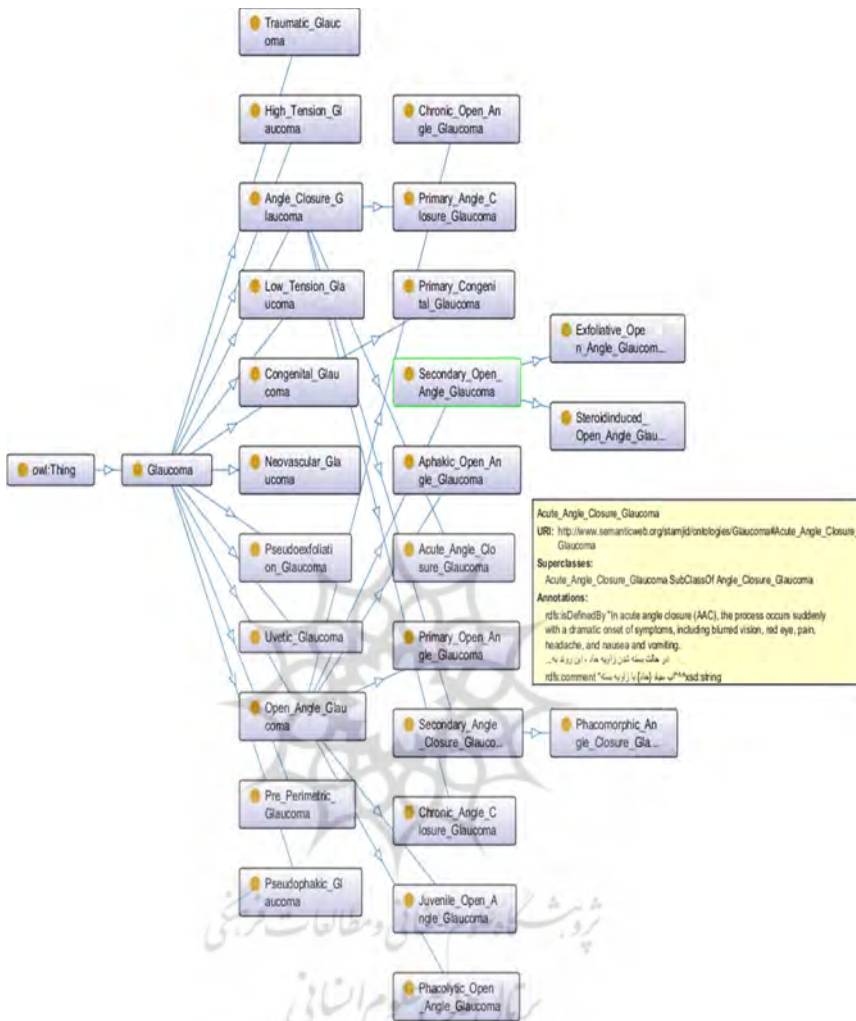
جدول ۲. الگوهای مضاعف LSP به کاررفته برای فاز مطالعه موردی

مثال	الگو
Pigmentary glaucoma is a type of secondary open-angle glaucoma...	NP1 is a type of NP
Normal-tension glaucoma (NTG), also known as low tension or normal pressure glaucoma.	NP1 also known as NP
Steroid-induced glaucoma is a form of secondary open angle glaucoma that results from the use of steroids.	NP1 is a form of NP
Exfoliative glaucoma is the most common type of secondary open-angle glaucoma worldwide.	NP1 is the most common type of NP
Childhood glaucoma, also referred to as congenital glaucoma, pediatric, or infantile glaucoma, occurs in babies and young children.	NP also referred to as NP1, NP2 or NP3...

جست‌وجوی الگوهای LSP از طریق زیربرنامه افزوده شده با همین نام به برنامه اصلی صورت می‌گیرد که این زیربرنامه کل ساختار فایل توکن شده را برای یافتن الگوهای مشابه با الگوی معرفی شده پالایش می‌کند.

۴-۳-۲. نهایی‌سازی ساختار واژگان

ساختار اولیه ایجاد شده نیازمند انجام بررسی توسط متخصص حوزه دانش است تا بتواند پس از پالایش و اعمال تغییرات لازم، به صورت نهایی تکمیل گردد. مفاهیم حاصل بر اساس روابط شناسایی شده شامل, Is a, also known, Is a form of, Is the most common type, also referred to,... در نرم‌افزار «پروتژ ۵» وارد گردید. به منظور قابلیت استفاده هستی‌شناسی در زبان فارسی، معادل فارسی مفاهیم در قسمت "comment" نرم‌افزار «پروتژ» وارد شد و همین مفاهیم برای قیاس با مراجع فارسی مورد استفاده قرار گرفت. تعاریف انگلیسی و فارسی هر مفهوم در قسمت "definition" وارد شد.



شکل ۴. ساختار نهایی شده هستی‌شناسی در «پروتز»

۵. تحلیل نتایج

رویکرد پیشنهادشده در مقاله حاضر توانسته است موضوع تنوع و پیچیدگی بسیار بالای واژگان و مفاهیم در پالایش حجم انبوه داده‌های متنی را رفع نماید. زمانی که منبع داده‌ها بانک مقالات علمی معتبر باشد، دشواری کار متن‌کاوی با توجه به ثقیل‌تر بودن ادبیات به کاررفته در چنین متن‌هایی به مراتب بیشتر می‌شود. با وجود این، در صورت غلبه بر این چالش می‌توان انتظار داشت که ساختار حاصل بتواند دامنه بیشتری

از واژگان تخصصی و علمی را پوشش داده و به کامل‌تر شدن هستی‌شناسی‌های موجود کمک مؤثری نماید. بر اساس نتایج حاصل از مرحله ارزیابی علمی، استفاده از مازول پردازش اولیه و به‌کارگیری فیلترهای مناسب در آن، کارایی بسیار مناسبی در جهت حذف محاسبات غیرضروری نشان داد که منجر به کاهش بسیار محسوس حجم و زمان محاسبات گردید. در خصوص مطالعه موردی صورت گرفته برای استخراج مفاهیم بیماری «گلوکوم»، حدود ۱۰،۰۰۰ خلاصه مقاله از بانک اطلاعاتی «پابمد» وارد تحلیل شد که مشتمل بر حدود ۳/۸ میلیون آیتم توکن شده بود. مازول پردازش اولیه این حجم از داده را به حدود ۱/۸۷ میلیون کاهش داد که به مفهوم کاهش ۵۰ درصدی حجم محاسبات غیرضروری است. نکته مهم این است که استفاده از چنین فیلتری می‌تواند به حذف عناصری منجر شود که در پردازش LSP کارایی دارند. بنابراین، در صورت به‌کارگیری الگوریتم LSP، با علم به افزایش زمان پردازش، لازم است از به‌کارگیری فیلترهایی که حروف اضافه را حذف می‌کنند، اجتناب شود. در شناسایی واژگان و روابط بین آن‌ها استفاده از رویکرد بسامد و رتبه‌بندی واژگان، در قیاس با روش LSP توانست کارایی بهتری داشته باشد و نیز توانست با شناسایی تعداد قابل قبولی از مفاهیم و روابط ساختار کلی هستی‌شناسی را توسعه دهد. این در حالی است که روش LSP صرفاً توانست لایه‌های اولیه و بالادستی هستی‌شناسی را پوشش دهد. اتکا به نیروی انسانی در مرحله نهایی‌سازی ساختار و بررسی شاخه‌های مترادف ضروری است و اگر بتوان این موضوع را با برنامه‌نویسی جامع‌تر (حداقل برای آشکارسازی و تلفیق شاخه‌های مترادف) مرتفع نمود، گام بلندی در جهت توسعه تمام‌خودکار هستی‌شناسی خواهد بود.

به‌منظور ارزیابی میزان انطباق نتایج به‌دست‌آمده با ساختارهای معتبر موجود، معیاری کمی برای قیاس با هستی‌شناسی‌های بازنمایی شده برای بیماری گلوکوم در بانک‌های هستی‌شناسی بیماری‌های انسانی توسعه داده شد. مقایسه با ساختار «سرعنوان‌های پزشکی مش»^۱ و «اصطلاحنامه و توصیفگرهای پزشکی فارسی» نیز انجام گردید. این قیاس بین مفاهیم، اصطلاحات و جایگاه سلسله‌مراتبی آن‌ها حاصل از برنامه گذشته، با مفاهیم، اصطلاحات و جایگاه هر اصطلاح در ساختار سلسله‌مراتبی هستی‌شناسی‌های معتبر موجود و اصطلاحنامه پزشکی انجام شده است. مطابق رابطه (۱)، درجه انطباق با به‌کارگیری

1. MESH Subject Headings

فرمولی ساده که بر اساس میزان انطباق و موقعیت انطباق، امتیازدهی می‌کند، به مقیاس کمی تبدیل شد. این قیاس هم تشابه در واژگان و هم تشابه در موقعیت واژه در ساختار هستی‌شناسی را مدنظر قرار داده است. اینکه دو واژه در دو ساختار درختی تحت قیاس، در موقعیت‌های همسان قرار گیرند، مؤید آن است که روابط نیز به درستی شناسایی شده‌اند. با ذکر این نکته که در خود ساختارهای ارائه‌شده در پایگاه‌های معتبر نیز اختلافاتی وجود دارد، قیاس انجام‌شده توانسته است کلیات ساختار را صحه‌گذاری کرده و تأیید نماید که رویکرد جدید توانسته گامی در مسیر سهولت ایجاد و توسعه نیمه‌خودکار هستی‌شناسی بردارد. با استفاده از این رابطه، میزان انطباق با «مش» ۸۰ درصد، انطباق و پوشش با «اصطلاحنامه پزشکی فارسی» ۱۰۰ درصد و در قیاس با ساختار واژگان «توصیفگرهای پزشکی فارسی» ۱۰۰ درصد محاسبه گردید. همچنین، در مقایسه با هستی‌شناسی بازنمایی‌شده در «هستی‌شناسی بیماری‌های بیماری‌های انسانی»^۱ میزان انطباق ۵۷ درصد و در قیاس با هستی‌شناسی نمایه‌شده در «هستی‌شناسی BioAssay»، ۹۱ درصد به دست آمد. افزون بر آن، رویکرد پیشنهادی توانست به‌طور میانگین حدود ۳۰ درصد به دامنه واژگان ساختار هستی‌شناسی موجود بیفزاید. این توسعه قابل توجه به دلیل توانمندی این شیوه برای پردازش حجم انبوه از متن‌های تخصصی است و با توجه به اینکه این متون در فضای ادبیاتی دانشگاهی و بر اساس به‌روزترین نتایج محققان تحریر می‌شود، می‌تواند نویدبخش قابلیت نسبی روش پیشنهادی برای بهبود هستی‌شناسی‌های موجود باشد. همچنین بر اساس نتایج به دست آمده، استفاده از این رویکرد با درجه مناسبی از اطمینان می‌تواند برای تهیه فرم نخست ساختار هستی‌شناسی در دامنه‌های متنوع و خاص مورد استفاده قرار گیرد. بر اساس نتایج حاصل از فاز تجربی، تهیه ساختار هستی‌شناسی گلوکوم با میزان انطباق ذکرشده صرفاً با صرف زمانی حدود پنج ساعت و توسط نیروی انسانی با سطح دانش و مهارت متوسط صورت پذیرفت که با توجه به حجم داده پردازش‌شده و کیفیت نتایج، می‌تواند راهبردی مناسب برای توسعه اولیه ساختار هستی‌شناسی و نیز نگاهداشت ساختارهای موجود با صرف حداقلی منابع باشد. این دستاورد در قیاس با کار مشابه (2019) (Chi, Jin and Hsieh) که به کارگیری ۱۶ نفر طی دو ساعت را گزارش داده‌اند، وابستگی به نیروی انسانی را تا شش برابر کاهش داده است. علت اصلی این بهبود، امکانات سطح

1. Human Disease Ontology (DOID)

بالای ماژول‌های پردازش متن به کار گرفته شده در نرم‌افزار TmbOnt_Alfa و عدم نیاز به بررسی متون ورودی به تحلیل توسط نیروی انسانی خبره است. سرانجام، با توجه به اینکه کار حاضر در بستر کدنویسی «متلب» توسعه داده شده که مجهز به ماژول‌های پردازشی بسیار گسترده در حوزه هوش مصنوعی و یادگیری ماشین است، می‌تواند بر خلاف کارهای مشابه که از بسترهای متن‌کاوی صرف استفاده نموده‌اند (مانند کار انجام شده در پژوهش Chaix et al. (2019) که از بستر مستقل متن‌کاوی Alvis استفاده شده)، انعطاف بسیار بالایی در توسعه‌های آتی داشته باشد.

$$\text{Compatibility Degree} = \frac{a * 3 + b * 2}{all * 3} * 100 \quad (1)$$

a: واژه‌های مشابه در موقعیت‌های همسان

b: واژه‌های مشابه در موقعیت‌های غیر همسان

شکل ۵. ساختار درختی واژه «گلوکوم» در «مش»

۶. نتیجه‌گیری

در این مقاله، راهکاری جهت توسعه هستی‌شناسی به روش نیمه‌خودکار ارائه گردید که بستری مناسب برای ایجاد یک هستی‌شناسی جدید و یا توسعه هستی‌شناسی موجود

بر مبنای انتشارات معتبر علمی نمایه‌شده در پایگاه‌های اطلاعات علمی فراهم می‌کند. در این روش، متن انبوه ورودی در قالب فایل متنی فراخوانی شده و بر اساس تنظیمات کاربر پردازش می‌شود. حاصل پردازش اولیه، به صورت اطلاعات میانی ذخیره می‌گردد و کاربر خبره می‌تواند بدون نیاز به صرف زمان برای پردازش مجدد، و صرفاً با وارد نمودن کلمه مورد نظر، نسبت به فراخوانی محتوای پیش‌پردازش‌شده و استخراج کلمات مرتبط بر اساس اولویت اقدام نماید. فرمت فایل ورودی برای متن می‌تواند به صورت pdf، txt و یا سایر فرمت‌های متنی بوده و کاربر انسانی صرفاً برای تصمیم‌گیری و مراحل نهایی‌سازی درگیر کار می‌شود. برتری رویکرد نیمه‌خودکار در توسعه و بهبود هستی‌شناسی‌ها نسبت به روش‌های دستی که صرفاً مبتنی بر عامل انسانی هستند، از نظر صرفه‌جویی در هزینه و زمان قابل توجه است. بر مبنای این رویکرد و در توسعه‌های آتی، این امکان وجود خواهد داشت که با کدنویسی فعالیت‌های کاربر انسانی، امکان ایجاد و بررسی رویکردهای کاملاً خودکار فراهم گردد.

بر اساس یافته‌های حاصل از مطالعه موردی توسعه هستی‌شناسی «گلوکوم»، هستی‌شناسی توسعه داده‌شده توسط رویکرد پیشنهادی به‌طور میانگین بیش از ۷۰ درصد با هستی‌شناسی‌های بازنمایی‌شده در پایگاه‌های معتبر هستی‌شناسی بیماری‌های انسانی^۱ انطباق داشت و به‌طور میانگین بیش از ۳۰ درصد واژگان جدید برای افزودن به دامنه واژگان آن‌ها را فراهم کرده است. بستر فراهم‌شده در پژوهش حاضر، صرفاً محدود به دامنه اطلاعاتی خاصی نبوده و برای طراحی و توسعه هستی‌شناسی در هر حوزه‌ای قابل استفاده است. همچنین، هستی‌شناسی به‌دست‌آمده بر خلاف هستی‌شناسی‌هایی که صرفاً برای یک کارکرد طراحی می‌شوند، قابلیت استفاده در کاربردهای مختلف مانند فهرست‌نویسی، نمایه‌سازی، آموزش و ... را دارد.

برای پژوهش‌های آتی، پیشنهاد می‌شود که ایجاد هستی‌شناسی در سایر حوزه‌های موضوعی با استفاده از رویکرد متن‌کاوی، همچنین قیاس هستی‌شناسی‌های به‌دست‌آمده از روش متن‌کاوی با هستی‌شناسی‌هایی که مفاهیم آن‌ها صرفاً از اصطلاحنامه‌ها حاصل شده‌اند، در دستور کار قرار گیرد. افزون بر آن، به نظر می‌رسد که بر مبنای رویکرد پیشنهادی در پژوهش حاضر، امکان حذف نقش ناظر انسانی (یا رساندن آن به حداقل

1. bio-ontologies

ممکن) می‌تواند بیشتر قابل اجرا باشد. این ایده می‌تواند گام مؤثری در راهبرد به‌سوی تولید کاملاً خودکار ساختارهای هستی‌شناسی در پژوهش‌های آینده باشد.

فهرست منابع

- احمدی، حمید. ترسیم و تحلیل شبکه مفهومی و هستی‌شناسی ساختار دانش حوزه علم سنجی ایران بر اساس رویکرد تحلیل حوزه. رساله جهت دریافت درجه دکتری. دانشگاه چمران اهواز. ۱۳۹۴.
- حسینی بهشتی، ملوک‌السادات. ۱۳۹۲. *ساخت‌واژه: اصطلاح‌شناسی و مهندسی دانش*. تهران: پژوهشگاه علوم و فناوری اطلاعات ایران؛ چاپار.
- فتحیان دستگردی، اکرم. ۱۳۸۹. مقایسه کارآمدی اصطلاحنامه و هستی‌شناسی در بازنمون دانش و بازیابی مفاهیم. پایان‌نامه جهت دریافت درجه کارشناسی ارشد، دانشگاه فردوسی مشهد، دانشکده علوم تربیتی و روان‌شناسی.
- معصومی، رحیم، امین معصومی گنجگاه، حبیب اوجاقی، عیسی بنزاده. ۱۳۹۱. توزیع فراوانی علل اختلالات بینایی در افراد بالای ۴۰ سال مراجعه‌کننده به درمانگاه چشم بیمارستان علوی طی سال‌های ۸۵-۱۳۸۴. *مجله دانشگاه علوم پزشکی و خدمات بهداشتی درمانی اردبیل* ۱۲ (۲): ۱۶۶-۱۷۲.

References

- Alfonseca, E., & S. Manandhar. 2002. An Unsupervised Method for General Named Entity Recognition and Automated Concept Discovery. *Proc 1st Int Conf Gen WordNet Mysore India* 69: 1-9. <https://doi.org/2682189>
- Babcock, S., J. Beverley, L. G. Cowell, & B. Smith. 2021. The Infectious Disease Ontology in the age of COVID-19. *J Biomed Semantics*. 2021 Jul 18; 12 (1):13. doi: 10.1186/s13326-021-00245-1. PMID: 34275487; PMCID: PMC8286442. <https://pubmed.ncbi.nlm.nih.gov/34275487/>
- Chaix, E., L. Deléger, R. Bossy, & C. Nédellec. 2019. Text mining tools for extracting information about microbial biodiversity in food. *Food Microbiol* 81: 63-75. <https://doi.org/10.1016/j.fm.2018.04.011>
- Charlet, J., B. Bachimont, & M. C. Jaulent. 2006. Building medical ontologies by terminology extraction from texts: An experiment for the intensive care units. *Comput Biol Med* 36: 857-870. <https://doi.org/10.1016/j.combiomed.2005.04.012>
- Chi, N. W., Y. H. Jin, & S. H. Hsieh. 2019. Developing base domain ontology from a reference collection to aid information retrieval. *Autom Constr* 100:180-189. <https://doi.org/10.1016/j.autcon.2019.01.001>
- Dutta, B., & M. DeBellis. 2020. CODO: an ontology for collection and analysis of COVID-19 data. *arXiv preprint arXiv:2009.01210*. <https://doi.org/10.48550/arXiv.2009.01210>
- Fabian, G., T. Wächter, & M. Schroeder. 2012. Extending ontologies by finding siblings using set expansion techniques. 28: 292-300. <https://doi.org/10.1093/bioinformatics/bts215>
- Foster, P. J., R. Buhrmann, H. A. Quigley, & GJ Johnson. 2002 The definition and classification of glaucoma in prevalence surveys. *Br J Ophthalmol* 86: 238-242. <https://doi.org/10.1136/bjo.86.2.238>
- Hearst, MA. 1992. Automatic acquisition of hyponyms from large text corpora. *Proc 14th Conf Comput Linguist* 23-28. <https://doi.org/https://doi.org/10.3115/992133.992154>

- Hsieh, SHang-hsieh, Lin Hsien-Tang, Chi NaiWen, kuang wu Chou, and Ken Yu Lin. 2011. Enabling the development of base domain ontology through extraction of knowledge from engineering domain handbooks. *Adv Eng Informatics* 25: 288–296. <https://doi.org/10.1016/j.aei.2010.08.004>
- Jiang, X, & AH Tan. 2005. Mining ontological knowledge from domain-specific text documents. *Proc - IEEE Int Conf Data Mining, ICDM* 665–668. <https://doi.org/10.1109/ICDM.2005.97>
- Kibbe, A. Warren, Cesar Arze, Victor Felix, Elvira Mitraka, Evan Bolton, Gang Fu, J. Christopher, Mungall Janos, X Binder, James Malone, Drashtti Vasant, Helen Parkinson, and Lynn M. Schriml. (2015) Disease Ontology 2015 update: An expanded and updated database of Human diseases for linking biomedical knowledge through disease data. *Nucleic Acids Res* 43: D1071–D1078. <https://doi.org/10.1093/nar/gku1011>
- Kless, Daniel, Ludger Jansen, & S. Milton. 2016. A content-focused method for re-engineering thesauri into semantically adequate ontologies using OWL. *Semantic Web*, vol. 7, no. 5, pp. 543-576. <https://content.iospress.com/articles/semantic-web/sw194>
- Liu, K, WR Hogan, & RS Crowley. 2011. Natural Language Processing methods and systems for biomedical ontology learning. *J Biomed Inform* 44: 163–79. <https://doi.org/10.1016/j.jbi.2010.07.006>
- Maedche, A., V. Pekar, & S. Staab. 2003. Ontology Learning Part One — on Discovering Taxonomic Relations from the Web. *Web Intell* 301–319. https://doi.org/10.1007/978-3-662-05320-1_14
- Missikoff, M., P. Velardi, & P. Fabriani. 2003. Text mining techniques to automatically enrich a domain ontology. *Appl Intell* 18: 323–340. <https://doi.org/10.1023/A:1023254205945>
- Morin, E., & C. Jacquemin. 2004 Automatic acquisition and expansion of hypernym links. *Comput Hum* 38: 363–396. <https://doi.org/10.1007/s10579-004-1926-2>
- Qian, Wang, Tao Lan, Zhu Lijun. 2007. Approach to ontology construction based on text mining, *New Zealand Journal of Agricultural Research* 50(5): 1383-1391. <https://doi.org/10.1080/00288230709510426>
- Soergel, D., Boris Lauser, Anita Liang, & Frehiwot Fisseha. 2004. Reengineering thesauri for New Application: the AGROVOC Example. *Journal of Digital Information*, vol 4: 4. Article No. 257. <https://www.fao.org/3/af234e/af234e.pdf> (accessed Nov. 19, 2021)
- Tsatsaronis, George, Petrova Alina, Kissa Maria, Yue Ma, Felix Distel, FranzBaader, and Michael Schroeder. 2013. Learning Formal Definitions for Biomedical Concepts. OWLED. https://www.researchgate.net/publication/244484656_Learning_Formal_Definitions_for_Biomedical_Concepts (accessed Aug. 11, 2019)
- Wächter T, & M. Schroeder. 2010. Semi-automated ontology generation within OBO-Edit. *Bioinformatics* 26: 88–96. <https://doi.org/10.1093/bioinformatics/btq188>
- Zhang Xiaodan, Liping Jing, Xiaohua Hu, and Xiaoua Zhou. 2007. A Comparative Study of Ontology Based Term Similarity Measures on PubMed Document Clustering. *Adv Databases Concepts, Syst Appl* 115–126. https://doi.org/10.1007/978-3-540-71703-4_12

سمیه تمجید

متولد ۱۳۶۱ دارای مدرک دکتری در رشته علم اطلاعات و دانش‌شناسی است. ایشان هم‌اکنون در کتابخانه مرکزی دانشگاه علوم پزشکی ایران مشغول فعالیت است. کتاب «راهنمای کاربردی جست‌وجو در پایگاه‌های اطلاعاتی» از جمله آثار مشترک ایشان است.

جست‌وجو و بازیابی اطلاعات، هستی‌شناسی، مدیریت فناوری اطلاعات از علایق پژوهشی ایشان است.



فاطمه نوشین فرد

دارای مدرک تحصیلی دکتری در رشته علم اطلاعات و دانش‌شناسی است. ایشان هم‌اکنون دانشیار گروه علم ارتباطات و دانش‌شناسی واحد علوم و تحقیقات تهران دانشگاه آزاد اسلامی است. سیرنیتیک، مدیریت دانش، رفتار اطلاع‌یابی، بازیابی اطلاعات، علم‌سنجی، کتابخانه‌های دیجیتال، و بازاریابی از جمله علایق پژوهشی ایشان است.



ملوک‌السادات حسینی بهشتی

دارای مدرک تحصیلی دکتری در رشته زبان‌شناسی همگانی از دانشگاه تهران است. ایشان هم‌اکنون دانشیار پژوهشکده علوم اطلاعات پژوهشگاه علوم و فناوری اطلاعات ایران (ایرنداک) است. کتاب «ساختواژه، اصطلاح‌شناسی و مهندسی دانش» از جمله آثار ایشان است. زبان‌شناسی، هستی‌شناسی، اصطلاح‌شناسی و مهندسی دانش از جمله علایق پژوهشی ایشان است.



نجلا حریری

دارای مدرک تحصیلی دکتری در رشته علم اطلاعات و دانش‌شناسی است. وی هم‌اکنون استاد گروه علم ارتباطات و دانش‌شناسی واحد علوم و تحقیقات تهران دانشگاه آزاد اسلامی است. کتاب «اصول و روش‌های پژوهش کیفی» از آثار ایشان است.



پژوهشگاه علوم انسانی و مطالعات فرهنگی
فهیمة باب الحوائجی
تالیفات جامع علوم انسانی

متولد سال ۱۳۵۵، دارای مدرک تحصیلی دکتری در رشته علم اطلاعات و دانش‌شناسی است. ایشان هم‌اکنون دانشیار گروه علم ارتباطات و دانش‌شناسی واحد علوم و تحقیقات تهران دانشگاه آزاد اسلامی است. اقتصاد اطلاعات، معماری اطلاعات و ذخیره و بازیابی اطلاعات از جمله علایق پژوهشی ایشان است.



پژوهش نامه
پژدازش و
مدیریت
اطلاعات

پژوهشگاه علوم انسانی و مطالعات فرهنگی
رتال جامع علوم انسانی