

علیرضا یاری مدرک لیسانس خود را در زمینه سیستم های کنترل در سال ۱۹۹۳ از دانشکده فنی دانشگاه تهران دریافت کرده است. در ادامه ایشان تحصیلات تکمیلی خود در زمینه مهندسی سیستم در مقاطعه فوق لیسانس و دکترا در دانشکده کامپیوتر دانشگاه فنی کیتامی کشور را بین ادامه داده که در سال



۲۰۰۰ موفق به دریافت مدرک دکترا از آن دانشگاه شده است. در حال حاضر زمینه تحقیقاتی ایشان در خصوص مراکزداده، سرویسهای وب و بسترسازی مناسب برای زبان فارسی در محیط رایانه ای می باشد. ایشان در حال حاضر مدیر گروه نرم افزار و سکوهای فناوری اطلاعات در مرکز تحقیقات مخابرات ایران می باشد.

ابوالفضل آل احمد در سال ۱۳۸۰ مدرک فوق دیپلم خود را در رشته نرم افزار کامپیوتر از مرکزآموزش عالی فنی مهندسی شهید باهنر شیراز، در سال ۱۳۸۲ مدرک کارشناسی خود را در رشته مهندسی کامپیوتر از دانشگاه شهید باهنر کرمان و در سال ۱۳۸۷ مدرک کارشناسی ارشد خود را در رشته



مهندسی فناوری اطلاعات از دانشگاه تهران دریافت کرده است. ایشان از اعضاء اصلی بروژه های تحقیقاتی موفقی در بتروشمی شیراز و پژوهشگاه صنعت نفت بوده است و از سال ۱۳۸۴ تا کنون عضو گروه تحقیقاتی پایگاه داده ها دانشگاه تهران بوده است. زمینه تحقیقاتی ایشان بازیابی اطلاعات، آنالیز وب و داده کاوی است.



- [16] B. Novak, "A Survey of Focused Web Crawling Algorithms", SIKDD 2004 Multi-Conference IS 2004, pp: 12–15, 2004.
- [17] K. Somboonviwat, T. Tamura, and M. Kitsuregawa, "Finding thai web pages in foreign web spaces", In ICDE Workshops, p. 135, 2006.
- [18] G. Botha and E. Barnards, "Two approaches to gathering text corpora from the World Wide Web", In Proceedings of the 16th Annual Symposium of the Pattern Recognition Association of South Africa, Langebaan, South-Africa, p. 194, November, 2005.
- [19] C. Castillo, "Effective Web crawling", Ph.D. Thesis, University of Chile, Department of Computer Science, 2004.
- [20] J. Kleinberg, "Authoritative sources in a hyperlinked environment", In Proceedings ACM-SIAM Symposium on Discrete Algorithms, 1998, also appears as IBM Research Report RJ 10076(91892) and online at <http://www.cs.cornell.edu/home/kleinber/auth.ps>.
- [21] G. Grefenstette, "Comparing two language identification schemes", In Proceedings of JADT 1995, 3rd International Conference on Statistical Analysis of Textual Data, 1995.
- [22] W.B. Cavnar, J. M. Trenkle, "N-gram-based text categorization", In Symposium on Document Analysis and Information Retrieval, Las Vegas, pp:161-175, 1994.
- [23] G. Churcher, "Distinctive character sequences", personal communication by Ted Dunning, 1994
- [24] A.H. Keyhanipour, A. Mohammad Zareh Bidoki, M. Mahmoudi, M. Azadnia, "Evaluation of Iran's Web Content from e-Government perspective", 12th International CSI conference, 2007.
- [25] E. Darrudi, M. R Hejazi, F. Oroumchian, "Assessment of a Modern Farsi Corpus", In Proceedings of the 2nd Workshop on Information Technology & its Disciplines (WITID) 2004, ITRC, Kish Island, Iran.
- [26] Ricardo Baeza-Yates, Carlos Castillo, Vicente Lopez, "Characteristics of the Web of Spain".
- [27] Daniel Gomes, Mario J. Silva, "A characterization of the portuguese web", University of Lisbon, Potugal, 2004.
- [28] R. Baeza-Yates and C. Castillo, Characterization of national web domains. Technical report, Universitat Pompeu Fabra, 2005.

[۲۹] امیر حسین کیهانی‌بور، علی‌محمد زارع بیدکی، مریم محمودی، محمد آزادنیا، "ازبایی محتوای وب ایران از منظر دولت الکترونیک"، دوازدهمین کنفرانس انجمن کامپیوتر ایران

[30] <http://ece.ut.ac.ir/dbrg/Hamshahri/>

مخصوصه عظیم‌زاده فارغ‌التحصیل رشته مهندسی کامپیوتر گرایش نرم‌افزار از دانشگاه تربیت معلم تهران در سال ۱۳۸۰ بوده و در سال ۱۳۸۵ مدرک کارشناسی ارشد خود را در همین رشته- گرایش از دانشگاه آزاد واحد تهران جنوب اخذ نموده است فعالیتهای پژوهشی ایشان از سال ۱۳۸۰ در مرکز تحقیقات مخابرات آغاز گردیده و هم‌اکنون نیز در زمینه بازبایی اطلاعات از وب و سترسازی مناسب برای زبان فارسی به همکاری خود با این مرکز ادامه می‌دهد.



واکنشی شده و کلمات تخصصی موجود در صفحه به عنوان داده‌های بازخوردی جهت توسعه دایره کلمات و یا منابع زبانی استفاده خواهد شد.

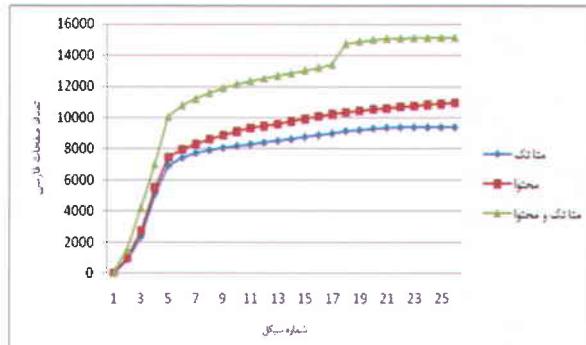
مراجع

- [1] G. Pant, P. Srinivasan, and F. Menczer, "Crawling the Web", In Web Dynamics: Adapting to Change in Content, Size, Topology and Use. Edited by M. Levene and A. Poulovassilis, Springer Verlag, pp: 153–178, 2004.
- [2] M.P.S.Bhatia, Divya Gupta, "Discussion on Web Crawlers of Search Engine", Proceedings of 2nd National Conference on Challenges & Opportunities in Information Technology (COIT-2008), Mandi Gobindgarh. March 29, pp: 227-230, 2008.
- [3] George Almanidis, Constantine Kotropoulos, Ioannis Pitas, "Combining text and link analysis for focused crawling - An application for vertical search engines", Information System 32(6), pp: 886-908, 2007.
- [4] A. Badia, T. Muezzinoglu, O. Nas-raoui, "Focused crawling: experiences in a real world project", In Proceedings of the 15International Conference on World Wide Web, Edinburgh S., pp: 1043-1044, 2006.
- [5] S. Chakrabarti, M. van den Berg, and B. Dom, "Focused crawling: a new approach to topic-specific Web resource discovery", In Proceedings of the 8th International World Wide Web Conference, Toronto, 1999.
- [6] F. Menczer, G. Pant, P. Srinivasan and M. Ruiz , "Evaluating Topic -Driven Web Crawlers", In Proceedings of the 24th Annual International ACM/SIGIR Conference, New Orleans, USA, 2001.
- [7] J. Cho, H. Garcia-Molina, and L. Page, "Efficient Crawling Through URL Ordering", In Proceedings of the 7th International World-Wide Web Conference, 1998
- [8] N. Angkawattanawit, A. Rungsawang, "Learnable Crawling: An Efficient Approach to Topic-specific Web Resource Discovery", In Proceedings of the 2nd International Symposium on Communications and Information Technology (ISCIT), 2002.
- [9] P. De Bra, G.-J. Houben, Y. Kornatzky, R. Post, "Information retrieval in distributed hypertexts", In Proceedings of RIAO'94, Intelligent Multimedia, Information Retrieval Systems and Management, New York, NY, 1994.
- [10] M. Hersovici, M. Jacovi, Y.S. Maarek, D. Pelleg, M. Shalheim, and S. Ur, "The Shark-Search algorithm - an application:tailored Web site mapping", In: 7th World-Wide, Web Conference, Brisbane, Australia, online, 1998.
- [11] M. Ehrig, A. Maedche, "Ontology-Focused Crawling of web documents", In Proceedings of the ACM Symposium on Applied Computing, 2003.
- [12] Yang, K., "Combining text- and link-based retrieval methods for Web IR", In The Ninth Text REtrieval Conf (TREC 9), 2001.
- [13] S. Raghavan and H. Garcia-Molina, "Crawling the hidden web", In Proceedings of VLDB '01, pp: 129–138, 2001.
- [14] L. Page, S. Brin, R. Motwani, and T. Winograd, "The pagerank citation ranking: Bringing order to the web", 1998.
- [15] K. Somboonviwat, M. Kitsuregawa, and T. Tamura, "Simulation Study of Language Specific Web Crawling", icde, 21st International Conference on Data Engineering (ICDE'05), p. 1254, 2005.

۶- نتیجه‌گیری

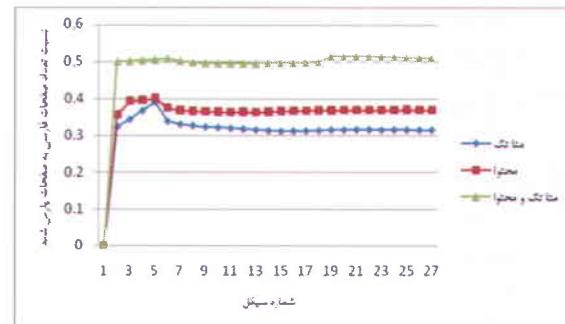
در این مقاله خزشگر فارسی با هدف بهبود خرش مستندات فارسی وب ارائه شد. با توجه به ارزیابی‌های صورت گرفته، مشخص شد بخش تشخیص زبان خزشگر فارسی با دقت بالایی عمل می‌کند. همچنین با تاثیر نتایج بخش شناسائی زبان بر اولویت‌دهی به واکنشی صفحات مرتبط، ضمن اینکه صفحات فارسی با دقت بالایی جمع‌آوری می‌شوند، امکان پوشش بیشتر صفحات فارسی نیز فراهم می‌شود. در واقع خزشگر فارسی با ایجاد امکان جمع‌آوری صفحات فارسی از سایر دامنه‌های وب، میزان پوشش صفحات فارسی را افزایش می‌دهد. ویژگی دیگر خزشگر ارائه شده سرعت جمع‌آوری صفحات فارسی می‌باشد. در واقع خزشگر فارسی دارای این قابلیت است که در بازه زمانی کوتاه‌تری تعداد صفحات مرتبط بیشتری جمع‌آوری نماید. بنابراین در مجموع استفاده از خزشگر فارسی منجر به بهبود خرش مستندات فارسی می‌گردد. همچنین در آزمایشات صورت پذیرفته تاثیر انتخاب دسته‌های متفاوت از URL‌های اولیه بر سرعت و پوشش خرش در سیکل‌های اولیه بررسی و تحلیل شده است. نتایج حاصله نشان‌دهنده آن است که در صورتیکه URL‌های اولیه حاوی پیوندهای خروجی مناسبی به سایر صفحات فارسی باشند، در مراحل اولیه نیز سازوکار وزندگی صفحات فارسی نتیجه مطلوبی برای خرش دربردارد. در مجموع در خزشگر فارسی بیشنهادی انتخاب URL‌های اولیه فارسی مناسب به تعداد زیاد، منجر به ایجاد شرایط اولیه مطلوبی هسته پوشش هرچه بیشتر صفحات فارسی وب می‌گردد. همچنین کارائی خزشگر فارسی بیشنهادی با روش‌های مبتنی بر محتوا، مبتنی بر فرابرچسب یا ترکیب این دو مورد بررسی قرار گرفته و عملکرد آنها مقایسه شده است. این آزمایش نشان می‌دهد که با ترکیب ویژگیها در مولفه تشخیص زبان عملکرد خزشگر به میزان قابل توجه ای بهبود یافته است. با توجه به بهبود جمع‌آوری اطلاعات فارسی وب، خزشگر فارسی دارای کاربردهای مختلفی در سامانه‌های جستجو و بازیابی اطلاعات است. به عنوان نمونه می‌توان از آن در موتورهای جستجو و بازیابی اطلاعات فارسی استفاده کرد. در موتورهای جستجوی فارسی، علاوه بر تمرکز بر مستندات فارسی، سرعت بالای این خزشگر در جمع‌آوری مستندات، بروز رسانی اطلاعات خرش شده را تسریع می‌بخشد. این امر در ارزیابی محتوای فارسی وب نیز صادق بوده و امکان استفاده از این ابزار جهت ارزیابی محتوای فارسی وب بسیار مفید خواهد بود.

این خزشگر هنوز امکان استفاده از هستان‌شناسی و یا گنجیه زبان فارسی را نداشته و امکان تشخیص ابعاد تخصصی و موضوعی صفحات فارسی را ندارد. از جمله فعالیت‌هایی که در کارهای آبتدی این تحقیق به آن خواهیم پرداخت، توسعه خزشگر موضوعی زبان فارسی می‌باشد. این خزشگر قادر خواهد بود که در ضمن تشخیص زبان صفحات موضوع مرتبط با عنوان پرس و جو را نیز تشخیص داده و جمع‌آوری کند. به این منظور اولین سازوکار مورد نیاز شناسایی صفحات فارسی وب می‌باشد که در این مقاله ارائه شده است. بعد از شناسائی زبان صفحه، در مرحله بعد می‌بایست کلمات بی‌محتوا و اضافه را دور ریخته و محتویات باقیمانده صفحه را با منابع زبانی موجود نظیر هستان‌شناس یا گنجینه و ازگان مقایسه شود. در صورت مرتبط بودن صفحه با موضوع مورد نظر پیوندهای موجود در آن



نمودار ۶- تعداد صفحات فارسی خرش شده با بکارگیری شاخصهای مختلف خرش زبانی

مطلوب با نمودار ۶ نرخ جمع‌آوری صفحات فارسی مبتنی بر اطلاعات فرابرچسب در مقایسه با بکارگیری سایر شاخصها از روند رشد کمتری برخوردار است. نتیجه حاصل از بکارگیری این شاخص مبتنی بر این واقعیت است که درصد قابل توجهی از صفحات فارسی فاقد فرابرچسب زبانی می‌باشند. استفاده از شاخص محتوا در منحنی دیگر این نمودار نشان‌دهنده بهبود نرخ جمع‌آوری صفحات فارسی نسبت به حالت قبل است. این موضوع نشان‌دهنده آن است که استفاده از ایستوازده‌ها می‌تواند منجر به شناسایی درصد بیشتری از صفحات فارسی شود.

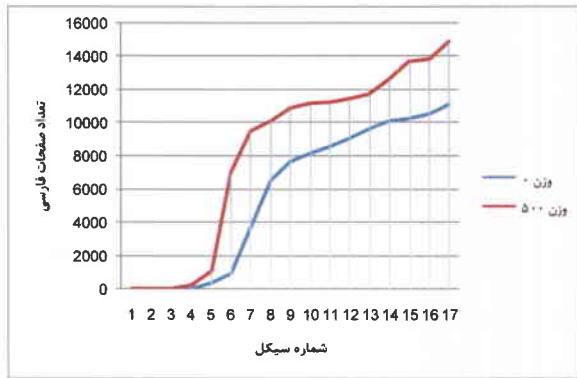


نمودار ۷- نسبت تعداد صفحات فارسی به صفحات پارس شده با بکارگیری شاخصهای مختلف خرش زبانی

اگر منحنی مربوط به بکارگیری ترکیبی از شاخصهای فرابرچسب و محتوا را در نمودار ۶ ملاحظه کنید نشان‌دهنده آن است که استفاده از ترکیبی از این دو شاخص تاثیر بسزایی بر افزایش نرخ صفحات جمع‌آوری شده در سیکلهای مختلف خرش دارد. از جمله دلایل این بهبود آن است که در این حالت صفحات فارسی که در روش مبتنی بر محتوا به دلیل غنی نبودن محتوا یا ذخیره‌سازی بر اساس کاراکترستی غیر از 1256-UTF-8 شناسایی نشده بودند، مورد شناسایی قرار می‌گیرند. همچنین نمودار ۷ نشان‌دهنده تاثیر شاخصهای مختلف خرش بر پوشش بیشتر صفحات فارسی می‌باشد. مطابق این نمودار استفاده از ترکیبی از شاخصهای خرش زبانی منجر به افزایش نسبت صفحات فارسی جمع‌آوری شده به کل صفحات خرش شده می‌گردد.



فارسی مورد استفاده در این وضعیت، منجر به دسترسی به سایتها با درجه پیوند خروجی قابل توجهی می‌شوند، این امر منجر شده که منحنی‌های موجود در این نمودار نسبت به منحنی‌های موجود در نمودار ۱ خصوصاً از سیکل ۶ به بعد روند رشد بسیار بیشتری داشته باشند. نکته‌ای که در رابطه با نمودار ۵ وجود دارد آن است که در سیکلهای اولیه روند رشد نمودار کندر از نمودار ۱ است که دلیل این موضوع را می‌توان در تعداد URL‌های اولیه مورد استفاده ذکر نمود.



نمودار ۵- تعداد صفحات فارسی خوش شده در سیکلهای مختلف بر اساس افزایش وزن

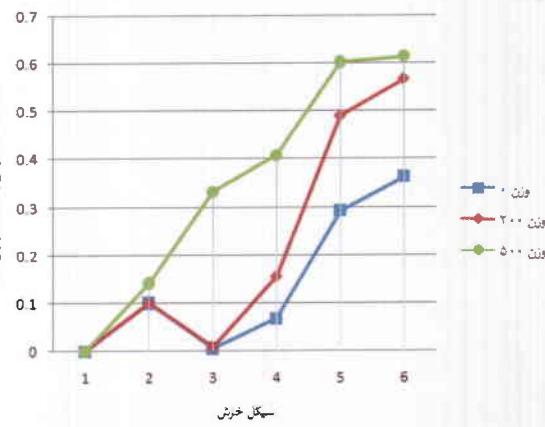
نکته مهم دیگری که در رابطه با نمودارهای ارائه شده قابل ذکر است آن است که ماهیت متغیر خوش مانع از رشد خطی نمودارها می‌گردد. در واقع تعداد صفحات فارسی واکنش شده در هر سیکل خوش بسته به عوامل مختلفی از جمله درصد پیوندهای واکنشی شده در سیکل مقابله دارد و همین موضوع موجب می‌شود که شب نمودار در هر سیکل دچار قدری افزایش یا کاهش شود. از طرفی با توجه به گستره شدن گراف جستجو همزمان با روند پیشرفت خوش در مجموع تعداد صفحات فارسی یافته شده افزایش می‌یابد.

جهت بررسی سرعت جمع‌آوری صفحات مرتبط توسط خوشگر فارسی، تعداد صفحات فارسی شناسائی شده توسط واير و خوشگر فارسی در دو اجرای مختلف مورد بررسی قرار گرفتند. نتایج اجرا با استفاده از خوشگر واير نشان داد که در انتهای سیکل ۳۲ از مجموع ۲۸۷۸۷ صفحه پارس شده، تعداد ۴۵۲۸۴ صفحه فارسی شناسائی شد. در حالیکه با استفاده از خوشگر فارسی در انتهای سیکل ۱۶ تعداد ۶۰۰۰ صفحه فارسی از بین صفحه شناسائی شدند. از نتایج حاصله می‌توان استنباط نمود که خوشگر ارائه شده جهت جمع‌آوری صفحات فارسی دارای سرعتی حداقل ۲ برابر خوشگر واير است. بنابراین خوشگر فارسی می‌تواند با سرعت بسیار بیشتری نسبت به خوشگر واير صفحات فارسی را جمع‌آوری نماید.

۴) بررسی تاثیر بکارگیری ویژگیهای مختلف صفحات وب بر کارائی خوش: جهت نشان دادن تاثیر استفاده از ویژگیهای مختلف صفحات وب بر کارائی خوش، خوشگر با بکارگیری URL‌های اولیه دسته دوم با وزن ۰ و به مدت ۲۵ دور اجرا گردیده است. شاخصهای زبانی مورد استفاده فرا برچسب، محتوا و ترکیبی از این دو شاخص بوده و نتایج در نمودارهای ۶ و ۷ ارائه شده است.

۲) نتایج اجرای خوشگر با استفاده از URL‌های دسته دوم:

نتایج اجرا با URL‌های دسته دوم در نمودار ۴ نشان‌دهنده آن است که نسبت صفحات فارسی خوش شده به کل صفحات پارس شده از روند رشد بسیار بیشتری نسبت به نمودار مشابه (نمودار ۲) مربوط به اجرای با URL‌های دسته اول برخوردار است. این مطلب تأثیر انتخاب URL‌های اولیه در نحوه خوش وب را نشان می‌دهد. به عنوان مثال در شرایطی که خوشگر با وزن ۵۰۰ برای صفحات فارسی اجرا گردیده است، در سیکل پنجم از نمودار ۲ نسبت صفحات فارسی به صفحات پارس شده حدود ۰.۲ می‌باشد در حالیکه این نسبت برای شرایط مشابه در نمودار ۴ حدود ۰.۶ می‌باشد.



نمودار ۶- نسبت تعداد صفحات فارسی به صفحات پارس شده براساس افزایش وزن در سیکلهای مختلف

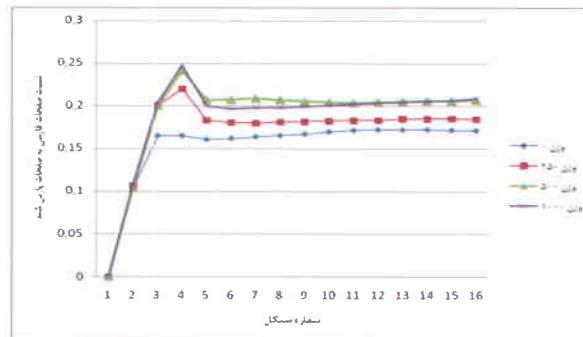
نمودار ۴- نسبت تعداد صفحات فارسی به صفحات پارس شده براساس افزایش وزن در سیکلهای مختلف مطابق نمودار ۴ خوشگر با وزن ۵۰۰ برای صفحات فارسی، از همان سیکلهای اولیه در مسیر صحیح خوش و یا در مسیر مناسب جهت خوش صفحات فارسی وب قرار گرفته است و همچنین از روند رشد بهتری نسبت به منحنی مشابه در نمودار ۲ برخوردار است. این امر بیانگر تأثیر انتخاب URL‌های اولیه بر روند و نتایج خوش می‌باشد. لازم به یادآوری است که نتایج ارائه شده در نمودار ۴ بیانگر شرایطی است که URL‌های اولیه از چندین دایرکتوری فارسی DMOZ انتخاب شده‌اند. همچنین در نمودار ۴ منحنی با وزن ۰ در سیکلهای اولیه خوش از روند رشد کمتری نسبت به منحنی با وزن ۲۰۰ در سیکلهای اولیه خوش از روند رشد کمتری نسبت به منحنی با وزن ۵۰۰ برخوردار است که از جمله دلایل این امر می‌توان به تأثیر پارامترهای دیگر در اولویت‌دهی صفحات در مراحل اولیه خوش اشاره نمود. پارامتر عمق از جمله پارامترهایی است که در اولویت‌دهی به صفحات نقش دارد که با پیمایش عمق بیشتری از صفحات تأثیر کمتری در دنبال نمودن آنها دارد. در واقع می‌توان گفت با توجه به اینکه در مراحل اولیه تعداد لینکهای غیرفارسی بیشتری در دایرکتوری وجود دارد، با تأثیر پارامتر عمق منحنی‌های موجود در نمودار برای وزنهای ۰ و ۲۰۰ در سیکلهای اولیه روند رشد مناسبی ندارند.

۳) نتایج اجرای خوشگر با استفاده از URL‌های دسته سوم:

در این اجرا بازهم URL‌های انتخاب شده باعث افزایش چشم‌گیر تعداد صفحات فارسی خوش شده نسبت به صفحات فارسی خوش شده در اجرای با URL‌های دسته اول می‌باشد. در نمودار ۵ نیز تأثیر انتخاب URL‌های اولیه در روند رشد منحنی‌ها دیده می‌شود. با توجه به اینکه URL‌های



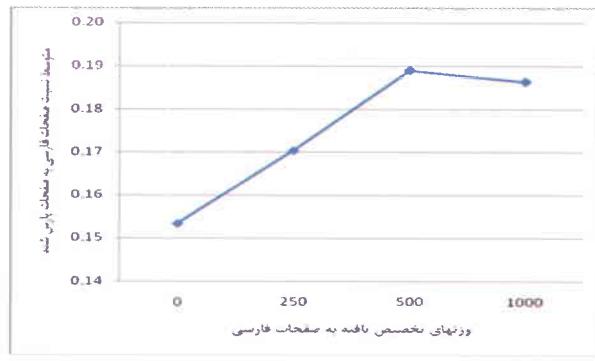
تعیین حداقل مقدار پارامتر وزنده‌ی به صفحات فارسی وابسته به شرایط اجرا و پارامترهای مختلفی است که در وزن دهنی به صفحه مؤثر هستند که از جمله این پارامترها می‌توان به پارامترهای عمق و رندوم اشاره نمود. بنابراین بسته به تطبیقات اولیه خزشگر و شرایط اجرای مختلف، حد اشباع رسیدن مقدار پارامتر وزن دهنی به صفحات متعلق به زبان متفاوت می‌باشد که در این مقاله این پارامتر با مقدار ۵۰۰ به حد اشباع رسیده است.



نمودار ۲- نسبت تعداد صفحات فارسی به صفحات پارس شده براساس افزایش وزن در سیکل‌های مختلف

همچنین نمودار ۲ نسبت صفحات فارسی به صفحات پارس شده را نشان می‌دهد. همانطور که از نمودار مشخص است، با افزایش وزن صفحات فارسی، میزان صفحات فارسی خزش شده به کل صفحات پارس شده نیز افزایش یافته است. این موضوع می‌تواند نشان‌دهنده افزایش پوشش صفحات فارسی وب باشد. همانطور که نمودارهای ۱ و ۲ نشان می‌دهند در اجرای خزشگر با URLهای دسته‌اول، منحنی‌ها در سیکل‌های اولیه از روند رشد نسبتاً یکسانی برخوردار هستند و پس از آن نحوه تاثیرگذاری وزن صفحات فارسی در بهبود خزش مستندات فارسی دیده می‌شود. دلیل این امر را می‌توان در تاثیر انتخاب URLهای اولیه در نتایج خزش دانست. در واقع با توجه به اینکه در این حالت URLهای ترکیبی از URLهای مربوط به دامنه کشورهای مختلف بوده و از نظر تعداد نیز قابل توجه هستند، تاثیر وزنده‌ی به صفحات بعد از پیمایش عمق مشخصی از صفحات وب دیده می‌شود.

نمودار ۳ متوسط نرخ پوشش صفحات فارسی را به ازای وزن‌های مختلف نشان می‌دهد. همانگونه که در نمودار ملاحظه می‌شود، از نقطه ۵۰۰ به بعد افزایش وزن صفحات تاثیر چندانی بر متوسط نرخ صفحات جمع‌آوری شده فارسی نسبت به کل صفحات پارس شده نداشته است.



نمودار ۳- متوسط پوشش صفحات فارسی به ازای وزن‌های مختلف

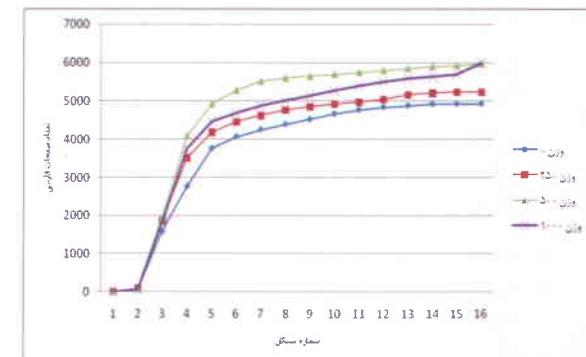
www.sku.ac.ir/fa/academic/members/golestanian/main.htm
www.honar.ac.ir/farhangestan/structure.htm
www.irandoc.ac.ir/library/copy.htm
www.iranculture.org/contactus/index.php
www.bawar.ir/

همانگونه که ملاحظه می‌شود URLهای ارائه شده از دامنه‌های مختلف .ir و .org جمع‌آوری شده‌اند که البته امکان شناسایی دامنه‌های دیگری مانند .net و .gov نیز وجود دارد. در حالیکه در کارهای مرتبطی که از خزشگر واب استفاده شده مانند [29]، ملاک شناسایی صفحات فارسی تنها تعلق آنها به دامنه ir. بوده است. این عامل موجب می‌شود که بخش قابل توجهی از صفحات فارسی وب نادیده گرفته شود. به عنوان مثال در نمونه بررسی شده تعداد ۴۸ پیوند متعلق به دامنه .com و ۵۰ پیوند متعلق به دامنه ir. بود. همچنین با توجه به اینکه خزشگر در حین عملیات خزش به صفحات فارسی وزن بیشتری می‌دهد، در مراحل بعدی خزش، لینکهای موجود در صفحات فارسی وزن بیشتری دریافت می‌کنند. بنابراین انتظار می‌رود با تاثیر این وزن در سیستم وزنده‌ی خزشگر واب، نرخ جمع‌آوری صفحات فارسی افزایش یابد که برای سه دسته URLهای فوق به شکل زیر مورد بررسی قرار گرفته است.

(۱) نتایج اجرای خزشگر با استفاده از URLهای دسته اول

نتایج اجرای خزشگر با دسته اول URLهای اولیه به ازای وزن‌های مختلف برای صفحات فارسی در نمودار ۱ نشان داده شده است. ستون افقی نشان دهنده شماره سیکل‌های خزش و ستون عمودی نشان دهنده میزان صفحات فارسی جمع‌آوری شده می‌باشد.

همانگونه که ملاحظه می‌شود در حالتیکه وزن صفحات فارسی با صفر مقداردهی شده است، نمودار رشد منحنی از روند کندتری برخوردار است و با افزایش وزن صفحات فارسی روند رشد نمودار بهبود می‌یابد. نکته قابل توجه در این نمودار این است که افزایش وزن صفحات متعلق به زبان تا حد مشخصی می‌تواند منجر به بهبود نتایج خزش گردد و زمانیکه مقدار آن از حد مشخصی تجاوز کند، با توجه به اینکه تاثیر آن در مقابل تاثیر سایر پارامترها در اولویت‌دهی به صفحات به صورت قابل توجهی بیشتر شده و به حد اشباع رسیده است، این افزایش وزن در نتایج نهایی خزش تغییر قابل توجهی ایجاد نخواهد کرد. به عنوان مثال همانگونه که در نمودار ملاحظه می‌گردد با افزایش وزن صفحات فارسی به مقدار ۱۰۰۰ نتایج بدست آمده در سیکل‌های مختلف خزش نسبت به وضعیتی است که خزشگر با وزنده‌ی شده است، دستخوش تغییراتی چندانی نگرددیده است.



نمودار ۱- تعداد صفحات فارسی خزش شده در سیکل‌های مختلف بر اساس افزایش وزن



■ آزمون براساس URLهای دسته دوم:
برای هر تست واپر تعداد ۶ دور سیکل خوش در نظر گرفته شده است. که هر تست حدود ۱:۳۰ ساعت زمان به خود اختصاص می‌دهد. مشخصات سه اجرای انجام گرفته به ازاء مقدایر مختلف برای پارامتر وزن صفحات فارسی در جدول ۳ بطور خلاصه آمده است.

جدول-۳ مشخصات اجراهای انجام گرفته برای ارزیابی خزشگر فارسی به ازاء دسته دوم seed

اجرای سوم	اجرای دوم	اجرای اول	وزن صفحات فارسی
۵۰۰	۲۵۰	-	وزن صفحات فارسی
۶	۶	۶	تعداد دور خوش
۱:۳۰ ساعت	۱:۳۰ ساعت	۱:۳۰ ساعت	مدت زمان اجرا
۱۵ فوریه	۱۵ فوریه	۱۵ فوریه	تاریخ
۱۰۰۶۶	۱۵۵۰	۱۲۷۵۰	تعداد صفحات پارس شده
۵۹۶۸	۸۰۰۸	۴۴۸۲	تعداد صفحات فارسی

■ آزمون براساس URLهای دسته سوم:
برای هر تست واپر تعداد ۱۶ دور سیکل خوش در نظر گرفته شده است. که هر تست حدود ۳ ساعت زمان به خود اختصاص می‌دهد. مشخصات دو اجرای انجام گرفته به ازاء وزن های ۰ و ۵۰۰ برای پارامتر وزن صفحات در جدول ۴ بطور خلاصه آمده است.

جدول-۴ مشخصات اجراهای انجام گرفته برای ارزیابی خزشگر فارسی به ازاء دسته سوم seed

اجرای دوم	اجرای اول	وزن صفحات فارسی
۵۰۰	-	وزن صفحات فارسی
۱۶	۱۶	تعداد دور خوش
۳ ساعت	۳ ساعت	مدت زمان اجرا
۱۹ فوریه	۱۹ فوریه	تاریخ
۲۲۷۸۳	۲۴۸۰۹	تعداد صفحات پارس
۱۴۸۸۶	۱۱۰۹۵	تعداد صفحات فارسی

جدول ۵ تعداد صفحات وب جمع آوری شده، درصد صفحات تکراری جمع-آوری شده و همچنین درصد صفحات بُویا و ایستای خوش شده را نشان می‌دهد. لازم به ذکر است صفحات بُویا به صفحاتی گفته می‌شود که توسط تکنولوژی‌های برنامه نویسی ایجاد می‌شوند و دارای پسوندهایی مانند .asp, .php, .aspx, .html. هستند اما صفحات ایستا دارای پسوندهایی مانند .html, .htm. هستند. پس از اجرای خزشگر در اجرای دوم از جدول ۲، در کل ۳۶۵۶۸ صفحه وب توسط خزشگر پردازش شد که از بین ۲۸۲۹۳ صفحه پارس شده تعداد ۵۲۳۳ صفحه فارسی بود.

جدول-۵ آمار صفحات وب خوش شده

تعداد کل صفحات	صفحه ۳۶,۵۶۸
تعداد صفحات یکتا	صفحه ۳۵,۳۲۷
تعداد صفحات تکراری	صفحه ۱,۴۴۱
تعداد صفحات استتا	صفحه ۲۲,۴۰۲
تعداد صفحات بُویا	صفحه ۱۴,۱۶۶

برای بررسی صحت عملکرد بخش تشخیص زبان فارسی، تعداد ۱۰۰ لینک از بین ۵۲۳۳ لینک فارسی شناسائی شده در اجرای دوم از جدول ۲ بصورت دستی بررسی شدند که همه به درستی شناسائی شده بودند. بخش کوچکی از لینکهایی که صفحات آنها به عنوان فارسی شناسائی شدند عبارتند از: www.farsnews.com/newstext.php?nn=8711290500 www.just.ac.ir/printme-1.3346.6386.fa.html

مقایسه با خزشگر واپر بدون تأثیر زبان، وزن صفحات فارسی را از صفر تا ۵۰۰ تغییر می‌دهیم. در حالتی که خزشگر فارسی با وزن صفر در واقع نشان‌دهنده عملکرد واپر بدون تأثیر اجزاء فارسی اضافه شده در این تحقیق می‌باشد. در ضمن برای پارامتر عمق وزن ۹۵، پارامتر رندوم وزن ۵ و برای پارامترها وزن ۰ در نظر گرفته شده است.

پارامتر دیگر حداکثر تعداد صفحات و سایتها قابل خوش می‌باشد که جهت اجرای خوش حداکثر ۶۰۰۰ صفحه به ازاء هر سایت و حداکثر ۶۰۰۰ سایت برای خوش در نظر گرفته شده است. یکی از پارامترهای مهم دیگری که جهت ارزیابی عملکرد خزشگر فارسی مقداردهی شده است عبارت است از دامنه‌هایی که خزشگر امکان خوش صفحات آنها را دارد. با توجه به اینکه خزشگر واپر جهت خوش دامنه کشورها طراحی شده است در کارهای انجام شده عموماً به دامنه اصلی هر کشور محدود می‌شود [26]. در برخی تحقیقات صورت پذیرفته نظر [28] به این موضوع پرداخته شده است که محدود بودن خزشگر به دامنه کشور منجر به از دست دادن بخش زیادی از صفحات مربوط به دامنه یک کشور می‌شود. با توجه به اینکه از اهداف این تحقیق ایجاد امکان پوشش صفحات فارسی موجود در سایر دامنه‌های وب می‌باشد، در پیکربندی خزشگر علاوه بر دامنه وب ایران (.ir)، دامنه‌های .com، .org، .net، .gov، نیز در نظر گرفته شده‌اند تا امکان بررسی عملکرد خزشگر در پوشش صفحات فارسی موجود در سایر دامنه‌ها فراهم گردد. همچنین جهت تعیین میزان تمایل خزشگر فارسی به سمت دامنه وب ایران، دامنه‌های کشورهای ژاپن (jp)، آلمان (de)، آرژانتین (ar)، انگلستان (uk) و آلمان (de). نیز در فایل پیکربندی مقداردهی شده است.

۳-۲-۵ ارزیابی خزشگر و نتایج به دست آمده

جهت ارزیابی عملکرد خزشگر سیستم عامل Linux Fedora Core 8 مورد استفاده قرار گرفت. این سیستم عامل روی سرور Intel با حجم ۳ ترابایت دیسک سخت و حافظه ۴ گیگابایت نصب شد.

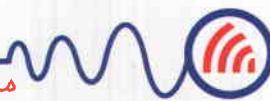
با توجه به اینکه عملکرد خزشگر تا حدود زیادی به لیست URLهای اولیه وابسته می‌شود، جهت ارزیابی عملکرد خزشگر فارسی از سه دسته URLهای اولیه استفاده شده است. در ادامه توضیح هریک از لیست URLهای اولیه آمده است.

■ آزمون براساس URLهای دسته اول:

برای هر تست واپر تعداد ۱۶ دور سیکل خوش در نظر گرفته شده است. که هر تست حدود ۳ ساعت زمان به خود اختصاص می‌دهد. مشخصات سه اجرای انجام گرفته به ازاء مقدایر مختلف برای پارامتر وزن صفحات فارسی در جدول ۲ بطور خلاصه آمده است.

جدول-۲ مشخصات اجراهای انجام گرفته برای ارزیابی خزشگر فارسی به ازاء دسته اول seed

اجرای سوم	اجرای دوم	اجرای اول	وزن صفحات فارسی
۵۰۰	۲۵۰	-	وزن صفحات فارسی
۱۶	۱۶	۱۶	تعداد دور خوش
۳ ساعت	۲:۳۰ ساعت	۲:۳۰ ساعت	مدت زمان اجرا
۷ فوریه	۷ فوریه	۷ فوریه	تاریخ
۲۸۸۳۱	۲۸۲۹۳	۲۸۷۷۷	تعداد صفحات پارس شده
۵۹۶۸	۵۲۳۳	۴۹۲۶	تعداد صفحات فارسی



- اضافه کردن یک ساختمندانه به ابردادهای مربوط به صفحات در بخش ذخیره سازی ابرداده.
 - تائیر وزندهی به صفحات فارسی در فرمول وزن دهی جهت اولویتدهی و اواکشی این صفحات. اضافه نمودن پارامتر وزندهی در فایل پیکربندی خرشنگر.
 - مقدار دهی به ضرایب وزندهی بر اساس معیارهای مختلف.

۲-۵ ارزیابی

جهت ارزیابی عملکرد خزشگر فارسی معیارهای صحت، پوشش و سرعت خزش مستندات فارسی وب مورد بررسی قرار گرفته‌اند. با توجه به اهمیت خش شناسائی زبان ابتدا باید صحت عملکرد آن در شناسائی صفحات فارسی تعیین شود که جهت ارزیابی آن صفحات وب شناسائی شده توسط خزشگر فارسی به صورت دستی مورد ارزیابی قرار گرفتند. معیار مهم دیگر معیار پوشش می‌باشد. از جمله شخصهای ارزیابی معیار پوشش میزان صفحات مرتبط خزش شده و شاخص دیگر امکان شناسائی و جمع‌آوری صفحات فارسی در سایر دامنه‌ها می‌باشد. سرعت خزش مستندات فارسی وب معیار سومی است که جهت ارزیابی عملکرد خزشگر مد نظر قرار گرفته است. در واقع به کمک این فاکتور زمان موردنیاز جهت جمع‌آوری صفحات فارسی از وب را می‌توان بر سیلکله‌های خزش مشخص می‌کنیم. در ادامه شرایط اولیه، نحوه آماده‌سازی، شرایط اجرا و ارزیابی نتایج خزشگر فارسی، آنچه خواهد شد.

۱-۲-۵ انتخاب URL های اولیه

جهت ارزیابی عملکرد خزشگر لازم است آن را تحت شرایط اولیه متفاوتی
رزیابی کرد؛ زیرا شرایط اولیه می‌تواند تاثیر زیادی در عملکرد خزشگر
داشته باشد. به همین دلیل چندین دسته URL^{۴۴} اولیه برای خزشگر در نظر
گرفته شده است. منبع انتخاب این URL‌ها دایرکتوری DMOZ^{۴۵} می‌باشد
که یک دایرکتوری مرجع حاوی دسته‌بندی‌های متنوع از آدرس و وبسایت-
های مربوط به کشورهای مختلف می‌باشد و بصورت دستی بروز رسانی

- URL‌های اولیه دسته اول: ارزیابی عملکرد خرشگر مستلزم ایجاد شرایط اولیه‌ای است که تمایل خرشگر به سمت صفحات فارسی را در مقایسه با صفحات زبانهای دیگر نشان دهد. به همین منظور URL‌های اولیه از پنج زبان مختلف ژاپنی، آلمانی، اسپانیائی، انگلیسی و فارسی و به ازای هر زبان ۱۵۰ URL انتخاب شد.
 - URL‌های اولیه دسته دوم: چندین دایرکتوری فارسی از DMOZ
 - URL‌های اولیه دسته سوم: ترکیبی از URL‌های فارسی و انگلیسی

۲-۲ پیکربندی خزشگر

خوشگر وایر دارای یک فایل پیکربندی تحت عنوان WIRE.conf می باشد که در این فایل تنظیمات اولیه خوشگر مقداردهی می شود. جهت وزندهی به صفحات فارسی به این فایل پارامتری به نام CONF_MANAGER_SCORE_PERSIAN_WEIGHT اضافه شده است که می توان جهت تعیین مقدار وزن مناسب برای صفحات فارسی آن را با مقادیر مختلف، ایش دهنده، نمود. جهت بررسی عملکرد خوشگر فارسی در

²⁴ DMOZ Open Directory: www.dmoz.org

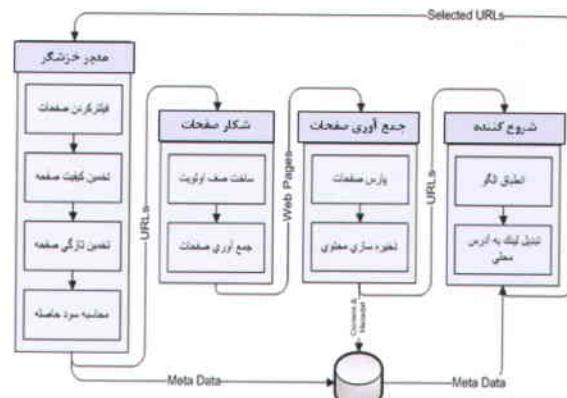
تاثیرگذاری زبان مستندات در وزن صفحات می‌باشد کلیه وزنها و به ویژه وزن زبان صفحات در فایل پکربندی خرشگر تنظیم گردد.

۵-پیاده‌سازی و ارزیابی خزشگر زبانی پیشنهادی

۱-۵ پیاده‌سازی

در این بخش شرح پیاده‌سازی چارچوب پیشنهادی در قالب یک خزشگر زبان فارسی آمده است. با توجه به اینکه طراحی یک خزشگر جهت خزش در ابعاد وسیع، دارای بی‌محدودیتها و چالشهای خاص خود می‌باشد، بنابراین بکارگیری یک خزشگر متن‌باز راه حل مناسبتری به نظر می‌رسد. در این راستا جهت پیاده‌سازی خزشگر زبان فارسی، از خزشگر متن‌باز واپر استفاده شد. این خزشگر از بین خزشگرهای متن‌باز موجود به دلیل دارا بودن ویژگی‌هایی نظیر مبنای علمی قوی، وجود مستندات کافی و بکارگیری آن در تحقیقات مختلف مانند [19]، انتخاب شده است.

معماری این خزشگر شامل چهار بخش اصلی مدیریت^{۲۰}، شکار^{۲۱}، جمع آوری^{۲۲} و شروع کننده^{۲۳} است که در یک سیکل بصورت گردشی اجرا می‌شوند و در هر سیکل نتایج حاصل از سیکل قبلی استفاده می‌شود. به عبارت دیگر بعد از هر سیکل مجموعه‌ای از پیوندها استخراج شده و به عنوان URL‌های قابل واکنش برای سیکل بعدی استفاده می‌شوند. ارتباط بخش‌های مختلف خزشگر و مولفه‌های تشکیل‌دهنده آنها در شکل ۳ نمایش داده شده است. با توجه به اینکه در خزشگر وایر ارتباط بین بخش‌های مختلف خزشگر توسط منبع ذخیره‌سازی متاداده برقرار می‌شود، بعد از تشخیص زبان صفحات مداری در بخش منبع ذخیره‌سازی متاداده ذخیره می‌شود تا این موضوع به اطلاع سایر بخش‌های خزشگر نیز برسد.



شکا. ۳- معماری خشک واب

در طرح پیشنهادی پس از تشخیص زبان صفحات در بخش جمع‌آوری صفحات، با توجه به معماری خرشگر وایر، اولویت دهی به صفحات فارسی از طریق تغییر فرمول وزندهی در بخش مدیریت صورت می‌یابد که در این اسکتا فعال‌سازی، نسبت بذوق فته است:

ایست واژه های^{۱۹} مورد استفاده در این مقاله از یک فهرست اولیه شامل ۱۰۰۰ ایست واژه از سایت پیکر همشهری [30] آزمایشگاه DBRG دانشگاه تهران واکنش شده و طی چندین مرحله پردازش نهایی شده است. در مرحله اول از پردازش بخشی از کلمات و واژگانی که تداول کثیر داشتند حذف گردیده اند. در مرحله بعد با توجه به اینکه زبان فارسی و عربی دارای اشتراک لغوی زیادی هستند بخشی از ایست واژه ها که بین دو زبان مشترک بودند حذف شده اند. سپس محتوای صفحاتی که به عنوان فارسی شناسایی شده بودند بررسی شده و در صورتیکه صفحه مذکور متعلق به زبان عربی بود و مبتنی بر ایست واژه های فراهم شده شناسایی شده باشد، ایست واژه های مذکور از لیست ایست واژه های مورد استفاده در این مقاله حذف گردیده است. بطور کلی جهت تسریع در شناسایی زبان صفحات سعی شده است که تعداد مناسبی ایست واژه که دارای فرکانس تکرار زیادی در صفحات فارسی باشند انتخاب شوند. لازم به ذکر است که لیست ایست واژه های مذکور با فرمت UTF-8 ذخیره سازی گردیده است. ساز و کار شناسایی صفحات فارسی وب در شکل ۲ ارائه گردیده است.

۴-۴ مولفه وزن دهی به صفحات

در طرح پیشنهادی فرض اولیه مبتنی بر این است که لینکهای موجود در صفحات مرتبط یا فارسی زبان به صفحات فارسی اشاره دارند. بنابراین بعد از شناسایی صفحات مرتبط یا فارسی زبان باید لینکهای موجود در این صفحات به لیست لینکهای قابل دانلود اضافه شوند. جایگاه مولفه وزن دهی در خرشگر در شکل ۱ نشان داده شده است. همانطور که در شکل مذکور مشاهده می گردد، پس از واکنشی صفحات می بایست تشخیص داده شود که آنها به زبان فارسی هستند یا خیر. سپس لینکهای موجود در کل صفحات جمع آوری شده را استخراج کرده و قبل از اینکه لینکهای استخراج شده از صفحات را به لیست لینکهای قابل واکنشی اضافه کند، وزن مناسبی به آنها اختصاص داده شود. به طوری که لینکهای استخراج شده از صفحات فارسی با احتمال بیشتری خوش شوند. بنابراین برای تغییر سازوکار اولویت بندی صفحات می بایست سازوکار وزن دهی به صفحات متعلق به زبان را در فرمول وزن دهی در نظر گرفت. فرمول وزن دهی مورد استفاده به شرح زیر است:

$$\text{URL Weight} = \text{PAGERANK_WEIGHT} * \text{pagerank} + \text{WLSCORE_WEIGHT} * \text{pagerank} + \text{HITS_HUB_WEIGHT} * \text{hubrank} + \text{HITS_AUTHORITY_WEIGHT} * \text{authrank} + \text{SITERANK_WEIGHT} * \text{siterank} + \text{QUEUESIZE_WEIGHT} * \text{queuesize} + \text{DEPTH_WEIGHT} * \text{depth_score} + \text{STATIC_WEIGHT} * \text{static_score} + \text{RANDOM_WEIGHT} * \text{random_score} + \text{PERSIAN_WEIGHT} * \text{ispersian};$$

در فرمول وزن دهی صفحات، وزن رتبه بندی صفحات نظیر pagerank و siterank authrank hubrank، وزن مربوط به صفحات، عمق، ایستایی و رندم صفحات در نظر گرفته شده است که از فایل تنظیمات خرشگر خوانده می شوند. جهت دخیل نمودن زبان صفحات در فرمول وزن دهی صفحات ضریبی تحت عنوان PERSIAN_WEIGHT و معیاری تحت عنوان ispersian که فارسی بودن یا نبودن زبان صفحات را مشخص می کند، در نظر گرفته شده اند. بنابراین برای تنظیمات لازم در خصوص میزان

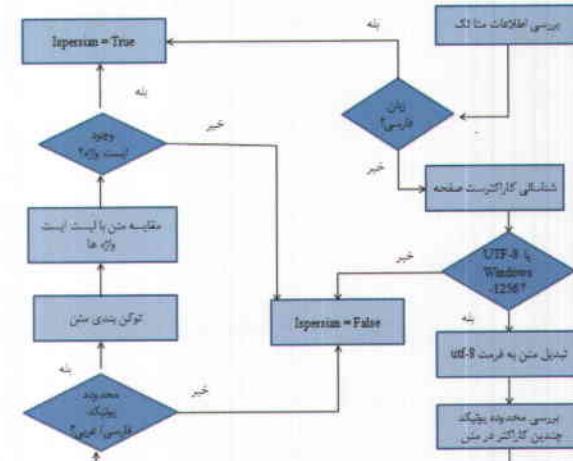
مستندات پس از تجزیه آنها در بخش پارس انجام می شود. ورودی مولفه شناسایی زبان محتوای صفحات وب است. در این بخش پس از بررسی زبان صفحات، آنها به عنوان فارسی یا غیرفارسی علامتگذاری می شوند. بخش های محتوای مورد استفاده در این طرح عبارتند از فرآبرچسب زبانی، کاراکترست و محتوای صفحات. عملیات تجزیه و پردازش محتوای صفحات در چند مرحله و به کمک شاخصهای زیر صورت می پذیرد:

۱. فرآبرچسب زبانی صفحه
۲. کاراکترست صفحات
۳. محدوده یونیکد صفحات
۴. ایست واژه ها

مولفه تشخیص زبان اول به بررسی فرآبرچسب زبانی صفحه می پردازد. در صورتیکه تگ زبانی صفحه فارسی باشد، صفحه به عنوان فارسی شناسایی می شود. اگر زبان صفحه به کمک اطلاعات فرآبرچسب شناسایی نشده باشد کاراکترست صفحه مورد بررسی قرار می گیرد. با توجه به اینکه کاراکترست صفحات فارسی عموما Windows-1256 یا utf-8 می باشد، می توان از این مشخصه به عنوان یک شاخص اولیه جهت شناسایی صفحات فارسی استفاده نمود. سپس محتوای صفحات با کاراکترست Windows-1256، به فرمت utf-8 تبدیل می شوند.

با توجه به اینکه ممکن است برخی از بخش های محتوای صفحه خالی از محتوا بوده و یا محتوا موجود در آنها به زبان فارسی نباشد، محدوده یونیکد کاراکتر موجود در متن مورد بررسی قرار می گیرد. در صورتیکه کاراکترها در محدوده یونیکد فارسی/عربی باشند، محتوای متنی توکن بندی شده و با ایست واژه های موردنظر مقایسه می شود. بر اساس وجود ایست واژه های فارسی صفحه به عنوان فارسی شناسایی می شود. برای توکن بندی متن از روش توکن بندی مبتنی بر مولفه موجود با معیار کاراکترهای جداگانه نظری^{۲۰}، ^{۲۱}، ^{۲۲}، ^{۲۳} و ^{۲۴} استفاده شده و کاراکترهای جداگانه فارسی نظیر ^{۲۵}، ^{۲۶} نیز در نظر گرفته شده است.

در نهایت جهت مشخص نمودن صفحات فارسی شناسایی شده یک متغیر بولین به نام ispersian مورد استفاده قرار گرفته است که در متغیر reserved متادیتا صفحه ذخیره می گردد تا امکان دسترسی به آن توسط سایر بخش های خرشگر وجود داشته باشد.

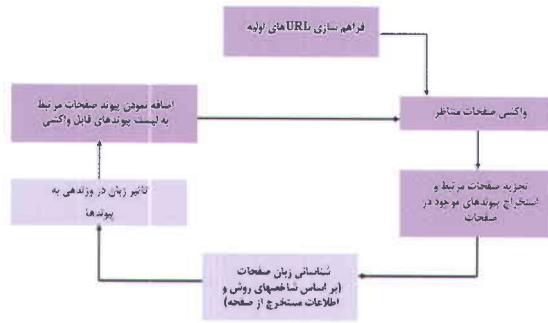


شکل ۲- سازوکار شناسایی زبان صفحات وب

^{۱۹} Stop Words

جدول ۱- مقایسه روش‌های خزش بر اساس ویژگی‌های خزش

استفاده از بلومنتر عملی جهت مواد ماصفحات نامنطبقه	اطلاعات مورد استفاده جهت بروز مرتبه بودن صفحات واکنش شده							اطلاعات مورد استفاده جهت بروز مرتبه بودن صفحات مورد آثاره بیوندها			شاخهای خرش روشنی خرش متوجه
	اطلاعات بازخورده	عنوان	كلمات کلیدی	شمارش بیوندها	مقایسه صفحه با صفحه پرس و جو	هستان شناس با کنجد و ازگان	ورانس	عن اطراف بیوندها	عن بیوند		
•			•				•			PageRank [14]	
•		•								Fish Search [9]	
•				•		•	•	•	•	Shark Search [10]	
•	•	•	•				•		•	Focused Crawling [4]	
•	•	•								Learnable Crawling [6]	
										Language Specific Web Crawling [15]	



شکل ۱- طرح کلی یک خزشگر زبانی

۱-۴- مولفه تشخیص زبان پیشنهادی

از آنجاییکه این مولفه در عملکرد صحیح خزشگر تائیر زیادی دارد، از مولفه‌های کلیدی خزشگر زبانی محسوب می‌شود. همانطور که در بخش فیل نیز اشاره شد، روش‌های مختلفی جهت شناسایی زبان صفحات وب وجود دارد. در این بخش با بکارگیری ترکیبی از شاخهای خزش زبانی روشی برای تشخیص زبان فارسی در خزشگرها ارائه گردیده است. با توجه به اینکه شناسایی زبان صفحات به صورت پویا صورت می‌پذیرد باید روش انتخاب شود که بتوان با سرعت بالایی زبان صفحات را شناسایی نمود. همچنین جهت شناسایی می‌توان از مشخصه‌هایی استفاده نمود که در بیشتر صفحات وب قابل دسترسی باشند. اطلاعات فراپرچسب به دلیل امکان تشخیص سریع زبان صفحه، به عنوان یک شاخص مورد استفاده در خزش زبانی پیشنهادی در نظر گرفته شده است، ولی از آنجاییکه صفحات فاقد فراپرچسب ربانی با این روش قابل تشخیص نمی‌باشند، بایراین شناسایی زبان صفحات بدون فراپرچسب نیاز به مکانیزم‌های مکملی دارد. دراین مقاله روش بکارگیری توکنها با کلمات منتخب زبان به عنوان روش مکمل برای شناسایی زبان صفحه پیشنهاد شده است. بخشی که خزشگرها محتوای صفحات را تجزیه می‌کنند پارس ناممده می‌شود. در طرح پیشنهادی، عملیات شناسایی زبان

بر اساس ویژگی‌های مستخرج از روش‌های خزش متوجه و نیازمندی‌های خزش زبانی، شاخهای مناسب در بردازه داده‌ها و سازوکار مناسب جهت شناسایی صفحات متعلق به زبان مورد نظر می‌باشد. شاخهای پایه مستخرج عبارتند از:

- اطلاعات فراپرچسب

- داده‌های محتوایی یا کلمات زبان

یکی از شاخهای مناسب جهت شناسایی زبان صفحات اطلاعات مستخرج از فراپرچسب‌ها می‌باشد. با توجه به اینکه فراپرچسب‌ها به سرعت قابل شناسایی هستند و در بسیاری از صفحات وب وجود دارند، به عنوان معیاری مناسب جهت شناسایی زبان محسوب می‌شوند. همچنین چون شناسایی زبان واسنگی زیادی به داده‌های زبانی دارد، استفاده از اطلاعات محتوایی صفحات شامل شاخهای متن بیوند و محتوای صفحات نیز روش مناسبی جهت شناسایی زبان صفحات می‌باشد. علاوه بر انتخاب شاخهای پایه جهت خزش زبانی، استفاده از شاخص بازخورده یا یادگیری نیز می‌تواند منجر به افزایش میزان بتوش خزش شود. در واقع می‌توان ترکیبی از شاخهای ذکر شده را برای خزش زبانی مورد استفاده فرارداد.

۴- چارچوبی برای طراحی یک خزشگر زبانی

در این بخش با توجه به سیاستگذاری واکنش صفحات وب، شاخهای و روش شناسایی زبان صفحات وب که در بخش‌های قبل شرح داده شد، روش برای خزش مبتنی بر زبان برای زبان فارسی پیشنهاد شده است. همانطوریکه در دیاگرام شکل ۱ آمده است، در خزشگرهای زبانی دو مولفه اصلی وجود دارد که در این مقاله آنها را مولفه تشخیص زبان و مولفه وزن دهنی به صفحات نامیده‌ایم. در واقع با استفاده از این دو مولفه یک خزشگر امکان اولویت دادن به صفحات با زبان خاص (در این مقاله زبان فارسی) را داشته و می‌تواند لینکهای مرتبط را شناسایی و دنبال کند.

در ادامه این مقاله، تحلیل و طراحی مولفه‌های مذکور برای یک خزشگر زبان فارسی آمده است.

» معیارهای شناسائی پیوندهای مرتبه

اطلاعاتی که به صورت متداول جهت شناسائی پیوندهای مرتبه مورد استفاده قرار می‌گیرد، شامل اطلاعات بررسی پیوند یا متن اطراف آن می‌باشد. همچنین در برخی از روش‌های مبتنی بر وراثت از امتیاز صفحه پدر نیز جهت امتیازدهی به پیوندهای موجود در آن استفاده می‌گردد. استفاده از متن پیوند جهت تعیین مرتبه بودن صفحاتی که به آن اشاره می‌کنند، منجر به تسريع عملیات شناسائی می‌شود. در مواردی که متن پیوند تهی باشد جهت شناسائی محتوای صفحات ارجاعی از متن اطراف پیوند کمک گرفته می‌شود. محدودیتی که در استفاده از این معیار وجود دارد آن است که گاهی شناسائی مرز اطلاعات مرتبه با پیوند دشوار است. در روش‌های وراثتی اهمیت صفحات در سطوح پائین تر با توجه به اهمیت صفحات در سطوح بالاتر افزایش یا کاهش می‌یابد [10,14]. مشکلی که در این گونه روشها وجود دارد انتشار خطای شناسائی صفحات مرتبه یا اولویت‌بندی صفحات، در یک مرحله به مراحل دیگر می‌باشد.

» معیارهای شناسائی صفحات مرتبه

مرتبه بودن محتوای صفحه مورد جستجو نیز می‌تواند مبتنی بر ویژگیهای متن یا پیوند و یا ترکیبی از هر دو ویژگی برسی گردد. در روش‌های مبتنی بر متن از جستجوی کلمات کلیدی، تطبیق محتوای صفحه با عبارت جستجو و تطبیق متن صفحه با هستان‌شناسی یا گنجینه واگان استفاده می‌گردد. روش‌های مبتنی بر پیوند نظری [14,20] به شناسائی صفحات محبوب یا پرمراجعه می‌بردارند.

در روش‌های مبتنی بر هستان‌شناسی [4,11]، بخش زیادی از صفحات مرتبه با توجه به ارتباط معنایی و مفهومی بازیابی می‌شوند که در روش‌های معمول بازیابی نمی‌باشند. روش‌های مبتنی بر فضای برداری [6,7] به عنوان دسته دیگری از روش‌های خوش‌نمترکر، به مقابله عنوان پرس و جو با صفحه می‌پردازند. روش مبتنی بر بازخورد [8] نیز با توجه به اینکه عملیات خوش را طی چندین مرحله اجرا می‌کند و نتایج هر مرحله را در اختیار مراحل بعدی قرار می‌دهد، روش زمانبری می‌باشد. از جمله مزیتهای این روش آن است که با توجه به بهبود حاصله در هر مرحله از خوش نسبت به مراحل قبلی با داشتن دانش از مرحله پیشین نتایج هر مرحله نسبت به مراحل قبلی قابل بهبود می‌باشد. در واقع در این حالت خوش‌گر از قابلیت یادگیری برخوردار است.

۳-۲ استخراج شاخصهای خوش زبانی

۳-۲-۱ نیازمندیهای خوش زبانی

با توجه به ایده خوش‌نمترکر که در بخش قبلی مورود شد، خوش‌گر زبانی نیز دارای دو نیازمندی می‌باشد: اولاً باید دارای روش مناسبی جهت شناسائی صفحات مرتبه یا صفحات متعلق به زبان باشد و ثانیاً باید از سیاست هوشمندانه‌ای در جهت دنبال نمودن صفحات مرتبه و برخورد با صفحات نامرتبط برخوردار باشد. علاوه بر این باید از روش اولویت‌بندی قابل قبولی نیز استفاده کند.

» شناسایی زبان

روشهای مطرح در شناسایی زبان از داده‌های محتوایی زیر استفاده می‌کنند:

- ایست واژه‌ها

▪ ۱- گرم‌های موجود در زبان

▪ ترکیب حروف انحصاری

استفاده از ایست واژه‌ها شامل جداگانه‌ها، حروف عطفی و حروف اضافه روش مناسبی جهت شناسایی زبان می‌باشد، زیرا احتمال وقوع آنها در انواع مختلف متون زیاد بوده و هر زبان تعداد محدود و مشخصی ایست واژه دارد. نتایج بکارگیری روش مبتنی بر ایست واژه [21] نشان می‌دهد که این روش جهت تشخیص زبان مستندات با تعداد مشخصی ایست واژه، روش مناسبی محسوب می‌شود. روش ۱- گرم [22] نیز روش مترکب و برکاربردی جهت شناسایی زبان مستندات می‌باشد. در این روش زبان مستنداتی که دارای محتوای متنی کمی باشند، قابل شناسایی نیست.

روش ترکیب حروف منحصر بفرد [23] روشی است که در مقابل روش مبتنی بر ایست واژه مطرح شده است. در این روش زیر رشته‌های خاص زبان به عنوان شناسه‌های آن زبان مورد استفاده قرار می‌گیرد. مشکل این روش احتمال کم وقوع رشته منحصر بفرد در متنهای کوتاه است. مشکل دیگر در این روش زمانی رخ می‌دهد که که زیررشته‌های منحصر بفرد بدستی انتخاب نشوند. زیرا بسیاری از زیررشته‌ها در زبانهای مختلف تکرار می‌شوند.

» سیاست واکنشی

خرشگر زبانی در برخورد با صفحات مرتبه دو سیاست را می‌تواند در پیش گیرد. سیاست اول آن است که پیوند موجود در صفحات مرتبه را واکنشی نموده و صفحات غیرمرتبه را دور بریزد. سیاست بعدی آن است که مسیر پیوندهای موجود در صفحه غیرمرتبه را تا عمق مشخصی دنبال نماید. زیرا احتمال دارد که بعد از پیمایش عمق مشخصی به یک صفحه مرتب ختم شود. یکی از روش‌های برخورد با این مسئله می‌تواند به این صورت باشد که برای هر صفحه یک پارامتر عمق در نظر گرفت و در صورتیکه پیوند موجود در آن غیرمرتب باشد، از پارامتر عمق یک واحد کم شده و برای صفحه مورد مراجعه توسعه پیوند غیرمرتب عمق جدید را در نظر گرفت. این روند ادامه می‌باید تا جاییکه به عمق صفر برسد و خوش در آن مسیر متوقف شود. همچنین در صورتیکه سیاست اولویت‌دهی واکنشی پیوندهای موجود در صفحات مرتب در مقابل صفحات غیرمرتب وجود داشته باشد، باید با مکانیزمی به پیوندهای موجود در صفحات مرتب وزن بیشتری اختصاص داد. بنابراین در یک خوش‌گر زبانی دو مولفه شناسایی زبان و الویت‌دهی واکنشی وابسته به زبان بوده و می‌بایست مد نظر قرار گیرد. این دو مولفه در شکل ۱ نشان داده شده است.

۴-۳ مقایسه روش‌های خوش و انتخاب شاخصهای

مناسب جهت خوش زبانی

روشهای مختلف خوش با توجه به شاخصها و ویژگیهای آنها که در بخش‌های قبلی آنها پرداخته شد، مقایسه و نتایج به طور خلاصه در جدول ۱ آمده است.



از اینترنت و یافتن مستندات متعلق به زبان مورد نظر می‌باشد. عملیات شناسایی زبان با جستجوی تعدادی از کلمات زبان توسط یک موتور جستجو صورت می‌گیرد. همچنین جهت محدود نمودن خزشگر، دامنه جستجوی آن به دامنه کشورهای آفریقایی جنوبی محدود می‌شود. در مراحل بعدی خرس با توجه به شناسایی کلمات بیشتری از زبان مورد نظر، تعداد مستندات بیشتری یافته می‌شوند. بعد از ایجاد این مجموعه زبانی، روش ^{۱۱}-گرم جهت شناسایی تعداد مستندات بیشتر متعلق به زبان هدف استفاده می‌شود.

۳- بررسی نیازمندیهای خزش متتمرکز و خرس

زبانی از دیدگاه مقایسه‌ای

با توجه به اینکه روش‌های خرس زبانی مبتنی بر ایده خزش متتمرکز مطرح شده‌اند، در این بخش جهت استخراج نیازمندیهای خرس زبانی ابتدا روش‌های خرس متتمرکز مورد بررسی قرار گرفته، سپس با تعیین سازوکار و ویژگیهای خرس متتمرکز و نیازمندیهای خرس زبانی، با مقایسه ویژگیهای مورد استفاده در روش‌های مختلف خرس، شاخصهای مناسب جهت خرس زبانی مشخص می‌شوند.

۱-۳ سازوکار خرس متتمرکز

با بررسی روش‌های متفاوت خرس [۱, ۱۹]، می‌توان دریافت که در عملیات خرس چهار موضوع کلیدی زیر در نظر گرفته می‌شود:

- شناسایی صفحات مرتبط
- شناسایی پیوندهای مرتبط در صفحات مرتبط
- سازوکار برخورد با صفحات نامرتبط
- سازوکار دنبال نمودن پیوندهای مرتبط

در واقع یک خزشگر عملیات خرس را از تعدادی URL اولیه آغاز می‌کند، صفحات متناظر با این URL‌ها را واکنش نموده و محتوای این صفحات را تجزیه می‌کند. سپس بر اساس شاخصهای مورد نظر، صفحات مرتبط را شناسایی کرده و اقدام به بررسی اطلاعات مستخرج می‌نمایند. بدیهی است که براساس بررسی‌های صورت گرفته در خصوص دنبال نمودن پیوندهای موجود در این صفحات تصمیم‌گیری می‌نماید. برای دنبال نمودن پیوندها، آنها به لیست URL‌های قابل واکنش اضافه می‌شوند. البته در روش‌های مبتنی بر اولویت، پیش از اینکه یک آدرس به لیست URL‌های قابل واکنش اضافه شود، بر اساس سازوکاری وزنده می‌شود.

۲- استخراج شاخصهای خرس متتمرکز

مهمنترین فرآیندی که در خرس صورت می‌گیرد، انتخاب صفحات مرتبط و متعاقباً پیوندهای مرتبط است. به این منظور برخی روش‌ها مانند روش ماهی پیوندهای موجود در صفحات مرتبط را مرتبط فرض نموده و آنها را به لیست URL‌های قابل واکنش اضافه می‌کنند. دسته‌ای دیگر از روش‌ها اعتبار خود پیوندها را نیز به کمک اطلاعاتی نظیر متن پیوند و متن اطراف پیوند بررسی می‌کنند. روش کوشه ماهی از جمله این روش‌ها می‌باشد.

قابلیت یادگیری شناخته می‌شوند. در این گونه سیستمها بتدربیج مسیرهای غیرمرتبط حذف شده و مسیرهای بهینه‌تری جهت جمع‌آوری صفحات مرتب مرتبط پیدا می‌شوند. همچنین روش‌های خزش متتمرکز وب می‌توانند از سطوح متفاوتی از دانش جهت یافتن صفحات مرتب استفاده نمایند. به عنوان مثال منبع اطلاعاتی جهت یافتن صفحات مرتب ممکن است کلمات کلیدی، عبارت پرس و جو [۶, ۷] یا دانش مکمل نظیر هستان‌شناسی^{۱۴} یا گنجینه واژگان^{۱۵} [۴, ۱۱] باشد. استفاده از دانش مکمل منجر به پوشش بیشتر و دقیق‌تر شدن نتایج خرس می‌شود. زیرا تنها به وجود کلمات کلیدی اکتفا نشده و ارتباط معنایی کلمات موجود در متن و موضوع مورد نظر نیز لحاظ می‌گردد. اطلاعات مورد استفاده توسط روش‌های خرس متتمرکز جهت تعیین سیاست واکنشی صفحات می‌تواند مبتنی بر وراثت باشد. این اطلاعات که با توجه به ارتباط پیر و فرزندی بین صفحات موجود در گراف وب حاصل می‌شود، منجر می‌شود که امتیاز فرزندان تحت تاثیر امتیاز صفحات اجاد باشد. برخی روش‌ها نظیر [۱۴] و الگوریتم کوشه ماهی [۱۰] از نمره وراثتی جهت امتیازدهی به فرزندان استفاده می‌کنند.

تحلیل صفحات وب نیز بسته به نیازها و اهداف خرس، در حین عملیات خرس یا بعد از آن صورت می‌گیرد [۱۳]. در واقع دو امکان تحلیل پویا^{۱۶} و ایستا^{۱۷} وجود دارد. در تحلیل ایستا ابتدا صفحات با حداقل مشخصات در یک مخزن جمع‌آوری شده و بعدها بسته به کاربرد مورد نظر تحلیل می‌شوند. در این حالت محدودیت‌های حافظه‌ای و منابع شبکه باید لحاظ شود. در حالیکه در تحلیل پویا امکان دسته‌بندی اطلاعات مناسب با کاربرد مورد نظر نیز وجود دارد.

مبحث دیگری که در ارتباط با خرس متتمرکز وجود دارد، خرس وب عمومی و وب پنهان است. وب عمومی به بخشی از وب اطلاق می‌شود که دسترسی به آن از طریق دنبال نمودن پیوندها امکان‌پذیر باشد و نیاز به تأیید یا ثبت‌نام نداشته باشد. در مقابل وب عمومی، بخش دیگری تحت عنوان وب پنهان وجود دارد که دسترسی به آن نیاز به اخذ برخی مجوزها و پرکردن برخی از فرمها دارد. این بخش از وب حجم وسیعی از اطلاعات مفید و با ارزش وب را دربردارد.

در ارتباط با خرس زبانی فعالیتهای محدودی صورت گرفته است [۱۵, ۱۷, ۱۸]. روشی تحت عنوان خرس زبانی وب (LSWC)^{۱۸} با هدف ایجاد آرشیوی بزرگ از وب کشور تایلند در [۱۷] پیش‌هاد شده است. روش LSWC جهت شناسایی زبان از اطلاعات فرابرچسب و روش ^{۱۱}-گرم استفاده می‌کند. در واقع شناسایی زبان در دو مرحله صورت می‌پذیرد. ابتدا از اطلاعات فرابرچسب کاراکترست جهت شناسایی زبان استفاده می‌شود و در صورت عدم شناسایی زبان صفحه، در مرحله بعد روش ^{۱۱}-گرم مورد استفاده قرار می‌گیرد. در روش معرفی شده در [۱۵] شناسایی زبان صفحه به کمک اطلاعات فرابرچسب کاراکترست یا به کمک ابزار شناسایی کاراکترست تحت عنوان Mozilla Charset Detector انجام می‌شود. به دلیل اینکه برخی از زبانها مانند زبان تایلندی توسط این ابزار شناسایی نمی‌شوند، از اطلاعات فرابرچسب آنها استفاده می‌شود. ایده اصلی روش [۱۸] خرس دامنه وسیعی

¹⁴ Ontology

¹⁵ Thesaurus

¹⁶ Dynamic

¹⁷ Static

¹⁸ Language Specific Web Crawling



این ویژگیها در روش پیشنهادی این روش را از سایر روش‌های مطرح در حوزه خزش زبانی تمایز نموده است.

همچنین با توجه به اهمیت جمع‌آوری صفحات هدف در روش‌های خزش متمرکز بطور کلی و روش‌های خزش زبانی به طور خاص، لازم است سیاستهای مورد نیاز جهت انتخاب و اولویت‌بندی واکنشی صفحات در حین عملیات خزش در نظر گرفته شود. در واقع با توجه به اینکه در خزشگرهای پیوندهای قابل واکنشی با توجه به اولویت‌بندی که برای خزشگر تعریف می‌گردد، چیده می‌شوند. لذا در صورتیکه سیاست واکنشی در خزشگر مبنی بر اولویت‌دهی به صفحات فارسی تعریف نگردد، علی‌رغم تشخیص زبان، خزشگر کلیه پیوندهای موجود در صفحات وب را واکنشی می‌نماید بنابراین بخش مهم دیگری که در راستای خزش زبانی در این مقاله به آن پرداخته شده است، سیاست واکنشی جهت دنبال نمودن صفحات متعلق به زبان می‌باشد.

در این مقاله جهت ارائه روشی برای خزش بهینه مستندات فارسی وب، ابتدا روش‌های مختلف خزش متمرکز مورد بررسی قرار گرفته و سپس دیدگاهی برای خزش زبانی به ویژه زبان فارسی ارائه شده است. براساس دیدگاه پیشنهادی و با استفاده از خزشگر متن باز وابر^{۱۲}، یک خزشگر زبانی برای زبان فارسی طراحی و پیاده سازی شده است. آزمایش‌های عمل آمده بر روی این خزشگر نشان می‌دهد که دیدگاه پیشنهادی در خزش موثر صفحات فارسی بسیار کارا عمل کرده است. در ادامه مقاله به شکل زیر سازماندهی شده است. در بخش ۲ کارهای مرتبط در زمینه خزش متمرکز موردنی گردد، در بخش ۳ شاخصهای خزش متمرکز استخراج شده و مبتنی بر این شاخصها و نیازمندی‌های خزش زبانی، دیدگاه خزش زبانی در بخش ۴ پیشنهاد می‌شود. در بخش ۵ بر اساس دیدگاه پیشنهادی نحوه پیاده‌سازی و ارزیابی خزشگر ارائه می‌گردد.

۲- مروری بر کارهای مرتبط

امروزه روش‌های خزش متمرکز به دلیل صرفه‌جویی در استفاده از منابع مورد توجه محققان قرار گرفته است. این روشها در بسیاری از موارد تحت عنوان خزش موضوعی^{۱۳} شناخته می‌شوند. در خزش موضوعی صفحات مرتبط با موضوع مورد نظر جمع‌آوری می‌شوند و بسته به روش آن ممکن است از کلمات کلیدی، عبارت پرس و جو یا واژگان تخصصی موضوع استفاده شود [۱,2,3,4,6,7,9,10,11,14]. در روش خزش متمرکز به جای جمع‌آوری و شاخص‌گذاری تمام صفحات وب (به منظور پاسخ‌گویی به تمام پرس و جوهای ممکن)، با شناسائی محدوده خزش، مسیرهای مرتبط مشخص شده و خزش در محدوده‌های غیرمرتبط متوقف می‌شود. این امر منجر به صرفه‌جویی قابل توجه در ساخت‌افزار و منابع شبکه شده و بروزرسانی مستندات واکنشی شده توسط خزشگر را نیز بهبود می‌بخشد.

عملیات خزش متمرکز می‌تواند به صورت یک مرحله‌ای یا چندمرحله‌ای انجام شود. در روش چندمرحله‌ای که به آنها روش‌های مبتنی بر بازخورد نیز گفته می‌شود، در هر مرحله خزش از داشت حاصل از مراحل پیشین استفاده می‌شود [۸]. به این ترتیب انتظار می‌رود نتایج هر مرحله نسبت به مراحل قبلی بهبود یابد. روش‌های مطرح در این دسته اغلب تحت عنوان روش‌های با

یکی از سازوکارهای موثر در بهبود عملکرد موتورهای جستجو در جهت بازیابی اطلاعات مطرح است.

برنامه‌ای که خزش وب توسط آن انجام می‌شود، خزشگر^۴ نامیده می‌شود. هدف اصلی از طراحی خزشگرهای وب بازیابی صفحات از وب و ذخیره نمودن آنها در مخازن محلی می‌باشد. چنین مخزنی بعدها برای کاربردهای مانند موتورهای جستجو مورد استفاده قرار می‌گیرد. خزش به دو صورت همه منظوره^۵ یا عمومی^۶ و خاص منظوره یا متمرکز^۷ [۳] امکان‌پذیر است. در خزش متمرکز موضوع یا حیطه خزش به صورت دقیق مشخص می‌شود در حالیکه یک خزشگر وب همه منظوره با شروع از مجموعه مشخصی URL‌ها^۸، هر تعداد صفحه که بتواند از وب واکنشی می‌کند. مهتمین ویژگی خزش متمرکز آن است که نیاز به جمع‌آوری همه صفحات وب ندارد بلکه تنها صفحات مرتبط را انتخاب و جمع‌آوری می‌کند. خزش زبانی به عنوان یکی از روش‌های خزش متمرکز به جمع‌آوری صفحات متعلق به یک زبان می‌پردازد. بسیاری از روش‌های خزش زبانی با توجه به ویژگی محلی زبانی در وب مطرح شده‌اند. این روشها عموماً بر اساس این واقعیت عمل می‌کنند که در گراف وب بین صفحات متعلق به یک زبان معمولاً پیوند وجود دارد. از جمله کاربردهای یک خزشگر زبانی استفاده از آن به عنوان یک مولفه مبنایی در جهت توسعه موتور جستجوی وب زبانی می‌باشد. همچنین پیکره ایجاد شده توسط خزشگر زبانی در جهت کاربردهای نظیر تشخیص الگو و پردازش زبان طبیعی قابل استفاده است. به عنوان مثال در [۲۵] از یک خزشگر غیرفارسی جهت ایجاد یک پیکره مبنی فارسی استفاده شده است.

بطور کلی در ارتباط با خزش زبانی فعالیتهای بسیار کمی صورت پذیرفته است که از جمله آنها می‌توان به خزش زبانی وب تایلند اشاره نمود [۱۵]، [۱۷] در روش‌های خزش زبانی مطرح مانند روش مبتنی بر استوازه^۹ [۲۱]، روش ترکیب حروف انحصاری^{۱۰} و روش مبتنی بر ۱۱-گرم [۲۲]، عموماً تاکید بر استفاده از کاراکترست زبان یا شناسائی زبان به کمک کلمات موجود در زبان است.

در ارتباط با خزش متمرکز مبتنی بر زبان فارسی نیز فعالیتهای زیادی صورت نگرفته است. در [۲۴] و کارهای مشابه در مورد وب سایر زبانها، گستره خزش محدود به دامنه کشورها شده است و همانطور که خود نویسنده‌گان گفته‌اند، این روش کارایی لازم را ندارد. با توجه به پراکندگی وب سایتها ایرانی روی سرورهای خارج از دامنه ایران (آ.) نیاز به ارائه روشی جهت بهبود خزش مستندات فارسی وب احساس می‌شود. در این مقاله خزشگر زبانی پیشنهاد شده با شناسائی صفحات فارسی وب در هر مرحله امکان پوشش صفحات فارسی بیشتری را برای مراحل بعدی فراهم می‌آورد. سازوکار خزش ارائه شده مبتنی بر ترکیب مناسبی از ویژگیهای فرابرجسب^{۱۱} و محتوای صفحات وب امکان شناسائی صفحات فارسی را با دقت و سرعت مناسبی فراهم می‌کند. با توجه به اینکه عموم روش‌های خزش زبانی مطرح تنها یکی از این دو ویژگی را بکارمی‌برند، بکارگیری ترکیبی از

⁴ Crawler

⁵ General Purpose

⁶ General

⁷ Focused

⁸ Uniform Resource Locator

⁹ Stop Word

¹⁰ Unique Letter Combination

¹¹ Meta-tag

¹² WIRE
¹³ Topic-Driven



یادداشت پژوهشی

طراحی و پیاده‌سازی یک خزشگر زبانی جهت بهبود سازوکار خزش در مستندات فارسی و وب

ابوالفضل آل‌احمد

علیرضا یاری

معصومه عظیم‌زاده

دانشگاه تهران
گروه تحقیقاتی پایگاه داده‌ها
a.aleahmad@ece.ut.ac.ir

مرکز تحقیقات مخابرات ایران
پژوهشکده فناوری اطلاعات
a_yari@itrc.ac.ir

مرکز تحقیقات مخابرات ایران
پژوهشکده فناوری اطلاعات
azim_ma@itrc.ac.ir

تاریخ دریافت: ۱۳۸۸/۱/۲۹ - تاریخ بذریش: ۱۳۸۸/۶/۲۴

چکیده- حجم زیاد، ماهیت پویا و غیرقابل کنترل وب چالش‌های زیادی را در خصوص خزش وب ایجاد نموده است. روش‌های خزش به طور کلی به دو دسته عمومی و متتمرکز قابل تقسیم هستند. در روش خزش عمومی همه صفحات وب جمع‌آوری می‌شوند و در روش خزش متتمرکز تنها بخشی از صفحات وب که با موضوع خاصی مرتبط هستند، جمع‌آوری می‌گردند. خزش زبانی به نوعی از خزش متتمرکز اطلاق می‌شود که صفحات نوشته شده به زبان مورد نظر را جمع‌آوری می‌کند. با توجه به اینکه وب حاوی گستره وسیعی از داده‌های بدون ساختار و نوشته شده به زبان‌های مختلف است، نحوه انجام خزش زبانی از جمله چالش‌های بازیابی اطلاعات در محیط وب است. در این مقاله برای بهبود خزش مستندات فارسی وب، یک خزشگر زبانی پیشنهاد گردیده و تشریح شده است. نتایج حاصل از پیاده‌سازی و تست این خزشگر نشان می‌دهد خزشگر زبانی در خزش صفحات فارسی وب با کارایی بهتری عمل می‌کند.

کلیدواژه‌ها: خزشگر فارسی، خزش متتمرکز، خزش زبانی، بازیابی اطلاعات.

لبه‌ها دسترسی از صفحات مبدأ به صفحات مقصد امکان‌پذیر می‌شود. ماهیت در حال تغییر وب علاوه بر تغییر محتوای صفحات و حذف و اضافه شدن آنها، منجر به تغییر ساختار ارتباطی فرایوندها یا صفحات نیز می‌گردد که این تغییرات وب منجر به عدم کنترل‌پذیری آن و طرح چالش‌های فراوانی در ارتباط با خزش و "بازیابی اطلاعات"^۲ از وب می‌شود. خزش^۳ وب به عنوان

۱- مقدمه

وب متشکل از مجموعه صفحاتی است که از طریق فرایوند^۱ به یکدیگر متصل شده‌اند که امکان دسترسی آنها به یکدیگر را تسهیل می‌کند. به طور خلاصه وب در قالب یک گراف قابل تصور است که صفحات آن گره‌های این گراف و فرایوندها لبه‌های آن را شکل می‌دهند. به کمک این فرایوندها یا

² Information Retrieval

³ Crawling

¹ Hyperlink

