



مروی بر شخصی سازی نتایج موتور جستجو با روش‌های هوشمند

دکتر مهدی یعقوبی^۱، ملیحه محمدزاده^۲

^۱ دانشکده فنی مهندسی، دانشگاه آزاد اسلامی واحد مشهد
yaghoubi@mshdiaum.ir

^۲ دانشکده فنی مهندسی، دانشگاه آزاد اسلامی واحد مشهد
malihemohamadzade@yahoo.com

چکیده

با افزایش روز افزون اطلاعات در وب ، یافتن مطالب مورد نیاز امر دشوار و وقت کیری شده است. شخصی سازی نتایج موتور جستجو هر تلاشی است که باعث شود نتایج جستجو متنا سب با علایق و دانش و نیاز کاربر نمایش داده شود. شده است که کارهای انجام شده در حیطه شخصی سازی نتایج موتور جستجو از دو منظر نحوه انجام کار و مدل جمع آوری اطلاعات بررسی شده اند. در این مقالات او دو دیدگاه فازی و قطعی با مساله برخورد شده است و روش‌های موجود بررسی شده است.

کارهای انجام شده در این محدوده به ویژه در سالهای اخیر به موفقیتها بسیاری رسیده اند و توانسته اند نتایج موتورهای جستجو را تا سطح زیادی بهبود بدهند. هر چند هنوز هیچ سیستم تمام اتوماتیک و تمام فازی که در آن از مفهوم کاوی صفحات وب استفاده شود، ایجاد نشده است.

کلمات کلیدی

شبکه های مفهومی فازی ، شخصی سازی ، مدل سازی کاربر ، موتور جستجو

دارد نتایج جستجو را مرتب میکند. بر این اساس ما کارهای

نجام شده را از دو منظر نحوه مدل

کاربر کاربر (شناختن کاربر) و نحوه

نجاد ای سیستم شخصی سازی بررسی کر

ده. در قسمت دوم کارهای انجام

شده را مبنظر نحوه مدل کردن

کاربر بررسی کرد. ایم. در قسمت سوم

نیز کارهای انجام شده را از منظر

نحوه اجرای سیستم کاربری

ایم. در قسمت سوم نتایج و معاایب

کارهای انجام شده را بررسی می

کنیم. در قسمت چهارم نتایج گیری و

در قسمت پنجم مراجع آورده شده

است.

۱ مقدمه

در شخصی سازی نتایج موتور جستجو سعی می شود تا نتایجی که موتورهای جستجو نمایش می دهند با علایق و دانش و نیازهای کاربر و منظور کاربر از جستجو مناسب باشد. یعنی می خواهیم ترتیب نتایجی که در یک موتور جستجو بر اساس پرس و جوی کاربر نمایش داده می شود به ترتیب علاقه کاربر باشد.

در این مقاله کارهایی که در حیطه شخصی سازی نتایج موتورهای جستجو انجام شده بررسی شده است روشنی که در همه این کارها مشترکاً انجام شده است این است که ابتدا کاربر شناخته می شود و نیازها و حیطه کاری او تشخیص داده می شود و سپس سیستم بر اساس شناختی که از کاربر

کند و بر این اساس مدل علایق کاربر به دست می آید.

در این سیستم ابتدا با کمک فایل result.htm که نتایج حاصل از موتور جستجو رالیست می کند صفحات لیست می شود سپس کاربر صفحاتی را که مورد علاقه اش است انتخاب می کند که به آن hit می گویند و صفحاتی که انتخاب نمی کند miss گفته می شود سپس در یک rerank صفحات hit و صفحاتی که بیشتر مشابه این صفحات هستند لیست می شوند. به این صورت سعی می شود نتایج نزدیک به علاقه و نیاز کاربر باشد. همانطور که دیده می شود در این روش نیز کاربر دخالت صریح دارد.

در [29] روند کار به این صورت است که ابتدا همه صفحات بازیابی شده در کلاس بد قرار میگیرند. پس از اینکه کاربر بر روی یکی از نتایج موجود در لیست بردا کلیک نمود، با این فرض که صفحه کلیک شده مورد علاقه کاربر می باشد، آن را به کلاس خوب منتقل مینماید و با استفاده از طبقه بندي کننده Naive Bayesian یک نمره به هر صفحه در لیست بردا اختصاص میدهد. سپس لیست بر اساس نمرات اختصاص یافته رتبه بندي شده و به کاربر ارائه می شود تا کاربر صفحه دیگری را به عنوان صفحه مورد علاقه انتخاب نماید و مراحل قبلی مجدداً تکرار می شود.

پروفایل کاربر را به صورت یک کوکی روی سیستم آه نخیره می کند . در این سیستم کلیک هون کاربر باید در حین جستجو هر بار صفحات مورد علاقه خود را انتخاب نماید و سیستم چندین بار با توجه به علاقه کاربر rerank می شود حالت ضمنی وجود ندارد.

کاربران معمولاً به انتشار علایقشان در اینترنت بدین هستند و همچنین این کاربرایشان وقت گیر و هزینه بر است و ممکن است که در بیان علایقشان دچار اشتباه شوند یا به خاطر عدم اعتمادی که به فضای وب

۲ کارهای انجام شده از منظر نحوه مدل سازی کاربر

در زمینه نحوه جمع آوری اطلاعات کاربر دو رویکرد وجود دارد.

۱-۲-مدل سازی صریح کاربر

در رویکرد اول اطلاعات شخصی کاربر به صورت صریح از او پرسیده می شود و بر اساس پاسخهایی که کاربر به سیستم اعلام می کند پروفایل شخصی او ساخته میشود. مثل کارهای انجام شده در [29],[12],[6],[5],[4],[3],[2],[14] و همچنین personal google.

در [2] از یک کلاسترینگ سلسله snaket مراتبی استفاده می کند که نامیده می شود که نتایج حاصل از ۱۶ موتور جستجوی کالا در پوشه های بر چسب زده سلسله مراتبی مرتب می شوند. حالت سلسله مراتبی دید کاملی از نتایج مرتب شده موتورهای جستجو ارائه می دهد. که کاربران با این حالت سلسه مراتبی به نیازهای جستجوی خودشان هدایت میشوند. ابتدا کاربر یک پرس وجو را به

SNAKET می فرستد و SNAKET نیز پس از خوش بندی سلسله مراتبی نتایج جستجو و برچسب گذاری خوش ها با جملات با طول متغیر، یک سلسله مراتب برچسب گذاری شده را به کاربر ارائه میدهد. کاربر نیز گروههایی را که برچسب آنها بیشترین تناسب با اطلاعاتی مورد نیازش را دارد انتخاب میکند . سپس SNAKET با فیلترکردن نتایج جستجوی متعلق به سایر خوش ها، نتایج جستجوی شخصی سازی شده را به کاربر ارائه می دهد.

همانطور که دیده می شود این روش به دخالت صریح کاربر نیاز دارد که روش مطلوبی نیست.

در [4] از یک پایگاه دانش استفاده شده که از رفتار کاربران به دست آمده است سیستم رفتار کاربران را مانیتور می کند سپس کاربر خودش حیطه علایقش را معین می

3- کارهای انجام شده از منظر نحوه اجرای کار

از نظر نحوه انجام کار نیز دو رویکرد وجود دارد.

3-1- روش قطعی

در رویکرد اول با این مساله به صورت یک مساله crisp برخورد شده است. در این سیستمها پس از ساختن پروفایل کاربر، کاربران با توجه به پروفایل شان به صورت قطعی و همیشگی عضو یک گروه می‌شوند و صفحات نیز پس از بررسی محتویاتشان به صورت قطعی گروهبندی می‌شوند و ارتباط بین مفاهیم و کاربران نیز به صورت قطعی برقرار می‌شود. مثل کارهای انجام شده در [3],[4],[16],[31],,,[17]

در [31] از الگوریتم BM25 برای مرتب کردن صفحات با توجه به پرس و جوی ارائه شده به سیستم استفاده می‌شود. در این روش با درنظر گرفتن یک مجموعه N تایی شامل کل اسناد موجود در وب و یک مجموعه R تایی شامل اسناد مورد علاقه کاربر، وزن هر مفهوم با استفاده از فرمول زیر محاسبه می‌شود:

فرمول N_i تعداد اسناد موجود در وب شامل مفهوم i و n_i تعداد اسناد مورد علاقه کاربر شامل مفهوم i می‌باشد. در [30] کلاسندی کننده بر اساس کلیدهای محتوا یکی از الگوریتم k نزدیکترین همسایه اسناده می‌کند. یک مجموعه از اسکرپتها که log file را پردازش می‌کند و کارایی آن را برای هر کاربر ارزیابی می‌کند استفاده می‌شود که آن را به مجموعه تست و آموزش تقسیم می‌کند. پروفایل کاربر به صورت یک سلسله مراتب وزن‌دار نمایش داده می‌شود. در این سیستم یک پروفایل کاربر داریم یک پروفایل اسناد که شباهت

دارند اطلاعات صحیحی را به سیستم ندهند، مجموع این دلایل باعث می‌شود که جمع آوری پروفایل کاربر به صورت اتومات نتایج بهتری را به دست آورد.

2-2- مدل سازی ضمنی کاربر

در رویکرد دوم سعی شده که دخالت کاربر به حداقل برسد و پروفایل او به صورت ضمنی از روی تعاملاتی که سیستم با کاربر دارد مثل مرورهای کاربر، پرس و جوهای گذشته او، ایمیلهایی که مشاهده کرده و... شناخته می‌شود. مثل کارهای انجام شده در [17],[16],[10],[30].

در [30] سعی شده که کمترین دخالت کاربر وجود داشته باشد. آنها پروفایل کاربران در موتور جستجو می‌سازند. در آن کار یک لفافه بند برای گوگل (google wrapper) اجرا شده است.

منابع مختلف پروفایل، پرس و جو ها و نتایج جستجو ها است. این پروفایل با کلاسندی اطلاعات در محتویات پروژه دایرکتوری باز ارثی ساخته می‌شود و سپس برای دوباره مرتب کردن نتایج جستجو استفاده می‌شود از بازخورد های کاربر استفاده می‌شود تا ترتیب نتایج گوگل با ترتیب جدید ما مقایسه شود و تا حد ۳۷ نتایج بهبود یافته است. پایه این کار بر ساختن پروفایل کاربر از تعاملات کاربر با یک موتور جستجوی خاص است. در این کار از google wrapper استفاده شده است یعنی یک لفافه حول موتور جستجوی گوگل برای ثبت کردن پرس و جو ها و نتایج جستجوها وکلیکها بر پایه هر کاربر قرار می‌گیرد. ایراد این کار این است که اولاً فقط یک موتور جستجوی خاص را بررسی می‌کند و ثانیاً سمت سرور است و به غیر از جستجوهای کاربر از وبگردی ها و رفتار کلی او در وب اطلاعاتی نخواهیم داشت.

رتیه بندی صفحات معتبر به دست آمده بر اساس علاقه کاربر استفاده می شود . در واقع استفاده از ساختار پیوند علاوه بر شبکه مفهومی فازی ، شخصی سازی نتایج جستجو بر اساس علاقه کاربران را بهبود می دهد . سیستم در یک شبکه مفهومی فازی مفاهیم و صفحات مرتبط به آنها را قرار می دهد . سپس نیاز کاربر را از پروفایل او استخراج می کند و سپس صفحات مرتبط با نیاز کاربر را با توجه به ساختار لینک و مفاهیم مورد علاقه کاربر پیدا کرده و در لیست مرتب شده قرار می دهد . اما این پروفایل با دخالت کاربر ساخته می شود و همچنین شبکه با نظارت فرد خبره ساخته می شود که امر مطلوبی نیست .

در [1] شخصی سازی نتایج موتور جستجو با استفاده از شبکه مفهومی فازی گسترش یافته انجام می شود . در این کار پروفایل کاربر و نتایج جستجو با یک شبکه مفهومی فازی گسترش یافته نمایش داده می شوند . در شبکه مفهومی فازی گسترش یافته 4 نوع ارتباط فازی که شامل ارتباط فازی مثبت ، ارتباط فازی منفی ، ارتباط فازی عمومی و ارتباط فازی خصوصی می باشد ، می تواند بین مفاهیم وجود داشته باشد که با یک ماتریس رابطه نمایش داده می شود . این شبکه نیز تحلیل فرد خبره ساخته می شود .

همانطور که دریده می شود کارهای انجام شده گذشته یا پروفایل کاربر را به صورت صریح می سازند و یا اگر پروفایل به روش صریح ایجاد شود در اجرای کارهای روشهای crisp استفاده می کند و سیستم وجود ندارد که در آن هم پروفایل کاربر به صورت ضمنی و اتوماتیک ساخته شود و هم از روشهای فازی استفاده شده باشد .

البته در [1],[9] تا حدودی سعی شده است که شبکه مفهومی فازی با پروفایل اتوماتیک راه اندازی شود اما ایراد این سیستم هم این است که در آن هر چند طبقه بندی مفاهیم

مفهومی بین آنها با فرمول زیر محاسبه می شود :

اما در رویکرد دوم با این مساله با دید فازی برخورد شده است . عدم قطعیت این مساله در پیدا کردن مفاهیم اسناد و هم در تشخیص نیاز کاربر از پرس وجودی او و هم در ارتباط بین سند و مفهوم مورد نیاز کاربر وجود دارد و همچنین ارتباط بین پرس وجودی ها و مفاهیم نیز دارای سطحی از عدم قطعیت است . ما نمی توانیم ادعا کنیم که کاربری به صورت قطعی عضو یک گروه است و به گروههای دیگر هیچ علاقه ای ندارد . این خوش بندی ماهیت غیر قطعی و فازی دارد و برخورد قطعی با آن نا مناسب است . سیستمهایی که از روش فازی برای دسته بندی کاربران و محتویات صفحات و مفاهیم استفاده کنند این عدم قطعیت را در نظر گرفته اند . مثل کارهای انجام شده در [9],[10],[1]

در [9] شبکه مفهومی فازی مسئول شخصی سازی لینکهای مرتبط به پرس وجودی کاربر است . این شبکه را می توان برای هر کاربر با توجه به پروفایل او ساخت . پس از ارسال پرس وجودی کاربر به یک موتور جستجوی متنی ، نتایج بازیابی شده در مجموعه ریشه قرار میگیرند ، سپس صفحات وبی که به صفحات موجود در مجموعه ریشه لینک دارند و صفحات وبی که این صفحات به آنها لینک دارند نیز به مجموعه ریشه اضافه می شوند . در مجموعه حاصل صفحاتی که بیشترین لینک به آنها وجود دارد ، یعنی صفحات معتبر ، و صفحاتی که بیشترین لینک را به صفحات دیگر دارند ، یعنی صفحات هاب ، تعیین می شوند که الگوریتم هایی برای تعیین آنها نیز معرفی می شود . حال از رویکرد شبکه مفهومی فازی برای

- [4] M. Radovanovic and M. Ivanovic, "CatS: A Classification-Powered Meta-Search Engine", *Advances in Web Intelligence and Data Mining*, Springer-Verlag, vol. 23, pp. 191–200, 2006
- [5] W. Kim, L. Kerschberg, A. Scime, "Learning for automatic personalization in a semantic taxonomybased meta-search agent", Elsevier, *Electronic Commerce Research and Applications* 1 (2002) 150–173
- [6] K. J. Kim, S. B. Cho, "Personalized mining of web documents using link structures and fuzzy concept networks", Elsevier, *Applied Soft Computing* 7 (2007) 398–410, 2005
- [7] Hong Zhang¹, Yanhong Ma², Qiuyu Zhang¹, Pengshou Xie¹, Zhongxian Bao "Personalized Intelligent Search Engine Basedon Web Data Mining"Proceedings of the 2009 International Workshop on Information Security and Application (IWISA 2009)Qingdao, China, November 21-22, 2009
- [8] Mariam Daoud , Lynda-Tamine Lechani,Mohand Boughanem,"Towards a graph-based user profile modeling f or a session-based personalized search ",Springer-Verlag London Limited 2009
- [9] C. H . Lee,Y . H . Kim,P . K . Rhee,"web personalization expert with combining collaborative filtering and association rule mining technique",Elsevier,expert system with application 21(2001) 131-137
- [10] P . Kazienko , M . Adamski,"Adrosa – adaptive personalization of web advertising",Elsevier,information sciences 177(2007) 2269-2295
- [11] B . P . C . Yen, R . C . W. Kong, "Personalization of information access for electronic catalogs on the web", Elsevier, *Electronic commerce research and applications* 1 (2002) 20 – 40
- [12] Y . F . kuo,L . S . Chen,"personalization technology application to internet content provider", Elsevier, *Expert Systems with Applications* 21 (2001) 203 – 215
- [13] M. Clements , A.P. de Vries , M.J.T. Reinders," The influence of personalization on tag query length in social media search", Elsevier, *Information Processing and Management* 46 (2010) 403–412
- [14] L. Yuen, M. Chang, Y. K. Lai and C. K. Poon, "Excalibur: a Personalized Metasearch

فازی است اما گروهندی کاربران به صورت crisp است که با عدم قطعیت مساله همواری ندارد . در مجموع در هیچ یک از کارهای گذشته تهیه پروفایل کاربر به صورت اتوماتیک و ضمنی و به روز رسانی پروفایل او و تولید شبکه فازی به صورت اتوماتیک و به روز رسانی آن و کلاسیندی فازی کاربران با هم انجام نشده است و در بعضی مقالات از روشهای فازی استفاده شده اما شبکه فازی توسط فرد خبره تولید می شود و این کار اتوماتیک نیست. فقط در یک مقاله این کار اتوماتیک انجام شده که در آن به روز رسانی شبکه فازی و به وجود ندارد. همچنین کلاسیندی فازی کاربران در هیچ کدام از کارهای قبلی وجود ندارد.در همه کارهای انجام شده حتی کارهایی که از روشهای فازی استفاده می کنند کلاسیندی به صورت crisp است که این امر اصل فازی بودن و غیر قطعی بودن موضوع را زیر سوال می برد. همچنین در هیچ کدام از کارهای انجام شده به روز رسانی شبکه فازی با گذشت زمان و با تغییر پروفایل کاربر وجود ندارد که این امر بسیار مهمی است،چون علیق و دانش کاربران در طول زمان تغییر می کند .

مراجع

- [1] S. M. Chen, Y. J. Horng and C. H. Lee, "Fuzzy Information Retrieval based on Multi-relationship Fuzzy Concept Networks", *Fuzzy Sets and Systems*, vol. 140, pp.183-205, 2003
- [2] P. Ferragina and A. Gulli, "A Personalized Search Engine Based on Web Snippet Hierarchical Clustering", Proceedings of the World Wide Web Conference (WWW) (The Tokio, The Japan), pp. 801-810, 2005
- [3] W. Alhalabi, M. Kubat and M. Tapia, "Search Engine Personalization Tool Using Linear Vector Algorithm", Proceedings of the 4th Saudi Technical Conference and Exhibition, pp. 336-344, 2006

Engine", Proc. 28th Annual International Computer Software and Applications Conference (COMPSAC'04), vol. 2, pp.49-50, 2004.

[15] M. Speretta and S. Gauch, "Personalizing Search Based on User Search History", Submitted to CIKM '04, (2004).

[16] J. Teevan, S. T. Dumais and E. Horvitz, "Personalizing Search via Automated Analysis of Interests and Activities", Proceedings of the ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 05), ACM Press, pp.449-456, 2005.

[17] S. Soulardatos, T. Dalamagas and T. Sellis, "Captain Nemo: A Meta-Search Engine with Personalized Hierarchical Search Space", INFORMATICA -LJUBLJANA- ,vol. 30, pp. 173-182, 2006