

امیرحسین کیهانی پور کارشناسی خود را در رشته مهندسی کامپیوتر-نرم افزار از دانشگاه صنعتی شریف و کارشناسی ارشد خود را در رشته مهندسی کامپیوتر-هوش ماشین و رباتیک از دانشگاه تهران دریافت نموده است. او هم اکنون مشغول به تحصیل در مقطع دکتری رشته مهندسی کامپیوتر-هوش ماشین و رباتیک در دانشگاه تهران می باشد. فعالیت های پژوهشی او در زمینه بازیابی اطلاعات وب، نظریه ترکیب اطلاعات و نیز الگوریتم های یادگیری ماشین می باشد.



## واژه نامه:

- <sup>1</sup> Adaptive
- <sup>2</sup> Term Frequency
- <sup>3</sup> In-Degree
- <sup>4</sup> Ordered Weighted Aggregation Operator
- <sup>5</sup> Benchmark
- <sup>6</sup> <http://research.microsoft.com/users/tyliu/LETOR/>
- <sup>7</sup> Ranking
- <sup>8</sup> Content
- <sup>9</sup> Connectivity
- <sup>10</sup> Boolean
- <sup>11</sup> Vector Space
- <sup>12</sup> Rank spamming
- <sup>13</sup> Web Information Retrieval
- <sup>14</sup> Links
- <sup>15</sup> Precision
- <sup>16</sup> Ordered Weighted Operator
- <sup>17</sup> Output links
- <sup>18</sup> Sinking pages
- <sup>19</sup> Damping Factor
- <sup>20</sup> Host
- <sup>21</sup> Domain
- <sup>22</sup> Super-node
- <sup>23</sup> Term Frequency
- <sup>24</sup> Inverse Document Frequency
- <sup>25</sup> Okapi
- <sup>26</sup> Relevance propagation
- <sup>27</sup> Score
- <sup>28</sup> Hyper-relevance
- <sup>29</sup> Sitemap-based term propagation
- <sup>30</sup> Aggregation
- <sup>31</sup> n-tuple
- <sup>32</sup> Constrained
- <sup>33</sup> Feature
- <sup>34</sup> Fine-grained
- <sup>35</sup> Coarse-grained
- <sup>36</sup> Popular
- <sup>37</sup> Rich-get-richer
- <sup>38</sup> Error Back Propagation
- <sup>39</sup> Average Precision
- <sup>40</sup> Mean Average Precision
- <sup>41</sup> Borda Voting
- <sup>42</sup> Off-line
- <sup>43</sup> Combination Rank Fine-Grained
- <sup>44</sup> Document Length



- [23] Borda, J. C. Memoire sur les elections au scrutin. Histoire de l'Academie Royale des Sciences, 1781.
- [24] Filev, D., and Yager, R. R., Learning OWA operator weights from data. Proceedings of the third IEEE Conference on Fuzzy Systems, Volume. 1, pp. 468-473, 1994.
- [25] LETOR, <http://research.microsoft.com/users/tyliu/LETOR2007>.
- [26] Jarvelin, K., & Kekalainen, J., IR evaluation methods for retrieving highly relevant documents. In Proceedings of the ACM Conference on Research and Development on Information Retrieval (SIGIR), 2000.
- [27] Xue, G. R., Yang, Q., Zeng, H. J., Yu, Y., and Chen, Z., Exploiting the hierarchical structure for link analysis, SIGIR Conference, 2005.
- [28] Qin, T., Liu, T.-Y., Zhang, X.-D., Chen, Z., & Ma, W.-Y., A study of relevance propagation for web search. In Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval, pp. 408-415, 2005.
- [29] Shakery, A., Zhai, C. X. Relevance Propagation for Topic Distillation UIUC TREC 2003 Web Track Experiments, Proceedings of TREC, 2003.
- [2] McBryan, O.A., GENVL and WWW: Tools for Taming the Web, In Proceedings of the First International World Wide Web Conference, pp. 79-90, 1994.
- [3] Gulli, A., & Signorini, A., "The Indexable web is more than 11.5 billion pages", In Proceedings of the 14th international conference on World Wide Web, ACM Press, pp. 902-903, 2005.
- [4] <http://www.worldwidewebsite.com/>, Jan. 2008.
- [5] <http://www.internetnews.com/stats/article.php/1363881>
- [6] [http://www.pewinternet.org/pdfs/PIP\\_Searchengine\\_users.pdf.2005](http://www.pewinternet.org/pdfs/PIP_Searchengine_users.pdf.2005).
- [7] Baeza-Yates, R., & Ribeiro-Neto, B., Modern Information Retrieval. ACM Press/ Addison-Wesley, 1999.
- [8] Salton, G., & Buckley, C. Term weighting approaches in automatic text retrieval. Information Processing and Management, Volume 24, Issue 5, pp. 513-523, 1988.
- [9] Robertson, S. E., Walker, S., Hancock-Beaulieu, M. M., Gatford, M., & Payne, A. Okapi at TREC-4. In NIST Special Publication. The Fourth Text REtrieval Conference (TREC-4), pp. 7396, 1995.
- [10] Henzinger, M., Motwani, R., & Silverstein, C., Challenges in web search engines, SIGIR Forum, Volume 36, Issue 2, pp. 11-22, 2002.
- [11] Henzinger, M., Hyperlink analysis for the web. IEEE Internet Computing, Volume 5, Issue 1, pp. 45-50, 2001.
- [12] Page, L., Brin, S., Motwani, R., & Winograd, T. The PageRank citation ranking: Bringing order to the web. Technical report, Computer Science Department, Stanford University, 1998.
- [13] Kleinberg, J. M. Authoritative sources in a hyperlinked environment. Journal of the ACM, Volume 46, Issue 5, pp.604-632, 1999.
- [14] Zareh Bidoki, A.M. and Yazdani, N., DistanceRank: An Intelligent Ranking Algorithm for Web Pages, Information Processing & Management, Volume 44, Issue 2, pp. 877-892, 2008.
- [15] Liu, T. Y., Xu, J., Qin, T., Xiong, W., and Li, H. LETOR: Benchmarking learning to rank for information retrieval. In Proceedings of SIGIR Workshop on Learning to Rank for Information Retrieval, 2007.
- [16] Najork, M., Zaragoza, H., & Taylor, M. J., Hits on the web: how does it compare? In Proceedings of SIGIR, pp. 471-478, 2007.
- [17] Yager, R.R., On ordered weighted averaging aggregation operators in multi-criteria decision making. IEEE Transactions on Systems, Man and Cybernetics, Volume 18, Issue 1, 1988.
- [18] TREC, <http://trec.nist.gov/data/webmain.html>, 2004.
- [19] Robertson, S.E., Walker, S., Microsoft Cambridge at TREC-9: Filtering track, *The Ninth Text REtrieval Conference (TREC-9)*, pp. 73-86, 2000.
- [20] O'Hagan, M., Aggregating template rule antecedents in real-time expert systems with fuzzy set logic, In Proceedings of 22th Annual IEEE Asilomar Conference on Signals, Systems and Computers, pp. 681-689, 1988.
- [21] Filev, D. & Yager, R. R., On the issue of obtaining OWA operator weights, Journal of Fuzzy Sets and Systems, Volume. 94, Issue. 2, pp. 157-169, 1998.
- [22] Cho, J., Roy, S., & E. Adams, R., Page Quality: In Search of an Unbiased Web Ranking. In Proceedings of ACM International Conference on Management of Data (SIGMOD), 2005.

**علی محمد زارع بیدکی** کارشناسی خود را در رشته مهندسی کامپیوتر - سخت افزار از دانشگاه صنعتی اصفهان و کارشناسی ارشد خود را در همین رشته از دانشگاه تهران دریافت نموده است. او همچنین دکتری خود را در رشته مهندسی کامپیوتر-ترم افزار در سال ۱۳۸۸ از دانشگاه تهران دریافت نموده است. او هم اکنون استادیار دانشکده مهندسی برق و کامپیوتر دانشگاه یزد می باشد. فعالیت های پژوهشی او در زمینه بازیابی اطلاعات وب، پردازش زبان های طبیعی و شبکه های کامپیوتری می باشد.

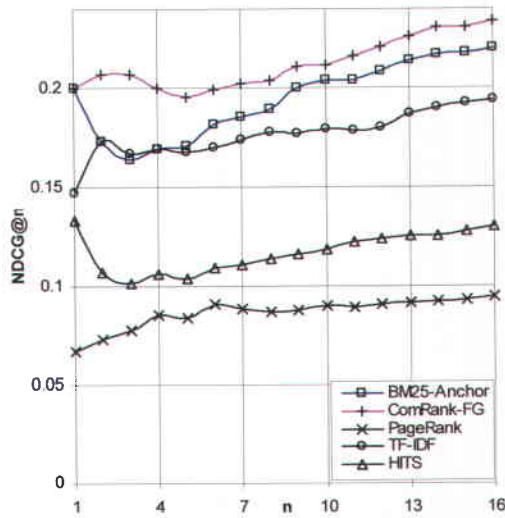


**محمد آزادنی** کارشناسی خود را در رشته مهندسی مخابرات از دانشگاه علم و صنعت و کارشناسی ارشد خود را در رشته مهندسی صنایع از دانشگاه صنعتی شریف دریافت نموده است. او از سال ۱۹۶۷ بعنوان پژوهشگر و مدیر پروژه و از سال ۱۳۷۹ به عنوان عضو هیئت علمی در مرکز تحقیقات مخابرات ایران فعالیت داشته است. فعالیت های پژوهشی او بیشتر در زمینه های سیستم های اطلاعاتی، مدیریت فناوری اطلاعات، پردازش زبان طبیعی و بازیابی اطلاعات وب می باشد.



**ناصر یزدانی** کارشناسی خود را در رشته مهندسی کامپیوتر از دانشگاه صنعتی شریف و دکتری خود را در همین رشته از دانشگاه کیس وسترن رزرو در آمریکا دریافت نموده است. او همچنین در شرکت ها و موسسه های تحقیقاتی مختلف آمریکا و ایران از جمله مرکز تحقیقات مخابرات ایران فعالیت داشته است. او هم اکنون دانشیار دانشکده مهندسی برق و کامپیوتر دانشگاه تهران می باشد. فعالیت های پژوهشی او در زمینه شبکه های کامپیوتری، بازیابی اطلاعات وب، سیستم های عامل و پایگاه های داده می باشد.





شکل (۹): مقایسه ComRank-FG با الگوریتم‌های دیگر در  $NDCG@n$

در این مقاله دو روش ترکیبی وقتی به نام‌های ComRank و ComRank-FG ارائه شده‌اند در مورد روش اول چندین الگوریتم رتبه بندی مشهور به عنوان ویژگی‌های درشت‌دانه و در روش دوم چندین ویژگی ریزدانه مانند TF و In-Degree با یکدیگر ترکیب شده‌اند. ترکیب روشها به دو صورت استفاده از رأی گیری وزن‌دار و عملگر تجمیع OWA انجام می‌گیرد. در روش رأی‌گیری وزن هر ویژگی (الگوریتم) متناظر با دقت آن تعیین می‌گردند. همچنین ویژگیها با استفاده از عملگر تجمیع OWA نیز ترکیب شده‌اند که در آن وزن‌ها به صورت پویا با استفاده از نظر افراد در رابطه با درجه‌ی ارتباط یک سند با پرسش یادگرفته می‌شوند.

برای ارزیابی و تست روش ارائه شده از بسته LETOR شامل مجموعه داده‌های محک کنفرانس TREC 2004 استفاده شده است. نتایج بدست آمده بهبود چشمگیری را در روش ترکیبی ارائه شده در مقایسه با بقیه الگوریتم‌ها نشان می‌دهد.

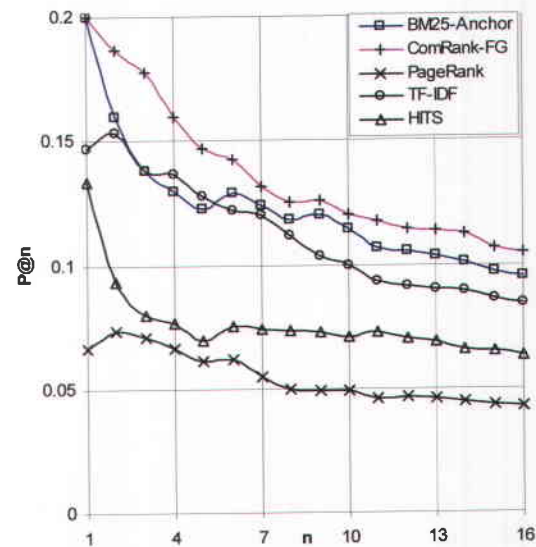
گرچه روش ترکیب‌کننده الگوریتم‌های مختلف دارای دقت بالایی می‌باشد ولی پیچیدگی آن نسبتاً زیاد است. در حالیکه الگوریتم ترکیبی ارائه شده مبتنی بر ویژگی‌های ریزدانه علاوه بر داشتن دقت مناسب دارای پیچیدگی پایینی است. علاوه بر آن روشهای ارائه شده به علت داشتن دقت بالا در چند نتیجه اول، برای محیط وب بسیار مناسب می‌باشند. نتایج فوق نشان دهنده اهمیت ترکیب ویژگی‌های محتوایی و ساختاری در روشهای رتبه بندی است. ارائه یک الگوریتم مناسب رتبه بندی با استفاده از شبکه‌های عصبی و الگوریتم‌های ژنتیک بر اساس ویژگی‌های ریزدانه، همچنین پیاده سازی الگوریتم رتبه‌بندی فوق در موتور جستجوی فارسی به عنوان کارهای آینده پیشنهاد می‌گردند.

## مراجع

- [۱] زارع‌بیدیکی ع. م.، آزادانیا م.، یزدانی ن. و کیهانی پور ا.ح. "الگوریتم ترکیبی وقتی جهت رتبه‌بندی صفحات وب"، چاپ شده در سیزدهمین کنفرانس سالانه انجمن کامپیوتر ایران، اسفند ۱۳۸۶.

ترکیب شده‌اند. لازم به ذکر است که هدف استفاده از In-Degree (درجه ورودی) صفحات به عنوان یک ویژگی می‌باشد، ولی با توجه در دسترس نبودن گراف، از این دو الگوریتم که به صورت برون‌خط با استفاده از درجه ورودی‌ها محاسبه شده‌اند، استفاده شده است. بنابراین هر سه ویژگی استفاده شده‌ی فوق ساده می‌باشند.

الگوریتم ترکیبی حاصل ComRank-FG نامیده شده و مانند روش ارائه شده در بخش ۴ عمل می‌نماید و نظر کاربران به عنوان تابع هدف در نظر گرفته شده است. بنابراین مانند بخش قبل فرآیند ارزیابی شامل دو بخش آموزش (۵۰ پرسش شامل ۴۰ هزار جفت پرسش و سند) و آزمون (۲۵ پرسش) می‌باشد.



شکل (۸): مقایسه الگوریتم ComRank-FG با الگوریتم‌های دیگر در  $P@n$

برای ارزیابی این الگوریتم با توجه به اینکه تمام ویژگی‌های الگوریتم BM25 برای متن در دست نیست از ویژگی‌های متن Anchor ها مانند TF استفاده شده است. همچنین با BM25 حاصل از متن Anchor عمل مقایسه انجام گرفته است.

در شکل‌های ۸ و ۹ الگوریتم ComRank-FG با بقیه الگوریتم‌ها مقایسه شده‌اند. همانطور که نشان داده شده است، الگوریتم ارائه شده از بقیه بهتر عمل کرده است. این در حالیست که دارای پیچیدگی پایینی نسبت به BM25 می‌باشد و وزن‌ها به صورت پویا محاسبه می‌شوند. علت اصلی این موفقیت ترکیب محتوا و ساختار وب به صورت مناسب بوده است. به عبارت دیگر با ترکیب ویژگی‌های محتوا و ساختار علاوه بر حل مشکلات موجود، جواب بهتری نیز ارائه گردیده است. بدین ترتیب ادعای فوق مبنی بر حصول یک الگوریتم دلخواه از طریق استفاده از ویژگی‌های ریزدانه ثابت می‌گردد.

## ۶- نتیجه گیری و کارهای آینده

مهمترین بخش یک موتور جستجو را واحد رتبه بندی تشکیل می‌دهد. در رتبه بندی، نتایج به ترتیب درجه ارتباط آنها با پرسش، مرتب شده و به کاربر ارائه می‌شوند. به عبارت دیگر هدف اصلی در رتبه‌بندی صفحات کشف کیفیت واقعی آنها می‌باشد. در حال حاضر روشهای مختلفی مبنی بر اتصال مانند PageRank و مبتنی بر محتوا مانند TF-IDF و BM25 ارائه شده اند.



پیچیدگی محاسباتی الگوریتم برابر  $O(mnr)$  می‌باشد که  $n$ ،  $m$  و  $r$  به ترتیب نشانگر تعداد جفت‌های سند و پرسش، تعداد ویژگی‌ها (۵ ویژگی) و تعداد تکرار هستند. مطابق آزمایشات انجام شده برای دقت  $0.001$  ده تکرار کافی می‌باشد. اگر چه محاسبه وزن‌ها به صورت بازگشتی انجام می‌شود ولیکن یک فرایند برون خط<sup>۴۳</sup> می‌باشد و در نهایت اثری روی زمان پاسخ به کاربران نخواهد داشت. به عبارت دیگر در ابتدا وزن‌ها با استفاده از داده‌های آموزش محاسبه شده و در سیستم جستجو با استفاده از وزن‌های محاسبه شده از قبل عمل ترکیب صورت می‌گیرد. جنبه‌ی دوم ترکیب چند الگوریتم مختلف با یکدیگر می‌باشد که در سیستم‌ها و موتورهای جستجوی مختلف نیز استفاده شده است [25]. گرچه در ظاهر پیچیده به نظر می‌رسد ولیکن با توجه به اینکه اکثر اعضای ترکیب از قبل محاسبه شده اند پیچیدگی این روش کمتر خواهد بود. برای مثال در روش  $ComRank$  که شامل ترکیب ۵ روش مختلف می‌باشد، الگوریتم‌های  $HostRank$  از قبل به صورت برون خط محاسبه شده‌است و همچنین اجزای بقیه‌ی روش‌ها که شامل  $TF$ ،  $IDF$  و درجه و روی می‌باشد نیز از قبل وجود دارند. بنابراین روش فوق به راحتی همانطور که در بقیه‌ی سیستم‌ها استفاده شده است قابل استفاده خواهد بود. همچنین روش ارائه شده در بخش بعد در مقایسه با  $ComRank$  پیچیدگی بسیار کمی دارد.

## ۵-۲- الگوریتم ترکیبی ریز دانه ( $ComRank-FG^{43}$ )

اگرچه با ترکیب روشهای مختلف رتبه بندی در بخش قبلی به نتایج مطلوبی رسیده‌ایم ولیکن این روش‌ها مشکلاتی هم دارند. مهمترین مشکل روش قبل پیچیده‌گی آن می‌باشد. به عبارت دیگری لازم است چندین روش رتبه بندی برای رسیدن به یک روش محاسبه شوند. علت اصلی آن درش دانه بودن ویژگیهای استفاده شده است. بنابراین هدف پیدا کردن ترکیبی از ویژگیهای ریزدانه برای رتبه بندی می‌باشد تا علاوه بر رتبه بندی مناسب دارای پیچیدگی کمتری باشند. بدین منظور از ویژگیهای ریزدانه را می‌توان فرکانس کلمه ( $TF$ )، (درجه‌ی ورودی یک صفحه‌ی وب)  $In-Degree$  و (طول سند)  $DL^{44}$  نام برد. در اینجا سعی داریم ترکیب مناسبی از آنها را ارائه نماییم.

هر تابع پیچیده را می‌توان به صورت یک تابع دو جمله‌ای نوشت. برای مثل قسمت اصلی الگوریتم  $BM25$  (بخش ۲)، عبارت است از

$$\frac{tf}{k + tf}$$

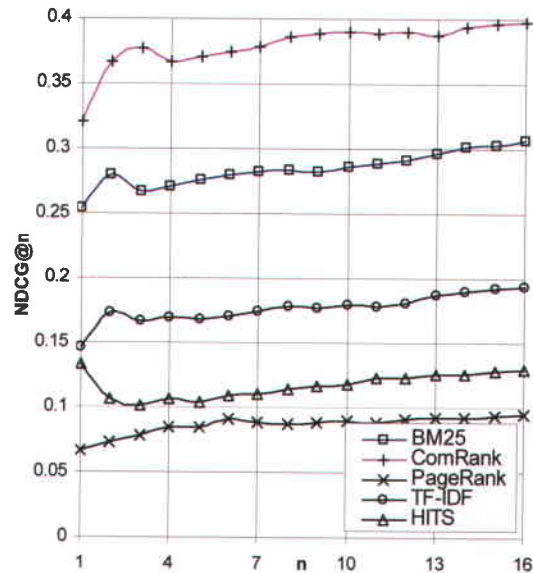
که با استفاده از سری تیلور میتوان آن را به صورت زیر بیان نمود:

$$f(x) = \frac{x}{k+x} = \sum_{i=1}^{\infty} (-1)^{i+1} \left(\frac{x}{k}\right)^i = x - \frac{x^2}{k^2} + \frac{x^3}{k^3} - \frac{x^4}{k^4} + \dots$$

به عبارت دیگر می‌توان با استفاده از ویژگیهای ریزدانه به یک الگوریتم رتبه بندی دلخواه با کارایی مناسب دست یافت.

در این بخش سه ویژگی  $TF$ ،  $PageRank$  و  $HostRank$  با استفاده از الگوریتم ارائه شده در شکل ۳ (که مبتنی بر  $OWA$  است) با یکدیگر

با تفاسیر فوق الگوریتم ترکیبی ارائه شده ( $ComRank$ ) با الگوریتم‌های  $PageRank$  و  $BM25$  مقایسه شده است. شکل‌های (۶) و (۷) به ترتیب مقایسه الگوریتم  $ComRank$  در مقایسه با دو الگوریتم  $BM25$  و  $PageRank$  را برای دو معیار  $P@n$  و  $NDCG$  را نشان می‌دهد. همانطور که در شکلها نشان داده شده است الگوریتم‌های ارائه شده در مقایسه با روشهای موجود کارکرد بهتری دارند. خصوصاً در  $n$  های کوچک مقدار آن حدود ۱۰٪ بهتر از  $BM25$  خصوصاً برای فاکتور  $NDCG$  می‌باشد.



شکل (۷): مقایسه  $ComRank$  با الگوریتم‌های دیگر در  $NDCG@n$

با توجه به اینکه کاربرهای وب معمولاً جوابهای اول را کلیک می‌کنند، الگوریتم فوق برای محیط وب بسیار مناسب خواهد بود.

مطابق شکل کارایی الگوریتم  $ComRank$  حدود ۴٪ کمتر از الگوریتم بوردا است. با توجه به اینکه روش بوردا به عنوان تابع هدف تعیین شده است، نتیجه فوق قابل قبول می‌باشد. بنابراین با داشتن یک مجموعه محک که برای هر جفت پرسش یک عدد بین صفر و یک را داشته باشد می‌توان به جواب‌های بهتری نیز دست یافت.

بنابراین از بحثهای فوق نتیجه می‌شود که از ترکیب چندین الگوریتم رتبه بندی می‌توان به یک روش جدید با کارایی بالایی دست پیدا کرد. جدول ۴ میانگین دقت ( $MAP$ ) را برای روشهای مختلف ارائه شده نشان می‌دهد.

جدول (۴): میانگین دقت ( $MAP$ ) الگوریتم‌های رتبه بندی.

نام الگوریتم	میانگین دقت برای همه پرسشها
Borda	۰/۲۶
BornaNorm	۰/۳۳
ComRank	۰/۳۰
BM25	۰/۲۴
PageRank	۰/۰۶

## ۵-۱- پیچیدگی الگوریتم ارائه شده

پیچیدگی الگوریتم ترکیبی ارائه شده را می‌توان از دو جنبه محاسبه‌ی وزن‌ها و استفاده از الگوریتم در یک سیستم جستجو بررسی کرد.

<sup>۱</sup> زمان اجرا روی یک  $PC$  با یک گیگابایت حافظه کمتر از سه دقیقه به طول انجامید.





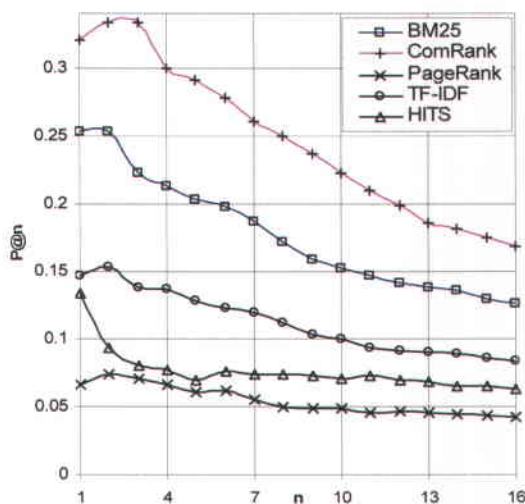
با توجه به اینکه محتوای وب داری پویایی بالایی می‌باشد، مهمترین مشکلی که روش فوق دارد ثابت بودن ضرایب ویژگیها (الگوریتم‌های رتبه بندی) می‌باشد (جدول ۳). بنابراین به ارزیابی روش ترکیبی ارائه شده در شکل (۳) که مبتنی بر یادگیری است و وزن ویژگی‌ها متغیر می‌باشند، می‌پردازیم.

مهمترین مسئله در یادگیری داشتن تابع هدف جهت پیدا کردن وزن‌های مناسب است. بهترین تابع درجه‌ی مرتبط بودن پرسش با سند از دید کاربر می‌باشد. با این تابع هدف با دو مشکل عمده مواجه می‌شویم. اولین مشکل از هم نوع نبودن درجه‌ی ارتباط با مقدار حاصل از ترکیب ویژگیها که همان الگوریتم‌های رتبه بندی هستند، به وجود می‌آید. به علاوه چونکه در مجموعه‌ی داده‌های موجود فقط دو سطح ۰ و ۱ توسط کاربر ارائه شده است تابع یادگیری با مشکل مواجه خواهد شد (جدول ۲). طبق آزمون‌های انجام شده، در حالتی که درجه ارتباط پرسشها و سند به عنوان تابع هدف در نظر گرفته شود، جواب بدست آمده از BM25 تا حدی بهتر خواهد بود. نتایج به علت نزدیک بودن به BM25 در اینجا ارائه نشده است.

با توجه به مشکل ذکر شده در بالا و موفقیت روش رأی‌گیری وزن‌دار ذکر شده (بوردا)، این روش (بوردا) به عنوان تابع هدف در نظر گرفته شده است. هدف از الگوریتم پیدا کردن وزن‌های متناظر به هر الگوریتم به صورت پویا می‌باشد. در این حالت الگوریتم ComRank نامیده می‌شود. در فعالیت قبلی انجام شده با توجه به اینکه روش نرمال‌سازی متفاوت بود، لذا در نظر گرفتن مقادیر بوردا به صورت مستقیم جواب مناسبی تولید نمود. بنابراین در آنجا تمام لیست‌ها مرتب شده، و ترتیب جواب‌های حاصل روش بوردا به عنوان تابع هدف در نظر گرفته شدند [1]. به عبارت دیگر چونکه مهمترین مسئله در بازایی اطلاعات ترتیب جواب‌ها هستند نه مقدار آنها، لذا جواب مناسبی در مقایسه با بقیه‌ی روشها حاصل گردید. خواهیم دید روش جدید ارائه شده در اینجا نسبت به حالتی که ترتیب مقادیر را در نظر بگیریم بهتر عمل می‌کند.

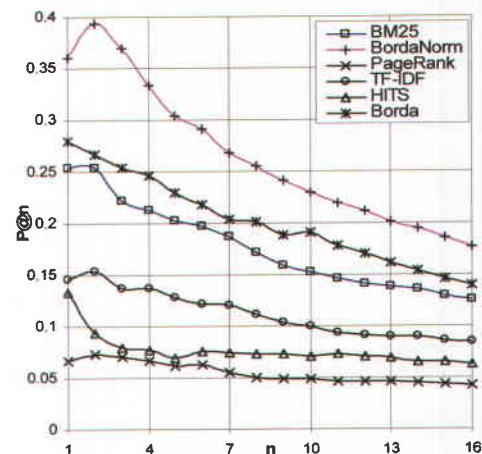
در ابتدا با توجه به اینکه الگوریتم ذکر شده مبتنی بر OWA است در ابتدا بردار W محاسبه می‌شود. ویژگیهای در نظر گرفته شده جهت ترکیب شامل الگوریتم‌های رتبه بندی HostRank، BM25، ComLink، و Sitmap، TF-IDF می‌باشند. بردار W بدست آمده با استفاده معادله ۱۳ (الگوریتم یادگیری وزن‌های OWA) برابر است با:

$$W = \{0.0013, 0.62, 0.082, 0.29, 0.002\}$$



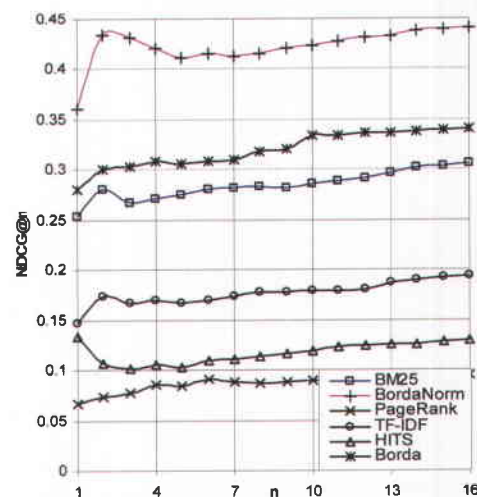
شکل (۶): مقایسه ComRank با الگوریتم‌های دیگر در P@n

در شکل‌های (۴) و (۵) الگوریتم ترکیبی بوردا در دو حالت بدون نرمال سازی و با نرمال سازی با بقیه روشها در دو معیار P@n و NDCG مقایسه گردیده‌اند.



شکل (۴): مقایسه الگوریتم بوردا با الگوریتم‌های دیگر در P@n در دو حالت نرمال سازی و بدون نرمال سازی

همانطور که دیده می‌شود روش رأی‌گیری بوردا که در آن الگوریتم‌های مختلف با استفاده از فرمول (۱۷) نرمال شده‌اند نسبت به بقیه افزایش چشمگیری داشته است. برای مثال در اولین نتیجه (n=1) و دومین نتیجه (n=2) به ترتیب حدود ۱۱٪ و ۱۵٪ نسبت به BM25 بهتر عمل کرده‌اند و نسبت به نتایج بدست آمده در مقاله قبلی [1] حدود ۵٪ بهتر است. همانطور که ذکر شد، علت آن همگن (نرمال) شدن آرگومانهای الگوریتم‌های استفاده می‌باشد. نکته‌ی دیگری که در شکل‌ها نمایان است افزایشی بودن معیار P@n (تقریباً) ثابت ماندن معیار NDCG با افزایش n است. شکل نمودار به صورت کلی به مجموعه‌ی داده‌ها و الگوریتم استفاده شده بستگی دارد. ولیکن در یک الگوریتم مؤثر با افزایش تعداد جواب‌ها (n) بایستی تعداد نتایج مرتبط کم شده و در نتیجه معیار دقت کاهشی و معیار NDCG (تقریباً) ثابت خواهند بود. با توجه به اینکه کاربران اغلب نتایج اول (۱۰ نتیجه) یک موتور جستجو را کلیک می‌کنند بنابراین هدف الگوریتم نیز داشتن دقت بالا در نتایج اولیه می‌باشد.



شکل (۵): مقایسه الگوریتم بوردا با الگوریتم‌های دیگر در NDCG@n در دو حالت نرمال سازی شده و بدون نرمال سازی



می‌شود و جمع وزن دار ویژگیها به عنوان خروجی در نظر گرفته خواهد شد. به عبارت دیگر نتایج بر حسب جمع وزن دار ویژگیها به صورت نزولی مرتب خواهند گردید.

جدول (۳): وزنهای متناظر هر الگوریتم جهت ترکیب که معادل دقت آنها در نظر گرفته شده است.

نام ویژگی	وزن (دقت)	توضیح
PageRank [12]	۰/۰۶	این الگوریتم مبتنی بر گراف وب می‌باشد و در گوگل استفاده شده است.
HostRank [27]	۰/۱۲	این الگوریتم رتبه بندی را با استفاده از ساختار سلسله مراتبی وب انجام می‌دهد.
TF-IDF [8]	۰/۱۴	این الگوریتم مبتنی بر متن Anchor می‌باشد (بهتر از متن جواب می‌دهد).
BM25 [9]	۰/۲۵	الگوریتم احتمالی مبتنی بر کل اسناد می‌باشد.
Sitemap[28]	۰/۲۵	الگوریتم ترکیبی اتصال و محتوا که بر روی سایتها عمل می‌کند.
ComLink [29]	۰/۲۸	الگوریتم ترکیبی اتصال و محتوا بر روی کل گراف وب می‌باشد.

حال مسئله‌ای که وجود دارد پیدا کردن وزن هر ویژگی جهت ترکیب می‌باشد. بدین منظور، مقدار  $P@n$  هر ویژگی در داده‌های آموزش (که از قبل محاسبه شده است) به عنوان وزن آن در نظر گرفته شده است. به عبارت دیگر با توجه به اینکه  $P@n$  درجه مناسب بودن هر ویژگی (الگوریتم) را نشان می‌دهد به نظر می‌رسد به عنوان وزن الگوریتم مورد نظر مناسب باشد. الگوریتمهای استفاده شده به عنوان ویژگی به همراه وزن نرمال شده آنها (در بازه صفر و یک) در جدول (۳) نشان داده شده است. در این جدول دقت نرمال شده هر الگوریتم (فرمول (۱۸)) به عنوان وزن در نظر گرفته شده است.

$$P@n\_weight_i = \frac{P@n_i}{\sum_{j=1}^6 P@n_j} \quad (18)$$

معادله ۱۹ ترکیب ویژگیهای ذکر شده را نشان می‌دهد که در آن اعداد فوق به عنوان وزن رأیگیری در نظر گرفته شده‌اند. ویژگی‌های در نظر گرفته شده متفاوت با ویژگیهای ارائه شده در مقاله [1] می‌باشند.

$$W = 0.06 * F_{PageRank} + 0.12 * F_{HostRank} + 0.14 * F_{TF-IDF} + 0.25 * F_{BM25} + 0.25 * F_{Sitemap} + 0.28 * F_{ComLink} \quad (19)$$

لازم به ذکر است که می‌توان ویژگیها را بدون نرمال سازی (Borda) و همچنین با نرمال سازی (BordaNorm) ترکیب کرد که هر کدام خصوصیات خاص خود را دارد. در حالت بدون نرمال سازی خصوصیات آرگومانهای هر الگوریتم مانند فاصله ذاتی آنها حفظ خواهد شد. در عین حال با توجه با اینکه محدوده‌ی الگوریتمهای مبتنی بر اتصال (بین ۰ و ۱) و محتوا (بین صفر تا ۳۰) متفاوت است، لذا ترکیب آنها بدون نرمال سازی مناسب نخواهد بود.

جدول (۲): نمونه‌ای از داده‌های LETOR شامل ویژگیها

Query	PageRank	TF-IDF	BM25	...	Relevancy
qid: 1	0.0023	9.23	19.31	...	1
qid: 1	0.0098	11.3	29.12	...	0
...	...	...	...	...	...
qid: 10	0.000	7.2	15.22	...	0
...	...	...	...	...	...
qid: 190	0.0076	12.23	30.1	...	1

برای ارزیابی و مقایسه دو معیار مناسب  $P@n$  و NDCG [26] استفاده شده است. معیار  $P@n$  نشانگر نسبت تعداد اسناد مرتبط در  $n$  سند ارائه شده به  $n$  است (معادله (۱۴)).

$$P@n = \frac{\# \text{ of relevant docs in top } n \text{ results}}{n} \quad (14)$$

با توجه به اینکه در  $P@n$  فقط مرتبط بودن یا نبودن یک سند در نظر گرفته می‌شود بنابراین از دقت کافی برخوردار نیست. بنابراین از معیار NDCG (معادله (۱۵)) به عنوان یک معیار تکمیلی نیز استفاده شده است. در این معادله که برای  $n$  نتیجه اول استفاده شده است  $r_j$  نشان دهنده درجه‌ی مرتبط بودن سند  $j$  با پرسش مربوطه می‌باشد.

$$NDCG(n) = \frac{\sum_{j=1}^n \frac{2^{r_j} - 1}{\log(1 + j)}}{\sum_{j=1}^n \frac{2^{r_j} - 1}{\log(1 + j)}} \quad (15)$$

پارامتر دیگری به نام میانگین دقت ( $AP^{rel}$ ) نیز در نظر گرفته شده و برای هر پرسش محاسبه می‌شود. مقدار آن با میانگین  $P@n$  برای تمام اسناد مرتبط با پرسش موردنظر (معادله (۱۶)) برابر می‌باشد. در این معادله اگر  $n$ -امین سند مرتبط باشد،  $rel(n)$  برابر ۱ و در غیر این صورت صفر خواهد بود.

$$AP = \frac{\sum_{n=1}^N (P@n * rel(n))}{\# \text{ total relevant docs for this query}} \quad (16)$$

همچنین پارامتری دیگری که برای ارزیابی استفاده می‌شود متوسط میانگین دقت ( $MAP^{rel}$ ) نام دارد و حاصل میانگین  $AP$ های تمام پرسشها می‌باشد.

با توجه به اینکه هر کدام از ویژگیها دارای حوزه مقادیر متفاوتی می‌باشند (جدول (۲))، برای یکدست شدن مقادیر ویژگیها به بازه صفر و یک عمل نرمال سازی به صورت زیر انجام شده است. در این رابطه  $f_k(q_i, d_j)$  نشان دهنده ویژگی  $k$  برای پرسش  $q_i$  و سند  $d_j$  می‌باشند. همچنین  $\min\{f_k(q_i, d_i)\}$

و  $\max\{f_k(q_i, d_i)\}$  به ترتیب نشان دهنده مقادیر کمینه و بیشینه ویژگی  $k$  و پرسش  $i$  در تمام اسناد می‌باشند.

$$f_k(q_i, d_j) = \frac{f_k(q_i, d_j) - \min\{f_k(q_i, d_i)\}}{\max\{f_k(q_i, d_i)\} - \min\{f_k(q_i, d_i)\}} \quad (17)$$

قبل از اینکه به چگونگی یادگیری وزنهای OWA بپردازیم، ابتدا روش رأیگیری ساده به نام بوردا<sup>۴</sup> را (ارائه شده در بخش قبل) مورد ارزیابی قرار می‌دهیم. در این روش به هر ویژگی یک وزن ثابت اختصاص داده



کار می‌رود. همچنین  $\hat{d}_k$  حاصل تجمیع داده‌های ردیف  $k$  با توجه به وزن‌های فعلی بدست آمده است.

$$(1) \text{ مساوی قرار دادن مقادیر اولیه تمام وزن‌ها } (w_i = \frac{1}{n}) \text{ و همچنین } t=0, \lambda_i(0) = 0, \beta = 0.3, \varepsilon = 0.001$$

$$(2) \text{ محاسبه لاندا در زمان } t+1: \lambda_i(t+1) = \lambda_i(t) - \beta w_i (b_{ki} - \hat{d}_k)(\hat{d}_k - d_k)$$

$$(3) \text{ محاسبه } w_i \text{ با استفاده از لاندا بدست آمده:}$$

$$w_i = \frac{e^{\lambda_i(t)}}{\sum_{j=1}^n e^{\lambda_j(t)}} \quad (13)$$

$$(4) \text{ اگر } (\hat{d}_k - d_k) > \varepsilon \text{ باشد به مرحله (2) می‌رویم.}$$

شکل (۳): الگوریتم یادگیری وزن‌های OWA

بنابراین پارامتر  $\lambda_i$  تعیین کننده وزن‌های OWA است که به وسیله پس‌خورد خطای  $(\hat{d}_k - d_k)$  بروزآوری می‌شود. الگوریتم فوق آنقدر تکرار می‌شود تا با خطای  $\varepsilon$  متوقف گردد. آزمایشات نشان داده شده است که نرخ یادگیری  $\beta = 0.3$  مناسب می‌باشد. همانطور که در شکل نشان داده شده است مقدار اولیه وزن‌ها برابر می‌باشند  $(1/n)$  که می‌توان به صورت تصادفی و یا دستی (با توجه به داش قبلی) نیز مقداردهی کرد که به کارهای آینده واگذار می‌گردد. با توجه به مطالب ذکر شده می‌توان خصوصیات الگوریتم فوق را در موارد زیر خلاصه کرد:

- مبتنی بر هر دو خصوصیت محتوایی و اتصالی
- مقیاس پذیری: به راحتی می‌توان الگوریتم جدید جهت ترکیب به آن اضافه کرد.
- وقتی بودن: ترکیب را بر اساس پرسش، اسناد و نظر کاربر انجام می‌شود. بنابراین بسته به محیط وزن‌ها نیز عوض می‌شوند.

## ۵- نتایج بدست آمده

برای ارزیابی و تست روش ارائه شده از بسته نرم‌افزاری LETOR [25] شامل مجموعه داده های محک برای یادگیری و ترکیب روشهای مختلف جهت رتبه بندی استفاده شده است. این مجموعه شامل ۴۴ ویژگی استخراج شده از Web TREC 2004, 2003 مانند PageRank, HITS, TF-IDF, BM25 و غیره برای ۷۵ پرسش مختلف می‌باشد (حدود ۶۵۰۰۰ ردیف داده برای جفت پرسش و سند). به علاوه نظر خبره (درجه‌ی مرتبط بودن) برای هر پرسش و سند نیز وجود دارد. نمونه‌ای از این داده در جدول (۲) نشان داده شده است. آزمایش انجام شده شامل دو مرحله‌ی آموزش (یادگیری) و آزمون می‌باشد. در مرحله آموزش از ۵۰ پرسش (شامل ۴۰ هزار جفت پرسش و سند)<sup>۱</sup> استفاده شده و مابقی (۲۵ پرسش) برای آزمون مورد استفاده قرار گرفته‌اند. بنابراین تمام معیارهای ارزیابی تعریف شده در زیر روی داده‌های آزمون انجام خواهد شد.

<sup>۱</sup> برای هر پرسش حدود ۸۰۰ سند وجود دارد.

ترکیب می‌باشد. با توجه به اینکه از اهداف مهم روشهای رتبه بندی داشتن دقت بالا است، ما برای هر الگوریتم دقت متناظرش را که توسط قضاوت‌های قبلی بدست آمده را به عنوان وزن آن در نظر گرفته‌ایم. معادله ۱۱ فرمول مربوطه را نشان می‌دهد که  $P_i$  و  $F_i$  به ترتیب نشان دهنده مقدار و دقت ویژگی  $i$  می‌باشند. در قسمت آزمون نشان داده خواهد شد که ترکیب فوق نتیجه چشمگیری داشته است و حتی از بهترین روشهای ارائه شده موجود نیز به مراتب بهتر عمل می‌کند.

$$W = P_1 * F_1 + P_2 * F_2 + \dots + P_n * F_n \quad (11)$$

بزرگترین مشکل روش بوردا ثابت بودن وزن‌ها می‌باشد. با توجه به اینکه محیط بازایابی اطلاعات کاملاً پویا و به محیط (هم بافت) وابسته است لذا لازم است وزن‌ها به صورت وقفی محاسبه شوند. یکی از روشهای ترکیب موفق پویا عملگر تجمیع OWA می‌باشد. بنابراین از عملگر تجمیع OWA برای ترکیب استفاده خواهد شد. برای بدست آوردن وزن‌های بردار OWA از مکانیزم یادگیری مناسب [24] بر مبنای قضاوت کاربر به صورت زیر استفاده شده است.

فرض کنید  $m$  ردیف اطلاعات در دسترس باشد که هر ردیف آن شامل یک بردار با  $n$  آرگومان به صورت  $A_k = (a_{k1}, a_{k2}, \dots, a_{kn})$  و بردار متناظر مرتب شده‌ی حاصل  $B_k = (b_{k1}, b_{k2}, \dots, b_{kn})$  است (آزمین عنصر  $b_{kj}$  متناظر با  $j$ -امین عنصر بزرگ در بردار  $A$  است) بعد از مرتب کردن  $A$ . در اینجا آرگومانهای هر ردیف متناظر با مقادیر الگوریتم‌های رتبه‌بندی مختلف به ازای یک پرسش و یک سند می‌باشند. بنابراین  $n$  الگوریتم رتبه‌بندی به عنوان ویژگی با  $m$  جفت پرسش و سند موجود خواهد بود. در بخش نتایج انواع الگوریتم‌های استفاده شده و چگونگی چینش در یک مجموعه آزمون توضیح داده شده است.

فرض کنید یک فرد خبره به هر ردیف اطلاعاتی (جفت سند و پرسش) یک مقداری به نام  $d_k$  اختصاص داده باشد (در اینجا نشانگر درجه مرتبط بودن هر پرسش با سند می‌باشد). این عدد در داده‌های استاندارد TREC [18] که در اینجا استفاده شده است شامل دو مقدار صفر و یک می‌باشد که صفر نشانگر نامرتب بودن سند و پرسش و یک نشانگر مرتبط بودن آنها است. این مقدار توسط کاربران خبره مقدار دهی شده است.

هدف ما پیدا کردن یک عملگر OWA یا بردار وزندار OWA به صورت  $W = (w_1, w_2, \dots, w_n)$  است که فرآیند تجمیع روی داده‌ها را مطابق نظر فرد خبره مدل کند. به عبارت دیگر یک تابع تجمیع با شرط زیر به ازای هر  $k$  مورد نیاز می‌باشد که هر ردیف اطلاعاتی مانند بردار  $A_k$  را به یک مقدار  $d_k$  نگاشت کند.

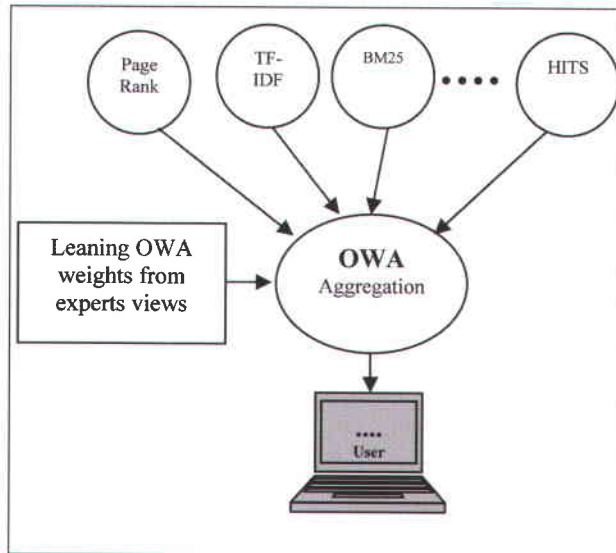
$$f(a_{k1}, a_{k2}, \dots, a_{kn}) = \sum_{i=1}^n w_i * b_{ki} = d_k \quad (12)$$

الگوریتم یادگیری وزن‌ها به صورت سناریوی تکراری در شکل (۳) نشان داده شده است. در [24] اثبات شده است که این الگوریتم با پیدا کردن مقادیر وزنی مناسب همگرا خواهد شد. مقدار  $0 \leq \beta \leq 1$  نشانگر نرخ یادگیری است و  $\lambda$  (لاندا) برای محاسبه  $w_i$  مطابق معادله (۱۳) به



دانه یک صفحه در نظر گرفت. در دو روش ارائه شده در این مقاله که به ترتیب ComRank و ComRank-FG نام دارند، روشهای درشت دانه و ریز دانه هر یک به صورت جداگانه برای رسیدن به راه حل مناسب تر ترکیب شده اند. در ادامه ابتدا ویژگیهای درشت دانه را با هم ترکیب کرده و بعد از ارزیابی آنها ترکیب ویژگیهای ریز دانه را مورد بررسی قرار می دهیم.

شکل (۲) نمای کلی الگوریتم جهت ترکیب چندین ویژگی را نشان می دهد.



شکل (۲): نمای کلی ترکیب روشهای رتبه بندی به عنوان ویژگیهای درشت دانه

همانطور که مشاهده می شود، الگوریتمهای مختلف مبتنی بر معیارهای محتوا و اتصال با استفاده از OWA که وزنهای آن با استفاده از نظر کاربر خبره بدست آورده شده است، ترکیب می شوند.

الگوریتمهای رتبه بندی موجود علاوه بر داشتن دقت پایین دارای مشکلات دیگری نیز می باشند. قرار گرفتن همیشگی صفحات محبوب<sup>۳۶</sup> در صدر لیست ارائه شده به کاربر، باعث می شود تا کاربر فقط صفحات خاصی را ببیند و در نتیجه صفحات تازه متولد شده با کیفیت بالا که کسی به آنها اشاره نمی کند نتوانند در دید کاربران قرار گیرند. این مشکل که "غنی تر شدن اغنیاء" نام دارد باعث می شود صفحات محبوب مرتباً محبوب تر شده و تعداد پیوند به آنها افزایش یابد. برای مثال الگوریتمهای مبتنی بر اتصال مانند PageRank از مشکل "غنی تر شدن اغنیاء"<sup>۳۷</sup> در رنج می برند و به پهنش رتبه بندی حساسیت کمتری دارند. در حالیکه الگوریتمهای مبتنی بر اتصال مانند BM25 و TF-IDF به مشکل فوق حساسیت نداشته ولی دارای مشکل پهنش رتبه می باشند. مسلماً ترکیب الگوریتمهای فوق و ارائه یک الگوریتم مبتنی بر محتوا و اتصال بصورت توأم باعث می شود تا مشکلات یکدیگر را تحت پوشش قرار دهند.

مهمترین مسئله باقیمانده چگونگی ترکیب ویژگیهای مختلف می باشد. جهت ترکیب از قضاوتهای قبلی کاربر با توجه به پرسشها و اسناد متناظر استفاده شده است.

ساده ترین روش، ترکیب وزن دار روشهای رتبه بندی مختلف می باشد که بوردا [23] دارد. در بوردا به هر الگوریتم یک وزن (حق رأی) داده می شود. بنابراین مهمترین مسئله در ترکیب، بدست آوردن وزن هر

به ازای  $\alpha = 1$  به بردار وزن  $W = [1 \ 0 \ \dots \ 0]^T$  می رسیم که دارای  $orness = 1$  است. به همین ترتیب به ازای  $\alpha = 0$  به بردار وزن  $W = [0 \ 0 \ \dots \ 1]^T$  می رسیم که دارای  $Orness = 0$  است. میزان Orness این عملگر میانگین گیری مرتب وزن دار به ازای مقادیر مختلف  $\alpha$  بصورت زیر می باشد:

$$Orness(W) = \frac{1}{n-1} \sum_{i=1}^n (n-i)w_i \quad (9)$$

یک نوع دیگر عملگر میانگین گیری مرتب وزن دار تحت عنوان بدبینانه را می توان با استفاده از وزنهای زیر تعریف نمود [17]:

$$\begin{aligned} w_1 &= \alpha^{n-1}, \\ w_2 &= (1-\alpha)\alpha^{n-2}, \\ w_3 &= (1-\alpha)\alpha^{n-3}, \\ &\dots \\ w_{n-1} &= \alpha(1-\alpha), \\ w_n &= (1-\alpha) \end{aligned} \quad (10)$$

این عملگر می تواند بصورت معکوس نیز تعریف گردد:

$$\begin{aligned} w_n &= (1-\alpha); \\ w_j &= w_{j+1}\alpha = w_j(1-w_n), \quad j = 2, \dots, n-1; \\ w_1 &= w_2(1-w_n)/w_n. \end{aligned}$$

می توان به فرآیند تهیه این وزنها بصورت دیگری نیز نگاه کرد. فرض کنید  $W$  برداری با بعد  $n$  از وزنها باشد که آنرا بصورت زیر به بردار  $V$  با بعد  $n+1$  تبدیل کنیم:

$$\begin{aligned} \forall i = 2, \dots, n : v_{i+1} &= w_i \\ v_2 &= (1-\alpha)w_1 = w_n w_1, \\ v_1 &= \alpha w_1 = (1-w_n)w_1. \end{aligned}$$

در مرجع [21] نشان داده شده است که این دو نوع عملگر، همه نیازمندیهای کلی عملگر میانگین گیری مرتب وزن دار را برآورده می سازند و نیز داریم:

$$\forall w_i \in [0,1], i = \{1, \dots, n\} \sum_{i=1}^n w_i = 1$$

#### ۴- الگوریتم ارائه شده (ComRank)

با توجه به مطالب ذکر شده در مقدمه، دقت الگوریتمهای فعلی کم می باشد و بسته به محتوا در شرایط خاص دارای جواب مناسب می باشند. بدین منظور سعی شده است با استفاده از ترکیب نتایج حاصل از خصوصیات یک صفحه به عنوان ویژگی<sup>۳۳</sup> با استفاده از عملگر تجمیع OWA الگوریتم ترکیبی جدیدی ارائه شود. با توجه به اینکه هدف، کشف صفحات با کیفیت بالا می باشد، بدیهی است که خصوصیات ذاتی صفحه نقش مهمی در کیفیت آن ایفا می کند. لذا این خصوصیات به عنوان ویژگیهای یک صفحه در نظر گرفته می شوند. خصوصیات صفحه را می توان در دو دسته ریز دانه<sup>۳۴</sup> و درشت دانه<sup>۳۵</sup> دسته بندی کرد. از خصوصیات ریزدانه صفحه می توان TF، IDF، تعداد پیوندهای ورودی (In-Degree) و اندازه صفحه (DL) را نام برد. همچنین الگوریتمهای رتبه بندی مختلف مبتنی بر اتصال و محتوا مانند PageRank، BM25 و TF-IDF را می توان به عنوان ویژگی درشت





مثل  $(x_1, x_2, \dots, x_n)$  را به یک عدد حقیقی چون  $y$ ، نسبت می‌دهد:

$$y = \text{Aggreg}(x_1, x_2, \dots, x_n)$$

در سال ۱۹۸۸ توسط آقای یانگر، انواع جدیدی از عملگرهای تجمیع موسوم به عملگرهای میانگین‌گیری مرتب وزن‌دار معرفی شد [17]. عملگر میانگین‌گیری مرتب وزن‌دار از بُعد  $n$ ، یک نگاشت  $F: \mathfrak{R}^n \rightarrow \mathfrak{R}$  با یک بردار وابسته  $n$

$$\text{تایی } W = [w_1, w_2, \dots, w_n] \text{ با شرط } \left\{ \begin{array}{l} w_j \in [0,1] \\ \sum_j w_j = 1 \end{array} \right\} \text{ می‌باشد}$$

بطوری که:

$$F(a_1, a_2, \dots, a_n) = \sum_{j=1}^n w_j b_j \quad (7)$$

که در آن،  $b_j$  برابر با  $j$  امین عنصر بزرگ در میان  $a_i$  ها می‌باشد. برای بدست آوردن وزن‌ها در عملگر میانگین‌گیری مرتب وزن‌دار، راه‌حل‌های متعددی از جمله توسط O'Hagan مطرح شده است که نیازمند حل یک مساله برنامه‌ریزی غیر خطی مقید<sup>۲۲</sup> است [20]. در کلاس عملگرهای میانگین‌گیری مرتب وزن‌دار نمایی، رابطه بسیار ساده‌ای بین درجه Orness و پارامتر تعیین‌کننده وزن‌های میانگین‌گیری مرتب وزن‌دار وجود دارد. درجه Orness نمایانگر این است که عملگر تجمیع تا چه حد به مقدار بیشینه یا کمینه نزدیک می‌شود. به عبارت دیگر، Orness نزدیک به یک نشانه نزدیکی تابع به عملگر MAX و بالعکس Orness نزدیک به صفر، نمایانگر نزدیکی تابع تجمیع به عملگر MIN است. بطور کلی دو روش برای محاسبه وزن‌ها وجود دارد. در روش خوش‌بینانه، وزن‌های میانگین‌گیری مرتب وزن‌دار بر اساس فرمول زیر محاسبه می‌شوند:

$$\begin{aligned} w_1 &= \alpha, \\ w_2 &= \alpha(1-\alpha), \\ w_3 &= \alpha(1-\alpha)^2, \\ &\dots \\ w_{n-1} &= \alpha(1-\alpha)^{n-2}, \\ w_n &= (1-\alpha)^{n-1} \end{aligned} \quad (8)$$

که در آن پارامتر  $\alpha \in [0,1]$  است. می‌توان این وزن‌ها را بصورت بازگشتی زیر نیز تعریف نمود:

$$\begin{aligned} w_1 &= \alpha \in [0,1] \\ w_j &= w_{j-1}(1-w_1), \quad j = 2, \dots, n-1 \\ w_n &= w_{n-1}(1-w_1)/w_1 \end{aligned}$$

می‌توان به فرآیند تهیه این وزن‌ها بصورت دیگری نیز نگاه کرد. فرض کنید برداری با بعد  $n$  از وزن‌ها باشد که آن را بصورت زیر به بردار  $V$  با بعد  $n+1$  تبدیل کنیم:

$$\begin{aligned} \forall i = 1, 2, \dots, n-1 : v_i &= w_i, \\ v_n &= \alpha w_n = w_1 w_n, \\ v_{n+1} &= (1-\alpha)w_n = (1-w_1)w_n. \end{aligned}$$

مساوی  $\gamma$  مقداردهی می‌شود. در آزمایشات انجام شده در این مقاله مقادیر  $k_1$ ،  $b$  با  $2/5$  و  $0/8$  و  $k_2$  و  $k_3$  با صفر مقداردهی شده‌اند [15]. همانطور که دیده می‌شود الگوریتم فوق از ویژگی‌های ریزدانه مانند فرکانس واژه (tf)، طول سند (dl) و غیره تشکیل شده است.

### الگوریتم ComLink

خانم شاکری در [29] یک روش رتبه‌بندی با استفاده از ترکیب پیوند و محتوا ارائه کرده است که در این مقاله ComLink نامیده شده است. در این مقاله یک مدل انتشار وابستگی<sup>۲۶</sup> بین صفحات ارائه شده است. در این مدل امتیازی<sup>۲۷</sup> به نام آبر وابستگی<sup>۲۸</sup>، برای هر سند  $p$  تعریف شده است که به سه پارامتر شباهت بین پرس‌وجو و سند  $p$  ( $S(p)$ )، جمع وزن‌دار آبر وابستگی صفحاتی که به  $p$  اشاره کرده و همچنین  $p$  به آنها اشاره می‌کند، بستگی دارد. ترکیب خطی این سه پارامتر به صورت زیر محاسبه می‌شود:

$$\begin{aligned} h(p) &= \alpha S(p) + \beta \sum_{p_i \rightarrow p} h(p_i) WI(p_i, p) + \gamma \sum_{p \rightarrow p_j} h(p_j) WO(p, p_j) \\ \alpha + \beta + \gamma &= 1 \end{aligned} \quad (5)$$

در معادله فوق،  $WI$  و  $WO$  به ترتیب وزن پیوندهای ورودی و خروجی صفحه  $p$  می‌باشند.

### الگوریتم SiteMap

خانم سانگ در مرجع [28] بر خلاف مدل فوق که گراف انتخاب شده مستقل از سایت است، یک روش ترکیبی انتشار مبتنی بر سایت به نام SiteMap ارائه کرده است. همچنین به جای انتشار امتیاز هر صفحه بر حسب میزان شباهت پرس‌وجو و صفحات، فرکانس تکرار واژه‌های هر پرس‌وجو در صفحات انتشار می‌یابد. ابتدا درخت یک وب سایت با استفاده از آنالیز URLها ساخته شده و سپس فرکانس واژه‌های پرس‌وجو از طریق پیوندها از پدر به فرزند (یا فرزند به پدر) در درخت سایت انتشار می‌یابد. معادله زیر روش محاسبه فرکانس تکرار واژه  $t$  را در صفحه  $p$  نشان می‌دهد.

$$f'_t(p) = (1+\alpha)f_t(p) + \frac{(1-\alpha)}{|Child(p)|} \sum_{q \in Child(p)} f_t(q) \quad (6)$$

در این معادله  $f_t(p)$  و  $f'_t(p)$  به ترتیب تعداد تکرار واژه  $t$  در صفحه  $p$  را قبل و بعد از انتشار نشان می‌دهند. این روش انتشار بر مبنای سایت<sup>۲۹</sup> نامیده می‌شود. روش فوق بازگشتی نبوده و فقط یکبار اجرا می‌شود (فقط اثر یک همسایه دیده می‌شود).

### ۳- عملگر تجمیع OWA

بصورت رسمی، تجمیع<sup>۳۰</sup> عبارتست از نگاشت یک  $n$  تایی<sup>۳۱</sup> از اشیاء عضو مجموعه‌ای خاص به یک شی از همان مجموعه. در خصوص عملگر ریاضی تجمیع، این مجموعه، همان مجموعه اعداد حقیقی است. در این حالت، این عملگر، تابعی است که یک  $n$  تایی از اعداد حقیقی



سند  $d$  می‌باشد). در فرمول (۳) TF-IDF واژه  $t$  را در سند  $d$  نشان می‌دهد.

$$TF-IDF_{t,d} = tf_{t,d} \cdot idf_t \quad (3)$$

$$tf_{t,d} = \frac{freq_{t,d}}{\max_{\ell} freq_{\ell,d}}$$

$$idf_t = \log \frac{N}{n_t}$$

متغیر IDF<sup>۲۴</sup> نشان‌دهنده عکس تکرار یک کلمه در کل اسناد می‌باشد. پارامتر  $N$  نشان‌دهنده تعداد اسناد و  $nt$  تعداد اسنادی که شامل کلمه  $t$  هستند را نشان می‌دهد.

### الگوریتم BM25

آقای رابرتسون [9] یک مکانیزم وزن گذاری مبتنی بر "۲-پواسون" به نام آکایی<sup>۲۵</sup> که دارای ویرایشهای متنوعی به صورت BMnn است را ارائه کرده است. نامگذاری این مکانیزم (۲-پواسون) به خاطر این است که توزیع هر واژه در یک مجموعه اسناد دارای توزیع پواسون می‌باشد (بدیهی است که تکرار هر واژه در هر مجموعه‌ای اسناد یک توزیع دوجمله‌ای است و در شرایطی که  $n$  زیاد و  $P$  کم باشد، توزیع دوجمله‌ای به پواسون تبدیل می‌شود). این روش احتمالی از بهترین روشهای رتبه‌بندی به شمار می‌رود که طبق آزمایشات انجام شده دارای دقت حدود ۲۵٪ می‌باشد [16].

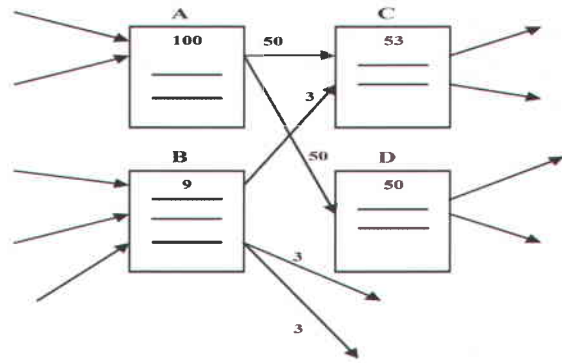
جدول (۱): متغیرهای استفاده شده در BM25

تعریف	نام متغیر
تعداد اسناد مرتبط به پرسش $Q$ شامل واژه $t$	$R$
تعداد اسناد مرتبط به پرسش $Q$	$R$
تعداد اسناد شامل واژه $t$	$n$
تعداد کل اسناد	$N$
فرکانس واژه $t$ در سند $d_i$	$tf$
فرکانس واژه $t$ در سند پرسش	$qtf$
میانگین طول اسناد	$avdl$
طول سند (تعداد واژه‌های در سند)	$dl$
تعداد واژه‌های در پرسش	$nq$
ثابت‌های قابل تنظیم (جهت tuning)	$K = k_1((1-b) + b(dl_i / dl_{avg}))$ $b, k_3, k_2, k_1, K$
شباهت میان سند $D$ و پرسش $Q$	$S(Q,D)$

در الگوریتم BM25 مقدار عددی حاصل برابر جمع وزن همه واژه‌های در پرسش در سند متناظر است و به شکل زیر محاسبه می‌شود (همه متغیرها در جدول (۱) توضیح داده شده‌اند).

$$S(Q,D) = \sum_{t \in Q} \left( \log \frac{(r+0.5)/(R-r+0.5)}{(n-r+0.5)/(N-n-R+r+0.5)} \right) \frac{(k_1+1)tf}{K+tf} \frac{(k_2+1)qtf}{k_3+qtf} + k_2 |Q| \frac{avdl-dl}{avdl+dl} \quad (4)$$

لازم است پارامترهای ثابت در نظر گرفته شده در فرمول فوق برای رسیدن به دقت بالا به صورت مناسب مقداردهی شوند. آزمایشات نشان می‌دهد که بسته به مجموعه‌ای داده استفاده شده بعضی از پارامترهای فوق متفاوت خواهند بود. طبق گفته‌ی آقای رابرتسون [19]  $k_1$  و  $b$  به صورت پیش فرض  $1/2$  و  $0.75$  مقداردهی می‌شوند.  $k_2$  و  $k_3$  اغلب صفر مقداردهی می‌شوند ولیکن برای پرسشهای با طول بزرگ،  $k_3$  اغلب



شکل (۱): مثالی از PageRank [Page98]

فرمول فوق برای حالتی که گراف کاملاً پیوسته باشد (هر گره به تمام گره‌ها دسترسی داشته باشد) مناسب است. در صورتیکه گراف وب پیوسته نبوده یا صفحات بدون ورودی یا خروجی<sup>۱۸</sup> موجود باشند، الگوریتم مذکور دچار اشکال می‌گردد (الگوریتم همگرا نخواهد شد). به عبارت دیگر بعد از اجرای کامل الگوریتم تعداد زیادی از صفحات دارای مقدار PageRank صفر خواهند بود. برای حل این مشکل از پارامتر  $d$  به نام ضریب استهلاک<sup>۱۹</sup> به صورت زیر استفاده شده‌است که در آن  $n$  نشان‌دهنده تعداد صفحات وب است.

$$r(i) = (1-d)/n + d * \sum_{j \in B(i)} r(j)/N(j) \quad (2)$$

بنابراین هر صفحه به تمام صفحات با احتمال  $(1-d)/n$  یک پیوند خواهد داشت. مقدار  $d$  در الگوریتم استفاده شده در گوگل  $0.85$  مقداردهی شده است. بدیهی است که رتبه‌بندی حاصل از PageRank رابطه مستقیمی با درجه‌ورودی صفحات دارد. قابل توجه است که درجه‌ورودی، یک ویژگی ریزدانه است.

### الگوریتم HostRank

الگوریتم HostRank [27] که مانند PageRank مبتنی بر اتصال می‌باشد در دو ساختار سلسله مراتبی و اتصالی وب را در رتبه‌بندی لحاظ می‌کند. در این روش صفحات ابتدا در یک ساختار سلسله مراتبی دایرکتوری، هاست<sup>۲۰</sup> و یا دامنه<sup>۲۱</sup> که گره برتر<sup>۲۲</sup> نامیده می‌شود قرار داده می‌شوند و عمل آنالیز اتصال بر روی گراف بدست آمده انجام می‌شود. سپس درجه‌ی اهمیت (ارزش) محاسبه شده‌ی هر گره بین صفحاتی که آن گره شامل آنها می‌شود توسط ساختار سلسله مراتبی توزیع می‌شود. نتایج آزمایشات روی TREC03 و TREC04 افزایش چشمگیری را نسبت به روش‌های دیگر مبتنی بر اتصال مانند PageRank و روش‌های دیگر مبتنی بر ساختار سلسله مراتبی نشان می‌دهند.

### الگوریتم TF-IDF

روش وزن‌دهی TF-IDF ارائه شده توسط آقای سالتون [8] از تکرار کلمات سند و پرس‌وجو برای محاسبه وزن استفاده می‌کند. هدف از TF<sup>۲۳</sup> یا تکرار کلمه این است که در صورتیکه یک کلمه چندین بار در یک سند ظاهر شود، آن کلمه سند را بهتر توصیف می‌کند. TF معمولاً با توجه به طول سند نرمال می‌شود و فرکانس کلمه بر فرکانس کلمه با بیشترین تکرار تقسیم خواهد شد ( $d$ , نشان‌دهنده TF کلمه  $t$  در

<sup>۱</sup> برای مثال yazduni.ac.ir, pweb.yazduni.ac.ir و yazduni.ac.ir/staff به ترتیب دامنه، هاست و دایرکتوری می‌باشند.



آزمایشات، افزایش قابل توجهی را در مقایسه با بقیه الگوریتم‌ها نشان می‌دهد. الگوریتم ارائه شده‌ی بر مبنای ویژگی‌های درشت‌دانه دارای میانگین دقت حدود ۳۰٪ می‌باشد که ۶٪ بهتر از بهترین الگوریتم فعلی یعنی BM25 می‌باشد. بزرگترین مشکل این الگوریتم پیچیدگی محاسباتی آن می‌باشد. در حالیکه الگوریتم ترکیبی مبتنی بر ویژگی‌های ریزدانه علاوه بر داشتن دقت مناسب، دارای پیچیدگی خیلی کمی می‌باشد.

این فعالیت در حقیقت ادامه کار انجام شده در [1] می‌باشد. در اینجا علاوه بر ارائه یک الگوریتم ترکیبی مبتنی بر ویژگی‌های درشت‌دانه کامل با جزئیات بیشتر و کارایی بهتر یک الگوریتم ریزدانه نیز ارائه شده است. ساختار مقاله بدین صورت می‌باشد که ابتدا در بخش بعدی مهمترین الگوریتم‌های پایه موجود بصورت خلاصه تشریح می‌شوند. در بخش (۳) کلیات الگوریتم OWA و خصوصیات آن شرح داده خواهد شد. جزئیات الگوریتم ارائه شده در بخش (۴) ارائه می‌شود و در بخش (۵) نتایج آزمایشات ارائه خواهد شد. در پایان بخش (۶) شامل نتیجه‌گیری و کارهای آینده می‌باشد.

## ۲- الگوریتم‌های رتبه‌بندی موجود

در این بخش الگوریتم‌های استفاده شده در مقاله را تشریح می‌کنیم.

### الگوریتم PageRank

الگوریتم PageRank یک الگوریتم مستقل از پرسش می‌باشد که در موتور جستجوی گوگل استفاده شده‌است و بر اساس اتصال بین صفحات عمل می‌کند. برای مثال اگر صفحه p1 به صفحه p2 اشاره کند، موضوع p2 برای ایجاد کننده p1 جذاب می‌باشد. بنابراین تعداد پیوندهای ورودی به یک صفحه درجه جذابیت آن صفحه برای دیگران را نشان می‌دهد. در نتیجه درجه جذابیت یک صفحه با تعداد پیوندهای ورودی آن افزایش می‌یابد. به‌علاوه وقتی به یک صفحه از صفحات مهم (با تعداد پیوند زیاد) اشاره شود، آن صفحه نیز رتبه‌ی بالایی خواهد داشت. به‌عبارت دیگر وزن هر صفحه در PageRank جمع وزن‌دار صفحاتی است که به آن اشاره می‌کنند. بنابراین الگوریتم PageRank بازگشتی بوده و با فرض برقرار بودن برخی شرایط می‌توان آن را با استفاده از زنجیره مارکف مدل کرد. فرض کنید  $N(i)$  و  $B(i)$  به ترتیب نشان‌دهنده تعداد پیوندهای خروجی<sup>۱۷</sup> و مجموعه صفحات ورودی صفحه  $i$  باشند. رتبه صفحه  $i$  با استفاده از PageRank به‌صورت زیر محاسبه می‌شود:

$$r(i) = \sum_{j \in B(i)} r(j) / N(j) \quad (1)$$

در نتیجه رتبه صفحه  $i$  مساوی جمع رتبه‌های صفحات ورودی تقسیم بر درجه خروجی آنها می‌باشد.

تقسیم ارزش یک صفحه بر درجه خروجی باعث می‌شود تا اولاً رتبه صفحه به صورت عادلانه بین فرزندان (خروجیها) تقسیم شود و ثانیاً جمع وزن‌های همه صفحات خروجی به عدد ثابت (یک) نرمال شود. شکل (۱) مثالی از PageRank را نشان می‌دهد. رتبه صفحات A که ۱۰۰ می‌باشد بین دو صفحه C و D تقسیم می‌شود. همچنین رتبه صفحه B نیز بین فرزندان که از جمله آنها C می‌باشد تقسیم می‌شود. بنابراین صفحه C، ۵۰ رتبه از صفحه A و ۳ رتبه از صفحه B دریافت می‌کند و در مجموعه دارای رتبه‌ی ۵۳ خواهد بود.

مهم از این مجموعه عظیم داده‌ای، موتورهای جستجو می‌باشند. لذا حدود ۸۰٪ از کاربران اینترنت سایتهای جدید را از طریق موتورهای جستجو کشف و بازدید می‌کنند [5]. به علاوه تعداد کاربرانی که از موتورهای جستجو استفاده می‌کنند نیز به طور دایم در حال افزایش است. برای مثال میانگین روزانه تعداد کاربران موتور جستجو در آمریکا از ۴۹/۳ میلیون در سپتامبر ۲۰۰۴ به ۶۰/۷ میلیون در سپتامبر ۲۰۰۵ (۲۳٪ رشد) افزایش یافته است [6].

بخش رتبه‌بندی<sup>۱۸</sup> یکی از مهمترین قسمت‌های موتور جستجو می‌باشد. "رتبه‌بندی" فرآیندی است که طی آن کیفیت صفحات از جنبه ارتباط با پرسش کاربر توسط موتور جستجو تخمین زده می‌شود. با توجه به اینکه به ازای هر پرسش کاربر معمولاً هزاران صفحه مرتبط وجود دارد لازم است آنها را اولویت بندی کرده و ۱۰ یا ۲۰ تای اول را به کاربر نشان دهد.

در حال حاضر دو روش عمده رتبه‌بندی مبتنی بر محتوا<sup>۱۹</sup> (استفاده شده در بازیابی اطلاعات سنتی) و مبتنی بر اتصال<sup>۲۰</sup> (استفاده شده در وب فعلی) وجود دارد.

در روشهای مبتنی بر محتوا که در بازیابی اطلاعات سنتی [7] استفاده می‌شوند، مدلهایی مانند بولی<sup>۲۱</sup>، احتمالی و فضای برداری<sup>۲۲</sup> جهت رتبه‌بندی اسناد بر مبنای محتوای آنها ارائه شده‌است. مهمترین این روشها در مدل برداری TF-IDF [8] و در مدل احتمالی BM25 [9] می‌باشند. به علت حجم زیاد محتوای وب مهمترین مشکل الگوریتم‌های بر مبنای محتوا پهنش رتبه<sup>۲۳</sup> [10] می‌باشد.

در حالیکه در حوزه بازیابی اطلاعات وب<sup>۲۴</sup> علاوه بر معیارهای فوق از ساختار اتصالی وب برای رتبه‌بندی استفاده می‌شود. پیوندها<sup>۲۵</sup> بیان‌کننده کیفیت محتوای یک صفحه از منظر صفحات بیرونی می‌باشد (بر خلاف محتوای متنی صفحه که کاملاً به ایجاد کننده آن وابسته است). به عبارت دیگر در رتبه‌بندی بر مبنای پیوند از محتوای صفحات دیگر برای ارزیابی یک صفحه استفاده می‌شود [11]. این خاصیت باعث می‌شود الگوریتم رتبه‌بندی با استفاده از اطلاعات استخراجی از پیوندها به مسائلی مانند پهنش حساسیت کمتری نشان دهد. از الگوریتم‌های مبتنی بر اتصال می‌توان PageRank [12] (اولین الگوریتم استفاده شده در گوگل)، HITS [13] و DistanceRank [14] را نام برد.

در مراجع [15, 16] نشان داده شده است که دو دسته الگوریتم‌های فوق دارای دقت<sup>۲۶</sup> پایین می‌باشند. همچنین نشان داده شده است که الگوریتم‌های مبتنی بر محتوا مانند BM25، الگوریتم‌های مبتنی بر اتصال و روشهای ترکیبی به ترتیب دارای دقت حدود ۲۵٪، ۱۴٪ و ۲۸٪ هستند.

با توجه با اینکه الگوریتم‌های ارائه شده در مجموع دارای کارایی پایینی می‌باشند و همچنین معمولاً هر الگوریتم در دامنه و محیط خاص خوب عمل می‌کند به نظر می‌رسد با استفاده از ترکیب الگوریتم‌های موجود بتوان علاوه بر حل مشکلات فعلی به کارایی بهتری نیز دست پیدا کرد.

در این مقاله هدف ارائه یک الگوریتم و فقهی ترکیبی برای دستیابی به دقت و کارایی بالاتر می‌باشد. این الگوریتم با استفاده از ترکیب الگوریتم‌های موجود به عنوان ویژگیهای درشت‌دانه، همچنین ترکیب ویژگیهای ریزدانه مانند TF، IDF و In-Degree به کمک فرآیند یادگیری سعی خواهد کرد به الگوریتم بهتری دست پیدا کند.

فرآیند یادگیری جهت ترکیب الگوریتم‌های مختلف با استفاده از OWA<sup>۲۷</sup> [17] انجام می‌شود. برای ارزیابی و مقایسه با سایر روشها از مجموعه داده محک TREC 2004 [18] استفاده شده است. نتایج



## الگوریتم ترکیبیِ وفقی جهت رتبه‌بندی صفحات وب با استفاده از ویژگیهای ریزدانه و درشت‌دانه

علی محمد زارع بیدکی	محمد آزادنیا	ناصر یزدانی	امیرحسین کیهانی‌پور
دانشگاه یزد	مرکز تحقیقات مخابرات ایران	دانشگاه تهران	دانشگاه تهران
دانشکده مهندسی برق و کامپیوتر	پژوهشکده فناوری اطلاعات	دانشکده مهندسی برق و کامپیوتر	دانشکده مهندسی برق و کامپیوتر
<a href="mailto:alizareh@yazduni.ac.ir">alizareh@yazduni.ac.ir</a>	<a href="mailto:azadnia@itrc.ac.ir">azadnia@itrc.ac.ir</a>	<a href="mailto:yazdani@ut.ac.ir">yazdani@ut.ac.ir</a>	<a href="mailto:keyhanipour@ut.ac.ir">keyhanipour@ut.ac.ir</a>

تاریخ دریافت: ۱۳۸۷/۱۲/۲۰ - تاریخ پذیرش: ۱۳۸۸/۶/۲۴

چکیده: حجم عظیم و پویا بودن اطلاعات وب، یکی از مهمترین چالش‌های بازیابی اطلاعات در پاسخ به پرسش کاربر می‌باشد. برای بهبود نتایج جستجو تاکنون الگوریتم‌های متنوعی مانند BM25 و PageRank ارائه شده‌اند. در این مقاله یک الگوریتم رتبه‌بندی وفقی<sup>۱</sup> ترکیبی برای دستیابی به دقت و کارایی بالاتر ارائه شده است. این الگوریتم با استفاده از ترکیب الگوریتم‌های رتبه‌بندی موجود به عنوان ویژگی درشت‌دانه مانند BM25 و TF-IDF و همچنین ترکیب ویژگیهای ریزدانه‌ی موجود مانند تکرار واژه‌ها<sup>۲</sup> و درجه ورودی صفحات<sup>۳</sup> به کمک فرآیند یادگیری به کارایی بهتری دست یافته است. در فرآیند یادگیری از عملگر تجمیع OWA<sup>۴</sup> و نظر افراد خبره در مورد درجه ارتباط پرسش و سند استفاده می‌شود. برای ارزیابی الگوریتم پیشنهادی از مجموعه داده‌های محک<sup>۵</sup> LETOR شامل داده‌های "WEB TREC 2004" استفاده گردید. آزمایشات، افزایش چشمگیری را در میزان دقت بازیابی نشان می‌دهند.

کلمات کلیدی: الگوریتم‌های رتبه‌بندی وب، رتبه‌بندی ترکیبی، عملگر تجمیع OWA، دقت

### ۱- مقدمه

۱۱/۵ میلیارد صفحه در سال ۲۰۰۵ [3] و ۲۶ میلیارد در ۲۰۰۸ رشد یافته‌است [4].  
برای اینکه کاربر به راحتی بتواند در این اقیانوس اطلاعاتی جستجو نموده و به اطلاعات مورد نظرش دسترسی پیدا کند، نیاز مبرمی به سامانه‌های بازیابی اطلاعات وجود دارد.  
در حال حاضر کاراترین ابزارها برای مدیریت، بازیابی و استخراج اطلاعات

توسعه و رشد نمایی وب باعث شده است تا با حجم عظیمی از اطلاعات (میلیاردها صفحه) شامل اسناد با فرمت‌های متفاوت (متن/صوت/تصویر...) در مکان‌های مختلف مواجه شویم. برای مثال حجم اطلاعات نمایه‌سازی شده در سال ۱۹۹۴ از ۱۱۰۰۰۰ صفحه [2] به

