

## سیستم پیشنهاد دهنده وب با استفاده از اطلاعات استفاده کاربران و پارتیشن بندی گراف

شهرزاد معتمدی مهر<sup>۱</sup>، مجید تاران<sup>۲</sup>، علی برادران هاشمی<sup>۳</sup>، محمدرضا میبدی<sup>۴</sup>

### چکیده

هدف سیستم‌های پیشنهاددهنده وب هدایت کاربران به سمت صفحاتی است که به بهترین وجه نیازها و علایق آنها را برآورده سازد. در این مقاله یک الگوریتم جدید مبتنی بر اتوماتای یادگیر توزیع شده و پارتیشن بندی گراف پیشنهاد می‌گردد. در الگوریتم پیشنهادی یک اتوماتای یادگیر توزیع شده بر اساس داده‌های استفاده کاربران از وب و گراف پیوند بین صفحات، شباهت صفحات یک سایت با یکدیگر را مشخص می‌کند. الگوریتم پیشنهادی همان الگوریتم HITS می باشد که در آن علاوه بر ساختار پیوند بین صفحات، رفتار کاربر در مشاهده این صفحات نیز در نظر گرفته شده است. برای این منظور از اتوماتای یادگیر توزیع شده برای یادگیری امتیازات Hub و Authority صفحات وب استفاده می گردد. الگوریتم پیشنهادی با ایجاد یک مدل مارکوف بر اساس اطلاعات فوق صفحات جدیدی را برای ادامه حرکت هر کاربر در سایت به وی پیشنهاد می‌دهد. نتایج آزمایشات انجام شده نشان می‌دهد که الگوریتم پیشنهادی در مقایسه با روش های گزارش شده مبتنی بر HITS و اتوماتای یادگیر توزیع شده از دقت بیشتری برخوردار است.

### کلمات کلیدی

اتوماتای یادگیر، داده کاوی استفاده از وب، سیستمهای پیشنهاد دهنده

## Web Recommendation using Web Usage Data and Graph Partitioning

Shahrzad Motamedi Mehr; Majid Taran; Ali B. Hashemi; M.R. Meybodi

### ABSTRACT

Recommendation systems aim at directing users toward the resources that best meet their needs and interests. One of the challenging tasks in improving web recommendation algorithms is the simultaneous use of users's activity log and hyperlink graph of the web site. In this paper, we propose a new recommendation algorithm based on web usage data and hyperlink graph of a web site. In the proposed algorithm, a distributed learning automata learns similarity between web pages of a web site using web usage data and hyperlink graph of the web site. The proposed approach is based on HITS algorithm in which in addition to link structure of web pages, the users' behavior in visiting these pages is also taken into consideration for Web Recommendation. A distributed learning automata is used to learn the hub and the authority scores of web pages. The proposed algorithm uses these information to build a Markov model which will be used to recommend new web pages for a user. Experiments show that the proposed method outperforms HITS algorithm and the only learning automata based method reported in the literature in terms of precision and coverage.

### KEYWORDS

Learning Automata, Web Usage Mining, Recommendations systems

<sup>۱</sup> دانشکده برق و مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه آزاد قزوین، قزوین، ایران، [motamedi@tmu.ac.ir](mailto:motamedi@tmu.ac.ir)

<sup>۲</sup> دانشکده برق و مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه آزاد قزوین، قزوین، ایران، [m\\_taran@isc.iranet.net](mailto:m_taran@isc.iranet.net)

<sup>۳</sup> دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیرکبیر، تهران، ایران، [a\\_hashemi@aut.ac.ir](mailto:a_hashemi@aut.ac.ir)

<sup>۴</sup> دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیرکبیر، تهران، ایران، [mmeybodi@aut.ac.ir](mailto:mmeybodi@aut.ac.ir)

## ۱. مقدمه

وب، محیطی وسیع، متنوع و پویا است که کاربران متعدد اسناد خود را در آن منتشر می کنند. وب طی یک فرآیند آشفته و غیر متمرکز رشد می کند و این روند منجر به تولید حجم وسیعی از مستندات متصل به یکدیگر گشته است که از هیچ گونه سازماندهی منطقی برخوردار نیستند. با توجه به حجم وسیع اطلاعات در وب، مدیریت آن با ابزارهای سنتی تقریباً غیر ممکن است و ابزارها و روش هایی نو برای مدیریت آن مورد نیاز است.

برای حل این مشکل، شخصی کردن وب به یک پدیده محبوب به منظور سفارشی کردن محیط های وب تبدیل شده است. هدف از سیستم های شخصی ساز فراهم کردن نیازهای کاربران، بدون اینکه به طور صریح آن ها را بیان کنند یا نشان بدهند، می باشد [۱۲]. شخصی سازی وب مجموعه ای از عملیات است که تجربه وب را برای یک کاربر خاص یا مجموعه ای از کاربران سازمان دهی می کند [۷].

روش های وب کاوی بر اساس آن که چه نوع داده ای را مورد کاوش قرار می دهند، به سه دسته داده کاوی محتوای وب<sup>۱</sup>، داده کاوی ساختار وب<sup>۲</sup> و داده کاوی استفاده از وب<sup>۳</sup> تقسیم می شوند [۱۴]. داده کاوی محتوای وب، فرآیند استخراج اطلاعات مفید از محتوای مستندات وب است. داده کاوی ساختار وب به کشف اطلاعات جدید با استفاده از پیوندهای<sup>۴</sup> بین صفحات وب می پردازد [۱۵] [۱۶] [۳۱]. داده کاوی استفاده از وب نیز داده های مربوط به استفاده کاربران از وب را مورد کاوش قرار می دهد و الگوهای استفاده از وب را به منظور درک و برآوردن بهتر نیازهای کاربران استخراج می کند. [۸] [۱۱] [۱۳]. بخش عمده ی فعالیت ها و تحقیقات انجام شده در وب کاوی به محتوای صفحات وب می پردازند. اما در سال های اخیر داده کاوی ساختار وب و داده کاوی استفاده از وب نیز مورد توجه قرار گرفته اند.

در [۷] از روش مبتنی بر قوانین انجمنی (AR)<sup>۵</sup> که با استفاده از کاوش آیت م های تکراری به دسته بندی صفحات می پردازد، استفاده شده است. الگوریتم ارائه شده در [۹] مبتنی بر آنالیز لینک ها می باشد که صفحات وب و کاربران سایت را به صورت گره و ابرپیوند مدل می کند و از الگوریتم HITS<sup>۶</sup> برای ارزیابی اهمیت آنها در گراف استفاده می کند و هدف آن اندازه گیری تخصص کاربران و اهمیت صفحات وب است. در [۱۰] دو متد مجزای رتبه بندی بر اساس آنالیز لینک ها ارائه داده شده است. Mobasher از درجه اتصالات بین صفحات سایت به عنوان فاکتوری تعیین کننده برای پیشنهاد بر اساس کاوش آیت م های تکرار شونده یا کشف الگوهای ترتیبی استفاده می کند [۱۲]. در سیستم های شخصی سازی که برای ارائه پیشنهادات فقط از رفتار کاربران استفاده می کنند با افزایش تعداد صفحات پیشنهادی، کارایی الگوریتم در حد قابل ملاحظه ای کاهش می یابد.

اتوماتای یادگیر توزیع شده قبلاً برای رتبه بندی صفحات وب [۳] [۴] و تعیین شباهت اسناد وب بکاربرده شده است [۵]. الگوریتم پیشنهادی مانند [۶] یک روش مبتنی بر اتوماتای یادگیر توزیعی است. نتایج بررسی ها نشان می دهد که الگوریتم های پیشین در تعداد صفحات بالا نسبت به الگوریتم پیشنهادی ما کارایی پایین تری دارد.

در [۱] و [۲] روش های جدیدی مبتنی بر اتوماتای یادگیر توزیع شده جهت تعیین شباهت بین صفحات ارائه شده است.

در [۱] برای محاسبه شباهت بین صفحات وب از اتوماتای یادگیر توزیع شده استفاده شده است که شباهت بین صفحات را با استفاده از فعل و انفعال کاربر و رابطه تراگذاری تعیین می کند. در [۲] اتوماتای یادگیر توزیع شده شباهت بین صفحات وب را با استفاده از اطلاعات پیمایش کاربران قبلی و پیوند گراف وب سایت تعیین می کند که نسبت به روش ارائه شده در [۱] از کارایی بالاتری برخوردار است در این روش با افزایش تعداد صفحات، مقدار دقت کم نمی شود. برای بالا بردن کارایی الگوریتم با تعداد صفحات زیاد از پارتیشن بندی گراف با استفاده از الگوریتم های چند سطحی<sup>۷</sup> استفاده شده است. در این روش به منظور کاهش اثر اطلاعات ناصحیح، کاربری که از پارتیشن خود خارج شود مسیر اشتباهی طی کرده و میزان شباهت محاسبه شده برای صفحات مسیر خارج از محدوده، با توجه به رابطه ای مشخص کاهش می یابد. جریمه دیگری که در روش پیشنهادی برای کاربر در نظر گرفته شده وجود دور در مسیر پیمایشی کاربر می باشد.

در این مقاله با استفاده از الگوریتم ارائه شده در [۲] شباهت صفحات وب تعیین می گردد. روش پیشنهادی با ترکیب داده های استفاده کاربران و داده های ساختاری صفحات وب الگوریتمی ترکیبی مبتنی بر اتوماتای یادگیر توزیع شده، پارتیشن بندی گراف و الگوریتم HITS به منظور رتبه بندی و پیشنهاد صفحات ارائه شده است. در روش پیشنهادی کاربر با صرف کمترین زمان به نتایج مطلوب خود دست می یابد.

در ادامه ابتدا در بخش ۲ اتوماتای یادگیر و اتوماتای یادگیر توزیع شده به اختصار معرفی می شوند. در بخش ۳ الگوریتم HITS و در بخش ۴ الگوریتم پیشنهادی ارائه می گردد. در بخش ۵ پس از معرفی مدل استفاده شده برای شبیه سازی، نتایج شبیه سازی ارائه و بررسی می گردد. بخش ۶ نتیجه گیری می باشد.

## ۲. اتوماتای یادگیر

اتوماتای یادگیر یک مدل انتزاعی است که بطور تصادفی یک اقدام از مجموعه متناهی اقدامهای خود را انتخاب کرده و بر محیط اعمال می‌کند. محیط اقدام انتخاب شده توسط اتوماتای یادگیر را ارزیابی کرده و نتیجه ارزیابی خود را توسط یک سیگنال تقویتی به اتوماتای یادگیر اطلاع می‌دهد. سپس اتوماتای یادگیر با اطلاع از اقدام انتخاب شده و سیگنال تقویتی، وضعیت داخلی خود را بروز کرده و اقدام بعدی خود را انتخاب می‌کند. شکل ۱ نحوه ارتباط بین اتوماتای یادگیر و محیط را نشان می‌دهد.



شکل ۱. ارتباط اتوماتای یادگیر با محیط

محیط را می‌توان توسط سه تایی  $E = \{\alpha, \beta, c\}$  نشان داد که در آن  $\alpha = \{\alpha_1, \alpha_2, \dots, \alpha_r\}$  مجموعه ورودیها،  $\beta = \{\beta_1, \beta_2, \dots, \beta_r\}$  مجموعه خروجیها و  $c = \{c_1, c_2, \dots, c_r\}$  مجموعه احتمالات جریمه می‌باشد. هرگاه  $\beta$  مجموعه دو عضوی باشد، محیط از نوع  $P$  می‌باشد. در چنین محیطی  $\beta_1 = 1$  به عنوان جریمه و  $\beta_2 = 0$  به عنوان پاداش در نظر گرفته می‌شود. در محیط از نوع  $Q$ ، مجموعه  $\beta$  دارای تعداد متناهی عضو می‌باشد و در محیط از نوع  $K$ ، تعداد اعضا مجموعه  $\beta$  نامتناهی است.  $c_i$  نشان دهنده احتمال نامطلوب بودن سیگنال تقویتی محیط در پاسخ به اقدام  $\alpha_i$  می‌باشد. در یک محیط ایستا<sup>۸</sup> مقادیر  $c_i$ ها ثابت هستند، حال آنکه در یک محیط غیر ایستا<sup>۹</sup> این مقادیر در طی زمان تغییر می‌کنند. بر اساس اینکه تابع بروز رسانی وضعیت اتوماتای یادگیر (که با اطلاع از اقدام انتخاب شده و سیگنال تقویتی  $\beta$ ، وضعیت بعدی اتوماتای یادگیر را محاسبه می‌کند) ثابت یا متغیر باشد، اتوماتای یادگیر به دو دسته اتوماتای یادگیر با ساختار ثابت و اتوماتای یادگیر با ساختار متغیر تقسیم می‌گردند [۱۷]. در این مقاله از اتوماتای یادگیر با ساختار متغیر استفاده شده است که در ادامه معرفی می‌شود.

اتوماتای یادگیر با ساختار متغیر توسط چهار تایی  $\{\alpha, \beta, p, T\}$  نشان داده می‌شود که در آن  $\alpha = \{\alpha_1, \alpha_2, \dots, \alpha_r\}$  مجموعه اقدامهای اتوماتای یادگیر،  $\beta = \{\beta_1, \beta_2, \dots, \beta_r\}$  مجموعه ورودیهای اتوماتای یادگیر،  $p = \{p_1, p_2, \dots, p_r\}$  بردار احتمال انتخاب هر یک از اقدامها و  $T, p(n+1) = T[\alpha(n), \beta(n), p(n)]$  الگوریتم یادگیری اتوماتای یادگیر می‌باشد. الگوریتمهای یادگیری متنوعی برای اتوماتای یادگیر ارائه شده است که در ادامه یک الگوریتم یادگیری خطی برای اتوماتای یادگیر بیان می‌گردد. فرض کنید اتوماتای یادگیر در مرحله  $n$ م اقدام  $\alpha_i$  خود را انتخاب نموده و محیط ارزیابی خود را توسط سیگنال تقویتی  $\beta(n)$  به اتوماتای یادگیر اعلام کند.

اتوماتای یادگیری که در بالا معرفی شد، دارای تعداد اقدامهای ثابتی می‌باشد. در بعضی از کاربردها به اتوماتای یادگیر با تعداد اقدام متغیر<sup>۱۰</sup> نیاز می‌باشد [۱۸]. یک اتوماتای یادگیر با تعداد اقدام متغیر، در لحظه  $n$  اقدام خود را از یک زیر مجموعه غیر تهی از اقدامها بنام مجموعه اقدامهای فعال  $V(n)$  انتخاب می‌کند. انتخاب مجموعه اقدامهای فعال اتوماتای یادگیر  $V(n)$  توسط یک عامل خارجی و بصورت تصادفی انجام می‌شود. نحوه فعالیت این اتوماتای یادگیر بصورت زیر است.

اتوماتای یادگیر برای انتخاب یک اقدام در زمان  $n$  ابتدا مجموع احتمال اقدامهای فعال خود  $K(n)$  را محاسبه و بردار  $\hat{p}(n)$  را مطابق رابطه (۱) ایجاد می‌کند. آنگاه اتوماتای یادگیر یک اقدام از مجموعه اقدامهای فعال خود را بصورت تصادفی و بر اساس بردار احتمال  $\hat{p}(n)$  انتخاب کرده و بر محیط اعمال می‌کند. در یک اتوماتای یادگیر با الگوریتم یادگیری خطی، اگر اقدام انتخاب شده  $\alpha_i$  باشد، اتوماتای یادگیر پس از دریافت پاسخ محیط، بردار احتمال  $\hat{p}(n)$  اقدامهای خود در صورت دریافت پاسخ مطلوب بر اساس رابطه (۲) و در صورت دریافت پاسخ نامطلوب طبق رابطه (۳) بروز می‌کند. سپس اتوماتای یادگیر بردار احتمال اقدامهای خود  $p(n)$  را با استفاده از بردار  $\hat{p}(n+1)$  و طبق رابطه (۴) بروز می‌کند.

$$K(n) = \sum_{\alpha_i \in V(n)} p_i(n)$$

$$\hat{p}_i(n) = \text{prob}[\alpha(n) = \alpha_i | \alpha_i \in V(n)] = \frac{p_i(n)}{K(n)} \quad (1)$$

$V(n)$  is the set of enabled actions

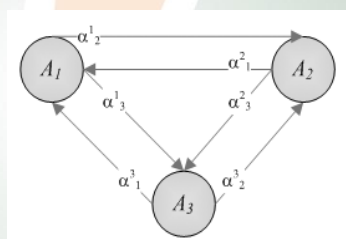
$$\begin{aligned} \hat{p}_i(n+1) &= \hat{p}_i(n) + a.(1 - \hat{p}_i(n)) \\ \hat{p}_j(n+1) &= \hat{p}_j(n) - a.\hat{p}_i(n) \quad \forall j \neq i \end{aligned} \quad (2)$$

$$\begin{aligned}\hat{p}_i(n+1) &= (1-b) \cdot \hat{p}_i(n) \\ \hat{p}_j(n+1) &= \frac{b}{\hat{r}-1} + (1-b) \hat{p}_j(n) \quad \forall j, j \neq i\end{aligned}\quad (3)$$

$$\begin{aligned}p_i(n+1) &= \hat{p}_i(n+1) \cdot K(n) & \text{for all } i, \alpha_i \in V(n) \\ p_j(n+1) &= p_j(n) & \text{for all } j, \alpha_j \notin V(n)\end{aligned}\quad (4)$$

## ۱.۲. اتوماتای یادگیر توزیع شده

اتوماتای یادگیر توزیع شده<sup>۱۱</sup> [۱۹] شبکه‌ای از چند اتوماتای یادگیر است که برای حل یک مساله مشخص با یکدیگر همکاری می‌کنند. یک اتوماتای یادگیر توزیع شده را می‌توان با یک گراف جهت‌دار مدل کرد. صورتی که مجموعه گره‌های آنرا مجموعه‌ای از اتوماتای یادگیر و یالهای خروجی هر گره مجموعه اقدامهای متناظر با اتوماتای یادگیر متناظر با آن گره است. هنگامی که اتوماتای یکی از اقدامهای خود را انتخاب می‌کند، اتوماتایی که در دیگر انتهای یال متناظر با آن اقدام قرار دارد، فعال می‌شود. بعنوان مثال در شکل ۲ هر اتوماتا ۲ اقدام دارد. اگر اتوماتای  $A_1$  اقدام  $\alpha_3$  خود را انتخاب کند، آنگاه اتوماتای  $A_3$  فعال خواهد شد. در گام بعد، اتوماتای  $A_3$  یکی از اقدامهای خود را انتخاب می‌کند که منجر به فعال شدن یکی از اتوماتاهای یادگیر متصل به  $A_3$  می‌شود. در هر لحظه فقط یک اتوماتای یادگیر در اتوماتای یادگیر توزیع شده فعال می‌باشد. بصورت رسمی، یک اتوماتای یادگیر توزیع شده با  $n$  اتوماتای یادگیر توسط یک گراف  $(A, E)$  تعریف می‌شود که  $A = \{A_1, A_2, \dots, A_n\}$  مجموعه اتوماتا و  $E \subset A \times A$  مجموعه لبه‌های گراف است بطوریکه لبه  $(i, j)$  متناظر با اقدام  $\alpha_j$  از اتوماتای  $A_i$  است. اگر بردار احتمال اقدامهای اتوماتای یادگیر  $A_j$  با  $p^j$  نشان داده شود، آنگاه  $p_m^j$  احتمال انتخاب اقدام  $\alpha_m$  از اتوماتای یادگیر  $A_j$  را نشان می‌دهد که احتمال انتخاب لبه خروجی  $(j, m)$  از میان لبه‌های خروجی گره  $j$  می‌باشد.



شکل ۲. اتوماتای یادگیر توزیع شده

## ۳. الگوریتم HITS

الگوریتم HITS، یکی از الگوریتم‌های رایج برای رتبه‌بندی صفحات وب بر اساس میزان ارتباط آنها با پرس وجوی کاربر است. HITS در سال ۱۹۹۹ توسط Kleinberg ارائه شد [۲۰]. این الگوریتم از دسته روشهای وابسته به پرس وجو<sup>۱۲</sup> است. در این نوع روشها، برای هر پرس وجو، تحلیل پیوند انجام می‌شود.

برای انجام تحلیل پیوند، ابتدا می‌بایست گراف خاص پرس وجو، به نام گراف همسایگی<sup>۱۳</sup> ساخته شود. سپس الگوریتم HITS، برای هر گره در گراف همسایگی، به طور تناوبی دو امتیاز، Authority و Hub را محاسبه می‌کند. در ادامه، گره‌ها با توجه به این امتیازات رتبه‌بندی می‌شوند. این الگوریتم فرض می‌کند، سندی که به اسناد بیشتری اشاره می‌کند Hub، خوبی است. همچنین، سندی که اسناد بیشتری به آن اشاره می‌کند Authority، خوبی می‌باشد.

## ۴. الگوریتم پیشنهادی

الگوریتم ارائه شده، صفحات وب جدید را بر اساس پیمایش کاربر، استفاده کاربران قبلی و اطلاعات پیوندی گراف سایت به کاربر جاری پیشنهاد می‌دهد. بدین منظور اتوماتای یادگیر توزیع شده شباهت بین صفحات وب را با استفاده از اطلاعات پیمایش کاربران قبلی و پیوند گراف وب سایت تعیین می‌کند [۲]. سپس یک رتبه‌بندی جدید صفحه مبتنی بر استفاده از وب، با بکاربردن شباهت تعیین شده محاسبه می‌شود. رتبه صفحه و شباهت بین صفحات محاسبه شده برای ایجاد مدل زنجیره مارکوف انتقال‌های کاربران در وب سایت استفاده می‌گردد. این مدل مارکوف برای پیش بینی احتمال ملاقات صفحات جدید برای پیشنهاد به کاربر استفاده می‌شود. جزئیات الگوریتم پیشنهادی در زیر آمده است.

## ۱.۴. بهبود الگوریتم HITS

در [۲۱] مجموعه‌ای از صفحات Hub و Authority به عنوان وب معرفی شده‌اند که تنها از پیوندهای بین صفحات وب برای تعیین صفحات Hub و Authority استفاده می‌کند، این روش از دقت کافی برخوردار نمی‌باشد. در [۲۲] علاوه بر پیوند بین صفحات وب از رفتار کاربران استفاده شده است که در نتیجه میزان تعداد صفحات نامرتبط کاهش می‌یابد. که با افزایش تعداد صفحات وب کارایی آن کاهش می‌یابد. در روش



پیشنهادی، محتوای صفحات وب در دو مرحله، ساخت مجموعه ریشه و اصلاح اجتماع وب مورد استفاده قرار می گیرد. یکی از مشکلات روش های کاوش استفاده از وب، اطلاعات ناصحیح می باشد. چرا که در برخی موارد، کاربران در وب سرگردان می شوند و بدون داشتن هدف مشخصی بر روی صفحات مختلف کلیک می کنند و گاهی اوقات کاربران به صفحه ای که پیمایش را آغاز کرده بودند بر می گردند. مراحل روش پیشنهادی به شرح زیر است:

- ایجاد مجموعه ریشه: در مرحله اول، موضوع اجتماع وب مورد نظر کاربر، به عنوان ورودی به الگوریتم ارائه می شود. سپس مجموعه ای از صفحات مرتبط با این موضوع انتخاب شده و مجموعه ریشه ساخته می شود.
- ایجاد مجموعه پایه: در این مرحله مجموعه ریشه که در مرحله قبل ایجاد شد، با استفاده از صفحاتی که اعضای مجموعه با آنها پیوند دارند، گسترش می یابد و مجموعه پایه را می سازند. برای این منظور ابتدا صفحاتی که صفحات مجموعه ریشه به آنها اشاره می کنند، به مجموعه پایه اضافه می شوند. سپس صفحاتی که به صفحات مجموعه ریشه اشاره می کنند، به این مجموعه اضافه می شوند.
- ایجاد اتوماتای یادگیر توزیع شده: در این مرحله برای هر یک از صفحات مجموعه پایه یک اتوماتای یادگیر ایجاد می شود. اعمال هر یک از این اتوماتاهای یادگیر، متناظر با صفحاتی است که صفحه جاری (صفحه مربوط به این اتوماتا) به آنها اشاره می کند. در ابتدا مولفه های بردار احتمال هر اتوماتا به صورت مساوی مقدار دهی اولیه می شوند.
- در این قسمت، الگوریتمی ترکیبی مبتنی بر اتوماتای یادگیر توزیع شده و پارتیشن بندی گراف به منظور تشخیص شباهت صفحات وب استفاده می گردد. در روش ارائه شده برای اسناد با تعداد بیشتر نتایج بهتری تولید می شود [۲].

## ۲.۴. پیشنهاد صفحات

در این مرحله پس از محاسبه امتیاز Hub و Authority صفحات وب مجموعه پایه، برای حل مشکل صفحات جدید و صفحات با فرکانس کم، به علت اینکه به آنها پیوندهای زیادی وجود نداشته، در فایل ثبت وقایع کاربران کمتر مشاهده شده اند و درجه اعتبار بالایی هم نخواهند داشت، بنابراین صفحات تنها بر اساس اهمیت آنها در ساختار سایت و با استفاده از معیار امتیاز Hub [۲۰] و امتیاز مرکزیت آنها رتبه بندی گردیده و هنگام ارائه پیشنهادها، با در نظر گرفتن صفحات بازدید شده توسط کاربر، فقط صفحات با بیشترین امتیاز مرکزیت (بسته به تعداد صفحات پیشنهادی از ۱ تا m صفحه) به عنوان صفحات پیشنهادی در نظر گرفته می شوند و به کاربر پیشنهاد می شوند [۲۳].

هدف از شخصی سازی بر اساس اطلاعات پیمایش کاربران محاسبه یک مجموعه پیشنهادی،  $rs$ ، برای نشست کاربر جاری می باشد [۲۴] [۲۵]. که بیشترین تطابق را با علایق کاربر داشته باشد. این جز تنها جز برخط سیستم بوده و باید از کارایی و دقت بالایی برخوردار باشد. فرض کنیم که کاربری که در حال گردش در سایت است و مسیر  $p_1 \rightarrow p_2 \rightarrow p_3 \rightarrow \dots \rightarrow p_k$  را پیموده است. تعداد آخرین صفحاتی را که توسط کاربر مشاهده شده و برای پیشنهاد صفحات جدید مورد استفاده قرار می گیرد را پنجره پیشنهاد می نامیم و اندازه آن را با  $rw$  نشان می دهیم که حداکثر برابر با تمام صفحات مشاهده شده و حداقل برابر با آخرین صفحه مشاهده شده می باشد. برای پیشنهاد صفحه  $P_{k+1}$  به کاربر از خاصیت مارکوف گراف استفاده می کنیم. طبق قاعده زنجیر مارکوف احتمال انتخاب مسیر در گراف دارای ویژگی مارکوف، از رابطه (۵) به دست می آید:

$$\Pr(p_1 \rightarrow p_2 \rightarrow p_3 \rightarrow \dots \rightarrow p_k) = \Pr(p_1) \times \prod_{i=2}^k \Pr(p_i | p_{i-1} \dots p_{i-m}) \quad (5)$$

به عنوان مثال، احتمال مسیر  $p_1 \rightarrow p_2 \rightarrow p_3$  برابر است با:

$$\Pr(p_1 \rightarrow p_2 \rightarrow p_3) = \Pr(p_1) \Pr(p_2 | p_1) \Pr(p_3 | p_2) = \Pr(p_1) \frac{\Pr(p_1 \rightarrow p_2) \Pr(p_2 \rightarrow p_3)}{\Pr(p_1) \Pr(p_2)} \quad (6)$$

که در آن  $\Pr(\bullet \rightarrow \bullet)$  برابر با احتمال گذار بین دو صفحه است و  $\Pr(\bullet)$  احتمال حالت پایدار صفحه متناظر می باشد که در دو بخش قبل به ترتیب در ماتریس  $p$  و بردار  $\vec{x}$  محاسبه شدند. برای پیشنهاد صفحه به کاربر، به ازای صفحات مختلف  $P_{k+1}$  که در مسیر  $p_1 \rightarrow p_2 \rightarrow p_3 \rightarrow \dots \rightarrow p_k$  ملاقات نشده اند، احتمال مسیر  $p_1 \rightarrow p_2 \rightarrow p_3 \rightarrow \dots \rightarrow p_k \rightarrow P_{k+1}$  را محاسبه می نماییم. احتمال هر مسیر امتیاز صفحه  $P_{k+1}$  را برای پیشنهاد به کاربر نشان می دهد. با مرتب کردن صفحات بر اساس امتیاز آنها، صفحاتی با بیشترین امتیاز به کاربر پیشنهاد می شود. برای هر یک از تعداد صفحات پیشنهادی دقیقاً فقط  $d$  صفحه (برابر با تعداد صفحات پیشنهادی) که بیشترین امتیاز را دارند، به کاربر پیشنهاد می شود.

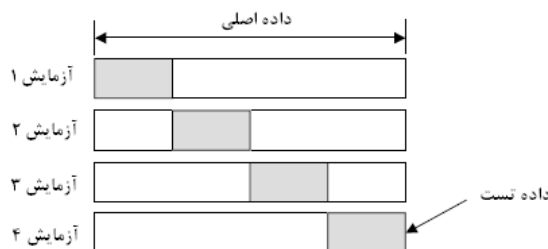
## ۵. ارزیابی الگوریتم پیشنهادی

در این قسمت ابتدا مدل بکار رفته برای تولید داده استفاده از وب و معیار ارزیابی تشریح می گردد. سپس نتایج آزمایشات الگوریتم پیشنهادی با الگوریتم های قبلی مقایسه می گردد.

### ۱.۵. مدل شبیه سازی

دو روش عمده برای ارزیابی الگوریتم‌هایی که از اطلاعات پیمایش کاربران استفاده می‌کنند وجود دارد. روش اول، استفاده از صفحات وب واقعی و داده‌های واقعی کاربران وب موجود در فایل‌های ثبت رخداد سایت‌ها می‌باشد. روش دوم مدل ارائه شده در [۲۶] می‌باشد. در این روش Liu و همکارانش نظم موجود در رفتارهای کاربران در محیط وب را با استفاده از یک مدل مبتنی بر عامل، مشخص و اعتبار مدل خود را با استفاده از چندین سایت وب بزرگ مانند مایکروسافت، تایید کرده‌اند. در این مقاله ما از داده‌های استاندارد سایت CTI DePaul استفاده می‌کنیم. این مجموعه داده اطلاعات نشست کاربران را به مدت ۲ هفته در سایت CTI DePaul در سال ۲۰۰۲ شامل می‌شود [۲۷]. این اطلاعات پیش‌پردازش شده و نشست‌های با اندازه ۱ و غیر استاندارد از آن حذف شده‌اند و در نهایت اطلاعات ۱۳۷۴۵ کاربر که از ۶۸۳ صفحه دیدن کرده‌اند در فایل‌های جداگانه قرار داده شده است.

برای انتخاب دو مجموعه آموزشی و تست از معیار  $K$ -Fold استفاده شده است [۲۸]. در این روش نمونه اصلی بطور تصادفی به  $k$  زیر نمونه تقسیم می‌گردد. برای هر  $k$  زیر نمونه، تنها یکی از زیر نمونه‌های بدست آمده به عنوان داده تست و  $k-1$  زیر نمونه باقیمانده به عنوان داده آموزشی استفاده می‌گردد. این پروسه به تعداد  $k$  بار تکرار می‌شود و به این ترتیب دقیقاً همه  $k$  زیر نمونه موجود، به عنوان داده تست استفاده می‌گردد. در نهایت برای تولید یک تخمین واحد، بین  $k$  نتیجه میانگین گرفته می‌شود. مزیت روش  $K$ -Fold این است که با تکرار زیر نمونه‌های تصادفی تمام مجموعه هم به عنوان داده آموزشی و هم داده تست استفاده می‌شوند. این روش در شکل ۳ نشان داده شده است.



شکل ۳. روش  $K$ -Fold

در اتوماتای یادگیر توزیع شده، احتمال انتخاب عمل  $z$  در اتوماتای  $i$ ، میزان ارتباط دو صفحه  $i$  ام و  $j$  ام را نشان می‌دهد. ساختار ارتباطی پیشنهادی با یک ماتریس  $n \times n$  به نام  $P$  بازنمایی می‌شود. در صورت فعال بودن عمل  $z$  در اتوماتای  $i$ ، درایه  $P_{ij}$  این ماتریس برابر با احتمال عمل  $z$  در اتوماتای  $i$  و در غیر اینصورت برابر با صفر قرار داده می‌شود. از آنجاییکه فرایند آموزش ترتیب توالی دسترسی‌ها را در نظر می‌گیرد نتایج فرایند یادگیری با یک ماتریس نامتقارن ( $p_{ij} \neq p_{ji}$ ) بازنمایی می‌شود. ماتریس نامتقارن تولید شده مبتنی بر اتوماتای یادگیر را ماتریس انتقال صفحات می‌نامیم. از حاصلضرب ماتریس نامتقارن  $P$  با ماتریس ترانهاده  $P^T$ ، آن، ماتریس متقارن جدیدی به نام ماتریس شباهت  $S$  طبق رابطه (۷) حاصل می‌شود.

$$S = P \cdot P^T$$

$$s_{ij} = \sum_k a_{ik} a_{kj}$$
(۷)

درایه  $s_{ij}$  در این ماتریس، درجه شباهت دو صفحه  $i$  ام و  $j$  ام را نشان می‌دهد.

## ۲.۵. معیار ارزیابی

برای ارزیابی، دو معیار "پوشش" <sup>۱۵</sup> و "دقت" <sup>۱۶</sup> معرفی می‌شوند. این دو معیار، بسیار شبیه به معیارهای متداول در بازیابی اطلاعات یعنی "فراخوانی" <sup>۱۷</sup> و "دقت" بازیابی اسناد هستند.

دقت پیشنهادها برابر با "نسبت پیشنهادهای درست به کل پیشنهادها" است. منظور از پیشنهاد درست، پیشنهادی است که با توجه به بخش دیده شده (پیشوند) یک جلسه کاربر تولید شده و در ادامه جلسه کاربر (پسوند) رخ دهد. اگر تعداد  $U$  جلسه کاربر را در نظر بگیریم، برای هر جلسه مثل  $u$  به ترتیب صفحات بازدید شده را، یک به یک، به مجموعه صفحات بازدید شده اضافه می‌کنیم. سپس، با دیدن هر صفحه  $p$  پیشنهادهایی تولید می‌کنیم. این مجموعه پیشنهاد را  $R(p)$  (مجموعه صفحات پیشنهاد شده پس از بازدید کاربر از صفحه  $p$ ) می‌نامیم. سپس  $R(p)$  با قسمت باقیمانده از جلسه کاربر، که آن را  $Tail(p)$  یا به اختصار  $T(p)$  می‌نامیم، مقایسه می‌شود. دقت پیشنهادات برابر با درصد اشتراک  $R(p)$  و  $T(p)$  خواهد بود و طبق رابطه (۸) محاسبه می‌گردد.

$$Precision = \frac{T(p) \cap R(p)}{R(p)}$$
(۸)

پوشش پیشنهادها، قدرت سیستم در پیش بینی تمام صفحاتی که ممکن است مورد نظر کاربران باشد را اندازه گیری می کند. این عدد برابر با "نسبت صفحات درست صفحات پیش بینی شده در ادامه جلسه یا  $T(p)$  به کل صفحات باقیمانده (تعداد صفحات  $T(p)$ ) در جلسه در هر قدم است" و طبق رابطه (۹) محاسبه می گردد.

$$Coverage = \frac{T(p) \cap R(p)}{T(p)} \quad (9)$$

هر چه مقدار دقت و پوشش بالاتر باشد، کارایی الگوریتم، مطلوب تر است. بررسی این وضعیت با استفاده از معیار  $F1$  ساده تر می باشد (رابطه (۱۰)).

$$F1 = \frac{2 \times Coverage \times Precision}{Coverage + Precision} \quad (10)$$

همانطور که رابطه (۱۰) نشان می دهد، هر چه مقدار دقت و پوشش بالاتر باشد، مقدار  $F1$  نیز افزایش می یابد.

### ۳.۵. تنظیمات

در آزمایشات، برای پارتیشن بندی گراف وب سایت از ابزار Metis [۲۹]، استفاده شده است که بر اساس الگوریتم های پارتیشن بندی چند سطحی [۳۰] کار می کند. تعداد پارتیشن ها  $k=25$  در نظر گرفته شده است. تعداد اعضای مجموعه ریشه و پایه به ترتیب ۱۰ و ۳۰ در نظر گرفته شده است.

همچنین برای آماده سازی داده تست و آموزشی مقدار پارامتر  $k$  در روش  $K-Fold$ ، ۲۰ در نظر گرفته شده است.

### ۴.۵. پارامترهای موثر در ارزیابی

اندازه طول پنجره پیشنهاد ( $rw$ ) و تعداد صفحات پیشنهادی که دو معیار "دقت" و "پوشش" برحسب آنها اندازه گیری می شوند، پارامترهای تاثیرگذار در کارایی الگوریتم هستند. در واقع براساس طول پنجره که مسیر پیمایش کاربر می باشد مسیر بعدی کاربر را پیش بینی می گردد. ابتدا با استفاده از مجموعه یادگیری الگوریتم اجرا شده و سپس بر اساس مقدار طول پنجره،  $rw$  صفحه متوالی را انتخاب کرده و به الگوریتم داده می شود. معیار ارزیابی رابطه معرفی شده در [۱۶] است. فرض کنیم مجموعه  $rp = \{x_{rw+1}, x_{rw+2}, \dots, x_{rw+|rs|}\}$  صفحات مشاهده شده توسط کاربر در ادامه نشست واقعی باشد. درجه شباهت مجموعه پیشنهادی و مجموعه صفحات واقعی از رابطه (۱۱) به دست می آید.

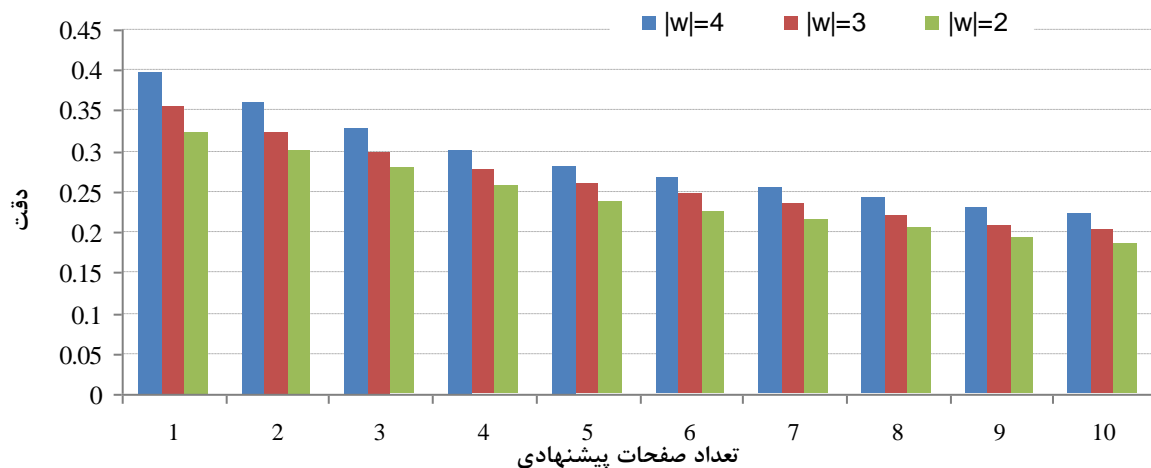
$$Sim(rs \cap rp) = \frac{rs \cap rp}{rp} \quad (11)$$

### ۵.۵. نتایج آزمایشات

الگوریتم پیشنهادی با روش های مبتنی بر اتوماتای یادگیر توزیع شده [۳][۴][۶][۱۲] که از یک روش آماری ساده نیز استفاده می کنند، مقایسه می گردد. در این روش آماری شباهت دو سند  $i$  و  $j$  بر اساس نسبت تعداد دفعاتی که کاربران از سند  $i$  به سند  $j$  حرکت کرده اند به تعداد دفعاتی که کاربران از سند  $i$  به هر سند دیگری مانند  $k$  حرکت نموده اند، طبق رابطه (۱۲) محاسبه می شود.

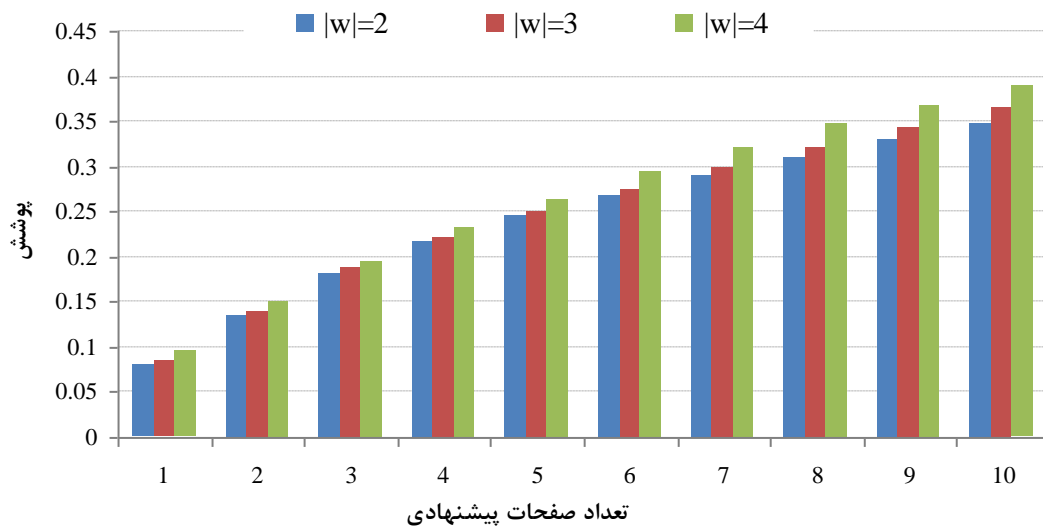
$$simpleSimilarity(i, j) = \frac{visited(i, j)}{\sum_{k=1}^n visited(i, k)} \quad (12)$$

کنفرانس داده کاوی ایران



شکل ۴. مقایسه دقت الگوریتم پیشنهادی با اندازه های متفاوت پنجره نسبت به تعداد صفحات پیشنهادی

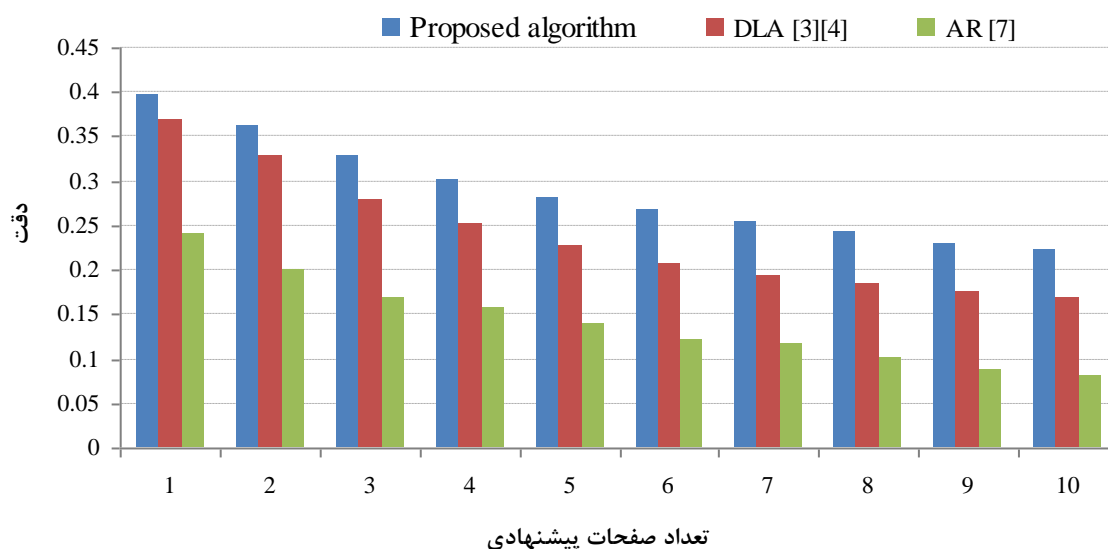
شکل ۴ نتایج دقت الگوریتم را نسبت به تعداد صفحات پیشنهادی مختلف برای ۳ حالت اندازه پنجره پیشنهاد نشان می دهد. در شکل مشخص است که دقت با تعداد صفحات پیشنهادی نسبت عکس دارد به طوری که با افزایش تعداد صفحات پیشنهادی، دقت کاهش می یابد.



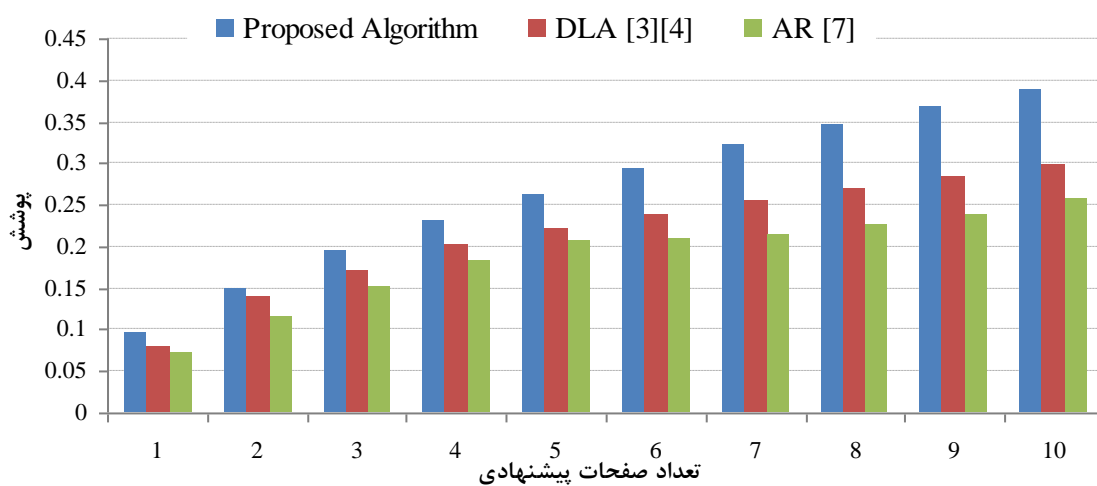
شکل ۵. مقایسه پوشش الگوریتم پیشنهادی با اندازه های متفاوت پنجره نسبت به تعداد صفحات پیشنهادی

شکل ۵ نتایج پوشش الگوریتم را نسبت به تعداد صفحات پیشنهادی مختلف برای ۳ حالت اندازه پنجره پیشنهاد نشان می دهد. پوشش با تعداد صفحات پیشنهادی نسبت مستقیم دارد به طوری که با افزایش تعداد صفحات پیشنهادی، پوشش افزایش می یابد. همانطور که در شکل های ۴ و ۵ مشاهده می شود، برای طول پنجره ۴ بالاترین دقت و پوشش را داریم. شکل ۶ و ۷ به ترتیب دقت و پوشش الگوریتم پیشنهادی با الگوریتم های مبتنی بر اتوماتای یادگیر توزیع شده [۳][۴] و قوانین انجمنی [۷] نسبت به تعداد صفحات پیشنهادی مختلف و با طول پنجره ۴، مقایسه شده است. همانطور که مشاهده می شود، الگوریتم پیشنهادی از دقت و پوشش بالاتری برخوردار است.





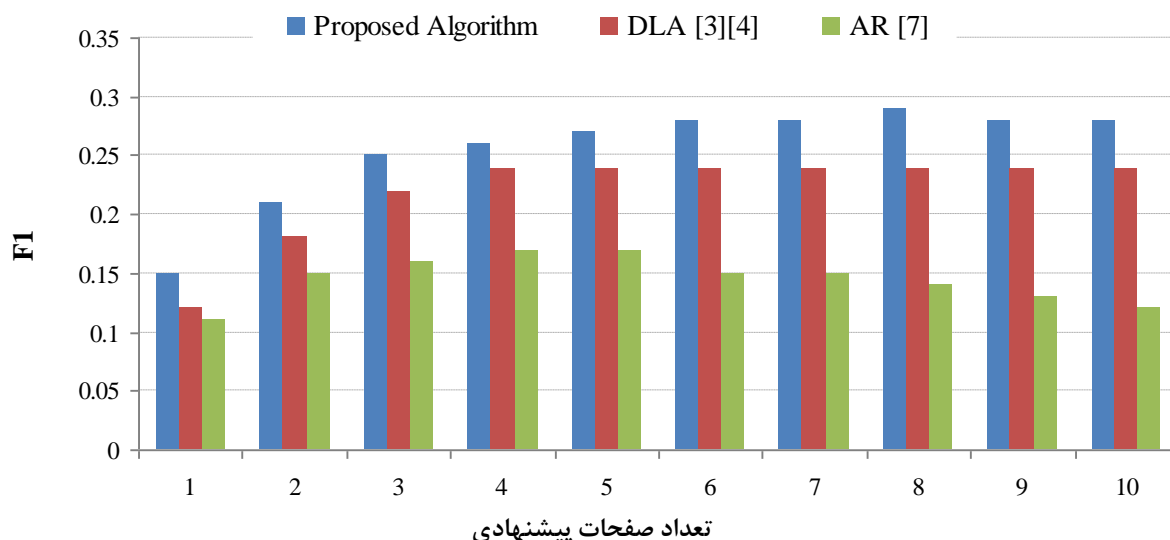
شکل ۶. مقایسه دقت الگوریتم پیشنهادی با الگوریتم های DLA [۳][۴] و AR [۷]



شکل ۷. مقایسه پوشش الگوریتم پیشنهادی با الگوریتم های DLA [۳][۴] و AR [۷]

شکل ۸ مقایسه الگوریتم پیشنهادی با الگوریتم های DLA [۳][۴] و AR [۷] براساس معیار  $F_1$  و با طول پنجره ۴ است و نشان می دهد که الگوریتم پیشنهادی نسبت به دو الگوریتم قبلی نتایج از دقت و پوشش بالاتری برخوردار است.

کنفرانس داده کاوی ایران



شکل ۸. مقایسه الگوریتم پیشنهادی با الگوریتم های DLA [۳][۴] و AR [۷] براساس معیار F1 و با طول پنجره ۴

## ۶. نتیجه گیری

در این مقاله الگوریتم ترکیبی جدیدی معرفی شده است که از اطلاعات پیمایش کاربران و پارتیشن بندی گراف به همراه الگوریتم HITS استفاده می کند. نتایج شبیه سازیها نشان داد که روش پیشنهادی در مقایسه با روش های گزارش شده مبتنی بر اتوماتای توزیع شده و قوانین انجمنی در تشخیص شباهت صفحات و پیشنهاد به کاربر از کارایی بالاتری برخوردار است. بصورتیکه دقت و پوشش بدست آمده با اندازه پنجره ۴، در الگوریتم پیشنهادی بیشتر از این مقدار در الگوریتم معرفی شده در [۳] [۴] و [۷] می باشد. الگوریتم ارائه شده همچنین مشکل صفحات جدید و صفات با فرکانس مشاهده کم که ارزش مشاهده شدن را دارند حل می کند و فرصت حضور در مجموعه صفحات پیشنهادی را فراهم می کند.

همچنین با بهبود الگوریتم HITS در الگوریتم ترکیبی پیشنهاد شده در این مقاله می توان برای بهبود ساختار استاتیک پیوندهای موجود بین صفحات سایت استفاده کرد.

## ۷. مراجع

- [۱]. SH. Motamedi Mehr, M. Taran, and M. R. Meybodi, "Web Usage Mining based on Distributed Learning Automata", Proceedings of ۱۰th Fuzzy Conference of Iran, Shahid Beheshti University, Tehran, Iran, ۲۰۱۰.
- [۲]. SH. Motamedi Mehr, M. Taran, A. B. Hashemi, and M. R. Meybodi, "Determining Web Pages Similarity Using Distributed Learning Automata and Graph Partitioning," International Symposium on Artificial Intelligence and Signal Processing (AISP) IEEE Iran Section, Tehran, Iran, ۲۰۱۱.
- [۳]. R. Forsati, M. R. Meybodi, "Web Personalization Using Distributed Learning Automata", Proceedings of ۳th Iran Information and knowledge Conference, Mashhad, Iran, ۲۰۰۷.
- [۴]. R. Forsati, M. R. Meybodi, "Effective page recommendation algorithms based on distributed learning automata and weighted association rules". Expert Systems with Applications: An International Journal, ۳۷, ۲ (۲۰۱۰), ۱۳۱۶-۱۳۳۰.
- [۵]. A. B. Hashemi, M. R. Meybodi, "Rating of document Using Distributed Learning Automata", Proceedings of ۱th Annual CSI Computer Conference of Iran, Shahid Beheshti University, Tehran, Iran, pp. ۵۵۳-۵۶۰, Feb. ۲۰-۲۲, ۲۰۰۶.
- [۶]. A. B. Hashemi, M. R. Meybodi, "Web Usage Mining Using Distributed Learning Automata", Proceedings of ۱th Annual CSI Computer Conference of Iran, Shahid Beheshti University, Tehran, Iran, pp. ۵۵۳-۵۶۰, Feb. ۲۰-۲۲, ۲۰۰۷.
- [۷]. B. Mobasher, H. Dai, T. Luo, M. Nakagawa, "Effective personalization based on association rule discovery from web usage data", Proceedings of the ۲rd ACM Workshop on Web Information and Data Management, ۲۰۰۱.
- [۸]. H. Dai, B. Mobasher, "Integrating Semantic Knowledge with Web Usage mining for Personalization", ۲۰۰۴.
- [۹]. J. Wang, Z. Chen, L. Tao, W. Ma, L. Wenyin, "Ranking User s Relevance to a Topic through Link Analysis on Web Logs", Proceeding of the WIDM ۰۳, ۲۰۰۳.
- [۱۰]. J. Borges, M. Levene, "Data Mining of User Navigation Patterns", in Revised Papers from the International Workshop on Web Usage Analysis and User Profiling, ۲۰۰۰, pp. ۹۲-۱۱۱.
- [۱۱]. L. Page, S. Brin, R. Motwani, T. Wingord, "The PageRank Citation Ranking: Bringing Order to the Web", Stanford University, ۱۹۹۸.
- [۱۲]. B. Mobasher, R. Cooley, J. Srivastava, "Automatic Personalization Based on Web Usage Mining", Communications of the ACM, Vol. ۴۳, No. ۸, ۲۰۰۰.
- [۱۳]. M. S. Aktas, M. A. Nacar, F. Menczer, "Personalizing PageRank Based on Domain Profiles", Proceeding of WEBKDD Workshop, Seattle, ۲۰۰۴.
- [۱۴]. B. Mobasher, H. Dai, Y. Sun, J. Zhu, "Integrating Web Usage and Content Mining for More Effective Personalization", Proceeding of the EC-WEB Conference, ۲۰۰۳.
- [۱۵]. M. Richardson, P. Domingos, "The Intelligent Surfer: Probabilistic Combination of Link and Content Information in PageRank", in Neural Information Processing System, ۲۰۰۲.
- [۱۶]. T. Haveliwalla, "Topic-Sensitive PageRank", in Proceedings of the ۱th International Conference on World Wide Web, New York: ACM Press, pp. ۵۱۷-۵۲۶, ۲۰۰۲.

- [14]. K. S. Narendra, M. A. L. Thathachar, "Learning automata: An introduction". Prentice Hall, 1989.
- [15]. M. A. L. Thathachar, R. Harita Bhaskar, "Learning Automata with Changing Number of Actions", IEEE Transactions on Systems Man and Cybernetics, vol. 19, no. 7, pp. 1092-1100, 1989.
- [16]. H. Beigy, M. R. Meybodi, "Utilizing Distributed Learning Automata to Solve Stochastic Shortest Path Problem". International Journal of Uncertainty, Fuzziness and Knowledge-based Systems, 14, 5 (2007), 591-614.
- [17]. J. Kleinberg, "Authoritative sources in a hyperlinked environment", Journal of the ACM, 47, 1999.
- [18]. S. Motiee, M. Meybodi, "Identification of Web Communities using Distributed Learning Automata", Proceedings of First Iran Data Mining Conference, IDMC'07, Tehran, Iran. 2007.
- [19]. D. Gibson, J. M. Kleinberg and P. Raghavan, "Inferring Web Communities from Link Topology", In Proc. of the 4th ACM Conference on Hypertext and Hypermedia. Pittsburgh, PA, pp. 220-224, 1994.
- [20]. T. Joachims, D. Freitag, T. M. Mitchell, "Web watcher: A tour guide for the world wide Web", Proceedings of International Joint Conference on Artificial Intelligence, 1997.
- [21]. B. Mobasher, H. Dai, T. Luo, M. Nakagawa, "Discovery and Evaluation of Aggregate Usage Profiles for Web Personalization", Data Mining and Knowledge Discovery, pp. 21-44, 2002.
- [22]. H. Liue, V. Keselj, "Combined mining of Web server logs and web contents for classifying user navigation patterns and predicting users' future requests", Data & Knowledge Engineering, 2007.
- [23]. J. Liu, S. Zhang, J. Yang, "Characterizing Web Usage Regularities with Information Foraging Agents", IEEE Transactions on Knowledge and Data Engineering, pp. 526-544, 2004.
- [24]. <http://maya.cs.depaul.edu/~classes/ect 544/data/cti-data.zip>
- [25]. G. J. McLachlan, K. A. Do, C. Ambroise, "Analyzing microarray gene expression data". Wiley, 2004.
- [26]. G. Karypis, V. Kumar, "METIS: A Software Package for Partitioning Unstructured Graphs, Partitioning Meshes, and Computing Fill-Reducing Orderings of Sparse Matrices". Minneapolis, City, 2004.
- [27]. G. Karypis, V. Kumar, "Multilevel k-way partitioning scheme for irregular graphs". Journal of Parallel and Distributed Computing, 48(1): 92-129, 1998.
- [28]. B. Mobasher, H. Dai, T. Luo, M. Nakagawa, "Using sequential and non-sequential patterns for predictive web usage mining tasks", Proceedings of the IEEE International Conference on Data Mining, Maebashi City, Japan, 2002.

- <sup>1</sup> Web Content Mining
- <sup>2</sup> Web Structure Mining
- <sup>3</sup> Web Usage Mining
- <sup>4</sup> Hyperlink
- <sup>5</sup> Association Rule
- <sup>6</sup> Hyperlink-Induced Topic Search
- <sup>7</sup> Multilevel
- <sup>8</sup> Stationary
- <sup>9</sup> Non-Stationary
- <sup>10</sup> Learning automata with changing number of actions
- <sup>11</sup> Distributed Learning Automata
- <sup>12</sup> Query Dependent Schemes
- <sup>13</sup> Neighborhood Graph
- <sup>14</sup> Transpose
- <sup>15</sup> Coverage
- <sup>16</sup> Precision
- <sup>17</sup> Recall

کنفرانس داده کاوی ایران