

۱۰۸۷

کاوش ساختار وب با استفاده از اتوماتای یادگیر و گرامر احتمالی

زهره اناری^۱؛ محمد رضا میبیدی^۲؛ بابک اناری^۳

چکیده

یکی از روش‌های داده کاوی در وب، کاوش استفاده از وب می‌باشد. هدف از این کاوش، استخراج اطلاعات مفید از داده‌هایی است که از تعامل کاربران در هنگام استفاده از وب به دست می‌آید. با استخراج این اطلاعات، می‌توان میزان ارتباط بین صفحات وب را تعیین و در نتیجه عملیاتی مثل خوشه بندی و رتبه بندی صفحات وب را انجام داد. هدف در این مقاله، تعیین میزان ارتباط بین صفحات وب با استفاده از گرامر احتمالی و اتوماتای یادگیر می‌باشد. در الگوریتم پیشنهادی ابتدا از مسیرهای حرکتی ذخیره شده کاربران در لاگ فایل‌ها، گرامر احتمالی تولید شده و قوانین انجمنی استخراج می‌گردند. قوانین انجمنی استخراجی توسط اتوماتای یادگیر مورد ارزیابی قرار گرفته و میزان ارتباط بین صفحات وب توسط اتوماتای یادگیر تعیین می‌گردد. نتایج شبیه‌سازی‌ها از طریق مقایسه با روشهای Bollen, AntWeb و E-DLA نشان می‌دهد که روش پیشنهادی از کارایی قابل ملاحظه‌ای برخوردار است.

کلمات کلیدی

کاوش استفاده از وب، اتوماتای یادگیر، گرامر احتمالی

^۱دستیار علمی، دانشکده کامپیوتر و فناوری اطلاعات، دانشگاه پیام نور، zanari323@yahoo.com

^۲استاد و عضو هیات علمی دانشگاه صنعتی امیرکبیر mmeybodi@aut.ac.ir

^۳عضو هیات علمی دانشگاه، دانشگاه آزاد اسلامی واحد شبستر، anari322@yahoo.com

Web Structure Mining using Learning Automata and Probabilistic Grammar

Zohreh Anari; Mohammad Reza Meybodi; Babak Anari

ABSTRACT

One of the data mining techniques on the web is web usage mining. The purpose of this mining, extracting useful information from the user interaction data that can be obtained when using the web. By extracting this information, we can determine the relationship between the web pages and then we can do operations such as clustering and ranking web pages. The purpose of this article is to determine the relationship between web pages, using probabilistic grammar and learning automata. In the proposed algorithm first probabilistic grammar generated and association rules are extracted from the user navigation paths stored in log files. Extracted association rules evaluated by learning automata and the relationship between web pages by learning automata is determined. The simulation results compared with Antweb, Bollen and E-DLA shows that the proposed method has a considerable efficiency.

KEYWORDS

Web usage mining, Learning Automata, Probabilistic Grammar

۱. مقدمه

استفاده از الگوهای حرکتی کاربران در بین صفحات وب، یک روش مهم برای کسب اطلاعات به منظور یاری رساندن کاربران وب در امر جستجو و حرکت در وب می باشد. استخراج اطلاعات از صفحات وب، به وسیله تکنیک های داده کاوی را کاوش وب می گویند. کاوش وب در سه سطح مطرح است. در سطح محتوا، در سطح ساختار و در سطح استفاده از وب. در سطح محتوا، هدف، کاوش محتوای وب می باشد. در سطح ساختار، هدف استفاده از توپولوژی ابر پیوندها برای تعیین ارتباط صفحات وب است. در سطح استفاده از وب، هدف، کشف اطلاعات مفید از داده - هایی است که از تعامل کاربران در هنگام استفاده از وب بدست می آید. بیشتر تکنیک های موجود برای کشف ساختار ارتباطی بین صفحات وب از اطلاعات موجود در لاگ فایل ها استفاده می کنند. لاگ فایل ها، فایل هایی هستند که مجموعه درخواست های کاربران به صفحات وب را در خود ذخیره می کنند. این اطلاعات به همان ترتیبی که به سرور می رسند، ذخیره شده و می توان دوباره آنها را بازیابی کرد. روش های موجود با کاوش در لاگ - فایل ها، یکسری اطلاعات آماری را از آنها استخراج می کنند که می توان از آنها به عنوان ابزاری برای کشف ساختار اسناد وب استفاده کرد. برخی از این روش ها عبارتند از: استفاده از قوانین یادگیری، مثل استفاده از قانون یادگیری هب در روش بولن [7]، استفاده از سیستم مورچه ها مثل [19] AntWeb، استفاده از اتوماتای یادگیر توزیع شده، [3,17,18]، استفاده از زنجیر مارکف [8] و استفاده از گرامرهای احتمالی ابرمتن [5]. ضعف روش گزارش شده مبتنی بر اتوماتای یادگیر توزیع شده (E-DLA) که بعداً در این مقاله مفصل تر به آن اشاره می شود، کارایی پایین آن در مقایسه با روش Ant Web است [2]. دو روش AntWeb و Bollen با وجود مزایایی که نسبت به روش های قدیمی تر دارند، به دلیل اینکه از ماتریس ارتباطات بین صفحات وب استفاده می کنند، برای استفاده در مجموعه های بزرگ و قابل گسترش مناسب نمی باشند. نقطه ضعف استفاده از زنجیره مارکف، نیاز به محاسبات بالا به واسطه محاسبه توان n ام ماتریس انتقال می باشد. هر چند با روش های فشرده سازی می توان هزینه محاسبه را تا حدی کاهش داد [8]. نقطه ضعف استفاده از گرامرهای احتمالی ابرمتن نیز داشتن پیچیدگی $O(n)$ به واسطه استفاده از الگوریتم پیمایش در عمق می باشد. در این مقاله یک روش جدید و مبتنی بر اتوماتای یادگیر توزیع شده که از گرامرهای احتمالی استفاده خواهد کرد، برای تعیین میزان ارتباط بین صفحات وب پیشنهاد می گردد. در این روش به هر صفحه وب یک اتوماتای یادگیر اختصاص داده می شود که وظیفه اش یادگیری ارتباطات آن صفحه با صفحات دیگر می باشد. کارایی الگوریتم پیشنهادی از طریق مقایسه با سه روش Ant web و Bollen و E-DLA مورد ارزیابی قرار خواهد گرفت. ادامه مقاله به این صورت سازماندهی شده است: در بخش ۲ اتوماتای یادگیر و اتوماتای یادگیر توزیع شده به اختصار شرح داده می شوند. در بخش ۳ گرامر احتمالی ابرمتن شرح داده می شوند. در بخش ۴ الگوریتم پیشنهادی و در بخش ۵ نتایج آزمایشات ارائه می - گردد. بخش پایانی نتیجه گیری می باشد.

۲. اتوماتای یادگیر

اتوماتای یادگیر یک مدل انتزاعی است که تعداد محدودی عمل را می تواند انجام دهد. هر عمل انتخاب شده توسط محیطی احتمالی ارزیابی شده و پاسخی به اتوماتای یادگیر داده می شود. اتوماتای یادگیر از این پاسخ استفاده نموده و عمل خود را برای مرحله بعد انتخاب می کند. شکل ۱ ارتباط بین اتوماتای یادگیر و محیط را نشان می دهد.



شکل ۱: ارتباط بین اتوماتای یادگیر و محیط.

محیط: محیط را می توان توسط سه تایی $E \equiv \{\alpha, \beta, c\}$ نشان داد که در آن $\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\}$ مجموعه ورودی ها، $\beta \equiv \{\beta_1, \beta_2, \dots, \beta_m\}$ مجموعه خروجی ها و $c \equiv \{c_1, c_2, \dots, c_r\}$ مجموعه احتمالات جریمه می باشد. هرگاه β مجموعه دو عضوی باشد، محیط از نوع P می باشد. در چنین محیطی $\beta_1 = 1$ به عنوان جریمه و $\beta_2 = 0$ به عنوان پاداش در نظر گرفته می شود. در محیط از نوع Q ، $\beta(n)$ می تواند به طور

گسسته یک مقدار از مقادیر محدود در فاصله [0,1] و در محیط از نوع S ، $\beta(n)$ متغیر تصادفی در فاصله [0,1] است. c_i احتمال اینکه عمل α_i نتیجه نامطلوب داشته باشد، می باشد. در محیط ایستا، مقادیر C_i بدون تغییر باقی می ماند حال آنکه در محیط غیرایستا این مقادیر در طی زمان تغییر می کنند. اتوماتاهای یادگیر به دو گروه با ساختار ثابت و متغیر تقسیم بندی می گردند. در ادامه به شرح مختصری درباره اتوماتای یادگیر با ساختار متغیر و اتوماتای یادگیر توزیع شده که در این مقاله از آنها استفاده شده است، می پردازیم.

۱.۲ اتوماتای یادگیر با ساختار متغیر

اتوماتای یادگیر با ساختار متغیر توسط ۴ تایی $\{\alpha, \beta, p, T\}$ نشان داده می شود که در آن $\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\}$ مجموعه عمل های اتوماتا، $\beta \equiv \{\beta_1, \beta_2, \dots, \beta_m\}$ مجموعه ورودی های اتوماتا، $p \equiv \{p_1, p_2, \dots, p_r\}$ بردار احتمال انتخاب هر یک از عمل ها و $p(n+1) = T[\alpha(n), \beta(n), p(n)]$ الگوریتم یادگیری می باشد. در این نوع از اتوماتاها، اگر عمل α_i در مرحله n انتخاب شود و پاسخ مطلوب از محیط دریافت نماید، احتمال $p_i(n)$ افزایش یافته و سایر احتمالات کاهش می یابند و برای پاسخ نامطلوب احتمال $p_i(n)$ کاهش یافته و سایر احتمالات افزایش می یابند. در هر حال تغییرات به گونه ای صورت می گیرد تا حاصل جمع $p_i(n)$ همواره مساوی یک باقی بماند. الگوریتم زیر نمونه ای از الگوریتم های یادگیری خطی در اتوماتای با ساختار متغیر است.

$$\begin{aligned} p_i(n+1) &= p_i(n) + a[1 - p_i(n)] \\ p_j(n+1) &= (1-a)p_j(n) \quad \forall j \neq i \end{aligned} \quad (1)$$

الف - پاسخ مطلوب

$$\begin{aligned} p_i(n+1) &= (1-b)p_i(n) \\ p_j(n+1) &= \frac{b}{r-1} + (1-b)p_j(n) \quad \forall j \neq i \end{aligned} \quad (2)$$

ب - پاسخ نامطلوب

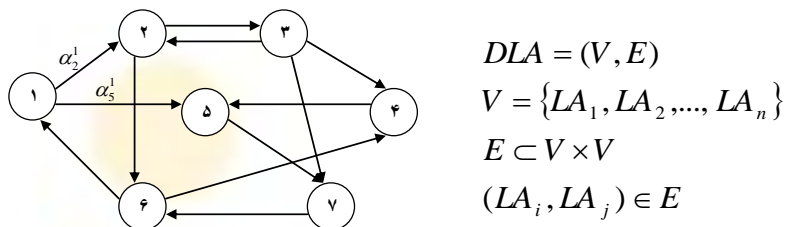
در روابط فوق a ، پارامتر پاداش و b پارامتر جریمه می باشد. با توجه به مقادیر a و b سه حالت را می توان در نظر گرفت. زمانی که a و b باهم برابر باشند، الگوریتم را LRP می نامیم. زمانی که b از a خیلی کوچک باشد، الگوریتم را $LREP$ می نامیم. زمانی که b مساوی صفر باشد، الگوریتم را LRI می نامیم. برای مطالعه بیشتر درباره اتوماتاهای یادگیر می توان به [9,12,13,16,20] مراجعه کرد.

۲.۲ اتوماتای یادگیر توزیع شده $(DLA)^4$

یک اتوماتای یادگیر توزیع شده، شبکه ای از اتوماتاهای یادگیر است که برای حل یک مسئله خاص با یکدیگر همکاری دارند. در این شبکه از اتوماتاهای یادگیر همکار در هر زمان، تنها یک اتوماتا فعال است. مدلی که برای شبکه DLA در نظر گرفته می شود یک گراف جهت دار است که هر یک از رئوس آن یک اتوماتای یادگیر است. وجود یال (LA_i, LA_j) در این گراف به این معناست که انتخاب عمل α_j توسط LA_i باعث فعال شدن LA_j می گردد. تعداد اعمال قابل انتخاب توسط LA_k بصورت $p^k = \{p_1^k, p_2^k, \dots, p_r^k\}$ نمایش داده می شود. در این مجموعه عدد p_m^k نشان

⁴ Distributed Learning Automata

دهنده احتمال مربوط به عمل α_m^k است. انتخاب عمل α_m^k توسط LA_k باعث فعال شدن LA_m می‌شود. تعداد اعمال قابل انجام توسط اتوماتای LA_k را نشان می‌دهد. شکل ۲ مثالی از یک اتوماتای یادگیر توزیع شده با هفت اتوماتای یادگیر را نشان می‌دهد. برای کسب اطلاعات بیشتر درباره اتوماتای یادگیر توزیع شده و کاربردهای آن می‌توان به [4,1,14,15] مراجعه کرد.



شکل ۲: یک اتوماتای یادگیر توزیع شده با هفت اتوماتای یادگیر

3. گرامر احتمالی ابرمتن^۵ (HPG)

یک جلسه کاری کاربر، دنباله‌ای از صفحات ملاقات شده توسط همان کاربر است، به طوری که فاصله زمانی بین ملاقات دو صفحه متوالی بیش از یک حد زمانی معین نباشد (پژوهشگران نشان داده اند که این زمان تقریباً ۲۵,۵ دقیقه یا ۳۰ دقیقه می‌باشد). می‌توان جلسات کاری کاربران استنتاج شده از لاگ داده را توسط زبانی به نام HPA[10] بیان کرد که توسط HPG[5] بدست می‌آید. HPG یک گرامر با قاعده احتمالی است به طوری که بین مجموعه غیرپایانه‌ها و پایانه‌ها یک نگاشت یک به یک وجود دارد. هر غیرپایانه نشان دهنده یک صفحه وب و هر قانون تولید نشان دهنده یال بین صفحات است. در این مدل، سیستم ابرمتن به صورت گراف جهت دار $G = (N, E)$ در نظر گرفته می‌شود. به طوری که N تعداد نودها و E یال‌های این گراف را نشان می‌دهد. اگر A, B دو صفحه وب باشند و کاربر از صفحه A به صفحه B حرکت کند، یک یال یک جفت $(A, B) \in E$ خواهد بود. تعداد زمان‌هایی که یک دنباله از دو صفحه در جلسه کاربری ظاهر شده (مانند دنباله $A \rightarrow B$) معادل تعداد دفعاتی است که لینک مورد نظر (یال) ملاقات شده است. این مدل بصورت $\Phi(H, M, S, F, \Gamma)$ تعریف می‌شود که در آن:

H : مجموعه نودها، $H = \{h_1, \dots, h_m\}$ که هر یک از اعضای آن معرف یک صفحه وب خواهد بود.

M : یک ماتریس $m \times m$ است به طوری که $M = \{p(h_i, h_j)\}$ ، $0 \leq p(h_i, h_j) \leq 1$ ، $(h_i, h_j) \in H$ ، $\sum_{k=1}^m p(h_i, h_k) = 1$.

S : حالت شروع جلسه کاری کاربر و F : حالت خاتمه جلسه کاری کاربر

Γ : تابعی از $H \times H$ به $[0,1]$ است. یعنی احتمال حرکت از صفحه h_i به صفحه h_j را نشان می‌دهد.

اگر صفحه h_i را به عنوان صفحه شروع و صفحه h_j را به عنوان صفحه مقصد در نظر بگیریم، یک قانون انجمنی $h_i \rightarrow h_j$ به این صورت تعریف می‌شود:

⁵ Hypertext Probabilistic Grammar

زمانی که کاربر به صفحه h_i می‌رسد لینک بعدی که او انتخاب می‌کند با یک احتمال معین h_j خواهد بود. احتمال لینک انتخاب شده به عنوان اطمینان⁶ توصیف شده و به این صورت تعریف می‌شود:

$$P(h_i, h_j) = \frac{R_j^i}{R^i} \quad (3)$$

R_j^i : تعداد دفعاتی که صفحه h_j بلافاصله بعد از صفحه h_i قرار گرفته است و R^i : تعداد تمامی صفحاتی که بعد از صفحه h_i ملاقات می‌شوند. احتمال رشته تولید شده توسط این گرامر، از حاصل ضرب اطمینان همه قوانین مورد استفاده برای اشتقاق این رشته به دست می‌آید. به عبارت دیگر اگر رشته تولید شده توسط این گرامر به صورت $\langle h_{i1}, \dots, h_{ir} \rangle$ باشد، احتمال رشته تولید شده به صورت رابطه زیر تعیین می‌شود:

$$P(\langle h_{i1}, \dots, h_{ir} \rangle) = \prod_{k=1}^{r-1} P(h_{ik}, h_{i(k+1)})$$

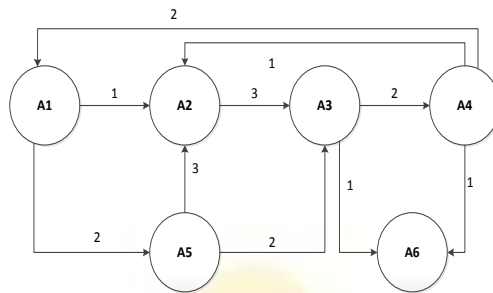
در این مدل، دنباله‌های حرکتی کاربران استخراج شده از لاگ داده به عنوان گرامر احتمالی ابرمتن مدل سازی می‌شوند، به طوری که رشته‌های تولید شده با احتمال بالا مسیرهای حرکتی را که مورد علاقه کاربران بوده را نشان می‌دهد [5]. برای مدل سازی جلسات کاربران در سیستم ابرمتن از مدل پیشنهادی [5] استفاده می‌کنیم.

مثال ۱: اگر حرکت‌های کاربر در یک جلسه کاری کاربر به صورت جدول ۱ باشد، در این مثال ما ۶ جلسه کاربری داریم. سیستم ابرمتن آن بصورت شکل ۳ خواهد بود. به عنوان مثال در شکل ۳ وزن یا $A_2 \rightarrow A_3$ برابر ۳ می‌باشد زیرا این دنباله حرکتی در جلسات کاربران ۳ بار ظاهر شده است. همچنین با توجه به شکل ۳ احتمال رشته $(A_1 \rightarrow A_2 \rightarrow A_3)$ به صورت زیر محاسبه می‌شود:

جدول ۱: دنباله‌ای از حرکت‌های کاربران

Session ID	User trail
1	$A_1 \rightarrow A_2 \rightarrow A_3 \rightarrow A_4$
2	$A_1 \rightarrow A_5 \rightarrow A_3 \rightarrow A_4 \rightarrow A_1$
3	$A_5 \rightarrow A_2 \rightarrow A_4 \rightarrow A_6$
4	$A_5 \rightarrow A_2 \rightarrow A_3$
5	$A_5 \rightarrow A_2 \rightarrow A_3 \rightarrow A_6$
6	$A_4 \rightarrow A_1 \rightarrow A_5 \rightarrow A_3$

⁶ Confidence



شکل ۳: ایجاد سیستم ابرمتن از حرکت‌های کاربر

ابتدا از حرکت‌های کاربر یک سیستم ابر متن می‌سازیم (شکل ۳). بعد از ایجاد سیستم ابرمتن می‌توان به تولید رشته‌ها از آن پرداخت. الگوریتمی وجود دارد که استخراج این رشته‌ها را انجام می‌دهد. الگوریتم استخراج رشته‌ها از سیستم ابرمتن را استخراج قوانین انجمنی می‌گویند [5]. این الگوریتم در ادامه بحث در بخش ۱.۳ آورده شده است.

۱.۳ الگوریتم استخراج قوانین انجمنی از سیستم ابرمتن

هدف از استخراج قوانین انجمنی [5] این است، که در گراف جهت دار موجود تمامی قوانینی را تولید کنیم که مقدار احتمال قوانین استخراجی از مقدار نقطه برش^۷ بیشتر باشد. این الگوریتم از هر نود این گراف شروع کرده و رشته‌هایی که مقدار احتمال آنها از نقطه برش بیشتر باشند را تولید می‌کند. اگر CR^A نشان دهنده قانون انتخاب باشد و x یک دنباله باشد. $forward-neigh(x)$ همسایه رو به جلو و ملاقات نشده را نشان می‌دهد. با یک دنباله x و یک لینک e ، $x+e$ الحاق دو رشته x و e را نشان می‌دهد. این الگوریتم حالت خاصی از الگوریتم DFS^8 می‌باشد. از هر یال این گراف درخت DFS ساخته شده و هر شاخه این درخت متناظر با یک قانون انتخاب (CR) است. هر درخت DFS مشخص کننده همه قوانین انتخابی است، که یال شروع شده از ریشه این درخت، اولین یال آن را نشان می‌دهد. کاوش یک شاخه زمانی پایان می‌پذیرد که به نودی برسیم که همه یال‌های آن از قبل ملاقات شده باشد یا مقدار اطمینان دنباله پایین تر از حد آستانه C باشد. بعلاوه کاوش هر شاخه مستقل از هر شاخه دیگر است. الگوریتم استخراج قوانین انجمنی در شکل ۴ آورده شده است.

⁷ Cut-Point

⁸ Candidate Rules

⁹ Depth-First-Search

شکل ۴: الگوریتم استخراج قوانین انجمنی از سیستم ابر متن

Algorithm 1.1(Modified_DFS (G, C))

```

1. begin
2. for each  $\{e \in E: C_e \geq C\}$ 
3. Explore ( $e, C_e$ )
4. end for
5. end.

```

Algorithm 1.2(Explore ($trail, C_{trail}$))

```

1. begin
    $CR = CR \cup trail$ 
3. for each ( $forw - neigh(trail)$ )
4. if  $(C_{fn} \times C_{trail} \geq C)$  then
5. Explore( $trail + f_n, C_{trail}$ )
6. end for
7. end

```

مثال ۲: در مثال ۱ با فرض این که مقدار نقطه برش $\lambda = 0.33$ باشد در این صورت نحوه استخراج قوانین انجمنی با شروع از A_1 به صورت زیر انجام می شود. از A_1 می توان به A_2 و A_5 رفت ولی $P(A_1 \rightarrow A_2) = \frac{1}{3} < \lambda$ است، پس مسیر $A_1 \rightarrow A_2$ را نادیده می گیریم. حال از صفحه A_5 کار را ادامه می دهیم. از این صفحه می توان به صفحات A_2 و A_3 رفت. مسیر $A_1 \rightarrow A_5 \rightarrow A_3$ نیز رد می شود زیرا خواهیم داشت $P(A_1 \rightarrow A_5 \rightarrow A_3) = \frac{2}{3} \times \frac{2}{5} < \lambda$ پس مسیر استخراجی با شروع از A_1 عبارت است از: $A_1 \rightarrow A_5 \rightarrow A_2$ جدول ۲ تمامی قوانین استخراجی را برای شکل ۳ به دست می دهد و برای هر کدام از نودها (صفحات وب) نیز داریم:

$$A_1: \begin{cases} p(A_1 \rightarrow \dots) \\ p(A_1 \rightarrow \dots) \end{cases}$$

جدول ۲: قوانین انجمنی استخراج شده از شکل ۳

Session ID	User Session
1	$A_4 \rightarrow A_1$
2	$A_5 \rightarrow A_2 \rightarrow A_3$
3	$A_5 \rightarrow A_3$
4	$A_3 \rightarrow A_4$
5	$A_2 \rightarrow A_3 \rightarrow A_4$
6	$A_1 \rightarrow A_5 \rightarrow A_2$

ل به نقطه برش معروف است. هر چه مقدار λ بزرگتر باشد، طول قوانین استخراج شده نیز کمتر خواهد بود و بالعکس هرچه این مقدار کمتر باشد، طول قوانین استخراجی نیز بیشتر خواهد بود. برای مقدار کمتر از ۰,۳، تعداد تکرارهای الگوریتم بالا است ولی برای مقدار بزرگتر از ۰,۳ این تکرارها کم می‌شود. پیچیدگی این الگوریتم خطی و به تعداد نودها وابسته است. به عبارت دیگر با افزایش تعداد نودها، تعداد تکرارهای الگوریتم نیز به صورت خطی رشد می‌کند. با افزایش تعداد نودها، زمان اجرای CPU به صورت توانی رشد می‌کند. همچنین با افزایش درجه گراف، تعداد قوانین تولید شده نیز بیشتر خواهد بود [6].

۴. الگوریتم پیشنهادی

صفحات وب و کاربران استفاده کننده از آن نقش یک محیط تصادفی را برای اتوماتای یادگیر موجود در DLA ایفا می‌کنند. خروجی DLA یک دنباله از صفحات وب مرور شده توسط کاربر هستند که مسیر حرکت کاربر را به سمت یک صفحه وب مورد نظر نشان می‌دهد. محیط با استفاده از این دنباله پاسخی برای DLA تولید می‌کند. با استفاده از این پاسخ ساختار داخلی اتوماتای یادگیر در اتوماتای یادگیر توزیع شده طبق الگوریتم یادگیری بروز می‌شود. در الگوریتم پیشنهادی از اتوماتای یادگیر با ساختار متغیر استفاده شده است. برای این کار در درروش پیشنهادی برای هر صفحه p_i یک اتوماتای یادگیر LA_i در نظر می‌گیریم. اگر n تعداد صفحات وب باشد در این حالت اتوماتا $n-1$ عمل دارد. زمانی که کاربر از صفحه p_i به صفحه p_j حرکت می‌کند j امین عمل از اتوماتای LA_i در DLA فعال شده و به محیط اعمال می‌شود. در صورتی که عمل انتخاب شده k امین عمل اتوماتای LA_i باشد (یعنی $j = \alpha_k^i$) احتمال متناظر این عمل یعنی p_k^i به عنوان میزان ارتباط صفحات i و j در نظر گرفته می‌شود. در این الگوریتم ابتدا یک لاگ فایل برای ذخیره حرکات کاربران در نظر می‌گیریم. سپس حرکات کاربران در این فایل ذخیره می‌شود، بعد از پر شدن این لاگ فایل، یک سیستم ابر متن از این فایل ساخته می‌شود و قوانین انجمنی با مقدار نقطه برش با استفاده از الگوریتم ۱.۲ و ۱.۱ (شکل ۴) از آن استخراج می‌شوند. این قوانین استخراجی، مسیریابی را نشان می‌دهند که شدت دنباله (مقدار احتمال) آنها بالاتر است، در نهایت فقط به این قوانین استخراجی توسط اتوماتای یادگیر پاداش داده می‌شود. دوباره لاگ فایل پاک شده حرکت های جدید در آن وارد شده و دوباره همین روال انجام می‌گیرد. الگوریتم پیشنهادی بصورت زیر است:

ورودی : دنباله های حرکتی کاربران

خروجی : استخراج قوانین انجمنی

۱- یک DLA متناظر با ساختار اسناد ایجاد کن.

۲- بردار احتمالات اتوماتاهای یادگیر در DLA را مقدار دهی اولیه کن.

۳- یک لاگ فایل جهت ذخیره حرکت های کاربران ایجاد کن.

۴- تا زمانی که لاگ فایل پر نشده، حرکت های کاربران را در آن ذخیره کن.

۵- اگر لاگ فایل پر شد، از اطلاعات آن یک سیستم ابرمتن بساز.

۶- با استفاده از الگوریتم ۱، ۱.2 و (شکل ۴) تمامی قوانین انجمنی با نقطه برش معلوم را از این سیستم ابرمتن استخراج کن .

۷- برای هر قانون استخراجی انجام بده

۷-۱- برای هر حرکت $D_m \rightarrow D_k$ (حرکت از صفحه D_k به صفحه D_m) کاربر در طول قانون استخراجی، انجام بده

۷-۲- بردار احتمال اعمال اتوماتای یادگیر LA_k را طبق الگوریتم یادگیری زیر بروز کن

$$\begin{aligned} p_m(n+1) &= p_m(n) + a[1 - p_m(n)] \\ p_j(n+1) &= (1-a)p_j(n) \quad j \neq m \quad \forall j \end{aligned}$$

۸- لاگ فایل را پاک کرده و برو به ۴

۹- پایان الگوریتم.

شکل ۵: الگوریتم پیشنهادی

۵. ارزیابی کارایی الگوریتم پیشنهادی و نتایج آزمایشات

برای انجام شبیه سازی از مدل ارایه شده در [11] برای تولید مجموعه صفحات وب و حرکات کاربران استفاده می شود. هر صفحه با بردار محتوای $C_n = [cw_n^1 \quad cw_n^2 \quad \dots \quad cw_n^M]$ نشان داده می شود که در آن M تعداد موضوعات در سیستم می باشد. هر عضو این بردار (cw_n^i) میزان ارتباط صفحه متناظر با آن بردار را با یکی از این موضوعات نشان می دهد. با استفاده از بردار محتوای هر کدام از صفحات، شباهت بین هر دو صفحه موجود در سیستم محاسبه می شود. ماتریس شباهت به دست آمده به عنوان ماتریس ارتباطات ایده ال بین صفحات وب در شبیه سازی ها به منظور ارزیابی کارایی روش پیشنهادی استفاده می شود. در این مدل توزیع موضوعات بین صفحات وب بصورت توزیع نرمال و پروفایل علاقه کاربران بصورت توزیع قانون توانی در نظر گرفته شده است که با تغییر پارامتر این توزیع و تعداد صفحات، سیستم های اطلاعاتی متفاوتی می توان ایجاد نمود.

جدول ۳: پارامترهای شبیه سازی

حد آستانه ایجاد اتصال	۰.۷
تعداد کاربران	۱۰۰۰۰ و ۲۰۰۰۰
تعداد صفحات وب	۳۰
تعداد موضوع ها	۵
T_c : مقدار ثابت صفحه اولیه سایت در موضوعات مختلف	۰.۲
ΔM_t^c : ضریب ثابت کاهش اشتیاق کاربر	-
ΔM_t^p : ضریب متغیر کاهش اشتیاق کاربر	-
α_u : پارامتر توزیع قانون- توانی (توزیع احتمال علائق کاربران)	۱
\emptyset : ضریب پاداش دریافتی از مشاهده یک صفحه	1.2
λ : ضریب جذب اطلاعات از یک صفحه توسط یک کاربر	۰.۵
μ_m : میانگین توزیع نرمال ΔM_t^p	5.79
σ_m : واریانس توزیع نرمال ΔM_t^p	0.25
μ_t : میانگین توزیع نرمال برای مقدار افزایش یک گره برای یک موضوع خاص	-
α_p : پارامتر توزیع قانون- توانی (توزیع احتمال های وزن های مطالب برای هر صفحه)	۳
σ : واریانس توزیع نرمال برای مقدار افزایش یک گره برای یک موضوع	0.25
θ : ضریب کاهش علاقه کاربر	۱
حداقل اشتیاق کاربر برای ادامه جستجو	0.2

در این مدل، انگیزه و استراتژی حرکتی کاربران نیز از طریق توزیع های آماری مدل شده اند. با تغییر پارامترهای این توزیع های آماری می توان کاربرانی با علائق، انگیزه ها و استراتژیهای متفاوت ایجاد نمود. در شبیه سازی ها، پارامترهای تعداد موضوع ها، تعداد کاربران، علائق، انگیزه ها و پارامترهای توزیع آنها را در طول شبیه سازی ثابت در نظر گرفته ایم. برای ارزیابی الگوریتم پیشنهادی از معیاری به نام کورولیشن استفاده خواهیم نمود. کورولیشن دو مجموعه داده مانند X, Y بصورت رابطه ۶ محاسبه می شود که در آن N تعداد داده ها است. مجموعه X ساختار ایجاد شده توسط مدل شبیه سازی [11] و Y نیز ساختار بدست آمده توسط الگوریتم پیشنهادی را نشان می دهد.

$$Corr(X, Y) = \frac{Cov(X, Y)}{\sigma_X \sigma_Y} = \frac{\sum XY - (\sum X \sum Y) / N}{\sqrt{(\sum X^2 - (\sum X)^2 / N)(\sum Y^2 - (\sum Y)^2 / N)}} \quad (5)$$

$$X = \{P_{ij} | i, j = 1, 2, \dots, n, \quad i \neq j\}$$

$$Y = \{P'_{ij} | i, j = 1, 2, \dots, n, \quad i \neq j\},$$

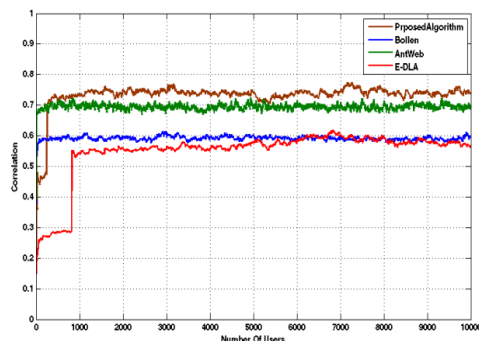
$$P_{ij} = \frac{(d_{ij})^{-1}}{\sum_{k=1}^n (d_{ik})^{-1}}$$

$$P'_{ij} = \text{probability of action } j \text{ from DLA } (i) \quad (6)$$

$$d_{ij} = \sqrt{\sum_{k=1}^M (cw_i^k - cw_j^k)^2}$$

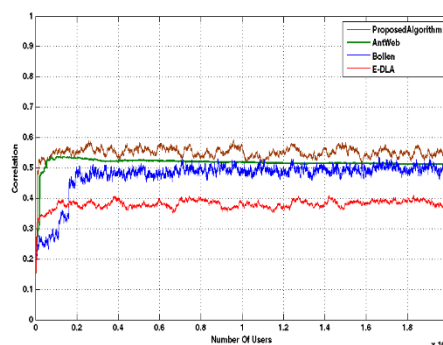
در رابطه فوق d_{ij} فاصله اقلیدسی بین دو صفحه i, j از یکدیگر در مدل شبیه‌سازی می‌باشد. پارامترهای شبیه‌سازی الگوریتم پیشنهادی مطابق با مقادیر جدول ۳ می‌باشد.

آزمایش ۱: شکل ۶ نتیجه مقایسه اجرای الگوریتم پیشنهادی (Proposed Algorithm) را با الگوریتم‌های AntWeb و Bollen و E-DLA نشان می‌دهد. در این شبیه‌سازی تعداد کاربران استفاده کننده از وب ۱۰۰۰۰، تعداد موضوعات مرتبط با هر سند ۵، تعداد صفحات وب ۳۰ و مقدار نقطه برش نیز $\lambda = 0.1$ در نظر گرفته شده است.



شکل ۶: مقایسه نتایج الگوریتم‌های موجود

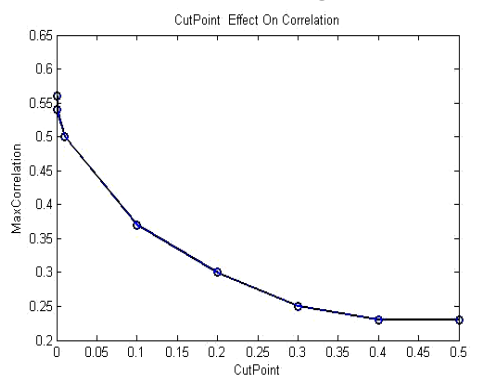
آزمایش ۲: شکل ۷ مقایسه نتیجه اجرای الگوریتم پیشنهادی را با الگوریتم‌های AntWeb و Bollen و E-DLA نشان می‌دهد. در این شبیه‌سازی تعداد کاربران استفاده کننده از وب ۲۰۰۰۰، تعداد موضوعات مرتبط با هر سند ۵، تعداد صفحات وب ۳۰ و مقدار نقطه برش نیز $\lambda = 0.1$ در نظر گرفته شده است.



شکل ۷: مقایسه نتایج الگوریتم‌های موجود

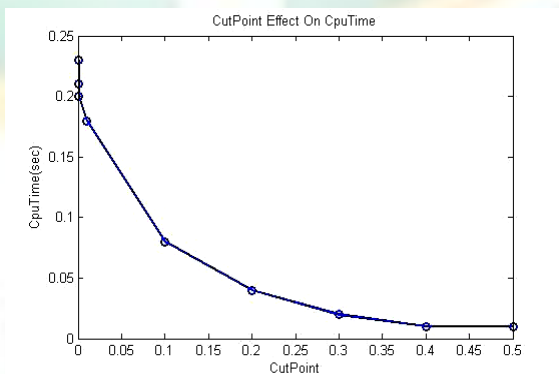
همان‌طوری که از شکل‌های ۶ و ۷ مشخص است، الگوریتم پیشنهادی دارای کارایی بالایی می‌باشد و می‌توان از این الگوریتم برای کشف ساختار ارتباطی بین صفحات وب بهره جست. با استفاده از مقدار نقطه برش λ ، در هر بار قوانین انجمنی استخراج شده و به این قوانین انجمنی پاداش داده شده است. نتایج آزمایشات ۱ و ۲ نشان می‌دهد که الگوریتم پیشنهادی در تعداد تکرارهای کمتری به کورولیشن مورد نظر می‌رسد.

آزمایش ۳ : در این آزمایش، تاثیر پارامتر نقطه برش بر روی کورولیشن نشان داده شده است. λ به نقطه برش معروف است. همان طوری که از شکل ۸ مشخص است، مقدار نقطه برش تاثیر زیادی در مقدار کورولیشن دارد. به عبارت دیگر با افزایش مقدار نقطه برش، مقدار کورولیشن نیز کمتر می شود. علت این امر را می توان در طول قوانین استخراجی دانست. هر چه مقدار λ بزرگتر باشد، طول قوانین استخراج شده نیز کمتر خواهد بود و بالعکس هر چه این مقدار کمتر باشد، طول قوانین استخراجی نیز بیشتر خواهد بود. نتایج این شبیه سازی در شکل ۸ آمده است.



شکل ۸: تاثیر نقطه برش بر روی کورولیشن

آزمایش ۴ : تاثیر نقطه برش بر روی زمان و کارایی الگوریتم پیشنهادی نیز انجام شد و تاثیر این پارامتر بر روی الگوریتم مشخص شد. نتایج این شبیه سازی در شکل ۹ آمده است. شکل ۹ تاثیر مقدار نقطه برش را در زمان اجرای الگوریتم نشان می دهد. همان طوری که از شکل ۹ مشخص است، مقدار نقطه برش تاثیر زیادی در زمان اجرا دارد. به عبارت دیگر با افزایش مقدار نقطه برش، مقدار زمان اجرا نیز کمتر می شود. علت این امر را می توان در طول قوانین استخراجی دانست. هر چه مقدار λ بزرگتر باشد، طول قوانین استخراج شده نیز کمتر خواهد بود و بالعکس هر چه این مقدار کمتر باشد، طول قوانین استخراجی نیز بیشتر خواهد بود.



شکل ۹: تاثیر نقطه برش بر روی زمان اجرا (استخراج قوانین انجمنی و پاداش به اتوماتا)

۶. نتیجه گیری

در این مقاله یک روش جدید برای تعیین میزان ارتباط بین صفحات وب با استفاده از اتوماتای یادگیر توزیع شده که از گرامرهای احتمالی ابر-متن استفاده می نماید، پیشنهاد گردید. کارایی الگوریتم پیشنهادی با روش های Ant Web و Bollen و E-DLA مورد مقایسه قرار گرفت و نشان داده شد که کورولیشن الگوریتم پیشنهادی بهتر از این روش ها است. مزیت الگوریتم پیشنهادی در این است که الگوریتم پیشنهادی به دلیل در نظر گرفتن قوانین انجمنی به حرکت هایی با کیفیت بالا پاداش می دهد. بنابراین الگوهایی با کیفیت پایین در نظر گرفته نشده و این امر باعث افزایش کارایی الگوریتم پیشنهادی نسبت به سایر الگوریتم ها می شود. الگوریتم پیشنهادی به دلیل استفاده از الگوریتم پیمایش در عمق (DFS)

دارای پیچیدگی زمانی $O(n)$ بود که در آن n تعداد نودها می‌باشد. بنابر این اگر تعداد نودهایی که در لاگ فایل قرار می‌گیرند، n باشد و K تعداد دفعات تکرار باشد، در آن صورت هزینه الگوریتم پیشنهادی برابر با $O(nK)$ خواهد بود. همچنین تاثیر پارامتر نقطه برش را در زمان اجراء و کورولیشن بحث کردیم و دیدیم که الگوریتم پیشنهادی ما تحت تاثیر این پارامتر قرار دارد به عبارت دیگر تعیین مقدار اولیه نقطه برش حائز اهمیت است. از کاربردهای الگوریتم پیشنهادی می‌توان به استخراج قوانین انجمنی، پیش بینی صفحات وب مورد علاقه کاربران، فراهم کردن وب سایتهای تطبیقی برای کاربران، خوشه بندی کاربران و رتبه بندی صفحات وب اشاره کرد.

۷. تقدیم و تشکر

از استاد بزرگوار جناب آقای دکتر محمد رضا میبیدی که با زحمات بی‌شائبه خویش ما را در کسب علم و معرفت یاری نمودند تشکر و قدردانی می‌نماییم.

این مقاله مستخرج از طرحی است که با حمایت مالی دانشگاه پیام نور با عنوان کاوش ساختار وب با استفاده از اتوماتای یادگیر و گرامر احتمالی به تصویب رسیده است.

کنفرانس داده‌کاوی ایران

- [1] M. Alipour; M. R. Meybodi, "Solving probabilistic traveling sales man problem using distributed learning automata", Proc. of 11th Annual CSI Computer Conference of Iran, Fundamental Science Research Center(IPM), Computer Science Research Lab, Tehran, Iran, pp.673-67, Jan. 24-26, 2006.
- [۲] B. Anari; M. R. Meybodi, "A new method based on distributed learning automata for determining web documents structure", Proc. of 12th Ann. Int. CSI Computer Conf. CSICC 2007, Tehran, Iran, pp. 2276- 2281, Feb. 20-22, 2007.
- [۳] A. Baradaran Hashemi ; M.R., Meybodi, "Web usage mining using distributed learning automata", Proc. Int. Conf. on Computer Engineering Department, Technical Report, 2005.
- [۴] H. Beigy; M. R. Meybodi, "Utilizing distributed learning automata to solve stochastic shortest Path Problem", Int. Jour. of Uncertainty, Fuzziness and Knowledge-based Systems, World Scientific Publishing Company, Vol. 14, No. 5, pp. 591-617, October 2006 .
- [۵] J. Borges ; M. Leven, "Data Mining of user navigation patterns", Proc. of the Web Usage Analysis and User Profiling, vol. 1, PP. 31-36, 1999.
- [۶] J. Borges; M. Levene, "Mining association rules in hypertext databases", Proc. of the 4th Int. Conf. On Knowledge Discovery and Data Mining, pp. 149-153, August 1998.
- [۷] F. Heylighen; J. Bollen, "Hebbian Algorithm for a Digital Library Recommendation System", Proc. Int. Conf. on Parallel Processing Workshops (ICPPW'02) IEEE, 2002.
- [۸] Z. Jianhan, "Mining web site link structures for adaptive web site navigation and search", Ph.D. Thesis, university of Ulster at jordanstown, October 2003.
- [۹] S. Lakshmivarahan, "Learning algorithms: theory and applications", New York: Springer- Verlag, 1981.
- [۱۰] M. Levene; G. Loizou, "A Probabilistic approach to navigation in hypertext", Information Sciences, pp. 1۱۴-1۲۵, 1999.
- [1۱] J. Liu ; S. Zhang ; J. Yang, "Characterizing web usage regularities with information foraging agents", IEEE Trans. in knowledge and data engineering, Vol. 16, No. 5, May 2004.
- [۱۲] P. Mars; J.R. Chen; R. Nambir, "Learning Algorithms: Theory and Applications in Signal Processing", Control, and Communication, CRC Press Inc, 1996.
- [۱۳] M. R. Meybodi; S. Lakshmivarahan, "On a class of Learning Algorithms which have Symmetric Behavior under Success and Failure", Lecture Notes in Statistics, Berlin: SpringerVerlag, pp. 145-155, 1984.
- [1۴] M. R. Meybodi; H. Beigy, "Solving stochastic shortest path problem using Monte Carlo sampling method: A distributed learning automata approach", Springer-verlag lecture notes in advances in Soft Computing: Neural Networks and Soft Computing, pp. 626-632, 2003. (ISBN: 3-7908-0005-8).
- [۱۵] M. R. Meybodi; H. Beigy, "Solving stochastic path problem using distributed learning automata", Proc. of 6th Annual International CSI Computer Conference (CSICC2001), Isfahan, Iran, pp. 70-86, Feb. 20- 22, 2001.
- [1۶] K. S. Narendra; M. A. L Thathachar, "Learning automata: An introduction," Prentice Hall, 1989.
- [1۷] S. Saati; M.R. Meybodi, "Document ranking using distributed learning automata," Proc. of 11th Annual CSI Computer Conf. of Iran, Fundamental Science Research Center(IPM), Computer Science Research Lab, Tehran, Iran, pp.467- 473, Iran, May 24-26 2006.
- [1۸] S. Saati; M.R Meybodi, "A self-organizing model for document structure using distributed learning automata", Proc. The Second Int. IEEE/WIC Conf. on Information and knowege Technology (IKT2005), Tehran, Iran, May 24-26, 2005.
- [1۹] W. Teles; L. Weigang; C. Ralha, "AntWeb-The adaptive web server based on the ants behavior", Proc. IEEE/WIC Int. Conf. on Web Intelligence (WI03), PP.558-564, 2003.
- [۲۰] M. A. L Thathachar Baskar; R. Harita, "Learning automata with changing number of actions", IEEE Trans. on System, Man and Cybernetic, vol. SMC-17, No. 6, Nov. 1987.