

پیمایش موضوعی وب با استفاده از اتوماتای یادگیر توزیع شده و پارتیشن بندی گراف

مجید تاران¹، شهرزاد معتمدی مهر²، علی برادران هاشمی³ و محمدرضا میبیدی⁴

¹ شرکت خدمات انفورماتیک، تهران، ایران، m_taran@isc.ir

² دانشکده فنی و مهندسی، دانشگاه خوارزمی، تهران، ایران، motamedi@tmu.ac.ir

³ دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیرکبیر، تهران، ایران، a_hashemi@aut.ac.ir

⁴ دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیرکبیر، تهران، ایران، mmeybodi@aut.ac.ir

چکیده - پیمایش وب جهت یافتن صفحاتی که توسط موتور جستجو شاخص گذاری شوند، از اهمیت بالایی برخوردار است. از آن جا که حجم صفحات وب بسیار بالا و همواره رو به افزایش است، موتورهای جستجو قادر به شاخص گذاری صفحات محدودی هستند. پیمایشگرهای موضوعی یا پیمایشگرهای متمرکز، در روند پیمایش خود به صورت انتخابگر عمل می کنند و صفحاتی را برای پیمایش انتخاب می کنند که تا حد ممکن در ارتباط با موضوعی خاص باشند. در این مقاله روشی ترکیبی مبتنی بر اتوماتای یادگیر توزیع شده و پارتیشن بندی گراف برای پیمایش موضوعی صفحات وب پیشنهاد می گردد. روش پیشنهادی با استفاده از الگوریتم *HITS* و ساختار پیوند بین صفحات که از طریق رفتار کاربر در مشاهده این صفحات بدست آمده است، صفحاتی را برای پیمایش انتخاب می کند. در این مقاله از پارتیشن بندی گراف وب برای بهبود کارایی استفاده شده است. به منظور ارزیابی، نتایج روش پیشنهادی با نتایج الگوریتم عرض اول، *Best First* و الگوریتمی دیگری مبتنی بر اتوماتای یادگیر توزیع شده مقایسه شده است. نتایج آزمایشها حاکی از کارایی روش پیشنهادی دارد.

کلید واژه- پیمایش موضوعی وب، اتوماتای یادگیر توزیع شده، پارتیشن بندی گراف، داده های استفاده از وب، الگوریتم *HITS*

مقدمه

پیمایشگرهای موضوعی، باید بتوانند مستقل از مجموعه اولیه صفحات (که کار خود را با آن آغاز می کنند) به مجموعه ای از صفحات مرتبط با موضوعی خاص دست یابند. پیمایشگرها با استفاده از پیوندهای موجود در صفحاتی که قبلاً بازیابی کرده اند، صفحات جدید را بازیابی می کنند. وقتی که صفحه ای بازیابی می شود، پیوندهای موجود در آن، به لیست صفحات بازیابی نشده اضافه می شود. این لیست، صف پیمایش (*Crawl Frontier*) نامیده می شود. چالش اصلی در پیمایش موضوعی آن است که مناسب ترین پیوند برای پیمایش از صف فوق الذکر انتخاب شود. یک پیمایشگر موضوعی باید از اطلاعات به دست آمده از پیمایش های قبلی برای تخمین میزان ارتباط یک پیوند جدید استفاده کند.

الگوریتمی که در انتخاب پیوند بعدی برای پیمایش به کار می رود، به هدف پیمایشگر بستگی دارد.

الگوریتم های پیمایش موضوعی را می توان به دو دسته کلی پیمایش بدون استفاده از دانش پیش زمینه (*Background Knowledge*) و پیمایش به کمک دانش پیش زمینه تقسیم کرد. در گروه اول، انتخاب پیوند برای پیمایش تنها بر اساس اطلاعات

با گسترش روز افزون شبکه جهانی اینترنت، بازیابی اسناد مرتبط با نیاز اطلاعاتی کاربران که در قالب یک پرس و جو مطرح می شود، روز به روز دشوارتر می گردد. موتورهای جستجو برای بازیابی اسناد مرتبط با یک پرس و جو از یک الگوریتم رتبه بندی استفاده می کنند که میزان ارتباط اسناد شاخص گذاری شده توسط موتور جستجو را با پرس و جوی کاربر محاسبه می کند. اما الگوریتم رتبه بندی هر اندازه که دقیق باشد، اگر صفحه ای توسط موتور جستجو شاخص گذاری نشده باشد، بازیابی نخواهد شد. از آن جا که حجم صفحات وب بسیار بالا و همواره رو به افزایش است، موتورهای جستجو قادر به شاخص گذاری صفحات محدودی هستند. به همین جهت از الگوریتم های پیمایش موضوعی جهت پیمایش صفحات مرتبط استفاده می شود. پیمایشگرهای موضوعی یا پیمایشگرهای متمرکز، در روند پیمایش خود به صورت انتخابگر عمل می کنند و صفحاتی را برای پیمایش انتخاب می کنند که تا حد ممکن در ارتباط با موضوعی خاص باشند.

موجود در صفحه وب انجام می‌شود. اما در گروه دوم برای تصمیم‌گیری از یک منبع خارجی مانند یک طبقه‌بندی از اسناد، هستان شناسی (Ontology)، ... استفاده می‌گردد. البته الگوریتم‌های دیگری نیز پیشنهاد شده‌اند که در هیچ یک از این دو گروه قرار نمی‌گیرند.

از جمله روشهایی که در دسته اول قرار دارند می‌توان به الگوریتم عرض اول [3]، Best First [4]، Shark Search [5] و PageRank [6] اشاره کرد.

پیمایش عرض اول (Breath First)، ساده‌ترین روش برای پیمایش می‌باشد. این الگوریتم در سال 1994 توسط Pinkerton ارائه شد. در این نوع پیمایش، صف پیمایش، یک صف FIFO است که با استفاده از آن پیوندها به همان ترتیبی که در صف قرار گرفته‌اند، پیمایش می‌شوند. زمانی که صف پر می‌شود، پیمایشگر می‌تواند تنها یکی از پیوندهای صفحه بازبایی شده را به صف اضافه کند. پیمایشگر عرض اول، به عنوان پیمایشگر پایه به کار می‌رود و از هیچ دانشی راجع به موضوع پیمایش استفاده نمی‌کند. الگوریتم BestFirst توسط cho و همکارانش در سال 1998 ارائه شد. ایده اصلی این الگوریتم آن است که از صف پیمایش، پیوندی که بهتر از سایرین معیاری را برآورده می‌کند، برای پیمایش انتخاب شود. در یک حالت ساده، معیار انتخاب، میزان شباهت صفحه حاوی پیوند (صفحه مبدا) با کلمات کلیدی موضوع مورد نظر می‌باشد. بنابراین از میزان شباهت صفحه مبدا با کلمات کلیدی موضوع برای تخمین مرتبط بودن صفحاتی که این پیوندها به آنها اشاره می‌کنند، استفاده می‌شود. الگوریتم SharkSearch در سال 1998 توسط Hersovici ارائه شد که نوعی از الگوریتم FishSearch [7] محسوب می‌شود. در FishSearch، پیمایشگر، در نقاطی از وب جستجوی خود را ادامه می‌دهد که صفحات مرتبط مشاهده شده‌اند و جستجو را در نقاطی که صفحات مرتبط یافته نشده‌اند، متوقف می‌کند. الگوریتم SharkSearch دو مزیت اصلی بر الگوریتم FishSearch دارد. تابعی که برای اندازه‌گیری میزان ارتباط صفحات به کار می‌رود، در الگوریتم FishSearch دودویی و در الگوریتم SharkSearch پیوسته و حقیقی است. همچنین امتیازی که به پیوندهای موجود در صف پیمایش اختصاص داده می‌شود، در الگوریتم SharkSearch پیچیده‌تر است. در الگوریتم FishSearch میزان ارتباط صفحه وب با موضوع مورد نظر بر اساس تطابق کلمات یا عبارات تعیین می‌شود. این امتیاز به متن پیوند، متن اطراف پیوند، امتیاز صفحات اشاره‌کننده به صفحه

حاوی پیوند بستگی دارد. این الگوریتم را می‌توان گونه پیچیده-تری از الگوریتم BestFirst دانست.

الگوریتم‌های ارائه شده [12][11][10][9][8] روش‌هایی هستند که از یک دانش پیش زمینه برای پیمایش موضوعی استفاده می‌کنند.

روش‌های پیمایش موضوعی وب که تاکنون گزارش شده است، از پیوندهای بین صفحات وب، متن صفحات وب و Token‌های موجود در پیوندها استفاده می‌کنند. استفاده از این اطلاعات به تنهایی نمی‌تواند به خوبی میزان ارتباط یک صفحه با موضوعی خاص را نشان دهد. به علاوه برخی پدیده‌ها در وب مانند سایت‌های بی‌انتهای (Bottomless Sites)، تله‌های Spiderها و پیوندهای تبلیغاتی و مصنوعی، پیمایشگرها را با مشکل مواجه می‌کند.

در [1] و [2] روش‌های جدیدی مبتنی بر اتوماتای یادگیر توزیع شده جهت تعیین شباهت بین صفحات ارائه شده است. در [1] برای محاسبه شباهت بین صفحات وب از اتوماتای یادگیر توزیع شده استفاده شده است که شباهت بین صفحات را با استفاده از فعل وانفعال کاربر و رابطه تراگذاری تعیین می‌کند. در [2] اتوماتای یادگیر توزیع شده شباهت بین صفحات وب را با استفاده از اطلاعات پیمایش کاربران قبلی و پیوند گراف وب سایت تعیین می‌کند که نسبت به روش ارائه شده در [1] از کارایی بالاتری برخوردار است در این روش با افزایش تعداد صفحات، مقدار دقت کم نمی‌شود. برای بالا بردن کارایی الگوریتم با تعداد صفحات زیاد از پارتیشن بندی گراف با استفاده از الگوریتم‌های چند سطحی (Multilevel) 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 رابطه‌ای مشخص کاهش می‌یابد. جریمه دیگری که در روش پیشنهادی برای کاربر در نظر گرفته شده وجود دور در مسیر پیمایشی کاربر می‌باشد.

در این مقاله با استفاده از الگوریتم ارائه شده در [2] شباهت صفحات وب تعیین می‌گردد.

روش پیشنهادی در گروه روش‌های پیمایش موضوعی وب قرار می‌گیرد که از یک دانش پیش زمینه (طبقه بندی از اسناد، هستان شناسی، ...) برای پیمایش استفاده می‌کنند. دانش پیش زمینه در این جا همان ساختار ارتباطی است که با استفاده از اتوماتای یادگیر توزیع شده و پارتیشن بندی گراف بدست آمده است. به علاوه از آن جا که در روش پیشنهادی از متن صفحات

وب استفاده نمی‌شود، این روش از نظر هزینه محاسباتی و زمان پردازش از روش‌های پیمایش مبتنی بر محتوای صفحات به صرفه‌تر می‌باشد.

ویژگی دیگر روش پیشنهادی آن است که پیمایشگر طراحی شده یک پیمایشگر یادگیر است که می‌تواند از دانش یاد گرفته شده در پیمایش‌های متفاوت استفاده کند.

همچنین یکی دیگر از مشکلات پیمایشگرهای موضوعی راكد شدن (Stagnation) می‌باشد به آن معنی که پیمایشگر دیگر قادر به جمع‌آوری صفحات مرتبط نمی‌باشد. این مشکل هیچ‌گاه برای پیمایشگرهای معمولی رخ نمی‌دهد. اما در پیمایشگرهای موضوعی در صورتی که موضوع مورد نظر خاص و محدود باشد، صفحات وب مرتبط به اندازه کافی با یکدیگر پیوند نداشته باشند و یا پیمایشگر در انتخاب پیوندها اشتباه کرده باشد، این مشکل ممکن است رخ دهد. در پیمایشگر طراحی شده در این بخش از آن جا که پیمایش تنها بر اساس محتوا و پیوندهای بین صفحات انجام نمی‌شود، این مشکل کاهش می‌یابد.

ادامه مقاله بدین صورت سازماندهی شده است. در بخش 2 اتوماتای یادگیر و اتوماتای یادگیر توزیع شده به اختصار معرفی می‌شوند. در بخش 3 الگوریتم پیشنهادی و در بخش 4 پس از معرفی مدل استفاده شده برای شبیه‌سازی، نتایج شبیه‌سازی ارائه می‌شود. بخش 5 نتیجه‌گیری می‌باشد.

1- اتوماتاهای 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.

محیط را می‌توان توسط سه‌تایی $E = \{a, b, c\}$ نشان داد که در آن $a = \{a_1, a_2, \dots, a_r\}$ مجموعه ورودیها، $b = \{b_1, b_2, \dots, b_r\}$ مجموعه خروجیها و $c = \{c_1, c_2, \dots, c_r\}$ مجموعه احتمالات جریمه می‌باشد. اتوماتای یادگیر به دو دسته اتوماتای یادگیر با ساختار ثابت و اتوماتای یادگیر با ساختار متغیر تقسیم می‌گردند. در این مقاله از اتوماتای یادگیر با ساختار متغیر استفاده شده است که در ادامه معرفی می‌شود.

اتوماتای یادگیر با ساختار متغیر توسط چهارتایی $a = \{a_1, a_2, \dots, a_r, L, a_r\}$ نشان داده می‌شود که در آن $b = \{b_1, b_2, \dots, b_r, L, b_r\}$ مجموعه اقدام‌های اتوماتای یادگیر، $p = \{p_1, p_2, \dots, p_r, L, p_r\}$ مجموعه ورودیهای اتوماتای یادگیر، احتمال انتخاب هر یک از اقدام‌ها و T ، $p(n+1) = T[a(n), b(n), p(n)]$ الگوریتم یادگیری اتوماتای یادگیر می‌باشد. الگوریتم‌های یادگیری متنوعی برای اتوماتای یادگیر ارائه شده است با استفاده از الگوریتم یادگیری خطی، اتوماتای یادگیر بردار احتمال انتخاب اقدام‌های خود را مطابق رابطه (1) تنظیم می‌کند.

$$\begin{aligned} p_i(n+1) &= p_i(n) + a_i(1-b(n))(1-p_i(n)) - b_i(n)p_i(n) \\ p_j(n+1) &= p_j(n) + a_j(1-b(n))p_j(n) + \frac{b_j(n)}{r-1} - b_j(n)p_j(n) \quad \text{if } j \neq i \end{aligned} \quad (1)$$

که a پارامتر پاداش و b پارامتر جریمه می‌باشد. اگر a و b با هم برابر باشند، الگوریتم $L_R - p$ ، اگر b از a خیلی کوچکتر باشد، الگوریتم L_{Rep} و اگر b صفر باشد، الگوریتم L_{R-I} نام دارد [14].

1-1 اتوماتای یادگیر توزیع شده

اتوماتای یادگیر توزیع شده 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. بصورتی 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. در هر لحظه فقط یک اتوماتای یادگیر در اتوماتای یادگیر توزیع شده فعال می‌باشد [15][13].

2- روش پیشنهادی

روش پیشنهادی برای پیمایش موضوعی از دو مرحله کلی تشکیل شده است. در مرحله اول 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 [2]. در مرحله دوم، با استفاده از ساختار ارتباطی به دست آمده از مرحله قبل، پیمایش موضوعی وب انجام می‌شود و صفحات مرتبط با یک موضوع خاص جمع‌آوری

می گردد. پس از آن که مرحله اول برای یک بار انجام شد، می توان مرحله دوم را برای دفعات متعدد انجام داد.

پس از اصلاح احتمال انتخاب اعمال اتوماتای یادگیر توزیع شده، امتیاز Hub و Authority هر یک از صفحات مجموعه پایه طبق فرمول (2) و (3) بروز می شود.

$$Authority(i) = \sum_{j \rightarrow i} r(j,i) \times Hub(j) \quad (2)$$

$$Hub(i) = \sum_{i \rightarrow j} r(i,j) \times Authority(j) \quad (3)$$

نحوه عملکرد بخش دوم به شرح زیر است:

ابتدا موضوع مورد نظر برای پیمایش به عنوان ورودی به الگوریتم داده می شود. سپس مجموعه ای به نام مجموعه آغازین که حاوی تعدادی صفحه وب می باشد، ایجاد می گردد و پیمایشگر کار خود را با این مجموعه آغاز می کند. این مجموعه را می توان به روش های مختلف ایجاد کرد. در پیمایشگر طراحی شده در این قسمت، از آنجا که در مدل شبیه سازی استفاده شده هر صفحه دارای برداری است که میزان ارتباط این صفحه با موضوعات موجود در مدل شبیه سازی را نشان می دهد، از این بردار برای ایجاد مجموعه آغازین استفاده می کنیم. برای این منظور صفحات را به پنج گروه تقسیم می کنیم به طوری که میزان ارتباط گروه اول با موضوع انتخاب شده بین 0,8 و 1، گروه دوم بین 0,6 و 0,8، ... و گروه پنجم بین 0,2 و 0 باشد. سپس از هر گروه 20% صفحات مجموعه آغازین را به صورت تصادفی انتخاب می کنیم.

سپس یک صف اولویت که صف پیمایش نامیده می شود، ایجاد شده و هر یک از صفحات مجموعه آغازین در این صف قرار می گیرد. این صف شامل صفحاتی است که تا مرحله جاری الگوریتم توسط پیمایشگر پردازش نشده اند. هر صفحه در این صف امتیازی دارد که اولویت آن برای پیمایش را تعیین می کند. صفحات بر اساس این امتیاز مرتب شده و در صف پیمایش قرار می گیرند. در ابتدا امتیاز کلیه صفحات مجموعه آغازین با یکدیگر برابر می باشد. همچنین مجموعه ای به نام مجموعه پیمایش ایجاد می گردد که شامل صفحاتی خواهد بود که پیمایش می شوند و در این مرحله تهی می باشد. سپس تا زمانی که تعداد صفحات پیمایش شده به حد آستانه ای برسد، عملیات زیر تکرار می شود:

صفحه ای که در ابتدای صف پیمایش قرار دارد، انتخاب و از صف حذف می گردد و به مجموعه صفحات پیمایش شده اضافه می شود. سپس صفحاتی که این صفحه در ساختار ارتباطی به آن ها اشاره می کند، استخراج می شوند و برای هر

یک از آنها امتیازی که اولویتشان را برای پیمایش تعیین می کند، محاسبه می شود. این امتیاز به دو عامل بستگی دارد و برای صفحه ای مانند Page_j طبق رابطه (4) محاسبه می شود.

$$crawl_score(Page_j) = r(i,j) \times hub(i) \quad (4)$$

که در آن $r(i,j)$ میزان ارتباط صفحه i (صفحه ای که به صفحه j اشاره می کند) و صفحه j می باشد و از ساختار ارتباطی محاسبه شده به دست می آید. $hub(i)$ امتیاز Hub صفحه i است که با استفاده از الگوریتم HITS محاسبه می شود. در این جا مجموعه پایه همان صفحاتی است که تا کنون پیمایش شده اند. الگوریتم HITS بر روی این صفحات اعمال می شود و امتیاز Hub صفحه i محاسبه می گردد.

پس از آن که امتیاز هر یک از صفحات استخراج شده محاسبه شد، این صفحات در صف پیمایش قرار می گیرد و صف مرتب می گردد. سپس صفحه ای که در ابتدای صف قرار دارد از صف انتخاب می شود و عملیات فوق الذکر برای آن تکرار می شود.

نکته ای که در طراحی این پیمایشگر در نظر گرفته شده آن است که برای پیمایش موضوعی، تنها به کارگیری میزان ارتباط صفحات کافی نیست و باید صفحاتی برای پیمایش انتخاب شوند که علاوه بر مرتبط بودن دارای پیوند به صفحات دیگر نیز باشند. برای رفع این مشکل علاوه بر در نظر گرفتن میزان ارتباط یک صفحه با موضوع مورد نظر (که بر اساس ساختار ارتباطی تعیین می شود)، معیاری دیگری نیز برای انتخاب یک صفحه برای پیمایش در نظر گرفته می شود. این معیار همان امتیاز Hub می باشد که با استفاده از الگوریتم HITS محاسبه می شود. با این تفاوت که در الگوریتم HITS امتیاز Hub برای مجموعه پایه محاسبه می شود، اما در این جا امتیاز Hub را با استفاده از مجموعه صفحاتی که تا کنون پیمایش شده اند، محاسبه می کنیم. صفحه ای که امتیاز Hub بالایی دارد، به صفحات مرتبط با یک موضوع اشاره می کند و به این ترتیب مشکل فوق الذکر را برطرف می نماید.

از آن جا که روش پیشنهادی از جمله روش هایی است که از یک دانش پیش زمینه استفاده می کند، مانند روش های این گروه وابسته به کیفیت دانش پیش زمینه می باشد. همان گونه که در پیمایش موضوعی با استفاده از طبقه بندی کننده صفحات و پیمایش موضوعی مبتنی بر هستان شناسی، کیفیت پیمایش به طبقه بندی کننده و هستان شناسی استفاده شده بستگی دارد، کارایی روش پیشنهادی نیز به کیفیت ساختار ارتباطی به دست آمده وابسته است.

شبه کد الگوریتم پیشنهادی در شکل 1 آمده است.

```
// DLA based Focused Crawling Algorithm
Procedure DLA_Crawling_algorithm
  create Initial Crawl Set based on Crawl Topic
  fill the Crawl Queue with pages of Initial Set
  set the priority of all pages in Crawl Queue equally
  Crawl Set = Null
  while (not termination) do // a definite number of pages have
    been crawled
      pagei = pick the page at the head of Crawl Queue
      add pagei to Crawl Set
      PointedPageSet = extract the pages pointed by pagei
      for each pagej in PointedPageSet do
        hub(i) = compute hub score of pagei by applying HITS
        algorithm on Crawl
        Set r(i,j) = retrieve the relevance score of pagei and
        pagej from Relationship
        Structure
        Crawl_Score(pagej) = r(i,j) * hub(i)
        add pagej to Crawl Queue
      end for
    sort Crawl Queue
  end while
  output Crawl Set
```

شکل 1. شبه کد الگوریتم پیمایش موضوعی مبتنی بر DLA

تصادفی واقعی از صفحات مرتبط وب باشد، اما تخمین مناسبی برای محاسبه میزان فراخوان فراهم می کند. این معیار نیز مشابه معیار فراخوان در کاربردهای بازیابی اطلاعات می باشد.

3-2- مدل شبیه سازی

برای شبیه سازی الگوریتم پیشنهادی و مقایسه آن با سایر روشها از مدل معرفی شده در [17] برای نشان دادن ساختار صفحات وب و چگونگی استفاده کاربران، استفاده شده است. اعتبار این مدل توسط Lui و همکاران [17] با استفاده از اطلاعات استفاده از وب چندین سایت وب بزرگ مانند مایکروسافت، تایید شده است.

جدول 1: پارامترهای شبیه سازی ها

0/7	حد آستانه ایجاد اتصال
50000	تعداد کاربران
500	تعداد اسناد
5	تعداد موضوعها
0/2	T_c مقدار ثابت سند اولیه (صفحه اولیه سایت) در موضوعات مختلف
-	ΔM_l^c ضریب ثابت کاهش اشتیاق کاربر
-	ΔM_l^v ضریب متغیر کاهش اشتیاق کاربر
1	a_u پارامتر توزیع قانون-توانی توزیع احتمال علائق کاربران
1/2	f ضریب پاداش دریافتی از مشاهده یک سند
0/5	I ضریب جذب اطلاعات از یک سند توسط یک کاربر
5/97	m_m میانگین توزیع نرمال ΔM_l^v
0/25	s_m واریانس توزیع نرمال ΔM_l^v
-	m_l میانگین توزیع نرمال برای مقدار افزایش یک گره برای یک موضوع خاص
3	a_p پارامتر توزیع قانون-توانی توزیع احتمال وزنه های مطالب برای هر سند
0/25	s_l واریانس توزیع نرمال برای مقدار افزایش یک گره برای یک موضوع خاص
1	q ضریب کاهش علاقه کاربر
0/2	حداقل اشتیاق کاربر برای ادامه جستجو
25	تعداد صفحات در هر پارتیشن

بر این اساس، در این مقاله مطابق با مدل رفتار کاربران، پروفایل علاقه کاربران بصورت توزیع قانون-توانی و توزیع محتوای صفحات وب بصورت توزیع نرمال در نظر گرفته شده است. سایر پارامترهای استفاده شده برای شبیه سازیهای انجام شده در این مقاله در جدول 1 نشان داده شده است. تعداد

3- نتایج شبیه سازی

در این بخش نتایج بدست آمده از شبیه سازی الگوریتم پیشنهادی ارائه می شود.

3-1- معیار ارزیابی پیمایشگر

دو شاخص اصلی که برای ارزیابی عملکرد پیمایشگر به کار می روند، عبارتند از: نرخ حاصل و فراخوان هدف.

- **نرخ حاصل:** این شاخص میانگین میزان ارتباط صفحات بازیابی شده را اندازه گیری می کند. نرخ حاصل با $H(t)$ نشان داده می شود و پس از پیمایش t صفحه از رابطه (5) محاسبه می شود.

$$H(t) = \frac{1}{t} \sum_{i=1}^t r_i \quad (5)$$

که r_i میزان ارتباط صفحه i با موضوع پیمایش می باشد. نرخ حاصل، مشابه با معیار ارزیابی دقت است که در کاربردهای بازیابی اطلاعات استفاده می شود.

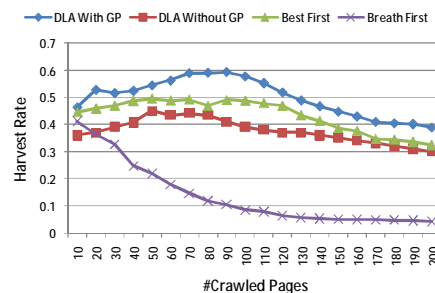
- **فراخوان هدف:** فراخوان هدف، میزان صفحات مرتبط بازیابی شده به کل صفحات مرتبط موجود در وب را اندازه گیری می کند. از آنجا که مجموعه کل صفحات مرتبط در وب، مشخص نیست، مجموعه ای به عنوان مجموعه هدف تعریف می شود که با استفاده از آن معیار فراخوان هدف به عنوان تخمینی از فراخوان واقعی محاسبه می شود. برای تعریف مجموعه هدف، می توان از یک نمونه تصادفی از صفحات موجود در وب استفاده کرد. هرچند، نمی توان انتظار داشت که این مجموعه یک نمونه

صفحات پیمایش شده در هر یک از این آزمایشات برابر 200 در نظر گرفته شده است.

همچنین پارتیشن بندی گراف توسط نرم افزار Metis انجام شده است [16]. در این نرم افزار پارتیشن بندی با استفاده از روشهای چند سطحی انجام می شود. این نرم افزار از الگوریتمهایی برای پارتیشن بندی گراف استفاده می کند که در سریعترین زمان، زیر گرافهایی با کیفیت بالا ایجاد کند.

3-3- مقایسه نتایج

مهمترین معیار ارزیابی پیمایشگر موضوعی نرخ حاصل است که میانگین میزان ارتباط صفحات بازیابی شده را نشان می دهد. برای ارزیابی این معیار سه موضوع از میان موضوعات موجود در مدل شبیه سازی انتخاب شده و پیمایش برای هر یک از آنها به چهار روش انجام شده است.



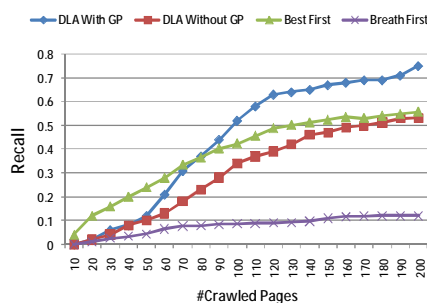
شکل 2. مقایسه نرخ حاصل برای چهار پیمایشگر

پیمایشگر اول، پیمایشگر عرض اول است و ضعیف ترین کارایی را در میان سایر پیمایشگرها دارد و به عنوان حد پایین برای ارزیابی عملکرد پیمایشگرها به کار می رود. پیمایشگر دوم، پیمایشگر BFS (Best First Search) است که نسبت به سایر پیمایشگرهای موضوعی بهتر عمل می کند. در این پیمایشگر اولویت پیوندها بر اساس میزان شباهت صفحه حاوی پیوند و مجموعه ای از صفحات مرتبط تعیین می شود. پیمایشگر سوم، پیمایشگر مبتنی بر DLA می باشد و پیمایشگر چهارم ترکیبی از DLA و پارتیشن بندی گراف می باشد.

همان طور که در شکل 2 دیده می شود، کارایی پیمایشگر عرض اول که یک پیمایشگر غیرموضوعی است، به سرعت کاهش می یابد و پس از پیمایش تعداد صفحات محدودی (200 صفحه)، میزان صفحات مرتبط بازیابی شده به صفر کاهش می یابد. در حالی که سه پیمایشگر موضوعی از کارایی بسیار بهتری برخوردارند و عملکرد پیمایشگرهای پیشنهادی بهتر از پیمایشگر BFS یا برابر با آن هستند.

شکل 3 فراخوان هدف نسبت تعداد صفحات مرتبط بازیابی شده به کل صفحات مرتبط را نشان می دهد.

در این آزمایش همچون آزمایش قبل پیمایش با استفاده از پیمایشگر عرض اول، پیمایشگر BFS، پیمایشگر مبتنی بر DLA و پیمایشگر ترکیبی مبتنی بر DLA و پارتیشن بندی گراف انجام شده و فراخوان هدف در مراحل مختلف پیمایش اندازه گیری شده است. همان طور که در شکل 3 مشاهده می شود، فراخوان هدف در پیمایشگر عرض اول ضعیف ترین و پیمایشگر مبتنی بر DLA و پارتیشن بندی گراف بهترین می باشد.

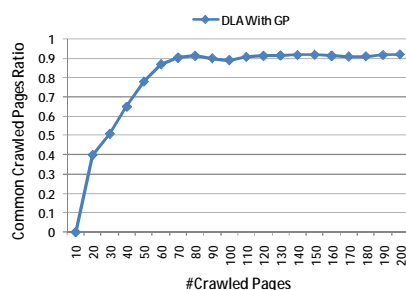


شکل 3. مقایسه فراخوان هدف برای چهار پیمایشگر

همانطور که در شکل 3 مشاهده می شود در هر چهار پیمایشگر پس از پیمایش تعداد مشخصی از صفحات فراخوان هدف به مقدار ثابتی میل می کند، چرا که تعداد صفحات مرتبط جدیدی که بازیابی می شود کاهش می یابد.

یکی از معیارهای دیگری که برای ارزیابی کارایی پیمایشگرهای موضوعی به کار می رود، میزان وابستگی آن به مجموعه آغازین می باشد.

برای ارزیابی به صورت تصادفی دو مجموعه از صفحات آغازین انتخاب می کنیم.



شکل 4. وابستگی به مجموعه آغازین

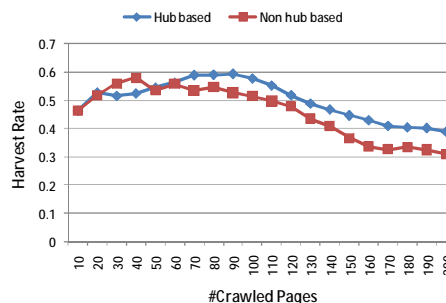
همان طور که در شکل 4 دیده می شود، تعداد صفحات مشترک پیمایش شده به سرعت به 90 درصد میل می کند و این نتیجه نشان می دهد که پیمایشگر مبتنی بر DLA وابسته به چگونگی انتخاب مجموعه آغازین نمی باشند.

پیوندهای بین صفحات در پیمایش موضوعی صفحات می تواند تاثیر بسزایی در بهبود نتایج داشته باشد.

مراجع

- [1] SH. Motamedi Mehr, M. Taran, and M. R. Meybodi, "Web Usage Mining based on Distributed Learning Automata", Proceedings of 10th Fuzzy Conference of Iran, Shahid Beheshti University, Tehran, Iran, 2010.
- [2] 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, 2011.
- [3] Pinkerton, B., "Finding What People Want: Experiences with the WebCrawler," Proceedings of the 2nd International World Wide Web Conference (Chicago), 1994.
- [4] Cho, J., Garacia-Molina, H., Page, L., "Efficient Crawling Through URL Ordering," Comput Netw. 30, pp. 161-172, 1998.
- [5] Hersovici, M., Javoci, M., Maarek, Y.S., Pelleg, D., Shtalham, M., "The Shark-Search Algorithm An Application: Tailored Web Site Mapping," Proceedings of the 7th International World-Wide Web Conference, 1998.
- [6] Page, L., Brin, S., Motwani, R., Winograd, T., "The PageRank Citation Ranking: Bringing Order to the Web," Stanford Publications, 1998.
- [7] Debra, P., Post, R., "Information Retrieval in the World Wide Web: Making Client-based Searching Feasible," Proceedings of the 1st International World Wide Web Conference (Geneva), 1994.
- [8] Chakrabarti, S., Van den Berg, M., Dom, B., "Focused Crawling: A New Approach to Topic-Specific Web Resource Discovery," Proceedings of the 8th International WWW Conference, Toronto, Canada, May 1999.
- [9] Diligenti, M., Coetzee, F., Lawrence, S., Giles, C., Gori, M., "Focused Crawling Using Context Graphs," Proceedings of the 26th International Conference on Very Large Databases (VLDB 2000), Cairo, Egypt, September 2000.
- [10] Su, C., et al, "An Efficient Adaptive Focused Crawler Based on Ontology Learning," Proceedings of the Fifth International Conference on Hybrid Intelligent Systems (HIS'05), 2005.
- [11] Aggarwal, C., Al-Garawi, F., Yu, P., "Intelligent Crawling on the World Wide Web with Arbitrary Predicates," Proceedings of the 10th International World Wide Web Conference, Hong Kong, May 2001.
- [12] Pant, G., Srinivasan, P., "Learning to Crawl: Comparing Classification Schemes," ACM Transactions on Information Systems, Vol. 23, No. 4, pp. 430-462, October 2005.
- [13] H. Beigy, and M. R. Meybodi, "Utilizing Distributed Learning Automata to Solve Stochastic Shortest Path Problem," International Journal of Uncertainty, Fuzziness and Knowledge-based Systems, vol. 14, no. 5, pp. 591-617, 2006.
- [14] K. S. Narendra and M.A.L. Thathachar, *Learning Automata: An Introduction*, Prentice Hall, 1989.
- [15] M. A. L. Thathachar and R. Harita Bhaskar, "Learning Automata with Changing Number of Actions," IEEE Transactions on Systems Man and Cybernetics, vol. 17, no. 6, Nov. 1987, pp. 1095-1100.
- [16] G. Karypis and V. Kumar, "METIS: A Software Package for Partitioning Unstructured Graphs, Partitioning Meshes, and Computing Fill-Reducing Orderings of Sparse Matrices," Version 5.0pre2. 2007: Minneapolis.
- [17] J. Liu, S. Zhang, and J. Yang, "Characterizing Web Usage Regularities with Information Foraging Agents," IEEE Transactions on Knowledge and Data Engineering, vol. 16, no. 4, April 2004, pp. 566-584.

یکی از ویژگی‌های پیمایشگرهای پیشنهاد شده آن است که در تخصیص امتیاز پیمایش به پیوندها، علاوه بر نظر گرفتن میزان ارتباط صفحه حاوی پیوند با موضوع پیمایش، امتیاز Hub صفحه حاوی پیوند نیز در نظر گرفته می‌شود. نتایج این آزمایش برای پیمایشگر مبتنی بر DLA در شکل 5 نشان داده شده است. برای ارزیابی این ویژگی، پیمایشگرهای موضوعی طراحی شده را در دو حالت مختلف با در نظر گرفتن امتیاز Hub و بدون آن اجرا می‌نماییم.



شکل 5. تاثیر به کارگیری امتیاز Hub در پیمایشگر مبتنی بر DLA

همان طور که در شکل دیده می‌شود، در ابتدای پیمایش نتایج بهتری در حالت بدون در نظر گرفتن امتیاز Hub حاصل می‌شود. در حالی که با احتساب امتیاز Hub در مراحل میانی نتایج از حالت اول بهتر می‌شود.

4- نتیجه گیری

روشهای پیشنهادی از جمله روش هایی است که از یک دانش پیش زمینه (ساختار ارتباطی) استفاده می کند. این پیمایشگرها، یک پیمایشگر یادگیر هستند که می توانند از دانش یاد گرفته شده (ساختار ارتباطی) در پیمایش های متفاوت استفاده کنند. بنابراین نسبت به سایر پیمایشگرهایی که محاسبات و پردازش های خود را در هر بار پیمایش تکرار می کنند، به صرفه تر می باشند. هر چند کارایی روش های پیشنهادی به کیفیت ساختار ارتباطی به دست آمده وابسته می باشند. همچنین این پیمایشگرها به دلیل به کارگیری امتیاز Hub در انتخاب صفحات، صفحاتی را انتخاب می کنند که حاوی پیوندهای بیشتری هستند و به این ترتیب دارای گزینه های بیشتری برای انتخاب می باشد.

نتایج ارزیابی پیمایشگرها پیشنهادی حاکی از بهبود عملکرد آن (نرخ حاصل و فراخوان هدف) نسبت به پیمایشگر Breath First یا Best First و عدم وابستگی آن به انتخاب مجموعه آغازین دارند. به این ترتیب استفاده از رفتار کاربران و