

تاریخ دریافت مقاله: ۹۷/۰۸/۰۶

تاریخ پذیرش مقاله: ۹۸/۰۷/۲۰

ترکیب الگوریتم HITS با الگوریتم Distance Rank برای بهبود نتایج در موتورهای جستجو

رعنا میلانی اباجلو

گروه مهندسی کامپیوتر، واحد ارومیه، دانشگاه آزاد اسلامی، ارومیه، ایران.

پست الکترونیکی: r.milani@abj@gmail.com

فرهاد سلیمانیان قره چیق*

گروه مهندسی کامپیوتر، واحد ارومیه، دانشگاه آزاد اسلامی، ارومیه، ایران.

پست الکترونیکی: bonab.farhad@gmail.com

چکیده:

واحد ارومیه انجام گرفته که نتایج نشان می‌دهد که روش پیشنهادی در مقایسه با الگوریتم‌های دیگر عملکرد بهتری دارد و توانسته است رتبه‌بندی متفاوت و بهتری نسبت به الگوریتم پایه HITS و سایر الگوریتم‌های رتبه‌بندی مانند Distance Rank و PR و WPR داشته باشد. همچنین الگوریتم‌های پیشنهادی بر مبنای معیارهای AP، P@n و NDC مورد ارزیابی قرار گرفت که نتایج نشان می‌دهد که روش پیشنهادی به ترتیب مقدار ۱ و ۱ و ۸/۱ را به دست آورده است.

واژه‌های کلیدی: موتورهای جستجو، الگوریتم HITS،

الگوریتم Distance Rank

۱. مقدمه

در دنیای اطلاعاتی امروز که هزاران گیگابایت اطلاعات در بین کاربران رد و بدل می‌شود، فضای وب یکی از بهترین فضاهای استفاده است. از سویی پیدایش و پیشرفت وب در دهه اخیر موجب تحول شگرفی در فراگیری اخبار

امروزه موتورهای جستجوگر از روش‌های وب‌کاوی برای نشان دادن نتایج بهتر استفاده می‌کنند که در لیست نتایج خود پیوندهای زیادی از صفحات وب را به کاربران نمایش می‌دهند و برای بهینه و محدود کردن لیست نتایج موتورهای جستجو از الگوریتم‌های رتبه‌بندی استفاده می‌شود. در این مقاله یک روش جدید که ترکیبی از الگوریتم HITS با الگوریتم Distance Rank است برای بهبود نتایج در موتورهای جستجو ارائه شده است که در روش پیشنهادی از فرایند اصلی الگوریتم Distance Rank برای بهبود الگوریتم HITS استفاده شده است. مشکل اصلی الگوریتم HITS این است که رتبه‌بندی صفحات وب بر اساس میزان ارتباط آن‌ها با پرس وجوی کاربر است. اما در الگوریتم Distance Rank از فاصله لگاریتمی میان صفحات به منظور رتبه‌بندی استفاده می‌شود. ارزیابی روش پیشنهادی بر روی سه مجموعه داده شامل گراف استاندارد، گراف تصادفی، گراف دانشگاه آزاد اسلامی

* نویسنده مسئول

و اطلاعات علمی در سراسر دنیا شده است؛ این امر باعث افزایش اطلاعات در جهان شده و مشکل پیدا کردن اطلاعات ارزشمند و معتبر از بین حجم وسیع اطلاعات را در پی دارد [۱-۲]. به همین خاطر جستجو مطالب امروزه به امر مهمی تبدیل شده است. موتورهای جستجو از جمله راه‌حل‌های مناسب برای حل این مسئله هستند که شرکت‌هایی همچون گوگل به‌طور مدام در حال یافتن بهترین و سریع‌ترین روش برای کاربران خود هستند. موتورهای جستجو برخلاف تصور کاربران که فکر می‌کنند این موتورها مطالب درخواستی‌شان را از بین تمام مطالب موجود در جهان جستجو می‌کند، در واقع آن‌ها این مطالب را از بین مطالب ذخیره‌شده در پایگاه خود جستجو کرده و در اختیار کاربران در آن لحظه جستجو قرار می‌دهند [۲-۴]. تعداد بینندگان یک وب‌گاه از مهم‌ترین عاملی که صفحات وب را دارای اهمیت می‌کند. در واقع موتورهای جستجو بینندگان را به وب‌گاه می‌آورند بدون توجه به این‌که وب‌گاه چه خدمتی ارائه می‌دهد یا چه چیزی می‌فروشد. اکثر مردم فقط ۱۰ وب‌گاه اول معرفی‌شده توسط موتور جستجو را مرور می‌کنند؛ بنابراین موتورهای جستجو و الگوریتم‌های بازبازی اطلاعات و الگوریتم رتبه‌بندی صفحات وب اطلاعات نقش حیاتی را در دنیای امروزی اطلاعات دارد [۵،۶]. الگوریتم‌های رتبه‌بندی صفحات وب، مجموعه‌ای از دستورالعمل‌ها است که موتور جستجوگر با اعمال آن‌ها بر پارامترهای صفحات موجود در پایگاه داده‌اش، تصمیم می‌گیرد که صفحات مرتبط را چگونه در نتایج جستجو مرتب کند. در واقع، رتبه‌بندی ارزش یک صفحه توسط موتور جستجو را مشخص کرده و صفحاتی با کیفیت بالا مبتنی بر درخواست و سلیقه کاربر را پیدا می‌کند [۲].

الگوریتم‌های زیادی برای رتبه‌بندی صفحات وب در سال‌های اخیر ارائه شده است [۷،۸-۹]. این الگوریتم‌ها یکی از اجزای اصلی موتورهای جستجو هستند و هدف آن‌ها تهیه یک رتبه برای هر صفحه وب است و همچنین

فضای جستجو را به میزان زیادی کاهش می‌دهند. یکی از الگوریتم‌های رایج برای رتبه‌بندی صفحات وب HITS است [۱] که برای هر پرس‌وجو تحلیل پیوندها انجام می‌شود. الگوریتم HITS یکی از الگوریتم‌های معروف برای رتبه‌بندی صفحات وب بر اساس میزان ارتباط آن‌ها با پرس‌وجوی کاربر است. این الگوریتم از دسته روش‌های وابسته به پرس‌وجو است و برای هر پرس‌وجو تحلیل پیوندها انجام می‌شود که می‌بایست گراف خاص پرس‌وجو به نام گراف همسایگی ساخته شود. از مسائل عمده رتبه‌بندی بر اساس پیوند که الگوریتم HITS یکی از این الگوریتم‌ها می‌باشد دارای مشکل غنی‌تر شدن اغنیاء^۱ می‌باشد. قرار گرفتن همیشه صفحات محبوب^۲ در صدر لیست ارائه‌شده به کاربر، باعث می‌شود تا کاربر فقط صفحات خاصی را ببیند و در نتیجه صفحات تازه متولدشده با کیفیت بالا که کسی به آن‌ها اشاره نمی‌کند نتوانند در دید کاربران قرار گیرند. این مشکل باعث می‌شود صفحات محبوب مرتباً محبوب‌تر شده و تعداد پیوند به آن‌ها افزایش یابد. لذا موتورهای جستجو با ارائه نکردن عادلانه اطلاعات به کاربران، باعث صرف وقت زیاد و در نتیجه کندی تولید علم و دانش خواهند شد. مهم‌ترین چالش در موتورهای جستجو رتبه‌بندی صفحات در پاسخ به پرس‌وجوی کاربر می‌باشد. از طرف دیگر الگوریتم Distance Rank مبتنی بر یادگیری تقویتی است [۱۰] و از فاصله لگاریتمی میان صفحات به‌عنوان تویخ دریافتی استفاده می‌کند. هدف آن کمینه کردن کل تویخ‌های دریافتی می‌باشد که باگذشت زمان آگاهی این الگوریتم بیشتر شده و عمل پیمایش را بهتر انجام می‌دهد. الگوریتم Distance Rank در رتبه‌بندی نسبت به بقیه مؤثرتر می‌باشد و حساسیت کمتری به مشکل غنی‌تر شدن اغنیاء نشان می‌دهد. در این مقاله از فرایند اصلی Distance Rank که مبتنی بر یادگیری تقویتی است برای پوشش نقطه ضعف الگوریتم HITS استفاده خواهیم کرد. روش پیشنهادی با سایر الگوریتم‌ها رتبه‌بندی از

1- Rich-get-richer
2- Popular

جمله Page Rank، Wighted Page Rank، Distance Rank، HITS، بر روی سه مجموعه شامل گراف استاندارد، گراف تصادفی، گراف دانشگاه آزاد اسلامی واحد ارومیه اجرا شده است که نتایج روش پیشنهادی امیدوار کننده است. ساختار کلی مقاله به شرح زیر سازماندهی شده است: در ادامه مقاله ابتدا در بخش دوم کارهای مرتبط با رتبه‌بندی صفحات وب را مورد بررسی قرار خواهیم داد. در بخش سوم مقاله در ابتدا الگوریتم HITS تشریح می‌شود سپس الگوریتم Distance rank تشریح می‌شود. در بخش چهارم، روش پیشنهادی بیان شده است. در بخش پنجم مقاله نتایج حاصل از بررسی و شبیه‌سازی و مجموعه داده‌ها و ارزیابی روش پیشنهادی و معیارهای ارزیابی آورده شده است. در بخش ششم مقاله نتیجه‌گیری کلی ارائه می‌شود.

۲. کارهای مرتبط

در این بخش به پژوهش‌های انجام‌شده در زمینه الگوریتم‌های رتبه‌بندی صفحات وب می‌پردازیم. رضا فتوح و همکاران در [۱۱] به ارزیابی کارایی الگوریتم‌های رتبه‌بندی صفحه برای استخراج صفحات وب پرداخته‌اند که برای حل مشکلات الگوریتم‌های رتبه‌بندی Page Rank(PR) و Weighted Page Rank (WPR) یک الگوریتم جدید به نام Weighted Page Content Rank ارائه شده است. نتایج تجزیه و تحلیل نشان می‌دهد که الگوریتم WPR دارای عملکرد بهتری نسبت به الگوریتم Page Rank می‌باشد. در پژوهش دیگر زارع بیدکی و همکاران در [۱۲] به الگوریتم ترکیبی افقی جهت رتبه‌بندی صفحات وب پرداخته‌اند که در این طرح روش‌های متعدد رتبه‌بندی بیان شده است که نتایج به دست آمده بهبود چشمگیری را در روش ترکیبی ارائه شده در مقایسه با بقیه الگوریتم‌ها نشان داد. فرصتی و همکاران یک الگوریتم مبتنی بر ساختار پیوندی صفحات و اطلاعات استفاده کاربران برای پیشنهاد صفحات وب در [۱۳] ارائه داده‌اند.

آن‌ها پس از معرفی معیار وزن‌دهی، الگوریتمی ترکیبی که از اطلاعات پیمایش کاربران و پیوند بین صفحات به منظور پیشنهاد صفحات به کاربران استفاده می‌کند، ارائه کردند. همچنین یک الگوریتم یادگیری مبتنی بر خصیصه‌های ذاتی جهت رتبه‌بندی صفحات وب در [۱۴] ارائه شده است. الگوریتم ارائه شده با اعمال وزن روی خصیصه‌های صفحات و مقایسه وزن اعمال شده با ارزش خصیصه صفحات، سعی در یافتن رتبه مناسب به منظور دسترسی راحت‌تر کاربران به صفحات با کیفیت از لحاظ محتویات داخلی شده است. جهت ارزیابی مجموعه داده‌ای LETOR به کار برده شده است. نتایج آزمایش‌ها حاکی از آن است که الگوریتم ارائه شده در مقایسه با الگوریتم‌های یادگیری که تاکنون ارائه شده است نتایج قابل توجهی را ارائه می‌دهد. در پژوهش دیگر کلانتری و همکاران در [۱۵] یک الگوریتم رتبه‌بندی صفحات وب مبتنی بر رویکرد یادگیری را ارائه داده‌اند. روش پیشنهادی با استفاده از PHP شبیه‌سازی و برای ارزیابی نتایج از معیار P@N استفاده شد. نتایج مثبت شده و ارزیابی نشان می‌دهد که روش پیشنهادی در این طرح دارای کارایی و دقت بالاتری نسبت به روش‌های موجود هستند و همچنین روش پیشنهادی بر مبنای رفتار و قضاوت کاربران و منصفانه است.

سورچی چاولا^۳ یک رویکرد جدید برای تولید مکان یکنواخت منبع‌های کلیک شده دارای رتبه بهینه با استفاده از الگوریتم ژنتیک بر اساس جلسات خوشه‌بندی شده جستجوی وب جهت جستجوی مؤثر در وبگاه‌های شخصی را ارائه داده است [۱۶] که روش پیشنهادی نتایج امیدوار کننده از خود نمایش گذاشته است. در پژوهش دیگر دویکا سینگ^۴ و دایا گاپاتا^۵ در [۱۷] یک الگوریتم جدید رتبه‌بندی به نام رتبه‌بندی صفحه بر اساس اولویت کاربر، ارائه و پیشنهاد دادند که از لحاظ ارتباطی کارآمد است زیرا آن عواملی را به کار می‌برد تا صفحات دارای مطالب مرتبط را تعیین کند و همچنین رفتار کاربر در نظر

3- Suruchi Chawla

4- Devika Singh

5- Daya Gupta

گرفته می‌شود. پاندا^۶ و سوت^۷ در [۱۸] بیان می‌کنند که امروزه، شبکه گسترده جهانی یک رسانه محبوب و جذاب برای انتشار اطلاعات است. در این مقاله، در مورد الگوریتم رتبه‌بندی صفحه بحث شده است که می‌تواند رتبه‌بندی بالاتری را برای صفحات مهم وب ایجاد کند. در [۱۹] جستجوی شخصی یک حوزه تحقیقی ضروری توسط رانی^۸ و سورانا^۹ را بیان کرده‌اند که هدف اصلیش این است که عدم اطمینان و شک یا تردید واژه‌های جستجو را تشخیص دهد. در این مقاله، طرحی پیشنهاد داده شده است که اولویت‌های مفهومی و ادراکی کاربر از طریق کلیک کاربران دیگر روی داده‌های حاصل از جستجو در وب مورد بررسی قرار می‌دهد.

در سال ۲۰۱۷، سن^{۱۰} و همکاران نسخه کارآمدی رتبه‌بندی ساده صفحه با زمان، فاکتور هدایت و مترادف ارائه دادند [۲۰]. این مقاله نسخه بازبینی شده‌ای از رتبه‌بندی صفحه را با استفاده از همان دستورالعمل کلاسیک و برخی از مواد دیگر اضافه شده مانند زمان، پیوندها، هدایت، مترادف ارائه می‌دهد که آن را برای کاربری که در الگوریتم پرکاربرد Rank Ranking محدودیت‌هایی دارد، و سوسه‌انگیزتر می‌کند. لاکشمی^{۱۱} و همکاران در سال ۲۰۱۸ مقاله‌ای تحت عنوان هدایت پویا حاصل از نتایج پرس وجو بر اساس نمایه‌گذاری درهم‌سازی با استفاده از الگوریتم بهبود یافته رتبه‌بندی PageRank ارائه دادند [۲۱]. در روش ارائه شده در این مقاله یک راه حل براساس دسته‌بندی پرس وجوهای وب به طور پویا و با استفاده از ساختار نمایه‌گذاری درهم‌سازی داده‌ها پیشنهاد شده است و صفحات وب حاصل با استفاده از الگوریتم بهبود یافته رتبه‌بندی PageRank رتبه‌بندی می‌شوند. با استفاده از این روش، نتایج بسیار غیرمرتبط و نتایج بسیار مهم بر اساس محتوا و تعداد پیوندها به عنوان نتایج برتر برای پرس

6- Panda
7- Sote
8- Rani
9- Sorana
10- Sen
11- Lakshmi

وجو داده شده، کاهش یافته‌اند. در سال ۲۰۱۹، ستی^{۱۲} و دیکسیت^{۱۳} یک سازوکار جدید رتبه‌بندی صفحه براساس الگوهای جستجوی کاربر ارائه دادند [۲۲]. در این مقاله، یک سازوکار جدید رتبه‌بندی صفحه مبتنی بر الگوهای جستجوی کاربر و پیوندهای بازدید شده، پیشنهاد شده است. نتایج شبیه‌سازی شده نشان می‌دهد که سازوکار رتبه‌بندی پیشنهادی بهتر از سازوکار PageRank معمولی در ارائه نتایج رضایت بخش به کاربر عمل می‌کند. در جدول (۱)، مزایا و معایب هر مدل نشان داده شده است. مدل‌های مختلفی برای رتبه‌بندی صفحات وب پیشنهاد شده است. اما هر مدل می‌تواند مزایا و معایبی داشته باشد.

۳. الگوریتم‌های پایه

در این بخش ما به معرفی دو الگوریتم HITS و Distance Rank خواهیم پرداخت. البته سعی شده است در این بخش هر یک از الگوریتم‌ها به صورت نظری و ریاضی مورد بررسی قرار گرفته و سپس در بخش بعدی (۴) به جزئیات روش پیشنهادی پرداخته خواهد شد.

الگوریتم HITS یک الگوریتم وابسته به پرس وجو می‌باشد که در سال ۱۹۹۸ توسط کلینبرگ^{۱۴} ارائه گردید [۱]. از آنجایی که این الگوریتم یک روش مبتنی بر پرس وجو است برای هر پرس وجو تحلیل پیوندها انجام می‌شود. برای تحلیل پیوندها، ابتدا گراف خاص پرس وجو به نام گراف همسایگی ساخته می‌شود که در حالت ایده‌آل شامل صفحات مرتبط یا موضوع پرس وجو می‌باشد. در شکل (۱) شبه کد الگوریتم HITS نشان داده شده است.

در این الگوریتم برای ساخت گراف همسایگی، ابتدا یک مجموعه از اسناد مرتبط با پرس وجو، به وسیله موتور جستجو واکشی می‌شوند. به این مجموعه، مجموعه ریشه گفته می‌شود. سپس مجموعه ریشه به وسیله

12- Sethi
13- Dixit
14- Kleinberg

جدول ۱: مزایا و معایب مدل‌های پیشنهاد شده برای رتبه‌بندی صفحات وب

مراجع	مدل	مزایا	معایب	مدل‌های مقایسه‌ای
[۱۱]	Weighted Page Content Rank	این الگوریتم اطلاعات مهم و مرتبط با پرس و جو را ارائه می‌دهد.	اختصاص وزن به پیوندهای نامرتب	Page Rank(PR) و Weighted Page Rank (WPR)
[۱۲]	الگوریتم ترکیبی افقی	کشف پیوندهایی با بیشترین بازدید و بیشترین تعداد مشاهده	افزایش زمان و عدم تشخیص بهترین پیوندها	-
[۱۳]	الگوریتم مبتنی بر وزندهی	الگوریتم بر پایه وزن هر پیوند و نرمال‌سازی آن‌ها کار می‌کند و مجموعه‌هایی از صفحات که مرتبط به پرس و جو هستند را جمع‌آوری می‌کنند.	اختصاص امتیاز به پیوندهایی که خود کار ایجاد شده‌اند و به پرس و جو کاربر مرتبط نیستند.	الگوریتم‌های مبتنی بر وزن
[۱۴]	الگوریتم مبتنی بر وزندهی	اختصاص وزن به پیوندهایی که مرتبط با کلمات کلیدی پرس و جو هستند.	نمایش صفحات نامرتب با پرس و جو کاربر به دلیل وزن یکسان صفحات	Page Rank(PR)
[۱۵]	الگوریتم رتبه‌بندی صفحات وب مبتنی بر رویکرد یادگیری	با استفاده از رفتار کاربران و استفاده از یادگیری از الگوریتم PageRank در حالت بهینه برای رتبه‌بندی صفحات استفاده می‌شود.	عدم تکرار ناکافی به منظور مرحله یادگیری و آموزش الگوریتم برای تشخیص پیوندهای مهم	Page Rank(PR)
[۱۶]	الگوریتم ژنتیک خوشه‌بندی	خوشه‌بندی پیوندهای بر مبنای نتایج اولیه و تعداد نمونه‌های مشابه در پیوندها انجام می‌شود.	عدم کشف مرکزیت خوشه به منظور یافتن پیوندهای مشابه	-
[۱۷]	الگوریتم اولویت	بازیابی پیوندها بر مبنای شاخص‌های مهم از قبیل رتبه و بگانه، تعداد کاربران و تعداد بازدید	عدم شباهت‌های ساختاری و میزان ارزشمندی امتیازدهی و سپس عمل رتبه‌بندی به تمامی صفحات	HITS
[۱۸]	الگوریتم اولویت	شناسایی کلمات مهم در پیوندها بر مبنای نتایج اولیه و تعداد مشاهدات و بگانه	یافتن صفحاتی که عدم ارتباط پیوندهای مفید میان آن‌ها پایین است.	Page Rank(PR)
[۱۹]	الگوریتم اولویت	در این الگوریتم همه صفحات به همراه سایر صفحاتی که به عنوان پرچسب یا پیوند به این صفحات پیوند داده شده‌اند و همچنین این صفحات به صفحات دیگری پیوند داده شده‌اند به عنوان یک مجموعه در نظر گرفته می‌شوند.	عدم پردازش نشانی پیوندها و نشانی‌های موجود در صفحات وب	Page Rank(PR)
[۲۰]	بهبود الگوریتم رتبه‌بندی	گروه‌ها یا دسته‌های حاصل را می‌توان بر اساس ویژگی‌ها به زیردسته‌های دیگری تقسیم نمود.	عدم شباهت‌های ساختاری رشته‌ای در پیوندها	-
[۲۱]	بهبود الگوریتم رتبه‌بندی	رتبه‌بندی پیوندها بر مبنای تعداد رجوع به لینک‌ها و تعداد مشاهدات	یافتن لینک‌های جعلی و نامرتب با پرس و جو کاربر	-
[۲۲]	الگوی جستجوی کاربران	بر مبنای پرس و جو کاربران، دنباله‌ای از ضروری‌ترین پیوندها نمایش داده می‌شود.	یافتن صفحاتی که عدم ارتباط پیوندهای مفید میان آن‌ها پایین است.	Page Rank(PR)

بزرگی شود، این عدد محدود و برای تعداد این اسناد حدی در نظر گرفته می‌شود. به این مجموعه جدید، مجموعه پایه یا گراف همسایگی گفته می‌شود. کلینبرگ [۱] صفحات وب را در دو فرم مرجع (Authorities) و قطب (Hub) تعریف می‌کند. به این صورت که یک صفحه مرجع دارای محتوای

همسایگانش تکمیل می‌گردد. همسایه‌ها، مجموعه‌ای از اسناد هستند که یا از اسناد موجود در مجموعه ریشه به آن‌ها پیوند داده شده است و یا به اسناد موجود در مجموعه ریشه پیوند داده‌اند. از آنجا که تعداد اسنادی که به اسناد موجود در مجموعه ریشه پیوند داده‌اند ممکن است عدد

```

//D is Distance Rank vector
D ← {∞, ∞, ∞, ...}
t ← 0; //iteration number
while δ > ε
    α ← e-β*t
    t ← t+1
    For every page j ∈ V
        Dt[j] ← (1-α)*Dt-1[j] + mini (γ*Dt-1[i] + log(O(i)))
        i ∈ B(j), α ≤ α ≤ 1, γ ≤ γ ≤ 1
    Normalize D
    δ ← ||Dt - Dt-1||
End while

```

شکل ۲: شبه کد الگوریتم Distance Rank

از فاصله لگاریتمی میان صفحات به منظور رتبه‌بندی استفاده می‌کند. منظور از فاصله میان دو صفحه A و Z وقتی A به Z اشاره کند، لگاریتم درجه خروجی A (تعداد پیوندها) می‌باشد. الگوریتم Distance Rank به صورت شکل (۲) محاسبه می‌شود که ϵ نشان‌دهنده خطا است. بیدوکی^{۱۰} در [۱۰] برای درک بهتر این مفهوم تعاریف زیر را ارائه کرده است:

(۱) اگر صفحه A به صفحه Z اشاره کند وزن پیوند میان A و Z برابر است با $\log O(i)$ که $O(i)$ نشان‌دهنده درجه خروجی صفحه A می‌باشد.

(۲) فاصله لگاریتمی میان A و Z عبارت است از وزن کوتاه‌ترین مسیر میان آن‌ها (جمع وزنه‌ای پیوندهای در مسیر) که با d_{ij} نشان داده می‌شود. بنابراین به جای فاصله عادی، تعریف جدیدی از فاصله به نام میانگین کلیک ارائه شده است. به عبارت دیگر وزن یک پیوند به جای $\log O(i)$ ، 1 می‌باشد.

تعاریف بالا نشان می‌دهند که اگر صفحه بارگذاری شده A دارای فاصله d_i باشد، فاصله هر کدام از بچه‌هایش از ریشه (صفحات اشاره شده توسط A در صورتی که ورودی دیگری نداشته باشند، با رابطه (۲) محاسبه می‌شود.

$$d_j = d_i + \log O(i) \quad (2)$$

که در رابطه (۲) d_i فاصله فرزند Z در صفحه A از ریشه می‌باشد و $\log O(i)$ ، به ازای دریافتی حاصل از انتقال از A به Z می‌باشد. در روش فاصله اگر یک صفحه دارای تعداد زیادی پیوند باشد فاصله کمتری نسبت به بقیه دارد و

```

authorities ap and hubs hp of the pages are stored in the vectors a and h
V ← collection of n pages
N ← number of iterations
z ← (1, 1, ..., 1) ∈ ℝ|V|
a0 ← z
h0 ← z
for i = 0 to N do
    {apply Eq. 5.27 to (ai-1; hi-1) and draw the new authority vector âi}
    {apply Eq. 5.28 to (âi; hi-1) and draw the new hub vector ĥi}
    ai ← âi
    hi ← ĥi
end for

```

شکل ۱: شبه کد الگوریتم HITS [۱]

مرتبط و یک صفحه قطب دارای تعدادی پیوند به صفحات مرجع می‌باشد؛ به عبارت دیگر یک صفحه با مرجع بالا به وسیله تعدادی از صفحات با قطب بالا اشاره شده است و همچنین یک صفحه با قطب بالا به تعدادی صفحه با مرجع بالا اشاره می‌کند. اگر فرض کنیم $(E, G=V)$ باشد الگوریتم HITS در دو مرحله زیر عمل می‌کند:

- نمونه‌گیری: در این مرحله مجموعه‌ای از صفحات مرتبط برای پرسش معین جمع‌آوری شده است.
- تکرار: در این مرحله مرجع و قطب با استفاده، از مرحله نمونه‌گیری پیدا می‌شوند.

$$H_p = \sum_{q \in I(p)} A_q, \quad A_p = \sum_{B(p)} H_q \quad (1)$$

که در رابطه (۱) H_q نمره قطب یک صفحه و A_q نمره یک صفحه می‌باشد. $I(p)$ مجموعه‌ای از صفحات مرجع p و $B(p)$ مجموعه‌ای از صفحات که مرجع به اشاره می‌کنند، می‌باشد. این روش مشکلاتی هم به همراه دارد که در زیر به چند مورد اشاره می‌کنیم:

- همه اسناد در همسایگی صفحات مرتبط از یک موضوع نبوده و در دسته‌های موضوعی مختلفی قرار دارند.
 - تشخیص میان قطب و مرجع آسان نیست، چون در HITS یک صفحه هم می‌تواند قطب خوب و همزمان مرجع خوبی باشد.
 - بعضی از وبگاه‌ها به خاطر بالا بردن رتبه مرتباً به یکدیگر اشاره می‌کنند؛ که به آن پهنش رتبه می‌گویند.
- الگوریتم Distance Rank مبتنی بر یادگیری تقویتی است که اولین بار ۲۰۰۸ در [۱۰] ارائه گردید. این الگوریتم

اگر صفحات اشاره‌کننده به آن فاصله کمی داشته باشند، فاصله این صفحه نیز کم خواهد بود. برای نشان دادن این مطلب در [۲۳] علاوه بر تعاریف ۱ و ۲ تعریف سومی را هم به شکل زیر ارائه کردند:

• اگر d_{ij} نشان‌دهنده فاصله لگاریتمی میان صفحه i و j باشد، میانگین فاصله صفحه j از بقیه صفحات با رابطه (۳) محاسبه می‌شود که N نشان‌دهنده تعداد صفحات وب است.

$$d_j = \frac{\sum_{i=1}^N d_{ij}}{N} \quad (3)$$

بعد از محاسبه بردار میانگین فاصله تمام صفحات، صفحات مطابق فاصله‌هایشان به صورت صعودی مرتب می‌شوند و لذا صفحات با فاصله کمتر از رتبه بیشتری برخوردار خواهند بود [۱۰] که جهت محاسبه میانگین فاصله برای هر صفحه، یک وابستگی میان فاصله آن صفحه و صفحاتی که به آن اشاره می‌کنند وجود دارد. با فرض این‌که فقط صفحه i به j اشاره می‌کند، بیدوکی برای محاسبه فاصله صفحه j (d_j) با استفاده از تعاریف ۱، ۲ و ۳ رابطه (۴) را ارائه کرده است. در رابطه (۴) چون N خیلی بزرگ است، از عبارت $\frac{d_{ij}}{N}$ چشم‌پوشی کردند.

$$d_j = \frac{\sum_{k=1}^N d_{kj}}{N} = \frac{\sum_{k \neq i}^N (d_{ki} + d_{ij}) + d_{ij}}{N} = \frac{\sum_{k \neq i}^N d_{ki}}{N} + d_{ij} = \quad (4)$$

$$\frac{\sum_{k=1}^N d_{ki} - d_{ii}}{N} + d_{ij} \xrightarrow{\text{Eq.}} d_j = d_i - \frac{d_{ii}}{N} + d_{ij} \approx d_i + d_{ij} = d_i + \log(O(i))$$

اگر $O(i)$ نشان‌دهنده تعداد پیوندهای خروجی از صفحه i و $B(j)$ نشان‌دهنده مجموعه صفحات اشاره‌کننده به j باشند، برای محاسبه مقدار Distance Rank رابطه (۵) را ارائه کرده‌اند.

$$d_j = \min_i (d_i + \log(O(i))) \quad i \in B(j) \quad (5)$$

و برای کامل‌تر شدن رابطه (۴)، آن را رابطه (۵) توسعه دادند.

$$d_{j,t+1} = (1 - \alpha) * d_{j,t} + \alpha * \min_j (\log(O(i)) + \gamma * d_{j,t}) \quad i \in B(j), 0 \leq \alpha \leq 1, 0 \leq \gamma \leq 1 \quad (6)$$

در رابطه (۶) α نرخ یادگیری و γ ضریب نزول می‌باشد

به $d_{j,t+1}$ فاصله صفحه j در زمان $t+1$ و $d_{j,t}$ به ترتیب نشان‌دهنده فاصله صفحه i و j در زمان t می‌باشند. ضریب نزول برای تنظیم اثرات صفحات قبلی که در مسیر رسیدن به j قرار دارند، استفاده می‌شود. در یک مسیر مانند، $j \rightarrow i \rightarrow k$ اثر روی با فاکتور γ اعمال می‌شود و نرخ یادگیری α با استفاده از رابطه (۷) محاسبه می‌شود که t زمان (شماره تکرار) است و ثابت β جهت تنظیم نرخ یادگیری استفاده می‌شود. در ابتدا به دلیل این‌که مقدار فاصله هیچ‌یک از صفحات در دسترس نیست، α مساوی یک بوده و با گذشت زمان به صورت نمایی و نزولی به صفر میل می‌کند.

$$\alpha = e^{-\beta * t} \quad (7)$$

رابطه (۷) نشان می‌دهد که فاصله در هر لحظه به فاصله خود و ورودی‌هایش در لحظه قبلی وابسته است. به عبارت دیگر عامل کاربر در هر لحظه انتخاب خود را بر اساس آگاهی‌های قبلی و وضعیت فعلی محیط انجام می‌دهد.

۴. روش پیشنهادی

از مسائل عمده رتبه‌بندی بر اساس پیوند که الگوریتم HITS یکی از این الگوریتم‌ها می‌باشد دارای مشکل غنی‌تر شدن اغنیا^{۱۶} می‌باشد. قرار گرفتن همیشه صفحات محبوب^{۱۷} در صدر لیست ارائه شده به کاربر، باعث می‌شود تا کاربر فقط صفحات خاصی را ببیند و در نتیجه صفحات تازه متولد شده با کیفیت بالا که کسی به آن‌ها اشاره نمی‌کند نتوانند در دید کاربران قرار گیرند. این مشکل باعث می‌شود صفحات محبوب مرتباً محبوب‌تر شده و تعداد پیوند به آن‌ها افزایش یابد. لذا موتورهای جستجو با ارائه نکردن عادلانه اطلاعات به کاربران، باعث صرف وقت زیاد و در نتیجه کندی تولید علم و دانش خواهند شد. مهم‌ترین چالش در موتورهای جستجو رتبه‌بندی صفحات در پاسخ به پرس‌وجوی کاربر می‌باشد. در حال حاضر روش‌های

16- Rich-get-richer
17- Popular

جدول ۲: نمونه نرمال شده قطب و مرجع HITS صفحات وب در ۲ تکرار در الگوریتم پیشنهادی

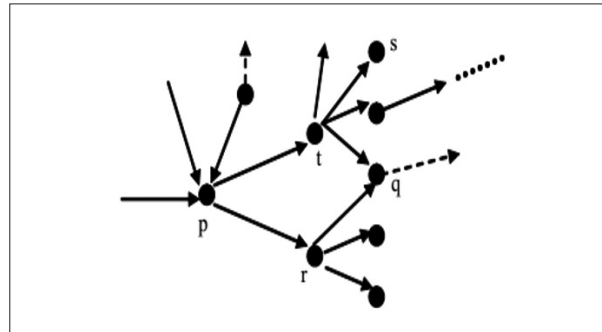
صفحه/وزن	H0	A0	H1	A1	H2	A2	NH2	NA2
A	۱	۱	۱	۵	۶	۲۵	۰,۱۹۳۵	۰,۳۵۷۱
B	۱	۱	۲	۶	۱۱	۳۱	۰,۳۵۴۸	۰,۴۴۲۹
C	۱	۱	۳	۳	۱۴	۱۴	۰,۴۵۱۶	۰,۲۰۰۰

که گفته شد الگوریتم HITS مشکل غنی تر شدن اغنیاء را دارد و باید نتایج این الگوریتم با الگوریتم Distance Rank که حساسیت کمتری به مشکل غنی تر شدن اغنیاء نشان می دهد ترکیب می شود. نحوه ترکیب این الگوریتم چنین خواهد بود میزان مرجع و قطب به دست آمده برای تمام صفحات نرمال می شود یعنی تمام مقادیر مرجع و قطب صفحات براساس میزان اعتبارشان نسبت مجموع کل صفحات مقداردهی می شود که در جدول (۲) نمونه نرمال شده برای سه صفحه نمایش داده شده است.

در جدول (۲) در ابتدا روش پیشنهادی نرمال کردن قطب و مرجع الگوریتم HITS را انجام می دهد و بعد از نرمال شدن نتایج برای ترکیب با نتایج الگوریتم DistanceRank آماده می شود و دلیل این نرمال کردن در روش پیشنهادی، ترکیب بهینه و درست براساس درصد مقادیر قطب و مرجع از مقادیر به دست آمده است با نتایج تولید شده در الگوریتم DistanceRank می باشد که برای این از رابطه جدید روش پیشنهادی به صورت رابطه (۸) ارائه می شود:

$$dj_{t+1} = (1-\alpha) * dj_t + \alpha * \min_j (\log(O(i)) + \gamma * dj_j) + 1 / (\text{Mean}(NH_t, NA_t)) \quad (8)$$

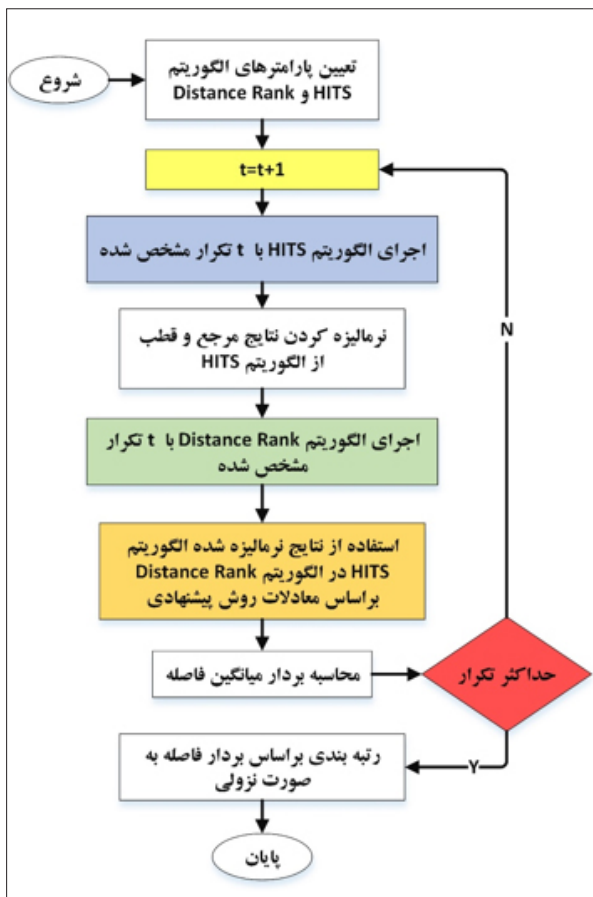
در رابطه (۸) در روش پیشنهادی ترکیبی از نتایج هر دو الگوریتم HITS و DistanceRank در رتبه بندی صفحات در روش پیشنهادی استفاده شده است. دلیل تقسیم در این رابطه به این دلیل می باشد که در الگوریتم DistanceRank بعد از محاسبه بردار میانگین فاصله تمام صفحات، صفحات مطابق فاصله هایشان به صورت صعودی مرتب می شوند و لذا صفحات با فاصله کمتر از رتبه بیشتری برخوردار خواهند بود که جهت محاسبه میانگین فاصله برای هر صفحه، یک وابستگی میان فاصله آن صفحه و



شکل ۳: الگوریتم Distance Rank به دست آوردن میانگین فاصله pq مساوی با $dp + \log 2 + \log 3$

رتبه بندی موجود مانند HITS از مشکلاتی مانند غنی تر شدن اغنیاء رنج می برند. در این بخش الگوریتمی ترکیبی از Distance Rank و HITS برای حل مشکلات فوق ارائه شده است. همان طور که بیان شد الگوریتم Distance Rank مبتنی بر یادگیری تقویتی است و از فاصله لگاریتمی میان صفحات به عنوان تویبخ دریافتی استفاده می کند. هدف آن کمینه کردن کل تویبخ های دریافتی می باشد. در این روش یک کاربر به صورت موج سوار واقعی مدل می شود. بدین معنی که در ابتدا کاربر آگاهی کمی از سیستم دارد و با گذشت زمان آگاهی او بیشتر شده و عمل پیمایش را بهتر انجام می دهد. الگوریتم Distance Rank در رتبه بندی نسبت به بقیه مؤثرتر می باشد و حساسیت کمتری به مشکل غنی تر شدن اغنیاء نشان می دهد. شکل (۳) نمونه صفحه جدید q و محاسبه میانگین فاصله توسط الگوریتم Distance Rank را نشان می دهد.

بنابراین در این بخش روش پیشنهادی ترکیبی از الگوریتم HITS و Distance Rank ارائه خواهیم داد. در روش پیشنهادی ابتدا الگوریتم HITS اجرا شود و بعد از اجرای الگوریتم HITS مقادیر صفحات وب را در دو فرم مرجع و قطب به دست خواهد آمد و همان طور



شکل ۴: روندنمای الگوریتم پیشنهادی

۶ گیگابایت حافظه موقت شبیه‌سازی شده است. برای ارزیابی روش پیشنهادی از سه مجموعه داده که اولین مجموعه مربوط گراف استاندارد که در [۱۱] معرفی شده است این مجموعه داده به صورت ماتریس و مجموعه پیوند در این گراف وارد سی‌شارپ شده است. دومین مجموعه داده مربوط به یک گراف تصادفی بزرگ است که این گراف به صورت تصادفی در محیط سی‌شارپ تولید شده است و لازم به ذکر است، داده‌های مورداستفاده در این مجموعه برای شبیه‌سازی رتبه‌بندی صفحات وب در این مقاله به صورت سیستم تصادفی تولید می‌شود و اعتبار و صحت داده‌های ورودی با توابع توزیع محیط سی‌شارپ بررسی می‌شود تا این‌که یک گراف تصادفی بزرگ مناسب برای رتبه‌بندی تولید شود و بتوانیم روش پیشنهادی بر روی یک گراف وب بزرگ مورد ارزیابی قرار دهیم. آخرین مجموعه داده مربوط صفحات دانشگاه آزاد اسلامی واحد

صفحاتی که به آن اشاره می‌کنند وجود دارد.

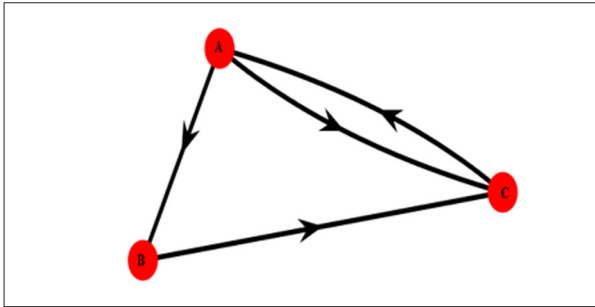
بنابراین هر چه صفحات در الگوریتم HITS دارای مرجع و قطب بالاتری باشند با تقسیم در بردار محاسبه مقادیر کمتری به دست خواهند آورد و صفحات رتبه‌های بالاتری کسب می‌کنند. یکی نکاتی در روش پیشنهادی وجود دارد مقادیر t یا تعداد تکرار می‌باشد که تعداد تکرار در هر الگوریتم یکسان در نظر گرفته شده است، بنابراین الگوریتم Distance Rank در تکرار $t+1$ از نتایج حاصل از الگوریتم HITS استفاده می‌کند. همچنین در رابطه (۸) α نرخ یادگیری و γ ضریب نزول می‌باشد به $d_{j,t+1}$ فاصله صفحه j در زمان $t+1$ و $d_{i,t}$ به ترتیب نشان‌دهنده فاصله صفحه و در زمان می‌باشند. ضریب نزول برای تنظیم اثرات صفحات قبلی که در مسیر رسیدن به z قرار دارند، استفاده می‌شود. در یک مسیر مانند، $z \rightarrow i \rightarrow k \rightarrow a$ اثر a روی z با فاکتور γ اعمال می‌شود و نرخ یادگیری α با استفاده از رابطه (۷) محاسبه می‌شود که t زمان (شماره تکرار) است و ثابت β جهت تنظیم نرخ یادگیری استفاده می‌شود. در ابتدا به دلیل این‌که مقدار فاصله هیچ‌یک از صفحات در دسترس نیست، α مساوی یک بوده و با گذشت زمان به صورت نمایی و نزولی به صفر میل می‌کند. روندنمای روش پیشنهادی به صورت شکل (۴) می‌باشد.

۵. ارزیابی و نتایج

در این بخش به ارزیابی و نتایج روش پیشنهادی خواهیم پرداخت. در این بخش ابتدا محیط شبیه‌سازی و مجموعه داده را معرفی خواهیم کرد. سپس به ارزیابی روش پیشنهادی خواهیم پرداخت و نهایتاً به مقایسه و ارزیابی روش پیشنهادی با مدل‌های دیگر خواهیم پرداخت.

۵.۱. محیط شبیه‌سازی و مجموعه داده

روش پیشنهادی که ترکیب الگوریتم HITS با الگوریتم Distance Rank برای بهبود نتایج در موتورهای جستجو، در محیط سی‌شارپ و بر روی سیستمی با پردازنده پنج هسته‌ای اینتل با قدرت پردازشی $2/30$ گیگاهرتز و



شکل ۵: گراف استاندارد و نحوه ارتباط بین صفحات

پیشنهادی و سایر الگوریتم‌ها و نحوه رتبه‌بندی صفحات را نشان می‌دهد؛ که گراف استاندارد به‌عنوان ورودی در نظر گرفته شده است و نحوه ارتباط صفحات به‌صورت یک ماتریس صفر و یک ذکر شده و صفحات موردنظر را رتبه‌بندی کرده و نمایش می‌دهد.

همان‌طور که در جدول (۳) مشاهده می‌شود. در این جدول نتایج مقایسه رتبه‌بندی صفحات روش پیشنهادی با سایر الگوریتم‌ها بر روی گراف استاندارد را نشان می‌دهد که نتایج حاصل نشان می‌دهد هر دو صفحه A و B در الگوریتم WPR و Pagerank یکسان هستند. صفحه A چون قطب مناسبی است در الگوریتم Hits بالاترین مقدار قطب را گرفته است و صفحه C چون یک مرجع مناسب هست بیشترین مقدار مرجع است. الگوریتم Distance بر اساس رتبه‌بندی که انجام داده کمترین مقدار را به صفحه A اختصاص داده است. روش پیشنهادی هم چون ترکیبی از روش Hits و Distance سعی کرده یک حد میانه بین این دو الگوریتم باشد و صفحه A در رتبه‌بندی اول قرار داده‌است و توانسته از الگوریتم Distance برای بهبود نتایج الگوریتم HITS در این مجموعه داده استفاده کند. البته روش پیشنهادی توانسته از بین دو صفحه C و B یک رتبه‌بندی مناسب (صفحه C یک مرجع مناسب است) انجام دهد چراکه صفحه C نیز با رتبه‌بندی بهتر نسبت به صفحه B در رتبه بعد از صفحه A قرار گرفته است.

همچنین در این بخش ما از مجموعه داده‌های گراف تصادفی بزرگ که در محیط سی‌شارپ تولید شده است، استفاده می‌کنیم؛ بنابراین به منظور ارزیابی روش

ارومیه است که برای نمایش ارتباط بین این صفحات و درک بهتر پیوندهای ورودی و خروجی بین این صفحات، گراف آن را با توجه به پیوندهای بین این صفحات رسم کردیم.

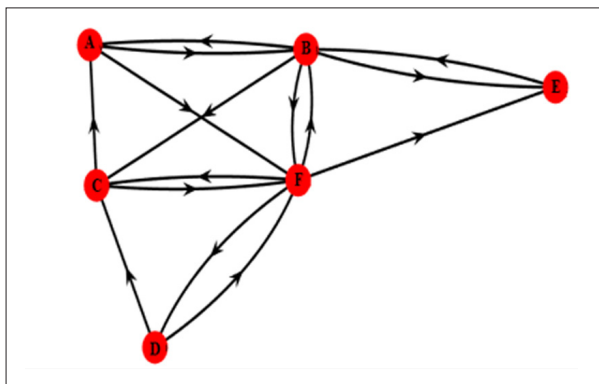
۲.۵. ارزیابی روش پیشنهادی

در این بخش، برای ارزیابی روش پیشنهادی از سه مجموعه داده که اولین مجموعه مربوط گراف استاندارد که در مقاله [۱۱] معرفی شده است این مجموعه داده به‌صورت ماتریس و مجموعه پیوند در این گراف وارد سی‌شارپ شده است. دومین مجموعه داده مربوط یک گراف تصادفی بزرگ است که این گراف به‌صورت تصادفی در محیط سی‌شارپ تولید شده است و لازم به ذکر است، داده‌های مورد استفاده در این مجموعه برای شبیه‌سازی رتبه‌بندی صفحات وب در این فصل به‌صورت سیستم تصادفی تولید می‌شود و اعتبار و صحت داده‌های ورودی با توابع توزیع محیط سی‌شارپ بررسی می‌شود تا این‌که یک گراف تصادفی بزرگ مناسب برای رتبه‌بندی تولید شود و بتوانیم روش پیشنهادی بر روی یک گراف وب بزرگ مورد ارزیابی قرار دهیم. آخرین مجموعه داده مربوط صفحات دانشگاه آزاد اسلامی ارومیه است که برای نمایش ارتباط بین این صفحات و درک بهتر پیوندهای ورودی و خروجی بین این صفحات؛ گراف آن را با توجه به پیوندهای بین این صفحات رسم شده است. بنابراین به منظور ارزیابی روش پیشنهادی ما صفحه وب گراف استاندارد را که در شکل (۳) آورده شده است را با روش پیشنهادی مورد بررسی قرار دادیم. برای نمایش ارتباط بین این صفحات و درک بهتر پیوندهای ورودی و خروجی بین این صفحات، گراف آن را با توجه به پیوندهای بین این صفحات به‌صورت شکل (۵) رسم کردیم.

بعد از حاصل شدن داده‌های و صفحات مطابق شکل (۵) شکل گراف استاندارد، صفحات وب این مجموعه داده را می‌توان با الگوریتم HITS و Distance Rank و الگوریتم ترکیبی پیشنهادی رتبه‌بندی کرد. جدول (۳) خروجی روش

جدول ۳: مقایسه رتبه‌بندی صفحات روش پیشنهادی با سایر الگوریتم‌ها بر روی گراف استاندارد

رتبه/صفحه	Pagerank[۱۱]	WPR [۱۱]	HITS		Distance	روش پیشنهادی
			قطب	مرجع		
A	۰,۵۹	۰,۵۹	۰,۵۷۴	۰,۱۱۱	۲,۵	۲,۷
B	۰,۴۲	۰,۴۲	۰,۳۵۷۱	۰,۳۳۳۳	۳,۱۹	۳,۳۷۱۸
C	۰,۵۹	۰,۷۶	۰,۰۷۱۴	۰,۵۵۵۶	۲,۸۵	۳,۰۳۱۸



شکل ۶: گراف تصادفی بزرگ و نحوه ارتباط بین صفحات

(۵) خروجی روش پیشنهادی نحوه رتبه‌بندی صفحات را نشان می‌دهد؛ که گراف تصادفی بزرگ به‌عنوان ورودی در نظر گرفته شده است و نحوه ارتباط صفحات به‌صورت یک ماتریس صفر و یک ذکر شده و صفحات موردنظر را رتبه‌بندی کرده و نمایش می‌دهد.

جدول (۶) مقایسه رتبه‌بندی صفحات روش پیشنهادی با سایر الگوریتم‌ها بر روی گراف تصادفی بزرگ را در حالت اجرا در محیط سی‌شارپ نشان می‌دهد. همان‌طور که در جدول (۶) مشاهده می‌شود. در این جدول نتایج مقایسه رتبه‌بندی صفحات روش پیشنهادی با سایر الگوریتم‌ها بر روی گراف تصادفی و نتایج حاصل را نشان می‌دهد. الگوریتم پیشنهاد در رتبه‌بندی صفحات، صفحه A را که هم قطب و هم مرجع خوبی دارد را به‌عنوان صفحه اول انتخاب کرده است و صفحه B را که از لحاظ هم مرجع دارای کمترین مقدار است به‌عنوان آخرین صفحه رتبه‌بندی قرار داده است. در نتیجه روش پیشنهادی در رتبه‌بندی دیگر صفحات مانند F از نظر الگوریتم HITS یک مرجع نسبتاً قوی و قطب ضعیف و از نظر الگوریتم

جدول ۴: تولید ماتریس تصادفی از صفر و یک برای ایجاد گراف تصادفی بزرگ در محیط اجرا

صفحه	A	B	C	D	E	F
A	۰	۱	۰	۰	۰	۱
B	۱	۰	۱	۰	۱	۱
C	۱	۰	۰	۰	۰	۱
D	۰	۰	۱	۰	۰	۱
E	۰	۱	۰	۰	۰	۰
F	۰	۱	۱	۱	۱	۰

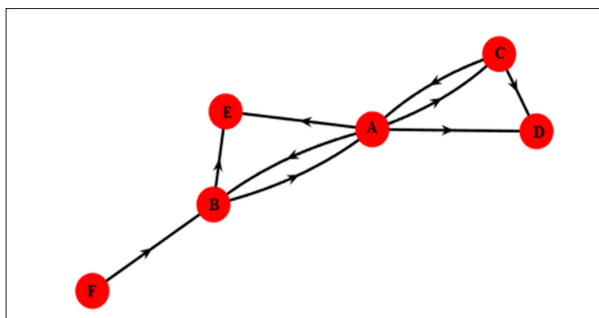
پیشنهادی ما صفحه وب گراف تصادفی بزرگ را که در جدول (۴) در محیط سی‌شارپ تولید شده است. ماتریس تصادفی در این نمونه شامل شش صفحه وب هست که بانام‌های مستعار نام‌گذاری شده است که اگر صفحه وب به صفحه وب دیگر اشاره داشته باشد مقدار آن برابر با یک و در غیر این صورت مساوی صفر خواهد بود. بنابراین این ماتریس به‌صورت تصادفی در محیط نرم‌افزار سی‌شارپ برای شش صفحه وب تولید شده که آن را در ادامه با روش پیشنهادی موردبررسی و ارزیابی قرار دادیم.

در جدول (۴) پیوندهای ورودی و خروجی صفحات وب به‌صورت تصادفی ایجاد شده است که برای نمایش ارتباط بین این صفحات و درک بهتر پیوندهای ورودی و خروجی بین این صفحات، گراف آن را با توجه به پیوندهای بین این صفحات به‌صورت شکل (۶) رسم کردیم.

بعد از حاصل شدن داده‌ها و ارتباط بین صفحات مطابق جدول (۴) و شکل (۶) گراف تصادفی بزرگ، صفحات وب این مجموعه داده را می‌توان با الگوریتم HITS و Distance Rank و الگوریتم ترکیبی پیشنهادی رتبه‌بندی کرد. جدول

جدول ۷: صفحات بررسی شده با روش پیشنهادی

عنوان صفحه	آدرس صفحه
دانشگاه آزاد اسلامی واحد ارومیه	http://www.iaurmia.ac.ir
دانشگاه آزاد اسلامی	http://www.iau.ac.ir
سیستم ثبت نام و اطلاع رسانی دانشگاه آزاد	http://amozesh.iaurmia.ac.ir/login.aspx
سامانه اتوماسیون پایان نامه های دانشجویان	http://thesis.iaurmia.ac.ir
سامانه نقل و انتقالات دانشجویان	http://transfers.stu...iau.ir
درگاه دانشجویی صندوق رفاه	http://bp.swf.ir



شکل ۷: گراف دانشگاه آزاد اسلامی ارومیه و نحوه ارتباط بین صفحات

و دانشگاه آزاد اسلامی و سیستم ثبت نام و اطلاع رسانی دانشگاه آزاد و سامانه اتوماسیون پایان نامه های دانشجویان و سامانه نقل و انتقالات دانشجویان و درگاه دانشجویی صندوق رفاه، برای نمایش ارتباط بین این صفحات و درک بهتر پیوندهای ورودی و خروجی بین این صفحات، گراف آن را با توجه به پیوندهای بین این صفحات به صورت شکل (۷) رسم کردیم.

جدول (۷) رتبه بندی روش پیشنهادی بر روی گراف دانشگاه آزاد اسلامی ارومیه را نشان می دهد. خروجی روش پیشنهادی نحوه رتبه بندی صفحات را نشان می دهد؛ که گراف دانشگاه آزاد اسلامی ارومیه به عنوان ورودی در نظر گرفته شده است و نحوه ارتباط صفحات به صورت

جدول ۵: رتبه بندی صفحات روش پیشنهادی بر روی گراف تصادفی بزرگ

تکرار ۵	تکرار ۴	تکرار ۳	تکرار ۲	تکرار ۱	صفحه/اجرا
۴,۳۸۴۵	۴,۶۸۷۹	۵,۲۲۵۶	۶,۲۸۴۵	۸,۲۱۸۸	A
۱۵,۰۲۰۶	۱۵,۲۶۲	۱۵,۷۳۳۴	۱۶,۳۲۶۵	۱۵,۰۳۱۴	B
۴,۷۵۸۱	۵,۰۲۹	۵,۵۶۴۴	۶,۵۶۸۵	۸,۸۳۶۶	C
۵,۷۶۱۶	۶,۰۳۲۴	۶,۵۷۰۴	۷,۵۷۸	۹,۸۵۲۳	D
۸,۲۱۶۹	۸,۴۸۶۹	۸,۹۷۱۳	۱۰,۰۳۶۴	۱۱,۸۷۲	E
۵,۶۰۸۳	۵,۸۷۸۳	۶,۳۶۶۵	۷,۴۵۵۱	۹,۸۵۵۲	F

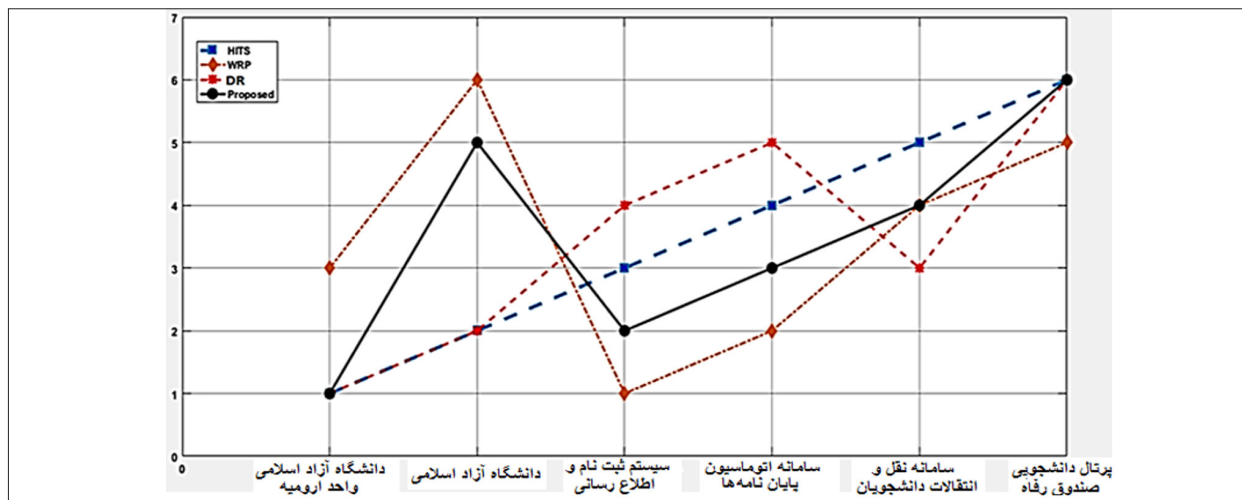
جدول ۶: مقایسه رتبه بندی صفحات روش پیشنهادی با سایر الگوریتم ها بر روی گراف تصادفی بزرگ

روش پیشنهادی	HITS		Distance	صفحه/رتبه
	مرجع	قطب		
۴,۴۳۸۵	۰,۲۷۱	۰,۲۳۷۹	۱,۸۸	A
۱۵,۰۲۰۶	۰,۰۲۲۹	۰,۱۰۶۴	۱,۸۸	B
۴,۷۵۸۱	۰,۲۳۶۵	۰,۱۹۴۱	۱,۷۶	C
۵,۷۶۱۶	۰,۱۷۷۶	۰,۱۴۴۶	۱,۷۶	D
۸,۲۱۶۹	۰,۰۴۳۹	۰,۲۳۹۷	۲,۱۱	E
۵,۶۰۸۳	۰,۲۴۸۱	۰,۰۷۵۵	۲,۱۱	F

Distance همین صفحه دارای رتبه متوسطی هست را در رتبه دوم قرار داده است و توانسته است یک حد میانه بین دو الگوریتم برقرار کند و بیشتر ضعف رتبه بندی الگوریتم Hits را بهبود دهد.

در این مقاله، ما از مجموعه داده های دانشگاه آزاد اسلامی ارومیه استفاده می کنیم. بنابراین به منظور ارزیابی روش پیشنهادی ما صفحه وب مربوط دانشگاه آزاد ارومیه را که در جدول (۷) آورده شده است را با روش پیشنهادی مورد بررسی قرار دادیم. در این جدول برای هر یک صفحات یک نام مستعار انتخاب کردیم و همچنین عنوان و نشانی دقیق هر صفحه را در مقابل هر اسم مستعار مشخص کردیم.

همچنین با توجه اطلاعات جدول (۷) و ارتباط مشخص شده بین صفحات دانشگاه آزاد اسلامی ارومیه



شکل ۸: مقایسه روش پیشنهادی رتبه‌بندی با پرس‌وجوی «دانشگاه آزاد اسلامی واحد ارومیه»

آن صفحه را به‌عنوان آخرین صفحه رتبه‌بندی کرده است.

۳.۵. مقایسه و ارزیابی

برای ارزیابی و مقایسه الگوریتم‌های پیشنهادی از معیارهای ارزیابی $P@n$, AP و NDC که توسط کالرو جارولین^{۱۸} و جانا ککالنن^{۱۹} در [۳۹] ارائه گردیده است استفاده کردیم. معیار $P@n$ در رابطه (۹) نشانگر نسبت تعداد اسناد مرتبط در n سند نخست ارائه‌شده به n است. هدف اصلی از این معیار آن است که دقت سیستم از دید کاربران محاسبه شود.

$$P@n = \frac{\text{#of relevant docs in top } n \text{ results}}{n} \quad (9)$$

معیار $P@n$ از دقت کافی برخوردار نیست. چراکه در این معیار فقط مرتبط بودن یا نبودن یک سند در نظر گرفته می‌شود. رابطه (۱۰) معیار $NDCG$ را نشان می‌دهد. در این معیار که برای n نتیجه نخست استفاده می‌شود، r_j نشان‌دهنده درجه مرتبط بودن سند j با پرس‌وجوی مربوطه می‌باشد.

$$NDCG@n = \frac{\sum_{j=1}^n \frac{r_j - 1}{\log(1+j)}}{\sum_{j=1}^n \frac{1}{\log(1+j)}} \quad (10)$$

AP یا میانگین دقت که برای هر پرس‌وجو محاسبه می‌شود، مقدار آن با میانگین $P@n$ برای تمام اسناد مرتبط با پرس‌وجوی موردنظر برابر می‌باشد. در رابطه (۱۱) اگر i -امین سند مرتبط با پرس‌وجو باشد، $rel(i)$ برابر ۱ و در

جدول ۸: رتبه‌بندی صفحات روش پیشنهادی بر روی گراف دانشگاه آزاد اسلامی ارومیه

صفحه/اجرا	تکرار ۱	تکرار ۲	تکرار ۳	تکرار ۴	تکرار ۵
A	۸,۰۲۵۹	۵,۷۰۶	۵,۳۸۴	۴,۷۷۵۷	۴,۶۹۵۲
B	۷,۹۴۲۵	۷,۴۸۲۶	۶,۲۳۴۷	۶,۰۲۸۳	۵,۷۰۸۴
C	۱۰,۷۶۱۶	۸,۰۱۸۶	۶,۶۹۴	۶,۴۸۸۲	۶,۱۷۳۴
D	۱۰,۱۹۰۵	۸,۰۷۲۴	۶,۹۶۸	۶,۳۳۹۷	۶,۲۵۱۵
E	۱۰,۱۹۰۵	۷,۳۷۲۴	۶,۹۶۸	۶,۳۳۹۷	۶,۲۵۱۵
F	۴۵,۰۰۶۸	۴۶,۸۸۵۷	۴۶,۵۵۲۹	۴۶,۲۶۴۸	۴۶,۱۳۲۴

یک ماتریس صفر و یک ذکرشده و صفحات موردنظر را رتبه‌بندی کرده و نمایش می‌دهد.

شکل (۸) نتایج به‌دست‌آمده در [۲۱] که با پرس‌وجوی دانشگاه ارومیه و نتایج رتبه‌بندی روش پیشنهادی و الگوریتم‌های WPR , $Distance Rank$ [۲۱] و $HITS$ را با پرس‌وجوی دانشگاه آزاد اسلامی واحد ارومیه مقایسه می‌کند. با توجه به شکل (۸) می‌توان به این نتیجه رسید که روش پیشنهادی با توجه به این که صفحه دانشگاه آزاد اسلامی واحد ارومیه در هر دو الگوریتم $Distance Rank$ و $HITS$ دارای رتبه بالای می‌باشد، این الگوریتم هم همانند دو الگوریتم دیگر عمل می‌کند و صفحه A را به‌عنوان رتبه یک در نظر می‌گیرد و همچنین صفحه F را که دارای کمترین رتبه در هر دو الگوریتم $Distance Rank$ و $HITS$ می‌باشد،

18- Kalervo Jarvelin
19- Jaana Kekalanien

جدول ۱۰: مقایسه روش پیشنهادی با معیارهای NDCG@n و AP

	NDCG@n	AP
روش پیشنهادی	۸,۱	۱
Distance Rank [10]	۶,۰۵	۰,۶۹۹
User attention time [23]	۶,۶۸۱	۰,۹۱۶
Page Rank [24]	۶,۰۵	۰,۶۹۹
Wighted Page Rank [25]	۶,۲۷۹	۰,۸۰۵
HITS [1]	۶,۴۳۱	۰,۸۰۵

داد که روش پیشنهادی در مقایسه با الگوریتم‌های دیگر عملکرد بهتری دارد و توانسته است رتبه‌بندی متفاوت و بهتری نسبت HITS و سایر الگوریتم‌های رتبه‌بندی مانند Distance Rank و PR و WPR داشته باشد. همچنین برای ارزیابی و مقایسه الگوریتم‌های پیشنهادی از معیارهای ارزیابی AP، P@n و NDC استفاده کردیم. می‌توان نتیجه گرفت که با ترکیب یک الگوریتم HITS با الگوریتم Distance Rank برای بهبود نتایج در موتورهای جستجو می‌توان یک الگوریتم قدرتمند برای رتبه‌بندی نتایج موتور جستجو داشت. که نتایج موتورهای جستجو را مطابق با حد میانه میان دو الگوریتم HITS با Distance Rank رتبه‌بندی می‌کند و تعداد خطاهای رتبه‌بندی را تا جای ممکن کاهش می‌دهد.

مراجع

1. Kleinberg, Jon M. «Authoritative sources in a hyper-linked environment.» ACM-SIAM Symposium on Discrete Algorithms. 1998.
2. Gomathi, C., M. Moorthi, and K. Duraiswamy. "Web access pattern algorithms in education domain." Computer and information science 1, no. 4 2008.
3. Joy Shalom Sona, Prof. Asha Ambhaikar "A Reconciling Website System to Enhance Efficiency with Web Mining Techniques: International Journal Of Scientific & Engineering Research Volume 3, 1 ISSN 2229-5518, Issue 2, February-2012.
4. Lin, Chun-Wei, and Tzung-Pei Hong. "A survey of fuzzy web mining." Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery 3, no. 3, 2013.
5. Suthar, Parth, and Bhavesh Oza. "A survey of web usage mining techniques." Int. J. Comput. Sci. Inf. Technol. (IJCSIT) 6, no. 6, 2015.
6. Herrouz, Abdelhakim, Chabane Khentout, and Mahied-

جدول ۹: مقایسه الگوریتم با معیار P@n

	n=2	n=3	n=4
روش پیشنهادی	۱	۱	۰,۸۲
Distance Rank [10]	۰,۵	۰,۳۳۳	۰,۵
User attention time [23]	۱	۰,۶۶۶	۰,۷۵
Page Rank [24]	۰,۵	۰,۳۳۳	۰,۵
Wighted Page Rank [25]	۰,۵	۰,۶۶۶	۰,۵
HITS [1]	۰,۵	۰,۶۶۶	۰,۷۵

غیر این صورت صفر خواهد بود و N نشان‌دهنده تعداد نتایج ارائه‌شده برای یک پرس‌وجو می‌باشد.

$$AP = \frac{\sum_{i=1}^N P@i \cdot rel(i)}{\#total\ relevant\ docs\ for\ this\ query} \quad (11)$$

با استفاده از معیارهای اشاره‌شده، روش پیشنهادی را با الگوریتم‌های Distance Rank، Weighted Page Rank، HITS، Rank page و الگوریتم مبتنی بر زمان توجه کاربران ارائه شده در [۶] مقایسه کردیم و که نتایج آن با معیار P@n که نمایانگر دقت الگوریتم در n نتیجه نخست می‌باشد به صورت جدول (۹) حاصل گردیده است.

با توجه به جدول (۹) دقت روش پیشنهادی در هر سه حالت (دقت در دو نتیجه نخست، دقت در سه نتیجه نخست و دقت در چهار نتیجه نخست) مناسب بوده است و با معیارهای NDCG و AP نتایج ارزیابی به صورت جدول (۱۰) حاصل گردیده است.

در جدول (۱۰) که نمایانگر میانگین دقت و درجه مرتبط بودن نتایج هر یک از الگوریتم‌ها می‌باشد، روش پیشنهادی از میانگین دقت (AP) بالایی نسبت به سایر الگوریتم‌ها عملکرد بهتری از خود به نمایش گذاشته است.

۶. نتیجه‌گیری

روش پیشنهادی برای ترکیب الگوریتم HITS با الگوریتم Distance Rank برای بهبود نتایج در موتورهای جستجو شبیه‌سازی شده است. نتایج مقایسه الگوریتم با سایر الگوریتم‌ها بر روی سه مجموعه از داده‌ها نشان

- pp. 27-32. IEEE, 2017.
21. Lakshmi, L., P. Bhaskara Reddy, and C. Shoba Bindu. "Dynamic Navigation of Query Results Based on Hash-Based Indexing Using Improved Distance PageRank Algorithm." In *Artificial Intelligence and Evolutionary Computations in Engineering Systems*, pp. 213-221. Springer, Singapore, 2018.
 22. Sethi, Shilpa, and Ashutosh Dixit. "A novel page ranking mechanism based on user browsing Patterns." In *Software Engineering*, pp. 37-49. Springer, Singapore, 2019.
 23. Xu, Songhua, Yi Zhu, Hao Jiang, and Francis CM Lau. "A User-Oriented Webpage Ranking Algorithm Based on User Attention Time." In *AAAI*, vol. 8, pp. 1255-1260. 2008.
 24. Page, Lawrence, Sergey Brin, Rajeev Motwani, and Terry Winograd. "The PageRank citation ranking: Bringing order to the web". Stanford InfoLab, 1999.
 25. Xing, Wenpu, and Ali Ghorbani. "Weighted pagerank algorithm." In *Communication Networks and Services Research*, 2004. Proceedings. Second Annual Conference on, pp. 305-314. IEEE, 2004.
 - dine Djoudi. "Overview of web content mining tools." arXiv preprint arXiv: 1307.1024, 2013.
 7. Du, Ranran, and Yongbin Lin. "Comment ranking by search engine." U.S. Patent Application 10/242,105, filed March 26, 2019.
 8. Ananthi, J. "A survey web content mining methods and applications for information extraction from online shopping sites." *International Journal of Computer Science and Information Technologies (IJCSIT)* 5, no. 3, 2014.
 9. Jia, Chen, Yunyan Du, Siying Wang, Tianyang Bai, and Teng Fei. "Measuring the vibrancy of urban neighborhoods using mobile phone data with an improved PageRank algorithm." *Transactions in GIS* 23, no. 2 (2019): 241-258.
 10. Bidoki, Ali Mohammad Zareh, and Nasser Yazdani. "DistanceRank: An intelligent ranking algorithm for web pages." *Information Processing & Management* 44, no. 2, 2008.
 ۱۱. فتوحی، رضا، و روح اله عبدی پور، «ارزیابی کارایی الگوریتم‌های رتبه‌بندی صفحه برای استخراج صفحات وب»، همایش ملی مهندسی کامپیوتر و توسعه پایدار با محوریت شبکه‌های کامپیوتری، مدل‌سازی و امنیت سیستم‌ها، مشهد، موسسه آموزش عالی خاوران، ۱۳۹۲.
 ۱۲. زارع بیدکی، علی محمد، آزادنیبا محمد، «الگوریتم ترکیبی وفقی جهت رتبه‌بندی صفحات وب»، کنفرانس ملی سالانه انجمن کامپیوتر ایران، ۱۳۸۶.
 ۱۳. فرصتی، رعنا، و محمدرضا میبیدی، «الگوریتمی مبتنی بر ساختار پیوندی صفحات و اطلاعات استفاده کاربران برای پیشنهاد صفحات وب»، دومین کنفرانس داده‌کاوی ایران، تهران، دانشگاه صنعتی امیرکبیر، موسسه پژوهشی داده‌پردازان گیتا، ۱۳۸۷.
 ۱۴. قاشللو، نگین، و فرهاد مردوخی، «ارائه یک الگوریتم یادگیر مبتنی بر خصیصه‌های ذاتی جهت رتبه‌بندی صفحات وب»، دومین کنفرانس ملی فناوری، انرژی و داده با رویکرد مهندسی برق و کامپیوتر، کرمانشاه، انجمن IEEE شاخه دانشجویی کردستان، ۱۳۹۵.
 ۱۵. کلانتری، میترا، «یک الگوریتم رتبه‌بندی صفحات وب مبتنی بر رویکرد یادگیری»، دومین همایش ملی مهندسی رایانه و مدیریت فناوری اطلاعات، تهران، گروه پژوهشی بوعلی، ۱۳۹۴.
 16. Chawla, Suruchi. «A novel approach of cluster based optimal ranking of clicked URLs using genetic algorithm for effective personalized web search.» *Applied Soft Computing* 46, 2016.
 17. Gupta, Daya, and Devika Singh. "User preference based page ranking algorithm." In *Computing, Communication and Automation (ICCCA)*, 2016 International Conference on, pp. 166-171. IEEE, 2016.
 18. Sote, A. M., and SR Pande. "Application of Page Ranking Algorithm in Web Mining." In *International Conference on Advances in Engineering & Technology-2014*. 2014.
 19. Rani, S. Geetha, and M. Sorana Mageswari. "A link-click-concept based ranking algorithm for ranking search results." *Indian Journal of Science and Technology* 7, no. 10, 2014.
 20. Sen, Tuhena, Dev Kumar Chaudhary, and Tanupriya Choudhury. "Modified Page Rank Algorithm: Efficient Version of Simple Page Rank with Time, Navigation and Synonym Factor." In *Computational Intelligence and Networks (CINE)*, 2017 3rd International Conference on,