

## روشی جهت بهینه سازی نتایج موتورهای جستجوگر با بهره مندی از الگوریتم های خوشه بندی و بردار پشتیبان

محمداحمدی

دانشجوی کارشناسی ارشد- کامپیوتر نرم افزار- موسسه آموزش عالی غیرانتفاعی- غیردولتی مقدس اردبیلی

عباس میرزایی

گروه مهندسی کامپیوتر- واحد اردبیل- دانشگاه آزاد اسلامی- اردبیل- ایران

### چکیده

امروزه دانش و اطلاعات وب در حال افزایش است، همچنین تعداد وبسایت ها نیز روز به روز بیشتر می شود؛ در نتیجه نیاز است شرایطی فراهم گردد تا در بحث هایی مثل تجارت الکترونیک، موتورهای جستجو، تبلیغات هدفمند و غیره بر اساس رویکرد کاربران مبتدی و حرفه ای از بین صفحات مختلف پیشنهاد مناسبی به آنها داده شود. اما این که بتوان نیازها و علایق کاربران مختلف را پیش بینی نموده و صفحات مورد علاقه آنان را پیشنهاد کرد با پیچیدگی هایی همراه است و باید بتوان با روش هایی بهینه و با استفاده از الگوهای حرکتی کاربران که در سرورهای وب ثبت گردیده است این کار را انجام داد. روش های مختلفی برای پیش بینی صفحات وب ارائه شده است. در این تحقیق روشی جهت بهبود عملکرد موتورهای جستجو با استفاده از خوشه بندی و بردار پشتیبان ارائه شده است که از چند گام مختلف اعم از پیش پردازش داده، ایجاد پروفایل ثبت کاربر، استخراج میزان علاقه مندی کاربر، خوشه بندی و بردار پشتیبان تشکیل شده است. ارائه روش پیشنهادی علاوه بر کاهش میزان حجم اطلاعات، پیچیدگی سیستم را کم نموده است، همچنین در مرحله پیش بینی، با ایجاد تغییراتی در مدل بردار پشتیبان میزان دقت سیستم نیز نسبت به مدل پایه افزایش یافته است. مدل بردار پشتیبان و روش خوشه بندی استفاده شده در این سیستم به صورت بهینه عمل کرده تا نتایج قابل قبولی از سیستم استخراج گردد. معیارهای ارزیابی در این تحقیق دقت، فراخوانی و پوشش می باشد. میزان دقت کلی سیستم برابر با  $83/5$  شده است که نسبت به کارهای مشابه  $1/8$  بهبود حاصل شده است. نتایج ارزیابی این سیستم به کمک ابزار متلب با دیگر روش های مشابه بیانگر عملکرد بهتر و بهینه است.

**واژگان کلیدی:** وب کاوی، خوشه بندی، الگوهای رفتاری کاربران، بردار پشتیبان

## مقدمه

موتورهای جستجو امروز همچون قلب تپنده اینترنت عمل می کنند و عملاً با افزایش سرسام آور تعداد صفحات موجود در وب، تنها راه کاوش و یافتن اطلاعات در این شبکه، استفاده از موتورهای جستجو می باشد. امروزه کمپانی هایی که موتورهای جستجوی خود را عرضه کرده اند رقابتی سنگین و فشرده را در عرصه عمومی آغاز کرده اند که نمونه آن رقابت شرکت های گوگل و یاهو می باشد که بتازگی مایکروسافت نیز موجودیت خود را در این عرصه شروع کرده است و قصد سرمایه گذاری وسیع در این حوزه را دارد. شرکت گوگل با الگوریتم انقلابی خود که البته پیشینه اش در کارهای موتورهای قبلی به چشم می خورد در حال حاضر محبوب ترین موتور جستجو می باشد و سهام آن ارزشی افسانه ای یافته است و این موفقیت مدیون نگرش جدید مؤسسين این موتور به ماهیت مسئله می باشد. موتورهای سنتی با الگوریتم هایی که مبتنی بر شمارش کلمات کلیدی است پیش می روند و این نوع نگرش به تنهایی جوابگوی نیازهای ذهنی مراجعه کنندگان به موتورهای جستجو نیست. در ارائه پاسخ ها الگوریتم های رتبه بندی<sup>۱</sup> مطابق خواست کاربر عمل نمی کنند و نمونه های کم و غیر تجاری ای از موتورهای جستجو وجود دارد که به مسئله شخصی سازی<sup>۲</sup> جستجو اهمیت داده اند. لذا با توجه به اهمیت مسئله، ضرورت مطرح گشتن روش هایی نوین در این حوزه آشکارتر می گردد. رشد رو به گسترش اینترنت و به تبع آن، توسعه کسب و کارهای الکترونیک در جهان باعث شده تا وبسایت ها از اهمیت بالایی برخوردار شده، نقش غیرقابل انکاری در برقراری ارتباط الکترونیکی بین سازمان ها و مؤسسات با مشتریان شان پیدا کنند؛ چرا که در چنین محیط الکترونیکی ای، وبسایت، پل ارتباطی بین سازمان ها و مشتریان شان بوده، مشتریان از طریق صفحات مختلف وبسایت به خدمات سازمان دسترسی پیدا می کنند. در سال های ابتدای توسعه اینترنت و کسب و کارهای الکترونیک، سازمان ها و شرکت ها درصدد توسعه هرچه بیشتر وبسایت، افزایش صفحات وب و جذب هرچه بیشتر بازدیدکنندگان به سایت خود بودند، لکن در سال های اخیر، بهبود کیفیت وبسایت، اولویت بالاتری پیدا کرده است. در توجیه این مسأله به طور ساده می توان گفت که افزایش صفحات وبسایت یک سازمان، یک کلاف بزرگ و پیچ در پیچ می سازد که حرکت در آن و یافتن خدمات یا محصولات مورد نظر را کاری وقت گیر، خسته کننده و حتی ناموفق می کند. این مسأله باعث می شود تا سازمان نتواند تمام محصولات و خدماتش را به نمایش بگذارد و در نتیجه احتمالاً مشتری خود را از دست بدهد (Jaiganesh, et al, 2020).

از این رو موتورهای جستجو در عصر حاضر نقش کلیدی برای دسترسی کاربران به اطلاعات موجود در فضای وب ایفا می کنند. مهم ترین وظیفه این سرویس ها یافتن بهترین نتایج در زمانی مناسب برای حجم بالایی از درخواست های همزمان است. یکی از مهم ترین چالش های موتورهای جستجو کاستن از زمان پاسخگویی به پرس و جوهای دریافتی است که این مسئله مستقیماً بر میزان رضایتمندی کاربران و درآمد موتور جستجو تاثیر گذار است. بنابراین استفاده از تکنیک هایی که بتواند با توجه به نیازهای کاربران داده های متناسب و مطلوبی ارائه دهد یک برتری و مزیت محسوب می شود لذا یکی از روش های پرکاربرد در زمینه افزایش رضایتمندی کاربران و بهبود نتایج مطلوب موتورهای جستجو استفاده از وب کاوی و آنالیز رفتار کاربران گذشته است.

## بیان مساله

در عصر نوین با انبوه درخواست ها و نیازهایی که کاربران به کمک موتورهای جستجو درصدد مرتفع کردنشان هستند استفاده از تکنیک های سریع و کارا یک چالش و ضرورت سرویس دهندگان است. بجهت رفع مشکل رتبه بندی در جواب های یک پرس و جوی ارائه شده به یک موتور جستجو می توان با استفاده از اطلاعات بدست آمده از کاربر به یکی از دو طریق استخراج از زمینه جستجو یا استفاده از پروفایل کاربر، جواب های اولیه را دوباره بردار بندی کرد و آنگاه نتایج این بردار بندی را با استفاده از یک الگوریتم خوشه بندی، دسته بندی کرد و آنگاه با اعمال الگوریتم های مارکوف به یک زیرمجموعه کوچک از جواب ها رسید که بهبود قابل توجهی نسبت به جواب هایی که موتورهای جستجوی متداول به کاربران می دهند را در بر دارد. رتبه بندی جدیدی از نتایج پرس و جو بدست آورد که در عمل بسیار به خواست کاربر نزدیکتر باشند و شرایطی جدید را در تنوع مشخصات ارائه شده نیز برآورده کند موتورهای جستجوی کنونی در عمل با

<sup>1</sup> Ranking

<sup>2</sup> personalization

دادن نتایج با شمار زیاد کاربر را سردرگم می کنند و هدف این الگوریتم و شمایی کلی آن است که به یک ساختار جدید در ارائه نتایج به کاربر دست یابد و او را از مواجه شدن با نتایجی با شمار بیش از اندازه ای که بتواند آن ها را بسنجد باز دارد (Jaiganesh, et al, 2020). تعداد نتایجی که یک کاربر تخصصی گوگل پس از تحویل پرس و جویش می آزماید بندرت از تعداد انگشتان دست فراتر می رود. و همین مسئله ضرورت موتورهای فراگیر و با این قابلیت جدید را بیش از پیش گوشزد می کند. پروفایل کاربر با روشی جدید که مبتنی بر یادگیری از روی تاریخچه جستجوهای کاربر است ارائه می گردد و این زمانی است که این پروفایل به روشی ساده تر مانند انتخاب مستقیم خود کاربر از مجموعه ای محدود از پروفایل ها، ممکن نباشد. نتایج اولیه از موتور جستجو با رتبه بندی بهبود نیافته نیز گرفته شده اند و اکنون باید ابتدا با استفاده از بردار پروفایل، بردار بندی بر روی نتایج موجود اولیه اعمال شود. روشی که در بردار بندی بکار می رود باید بتواند در اعمال علایق کاربر به نتایج، کارآ باشد و نتیجه کلی الگوریتم را توجیه کند. روش خوشه بندی نیز باید بتواند در توجیه خود بنابر معیارهای کیفیت خوشه بندی به مقدار قابل قبولی دست یابد و سرانجام اینکه رتبه بندی ثانویه باید بهبود قابل توجهی نسبت به رتبه بندی اولیه از خود نشان دهد که این بهبود باید در راستای علایق کاربر و شخصی سازی نتایج جستجو باشد و در ضمن کاربر را از افتادن در دام جواب هایی با رتبه بندی بالا اما یکدست و شبیه به هم برهاند و یک تنوع و کیفیت را در دادن زیرمجموعه کوچکی از جواب های اولیه رعایت کند (Robertson, et al, 2019)

هر موتور جستجوی ناگزیر از آن است که قضاوت و رتبه بندی نهایی خود را نه بر اساس چند کلمه محدود یا براساس ارجاعات صفحات به همدیگر (روش مورد استفاده در گوگل)، که بر اساس نیازهای کاربری که به اطلاعات نیاز دارد، انجام دهد و قطعاً به پروفایلی از کاربر نیاز دارد. با توجه به پیشرفت های روزافزون در فناوری دیسک های سخت قاعداً می توان برای هر کاربر یک یا چند پروفایل بسته نوع کاربر ذخیره کرد و آنگاه در هر بار ورود کاربر به محیط جستجو بر اساس پروفایلش که پویا بروز می شود جستجو را انجام داد. امکان بروز رسانی پویای کاربر بنظر ناگزیر می رسد تا بتوان همواره علایق کاربر را ردگیری کرد. به نظر می رسد که استفاده از تکنیک های وب-کاوی در این فرایند ناگزیر باشد و به جهت دستیابی به سرعت و دقتی قابل قبول این تکنیک ها می توانند کمک شایانی در رتبه بندی جواب های ارائه شده به کاربر داشته باشند. در این میان الگوریتم های خوشه بندی، بردار پشتیبان و رتبه بندی صفحات با کاربرد چند گانه می توانند نقش اصلی را ایفا کنند. در روش ارائه شده در این تحقیق از مدل بردار پشتیبان استفاده شده است و برای کاهش پیچیدگی از پروفایل کاربران و ایجاد انجمن ها استفاده شده است. در حقیقت پروفایل کاربران شامل اطلاعات مربوط به زمان و توالی بازدید صفحات کاربران و میزان علاقه مندی کاربران در یک نشست می باشد. برای افزایش دقت سیستم مورد ارائه همچنین از الگوریتم خوشه بندی  $k$ -means به صورت تلفیقی با مدل بردار پشتیبان استفاده شده است.

## اهمیت موضوع

با توجه به افزایش حجم وب، پیمایش و جستجوی آن از اهمیت بالایی برخوردار است. در پیمایش این حجم وسیع از صفحات بهتر آن است، صفحاتی ابتدا پیمایش شوند که مرتبط با موضوع مورد نظر می باشند. "پیمایش هوشمند" روشی است که برای پیمایش صفحات مرتبط با یک موضوع به کار می رود. در این روش سعی بر آن است که در هنگام پیمایش، صفحات هاب<sup>۱</sup> (صفحاتی است که به صفحات مهمی اشاره می کنند) خوب تشخیص داده شوند تا از آن ها به عنوان منبعی برای رسیدن به اعتبار صفحه یا اتوریته<sup>۲</sup> (صفحاتی هستند که محتوای مهمی دارند) استفاده شود (Okabe and Seiji, 2017). وب طی یک فرآیند آشفته و غیر متمرکز در حال رشد است و این روند منجر به تولید حجم وسیعی از مستندات متصل به یکدیگر گشته است که از هیچ گونه سازماندهی منطقی برخوردار نیستند. در واقع وب به مجموعه بزرگی از داده های ساخت یافته و نیمه ساخت یافته تبدیل شده است که کاربران آن از همپوشانی داده ها رنج می برند. بنابراین تحلیل رفتارهای کاوشی کاربران وب و بررسی واقعی علایق کاربران اهمیت خاصی پیدا کرده است.

<sup>1</sup> Hub

<sup>2</sup> Page Authority

رشد چشم گیر اطلاعات در وب جهان گستر به طور ناگهانی منجر به اضافه بار اطلاعاتی شده است و از این رو پیدا کردن یک بخش خاص از اطلاعات دشوار و وقت گیر است. سیستم های شخصی سازی در سال های اخیر به منظور مقابله با این مشکل و با هدف ارائه یک تجربه شخصی به کاربران بر اساس اولویت ها، و نیازهای فردی آنها ایجاد شده اند. با رشد نمایی صفحات وب و پدیده گردآوری نمودن اطلاعات ضرورت بهره گیری از یک موتور جستجوی قدرتمند که به کاربر بتواند پیشنهادهای مناسب را ارائه بدهد، ملاحظه می شود. در این راستا کارهای متعددی انجام پذیرفته است. برای این منظور بعضا از الگوریتم های فراابتکاری و برخی دیگر از مدل های تصادفی بهره برده اند. بنابراین ارائه روش های بهینه سازی جهت بهبود نتایج موتورهای جستجوگر یک الزام در این حوزه محسوب می شود.

#### ادبیات و پیشینه

کاظمی در سال ۱۳۹۳، برای افزایش دقت در حذف کردن حالات و داشتن دقت بهتر در پیش بینی مدل، روشی را برای حذف بر اساس خطا ارایه می کند که براساس این روش برای حذف یک حالت با مرتبه بالاتر به تخمین خطای حالت های مرتبه پایین تر با مجموعه ای از جلسه هایی که توسط مرتبه بالاتر هم پوشش داده می شود، می پردازد. در واقع در این روش برای هر حالت به تعداد حالات مرتبه بالاتر آن خطا محاسبه می شود.

مقیم و همکارانش در سال ۱۳۹۳، ابتدا فایل رخداد سرویس دهنده وب را به کمک روش های پیش پردازش به فایلی حاوی نشست های کاربر تبدیل می کند و خروجی مرحله اول را که با استفاده از مدل خوشه بندی کامیانه بهینه شده، خوشه بندی می کند، سپس این مجموعه را توسط مدل مارکوف همه مراتب براساس دو روش متفاوت مورد بررسی قرار داده و مرتبه ای که بالاترین صحت کار را می دهد بعنوان برچسب برای آن خوشه انتخاب می کند. آخرین مرحله هم ورود فرد جدید به سایت و پیش بینی بهترین گزینه برای وی می باشد. Chintan و همکارانش در سال ۲۰۱۹، مروری بر آنچه در داده کاوی وب وجود دارد صورت دادند که شامل فرآیند وب کاوی کاربردی ۱ و بررسی عمیق کارهایی که تا کنون در روش پیش پردازش داده ها برای وب کاوی کاربردی انجام شده است. هنگامی که کاربران با وب سایت ارتباط برقرار می کنند آثاری با فرمت های مختلف در مکان های مختلف تولید و برجای می گذارند. این آثار در راه مناسب ثبت و ضبط می شوند، سیاهه ها ممکن است از ناخالصی ها و اختلالات جمع آوری می شوند، از اینرو تکنیک های مختلف داده کاوی نمی تواند به طور مستقیم بر آنها اعمال شود. بنابراین نیازمندی ها، مراحل و روش های پیش پردازش داده ها شرح داده شده است. تکنیک های مختلف الگوی کشف که می تواند در پیش پردازش سیاهه ها جمع آوری شود در مرحله قبل به کار برده شود و به استخراج دانش از آن منجر گردد، شرح داده می شود.

Priyanka و همکارانش در سال ۲۰۲۰، فرایند استخراج از فایل های سرور وب به منظور استخراج الگوهای استفاده شده در پیش بینی لینک های وب با کمک مدل مارکوف ارائه شده است. نتیجه این روش پیش بینی صفحه وب محبوب و یا جهت دهی رفتار کاربر می باشد. روش پیشنهادی جهت دهی رفتار کاربر تشخیص جامعه می شود بر اساس اندازه گیری شباهت دو به دو آن ها همراه با مدل مارکوف با مفهوم الگوریتم آپریوری که برای لینک وب استفاده می شود.

Meera و همکارش در سال ۲۰۱۵، یک چارچوب دو لایه ای را ارائه کرده اند. این چارچوب دو لایه ای از دو لایه به صورت آفلاین و برخط تشکیل شده است. بدین صورت که لایه ای اول مربوط به مرحله آموزش می باشد و از اطلاعات کاربران در فاز آفلاین بهره گرفته می شود و لایه ای دوم مربوط به فاز برخط می باشد که عملیات پیش بینی در این فاز به نمایش گذاشته می شود. در این مقاله از مدل های مارکوف و مارکوف مخفی به منظور مدل های پیش بینی کننده استفاده شده است نتایج هر کدام از مدل ها با استفاده از تئوری ترکیب شواهد دمپستر - شافر ترکیب شده و صفحاتی که احتمال بازدید آن ها بیشتر است نمایش داده می شود. این چارچوب قابل تغییر می باشد و از مدل های دیگر همچون خوشه بندی، قوانین انجمنی و غیره می توان در آن استفاده کرد و به آن افزود. خروجی مبتنی بر تلفیق مارکوف و مارکوف مخفی نشان می دهد که سیستم از دقت بالاتری نسبت به تک تک هر کدام از مدل ها دارا می باشد.

<sup>1</sup> Web Usage Mining (WUM)

Madhuri و همکارانش در سال ۲۰۱۹، روشی جهت رشد نامحدود خدمات وب و سیستم های اطلاعاتی مبتنی بر وب، حجم داده کلیک-استریم<sup>۱</sup> و داده کاربر جمع آوری شده توسط سازمان های مبتنی بر وب در فعالیت های روزانه خود به حجم بالایی رسیده است. تحلیل چنین حجم بالایی از داده می تواند به ارزیابی اثربخشی کمپین های تبلیغاتی، بهینه سازی عملکرد کاربردهای مبتنی بر وب کمک کند و محتوای سفارش گیری را برای ویزیتورها فراهم کند. در کار قبلی، ما روشی پیشنهاد کردیم، تحلیل الگوی رابطه ای گری با استفاده از زنجیره های مارکوف، که شامل کشف الگوهای معنی دار و ارتباطات از یک مجموعه بزرگ داده است، که اغلب در لاگ های دسترسی به سرور برنامه کاربردی و وب، لاگ های پروکسی<sup>۲</sup> و غیره ذخیره شده اند. در اینجا ما یک رویکرد جدید برای تحلیل رفتار گشت وگذار کار با استفاده از GRPA<sup>۳</sup> با زنجیره های مارکوف با طول متغیر پیشنهاد می کنیم. یک VLMLC<sup>۴</sup> یک گسترش مدل است که به تاریخچه طول متغیر اجازه ضبط شدن می دهد GPRPA با زنجیره های مارکوف با طول متغیر بطور موثر و اثربخشی نشان دهنده اطلاعات متوالی در داده استفاده از وب است و می تواند بسط داده شود تا باعث ادغام با یک مدل پیش بینی رفتار گشت وگذار کاربر وب برای کاربردهای استخراج استفاده از وب بهتر شود. تعداد زیادی از صفحات وب روی بسیاری از وبسایت ها باعث مشکلات گشت وگذار می شوند. زنجیره های مارکوف اخیرا برای مدلسازی رفتار حرکتی کاربر روی وب جهان گستر استفاده می شوند.

Zhu و همکارانش در سال ۲۰۱۸، روشی را برای ایجاد یک مدل مارکوف از وبسایت های مبتنی بر رفتار ویزیتوری گذشته پیشنهاد کرده-اند که از مدل مارکوف برای ایجاد پیش بینی های لینک استفاده شده که به کاربران جدید برای گشت وگذار در وبسایت ها کمک می کنند. یک الگوریتم برای فشرده سازی ماتریس احتمال تراکنش استفاده می شود تا صفحات وب با رفتارهای تراکنشی مشابه را کلاستر بندی کند و ماتریس تراکنش را به یک اندازه بهینه برای محاسبه احتمال اثربخشی در پیش بینی لینک فشرده کند. یک روش فوروارد مسیر حداکثری برای اصلاحات بیشتر در مورد اثربخشی پیش بینی لینک استفاده می شود. پیش بینی لینک در یک سیستم برخط به نام جستجوگر مرور برخط اجرا می شود تا به حرکت کاربر در وبسایت انطباقی کمک کند.

Montgomery و همکارانش در سال ۲۰۲۰، نشان داده اند که چطور اطلاعات مسیر می تواند با استفاده از مدل احتمالی چندگانه مرور وب گروه بندی و تحلیل شود. در روش ارائه شده با استفاده از داده ها از فروشگاه برخط اصلی ارزیابی شده است. نتایج نشان می دهند که مولفه حافظه این مدل در پیشگویی صحیح یک مسیر ضروری است. در مقایسه، مدل های ماکوف مرتبه اول و احتمالی چندگانه سنتی مسیرها را به طرز ضعیفی پیش بینی می کنند. این نتایج پیشنهاد می کنند که مسیرها ممکن است اهداف یک کاربر را منعکس کنند، که می تواند در پیش بینی حرکات آینده در یک وبسایت مفید باشد. یک کاربرد احتمالی مدل تحقیق پیش بینی تبدیل خرید است. ما متوجه می شویم که تنها بعد از شش مشاهده می توان با دقت بیش از 40% خریداران را پیش بینی کرد که بسیار بهتر از نرخ پیش بینی تبدیل خرید 7% مارک انجام شده بدون اطلاعات مسیر است. این روش می تواند برای شخصی سازی طراحی وب و پیشنهاد محصول بر اساس مسیر کاربر استفاده شود.

Kalbhori و همکارانش در سال ۲۰۱۹، تلاش کرده اند تا با استفاده از اطلاعات در خصوص صفحات از قبل بازدید شده کاربر و مرور سوابق آن، مجموعه بعدی صفحات وبی که ممکن است مشاهده نماید را پیش بینی نموده و مشکل طبقه بندی را برطرف نمایند. این برای پیش بینی رفتار کاربر بسیار مفید است در حالی که او از اینترنت به دلایل بسیاری مانند افزایش سرعت مرور و یا به حداقل رساندن بار سرور و غیره استفاده می نماید. در این مرجع بر اساس مدل پنهان مارکوف فازی پیش بینی روند بعدی کاربر پیشنهاد شده است که از نسخه k امین مرتبه مارکوف استفاده می نماید. این مرجع همچنین گزارشی از مقایسه روش های مختلف پیش بینی درخواست آینده با روش مناسب خود ارائه می نماید.

Yang و همکاران در سال ۲۰۱۸، به بررسی الگوریتم های فیلترینگ مشارکتی به کار گرفته شده در برنامه های کاربردی اینترنتی تلفن همراه پرداخته شده است. در ابتدا یک چارچوب برای سیستم پیشنهاد دهنده بر اساس داده های کاربران مختلف از جمله رتبه بندی کاربر

<sup>1</sup> Click Stream

<sup>2</sup> Proxy log

<sup>3</sup> Grey Relational Pattern Analysis

<sup>4</sup> Variable Length Markov Chains

و رفتارهای کاربران پیشنهاد می شود. سپس ویژگی های کلیدی این دو نوع داده ها بحث شده است. علاوه بر این، چندین نوع الگوریتم به عنوان روش های مبتنی بر حافظه و روش های مبتنی بر مدل طبقه بندی و مقایسه شده است و در پایان دو مطالعه موردی در مجموعه داده های مووی لینز بر اساس رتبه بندی کاربر بر اساس رفتار کاربر برای تعیین اعتبار این چارچوب ارائه شده است. Katarya و همکارش در سال ۲۰۱۷، یک سیستم توصیه گر ارائه داده اند که علاقه کاربران را نسبت به مجموعه ای از آیتم ها (کتاب، فیلم، موزیک، نرم افزار) جمع آوری می کند. این اطلاعات به صورت صریح (معمولاً به صورت جمع آوری امتیازان) و یا ضمنی (به وسیله رصد فعالیت های کاربران از قبیل آهنگ های گوش داده، نرم افزارهای دانلود شده، صفحات وب بازدید شده و غیره) به دست می آیند. این سیستم های توصیه گر تلاش می کند که به وسیله پیشنهاد اطلاعاتی که مورد علاقه کاربر باشد بر مشکل سرریز اطلاعات غلبه کند. با افزایش محبوبیت شبکه های اجتماعی، اطلاعات شبکه اجتماعی رشد سریعی داشته است و باعث شده است که کاربران با سرریز اطلاعات مواجه می شوند.

### اهداف و فرضیه های پژوهش

از جمله اهداف مهمی که در این تحقیق می توان به آن ها اشاره کرد عبارتند از:

هدف اصلی:

- بهینه سازی نتایج موتورهای جستجوگر با بهره مندی از الگوریتم های خوشه بندی و بردار پشتیبان

اهداف فرعی:

- ۱) افزایش دقت موتورهای جستجو با بهره گیری از الگوی رفتاری کاربران
- ۲) بهره گیری از بردار پشتیبان در داخل خوشه ها به منظور افزایش دقت سیستم پیشنهادی
- ۳) ایجاد پایگاه داده آموزشی قوی از روی الگوی رفتاری کاربران جهت تخمین رفتارهای آتی کاربران
- ۴) ارائه یک موتور جستجو بر وب براساس الگوهای رفتاری ایجاد شده از کاربران و الگوریتم بردار پشتیبان

فرضیه های تحقیق عبارتند از:

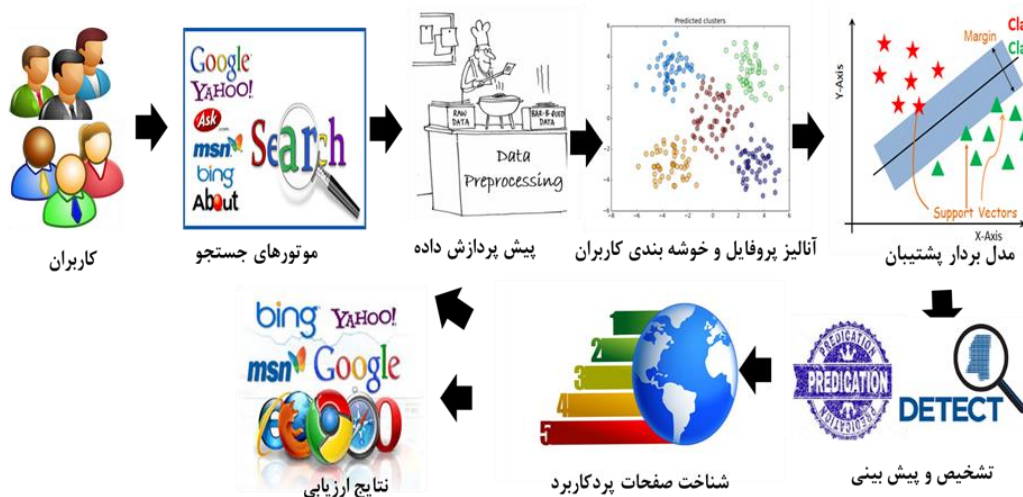
- ۱) می توان با بهره گیری از بردار پشتیبان و خوشه بندی اقدام به ارایه توصیه مناسب در موتورهای جستجو نمود.
- ۲) می توان با بهره گیری از فایل های ثبت کاربران امکان ایجاد یک پروفایل برای آنها براساس الگوی رفتاری آنها ایجاد نمود.
- ۳) بهره گیری از بردار پشتیبان منجر به افزایش دقت و کاهش زمان پاسخ سیستم پیشنهادی می گردد.

### روش تحقیق

روش پیشنهادی بر اساس وب کاوی مبتنی بر کاربرد ارائه شده است که می تواند بدون نیاز به دخالت کاربران و به صورت کاملاً صریح، نیازها و علاقه آنان را تشخیص داده و پیشنهادهای مناسب را به آن ها ارائه نماید. تشخیص الگوهای رفتاری کاربر و استخراج اطلاعات با ارزش از بین این میزان داده گسترده شامل چندین فاز می باشد که عبارتند از: پالایش داده و پیش پردازش، تبدیل داده، استفاده از تکنیک های داده کاوی و تفسیر نتایج. در این تحقیق بر روی ایجاد پروفایل کاربران تمرکز شده است. در ادامه بر روی یادگیری بردار پشتیبان تمرکز شده است و برای مدل کردن پیش بینی رفتار پیمایشی کاربر بر روی وب مناسب عمل می نمایند. به صورت کلی ورودی برای این مسائل دنباله ای از صفحات وب است که به وسیله کاربر ملاقات شده و نتیجه ساختن داده آموزشی است که بتواند محتمل ترین صفحه ای که کاربر به عنوان صفحه بعدی خواهد دید را پیش بینی کند. در این تحقیق از بردار پشتیبان استفاده شده و برای کاهش پیچیدگی از پروفایل کاربران استفاده شده است. در حقیقت پروفایل کاربران شامل اطلاعات مربوط به زمان و توالی بازدید صفحات کاربران و میزان علاقه مندی کاربران در یک نشست می باشد. برای افزایش دقت سیستم مورد ارائه همچنین از الگوریتم خوشه بندی به صورت تلفیقی با SVM استفاده شده است.



در این روش، تعاملات کاربران وبسایت مورد نظر که شامل پروفایل کاربر، صفحات وب و اطلاعات ثبت شده مربوط به استفاده کاربران می باشد، مورد پردازش قرار می گیرند. کاربران با جست و جویهای مختلف از طریق موتورهای جستجو رفتارشان ثبت و ذخیره می شود. این رفتار و تعاملات در قالب پروفایل کاربری که برگرفته از شاخص های مهمی مانند: رفتار، زمان، علائق و غیره است ذخیره می شود. این پروفایل ها باید مورد آنالیز و بررسی قرار بگیرند تا به کمک تکنیک های مختلف صفحات مورد بازدید و نیاز کاربران شناخته شوند. در ادامه ابتدا معماری روش پیشنهادی مطابق با شکل (۱) تشریح شده و سپس گام های تعریف شده برای رسیدن به هدف تحقیق با جزئیات لازم تشریح شده است همچنین به چگونگی ارزیابی و شبیه سازی اشاره شده است.



شکل ۱. بلوک دیاگرام روش پیشنهادی

### پیش پردازش داده

در این مرحله، ساختار وبسایت و محتوای آن مورد پردازش قرار گرفته و یک مدل از آن تولید می گردد. یک وبسایت، شامل تعداد زیادی صفحات وب است که در آن، صفحات وب از طریق پیوندها به هم مرتبط اند. ابتدا نیاز است که داده را از فایل های ثبت وقایع متعدد یکسان سازی کنیم؛ تا کاربران یکسان را در فایل های ثبت وقایع متفاوت تعیین شوند. این فایل ها می توانند به فرمت های معمولی<sup>۱</sup> (آدرس IP کاربر، زمان و تاریخ دستیابی، روش درخواست (Get, Post, ...))، آدرس صفحه دستیابی شده، پروتکل (HTTP, HTTP1.1)، کد خروجی و تعداد بایت های جا به جا شده) یا توسعه یافته<sup>۲</sup> ذخیره شوند. در فرمت توسعه یافته فیلد ارجاع دهنده<sup>۳</sup> و فیلدهای عامل کاربر<sup>۴</sup> وجود دارند. اولین گام در پیش پردازش داده حذف محتوای نامرتب است. برای مثال فایل های تصویری مانند gif, jpg, با چک کردن رشته ی رکوردها در فایل ثبت وقایع حذف می شوند. اگر این رکوردها توسط جست و گو وب ایجاد شود با مقایسه با robot.txt حذف می شود.

### ایجاد پروفایل برای کاربران و ترسیم ماتریس کاربر-صفحه

مرحله ی تشخیص کاربران پس از پاک سازی داده ها آغاز می شود. تشخیص کاربران یک مرحله ی بسار پیچیده است. از آنجا که ممکن است فایل ثبت وقایع کاربران توسط یک سرور وب تنها ضبط شود و یا ممکن است توسط پروکسی در ترکیب با سرورهای وب پیچیده ضبط شود. در این گام، آدرس یا DNS آدرس ارجاع دهنده، آدرس الکترونیک و عامل کاربرد، زمینه های مرتبطی در فایل ثبت وقایع

<sup>1</sup> CLF (common log file)

<sup>2</sup> ELF(extend log file)

<sup>3</sup> Referrer url

<sup>4</sup> User agent

هستند. آدرس DNS که همان IP است در  $CLF^1$  قرار دارد. آدرس ارجاع دهنده<sup>۲</sup>، کاربر تایید شده<sup>۳</sup>، عامل کاربردی<sup>۴</sup> در فایل  $ELF^5$  قرار دارند. در طول تعاملات کاربران با صفحات وبسایت، تمامی رفتارهای کاربران که شامل مرور صفحات وبسایت می باشد در یک فایل متنی ذخیره می گردند. به منظور رسیدن به الگوهای استفاده از وب سایت، باید فایل های ثبت رفتارهای کاربران مورد پردازش قرار گیرد. برای شناسایی کاربران و تشخیص الگوهای استفاده از وب مربوط به آن کاربر با توجه به ساختار فایل های ثبت رخداد سمت سرور، راه کاری ارائه شده که کاربران مختلف را از هم متمایز می نماید همچنین الگوهای استفاده از صفحات وب را تشخیص می دهد. جهت انجام این کار، در این بخش از مشخصه شناسه نشست استفاده می شود اما مشکلی که در این مورد می تواند ایجاد شود این است که یک کاربر باید به یک مرورگر و محدوده زمانی خاصی محدود شود و همچنین نمی تواند استفاده های مکرر یک کاربر را تشخیص دهد یا به عبارت دیگر، درخواست های مشاهده صفحات مربوط به یک کاربر، محدود می شود به درخواست های که در طول مدت یک نشست ارائه شده است. روشی که در این تحقیق در نظر گرفته ایم استفاده از آدرس IP به همراه شناسه Agent که دسترسی به هر دو خیلی راحت تر بوده و در فایل های ثبت رخداد سمت سرور نیز ثبت می شوند و مشکلی که در این روش وجود دارد این است که نمی تواند منحصر بودن شناسه ترکیبی IP و Agent را تضمین نماید. برای حل این مشکل نیز فرض می کنیم که از هر کامپیوتر تنها یک کاربر با یک IP معتبر می تواند استفاده نماید. هدف از این مرحله تعیین راه کاری است که بتوان صفحات مشاهده شده توسط کاربران را نشان داد. جهت انجام این کار مجموعه هایی بصورت زیر تعریف می شود:

الف) مجموعه آدرس صفحات وب سایت  $P = \{P_1, P_2, \dots, P_n\}$

ب) مجموعه تراکنش های کاربر  $T = \{t_1, t_2, \dots, t_m\}$

در اینجا هر تراکنش شامل یک زیر مجموعه از مجموعه P (مجموعه آدرس صفحات) می باشد. مطابق با مدل ارائه شده در روش پیشنهادی خوشه بندی صفحات وب بر اساس وقوع مشترک در بین نشست های کاربران محاسبه می شود. برای M نشست و N صفحه یک نشست  $S_i$  توسط آمین سطر ماتریس صفحه-نشست می تواند نشان داده شود. هر ورودی  $X_{ik}$  وزن صفحه  $P_k$  در نشست را نشان می دهد که مقدار ۰ نشان دهنده این است که صفحه  $P_k$  در نشست حاضر نیست. بنابراین، در ماتریس صفحه-نشست X هر سطر نشان دهنده یک نشست بر حسب صفحات درخواست شده در آن نشست است.

#### خوشه بندی بر اساس الگوریتم K-means

داده های مهمی در پردازش استفاده از وب وجود دارد که همه آنها ممکن است در دسترس نباشند، که به برخی از آنها مانند فایل های ثبت رخداد سمت سرور، محتویات سایت، داده های مرتبط با کاربر و غیره می توان اشاره کرد. اما مرحله ی تشخیص کاربران پس از پاک سازی داده ها آغاز می شود. تشخیص کاربران یک مرحله ی بسار پیچیده است. از آنجا که ممکن است فایل ثبت وقایع کاربران توسط یک سرور وب تنها ضبط شود و یا ممکن است توسط پروکسی در ترکیب با سرورهای وب پیچیده ضبط شود. در این گام به کمک آنالیز پروفایل کاربران با توجه به میزان علاقه مندی در هر نشست، صفحات محبوب انتخاب و در نهایت بر اساس علاقه موجود به صفحات توسط کاربران خوشه بندی انجام می گردد (Liu, et al, 2017).

برای استخراج میزان علاقه مندی کاربران به صفحات در یک نشست از ماتریس میانگین هارمونی به کمک دو معیار فرکانس بازدید صفحه و مدت زمان حضور کاربر در صفحه استفاده می شود. رابطه (۱) این مهم را نشان می دهد.

$$Interest(Page) = \frac{frequency(p) * Duration(p)}{frequency(p) * Duration(p)} = z \quad (1)$$

<sup>1</sup> Common log file

<sup>2</sup> Referrer URL

<sup>3</sup> Authuser

<sup>4</sup> User agent

<sup>5</sup> extend log file



در حقیقت  $a$  بیانگر تعداد مراجعات کاربر به یک صفحه وب در یک نشست و  $b$  مدت زمان سپری شده بر روی یک صفحه است. این داده ها از طریق پروفایل کاربران و بر اساس رفتار کاربران ثبت و ذخیره می شود. خروجی ماتریس میانگین هارمونیک بصورت ماتریسی از  $U \times P$  است که سطرهای ماتریس، کاربرها و ستون های آن صفحات بازدید شده توسط کاربران را نشان می دهد.

$$A_{U \times P} = \begin{matrix} U1 \\ U2 \\ \vdots \\ Un \end{matrix} \begin{bmatrix} P1 & P2 & \dots & Pn \\ Z_{11} & Z_{12} & \dots & Z_{1n} \\ Z_{21} & Z_{22} & \dots & Z_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ Z_{n1} & Z_{n2} & \dots & Z_{nn} \end{bmatrix}_{4 \times 4}$$

با توجه به تشریح انجام شده برای محاسبه فرکانس صفحه از رابطه (۲) و مدت زمان حضور کاربر در صفحه از رابطه (۳) استفاده می شود.

$$a = \text{Frequency}(u, p) = \frac{\text{number of visits (page)}}{\sum_{p \in \text{visited}} p(\text{number of visits}(p))} \quad (2)$$

$$b = \text{Duration}(u, p) = \frac{\text{total duration}(p)}{\max_{p \in \text{visited}(p)} (\text{total size}(p))} \quad (3)$$

در رابطه (۳) Size بیانگر اندازه صفحه بر حسب بایت است. مطابق با روابط تعریف شده هر سطر از ماتریس را می توان به صورت برداری بر اساس میزان علاقه هر کاربر به صفحه به صورت  $U_1(a, b, c, \dots)$  و  $U_2(a', b', c', \dots)$  نشان داد.

#### تشخیص با الگوریتم یادگیری بردار ماشین

از این الگوریتم برای تشخیص صفحات پرکاربرد و پیشنهاد به کاربر استفاده می شود در حقیقت به کمک شاخص های آموزشی این وضعیت آنالیز و بررسی می شود و در صورت قابلیت توانایی انجام دادن توسط منبع با اتکا بر حدا آستانه بالا و بهبود جستجوی موتور جستجو انجام می شود. ماشین بردار پشتیبان، الگوریتم طبقه بندی<sup>۱</sup> بوده و به عنوان یکی از بهترین تکنیک های دسته بندی و پیش بینی و تشخیص outlier شناخته می شود و برخلاف الگوریتم های خوشه بندی در دسته یادگیری با نظارت محسوب می شود و دو فاز آموزش و تست دارد. در واقع SVM برای اینکه داده های غیر خطی را از هم تفکیک کند باید از کرنل های مختلف استفاده کند. برای این کار دیگر در فضای دو بعدی کار نمی کند بلکه داده ها به فضایی با ابعاد بیشتر نگاشت داده می شوند تا بتوان آن ها را در این فضای جدید بصورت خطی تفکیک نمود. در حقیقت به کمک این الگوریتم که در دو فاز یادگیری و تشخیص استفاده می شود داده های آموزشی از پیش به دو کلاس پربازدید و کم بازدید برچسب شده و تقسیم می گردند، این دیتاست آموزشی با افزایش آزمایشات تکمیل تر از قبل می شود و نتایج بدست آمده پشیمانی تخصیص را کاهش می دهد.

#### یافته ها

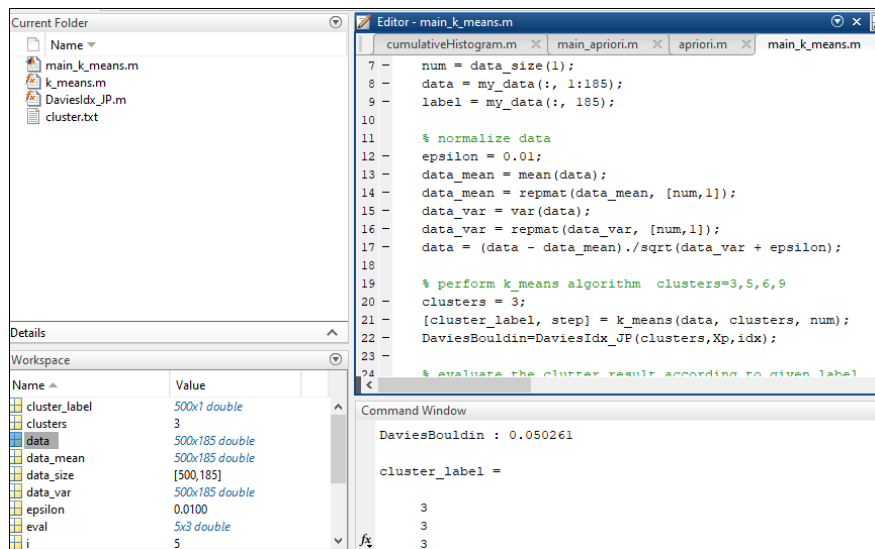
در این تحقیق برای مقایسه و ارزیابی سیستم پیشنهاد دهنده با سایر سیستم ها از دیتاست ناسا استفاده شد. این دیتاست در <https://www.kaggle.com/datasets/souhagaa/nasa-access-log-dataset> در دسترس قرار دارد. ویژگی های زیادی در هر رکورد از این دیتاست وجود دارد که شامل آدرس آی پی یا نام دی ان اسی درخواست دهنده تاریخ و ساعت درخواست صفحه آدرس و نام صفحه نوع ارسال صفحه مثلاً گت یا پست و غیره بایت ارسالی و دیگر دارد. در اینجا با توجه به حجم زیاد دیتاست نمونه آزمایش شده به تعداد ۵۰۰ رکورد مورد انتخابی قرار گرفته است. رکوردهای انتخاب شده از نشست های مختلف ثبت شده در زمان های مختلف به صورت تصادفی انتخاب شده است. شکل (۲) نمونه ای از این مجموعه داده ها را نشان می دهد.

<sup>1</sup> classifier

```
128.158.41.74 - [06/Jul/1995:09:37:57 -0400] "GET /icons/menu.xbm HTTP/1.0" 200 527
128.158.41.74 - [06/Jul/1995:09:37:57 -0400] "GET /icons/image.xbm HTTP/1.0" 200 509
192.189.46.188 - [06/Jul/1995:09:37:57 -0400] "GET /images/ksclogo-medium.gif HTTP/1.0" 200 5866
192.189.46.188 - [06/Jul/1995:09:37:57 -0400] "GET /images/NASA-logomall.gif HTTP/1.0" 200 786
192.189.46.188 - [06/Jul/1995:09:37:57 -0400] "GET /images/MOSAIC-logomall.gif HTTP/1.0" 200 363
192.189.46.188 - [06/Jul/1995:09:37:58 -0400] "GET /images/USA-logomall.gif HTTP/1.0" 200 234
herman.and.nl - [06/Jul/1995:09:37:58 -0400] "GET /pub/wiavn/readme.txt HTTP/1.0" 404 -
199.10.135.154 - [06/Jul/1995:09:37:59 -0400] "GET /icons/menu.xbm HTTP/1.0" 200 527
192.189.46.188 - [06/Jul/1995:09:37:59 -0400] "GET /images/WORLD-logomall.gif HTTP/1.0" 200 669
aem003.aerlepa.gov - [06/Jul/1995:09:38:00 -0400] "GET /images/launch-logo.gif HTTP/1.0" 200 1713
host62.ascend.interop.eunet.de - [06/Jul/1995:09:38:00 -0400] "GET /history/history.html HTTP/1.0" 200 1602
www-42.proxy.aol.com - [06/Jul/1995:09:38:02 -0400] "GET /shuttle/missions/missions.html HTTP/1.0" 200 8677
199.10.135.154 - [06/Jul/1995:09:38:04 -0400] "GET /icons/sound.xbm HTTP/1.0" 200 530
jeuro.dia.net - [06/Jul/1995:09:38:05 -0400] "GET /shuttle/countdown/shot.html HTTP/1.0" 200 4524
163.205.3.60 - [06/Jul/1995:09:38:06 -0400] "GET /kx.html HTTP/1.0" 200 7067
163.205.137.21 - [06/Jul/1995:09:38:07 -0400] "HEAD /shuttle/missions/sts-71/images/KSC-95EC-0874.jpg HTTP/1.0" 200 0
lunatic.bangor.ac.uk - [06/Jul/1995:09:38:07 -0400] "GET /shuttle/missions/sts-71/images/KSC-95EC-0544.jpg HTTP/1.0" 200 70128
134.217.85.114 - [06/Jul/1995:09:38:08 -0400] "GET /shuttle/countdown/ HTTP/1.0" 200 3985
dd13-031.compuerve.com - [06/Jul/1995:09:38:09 -0400] "GET /shuttle/missions/sts-71/images/images.html HTTP/1.0" 200 7634
163.205.3.60 - [06/Jul/1995:09:38:10 -0400] "GET /images/ksclogo-medium.gif HTTP/1.0" 200 5866
herman.and.nl - [06/Jul/1995:09:38:10 -0400] "GET /pub/wiavn/readme.txt HTTP/1.0" 404 -
163.205.137.21 - [06/Jul/1995:09:38:11 -0400] "GET /shuttle/missions/sts-71/images/KSC-95EC-0875.jpg HTTP/1.0" 200 41160
192.189.46.188 - [06/Jul/1995:09:38:11 -0400] "GET /shuttle/countdown/ HTTP/1.0" 200 3985
aem003.aerlepa.gov - [06/Jul/1995:09:38:12 -0400] "GET /shuttle/resources/orbiters/orbiters.html HTTP/1.0" 200 7025
128.158.41.74 - [06/Jul/1995:09:38:12 -0400] "GET /history/apollo/apollo-13/images/70HC467.gif HTTP/1.0" 200 159813
192.189.46.188 - [06/Jul/1995:09:38:12 -0400] "GET /shuttle/countdown/count.gif HTTP/1.0" 200 40310
aem003.aerlepa.gov - [06/Jul/1995:09:38:13 -0400] "GET /shuttle/resources/orbiters/orbiters-logo.gif HTTP/1.0" 200 4179
163.205.3.60 - [06/Jul/1995:09:38:13 -0400] "GET /images/NASA-logomall.gif HTTP/1.0" 200 786
nbl-10.larc.nasa.gov - [06/Jul/1995:09:38:14 -0400] "GET /shuttle/technology/sts-newsref/stsref-toc.html HTTP/1.0" 200 81920
163.205.3.60 - [06/Jul/1995:09:38:14 -0400] "GET /images/MOSAIC-logomall.gif HTTP/1.0" 200 363
163.205.3.60 - [06/Jul/1995:09:38:15 -0400] "GET /images/USA-logomall.gif HTTP/1.0" 200 234
163.205.3.60 - [06/Jul/1995:09:38:15 -0400] "GET /images/WORLD-logomall.gif HTTP/1.0" 200 669
128.159.124.51 - [06/Jul/1995:09:38:16 -0400] "GET /kx.html HTTP/1.0" 200 7067
aem003.aerlepa.gov - [06/Jul/1995:09:38:16 -0400] "GET /shuttle/resources/orbiters/orbiters-logo.gif HTTP/1.0" 200 1932
```

شکل ۲. بخشی از مجموعه داده ناسا مورد مطالعه

با توجه به روش پیشنهادی ارائه شده پس از پیش پردازش داده، ایجاد پروفایل ثبت کاربران و استخراج علاقه‌مندی به صفحات در مرحله بعد خوشه‌بندی کاربران جهت کاهش پیچیدگی فضای مساله انجام می‌گردد. قابل ذکر است که مقادیری که صفر نشان داده شده‌اند بیانگر این موضوع است که هیچ نشست در این صفحه توسط کاربر انجام نشده است. در ادامه با توجه به ۵۰۰ رکورد مورد بررسی از نشست‌های انجام شده ۱۸۵ صفحه مورد آنالیز و بررسی قرار گرفته شده است. عبارتی کاربران در نشست‌های مختلف میزان علاقه‌مندی متفاوتی به صفحات مورد بازدید داشته‌اند که در اینجا نشان داده شده است. تعداد خوشه‌ها با بررسی چندین وضعیت و انتخاب  $k$  بهینه انجام شده است و این امکان به کمک شاخص دیوس بولدین محقق شده است. تعداد خوشه‌های مورد بررسی ۳، ۵، ۷ و ۹ خوشه می‌باشد که از بین سناریوهای مختلف مقدار  $k=9$  انتخاب شده است و این مهم با توجه به کم بودن مقدار شاخص بولدین در مقایسه با سناریوهای دیگر بوده است. در شکل (۳) بخشی از شبیه‌سازی انجام شده برای خوشه‌بندی توسط مدل مورد توسعه در ابزار متلب نشان داده شده است.



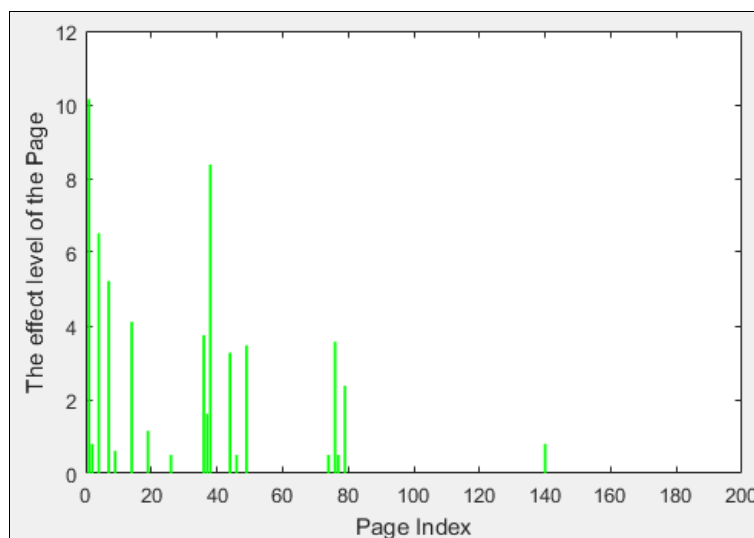
شکل ۳. بخشی از شبیه‌سازی انجام شده برای خوشه‌بندی نشست‌های انجام شده توسط کاربران در ابزار متلب

با توجه به روش پیشنهادی ارائه شده پس از پیش پردازش داده، ایجاد پروفایل ثبت کاربران و استخراج علاقه‌مندی به صفحات در مرحله بعد خوشه‌بندی کاربران جهت کاهش پیچیدگی فضای مساله انجام شده است.

جدول ۱. خروجی حاصل از پیش‌بینی کاربران جدید از تعلق به خوشه های تعریف شده

خوشه مورد انتظار	پیش بینی شده	خوشه مورد انتظار	پیش بینی شده
6	'6'	8	'8'
6	'6'	1	'1'
1	'1'	6	'6'
2	'۱'	8	'8'
7	'7'	4	'4'
7	'7'	2	'2'
1	'1'	1	'1'

شکل (۴) میزان علاقه‌مندی کاربران به صفحات در دیتاست مورد مطالعه را نشان می‌دهد. در ادامه هر خوشه پیشنهاد شده با توجه به میزان علاقه‌مندی صفحات موجود در آن خوشه به ترتیب به کاربران پیشنهاد داده می‌شود. در آزمایش انجام شده از بین ۹ خوشه موجود ۶ خوشه انتخاب شده‌اند در حقیقت صفحات موجود در این خوشه‌ها به ترتیب میزان علاقه‌مندی به کاربران مورد تست واقع شده پیشنهاد می‌شود. با توجه به اهمیت صفحات (میزان علاقه‌مندی) در هر خوشه امکان پیشنهاد صفحات مختلفی وجود دارد اما مطابق با جدول (۲) در این آزمایش با ۵۰۰ کاربر آموزش داده شده در اکثر خوشه‌های منتخب ۱۰ صفحه نشان داده شده ارجح‌تر هستند. اعتبارسنجی در سیستم‌های توصیه‌گر غالباً بر روی موارد مرتبط با دقت، فراخوانی<sup>۱</sup> و پوشش<sup>۲</sup> استفاده می‌گردد.



شکل ۴. میزان علاقه‌مندی کاربران به صفحات در پایگاه داده مورد مطالعه

جدول ۲. ده صفحه پیشنهادی مهم با توجه به میزان علاقه‌مندی

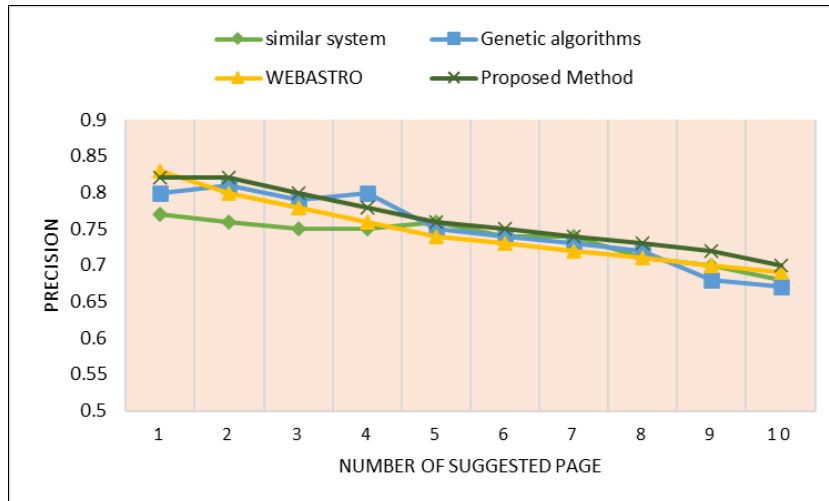
اولویت	نام صفحه	اولویت	نام صفحه
1	ksc	6	lcc
2	mission-sts-71	7	apollo
3	tour	8	apollo-15-info
4	history	9	procurement
5	mission-sts-70	10	sts-gnnc

برای ارزیابی و مقایسه سیستم از سیستم مشابه در (Chidambaram, et al, 2018; Le, et al, 2024; Thwe, 2021) ارائه گردیده،

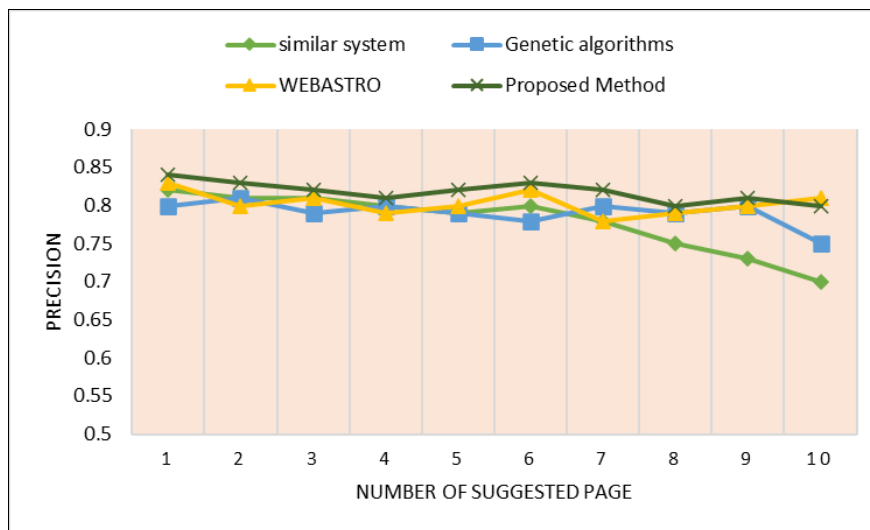
<sup>1</sup> recall

<sup>2</sup> coverage

استفاده شده است. شکل (۵) اعمال مراتب مختلف مدل بردار پشتیبان بر روی سیستم را نشان می دهد. همانطور که مشاهده می شود میزان دقت سیستم در مراتب بالاتر مدل تغییری نکرده است؛ به همین دلیل در این تحقیق برای کاهش میزان پیچیدگی سیستم، مدل بردار پشتیبان برای اعمال بر روی سیستم انتخاب شده است تا علاوه بر دارا بودن دقت مناسب، از افزایش پیچیدگی سیستم جلوگیری گردد. نتایج حاصل از ارزیابی پیاده سازی در ادامه نشان داده شده است.

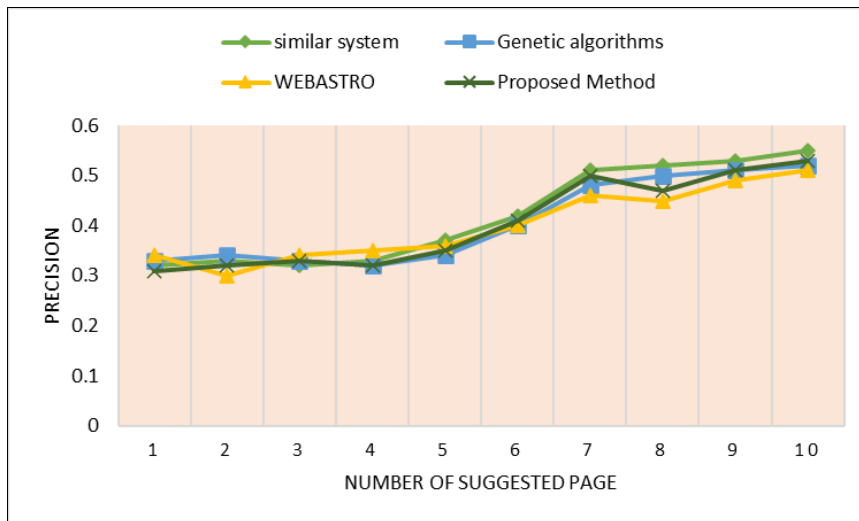


شکل ۵. ارزیابی دقت (Chidambaram, et al, 2018; Le, et al, 2024; Thwe, 2021)



شکل ۶. ارزیابی فراخوانی (Chidambaram, et al, 2018; Le, et al, 2024; Thwe, 2021)

با بیشتر شدن تعداد صفحات پیشنهادی دقت سیستم کاهش می یابد، یعنی تعداد صفحات پیشنهادی درست کمتری نسبت کل صفحات، پیشنهاد شده است. طبق نمودارهای به دست آمده و همانطور که انتظار می رفت، میزان دقت سیستم پیشنهادی نسبت به سیستم مشابه بهبود یافته است. این بهبود دقت همراه با کاهش زیاد حجم داده ها ورودی به سیستم می توان مزیت بسیار خوبی برای این سیستم نسبت سیستم مشابه باشد.



شکل ۷. ارزیابی پوشش (Chidambaram, et al, 2018; Le, et al, 2024; Thwe, 2021)

### بحث و نتیجه گیری

روش پیشنهادی ارائه شده توانسته است حجم بسیار زیادی از اطلاعاتی را که میزان اهمیت آن‌ها پایین می‌باشد را حذف نماید و با این کار میزان زمان پاسخگویی سیستم را کاهش داد. زمان پاسخگویی سریع یکی از ارکان اصلی هر سیستمی می‌باشد که پایین بودن آن نشان از کارایی سیستم مذکور می‌باشد. استفاده از سیستم خوشه‌بندی بهینه شده یکی دیگر از مزایای این سیستم در کاهش میزان پیچیدگی آن می‌باشد. به طوری که با کوچک کردن ماتریس‌های به کار رفته در آن با حذف موارد غیر لازم میزان فضای مورد استفاده برای ذخیره‌سازی آن‌ها را نیز کاهش می‌دهد. این سیستم با حذف اشتباهات کاربران در انتخاب پیوندها و استفاده از مدل پیش‌بینی کارآمد می‌تواند باعث افزایش میزان دقت کلی سیستم شود و به این ترتیب پیشنهادهای مناسب‌تری را به کاربران ارائه دهد. روش ارائه شده در این تحقیق بر اساس وب‌کاوی مبتنی بر کاربرد ارائه گردیده است که می‌تواند بدون نیاز به دخالت کاربران و به صورت کاملاً صریح، نیازها و علاقه آنان را تشخیص داده و پیشنهادهای مناسب را به آن‌ها ارائه نماید. هدف از این تحقیق ارائه روشی جهت بهبود عملکرد موتورهای جستجو با استفاده از وب‌کاوی مبتنی بر کاربرد و بردار پشتیبان می‌باشد که نتایج به دست آمده از این تحقیق حاکی از کارا بودن این سیستم در کاربردی بودن آن در موارد می‌باشد. در حقیقت ارائه صفحات متناسب با نیازمندی کاربران به صورت خودکار و با دقت و کارایی بالا علاوه بر کاهش زمان پاسخ‌دهی و افزایش پوشش‌دهی نیازهای کاربران، سرویس‌دهنده نیز مشمول سوددهی واقع می‌شود، بعبارتی سرویس‌دهنده با ارائه محصولات و سرویس‌های خود در این امر نیز موفق می‌شود که با شناخت علایق کاربران بهترین تصمیم‌ها را در راستای پیشرفت و سوددهی خود انجام دهد.

پیشنهادهایی که برای کارهای آینده در این زمینه می‌باشند عبارتند از:

- یکی از موارد اطلاعاتی موجود در فایل‌های ثبت وب آدرس IP کاربران می‌باشد. با استفاده از این اطلاعات می‌توان از موقعیت جغرافیایی کاربران اطلاع حاصل کرد. از این اطلاعات می‌توان در سیستم پیشنهاد دهنده برای ارائه پیشنهادها مرتبط به آن موقعیت جغرافیایی استفاده نمود.
- استفاده از سیستم ردیابی حرکات چشم یکی دیگر از روش‌هایی می‌باشد که می‌تواند در این گونه سیستم‌های پیشنهاد دهنده به کار آید. از آنجا که این سیستم‌ها به صورت صریح و بدون دخالت کاربر عمل می‌نمایند این تکنیک می‌تواند بسیار مفید واقع شده و اطلاعات بسیار مهمی را از علاقه کاربران در اختیار سیستم قرار دهد.

با توجه به مدت زمان مشاهده صفحات توسط کاربران، می توان پی برد که مدت زمان بازدید از یک صفحه خاص توسط کاربران مختلف، متفاوت می باشد. به نظر می رسد با جداسازی این کاربران با توجه به مدت زمان بازدید آن صفحه و قرار دادن کاربران با مدت زمان بازدید مشابه در یک خوشه، می توان به دقت بالاتری از پیش بینی در سیستم پیشنهادی دست یافت.

## منابع

- کاظمی، شهره، (۱۳۹۳) " ترکیب مدل مارکوف به همراه روش خوشه بندی دوره های کاربر برای پیش بینی صفحه بعدی"، همایش ملی مهندسی رایانه و مدیریت فناوری اطلاعات، تهران.
- مقیم، مهدی و همکاران، (۱۳۹۳) " تلفیق خوشه بندی و مدل مارکوف در یک چارچوب جدید برای پیش بینی صفحه بعدی انتخابی توسط کاربر"، نوزدهمین کنفرانس ملی سالانه انجمن کامپیوتر ایران، تهران.
- Jaiganesh, M., Ramadoss, B., Vincent, A. and Mercy, S. (2020). Performance Evaluation of Cloud Services with Profit Optimization, *Procedia Computer Science*, vol. 54, pp. 24.
- Robertson, A. M., & Willett, P. (2019), "An upperbound to the performance of ranked output searching: optimal weighting of query terms using a genetic algorithm". *Journal of Documentation*, Vol 4, pp.405–420.
- Okabe, M., Seiji, Y., (2017), "Learning Filtering rulesets for ranking refinement in relevance Feedback", *Knowledge-Based Systems* 18 , pp117-124.
- Chintan. R., Varnagar. N, Madhak. T, Kodinariya. M, NRathod. J., (2019), "Web Usage Mining: A Review on Process, Methods and Techniques", *IEEE ,Information Communication and Embedded Systems (ICICES)*, International Conference , pp: 40 .
- Priyanka S. Panchal, Prof. Urmi D. Agravat., (2020), "Hybrid Technique for User's Web Page Access Prediction based on Markov Model", *IEEE 4th ICCCNT, Computing, Communications and Networking Technologies*, pp. 1 – 8.
- Meera. N, and Sakina Banu. SH., (2015), "Predicting user's web navigation behavior using hybrid approach", *Procedia Computer Science* 45, pp.3-12.
- Madhuri, B. C., Chandulal, A. J., Ramya, K., & Phanindra, M. (2019). "Analysis of users' web navigation behavior using grpa with variable length markov chains", *International Journal of Data Mining & Knowledge Management Process*, 1(2), pp.1-20.
- Zhu. J, Hong. J, and G. Hughes. J., (2018), "Using Markov Chains for Link Prediction in Adaptive Web Sites", *Soft-Ware, LNCS* 2311, pp. 60–73.
- Montgomery., L. A, Li. SH, Srinivasan. K, and C. Liechty. J., (2020), "Modeling Online Browsing and Path Analysis Using Clickstream Data", *Marketing science*, 23(4), pp.579.
- Kalbhor, M., & Jain, K. (2019), "Fuzzy based hybrid approach for user request prediction using Markov model", In *Computer, Communication and Control (IC4)*, vol 3, pp501.
- Yang, Z., Wu, B., Zheng, K., Wang, X., & Lei, L., (2018), "A survey of collaborative filtering-based recommender systems for mobile internet applications". *IEEE Access*, Vol 4, pp.3273-3287.
- Katarya, R., & Verma, O. P. (2017). "An effective collaborative movie recommender system with cuckoo search", *Egyptian Informatics Journal*, 18(2), pp.105-112.
- Liu, H., Wu, J., Liu, T., Tao, D., Fu, Y. (2017). "Spectral Ensemble Clustering via Weighted K-Means: Theoretical and Practical Evidence", *IEEE Transactions on Knowledge and Data Engineering*, Vol 29, pp.1129 – 1143.
- Chidambaram, S., Srinivasagan, K.G., (2018), "Performance evaluation of support vector machine classification approaches in data mining", *Springer in Cluster Computing* , Vol3, pp.1–8.
- Le, T. D., Le-Dinh, T., & Uwizeyemungu, S. (2024). "Search engine optimization poisoning: A cybersecurity threat analysis and mitigation strategies for small and medium-sized enterprises", *Technology in Society*, 76, 102470.
- Thwe, P. (2021). "Web page access prediction based on integrated approach", *International Journal of Computer Science and Business Informatics*, 12(1), pp.55-64.





## A Method To Optimize Search Engine Results Using Clustering And Support Vector Machine (SVM) Algorithms

### Abstract

Nowadays, the knowledge and information of the web is increasing, and the number of websites is also increasing day by day; As a result, there is a need to provide conditions in discussions such as e-commerce, search engines, targeted advertising, etc. based on the approach of beginner and professional users from different pages. But to be able to predict the needs and interests of different users and to suggest their favorite pages is associated with complications and it should be possible to do this with optimal methods and by using the movement patterns of users recorded in the web servers. Various methods for predicting web pages have been presented, but the measures taken in this field have not yet reached a level that provides full user satisfaction, because each of the techniques available in this field has its advantages and disadvantages. Times and conditions are different, and for this reason, one technique cannot be determined to be better than other techniques in all aspects. Therefore, in order to achieve the desired quality, more research should be done so that the obtained results have high and acceptable accuracy and the way of doing the work is also a little complicated. In this research, a method to improve the performance of search engines using clustering and support vector has been presented, which consists of several different steps, including data pre-processing, creating a user registration profile, extracting the level of interest of the user, clustering and support vector. In addition to reducing the volume of information, the presentation of the proposed method has reduced the complexity of the system, and in the forecasting stage, by making changes in the support vector model, the accuracy of the system has also increased compared to the base model. The support vector model and the clustering method used in this system worked optimally to obtain acceptable results from the system. The evaluation criteria in this research are precision, recall and coverage. The overall accuracy of the system is equal to 83.5. The evaluation results of this system with the help of MATLAB tool with other similar tools show better and optimal performance.

**Keywords** web mining, clustering, user behavior patterns, support vector