





دانشگاه صنعتی اصفهان

دانشکده برق و کامپیوتر

خوشه بندی کاربران وب و واکنشی اولیه صفحات وب

با استفاده از آنالیز معنایی پنهان احتمالاتی

پایان نامه کارشناسی ارشد مهندسی کامپیوتر- نرم افزار

فاطمه زهرا حیدری

استاد راهنما

دکتر محمدعلی منتظری

دکتر احمد برآنی



دانشگاه صنعتی اصفهان

دانشکده برق و کامپیوتر

پایان نامه کارشناسی ارشد مهندسی کامپیوتر - نرم افزار خانم فاطمه زهرا حیدری

تحت عنوان

خوشه بندی کاربران وب و واکشی اولیه صفحات وب با استفاده از آنالیز معنایی پنهان احتمالاتی

در تاریخ ۱۳۹۳/۰۶/۲۶ توسط کمیته تخصصی زیر مورد بررسی و تصویب نهایی قرار گرفت

- | | |
|------------------------|----------------------------------|
| دکتر محمدعلی منتظری | ۱- استاد راهنمای پایان نامه |
| دکتر احمد برآنی | ۲- استاد راهنمای پایان نامه |
| دکتر عبدالرضا میرزایی | ۳- استاد داور |
| دکتر محمدعلی خسروی فرد | ۴- سرپرست تحصیلات تکمیلی دانشکده |

تقدیم به

پدرم به استواری کوه

مادرم به زلالی چشمه

همسرم به صمیمیت باران

کلیه حقوق مادی مترتب بر نتایج
مطالعات، ابتکارات و نوآوریهای ناشی
از تحقیق موضوع این پایان نامه متعلق
به دانشگاه صنعتی اصفهان است.

فهرست مطالب

<u>صفحه</u>	<u>عنوان</u>
	فصل اول : مقدمه
۱-۱	انگیزه و پیش زمینه ها ۱
۲-۱	هدف پژوهش ۲
۳-۱	ساختار پژوهش ۳
	فصل دوم: استخراج اطلاعات از داده های کاربرد وب
۱-۲	مقدمه ۵
۲-۲	داده کاوی و استخراج دانش ۵
۳-۲	وب کاوی ۷
۴-۲	چالش های وب کاوی ۸
۵-۲	اجزای اصلی وب کاوی ۸
۱-۵-۲	محتوا کاوی وب ۹
۲-۵-۲	ساختار کاوی وب ۹
۳-۵-۲	کاربرد کاوی وب ۹
۶-۲	آماده سازی و مدل سازی داده ۱۱
۱-۶-۲	منابع داده ۱۱
۲-۶-۲	آماده سازی و پیش پردازش داده ها ۱۴
۷-۲	کشف الگو از داده های کاربرد وب ۲۰
۸-۲	سطوح و انواع تحلیل ۲۲
۹-۲	کاربرد های کاربرد کاوی وب ۲۳
۱۰-۲	جمع بندی ۲۴
	فصل سوم : خوشه بندی کاربران وب
۱-۳	مقدمه ۲۵
۲-۳	خوشه بندی ۲۵
۳-۳	مروری بر کارهای انجام شده در زمینه خوشه بندی کاربران وب ۲۸
۴-۳	روش های براساس عامل های پنهان ۳۱

۳۲.....	تکنیک تحلیل معنایی پنهان	۱-۴-۳
۳۵.....	تکنیک نمایه سازی تصادفی	۲-۴-۳
۳۷.....	تکنیک تحلیل معنایی پنهان احتمالاتی	۳-۴-۳
۴۱.....	کشف و آنالیز الگوهای کاربردی با استفاده از تحلیل معنایی پنهان احتمالاتی	۴-۴-۳
۴۶.....	مروری بر کارهای انجام شده در زمینه خوشه بندی کاربران وب براساس عامل های پنهان	۵-۳
۴۸.....	برآورد تعداد عامل های پنهان	۶-۳
۴۹.....	ارزیابی روش های خوشه بندی	۷-۳
۵۲.....	جمع بندی	۸-۳

فصل چهارم: روش پیشنهادی برای خوشه بندی کاربران وب براساس تحلیل معنایی پنهان احتمالاتی

۵۳.....	مقدمه	۱-۴
۵۴.....	خوشه بندی کاربران وب براساس آنالیز احتمالاتی معنایی پنهان	۲-۴
۵۴.....	پیش پردازش داده	۱-۲-۴
۵۵.....	مدل سازی کاربر براساس تحلیل معنایی پنهان احتمالاتی	۲-۲-۴
۵۵.....	مجموعه گردشگری کاربران شخصی	
۵۶.....	تقسیم بندی URLها	
۵۷.....	اعمال تحلیل معنایی پنهان احتمالاتی با استفاده از توابع وزن دهی مختلف برای هر کاربر	۳-۴
۶۰.....	خوشه بندی الگوهای کاربری ایجاد شده	۴-۴
۶۲.....	ایجاد پروفایل عمومی کاربران	۵-۴
۶۳.....	معرفی روش ها به منظور مقایسه	۶-۴
۶۳.....	جمع بندی	۷-۴

فصل پنجم: ارزیابی روش پیشنهادی و کاربرد آن

۶۵.....	مقدمه	۱-۵
۶۵.....	ارزیابی نتایج	۲-۵
۶۵.....	پیش پردازش و منبع داده	۱-۲-۵
۶۶.....	بررسی تعداد عامل های پنهان با استفاده از توابع وزنی مختلف	۱-۲-۵
۶۸.....	اندازه گیری صحت خوشه بندی	۲-۲-۵
۷۲.....	ایجاد پروفایل های عمومی کاربر	۳-۲-۵

۷۳.....	کاربرد : واكشی اولیه صفحات	۳-۵
۷۴.....	قوانین واكشی اولیه	۱-۳-۵
۷۵.....	نتایج	۲-۳-۵
۷۶.....	جمع بندی	۴-۵
فصل ششم: نتیجه گیری		
۸۱.....	مقدمه	۱-۶
۸۱.....	ارزیابی الگوریتم پیشنهادی	۲-۶
۸۲.....	پیشنهاداتی برای کارهای آینده	۳-۶
۸۲.....	جمع بندی	۴-۶
۸۳.....	مراجع	

چکیده

در سالهای اخیر با توجه به رشد سریع وب جهانی، تحقیقات وسیعی پیرامون مدل کردن رفتار پیمایشی کاربران در وبسایتها انجام گرفته است. در این راستا کاربرد کاوی وب با هدف به دست آوردن الگوهای رفتار پیمایشی کاربران وب، مورد استفاده بسیاری از محققان قرار گرفته است.

به طور کلی کاربران وب رفتارهای متفاوتی متناسب با نیازهای اطلاعاتی و وظایف مورد علاقه خود در بازدید از وبسایتها از خود نشان می دهند، تمامی رفتارهای پیمایشی کاربران در فایل های ثبت وب قابل ردیابی است. یکی از تکنیک های مورد استفاده در کاربرد کاوی وب خوشه بندی کاربران وب می باشد. در تکنیک خوشه بندی کاربرانی که رفتار پیمایشی مشابهی دارند در یک خوشه قرار می گیرند. هر خوشه منجر به ایجاد پروفایل های کاربری می شود که در برنامه های کاربردی مانند واکنشی اولیه و حافظه نهان مورد استفاده قرار می گیرد. تکنیک های متداول و استاندارد کاربرد کاوی وب برای خوشه بندی کاربران وب می تواند الگوهای کاربردی را مستقیماً کشف کند، اما این تکنیک ها به طور خودکار نمی توانند مشخصات یا کیفیت عامل های پنهانی که منجر به کشف الگوهای پیمایشی مشترک می شوند را تعیین کنند. بنابراین نیاز به گسترش تکنیک هایی می باشد تا بتوان بصورت خودکار اهداف اساسی پیمایشی کاربران را شناسایی و رابطه معنایی پنهان میان کاربران وب و همچنین رابطه معنایی پنهان بین کاربران وب و اشیاء وب را استخراج کرد.

در این پژوهش، روشی براساس آنالیز معنایی پنهان احتمالاتی پیشنهاد می شود که مشخصات ذاتی رفتار پیمایشی کاربران را مشخص می کند. روش پیشنهادی فاکتورهای پنهان به دست آمده را جهت خوشه بندی الگوهای پیمایشی کاربران مورد استفاده قرار می دهد و پروفایل های کاربری را ایجاد می کند. نتایج خوشه بندی برای پیش بینی و واکنشی درخواست های وب گروه های کاربران مورد استفاده و ارزیابی قرار می گیرد. کارایی و برتری روش خوشه بندی کاربران از طریق آزمایش بر روی فایل های ثبت واقعی نشان داده می شود. روش پیشنهادی برای گروه بندی کاربران وب و واکنشی اولیه با کارهای قبلی مقایسه شده و نتایج به دست آمده کارایی بهتر و میزان دقت بالای روش پیشنهادی را در مقایسه با روش های دیگر از خود نشان می دهد.

کلمات کلیدی: کاربرد کاوی وب ، خوشه بندی کاربران وب ، آنالیز معنایی پنهان احتمالاتی، واکنشی اولیه صفحات

فصل اول

مقدمه

۱-۱ انگیزه و پیش‌زمینه‌ها

در دنیای امروز، وب رسانه‌ای محاوره‌ای و محبوب برای انتشار اطلاعات می‌باشد. در نتیجه، وب به بزرگترین بازار بدون مرز برای ارائه خدمات و اطلاعات و یا خرید و فروش کالا تبدیل شده است. وجود انبوه عرضه‌کنندگان سبب شده است که شرکت‌ها بر سر یک یک کاربرانی (مشتریانی) که از وب‌سایت آنها دیدار می‌کنند و شاید در مواقعی کالا و خدمات آنها را خریداری می‌کنند، با یکدیگر رقابت کنند. راز پیروزی در چنین رقابت سرسختانه‌ای اطلاع از نیازهای کاربران وب‌سایت و ارائه خدماتی که بتواند نیازها را برآورده سازد. در بیشتر موارد طراحی یک سایت براساس درک و تصور طراح از الگوهای کاربردی کاربران از وب‌سایت صورت می‌پذیرد. بنابراین اطلاع از الگوهای کاربردی کاربران منجر به انطباق هر چه بهتر خدمات ارائه شده توسط وب‌سایت مطابق با نیازهای کاربران و به تبع آن افزایش احتمال جذب کاربر می‌گردد. از سوی دیگر کاربران وب به دلیل رشد سریع اطلاعات و تعداد زیاد کاربران با مشکل سربار اطلاعات مواجه هستند. در نتیجه اینکه چگونه اطلاعات مورد نیاز کاربران در اختیار آنها قرار داده شود نیز به موضوعی بحرانی در بازیابی اطلاعات و مدیریت داده‌ها براساس وب تبدیل شده است.

راه‌حل‌های فراوانی براساس راه‌حل‌های تکنیکی برای حل مشکلات بیان شده ارائه شده است. در میان راه‌حل‌های ارائه شده، وب کاوی^۱ به عنوان راه‌حلی کارآمد برای کشف ساختار محتوای وب، رفتار گردشگری کاربران براساس علاقه و تعاملات اساسی بین کاربران وب و اشیاء وب مورد استفاده قرار می‌گیرد.

وب کاوی فرایند کشف ارتباط ذاتی میان داده‌های وب می‌باشد، داده‌های وب به صورت اطلاعات متنی، پیوندی^۲ و کاربردی است که از طریق آنالیز ویژگی‌های داده‌های وب و با استفاده از تکنیک‌های داده کاوی الگوهای جالبی را می‌توان کشف و استخراج کرد. با توجه به داده‌های مورد استفاده، وب کاوی شامل محتوا کاوی وب، ساختار کاوی وب و کاربرد کاوی وب می‌باشد. در این پروژه، تمرکز اصلی بر روی کاربرد کاوی وب می‌باشد. همانطور که در فصل دو توضیح داده می‌شود، کاربرد کاوی وب بکارگیری روش‌های داده کاوی بر روی داده‌های کاربردی که نحوه استفاده کاربران از وب‌سایت را بیان می‌کند، برای کشف دانش از رفتار پیمایشی کاربران می‌باشد. گام‌های اصلی در کاربرد کاوی وب پیش‌پردازش داده‌ها، کشف الگو و آنالیز الگو می‌باشد. در مرحله پیش‌پردازش، داده‌های جمع‌آوری شده به فرم داده مورد نیاز برای مرحله کشف الگو تبدیل می‌شود. در مرحله کشف الگو از تکنیک‌های مختلف داده کاوی مانند خوشه‌بندی، قوانین انجمنی، تحلیل آماری و الگوهای دنباله‌ای برای کشف الگوهای رفتاری کاربران و بدست آوردن اولویت‌ها و علائق کاربر استفاده می‌شود. در مرحله آنالیز الگو، الگوهای نامربوط از میان الگوهای بدست آمده حذف

¹ Web Mining

² Linkage

می‌گردد. الگوهای بدست آمده از فرایند کاربرد کاوی می‌تواند در حیطه‌های مختلف مانند شخصی‌سازی، تجارت الکترونیکی، بهبود وب‌سایت و واکنشی اولیه و حافظه نهان بکار گرفته شود. در سال‌های اخیر، تحقیقات در حوزه وب‌کاوی و به خصوص در زمینه کاربرد کاوی وب گسترش یافته است. از اواسط ۱۹۹۰ تا به حال هزاران مقاله در زمینه وب‌کاوی منتشر شده است که حدوداً پنجاه درصد از این مقالات درباره‌ی کاربرد کاوی وب بوده است.

هنگام مواجهه با مسئله‌ای در کاربرد کاوی وب، مشکل اصلی انتخاب مناسب‌ترین تکنیک برای مسئله مورد بحث است. مانند همه زمینه‌های دیگر، برای مسائل کاربرد کاوی وب پاسخ یکتا وجود ندارد. در مطالعات صورت گرفته، از میان تکنیک‌های مختلف داده‌کاوی، تمرکز اصلی بر تکنیک خوشه‌بندی می‌باشد. در کاربرد کاوی وب خوشه‌بندی برای گروه‌بندی اشیائی که مشخصات یکسان دارند، مورد استفاده قرار می‌گیرد. خوشه‌بندی در کاربرد کاوی وب به خوشه‌بندی صفحات وب و خوشه‌بندی کاربران وب تقسیم می‌شود. خوشه‌بندی صفحات وب صفحات وب اصول ذاتی صفحات وب را انعکاس می‌دهد، در حالیکه خوشه‌بندی کاربران وب علاقه‌گردشگری کاربران را نشان می‌دهد. هدف اصلی در این پژوهش مدل کردن رفتار گردشگری کاربران براساس خوشه‌بندی کاربران علاقه‌مند و نمایش فضای وظایف گردشگری در سطح معنایی، از طریق وب‌کاوی و تحلیل معنایی پنهان احتمالاتی می‌باشد.

در خوشه‌بندی کاربران وب، کاربرانی که هنگام گردشگری در وب‌سایت رفتارهای مشابهی انجام می‌دهند در یک گروه قرار می‌گیرند. هدف نهایی خوشه‌بندی کاربران همانند دیگر تکنیک‌های داده‌کاوی، فراهم کردن قابلیت تحلیل خوشه‌ها به منظور استخراج هوش تجاری یا استفاده از آن‌ها برای اعمالی نظیر شخصی‌سازی وب و بهبود وب‌سایت‌ها و واکنشی اولیه صفحات و غیره می‌باشد. بسیاری از مطالعات صورت گرفته در زمینه خوشه‌بندی کاربران وب، برای مدل کردن و گروه‌بندی کاربران وب از تکنیک‌های استاندارد خوشه‌بندی استفاده می‌کند. اما الگوهای کاربردی بدست آمده از این تکنیک‌ها، الگوهایی در سطح صفحات وب می‌باشند. این الگوها نه تنها مشخصات ذاتی فعالیت‌های کاربران را بدست نمی‌آورند، بلکه فاکتورهای اساسی و غیرقابل مشاهده‌ای را که منجر به الگوهای گردشگری خاص می‌شود تعیین نمی‌کنند.

بنابراین برای درک بهتر عواملی که منجر به الگوهای گردشگری می‌شود نیاز به توسعه و گسترش تکنیک‌ها می‌باشد، بطوریکه تکنیک‌های جدید بتوانند بصورت خودکار هدف‌های اصلی گردشگری کاربران را مشخص کند و رابطه معنایی میان کاربران و همچنین میان اشیاء وب را تعیین کنند. مدل‌های متغیر پنهان، همانند تحلیل معنایی پنهان احتمالاتی^۱ بصورت گسترده‌ای برای کشف رابطه پنهان میان هم‌رخدادی اشیاء بکار می‌رود. تحلیل معنایی پنهان احتمالاتی روش احتمالاتی را برای کشف متغیرهای پنهان ارائه می‌دهد که روشی بسیار انعطاف‌پذیر و دارای پایه آماری بسیار قوی می‌باشد. اساس تکنیک تحلیل معنایی پنهان احتمالاتی مدلی به نام مدل aspect می‌باشد. فرض کنید مجموعه‌ای از عامل‌های پنهان وجود دارد که اساس هم‌رخدادی میان دو مجموعه از اشیاء می‌باشد، تکنیک تحلیل معنایی پنهان احتمالاتی از الگوریتم ماکزیمم انتظار وقوع برای برآورد مقادیر احتمالاتی که رابطه بین عامل‌های پنهان و دو مجموعه از اشیاء اندازه‌گیری می‌کند، استفاده می‌کند. دانش پس‌زمینه در تکنیک تحلیل معنایی پنهان احتمالاتی مفهوم عامل پنهان می‌باشد که فرض می‌کند عامل‌های پنهان قطعیت هم‌رخدادی کاربر و صفحات را کنترل کند. تحلیل معنایی پنهان بطور

¹ Probabilistic Semantic Latent Analysis

موفقیت آمیزی در کاربرد کاوی وب معرفی شده است و موفقیت بزرگی را در مطالعات مربوطه مانند مدلسازی کاربر، سیستم‌های پیشنهاددهنده و غیره بدست آورده است.

در این پژوهش، کاربرد کاوی وب براساس تحلیل معنایی پنهان احتمالاتی ارائه می‌شود. در کاربرد کاوی وب داده‌های هم‌رخداد کاربران وب و اشیاء وب می‌باشند. در این پروژه فاکتورهای پنهان که رابطه میان موجودیت‌ها را نشان می‌دهد به عنوان وظیفه^۱ شناخته می‌شود. این مطلب برای تایید این حقیقت است که فاکتورهای پنهان اهداف گردشگری کاربران در وب‌سایت‌ها و همچنین تعاملاتشان با اشیاء وب را نشان می‌دهد. با اعمال تکنیک تحلیل معنایی پنهان احتمالاتی، می‌توان بطور موثر مشخصات عامل‌های پنهان را شناسایی کرد، سپس از لحاظ کمی رابطه بین کاربران وب و وظایف همچنین رابطه بین اشیاء وب و کاربران وب را اندازه‌گیری کرد. این روابط بصورت احتمالاتی اندازه‌گیری می‌شود که به نوبه خود به کشف الگوهای کاربردی مختلف با استفاده از استنتاج احتمالاتی منجر می‌شود.

۲-۱ هدف پژوهش

در این پروژه مطالعات صورت گرفته بر کشف الگوهای کاربردی از طریق کاربرد کاوی وب و سپس استفاده از الگوهای کشف شده برای ایجاد گروه‌های کاربران علاقه‌مند و استفاده از گروه‌های بدست آمده برای کاربرد واکنشی اولیه صفحات می‌باشد. هدف اصلی در این پژوهش در روشی که از ساختار سلسله مراتبی وب‌سایت بجای صفحات وب استفاده می‌شود، بهبود در مدل کاربر و در نتیجه بهبود در گروه‌بندی کاربران وب می‌باشد که از تکنیک تحلیل معنایی پنهان احتمالاتی در کاربرد کاوی وب استفاده می‌کند. همچنین کاربرد گروه‌های ایجاد شده برای واکنشی اولیه صفحات ارائه می‌شود.

۳-۱ ساختار پژوهش

ساختار پژوهش به صورت زیر می‌باشد. فصل دو به ارائه‌ی اطلاعات پیش‌زمینه‌ای لازم در مورد کاربرد کاوی وب تخصیص پیدا می‌کند. مرور بر ادبیاتی که در این فصل آمده است برای دادن اطلاعات به خواننده برای درک موضوع است.

فصل دو به جزئیات وب کاوی و بطور اخص، کاربرد کاوی وب می‌پردازد. با استفاده از رویکرد بالا به پایین از وب کاوی شروع کرده و سپس بحث را اختصاصی کرده و به کاربرد کاوی وب پرداخته می‌شود.

فصل سه، به معرفی تکنیک خوشه‌بندی کاربران وب پرداخته می‌شود. این فصل شامل دو بخش است، در بخش اول کارهای انجام در زمینه خوشه‌بندی کاربران وب با استفاده از تکنیک‌های داده کاوی بررسی می‌شود. در بخش دوم، کارهای صورت گرفته با استفاده از تکنیک‌های آنالیز معنایی در زمینه خوشه‌بندی کاربران وب معرفی و بررسی می‌شود. فصل چهارم به ارائه روش پیشنهادی برای خوشه‌بندی کاربران وب بر اساس تحلیل معنایی پنهان احتمالاتی پرداخته می‌شود که فصل اصلی این در پژوهش می‌باشد. بخش بعد جزئیات مربوط به پیاده سازی، بررسی نتایج آزمایشات و مقایسه و تحلیل آنها با سایر روش‌ها مطرح می‌شود.

¹ Task

فصل پنجم به ارزیابی روش پیشنهادی اختصاص دارد. در این فصل نتایج آزمایشاتی که برای ارزیابی سیستم انجام شده است، ارائه خواهد شد و رفتار الگوریتم تفسیر می‌شود. در بخش دوم این فصل نیز الگوریتم پیشنهادی برای وظیفه حقیقی واکنشی اولیه صفحات اعمال می‌شود تا درخواست‌های آینده کاربران گروه‌بندی شده بر طبق صفحات عمومی مشاهده شود توسط کاربران، پیش‌بینی شود و نتایج بدست آمده با روش‌های دیگر مقایسه شده می‌شود. و نهایتاً فصل ششم که به جمع‌بندی و خلاصه کردن نتایج فصل‌های پیشین اختصاص دارد، و دیدی کلی بر مسئله ایجاد می‌کند.

فصل دوم

استخراج اطلاعات از داده‌های کاربرد وب

۱-۲ مقدمه

در این فصل به تعریف مفاهیم در زمینه وب کاوی خواهیم پرداخت، بعد از ارائه دیدی کلی بر مسئله‌ی وب کاوی، به کاربرد کاوی وب خواهیم پرداخت. روش‌های موجود برای این نوع کاوش در ادبیات این مسئله، در این فصل عنوان شده و کاربردهای آن ذکر می‌شود.

۲-۲ داده کاوی و استخراج دانش

با گسترش سیستم‌های پایگاه داده‌ای و حجم بالای داده‌های ذخیره شده در این سیستم‌ها، برای اینکه بتوان داده‌های ذخیره شده را پردازش کرد نیاز به ابزارهایی است که بتواند اطلاعات حاصل از این پردازش را در اختیار کاربران قرار دهد. با استفاده از پرس و جوهای ساده در SQL و ابزارهای گوناگون گزارش‌گیری معمولی، می‌توان اطلاعاتی را در اختیار کاربران قرارداد تا با استفاده از آن بتوانند به نتیجه‌گیری در مورد داده‌ها و روابط منطقی میان آنها پردازند. اما وقتی که حجم داده‌ها بالا باشد کاربران هرچند زبردست و مجرب باشند نمی‌توانند الگوهای مفید را در میان حجم انبوه داده‌ها تشخیص دهند و حتی اگر قادر به اینکار هم باشند هزینه عملیات از نظر نیروی انسانی و مادی بسیار بالا است.

علاوه بر این، با وجود سیستم‌های یکپارچه اطلاعاتی، سیستم‌های یکپارچه بانکی و تجارت الکترونیک، لحظه به لحظه به حجم داده‌ها به پایگاه داده‌های مربوطه اضافه شده و باعث بوجود آمدن انبارهای عظیمی از داده‌ها شده‌است بطوری که ضرورت کشف و استخراج سریع و دقیق دانش از این پایگاه داده‌ها را بیش از پیش نمایان کرده‌است.

داده کاوی بعنوان یک راه حل برای این مسائل مطرح می‌باشد در متون علمی تعاریف گوناگونی برای داده کاوی ارائه شده است برخی از این تعاریف عبارتند از:

- داده کاوی عبارت است از فرایند استخراج اطلاعات معتبر، از پیش ناشناخته، قابل فهم و قابل اعتماد از پایگاه داده‌های بزرگ و استفاده از آن در تصمیم‌گیری در فعالیتهای تجاری مهم.
- اصطلاح داده کاوی به فرایند نیمه خودکار تجزیه و تحلیل پایگاه داده‌های بزرگ به منظور یافتن الگوهای مفید اطلاق می‌شود.
- داده کاوی یعنی جستجو در پایگاه داده‌ها برای یافتن الگوهای میان داده‌ها.
- داده کاوی یعنی تجزیه و تحلیل مجموعه داده‌های قابل مشاهده برای یافتن روابط مطمئن بین داده‌ها.

کاوش اطلاعات، حجم عظیمی از داده‌های خام را به فرمی تغییر می‌دهد که انسان بتواند آنها را به راحتی بفهمد و برای تصمیم‌گیری بتواند از این اطلاعات استفاده کند. در کل، داده‌کاوی به دو شاخه‌ی مرتبط مدل‌کردن و کشف الگو تقسیم می‌شود. مدل‌کردن خلاصه‌کردن ساختارهای عظیم داده‌است. مثال‌هایی از چنین مدل‌کردنی عبارتند از:

- تحلیل جزءبندی
- تحلیل رگرسیون
- تجزیه‌ی سری‌های زمانی

و فرآیند کشف دانش، می‌تواند در قدم‌های زیر خلاصه شود [۱]:

- پاک‌سازی داده
- انتخاب و تغییر شکل داده
- کاوش داده
- ارزیابی
- ارائه‌ی دانش

دو قدم اول، فاز پیش‌پردازش و دو قدم آخر فاز پس‌پردازش هستند. قدم کاوش داده، قسمتی از فرآیند کشف دانش است که از داده‌های خام دانش استخراج می‌کند. در این قدم از استراتژی‌های یادگیری ماشین، آمار و ... برای حل مسائل داده‌کاوی استفاده می‌شود. دانش استخراج شده باید دارای صفات دقیق، قابل درک، جالب باشد [۲].

البته درجه‌ی اهمیت هر کدام از این صفات بسته به کاربر و کاربرد دارد. برای مثال اگر قرار باشد کاربر از دانش استخراج شده استفاده کند و تصمیم‌گیری‌های استراتژیک توسط کاربر انجام شود، قابل درک بودن دانش استخراج شده اهمیت بیشتری پیدا می‌کند. و همچنین اگر دانش استخراج شده قابل فهم نباشد، هدف تحقق پیدا نخواهد کرد.

در مسائل داده‌کاوی، هر چه حجم داده‌ها بیشتر می‌شود، میل بیشتری برای کشف الگوهای مخفی در داده‌ها برای به دست آوردن سود تجاری و ... به وجود می‌آید. در قدم اصلی داده‌کاوی ممکن است از چندین الگوریتم داده‌کاوی استفاده شود. کار اصلی الگوریتم داده‌کاوی با توجه به نوع مسئله‌ی کشف دانش تغییر می‌کند.

دو نوع اصلی الگوریتم‌های داده‌کاوی، کلاس‌بندی و خوشه‌بندی می‌باشد. خوشه‌بندی به فرایند تقسیم‌بندی داده به یک یا چند گروه به طوری که فاصله بین خوشه‌ها حداکثر و فاصله درون خوشه‌ها حداقل باشد، اطلاق می‌شود. با توجه تعریف خوشه‌بندی، در این پروژه الگوریتم مورد استفاده و پیشنهادی در دسته‌ی خوشه‌بندی قرار می‌گیرد، چرا که در نهایت کاربران وب به تعداد مشخصی گروه تقسیم‌بندی می‌شوند به این ترتیب که کاربرانی که در مواجهه با وب‌سایت رفتار گردشگری مشابهی دارند، در یک گروه قرار می‌گیرند.

وب‌سازان برای رسیدن به اطلاعات زیاد و متنوعی که در سایت‌های مختلف در شبکه گسترده جهانی ذخیره شده‌است. رشد گسترده وب مسبب رشد ثابت اطلاعات شده‌است که چندین مشکل را به همراه دارد، بدست آوردن اطلاعات مرتبط، استخراج دانش مفید و یادگیری در مورد عملکردهای مشتریان یا کاربران، نمونه‌هایی از این مشکلات

است. بنابراین در داده کاوی حوزه‌ای به نام وب کاوی بوجود آمده است که سعی بر حل این مشکلات دارد. در بخش بعدی در مورد وب کاوی و انواع آن بحث می‌شود.

۳-۲ وب کاوی^۱

وب، اکنون بزرگترین انبار داده‌ی بشر است. ایجاد وب برمی‌گردد به سال ۱۹۹۰ که یک کارمند CERN اولین برنامه برای حرکت روی اتصالات دو طرفه بین مجموعه‌های از اسناد را با رابط کاربر گرافیکی^۲ نوشت. نام این نرم‌افزار را وب گسترده‌ی جهانی^۳ گذاشتند. اگرچه عمل هدایت بین اسناد به این طریق، جدید نبود ولی اضافه شدن رابط کاربر کار جدیدی بود. این رابط کاربر جدید با ابداع دیگری در CERN به هم آمیخت. این ابداع HTTP^۴ بود. در سال ۱۹۹۳ حجم ترافیک HTTP ده برابر شد [۳] و بعد از آن وب گسترده‌ی جهانی از اسم یک نرم‌افزار تبدیل شد به یک وب واقعی که با ارتباطات بین اسناد ابرمتن به وجود آمده‌است.

وب کاوی استفاده از تکنیک‌های داده کاوی برای کشف و استخراج خودکار اطلاعات از مستندات و سرویس‌های وب می‌باشد. در واقع وب کاوی به تمام فرایند کشف اطلاعات گفته می‌شود و نه صرفاً نرم‌افزارهای کاربردی که از ابزارهای استاندارد داده کاوی استفاده می‌کنند. برخی از صاحب‌نظران در این حوزه وظایف وب کاوی را به چهار وظیفه زیر تقسیم می‌کنند که در ادامه آورده می‌شود:

- پیدا کردن منبع: یعنی فرایند بازیابی داده‌ها از اسناد وب مورد نظر.
 - انتخاب و پیش‌پردازش اطلاعات: در این مرحله به صورت خودکار اطلاعات خاصی از اسناد بازیابی شده، انتخاب و پیش‌پردازش می‌شوند.
 - تعمیم [۴]: در این مرحله به صورت خودکار الگوهای عام در یک یا چندین سایت وب کشف می‌شود. برای این عمل تکنیک‌های مختلف داده کاوی، یادگیری ماشین و متدهای مبتنی بر وب استفاده می‌شود.
 - تحلیل: در این مرحله الگوهای به دست آمده در مرحله قبل اعتبار سنجی [۵] و تفسیر می‌شوند.
- از طرفی عواملی که بر درک و ارزیابی کاربر از صفحات وب در هنگام فرایند داده کاوی تاثیر می‌گذارد عبارتند از:

الف) محتوای صفحه وب

ب) طراحی صفحه وب

ج) طراحی کل سایت که شامل ساختار آن می‌باشد.

عامل اول در ارتباط با داده‌هایی است که توسط یک سایت ارائه می‌شود. عامل دوم در ارتباط با روشی است که محتویات سایت برای کاربران در دسترس و قابل درک می‌شود. توجه شود که میان طراحی یک صفحه‌ی مستقل و طراحی کل یک سایت فرق وجود دارد، زیرا یک سایت شبکه‌ای است از صفحات مرتبط است و تا وقتی که طراحی یک سایت ساده و قابل درک نباشد، کاربران تمایل به استفاده از وب سایت ندارند.

¹ Web mining

² GUI

³ World Wide Web

⁴ HyperText Transfer Protocol

۴-۲ چالش‌های وب‌کاوی

وب‌کاوی با چالش‌ها و محدودیت‌های متنوعی روبه‌رو است. از یک دیدگاه می‌توان این محدودیت‌ها را به دو گروه تکنیکی و غیرتکنیکی تقسیم کرد. از محدودیت‌های غیرتکنیکی می‌توان به عدم پشتیبانی مدیریت، کافی نبودن بودجه و عدم وجود منابع مورد نیاز مانند نیروی انسانی متخصص اشاره کرد. اما مشکلات تکنیکی بسیار است که به برخی از آنها در این جا اشاره می‌شود:

۱. **داده‌های ناصحیح و نادقیق:** برای آن که فرآیند وب‌کاوی با موفقیت انجام شود، لازم است داده‌های جمع‌آوری شده صحیح و در قالب مناسب باشند. اما معمولاً مشکلات زیادی در این زمینه وجود دارد، چرا که ممکن است داده‌ها دقیق نباشند و از طرفی داده‌ها می‌توانند ناکامل بوده و برخی مقادیر موجود نباشد. همچنین تخمین میزان اطمینان درباره صحت و دقت داده‌ها به سادگی امکان‌پذیر نیست.
۲. **عدم وجود ابزارها:** محدودیت دیگر وب‌کاوی، عدم وجود ابزارهای مناسب و کامل برای آن می‌باشد. متخصصان باید تصمیم بگیرند آیا برای یک کاربرد وب‌کاوی، ابزار خاص آن را توسعه دهند و یا از ابزارهای موجود استفاده کنند.
۳. **ابزارهای سفارشی:** ابزارهای موجود تنها یکی از انواع وب‌کاوی مانند طبقه‌بندی یا خوشه‌بندی را پشتیبانی می‌کنند. اما بهتر آن است که یک ابزار قادر به انجام چندین تکنیک وب‌کاوی باشد تا کاربران بتوانند با توجه به نیازمندی‌های خود از تکنیک مناسب استفاده کنند.

البته در حال حاضر تحقیقات بسیاری در زمینه وب‌کاوی در حال انجام است که هدف آن‌ها حل این مشکلات می‌باشد.

۵-۲ اجزای اصلی وب‌کاوی

همان‌طور که گفته شد وب منبع عظیمی از داده است. این داده یا محتویات وب است که از طریق میلیون‌ها صفحه‌ی وب به دست می‌آید و یا اطلاعات اتصالات بین صفحات وب و یا داده‌های دسترسی روزانه کاربران به صفحات وب که در سرورهای وب در فایل‌های متنی به نام ثبت^۱ ذخیره می‌شوند. براساس داده‌های مورد استفاده در کاوش وب، سه شاخه‌ی مجزا در وب‌کاوی به وجود آمده‌است. این سه شاخه، عبارتند از:

- ۱- محتوا کاوی وب^۲
- ۲- ساختار کاوی وب^۳
- ۳- کاربرد کاوی وب^۴

^۱ Log

^۲ Web content mining

^۳ Wes structure mining

^۴ Web usage mining

۲-۵-۱ محتوا کاوی وب

استفاده از تکنیک‌های داده کاوی به منظور استخراج اطلاعات مفید از محتوای صفحات وب، محتوا کاوی وب نامیده می‌شود. این اطلاعات عبارتند از متن، تصویر، ویدئو و ... همچنین استخراج منابع وب، طبقه‌بندی مستندات، دسته‌بندی و استخراج اطلاعات از صفحات وب بخشی از وظایف محتوا کاوی وب است. در اینترنت به دلیل عدم وجود ساختاری برای مستندات، عملیات مکانیزه کردن استخراج، سازماندهی و مدیریت اطلاعات وب مشکل است. کاربردهای معمول این نوع کاوش وب، سازماندهی وب بر اساس محتوا و درجه بندی^۱ صفحات وب بر اساس محتوا است.

محتوا کاوی وب، از جهاتی شبیه به داده کاوی و متن کاوی است. به داده کاوی شبیه است زیرا تکنیک‌های داده کاوی زیادی در مسئله‌ی محتوا کاوی وب استفاده می‌شوند و به متن کاوی شبیه است زیرا اکثر محتویات وب، متن هستند. ولی از جهاتی هم با داده کاوی متفاوت است زیرا داده‌های وب اکثراً نیمه ساختاریافته هستند در حالیکه داده کاوی با داده‌های ساختاریافته سروکار دارد. با متن کاوی نیز از آن جهت متفاوت است که متن موجود در وب نیمه ساخت یافته است ولی متن کاوی با داده‌های متنی بدون ساختار سروکار دارد.

۲-۵-۲ ساختار کاوی وب

ساختار کاوی وب ساختار فرامتن را پیدا می‌کند. در واقع ساختار کاوی وب به کشف مدل پس زمینه حاکم بر ساختار فرایوندهای وب می‌پردازد و هدف آن ایجاد اطلاعاتی همچون تشابه یا ارتباط بین سایت‌های مختلف وب است. ساختار کاوی وب درگیر تحلیل لینک‌های داخلی و خارجی صفحات وب است و برای رتبه‌بندی موتور جستجو به کار می‌رود. این بخش سعی در کشف ساختارهای لینک‌های پشت صفحات وب می‌کند که اطلاعات بیشتری از مستندات مفید را بیان می‌کنند. همچنین ساختار کاوی وب، عمل کشف اعتبار سایتها، برای موضوعات و بازنگری سایت‌هایی است که به چندین شخص حقیقی یا حقوقی وابسته هستند. لینک‌های درون یک صفحه، عمومی بودن یک صفحه را مشخص می‌کند در حالی که لینک‌های بیرونی یک صفحه توانگری صفحه را نشان می‌دهد و یا یک صفحه‌ای که بصورت دائم تکرار می‌شود می‌تواند یک صفحه مهم باشد. بنابراین می‌توان از میان ساختارهای لینک، وب را کاوش کرد. Page Rank، CLEVER مثالهایی برای ساختار کاوی وب است.

۲-۵-۳ کاربرد کاوی وب

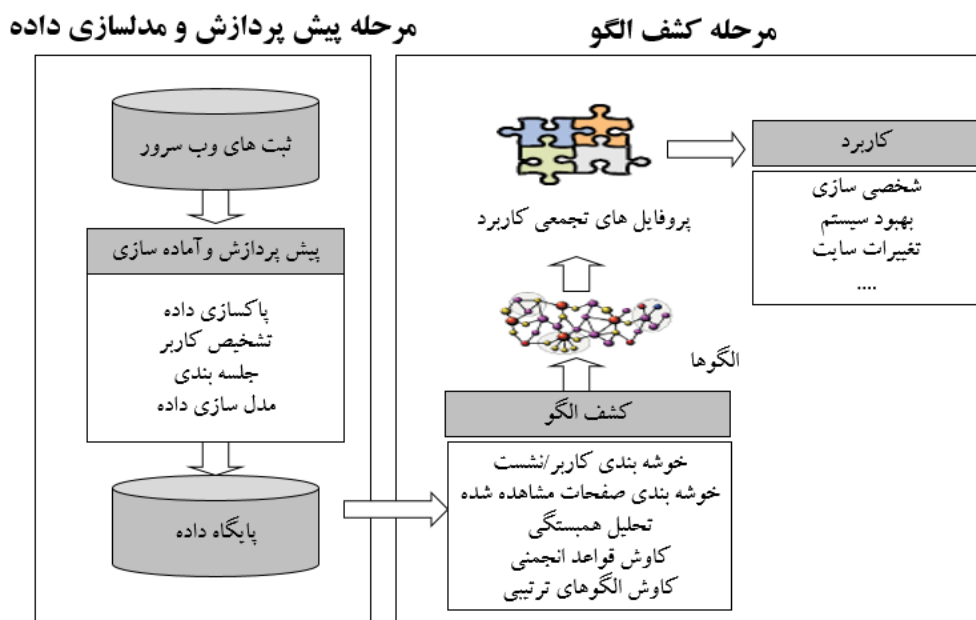
کاربرد کاوی وب بخشی از وب کاوی است که استخراج مکانیزه الگوی دسترسی کاربر از یک یا چند سرویس دهنده وب را بر عهده دارد. کاربرد کاوی وب سعی می‌کند از تعاملات کاربر با وب اطلاعاتی کسب کند و از آنها بصورت سابقه‌ای در مراجعات بعدی کاربر سود ببرد. داده‌های مورد استفاده وب شامل داده‌های ثبت شده، User Session، Browser Log، User Profile، Proxy Server Log، Web Server usage Log، User Query، Cookies، User، Bookmark Log، Transaction و ... می‌شوند. سازمان‌ها اغلب در عملیات روزانه‌شان حجم زیادی از اطلاعات را جمع‌آوری و تولید می‌کنند. بیشتر این اطلاعات معمولاً بطور خودکار توسط سرویس دهنده‌های وب تولید و توسعه Server Access Log، جمع‌آوری می‌شوند. سایر اطلاعات کاربر شامل Refer Logها است که شامل اطلاعاتی درباره صفحات استنادی برای منشأ صفحه و ثبت نام کاربر با ابزارهای کلی جمع‌آوری داده مانند CGI است. تحلیل Server

^۱ Ranking

Access Log و اطلاعات ثبت شده کاربران می توانند اطلاعات ارزشمندی در مورد اینکه چگونه بایستی یک وب سایت را در یک سازمان ایجاد کرد فراهم کنند.

از آنجایی که موضوع اصلی پروژه کاربرد کاوی وب می باشد، در ادامه این بخش به ارائه مسائل مرتبط با این نوع کاوش وب پرداخته می شود. فرآیند کلی کاربرد کاوی وب شامل سه مرحله است: آماده سازی و مدلسازی داده، کشف الگو از داده های کاربرد وب و استفاده از الگوهای کشف شده برای کاربردهایی مانند شخصی سازی وب، واکنشی اولیه صفحات، تجارت الکترونیک. در مرحله آماده سازی داده، پیش پردازش های لازم بر روی ثبت های جمع آوری شده، بیان می شود. ثبت های وب برای مقاصد مختلف حذف می شوند، مثلا برای مرتب کردن داده های غیر لازم (مثلا دسترسی از طریق web spider)، برای تعیین نشست های کاربران^۱ (به وسیله کوکی^۲ها)، برای ذخیره ی داده ها در یک پایگاه داده ی رابطه ای و در نهایت ثبت های خام وب به فرمتی تبدیل می شود که می تواند در داده کاوی مورد استفاده قرار گیرد. در مرحله کشف الگو وظیفه های مختلف داده کاوی مانند خوشه بندی، کاوش قواعد انجمنی و کشف الگوهای ترتیبی را می توان بر روی داده ی پیش پردازش شده اجرا کرد. نتایج فاز کاوش به پروفایل های تجمعی کاربرد تبدیل می شوند که پروفایل های بدست آمده در این مرحله برای کاربردهایی مانند واکنشی اولیه صفحات، شخصی سازی وب، تغییر در وب سایت و تجارت الکترونیک.

در شکل ۱-۲ یک نمای جامع از فرآیند کاربرد کاوی وب ارائه داده می شود. در ادامه ی این بخش از این چارچوب نشان داده شده از شکل ۱-۲ به عنوان راهنما استفاده خواهد شد. و در ادامه بحث مفصلی از فعالیت های کاربرد کاوی وب لازم برای این فرآیند شامل پیش پردازش، تکنیک های متداول کشف الگو که بر روی این داده های کاربرد بکار می روند ارائه خواهد شد.



شکل ۱-۲ فرآیند کاربرد کاوی وب [۸]

^۱ User sessions

^۲ Cookie

۶-۲ آماده‌سازی و مدل‌سازی داده

یک مرحله‌ی مهم در هر عملیات داده‌کاوی ایجاد مجموعه داده‌ی مناسبی است که بر روی آن الگوریتم‌های داده‌کاوی اجرا می‌شوند. این فرایند می‌تواند شامل پیش‌پردازش داده‌های اولیه، تبدیل داده‌ها به فرمتی مناسب به عنوان ورودی عملیات داده‌کاوی، باشد. به مجموعه‌ی این عملیات آماده‌سازی داده گفته می‌شود.

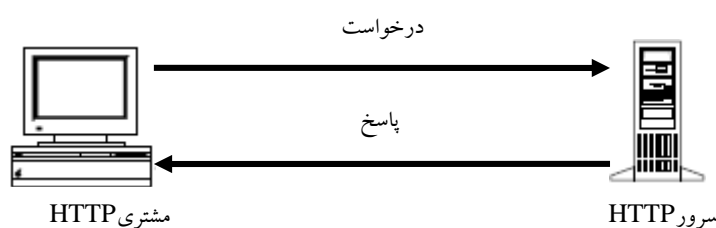
فرایند آماده‌سازی داده اغلب پرهزینه‌ترین مرحله از نظر زمان و محاسبات در فرآیند کشف دانش محسوب می‌شود [۶] و در مورد کاربرد کاوی وب نیز این موضوع صادق است. مرحله‌ی آماده‌سازی داده در وب‌کاوی کاربرد وب اغلب نیازمند استفاده از الگوریتم‌های خاص و مکاشفه‌ای است که در حوزه‌های دیگر داده‌کاوی در این مرحله بکار برده نمی‌شوند. این مرحله در استخراج موفقیت‌آمیز الگوهای مفید از داده‌ها بسیار حیاتی است. در ادامه موضوعات و مفاهیم مرتبط با آماده‌سازی و مدل‌سازی داده‌ها در کاربرد کاوی وب توضیح داده می‌شود.

۱-۶-۲ منابع داده

منابع عمده‌ی داده در کاربرد کاوی وب فایل‌های ثبت وب‌سرور می‌باشد که شامل ثبت‌های دسترسی وب‌سرور است.

داده‌های کاربرد

پیش از آنکه به داده‌های کاربرد و منابع آن پرداخته شود ابتدا تعاریفی از متاداده‌هایی که توسط وب‌سرورها تولید و استفاده می‌شوند ارائه می‌گردد [۷]. شکل ۲-۲ یک تراکنش HTTP^۱ را بین یک مشتری HTTP و یک سرور HTTP نشان می‌دهد. برای سادگی فرض کنید که مشتری HTTP یک مشتری وب است و یک سرور HTTP نیز یک وب‌سرور می‌باشد. یک مشتری وب که برای کاربران انسانی طراحی شده است یک مرورگر وب نامیده می‌شود مانند Mozilla Netscape Navigator، Firefox و Microsoft Internet Explorer. مثال‌هایی از وب‌سرور عبارتند از IBM HTTP Server، Apache HTTP Server و Microsoft Internet Information Server (IIS).



شکل ۲-۲ تراکنش HTTP

در یک تراکنش HTTP داده‌های کاربرد^۲ با متاداده‌های زیر تعریف می‌شوند:

- آدرس IP ماشین مشتری
- شناسه کاربر
- زمانی که سرور پردازش درخواست را انجام می‌دهد.

^۱ HyperText Transfer Protocol

^۲ Usage data