

## چارچوب سیستم حمایت از پژوهش برای وب کاوی

حمیرا عزت پناه<sup>۱\*</sup>

۱- گروه مهندسی کامپیوتر، دانشگاه گیلان، رشت، ایران.

## چکیده

امروزه طراحی و پیاده سازی یک سیستم پشتیبانی تحقیق و پژوهش، به چالشی برای محققان مایل به استفاده از اطلاعات مفید، تبدیل شده است. این مقاله چارچوبی برای پشتیبانی سیستم های وب کاوی پیشنهاد می کند. این سیستم ها به منظور شناسایی، استخراج، فیلترسازی و تجزیه و تحلیل داده ها از منابع وب طراحی شده اند. همچنین این سیستم ها، بازیابی وب و تکنیک های داده کاوی را با هم جهت ارائه زیر ساختی کارآمد به منظور پشتیبانی اطلاعات وب کاوی برای تحقیق ترکیب می کنند.

**کلمات کلیدی:** داده کاوی، بازیابی وب، خوشه، طبقه بندی، قوانین انجمن، سیستم پشتیبانی پژوهش.

## ۱- مقدمه

تکامل شبکه جهانی وب برای ما رشد مقدار زیادی از داده ها و اطلاعات را همراه داده های فراوان ارائه شده توسط وب، به ارمان آورده است. اینجاست که شبکه جهانی وب به یک منبع مهم برای تحقیق تبدیل می شود. با این حال، اطلاعات مرسوم و تکنیک های استخراج داده ها، به دلیل نیمه ساختار یافته بودن آن و یا حتی به دلیل طبیعت بدون ساختارشان نمی توانند به طور مستقیم به وب اعمال شوند. صفحات وب اسناد ابر متن است که شامل هر دوی متن و لینک به دیگر اسناد می باشد. علاوه بر این، اطلاعات وب ناهمگن و پویا هستند. بنابراین طراحی و پیاده سازی سیستم پشتیبانی پژوهش وب کاوی به یک چالش برای محققان به منظور استفاده از اطلاعات مفید وب، تبدیل شده است. معمولا وب کاوی، به عنوان محتوای وب کاوی و کاربرد وب کاوی طبقه بندی می شود. استخراج محتوای وب، مطالعه محتوای وب کاوی و بازیابی از اطلاعات وب است، در حالی که کاربرد وب کاوی، الگوی دسترسی کاربر را کشف و تجزیه و تحلیل می کند. [۱]

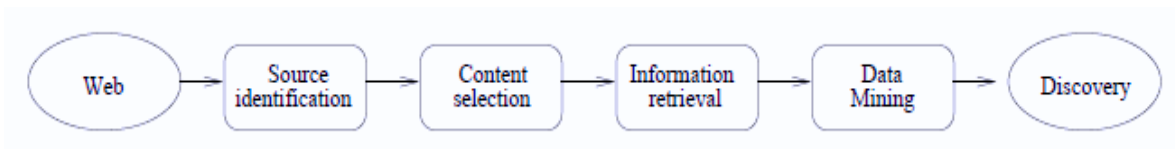
این مقاله چارچوبی برای طراحی وب در سیستم های پژوهش داده کاوی ارائه می کند. این سیستم ها برای شناسایی، استخراج، فیلترسازی و تجزیه و تحلیل داده ها از منابع وب طراحی شده اند و از تکنیک های بازیابی وب و داده کاوی با هم جهت ارائه زیر ساخت کارآمد به منظور حمایت از داده کاوی وب برای تحقیق ترکیب می کنند. این چارچوب از چند مرحله تشکیل شده است. ویژگی های هر مرحله در حال کاوش هستند و تکنیک های پیاده سازی نیز ارائه شده اند. مطالعه موردی روی داده کاوی از یک سایت توسعه یافته نرم افزار بزرگ، به طور مثال از نحوه استفاده از این چارچوب، ارائه شده است. کار ما یک راه حل عمومی فراهم می کند که محققان می توانند استفاده از منابع وب در پژوهش خود را دنبال

\*کارشناسی نرم افزار از دانشگاه گیلان  
h.ezatpanah@yahoo.com

کنند. بقیه این مقاله به شرح زیر است: بخش ۲ یک چارچوب برای حمایت پژوهش وب کاوی سیستم‌ها و بازبینی مختصری از اجزای آن ارائه می‌کند. بخش ۳ طراحی و پیاده‌سازی بازبازی اطلاعات وب را توصیف می‌کند. بخش ۴ در پردازش و تجزیه و تحلیل داده‌های وب تمرکز دارد. بخش ۵ یک مطالعه موردی بر اساس وب کاوی سیستم پشتیبانی پژوهش ارائه می‌کند. نتیجه‌گیری و کار آینده در بخش ۶ ارائه شده است.

## ۲- نمای کلی چارچوب

به منظور کشف اطلاعات وب، ما پژوهش چارچوب سیستم پشتیبانی برای داده‌های وب کاوی را به عنوان نشان داده شده در شکل ۱، متشکل از ۴ مرحله، می‌سازیم. شناسایی منبع، انتخاب محتوا، بازبازی اطلاعات و داده کاوی. در مرحله اول، وب سایت‌های مناسب باید با توجه به نیازهای پژوهشی انتخاب شوند. این شامل در دسترس بودن شناسایی، ارتباط و اهمیت وب سایت است. جستجوی واژه‌های کلیدی با استفاده از موتور جستجو می‌تواند برای پیدا کردن وب سایت مناسب، مورد استفاده قرار گیرد. بعد از پیدا کردن تمام وب سایت‌های شناسایی شده در فاز اول، مرحله دوم انتخاب مطالب مناسب در آن وب سایت است. مانند اسناد، گروه‌های خبری، انجمن‌ها، لیست پستی و ... معمولاً یک وب سایت حاوی بسیاری از صفحات وب، از جمله اطلاعات مربوط و نامربوط است. این فاز مهم است به این دلیل که تصمیم می‌گیرد که اطلاعات وب باید استخراج شود. انتخاب صفحات وب بر اساس هدف و تجربه محقق است. در مرحله بازبازی اطلاعات، یک خزنده طراحی شده است. در مرحله بازبازی اطلاعات، خزنده به صورت خودکار اطلاعات انتخاب شده را در طول انتخاب فاز استخراج می‌کند.



شکل ۱- چارچوب سیستم پشتیبانی پژوهش برای داده‌های وب کاوی

ابزارها و تکنیک‌های خاص به طور موثر برای بازبازی مفید دانش/اطلاعات از منابع وب، به کار گرفته شدند. تلاش اضافی ممکن است برای بازبازی محتوای پویا و منابع داده مشخص مانند گروه خبری، انجمن و ... مورد نیاز باشد. مرحله نهایی، انجام داده کاوی در استخراج اطلاعات وب است. این شامل آماده‌سازی داده‌ها برای تجزیه و تحلیل است. استخراج صفحه وب ممکن است حاوی اطلاعات از دست رفته، داده‌های اضافی، فرمت اشتباه و کاراکترهای غیر ضروری باشد. علاوه بر این، برخی از اطلاعات باید به منظور حفاظت از حریم شخصی پردازش شوند. تکنیک‌های داده کاوی پیشرفته در اینجا برای کمک به تجزیه و تحلیل اطلاعات به کار گرفته می‌شوند.

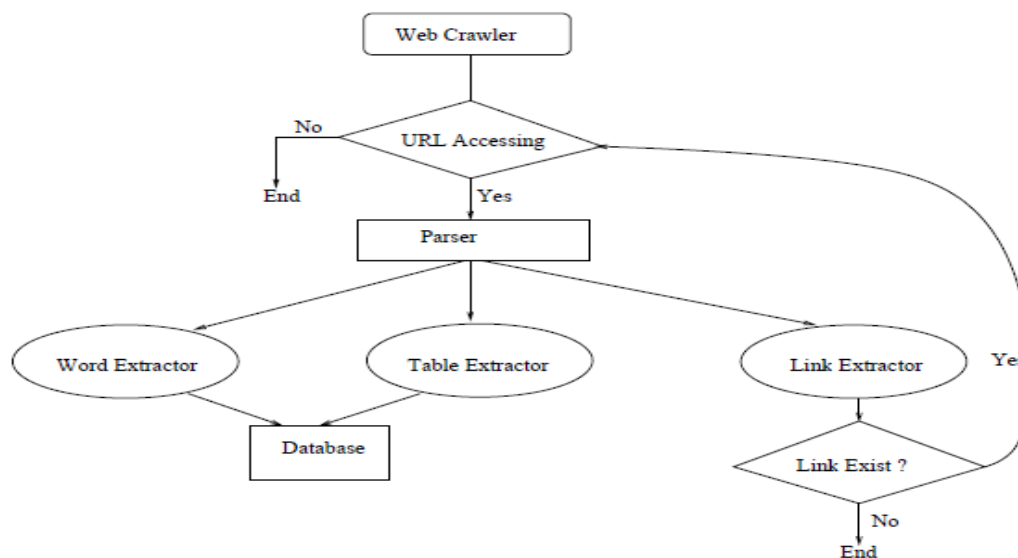
## ۳- بازبازی اطلاعات

بازبازی اطلاعات برای دسترسی به داده‌ها در وب استفاده می‌شود. به این معنا که سیستم پشتیبانی پژوهش وب کاوی باید قادر به جستجو و بازبازی محتویات خاص وب، به طور موثر و کارآمد باشد. دو دسته عمده از ابزار جستجو وجود دارد:

دایرکتوری (پاهو، نت اسکپ و ...) و موتورهای جستجو (لیکس، گوگل و ...) برای استفاده از دایرکتوری‌ها همراه افزایش وب سایت این کار سختی است. موتورهای جستجو نمی‌توانند هر نیاز مورد جستجویی را بر آورده کنند. در سیستم ما، خزنده وب بر اساس ابزارهای پیشرفته و تکنیک‌ها برای کمک به پیدا کردن اطلاعات مفید از منابع، توسعه داده شده است. خزنده وب عنکبوت، ربات، کرم و ... نامیده می‌شود. خزنده وب شامل برنامه‌ای است که به طور خودکار از سایت‌ها عبور می‌کند، اسناد را دانلود می‌کند و لینک‌ها را به صفحات دیگر می‌فرستد. [۴]

خزنده وب یک کپی از تمام بازدیدهای صفحات برای استفاده‌های بعدی نگه می‌دارد. بسیاری از موتورهای جستجوی وب از خزنده وب برای ایجاد مطالب جهت نمایه‌سازی استفاده می‌کنند. آن‌ها همچنین می‌توانند در دیگر برنامه‌های کاربردی مانند اعتبار سنجی صفحه، تجزیه و تحلیل ساختاری و تجسم، به روز رسانی اخبار، معکوس‌سازی و ... استفاده شوند. [۲]

موتورهای جستجو برای وب‌کاوی یک پروژه تحقیقاتی کافی نیستند. برای طراحی، یک خزنده وب لازم است که شامل متودی برای پیدا کردن و جمع‌آوری تحقیقات اطلاعات مرتبط از وب باشد. اگرچه پروژه‌های تحقیقاتی، اطلاعات وب سایت‌های متفاوتی دارند که منجر به خزنده‌های وب می‌شود، این خزنده‌ها هنوز برخی طرح‌های مشترکی دارند که در شکل ۲ نشان داده شده است. آن‌ها می‌توانند با جاوا، پیل، پایتون و ... پیاده‌سازی و اجرا شوند. خزنده وب باید همراه یک <sup>۲</sup> (URL) شروع شود، دیگر لینک‌ها روی صفحه <sup>۳</sup> (HTML) شناسایی شوند، این لینک‌ها بازدید شوند، اطلاعات از صفحات وب استخراج شوند و اطلاعات در پایگاه داده ذخیره شوند. بنابراین، یک خزنده وب شامل یک روش دسترسی URL، تجزیه‌کننده صفحه وب همراه با برخی استخراج‌کننده‌ها و پایگاه داده‌ها است.



شکل ۲- خزنده وب

تابع دسترسی خزنده باید بارها و بارها از دسترسی به آدرس یکسان وب جلوگیری کند و باید لینک‌های مرده را شناسایی کند. تجزیه‌کننده برچسب ابتدایی، برچسب پایانی، متن و نظرات را به رسمیت می‌شناسد. پایگاه‌های داده ذخیره‌سازی را برای استخراج اطلاعات فراهم می‌کنند. اجزای کلیدی خزنده وب شامل پارسر، که خود شامل استخراج کلمه، استخراج جدول و استخراج لینک است می‌باشد. استخراج‌کننده کلمه برای استخراج اطلاعات کلمه استفاده می‌شود. آن

<sup>۲</sup> Uniform Resource Locator

<sup>۳</sup> HyperText Markup Language

باید توابع چک کردن رشته را فراهم کند. جداول معمولاً در صفحات وب به منظور چینش و هم تراز کردن اطلاعات استفاده می‌شوند. استخراج کننده جدول محل داده‌ها در یک جدول را شناسایی می‌کند. استخراج کننده لینک، لینک‌های موجود در یک صفحه وب را بازیابی می‌کند. دو نوع از لینک‌ها- لینک‌های مطلق و لینک‌های نسبی وجود دارند. لینک مطلق آدرس کامل یک صفحه وب را می‌دهد، در حالی که لینک وابسته به اضافه کردن یک آدرس کامل با اضافه کردن یک پیشوند نیاز دارد.

#### ۴- تکنولوژی داده کاوی

به منظور تسهیل در وب کاوی، الگوریتم‌های داده کاوی زیر می‌توانند پیدا کردن الگوها و روند‌هایی در اطلاعات جمع‌آوری شده از وب سایت اعمال کنند: خوشه بندی، طبقه بندی، قوانین انجمنی. در بخش‌های زیر، ما هر یک از الگوریتم‌ها و برنامه‌های کاربردی خود را توضیح خواهیم داد.

#### ۴-۱- مشاهده قوانین انجمن

قوانین انجمن برای پیدا کردن ارتباط جالب و یا رابطه همبستگی میان یک مجموعه بزرگی از داده‌ها تلاش می‌کند. نمونه‌ای از قوانین انجمن کاوی تجزیه و تحلیل سبد خرید است. قوانین انجمن چیزی است مانند ۸۰٪ افرادی که آبجو، مرغ سرخ شده در روغن، خرید می‌کنند. قوانین انجمن همچنین می‌توانند برای پیش‌بینی الگوهای دسترسی به وب شخصی استفاده شوند. به طور مثال، ما ممکن است که ۸۰٪ افرادی که به صفحه A و صفحه B و همچنین صفحه C دسترسی دارند را کشف کنیم. صفحه C ممکن است یک لینک مستقیم از صفحه A و یا صفحه B نداشته باشد. اطلاعات کشف شده ممکن است برای ایجاد یک لینک به صفحه C و یا صفحه B مورد استفاده قرار گیرد. یک نمونه از این نرم افزار amazon.com است. ما اغلب چیزی شبیه به "مشتریانی که این کتاب و همچنین کتاب A را می‌خرند" می‌بینیم. قوانین انجمن می‌تواند به وب داده برای کشف رفتار کاربران وب و پیدا کردن الگوهای رفتاری خود ارائه شود.

#### ۴-۲- طبقه بندی

هدف طبقه بندی، پیش‌بینی یک مورد یا چندین کلاس (یا مشاهده) است. هر مورد شامل  $n$  صفت است، که یکی از این موجودیت‌ها هدف است، و بقیه پیش‌بینی ویژگی‌هاست. هر یک از اهداف ارزش ویژگی‌های یک کلاس پیش‌بینی بر اساس  $n-1$  خواص و ویژگی‌های پیش‌بینی کننده است. طبقه بندی یک فرآیند دو مرحله‌ای است. اول، یک طبقه بندی مدل بر اساس مجموعه‌ای از آموزش‌های داده‌ای ساخته شده است. دوم، مدل به منظور داده‌های جدید برای طبقه بندی ارائه شده است. در بین این دو مرحله، برخی مراحل دیگر ممکن است به کار گرفته شود، از جمله محاسبات بلند. محاسبات بلند برای تایید این است که آیا یک طبقه بندی، مدل با ارزشی است. مقدار بزرگتر از ۱ به طور معمول خوب است. مدل‌های طبقه بندی می‌توانند جهت تصمیم‌گیری کسب و کار به کار گرفته شوند. کاربردها شامل طبقه بندی عنوان پیام‌ها ایمیل‌های ناخواسته، تشخیص تقلب در کارت‌های اعتباری، تشخیص شبکه نفوذ و ... می‌باشند.

## ۴-۳- خوشه بندی

خوشه برای گروه‌بندی طبیعی از داده‌ها استفاده می‌شود. این‌ها گروه‌بندی طبیعی خوشه هستند. خوشه مجموعه داده‌هایی است که شبیه یکدیگرند. یک الگوریتم خوشه‌بندی خوب، خوشه‌هایی تولید می‌کند، به طوری که شباهت بین خوشه کم است و شباهت درون خوشه بالاست. خوشه می‌تواند جهت مشتریان گروه همراه رفتار مشابه و تصمیم‌گیری کسب و کار در صنعت استفاده شود.

## ۵- مطالعه موردی: منبع باز نرم‌افزار توسعه یافته (OSS)

جامعه<sup>۴</sup> (OSS) مقدار قابل توجهی از زیرساخت‌های اینترنت را توسعه داده است و دارای چندین دستاورد برجسته فنی، از جمله آپاچی، پرل، لینوکس و ... است. این برنامه‌ها نوشته شده بودند همچنین توسعه یافته شده بودند، و تا حد زیادی توسط شرکت‌کنندگان در بخشی از زمان اشکال زدایی شده بودند. که در اغلب موارد برای کار پرداخته نشده بودند، و از هیچ تکنیک مدیریت پروژه‌ای بهره‌نگرفته بودند. پژوهش چگونه ممکن است از تابع جامعه OSS به برنامه ریزان (IT) کمک کند و منجر به تصمیم‌گیری آگاهانه‌تر و توسعه موثرتر استراتژی برای استفاده از نرم‌افزار OSS شود. جامعه توسعه داده شده OSS یک انجمن جهانی مجازی است. بنابراین، در حال حاضر ما مزیتی در تعامل دیجیتال داریم که بایگانی شده هستند و داده‌ها می‌توانند داده‌کاوی و استخراج شوند. با حدود ۷۰۰۰۰ پروژه، ۹۰۰۰۰ توسعه‌دهنده، و ۷۰۰۰۰۰ کاربران ثبت نام شده، SourceForge.net با حمایت مالی نرم‌افزار VA بزرگترین توسعه و همکاری را با OSS داشته است. این سایت اطلاعات دقیق در مورد پروژه و توسعه‌دهندگان، از جمله ویژگی‌های پروژه، فعال‌ترین پروژه‌ها و رتبه‌بالی توسعه‌دهندگان فراهم می‌کند.

## ۵-۱- جمع‌آوری داده‌ها

پس از اطلاع‌رسانی به سورس برنامه‌های ما و دریافت اجازه، ما ماهنامه داده‌ها را در سورس جمع کردیم. ابزار سورس مدیریت پروژه، ردیابی اشکال، فهرست خدمات پست الکترونیکی، بحث و تبادل نظر، کنترل نسخه نرم‌افزار برای پروژه‌های میزبان را فراهم می‌کند. اطلاعات در جامعه جمع‌آوری شده، پروژه و سطح توسعه‌کننده، مشخص کل پدیده OSS است، در سرشماری از چند پروژه، رفتار و مکانیزم‌هایی در محل کار و سطح پروژه و توسعه‌کننده بررسی می‌شوند. اطلاعات اولیه مورد نیاز برای این تحقیق شامل دو جدول آمار پروژه و توسعه‌دهندگان می‌باشد. جدول آمار پروژه در شکل ۱ نشان داده شده است. شامل رکوردهایی با ۹ فیلد: شناسه پروژه، نمایش‌های طول عمر، رتبه، صفحه، دانلود، اشکالات، پشتیبانی، تکه و<sup>۱</sup> (CVS).

جدول توسعه‌دهندگان دو فیلد دارد: شناسه پروژه و شناسه توسعه‌دهنده. به این دلیل که پروژه‌ها می‌توانند بسیاری توسعه‌دهنده داشته باشند و توسعه‌دهندگان نیز می‌توانند بسیاری پروژه داشته باشند، هیچ یک از فیلدها منحصر به فرد برای کلید اصلی نیستند. بنابراین کلید کامپوزیت متشکل از هر دو ویژگی به عنوان یک کلید اصلی در خدمت است. هر پروژه در سورس دارای یک ID (شناسه) منحصر به فرد در سورس خواهد بود. خزنده وب توسط پرل و

<sup>4</sup> Office of Strategic Services

<sup>5</sup> Information technology

<sup>6</sup> Concurrent Versions System

<sup>۷</sup> (CPAN) (جامع بایگانی پرل - مخزن ماژول پرل/ کتابخانه)، گذشته از سورس برای وب سرور جهت جمع آوری اطلاعات لازم، ماژول سازی شده است. همه صفحات اصلی پروژه ها در سورس، یک طراحی سطح بالای مشابه دارند. بسیاری از این صفحات به صورت پویا از یک پایگاه داده ایجاد شدند. خزنده وب از کتابخانه-libwww- پرل، به بهانه هر پروژه صفحه خانگی استفاده می کند. CPAN دارای یک تجزیه گر HTML عمومی جهت شناسایی و به رسمیت شناختن شروع برچسب ها، پایان برچسب ها، متن و نظرات و ... است. به این دلیل که هر دو اطلاعات آماری و عضو در جداول ذخیره شده اند، خزنده وب از ماژول پرل موجود به نام HTML:TableExtract و رشته مقایسه بارانه پرل برای استخراج اطلاعات استفاده می کند. اگر این ها بیشتر از اعضای یک صفحه باشند، استخراج لینک ها استفاده می شوند.

### ۵-۲- داده کاوی

چندین الگوریتم داده کاوی وجود دارد که می تواند برای وب داده استفاده شود. که در میان آن ها ساده و بی تکلف، طبقه بندی شده، و دارای رگرسیون درخت (سبب خرید) وجود دارد. بیا بیا ابتدا روی الگوریتم های طبقه بندی ساده و بی تکلف بیز و سبب خرید تمرکز کنیم.

#### جدول ۱- آمار پروژه

project ID	lifespan	rank	page views	downloads	bugs	support	patches	all trackers	tasks	cvs
1	1355 days	31	12,163,712	71,478	4,160	46,811	277	52,732	44	0
2	1355 days	226	4,399,238	662,961	0	53	0	62	0	14,979
3	1355 days	301	1,656,020	1,064,236	364	35	15	421	0	12,263
7	1355 days	3322	50,257	20,091	0	0	0	12	0	0
8	1355 days	2849	6,541,480	0	17	1	1	26	0	13,896

### ۵-۲-۱ طبقه بندی

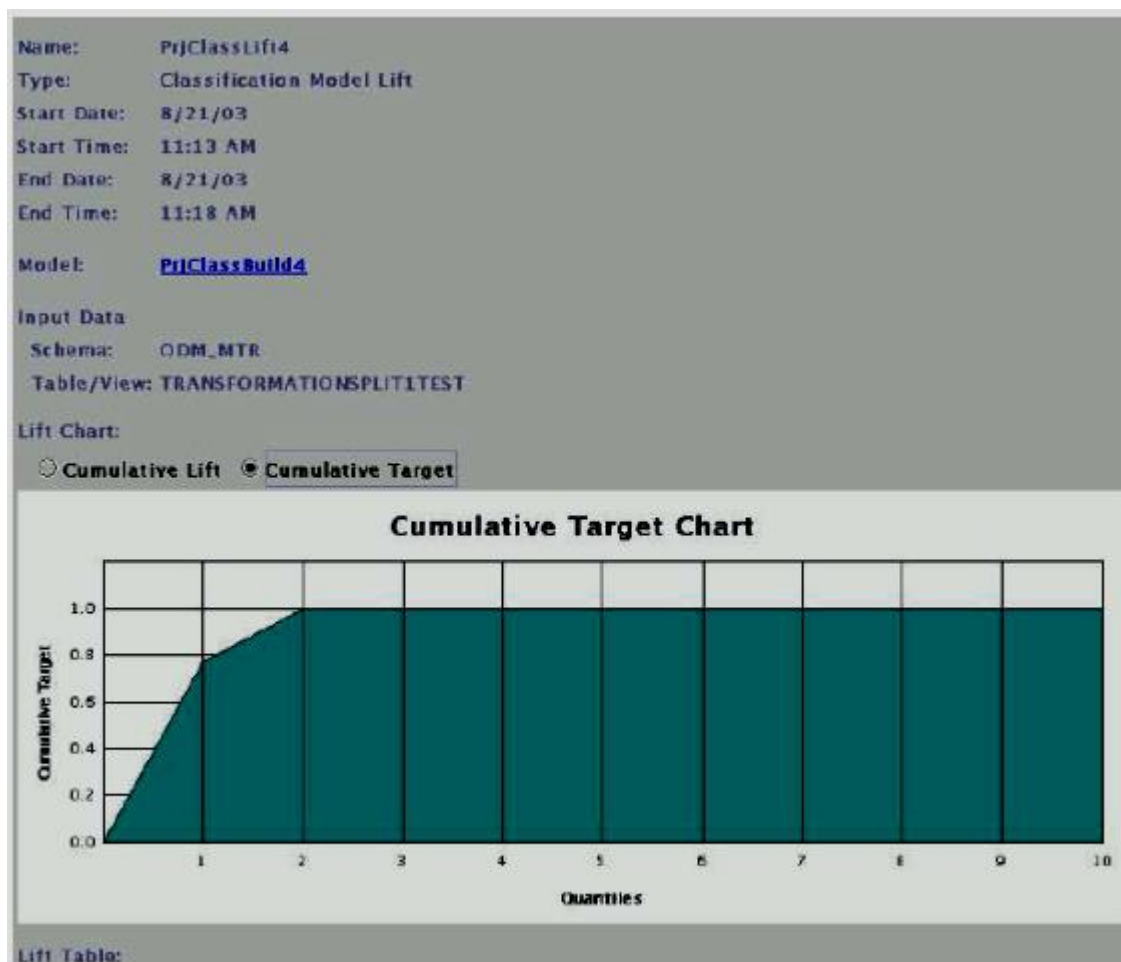
الگوریتم ساده و بی تکلف بیز با استفاده از قضیه بیز پیش بینی می کند. ساده و بی تکلف بیز فرض می کند که هر صفتی مستقل از دیگران است. در این مطالعه، که موردی نیست، به عنوان مثال، "دانلود" ویژگی ای نزدیک به ویژگی "CVS" است، ویژگی "رتبه" از نزدیک به ویژگی های دیگر مرتبط است و از دیگر ویژگی ها محاسبه می شود. سبب خرید ایجاد طبقه بندی و درخت رگرسیون برای پیش بینی متغیر وابسته (رگرسیون) و طبقه بندی متغیرهای پیش بینی کننده (طبقه بندی) است. ما تنها زمانی علاقمند به نوع طبقه بندی از مشکلات هستیم که ما تنها از دسته این نوع مشکلات در حال استفاده باشیم. اجرا و پیاده سازی اوراکل از سبب خرید به نام تطبیقی شبکه بیز<sup>۸</sup> (ABN) نامیده می شود. ABN پیش بینی اهداف باینری و طبقه بندی چندگانه است. بنابراین جداسازی موجودیت هدف مطلوب است. در مطالعه این مورد، ما سعی می کنیم دانه‌ها (بارگیری) را از ویژگی های دیگر پیش بینی کنیم.

همانطور که قبلاً گفته شد، ویژگی "دانلود" نه برابر داخل باکت است. ما پیش بینی می کنیم که دریافت دانه‌های باکت بر اساس ارزش دیگر ویژگی هاست. به عنوان انتظار، الگوریتم بیز ساده و بی تکلف است و مناسب برای پیش بینی "دانلود" و آن مربوط به دیگر ویژگی ها مانند "CVS" است. دقت در بیز 10% است. در حالی که بیز روی پیش بینی

<sup>7</sup> Comprehensive Perl Archive Network

<sup>8</sup> Australian Business Number

"دانلود" به بدی عمل می‌کند. الگوریتم ABN می‌تواند "دانلود" را با کمال دقت ۶۳٪ پیش‌بینی کند. در ابتدای دید، دقت ۶۳٪ جذاب نیست، اما یک پیش‌بینی خوب است زمانی که ما می‌توانیم ۱۰ درصد پیش‌بینی درست و بدون طبقه بندی داشته باشیم. در نتیجه محاسبات بلند تایید می‌کند که طبقه بندی کاملاً خوب است، همانطور که در شکل ۳ نشان داده شده است.



شکل ۳- نمودار بلند

این ارقام نشان می‌دهد که ما همه رکوردهای مربوط به ویژگی "دانلود" را می‌یابیم، ویژگی ای که فقط در ۲۰٪ برای اولین بار از سوابق ۱ است. قوانین ساخته شده توسط طبقه بندی ABN نشان می‌دهد که "دانلود" نزدیک و مرتبط "CVS" است. ما نتیجه می‌گیریم که الگوریتم ABN، مناسب برای پیش‌بینی دانلود در مطالعه ماست. جدول زیر مقایسه دو الگوریتم یعنی ABN و Navis bayes است. از جدول، ما می‌بینیم که ABN جهت ساخت یک مدل طبقه بندی، بسیار طولانی طول می‌کشد، اما در نتیجه مدلی بسیار دقیق تر است.

## جدول ۲. مقایسه بین ABN و Naive Bayes

Name	Build Time	Accuracy
ABN	0:19:56	63%
Naive Bayes	0:0:30	9%

## ۵-۲-۲ مشاهده قوانین انجمن

مشکل قوانین انجمن را می‌توان به دو مسئله زیر تجزیه کرد:

یافتن همه ترکیبات از آیت‌ها، نام مجموعه اقلام مکرر، که پشتیبانی پس از آن حداقل حمایت است. استفاده از مجموعه اقلام مکرر برای تولید انجمن، برای مثال، اگر AB مجموعه اقلام مکرر هستند، پس از آن نقش  $A \rightarrow B$  را حفظ می‌کند اگر نسبت پشتیبانی (AB) برای حمایت از A بیشتر از حداقل اعتماد باشد. یکی از معروف ترین قوانین الگوریتم های کاوش به نام Priori است. اوراکل این الگوریتم را با استفاده از SQL پیاده سازی می‌کند. ما سعی می‌کنیم از این الگوریتم برای پیدا کردن ارتباط بین ویژگی های پروژه استفاده کنیم. الگوریتم ۲ ورودی دارد، یعنی حداقل حمایت و حداقل اعتماد به نفس. ما (۰.۰۱) را برای حداقل حمایت و (۰.۵) را برای حداقل اعتماد به نفس انتخاب می‌کنیم. ما همه ویژگی های CVS و دانلود را می‌یابیم. قوانین بیشتر را می‌توان در شکل ۵ دید. هیچ یک از قوانین بر اساس علاقه کشف نشده اند.

## ۵-۲-۳ خوشه بندی

ما علاقمند به قرار دادن پروژه با ویژگی های مشابه به صورت خوشه هستیم. دو الگوریتم را به این منظور می‌توان مورد استفاده قرار داد: k-means و o-cluster.

k-means یک الگوریتم مبتنی بر راه دور است، که داده ها را به تعدادی خوشه از پیش تعریف شده پارتیشن بندی می‌کند. o-cluster الگوریتمی سلسله مراتبی و مبتنی بر شبکه است. خوشه منجر به تعریف مناطق متراکم در ویژگی فضا است. بعد از فضای ویژگی، شماری از تعداد صفات، درگیر در الگوریتم خوشه بندی هستند. ما دو الگوریتم خوشه بندی پروژه را در مطالعه این مورد به کار می‌بریم. شکل ۴ و ۵ نتایج خوشه بندی و قوانینی که خوشه را تعریف می‌کند نشان می‌دهد.



Settings Clusters Rules Results

Leaf Clusters: 10  
Cluster Levels: 5  
Cases: 50000

Clusters:  Show Leaves Only

Cluster ID:	Cases	Split Rule
1	50000	SUPPORT in (4)
3	23828	LIFESPAN in (4)
5	8534	SUPPORT in (6)
18	3590	n/a
19	4944	n/a
7	15294	PAGE_VIEWS in (7)
8	10449	PAGE_VIEWS in (4)
14	5412	n/a
15	5037	n/a
9	4845	n/a
2	26172	ALL_TRKS in (4)
4	10993	SUPPORT equal (1)
12	3035	n/a
13	7958	n/a
5	15179	BUGS in (6)
16	5043	BUGS in (6)

Detail  
Expand All  
Collapse All

شکل ۴- خوشه بندی

Cluster Rule Display Criteria:

Only Show Rules for Leaf Clusters

Only Show Attributes with Minimum Relevance Rank: 10 Refresh

Rules

If (condition)	Then (cluster)	Confidence	Support
ALL_TRKS in (10, 3, 4, 5, 8, 9) and BUGS in (1, 10, 3, 8, 9) and CVS in (1, 10, 2, 5, 6, 8, 9) and DOWNLOADS in (10, 5, 6, 7, 8, 9) and LIFESPAN in (10, 5, 6, 7, 8, 9) and PAGE_VIEWS in (10, 8, 9) and PATCHES in (1, 10, 5, 6, 9) and RANK in (2, 3, 4, 5, 6, 7, 9) and SUPPORT in (10, 5, 6, 7, 8) and TASKS in (1, 10, 2, 3, 8, 9)	CLUSTER equal (9)	1.0	1.0
ALL_TRKS in (10, 5, 6, 8, 9) and BUGS in (1, 10, 3, 8, 9) and CVS in (1, 10, 2, 5, 6, 8, 9) and DOWNLOADS in (10, 5, 6, 7, 8, 9) and LIFESPAN in (10, 5, 6, 7, 8, 9) and PAGE_VIEWS in (10, 8, 9) and PATCHES in (1, 10, 5, 6, 9) and RANK in (2, 3, 4, 5, 6, 7, 9) and SUPPORT in (10, 5, 6, 7, 8) and TASKS in (1, 10, 2, 3, 8, 9)	CLUSTER equal (11)	1.0	1.0
ALL_TRKS in (1, 2, 3) and BUGS in (2, 3, 4) and CVS in (1, 10, 2, 5, 6, 8, 9) and DOWNLOADS in (10, 5, 6, 7, 8, 9) and LIFESPAN in (10, 5, 6, 7, 8, 9) and PAGE_VIEWS in (10, 8, 9) and PATCHES in (1, 10, 5, 6, 9) and RANK in (2, 3, 4, 5, 6, 7, 9) and SUPPORT in (10, 5, 6, 7, 8) and TASKS in (1, 10, 2, 3, 8, 9)	CLUSTER equal (12)	1.0	1.0
ALL_TRKS in (1, 2, 3, 4) and BUGS in (1, 2, 3, 4) and CVS in (1, 10, 2, 5, 6, 8, 9) and DOWNLOADS in (10, 5, 6, 7, 8, 9) and LIFESPAN in (10, 5, 6, 7, 8, 9) and PAGE_VIEWS in (10, 8, 9) and PATCHES in (1, 10, 5, 6, 9) and RANK in (2, 3, 4, 5, 6, 7, 9) and SUPPORT in (10, 5, 6, 7, 8) and TASKS in (1, 10, 2, 3, 8, 9)	CLUSTER equal (13)	1.0	1.0
ALL_TRKS in (3, 4, 5, 6, 8, 9) and BUGS in (1, 10, 3, 8, 9) and CVS in (1, 10, 2, 5, 6, 8, 9) and DOWNLOADS in (10, 5, 6, 7, 8, 9) and LIFESPAN in (10, 5, 6, 7, 8, 9) and PAGE_VIEWS in (10, 8, 9) and PATCHES in (1, 10, 5, 6, 9) and RANK in (2, 3, 4, 5, 6, 7, 9) and SUPPORT in (10, 5, 6, 7, 8) and TASKS in (1, 10, 2, 3, 8, 9)	CLUSTER equal (14)	1.0	1.0
ALL_TRKS in (3, 4, 5, 6, 9) and BUGS in (1, 10, 3, 8, 9) and CVS in (1, 10, 2, 5, 6, 8, 9) and DOWNLOADS in (10, 5, 6, 7, 8, 9) and LIFESPAN in (10, 5, 6, 7, 8, 9) and PAGE_VIEWS in (10, 8, 9) and PATCHES in (1, 10, 5, 6, 9) and RANK in (2, 3, 4, 5, 6, 7, 9) and SUPPORT in (10, 5, 6, 7, 8) and TASKS in (1, 10, 2, 3, 8, 9)	CLUSTER equal (15)	1.0	1.0

Rule Detail

IF:  
ALL\_TRKS in (10, 3, 4, 5, 8, 9) and BUGS in (1, 10, 3, 8, 9) and CVS in (1, 10, 2, 5, 6, 8, 9) and DOWNLOADS in (10, 5, 6, 7, 8, 9) and LIFESPAN in (10, 5, 6, 7, 8, 9) and PAGE\_VIEWS in (10, 8, 9) and PATCHES in (1, 10, 5, 6, 9) and RANK in (2, 3, 4, 5, 6, 7, 9) and SUPPORT in (10, 5, 6, 7, 8) and TASKS in (1, 10, 2, 3, 8, 9)

THEN  
CLUSTER equal (9)

Confidence=1.0  
Confidence=1.0

شکل ۵- قوانینی که خوشه را تعریف می کند

## ۶- نتیجه گیری و کار آینده

این مقاله یک چارچوب برای سیستم پشتیبانی وب کاوی پژوهش ارائه می کند و روش های آن را توصیف می کند. سپس تکنیک های پیاده سازی در استخراج و تجزیه و تحلیل داده های وب بحث می شود. سورتس وب کاوی موردی است که ارائه شده است، به طور مثال، چگونگی کاربرد این چارچوب، یک مطالعه اکتشافی از بازیابی و داده کاوی داده های وب است. ما سعی می کنیم اطلاعات استخراج شده و نرم افزار داده کاوی که می تواند برای کشف دانش در وب داده استفاده شود را بررسی کنیم. نوآوری های واقعی و جالب هنوز در حال پیشرفت هستند. ما انتظار به کشف الگوهای جالب از داده ها داریم. همچنین لازم به ذکر است که این پژوهش بنیاد علوم ایالت متحده شد.

1. R. Cooley, B. Mobasher, and J. Srivastava. Web mining: Information and pattern discovery on the world wide web. In International Conference on Tools with Artificial Intelligence, pages 558–567, Newport Beach, 1997.
2. Francis Crimmins. Web crawler review. [http://dev.funnelback.com/crawler review.html](http://dev.funnelback.com/crawler%20review.html), 2001.
3. Yao J.T. and Yao Y.Y. Web-based information retrieval support systems: building research tools for scientists in the new information age. In Proceedings of the IEEE/WIC International Conference on Web Intelligence, Halifax, Canada, 2003.
4. M. Koster. The web robots pages. <http://info.webcrawler.com/mak/projects/robots/robots.html>, 1999.
5. Bamshad Mobasher, Honghua Dai, Tao Luo, Yuqing Sun, and Jiang Zhu. Integrating web usage and content mining for more effective personalization. In EC-Web, pages 165–176, 2000.
6. Yao Y.Y. Information retrieval support systems. In FUZZIEEE' 02 in The 2002 IEEE World Congress on Computational Intelligence, Honolulu, Hawaii, USA, 2002.
7. Yao Y.Y. A framework for web-based research support systems. In Computer Software and Application Conference, (COMPOSAC 2003), Dallas, Texas, 2003.
8. Osmar R. Zaiane, Man Xin, and Jiawei Han. Discovering web access patterns and trends by applying OLAP and data mining technology on web logs. In Advances in Digital Libraries, pages 19–29, 1998.