

تحلیل رفتار کاربران برای بهبود پالایش دسترسی به سایتها

مهدی نقوی^{۱*}، وحید نقوی^۲، بهروز مینایی^۳

۱- دانشجوی دکتری، ۳- استادیار، دانشگاه علم و صنعت ایران، دانشکده کامپیوتر

۲- دانشجوی کارشناسی ارشد، دانشگاه آزاد اسلامی، واحد اراک

E-mail: naghavi@iust.ac.ir

(دریافت: ۸۸/۰۹/۰۴، پذیرش: ۸۹/۰۶/۲۹)

چکیده

در این مقاله روش نوینی جهت بهبود پالایش دسترسی به سایتها ارائه شده است. این روش مبتنی بر تحلیل رفتار بازدیدکنندگان صفحات وب می باشد. برای تحلیل رفتار کاربران، معماری ارائه شده است که در سه مرحله بر اساس اطلاعات موجود در فایل های گزارش سرورهای پروکسی، الگوی رفتار کاربران وب براساس یک موتور داده کاو، استخراج شده و در بانک اطلاعات ذخیره می گردد. هنگام درخواست کاربر برای یک صفحه جدید، با بررسی و تحلیل سابقه رفتار کاربر، دسترسی مجاز و یا غیر مجاز به صفحه مورد نظر، اعلام می گردد. معماری ارائه شده با استفاده از داده های یکی از سایت های معتبر جهانی آزمایش شده و نتایج آن مورد بررسی و تحلیل قرار گرفته است.

کلیدواژه ها: پالایش دسترسی به وب؛ تحلیل رفتار؛ فیلترینگ؛ کاوش وب؛ مسدود کردن سایتها

Analyzing Web Users Behavior to Improve Web Filtering

M. Naghavi^{1*}, V. Naghavi², B. Minaei³

Department of Computer, Iran University of Science and Technology

E-mail: naghavi@iust.ac.ir

Abstract

In this paper, a new method to improve sites access filtering is presented. This method is based on web visitor behavior analysis. In order to analyze the behavior of users, A procedure is proposed following which, based on the information in the log files of proxy servers, user's behavior history was extracted and stocked in a data bank by a data-mining engine. Authorized or unauthorized access to the desired page, when a user requests a new page, is reported after investigation of the user's behavior history. The above architecture was tested, using data from one of the globally credible sites, and the results were analyzed.

Keywords: Behavior Analysis; Web Mining; Web Usage Mining; Web Filtering; Web Blocking

۱. مقدمه

گسترش روزافزون دنیای وب باعث بوجود آمدن نیازهای جدیدی شده است که یکی از آنها، پالایش سایتها بر اساس اهداف مختلف است. طبق گزارش WebSpy در سال ۲۰۰۸، روزانه حدود ۴۰ میلیون کاربر فقط به منظور سرگرمی و گذراندن وقت وارد شبکه اینترنت می‌شوند [۱]. یک سوم وقت کارمندان شرکت‌های آمریکایی در امور غیر مرتبط با کار در اینترنت صرف شده و هزینه‌های ناشی از آن بالغ بر ۸۵ میلیارد دلار در سال برآورد شده است [۲]. طبق گزارش Covenant Eyes^۱ از هر ۵ نفر، ۲ نفر جهت بازدید از سایت‌های غیر اخلاقی به اینترنت متصل می‌شوند. کلمه sex دارای رتبه اول کلمات جستجو در موتورهای جستجو بوده و سایر کلمات مرتبط نیز از رتبه‌های بالا برخوردار می‌باشند [۳].

این مقاله با رویکرد تحلیل رفتار بازدیدکنندگان وب، نیاز فوق را مورد توجه قرار داده و در راستای بهبود فیلترینگ سایتها ارائه شده است. در ابتدا موضوع پالایش دسترسی به وب مورد بررسی قرار گرفته و سپس مبانی مربوط به کاوش رفتار در وب بیان شده است. در ادامه، کارهای مرتبط با موضوع کاوش رفتار در وب توضیح داده شده و به بیان طرح مسئله و ارائه راه‌حل پرداخته و در نهایت چگونگی معماری استخراج الگو و تحلیل رفتار کاربر ارائه شده است. در بخش‌های بعدی به عنوان نمونه به بررسی آزمایشات صورت گرفته بر روی داده‌های پایگاه AOL پرداخته شده و در نهایت چالش‌های پیش‌رو و پیشنهاد کارهای آینده بیان شده است.

۲. پالایش دسترسی به وب

عوامل فوق، سازمان‌ها و شرکت‌ها را بر آن داشت که جهت کم کردن سطح استفاده نادرست از اینترنت، با استفاده از روش‌های فیلترینگ، سایت‌های غیرمرتبط با موضوعات کاری و یا اهداف از پیش تعریف شده را مسدود نمایند. بسیاری از فیلترهای وب از الگوریتم ساده کنترل وجود کلمات کلیدی^۲ ممنوعه در کلمات تشکیل دهنده درخواست کاربر و

URL (Uniform Resource Locator) استفاده می‌کنند. اگر چه این روش دارای مزایایی از جمله سادگی و سرعت مناسب می‌باشد، ولی قدرت تشخیص مفهوم یک کلمه در جمله را ندارد و همین امر سبب می‌شود تا تعداد زیادی از صفحات مجاز مسدود شوند. به طور مثال کلمه "essex" که نام یک دانشگاه و چندین شهر در نقاط مختلف دنیا است، به سبب داشتن حروف کلمه sex، ممنوعه منظور شده و فیلتر، آن را غیرمجاز تشخیص داده و صفحه مربوطه را مسدود می‌کند و یا ممکن است کاربر با یک هدف معتبر از چند کلمه مربوطه برای جستجوی موضوع مورد نظر خود استفاده کند. به طور مثال بسیار طبیعی است که یک پزشک برای تحقیق در موضوع تاثیر جنسیت در بیماری‌های قلبی از واژه sex در جستجوهایش استفاده نماید، در صورتیکه فیلتر، آن را غیرمجاز تشخیص داده و صفحه مربوطه را مسدود می‌نماید. از طرف دیگر ممکن است به دلیل عدم وجود کلمات کلیدی ممنوعه در درخواست کاربر و URL مربوط به یک صفحه غیر مجاز، فیلتر این صفحه را مجاز شناخته و اجازه عبور دهد.

در این مقاله روشی ارائه شده است که با تحلیل رفتار بازدیدکنندگان وب و با استفاده از تکنیک‌های کاوش در وب، پالایش دسترسی به وب را بهبود می‌بخشد. با توجه به موارد فوق، روش ارائه شده، برای فیلترهایی کارایی خواهد داشت که از روش کنترل کلمات ممنوعه استفاده می‌کنند.

۳. کاوش رفتار در وب

رشد روزافزون صفحات وب باعث پدید آمدن حجم وسیعی از اطلاعات غیر ساخت یافته در شبکه اینترنت شده است [۴] که تا آوریل ۲۰۰۸ بیش از ۱۶۵/۷ میلیون سایت وب به صورت فعال مشغول به کار شده‌اند [۵]. در چنین وضعیتی کاوش هدفمند در وب به‌عنوان یک مسئله اساسی در دنیای وب مطرح و به یکی از محورهای تحقیقاتی تبدیل شده است.

کاوش در وب^۳ را می‌توان به سه محور کاوش محتوای وب^۴، کاوش ساختار وب و کاوش رفتار در وب^۵ دسته‌بندی کرد [۶]. در این مقاله بر روش کاوش رفتار در وب تمرکز شده و در ادامه به طور مختصر تشریح می‌گردد.

در تکنیک کاوش رفتار در وب با استفاده از داده‌کاوی در

۱. یکی از شرکت‌های نرم‌افزارهای فیلترینگ که اغلب فعالیت‌های اینترنتی را گزارش می‌کند.

۲. کلمه کلیدی به کلماتی می‌گوییم که وجود آن‌ها در درخواست کاربر و URL احتمالاً نشانگر سمت و سوی کاربر به یک موضوع ممنوعه (غیر مرتبط با اهداف سازمان‌ها و شرکت‌های مربوطه) می‌باشد.

3. Web Mining
4. Web Content Mining
5. Web Usage Mining

آماده‌سازی داده‌ها است. در این مرحله داده‌های فایل گزارش از اطلاعات غیر ضروری و غیر مفید پاکسازی شده و اطلاعات مورد نیاز به صورت یکپارچه جمع‌آوری و در پایگاه داده مناسبی ذخیره می‌گردد [۸].

مرحله دوم که اصلی‌ترین مرحله کاوش رفتار در وب می‌باشد، استخراج الگوی رفتار کاربر است [۹]. در این مرحله براساس اطلاعات دسترسی کاربر و با توجه به اهداف مربوط به تحلیل، اطلاعات مورد نظر به‌طور دقیق مورد بررسی قرار گرفته و الگوی رفتاری کاربر استخراج می‌شود. از مشکل‌ترین قسمت‌های این مرحله شناخت نشست کاربر است. معمولاً اتصال کاربران به سایت‌های عمومی بدون نام و مشخصات است. همچنین نمی‌توان به آدرس IP او نیز متکی بود. زیرا بخشی از کاربران به واسطه ISPها و یا ICPها به اینترنت متصل شده و ممکن است آدرس IP این واسطه‌ها به جای آدرس مبدأ ارسال شده باشد. در این صورت جهت شناخت کاربر باید از

فایل‌های گزارش^۱ سرورهای پروکسی و یا میزبان، سعی در درک رفتار کاربر و پیش‌بینی آن در آینده می‌شود. تحلیل اطلاعات فایل‌های گزارش و داده‌های ثبت شده کاربران، می‌تواند اطلاعات مهم و قابل توجهی برای صاحبان سایت‌ها و سازمان‌های مربوطه فراهم کند.

۳-۱. ساختار فایل گزارش

جدول (۱) نمونه‌ای از فایل گزارش سرورهای پروکسی را نشان می‌دهد [۷].

هر سطر فایل گزارش، مربوط به دسترسی یک کاربر می‌باشد. مشخصه‌های معمول این فایل در جدول (۲) آمده است که برای آگاهی بیشتر از توضیحات این مشخصه‌ها می‌توان به [۵] رجوع نمود.

۳-۲. مراحل کاوش رفتار در وب

کاوش رفتار در وب شامل سه مرحله اصلی است. مرحله اول،

جدول ۱. نمونه‌ای از فایل گزارش دسترسی

2006-02-01 00:08:43 1.2.3.4 - GET /classes/cs589/papers.html - 200 9221 HTTP/1.1 maya.cs.depaul.edu Mozilla/4.0+(compatible;+MSIE+6.0;+Windows+NT+5.1;+SV1;+.NET+CLR+2.0.50727) http://dataminingresources.blogspot.com/
2006-02-01 00:08:46 1.2.3.4 - GET /classes/cs589/papers/cms-tai.pdf - 200 4096 HTTP/1.1 maya.cs.depaul.edu Mozilla/4.0+(compatible;+MSIE+6.0;+Windows+NT+5.1;+SV1;+.NET+CLR+2.0.50727) http://maya.cs.depaul.edu/~classes/cs589/papers.html
2006-02-01 08:01:28 2.3.4.5 - GET /classes/ds575/papers/hyperlink.pdf - 200 318814 HTTP/1.1 maya.cs.depaul.edu Mozilla/4.0+(compatible;+MSIE+6.0;+Windows+NT+5.1) http://www.google.com/search?hl=en&lr=&q=hyperlink+analysis+for+the+web+survey
2006-02-02 19:34:45 3.4.5.6 - GET /classes/cs480/announce.html - 200 3794 HTTP/1.1 maya.cs.depaul.edu Mozilla/4.0+(compatible;+MSIE+6.0;+Windows+NT+5.1;+SV1) http://maya.cs.depaul.edu/~classes/cs480/
2006-02-02 19:34:45 3.4.5.6 - GET /classes/cs480/styles2.css - 200 1636 HTTP/1.1 maya.cs.depaul.edu Mozilla/4.0+(compatible;+MSIE+6.0;+Windows+NT+5.1;+SV1) http://maya.cs.depaul.edu/~classes/cs480/announce.html
2006-02-02 19:34:45 3.4.5.6 - GET /classes/cs480/header.gif - 200 6027 HTTP/1.1 maya.cs.depaul.edu Mozilla/4.0+(compatible;+MSIE+6.0;+Windows+NT+5.1;+SV1) http://maya.cs.depaul.edu/~classes/cs480/announce.html

پایانی^۴ و زمان دیدن^۵ استفاده شده است که در این روش، صفحه پایانی فراوانی صفحات مورد دستیابی، محاسبه و در یک جدول قرار می‌گیرد. جدول (۳) یک نمونه از میزان فراوانی صفحات را نشان می‌دهد. عامل EDF^۴ تعداد مراجعات به صفحه مربوطه در آخرین نشست کاربر را نشان می‌دهد.

جدول ۳. محاسبه فراوانی صفحات پایانی

ID.P	EDF	Page Name
P _۱	۱۲	/staff/dc/hcow/menu.html
P _{۲۴}	۲۷	/welcome.ht
P _{۲۹}	۵	/courses/Undergraduate/elective.htm
P _{۳۶}	۱۸	/staff/dc/hcow/links.html
P _{۳۹}	۲۹	/courses/postgradindex.htm
P _{۵۹}	۴	/staff/khui/teaching/cm3016/
...

عامل دیگری به نام IV^۶ تعریف شده که تعداد دفعاتی که یک صفحه دیده شده است را نشان می‌دهد. با تعریف این عامل، شرط $EDF/IV \geq$ برای پذیرفتن صفحات، تعیین شده است. این شرط اجبار می‌کند که حداقل نیمی از صفحات مورد دستیابی در آخرین نشست، مربوط به صفحه پایانی باشد که در جدول (۴)، مقادیر عامل‌ها و نتیجه شرط مشخص شده است. با توجه به جدول (۴) صفحات P36 و P59 به عنوان صفحات هدف گزینش می‌شوند.

جدول ۴. محاسبه نتایج شروط

ID.P	EDF	IV	EDF/IV	Con.
P _۱	۱۲	۱۰۲	۰/۱۲	No
P _{۲۴}	۲۷	۱۲۸	۰/۰۲	No
P _{۲۹}	۵	۱۱	۰/۰۵	No
P _{۳۶}	۱۸	۳۱	۰/۵۸	Yes
P _{۳۹}	۲۹	۶۷	۰/۴۳	No
P _{۵۹}	۴	۸	۰/۵۰	Yes
...

پس از آنکه صفحات هدف مشخص شدند، باید مسیره‌ها استخراج شوند. منظور از مسیره‌ها، استخراج الگوی دسترسی به صفحات است.

تکنیک‌های دیگری مانند دسته‌بندی نشست‌های کاربران استفاده کرد [۱۰].

مرحله سوم، تحلیل الگوهای استخراج شده است. در این مرحله الگوهای استخراج شده، مورد بازبینی قرار گرفته و رفتار کاربر مورد تحلیل قرار می‌گیرد. به‌عنوان مثال درمی‌یابیم که کاربر پیگیر موسیقی بوده و همچنین به چه نوع موسیقی علاقه‌مند است [۱۱].

جدول ۲. مشخصه‌های اصلی فایل گزارش دسترسی

Remote Host	ng21.exabot.com
Login Name	-
Authorised User	-
Date / Time	[01/Oct/2004:22:27:46 +0100]
Request	"Get /staff/asga/context_learn HTTP/1.1"
Status	301
Bytes Sent	388

۴. کارهای انجام شده مرتبط

گانزن باتامالای، مقاله وبسایت‌های خودسازگار را بر اساس تحقیقات مربوط به پایان‌نامه کارشناسی ارشد خود در دانشگاه رابرت گاردن^۲ را در سال ۲۰۰۸ ارائه کرده است [۵]. او با استفاده از کاوش گزارش‌های دسترسی، الگوی رفتار کاربر را استخراج کرده و بر اساس آن سعی در یافتن سریع‌تر صفحات هدف مورد نظر کاربر کرده است.

در این مقاله روشی ارائه شده که راه‌های کوتاه‌تری برای دسترسی به صفحات هدف شناسایی کرده و ضمن پاکسازی فایل گزارش از اطلاعات ناخواسته، جلسه‌بندی^۳ کاربران صورت می‌گیرد. منظور از جلسه‌بندی، شناسایی و جداسازی آن‌دسته از صفحاتی است که کاربر در یک نشست در طی اتصال به وبسایت، به آن صفحات دسترسی می‌یابد.

در این مقاله، زمان متوسط صرف شده یک کاربر بر روی یک صفحه حدود ۴۰ ثانیه محاسبه شده و بر اساس نمودارهای ارائه شده، زمان متوسط یک نشست ۳ دقیقه برآورد شده است. این زمان‌ها برای بدست آوردن مجموعه صفحات مرتبط به هم در یک نشست کاربر لازم است.

جهت شناسایی و یافتن صفحه هدف کاربر، از روش‌های مدرک

4. End Document Method
5. Browsing Time Method
6. End Document Frequency
7. Incoming Visits

1. Session Clustering
2. Robert Gordon University
3. Sessionization

P5 در مثال اول و صفحات P15 و P9 در مثال دوم به عنوان نقاط دسترسی سریع منظور شده و ایجاد می‌شوند.

۵. طرح مسئله و ارائه راه حل

پالایش صفحات وب بر اساس کلمه‌های کلیدی ممنوعه، یک راه حل افراطی بوده و باعث پایین آمدن کارایی و نارضایتی کاربران می‌گردد. ممکن است هدف کاربران از به کار بردن کلمات کلیدی ممنوعه، ورود به صفحات ممنوعه نباشد، ولی چون ملاک ما فقط کلمات ممنوعه است، به صورت افراطی از دسترسی کاربر به صفحات مورد نظرش جلوگیری به عمل می‌آید. اگر بتوان سابقه رفتاری کاربر را استخراج کرده و با بررسی درخواست جاری وی، قصد او را از این درخواست تخمین زده و عکس‌العمل مناسب‌تری نشان دهیم، توانسته‌ایم یک روش عادلانه‌تری را ارائه دهیم. به همین منظور در ادامه معماری تحلیل رفتار کاربر ارائه شده و تشریح می‌گردد.

۶. معماری تحلیل رفتار کاربر

شکل (۳)، معماری پالایش درخواست‌های کاربر، بر اساس رفتار او را نشان می‌دهد. بر اساس این معماری، در سه مرحله عمل پالایش صورت می‌پذیرد:

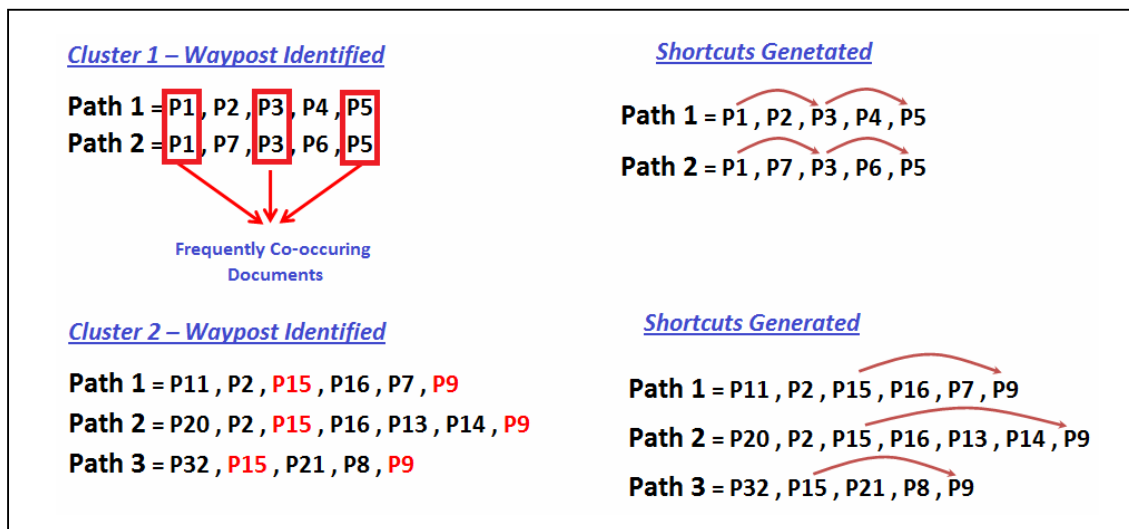
Target Documents Identified = P1,P4,P8,P20

User ID	Browsing Session
U1S1	P1 , P2 , P3 , P4 , P5 , P6
U2S1	P5 , P9 , P8 , P13 , P14 , P17 , P20
U3S1	P16 , P7 , P21 , P22 , P101 , P102 , P60

Path No	Path Sequence
1	P1 , P2 , P3 , P4
2	P5 , P9 , P8
3	P8 , P13 , P14 , P17 , P20

شکل ۱. مثالی از استخراج مسیرها

برای این منظور از میان صفحات دسترسی شده توسط هر کاربر، مسیرهایی انتخاب می‌شوند که ابتدا و انتهای آن‌ها جزء صفحات هدف باشند. شکل (۱) مثالی از انتخاب مسیرها است. در نهایت با استفاده از الگوریتم k-means مسیرها دسته‌بندی می‌شوند. سپس با استفاده از آن، گره‌هایی برای دسترسی‌های سریع استخراج می‌شود که شکل (۲) مثال‌هایی از ایجاد این نقطه‌های دسترسی‌های سریع را نشان می‌دهد. همان‌طور که در شکل (۲) دیده می‌شود، صفحات P1، P3 و



شکل ۲. مثالی از استخراج نقاط دسترسی سریع

کاربرد آن‌ها در جستجوهای ممنوعه، صورت پذیرد. وزن‌ها در جدولی مانند جدول (۵) ذخیره می‌شوند.

جدول ۵. نمونه‌ای از جدول وزن‌دهی کلمات کلیدی

ردیف	کلمه	وزن
۱	sensation	۳۰
۲	bosom	۳۰
۳	Sex	۷۰
۴	Bare	۷۰
۵	couple	۵۰
۶	plump	۵۰
۷	whore	۹۰
۸	girl	۳۰
۹	porn	۸۰
۱۰	adult	۶۰

همچنین در این مرحله باید توسط مدیر سیستم، قوانین مورد نیاز در بانک قوانین ثبت شود. این بانک، قوانینی مانند قانون وزن‌دهی به کاربران در گروه‌های مختلف و یا قوانین مربوط به تعیین مجوز دستیابی با توجه به رفتار کاربر را در خود ذخیره می‌نماید. از دیگر قوانینی که می‌توان در بانک قوانین ثبت کرد، اطلاعات مربوط به مرزبندی گروه‌های بازدیدکننده وب و نحوه تعیین آن‌ها خواهد بود که در بخش بعدی توضیح داده می‌شوند.

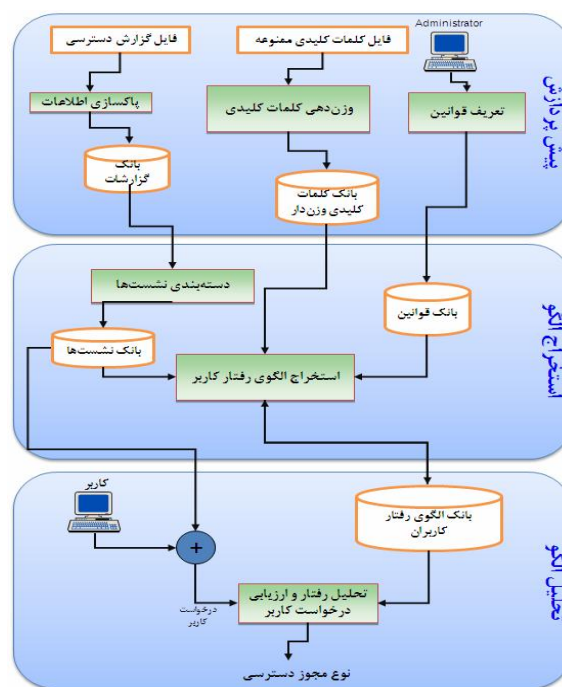
مرحله دوم، مرحله استخراج الگوی رفتار است. در این مرحله باید برای شناسایی کاربران و انتساب کد شناسایی، نشست‌های ایشان را دسته‌بندی کرده و در پایگاه داده مربوطه ذخیره نمود. بر اساس اطلاعات موجود در این پایگاه و پایگاه‌های کلمه‌های کلیدی وزن‌دار، قوانین و همچنین بر اساس سابقه رفتار قبلی کاربر، رفتار وی مورد بررسی قرار گرفته و پس از استخراج الگوی رفتاری، در بانک مربوطه ذخیره می‌گردد. چگونگی استخراج الگوی رفتار کاربران در بخش ۳ توضیح داده شده است.

در مرحله تحلیل الگو، پس از استخراج الگوی رفتاری کاربر، درخواست‌های او بر اساس رفتار گذشته ارزیابی شده و اجازه دسترسی به صفحه مورد درخواست مشخص می‌شود.

۷. تحلیل رفتار بازدیدکنندگان

هدف این بخش تبیین چگونگی تحلیل رفتار بازدیدکنندگان

مرحله اول پیش پردازش است که در این مرحله باید سه پردازش اساسی صورت بگیرد؛ فایل گزارش^۱ باید به صورت یک بانک اطلاعات آماده‌سازی و قابل استفاده گردد. به همین منظور باید داده‌ها را از داده‌های غیر مفید پاکسازی نمود. در این پالایش ضمن تبدیل اطلاعات قبلی به شکل ساخت‌یافته، فیلدهای مورد نیاز ذخیره شده و فیلدهای غیر ضروری دور ریخته می‌شود.



شکل ۳. معماری پالایش درخواست‌ها بر اساس رفتار کاربر

هنگامی که کاربر یک صفحه خاص را دریافت می‌کند، ممکن است در فایل گزارش سرور چندین سطر گزارش مربوط به این درخواست ثبت شود و در بیشتر مواقع ممکن است فقط یک سطر این اطلاعات مورد نیاز باشد که در این صورت باید این اطلاعات اضافه نیز حذف گردد. به‌عنوان مثال اطلاعات مربوط به فایل‌های غیر مرتبط مانند فایل‌های تصویری می‌تواند حذف گردد.

در همین مرحله باید به کلمه‌های کلیدی وزن خاصی را نسبت داده و در پایگاه داده مربوطه نگهداری نمود. این وزن‌دهی توسط مدیر سیستم فیلترینگ انجام می‌شود و باید متناسب

1. Log File

گروه‌های فوق نسبت داده می‌شود. این عدد که وزن رفتاری نامیده می‌شود، در تحلیل رفتار کاربر به ما کمک خواهد کرد. برای اختصاص وزن باید از قوانینی که قبلاً تعریف شده و در بانک قوانین موجود است، استفاده کنیم.

برای تعیین رفتار کاربر باید به چند نکته توجه کرد؛ نخست اینکه یک کاربر می‌تواند رفتار کاملاً هنجار و یا رفتار کاملاً ناهنجار داشته باشد (رفتار کاملاً هنجار را معادل وزن صفر در نظر گرفته و رفتار کاملاً ناهنجار معادل وزن صد تعریف می‌گردد).

دوم اینکه برای تعیین رفتار (W) یک فرد در زمان t باید رفتار گذشته (W) او را نیز لحاظ کنیم.

سوم اینکه هر چند سابقه رفتار فرد در تعیین رفتار جدید فرد تاثیر دارد، اما تغییر رفتار او نیز باید بتواند به سرعت وضعیت رفتاری او را اصلاح نماید. باید توجه داشت که ما قصد نداریم کاربر را به لحاظ رفتاری ارزیابی نماییم، بلکه به دنبال این هستیم که رفتار زمان حال او را با توجه به گذشته‌اش بهتر درک کنیم. منظور از رفتار گذشته، گذشته نزدیک در همان نشست بوده و به رفتار کاربر در نشست‌های قبلی توجه نخواهیم کرد.

با توجه به موارد فوق، فرمول (۱) پیشنهادی برای محاسبه وزن رفتاری کاربر در زمان t می‌باشد:

$$W(t) = aW(t-1) + b \cdot \sum_i^n (1) \quad (1)$$

$W(t)$ رفتار جاری کاربر و وزن کلمات کلیدی ممنوعه بوده که کاربر در درخواست جاری خود استفاده کرده و از جدول مربوط به وزن کلمات کلیدی استخراج می‌گردد و $W(t-1)$ رفتار محاسبه شده او در آخرین دسترسی‌اش به صفحات وب می‌باشد. ضرایب a و b عددی بین ۰ و ۱ است. ضریب a میزان تأثیر رفتار گذشته فرد در محاسبه رفتار جاری او و ضریب b سهم تأثیر کلمات کلیدی به‌کار رفته در درخواست جاری او را در تعیین رفتار جاری فرد نشان می‌دهد. مقدار n نشان دهنده تعداد کلمات کلیدی منظور شده یک درخواست می‌باشد. همان‌طور که قبلاً اشاره شد، مقدار W حداکثر ۱۰۰ می‌باشد و اگر مقدار محاسبه شده در مواردی بیش از عدد ۱۰۰ به‌دست آید، باید آن را ۱۰۰ در نظر گرفته و در محاسبات مرحله بعد لحاظ کرد.

وب بر اساس اطلاعات گزارش دسترسی کاربران می‌باشد. جهت بهبود فیلتر کردن سایت‌ها از سابقه رفتاری هر کاربر کمک گرفته و معین می‌گردد که دسترسی جاری کاربر با چه هدفی صورت گرفته و مجاز و یا غیر مجاز بودن دسترسی اعلام می‌شود. به همین منظور کاربران بر اساس رفتار کاربران در مراجعه به صفحات وب در سه دسته‌بندی تعریف می‌گردند:

گروه کاربران با رفتارهای هنجار، ناهنجار و نامشخص. رفتارهای هنجار به آن دسته از رفتارهای کاربر گفته می‌شود که از دید سیستم فیلترینگ، جستجو و گشت و گذار عادی در وب باشد. اگر کاربر از روال عادی خارج شده و سعی در دسترسی به صفحات ممنوعه را داشته باشد به‌عنوان کاربر با رفتار ناهنجار تعریف می‌شود. به‌عنوان مثال کاربر به صفحات مربوط به معاملات مواد مخدر و یا در ساعات اداری به سایت‌های غیر مرتبط با کار متصل گردد.

گروه کاربران با رفتار نامشخص به کاربرانی گفته می‌شود که رفتار آنان به روشنی مشخص نشده و در حال بررسی است. این کاربران گاهی دارای رفتار عادی و گاهی نیز دارای رفتار غیر عادی بوده و نمی‌توانیم تشخیص دهیم این کاربران دارای رفتار هنجار و یا ناهنجار هستند.

۷-۱. گروه‌بندی کاربران

در ابتدای ورود کاربر به سیستم، برای او وزن رفتاری اولیه منظور شده به نحوی که در گروه هنجار دسته‌بندی شود. منظور از وزن رفتاری، عددی است بین صفر تا صد که نشانگر رفتار کاربر می‌باشد. هرچه این عدد به صد نزدیک‌تر باشد نشانگر رفتار ناهنجار بیشتر و هر چه به صفر نزدیک‌تر باشد، نشانگر ناهنجاری کمتر و در واقع رفتار سالم‌تر و هنجارتر می‌باشد.

در صورتی که کاربری از گروه هنجار صفحه ممنوعه‌ای را تقاضا بکند، آنگاه به او اجازه دیدن صفحه درخواستی داده می‌شود، ولی رفتار او به صورت دقیق‌تری مورد بررسی قرار می‌گیرد. در صورتی که کاربر، اصرار به ادامه رفتارهای ناهنجار را داشته باشد، وزن رفتاری او در گروه ناهنجار افزوده شده و دسترسی‌ها برای وی محدودتر می‌گردد.

۷-۲. الگوریتم اختصاص وزن رفتاری

با توجه به رفتار کاربر، یک عدد متناسب با حضور او در

فقط یک کلمه ممنوعه با وزن ۳۰ وجود داشته باشد، وزن او ابتدا به ۸۰ و سپس به ۶۵ کاهش می یابد.

۸. استفاده از تحلیل رفتار برای پالایش

دسترسی کاربران به وب

همان طور که در فصل قبل بیان گردید، می توان با استفاده از فایل های گزارش دسترسی کاربران به صفحات مختلف، رفتار آنان را مورد ارزیابی قرار داده و جهت استفاده در تصمیم گیری های بعدی نگهداری کرد. با توجه به دسته بندی کاربر و وزن رفتاری او در هر گروه، می توان براساس معماری ارائه شده، درخواست او را مورد ارزیابی قرار داده و مجاز و یا غیر مجاز بودن درخواست او را مشخص نمود.

در تحلیل رفتار کاربر باید به مواردی که باعث می شود یک دسترسی، به اشتباه مجاز و یا غیر مجاز اعلام شود، توجه کرد. وزن دهی کلمات تا حدی باعث می شود درخواست های هنجاری که دارای کلمات کلیدی ممنوعه با وزن خفیف می باشند، وزن رفتاری کاربر را در گروه ناهنجار بالا نبرده و حتی رفتار آرام او، وزن را کاهش می دهد. باید یک حد آستانه برای مجاز بودن دسترسی، در بانک قوانین تعریف کرد که براساس آن، نوع مجوز دستیابی بر اساس اطلاعات تحلیل رفتار کاربر تعیین می گردد.

۹. بررسی نتایج یک آزمایش

برای ارزیابی معماری ارائه شده و بررسی نتایج آن، ابتدا باید داده های اولیه مناسبی تهیه گردد. به همین منظور داده هایی به سه صورت منتخب، تصادفی و واقعی فراهم، مورد آزمایش و نتایج آن، بررسی و مورد تحلیل قرار گرفت. در ادامه، این نتایج بیان می شوند.

برای بررسی بهتر ابتدا برای پنج کاربر فرضی داده هایی تولید شده و بر اساس آن، رفتار کاربر محاسبه گردید. جدول (۶) این داده های فرضی را نمایش داده است. کاربر U1 یک کاربر فرضی است که در ۱۵ درخواست پیاپی خود از هیچ کلمه ممنوعه ای استفاده نکرده است. لذا وزن کلمه های کلیدی مورد استفاده، صفر منظور می شود.

کاربر فرضی U2 کاربری است که در همه ۱۵ درخواستش از حداقل ۳ کلمه کلیدی ممنوعه با وزن ۱۰۰ استفاده نموده است. کاربر U3 کاربری است که داده های آن به صورت دستی و

به عنوان مثال اگر تعداد کلمات کلیدی در یک درخواست که مؤثر در تعیین رفتار جاری کاربر باشد را ۳ فرض کرده و ضرایب a و b را به ترتیب $\frac{3}{4}$ و $\frac{1}{6}$ منظور نماییم، فرمول فوق به صورت فرمول (۲) در می آید.

$$W(t) = \left(\frac{3}{4}W(t-1) + \frac{1}{6}(W1+W2+W3) \right) \quad (2)$$

در فرمول (۲)، W نشانگر آخرین وزن به دست آمده برای هر کاربر در زمان t و در نشست جاری می باشد و فقط ۳ کلمه از کلمات کلیدی ممنوعه که کاربر به کار می برد محاسبه می شود. یعنی اگر تعداد کلمات کلیدی ممنوعه که کاربر به کار می برد بیش از سه کلمه باشد، فقط سه کلمه از آن انتخاب می شود. این سه کلمه ممنوعه می تواند سه کلمه اول ممنوعه بوده و یا آن کلماتی انتخاب شوند که بیشترین وزن را دارا می باشند. در ابتدای کار، مقدار ۵۰ به عنوان وزن اولیه به رفتار کاربر اختصاص داده می شود. کاربری که در اولین و دومین درخواست او هیچ واژه ممنوعه ای وجود نداشته باشد، وزن او به صورت فرمول های (۳) و (۴) محاسبه می گردد.

$$W = TRUNC \left(\frac{3 \times 50}{4} + \frac{0+0+0}{6} \right) = 37 \quad (3)$$

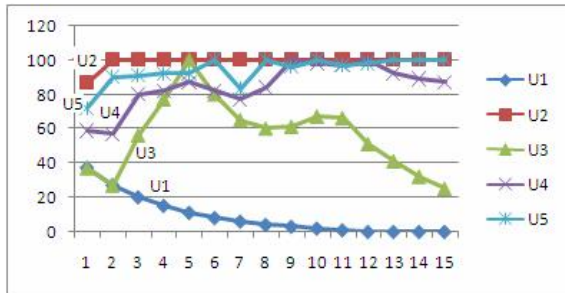
$$W = TRUNC \left(\frac{3 \times 37}{4} + \frac{0+0+0}{6} \right) = 27 \quad (4)$$

و اگر در درخواست سوم او سه کلمه ممنوعه با وزن های ۹۰، ۶۰ و ۷۰ وجود داشته باشد، وزن او در قسمت کلمات ممنوعه ۵۶ محاسبه می شود. اگر در درخواست چهارم، کلمات ممنوعه با وزن های ۶۰، ۷۰ و ۸۰ و در درخواست پنجم، کلمات ممنوعه با وزن های ۹۰، ۹۰ و ۹۰ را به کار برند، وزن رفتاری کاربر ۷۷ محاسبه می گردد.

لازم به ذکر است در محاسبه آخر، مقدار بدست آمده بزرگ تر از ۱۰۰ بوده که باید این مقدار کوچک تر و یا مساوی ۱۰۰ شود.

همان طور که دیده می شود سعی بیشتر برای جستجوی کلمات ممنوعه و یا آدرس دهی آن ها، وزن رفتار کاربر را در گروه کلمات ممنوعه بالا می برد. اگر کاربر رفتار خود را اصلاح کرده و کلمات ممنوعه کمتری به کار برده و یا اصلا به کار نبرد، وزن او نیز رو به کاهش می رود. به طور مثال اگر در دو درخواست بعدی

مقدار تصادفی فرض کرده که عددی بین ۰ تا ۱۰۰ است. بنابراین کاربرانی که رفتار ناهنجار خود را به هر شکل ادامه دهند، شناسایی شده و نمودار آنها بر خط ۱۰۰ منطبق می‌شود.



نمودار ۱. نمودار رفتاری ۵ کاربر فرضی

جدول ۶. داده‌های ایجاد شده برای ۵ کاربر فرضی

U1	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰
U2	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰
U3	۰	۰	۹۰	۶۰	۹۰	۳۰	۳۰	۵۰	۶۰	۴۰	۱۰۰	۱۰۰	۹۰	۹۰	۸۰	۹۰	۹۰	۸۰	۸۰
U4	۱۲	۱۶	۶۸	۲۰	۱۷	۲۵	۵۵	۳۳	۷۶	۱۲	۳۱	۹۰	۳۹	۴۶	۲	۵۷	۶۰	۵۶	۷۵
U5	۹۲	۹۰	۱۲	۵۵	۶۴	۵۱	۱۷	۸۵	۸۵	۸۳	۵۴	۴۹	۳۷	۹۲	۴۵	۴۴	۳۶	۵۶	۷۷

انتخابی تولید شده و کاربرهای U4 و U5 کاربرهایی هستند که به هر سه درخواست ایشان یک مقدار تصادفی نسبت داده شده است. محاسبه وزن رفتاری این ۵ کاربر بر اساس فرمول (۲) در جدول (۷) آورده شده است.

جدول ۷. رفتار محاسبه شده برای ۵ کاربر فرضی

U1	۳۷	۲۷	۳۰	۱۵	۱۱	۸	۶	۴	۳	۲	۱	۰	۰	۰	۰	۰	۰	۰	۰
U2	۸۷	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰
U3	۳۷	۲۷	۵۶	۷۷	۱۰۰	۸۰	۶۵	۶۰	۶۱	۶۷	۶۶	۶۶	۷۷	۹۶	۱۰۰	۵۹	۵۷	۸۰	۸۲
U4	۵۹	۵۷	۸۰	۸۲	۸۷	۸۲	۷۷	۸۴	۱۰۰	۹۸	۹۷	۱۰۰	۹۲	۸۹	۸۷	۵۹	۵۷	۸۰	۸۲
U5	۷۲	۹۰	۹۱	۹۲	۹۲	۱۰۰	۸۳	۱۰۰	۹۶	۱۰۰	۹۷	۹۸	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰

جدول ۸. داده‌های ۵ کاربر واقعی

U6	۷۰	۷۰	۷۰	۷۰	۷۰	۲۰	۲۰	۲۰	۲۰	۳۰	۲۰	۰	۰	۰	۰	۶۰	۶۰	۷۰	۷۰
U7	۶۰	۶۵	۶۵	۶۵	۶۵	۶۵	۶۵	۰	۶۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰
U8	۰	۰	۱۰۰	۸۰	۱۰۰	۰	۰	۰	۶۵	۶۵	۵۵	۶۵	۶۵	۶۵	۰	۰	۰	۰	۰
U9	۰	۵۰	۰	۰	۵۰	۵۰	۰	۵۰	۵۰	۵۰	۵۰	۵۰	۵۰	۵۰	۰	۷۰	۰	۰	۰
U10	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰

برای بررسی بیشتر، اطلاعات دسترسی به صفحات وب مربوط به ۵ کاربر از فایل گزارش‌های سایت AOL انتخاب شده و داده‌های آنها در جدول (۸) آورده شده است. محاسبه رفتار ایشان در جدول (۹) و نمودار رفتاری آنها در نمودار (۲) آمده است.

جدول ۹. وزن رفتار محاسبه شده برای ۵ کاربر واقعی

U6	۴۹	۴۸	۴۷	۴۶	۴۶	۴۹	۵۱	۵۳	۵۴	۴۵	۵۳	۳۹	۲۹	۴۳	۵۳	۰	۰	۰	۰
U7	۵۹	۵۵	۵۲	۴۹	۴۷	۴۶	۴۵	۳۳	۴۶	۵۶	۴۲	۳۱	۲۳	۱۷	۱۲	۰	۰	۰	۰
U8	۳۷	۲۷	۳۶	۴۰	۴۶	۳۴	۲۵	۱۸	۱۳	۲۰	۲۵	۲۷	۴۱	۴۱	۰	۰	۰	۰	۰
U9	۳۷	۳۶	۳۷	۲۰	۲۳	۲۵	۱۸	۲۱	۲۴	۲۶	۲۷	۲۸	۲۹	۲۱	۲۷	۰	۰	۰	۰
U10	۳۷	۳۲	۲۹	۲۱	۱۵	۱۱	۸	۱۱	۱۳	۱۴	۱۵	۱۶	۱۲	۹	۶	۰	۰	۰	۰

نمودار (۱)، نمودار این داده‌ها را بر اساس وزن رفتاری کاربر نسبت به زمان، مشخص کرده است. منحنی‌های پایینی و بالایی نمودار (۱) در واقع حد پایینی و حد بالایی داده‌ها را مشخص ساخته است. منحنی U1 مربوط به کاربری است که هیچ کلمه ممنوعه‌ای به کار نبرده است و لذا منحنی کم‌کم به سمت صفر میل کرده و منطبق بر خط صفر ادامه یافته است. یعنی U1 یک کاربر با رفتار کاملاً هنجار می‌باشد. منحنی U2 مربوط به کاربری است که بر خلاف کاربر اول، رفتاری کاملاً ناهنجار داشته و به سرعت بر خط ۱۰۰ منطبق شده است. نمودار U3 کاربری است که پس از ارتباط اولیه، رفتارهای ناهنجاری داشته و پس از چند دسترسی رفتارش اصلاح شده و به صورت متعادل درآمد است. منحنی‌های U4 و U5 مربوط به کاربران با رفتار تصادفی است. البته بهتر است گفته شود با رفتار ناهنجار تصادفی، زیرا برای هر یک، حتماً سه

آن و تعیین مجوز دستیابی کاربر، زمان زیادی صرف نموده و باعث تاخیر چشمگیر در پاسخ به کاربران می شود. اگر در هنگام پیاده سازی به مسایل فوق توجه نگردد، ممکن است سیستم شبکه، ناکارآمد شود. طبق معماری ارائه شده، مجوز عبور درخواستی کاربر، بر اساس اطلاعات بانک رفتار وی صادر می گردد. با توجه به اینکه هنگام درخواست، استخراج الگو و کلاس بندی کاربران صورت نمی پذیرد، پس این تاخیرها شامل زمان پاسخ نمی شود. اما تطبیق کاربر با نشست مربوطه و دسترسی به بانک رفتار همچنان جزء مسایلی خواهد بود که باید در پیاده سازی، آن را به حداقل رساند.

۱۰-۲. خطا در گذر از فیلتر

از دیگر چالش های مطرح، اشتباه در مسدود ساختن یک صفحه و یا صدور مجوز عبور می باشد. در روش ارائه شده، رفتار کاربر با مجموعه ای از عوامل در حین درخواست و سابقه رفتار و نیز وضعیت او در درخواست جدید سنجیده می شود. ممکن است این نتیجه بنا به دلایلی از جمله به روز نبودن بانک کلمات کلیدی ممنوعه و یا عدم دقت در تعیین ضرایب مربوط به فرمول رفتار کاربر، با واقعیت تفاوت داشته باشد.

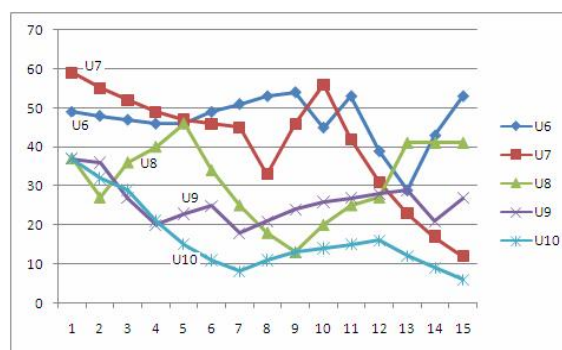
۱۰-۳. وزن دهی کلمات کلیدی ممنوعه

در وزن دهی کلمات ممنوعه، مشکل اصلی، چند ماهیتی بودن کلمات می باشد. یعنی یک کلمه ممکن است به عنوان یک واژه ممنوعه به کار برده شده ولی به عنوان یک کلمه عادی استفاده شود. وزن دهی کلمات می تواند تا حدی ماهیت کلمات را مشخص نماید. اگر کلمات را به صورت صحیحی وزن دهی نماییم و با توجه به اینکه نوعی میانگین یابی در تعیین وزن ها وجود دارد، تا حدی باعث تاثیر ماهیت کلمات در وزن های به دست آمده می شود. به هر حال، راه حل فوق یک راه حل نسبی بوده و این موضوع می تواند به عنوان یک موضوع مستقل در تحقیقات آتی پیگیری شود.

۱۱. نتیجه گیری

هدف در این مقاله استفاده از تحلیل رفتار کاربران برای بهبود مسدودسازی درخواست های ممنوعه است. به همین منظور با استفاده از اطلاعات فایل گزارش دسترسی کاربران و داده کاوی آن، الگوی رفتار ایشان استخراج شده و در یک بانک اطلاعات

بر طبق نمودار رفتار کاربران، کاربر U6 از ابتدا رفتار ناهنجار داشته و سعی در ادامه آن نیز دارد. اگر در زمان هایی نیز به سمت هنجار تمایل پیدا کرده، سریعاً به حالت ناهنجار باز گشته است. کاربر U7 با رفتار ناهنجار شروع کرده و پس از مدتی دارای رفتار هنجار شده است. کاربر U8 رفتار نامشخصی داشته و به صورت شبه سینوسی بین هنجار و ناهنجار نوسان می کند. کاربر U9 رفتاری متعادل داشته و کاربر U10 کاربری آرام بوده است. بنابراین می توان گفت چنین کاربرانی همزمان دو یا چند موضوع را پیگیری می کنند و لذا رفتار ایشان نامشخص است.



نمودار ۲. نمودار رفتاری ۵ کاربر واقعی

در این آزمایش از داده های کاربران با رفتارهای افراطی و یا بسیار آرام صرف نظر شده است. در این مقاله مجال بررسی داده های بیشتر برای تحلیل مسایلی از قبیل کاربران چند رفتاری نمی باشد. کاربران چند رفتاری می تواند به عنوان یک موضوع جداگانه در تحقیقات آتی مورد بررسی قرار گیرد.

۱۰. چالش های فرارو

اگر چه پلایش وب با استفاده از تحلیل رفتار کاربران می تواند کارایی فیلترها را بهبود بخشد، ولی با چالش هایی روبرو می باشد که در ادامه تشریح می شود.

۱-۱۰. سرعت دستیابی

با توجه به اینکه سیستم فیلترینگ در مسیر اصلی ارتباط با اینترنت قرار داده می شود و درخواست های همه کاربران باید با مجوز این سیستم پاسخ داده شود، سرعت دستیابی به صفحه مورد نظر از اهمیت ویژه ای برخوردار خواهد بود. پردازش های مربوط به کلاس بندی کاربران، استخراج الگوی رفتار و تحلیل

برای بالا بردن کارایی بیشتر فیلترها، می‌توان به جای استفاده از مجوز عبور و یا عدم عبور، بر اساس درخواست کاربر و سابقه رفتاری وی مجوزهای خاصی را صادر نمود. یعنی براساس طیف اخلاقی کاربران، مجوزهای معنی‌داری را صادر کرد. به‌عنوان مثال کاربرانی با سابقه رفتاری ناهنجار با نمره ۸۰ الی ۹۰ اجازه دیدن متن را داشته ولی مجوز دیدن فایل‌های چند رسانه‌ای را نداشته باشند. این امر کارایی فیلترها را بیشتر کرده و هدفمندتر می‌گردد.

۱۳. مراجع

- [1] WebSpy Ltd, "Internet Use Statistic.", Internet: <http://www.webspy.com/resources/whitepapers/2008%20WebSpy%20Ltd%20%20Internet%20Use%20Statistics.pdf>, 2008.
- [2] Connect World, "Employee Internet Usage."; Internet: <http://www.connectworld.net/cc/employee-internet-usage.html>, It is a snapshot of the page as it appeared on 14 Dec 2010.
- [3] Covenant Eyes, "Internet Pornography Statistics."; Internet: <http://www.covenanteyes.com/support/pdfs/Covenant%20Eyes%20Pornography%20Statistics.pdf>, 2006
- [4] Yaltaghian, Behnak "Improving the ranking of Search Engine Output: A Network Analysis Approach."; In IBM Centre for Advanced Studies Conferenc, University of Toronto, Toronto, Ontario, M5S 2E4, 2002.
- [5] Bathumalai, Ganesan "Self Adapting Websites: Mining User Access Logs."; M. Sc. Thesis, The Robert Gordon University, 2008.
- [6] Kosala, Raymond; Blockeel, Hendrik "Web Mining Research: A Survey."; In ACM SIGKDD Explorations, 2000, 2(1), 1-15.
- [7] Mobasher, Bamshad "Web usage mining."; In Bing Liu, Web Data Mining: Exploring Hyperlinks, Content and Usage Data, Berlin-Heidelberg, Springer-Verlag 2007, 449-483.
- [8] Hofgesang, I. Peter "Web Usage Mining Structuring Semantically Enriched Clickstream Data."; M. Sc. Thesis, Vrije University Amsterdam, The Netherlands, 2004.
- [9] Mihara, Koichiro; Terabe, Masahiro; Hashimoto, Kazuo "A Novel Web Usage Mining Method."; In Proceedings of the Fourth International Conference on Web Information Systems and Technologies 2008.
- [10] Xu, Guandong "Web Mining Techniques for Recommendation and Personalization."; Ph. D. Thesis, Victoria University, Australia, 2008.
- [11] Liu, Ning-Han; Lai, Szu-Wei; Chen, Chien-Yi; Hsieh, Shu-Ju "Adaptive Music Recommendation Based on User Behavior in Time Slot."; In International Journal of Computer Science and Network Security 2009, 9(2), 219-227, Available on: http://paper.ijcsns.org/07_book/200902/20090229.pdf

ذخیره می‌گردد. برای کمی کردن رفتار کاربر، قوانینی ارائه شد که در آن، عواملی مانند وزن کلمات کلیدی مورد استفاده و همچنین سابقه رفتار کاربر تعیین‌کننده هستند. بر اساس نتایج بدست آمده از این محاسبات، وزن رفتاری کاربر تعیین شده و در بانک اطلاعات مربوطه ذخیره می‌گردد. در هر دسترسی کاربر، ضمن استفاده از این اطلاعات به عنوان سابقه رفتار وی، بر اساس دسترسی جدید، اطلاعات فوق به‌روز شده و ذخیره می‌گردد. با مراجعه به این بانک و بررسی واژه‌های به کار رفته توسط کاربر، تقاضای او تحلیل و ارزیابی گشته و مجوز عبور صادر شده و یا صفحه درخواست شده مسدود می‌گردد.

دوری از مسدود کردن افراطی صفحات وب، مزیت اصلی روش ارائه شده، نسبت به روش معمول مورد استفاده در فیلترها است.

بررسی چگونگی جلوگیری از اشتباهات ناشی از ماهیت کلمات کلیدی و اشتباه در گذر از فیلتر و ارائه راه‌کار مناسب، می‌تواند جزء کارهایی باشد که در آینده به بهبود پالایش صفحات منجر شود. اصلاح وزن‌دهی بر اساس اطلاعات جانبی و به‌روز کردن خودکار آن می‌تواند یک رویکرد باشد.

ارائه یک معماری که در آن با استفاده از اطلاعات دیگری مانند کلمات غیر ممنوعه که به نحوی در وزن‌دهی اثر منفی داشته و با استفاده از آن بتوان رفتارهای حرفه‌ای را، مانند آنچه در مورد جستجوی یک پزشک گفته شد، تشخیص داده و باعث بهبود درک رفتار شود، می‌تواند جزء کارهای آتی باشد.

اصلاح هوشمندانه وزن‌دهی کلمات کلیدی بر اساس رفتار کاربر می‌تواند به‌عنوان محور تحقیقی جدیدی مطرح و پیگیری گردد. همان‌طور که قبلاً بیان شد، در ابتدا باید در جدول وزن‌دهی کلمات، به هر کلمه کلیدی وزنی متناسب کاربرد آن در جستجوهای ممنوعه اختصاص یابد که این کار به صورت دستی و بر اساس بررسی و تحلیل مدیر سیستم فیلتر صورت می‌پذیرد و می‌توان ضرایب کلمات کلیدی را بر اساس رفتار کاربر و رهگیری او اصلاح کرده و کم و یا زیاد نمود.

برای بالا بردن دقت تعیین الگوی رفتار کاربر می‌توان پارامتر زمان را نیز منظور نمود. تناوب درخواست‌ها و یا فاصله زمانی آن‌ها می‌تواند ملاکی بر اصرار و یا کم‌توجهی کاربر بر ادامه رفتار خود باشد. با استفاده از فیلد زمان درخواست کاربر در فایل گزارش، می‌توان سنجه‌های جدیدی برای اصلاح تعیین الگوی رفتار طراحی کرده و استفاده نمود.

