

A Survey of Threats and Solution to Privacy in Web Search Engines

Hadi Sabouhi¹, Maryam Fathi Ahmad Sharaei²

1 Post Doc. and Assistant Professor of Islamic Azad University, Karaj Branch, Faculty of Mechatronics, Computer Department, Alborz, Iran

saboohi@kiau.ac.ir

2 PhD Student of Islamic Azad University, Karaj Branch, Faculty of Mechatronics, Computer Department, Alborz, Iran

fathimaryam2009@gmail.com

Abstract

We use search engines to find our needed information through amount of data. Most of the time people use Internet to get information about every job. Web search engines often confuse people by presenting different results for same queries. Search engines create profiles for users to solve this problem. According to each user's profile, search engines represent results. Often, users' privacy are in danger by using these tools. In web search engines results, we need to have quality and privacy protecting both together. In this article we present problems that threaten the users' privacy. We also review how to prevent web search engines from violating the users' privacy.

key words: Privacy, personalization, users' profile in search engines.

Archive of SID

بررسی تهدیدها و راه حل های حفظ حریم خصوصی کاربران

در موتورهای جستجوی وب

هادی صبحی^۱، مریم فتحی احمدسرائی^۲

۱ فوق دکتری کامپیوتر، استادیار دانشگاه آزاد اسلامی، واحد کرج، دانشکده مکترونیک، گروه کامپیوتر، البرز، ایران
saboohi@kiauo.ac.ir

۲ دانشجوی دکتری کامپیوتر، دانشگاه آزاد اسلامی، واحد کرج، دانشکده مکترونیک، گروه کامپیوتر، البرز، ایران
fathimaryam2009@gmail.com

چکیده

موتورهای جستجوی وب برای پیدا کردن اطلاعات خاصی در میان حجم عظیمی از داده‌ها در زمان کم مورد استفاده قرار می‌گیرند. افراد در هر کاری تقریباً به اینترنت متکی هستند و موتورهای جستجوی وب گاهی اوقات با ارائه نتایج مختلف ممکن است مردم را سردرگم کنند. برای حل این مشکل از ساخت پروفایل استفاده می‌شود. نتایج با توجه به پروفایل هر کاربر نمایش داده می‌شود. این ابزارها معمولاً باعث به خطر افتادن حریم خصوصی کاربران می‌شود. محافظت از حریم خصوصی و داشتن کیفیت در نتایج جستجو را باید بتوانیم در کنار هم داشته باشیم. در این مقاله به مشکلاتی که حریم خصوصی کاربران را تهدید می‌کند، اشاره می‌شود. همچنین برای جلوگیری از نقض حریم خصوصی کاربران توسط موتورهای جستجوی وب راه حل‌های ارائه شده اخیر بررسی خواهد شد.

کلمات کلیدی

حریم خصوصی، شخصی سازی، پروفایل کاربران موتورهای جستجو

۱- مقدمه

کاربر که در پروفایل او وجود دارد می‌تواند کیفیت نتایج ارائه شده را بهبود بخشید. این روش‌ها موجب تجاوز به حریم خصوصی کاربران می‌شوند. حمله کننده نیز می‌تواند از این اطلاعات که بین کاربر و موتور جستجوی وب رد و بدل می‌شود، به دلیل فقدان سیاست‌های حریم خصوصی سوء استفاده کند [۱]. در ادامه در بخش ۲ به بررسی تهدیدهای حریم خصوصی می‌پردازیم و در بخش ۳ راه حل‌های حفظ حریم خصوصی را مورد بررسی قرار می‌دهیم. در آخر در بخش ۴، نتیجه گیری و کارهای آینده را ارائه می‌دهیم.

۲ - بررسی تهدیدهای حفظ حریم خصوصی کاربران

در بحث حریم خصوصی، باید بررسی شود که موتورهای جستجو چقدر قابل نفوذ هستند و می‌توانند حریم خصوصی کاربران را تهدید کنند [۲]. زمانی که کاربر برای اولین بار به موتور جستجو متصل می‌شود، ابتدا IP، زمان، تاریخ و

موتورهای جستجوی وب، کار دشوار استخراج اطلاعات مورد نیاز از حجم انبوهی از داده‌ها را انجام می‌دهند [۱]. بعد از ایمیل، موتورهای جستجوی وب دومین ابزاری هستند که کاربران اینترنت از آن استفاده می‌کنند [۲]. موتورهای جستجوی وب در کمترین زمان، اطلاعات را بازیابی می‌کنند. استفاده از موتورهای جستجوی وب باعث تولید گزارش^۱ از پرس‌وجوهای^۲ کاربران می‌شود [۳]. اگر به گزارش‌ها و تاریخچه جستجوی کاربران توجهی نشود، ممکن است نتایج نادرستی بازگردانده شود [۴و۵]. شخصی‌سازی^۳ جستجوی وب یکی از معمول‌ترین روش‌ها، برای جمع‌آوری و تحلیل اهداف کاربر است. شخصی سازی جستجوی وب دو نوع است [۴]: الف) روش‌های مبتنی بر گزارش کلیک‌های کاربران: که از تاریخچه صفحه‌هایی که کاربر کلیک کرده است، استفاده می‌کنند. ب) روش‌های مبتنی بر پروفایل کاربران: که با تولید مدل‌های علائق کاربران می‌توان جستجوی بهتری را تجربه کرد. با استفاده از اطلاعات رفتاری

که با کلیک کردن کاربر بر روی URLها معمولا خاتمه می‌یابند با این فرض که آن پرس‌وجوها را وابسته در نظر می‌گیرند، یعنی شبه‌شناسه‌ای [۱۰]. (ب) کنترل افشای آماری؛ پرس‌وجوی کاربران با یکدیگر گروه‌بندی و با یک پرس‌وجوی نماینده جایگزین می‌شود. نمونه‌هایی از این رویکرد در ادامه بحث شده است: به دلیل حجم بالای گزارش‌های پرس‌وجوها، در [۱۱] پیشنهاد داده‌اند تا مجموعه پرس‌وجوهای یک کاربر را به وسیله عمومیت بخشیدن پرس‌وجوهای او ناشناس سازی کنند. مشکل این رویکرد این است که یک پرس‌وجو می‌تواند در فرهنگ لغت عمومی بی‌معنی باشد، ولی در یک فرهنگ لغت تخصصی‌تر خطرناک باشد.

۳-۱-۲- پردازش زمان واقعی^۷

روش‌های بررسی شده تنها برای داده‌های ایستا^۸ بود. این روش‌ها برای جریان داده‌های^۹ آنلاین مناسب است و دو زیر روش دارد:

الف) ناشناس کردن همه پرس‌وجوها از لحظه صفر تاکنون که با مقدار زیادی داده روبرو هستیم. نمونه‌هایی از این رویکرد: استفاده از تکنیک‌های معروف SDC در این رویکرد به هنگام پارتیشن کردن پایگاه داده، باعث تحمیل هزینه‌های زیادی می‌شود [۱۲]. به صورت دقیق‌تر، مولفین ناشناس‌سازی گزارش‌های پرس و جو را اعمال کردند و برداشت آنها این بود که این تکنیک‌ها با مجموعه‌های بزرگ داده‌ها، که نیازمند آپدیت پیوسته و ناشناس‌سازی داده‌های جدید است که همواره افزوده می‌شوند، تناسب خوبی ندارند.

ب) ناشناس کردن پرس‌وجوها در یک زمان محدود که در دوره‌های سه یا شش ماهه انجام می‌شود. این روش، هزینه، دقت و حجم مجموعه داده را کاهش می‌دهد. نمونه‌هایی از این رویکرد: این رویکردها، داده‌ها را به‌عنوان یک جریان در نظر می‌گیرند و اطلاعات را به محض دریافت و با ارائه حداقل تاخیر، پردازش می‌کنند [۱۳]. جریان‌های داده می‌توانند نامتناهی باشند و حجم عظیمی از داده‌ها را تولید نمایند. بنابراین، بهره‌برداری از جریان‌های عظیم داده، یک کار چالش برانگیز است. اکثر رویکردهای موجود می‌خواهند تاخیر بین پردازش جریان‌های داده را کاهش دهند. این روش‌ها عمدتاً بر مبنای خوشه بندی جریان‌های داده ورودی هستند و همگی نیاز دارند تا جهت ایجاد خوشه‌های ناشناس منتظر داده‌های جدید باشند.

۳-۲-۲- دسته بندی دوم روش‌های حفظ حریم خصوصی

رویکردهای حفظ حریم خصوصی را به دو گروه کلی تقسیم می‌کنیم: الف) رویکردهای سمت کاربر، ب) رویکردهای سمت سرور.

۳-۲-۱- رویکردهای سمت کاربر

رویکردهای سمت کاربر که کاربر به‌طور مستقیم بر روی پروفایل خود کنترل دارد.

تنظیمات مرورگر کاربر گزارش می‌شود. طرفداران حریم خصوصی ادعا می‌کنند که از جستجوهای یک فرد در یک مدت زمان می‌توان به زندگی شخصی و علایق آن فرد رسید که این مسئله خود، نقض حریم خصوصی است. موتورهای جستجو از کوکی‌ها استفاده می‌کنند [۲]. کوکی‌ها می‌توانند حریم خصوصی کاربران را در معرض خطر قرار دهند. موتورهای جستجوی وب می‌توانند این اطلاعات را رمزنگاری کنند و با این کار می‌توانند هر اطلاعات دیگری را از کاربران ذخیره کنند. کاربران باید توافقنامه‌های موتورهای جستجوی وب را قبول کنند. در این توافقنامه‌ها، حق مدیریت و فروش اطلاعات به شرکت‌های دیگر ذکر شده است. اینکه چه اطلاعاتی از کاربران آنلاین است و برای شاخص گذاری در اختیار موتورهای جستجو قرار می‌گیرد به میزان مسئولیت‌پذیری مدیر وب سایت بستگی دارد [۲]. یک راه ساده استفاده از پراکسی است. به این ترتیب موتورهای جستجو نمی‌توانند IP کاربران را به‌دست آورند [۷].

۳- بررسی راه حل‌های ارائه شده جهت برقراری

حفظ حریم خصوصی کاربران

روش‌های حفظ حریم خصوصی را می‌توان از چند جهت دسته بندی کرد. در این مقاله ما دو روش دسته بندی را بررسی می‌کنیم. یک روش دسته بندی روش‌ها را به دو دسته پردازش دسته‌ای و پردازش زمان واقعی تقسیم بندی می‌کند. در روش دسته بندی دیگر، روش‌ها به دو دسته رویکردهای سمت کاربر و سمت سرور دسته بندی می‌شوند.

۳-۱- دسته بندی اول روش‌های حفظ حریم خصوصی

در دسته بندی اول، روش‌های حفظ حریم خصوصی عبارتند از [۳]:

۳-۱-۱- پردازش دسته‌ای^۴

الف) حذف پرس‌وجو؛ اطلاعات شناسه‌ای مانند شماره ملی یا شماره بیمه و غیره را حذف کند. نمونه‌هایی از این رویکرد: در [۸]، مجموعه پرس‌وجوی قدیمی را حذف می‌کند، در حالی که فرض می‌کند گزارش پرس‌وجو به اندازه کافی بزرگ نخواهند بود که قادر به افشای هویت شود. در [۹] پرس‌وجوهای نامعمول حذف می‌شوند و فرض می‌کند این پرس‌وجوها به احتمال بیشتر به اطلاعات شناسه‌ای اشاره می‌کنند. انتخاب حدآستانه حذف مناسب می‌تواند کاملاً چالش برانگیز باشد. به‌عنوان یک نتیجه، یک روش ممکن است موجب حذف مقدار زیادی از پرس‌وجوهای غیرشناساگر ولی مفید شود. سایر روش‌های مبتنی بر حذف به حذف داده‌های شناسه‌ای مرتبط با پرس‌وجوها و اطلاعات خصوصی مانند شماره‌های تامین اجتماعی، کارت‌های اعتباری، آدرس‌ها و غیره تمرکز دارند. رویکردهای پیچیده‌تر در این زمینه بر حذف آن پرس‌وجوهایی تمرکز دارند

روش اول از رویکردهای سمت کاربر

در این سیستم هر کاربر یک حساب کاربری ایجاد می کند و اطلاعاتی مانند آدرس، شغل، و غیره را ثبت می کند. سپس شروع به جستجو می کند و پروفایل کاربر در سرور ذخیره می شود. اگر کاربر بر روی خصوصی^{۱۱} کلیک کند، لینکها تنها برای خود کاربر قابل دسترس هستند. اگر کاربر بر روی گزینه عمومی^{۱۱} کلیک کند، پروفایل او برای همه قابل دسترس است. جریان کاری این روش در شکل ۱ نمایش داده شده است. در این روش از الگوریتم های Spy Naïve و Bayes استفاده می شود [۴].

روش دوم از رویکردهای سمت کاربر

تاریخچه مرورگر و ایمیل های هر فرد می تواند منبع اطلاعاتی خوبی برای پروفایل کاربر باشد. فرض ما این است که اگر عبارتی در این اسناد^{۱۲} تکرار شود جزء علائق کاربر است. در این رویکرد پروفایل کاربر را به صورت سلسله مراتبی بر مبنای تکرار عبارت ها ایجاد می کنیم. در این سلسله مراتب عبارات با تکرار بیشتر در سطوح بالاتر قرار می گیرند. $D(t)$ شامل همه اسنادی است که عبارت t در آنها وجود دارد. برای سازماندهی سلسله مراتبی عبارات از دو قانون هیوربستیک زیر استفاده می کنیم: *الف) عبارت مشابه: دو عبارتی که مجموعه اسناد یکسانی را با همپوشانی^{۱۳} زیادی پوشش می دهند، معمولاً علائق یکسانی را نشان می دهند. ب) عبارات والد و فرزندی^{۱۴}: معمولاً عبارات خاص با عبارات کلی ظاهر می شوند ولی عکس این قضیه درست نیست. بنابراین t_B فرزند t_A است اگر احتمال شرطی $p(t_A | t_B) > \delta$ اتفاق بیفتد.*

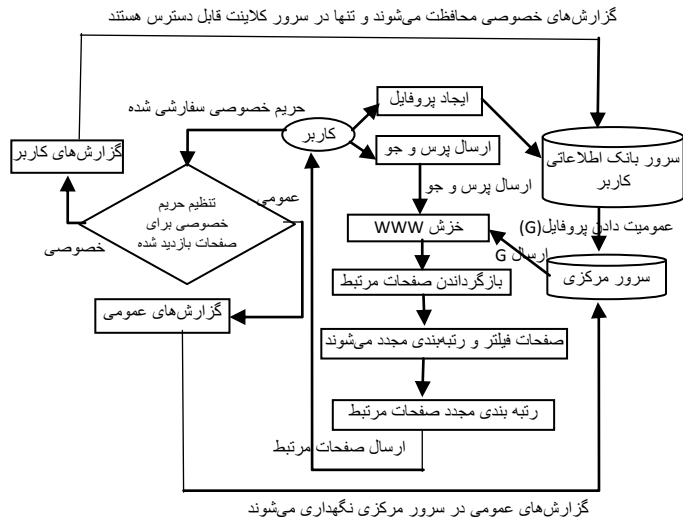
الگوریتم بررسی شده را به صورت شبه کد در الگوریتم ۱ مشاهده می کنیم. پروفایل کاربر کاملاً تحت کنترل کاربر است و می تواند عبارتهایی را که تمایل دارد در ساختار سلسله مراتبی پنهان کند. کاربر دقیقاً می داند که کدام بخش از پروفایل او محافظت شده است. برای هر پروفایل یک لیست از زوج مرتب های (t, w_t) داریم. به طوری که t عبارت و w_t وزن عبارت است.

۳-۲-۲- رویکردهای سمت سرور

در این رویکردها کاربر به طور مستقیم بر محافظت از حریم خصوصی خود دخالت ندارد بلکه به واسطه یک سرور میانی این کار انجام می شود.

روش اول از رویکردهای سمت سرور

پروتکل UUP از دریافت پروفایل معتبر کاربران توسط موتورهای جستجو جلوگیری می کند. منظور از U کاربران، C نود مرکزی که الگوریتم UUP را اجرا می کند و W همان موتورهای جستجو هستند.



شکل ۱ - نمایش جریان کار سیستمی از رویکردهای سمت کاربر

الگوریتم: $Split(n, S(t), minsup, \delta)$

ورودی: به یک گره n برچسب عبارت t داده می شود، اسناد پشتیبانی $S(t)$ ، حدآستانه $minsup$ و δ

- لیستی از کلمات پرتکرار $\{t_i\}$ تولید می شود به طوری که $D(t_i) \geq minsup$ و بر اساس تعداد تکرار به صورت نزولی مرتب شده است.
- برای هر عبارت t_i :
- اگر $k < i$ و $sim(t_i, t_k) > \delta$
- برچسب گره t_i/t_k تنظیم می شود و $S(t_i/t_k) = S(t_k) \cup D(t_i)$
- در غیراینصورت اگر $k < i$ و $P(t_k | t_i) > \delta$
- برچسب گره همان t_k می ماند و $S(t_k) = S(t_k) \cup D(t_i)$
- در غیراینصورت
- یک گره جدید با برچسب t_i ایجاد می شود و $S(t_i) = D(t_i)$
- $Sup(t_i)$ برای هر گره با برچسب t_i محاسبه می شود و کلمات به صورت نزولی مرتب می شوند

الگوریتم: $BuildUP(n, D, minsup, \delta)$

ورودی: گره n ، اسناد پشتیبانی D ، حدآستانه $minsup$ و δ

خروجی: پروفایل یک کاربر U

- $Split(n, D, minsup, \delta)$
- برای هر فرزند c_i با برچسب t_i از گره n :
- $BuildUP(c_i, S(t_i), minsup, \delta)$

الگوریتم ۱- الگوریتم $Split$ کردن مجموعه ای از داکومننت ها

کنید $\{C_1^D, \dots, C_n^D\}$. مجموعه پرس و جوهای رمز شده پس از ارسال در گروه توسط U_i است. برای هر کاربر U_i از 1 تا n مراحل زیر اجرا می شود:

(a) کاربر U_i پرس و جوهای رمز شده $\{C_1^{i-1}, \dots, C_n^{i-1}\}$ از کاربر U_{i-1} را دریافت می کند و مجدداً با الگوریتم ElGamal رمز می کند و $\{C_1^i, \dots, C_n^i\}$ را به دست می آورد. (b) سپس ترتیب پرس و جوهای رمز شده را عوض می کند و مجموعه $\{C_{S(1)}^i, \dots, C_{S(2)}^i\}$ را به دست می آورد. (c) U_i مجموعه پرس و جوهای رمز شده تا آن مرحله $\{C_1^i, \dots, C_n^i\}$ را به U_{i+1} ارسال می کند. در نهایت کاربر U_n مجموعه $\{C_1^n, \dots, C_n^n\}$ را به همه اعضای گروه به صورت همه پختی ارسال می کند. در این مرحله $\{C_1, \dots, C_n\} = \{C_1^n, \dots, C_n^n\}$ را در نظر می گیریم. U_i پرس و جوی $(C_{1i})^{a_i}$ را از U_j دریافت می کند به طوری که $j=1, \dots, n$ و $i \neq j$ است. سپس U_i پرس و جوی $(C_{1j})^{a_j}$ را با کلید خصوصی خود محاسبه می کند. کاربر U_i پرس و جوی رمز شده m_i را به صورت رابطه (۲) به دست می آورد:

$$m_i = \frac{C_{2i}}{(C_{1i})^{a_i} \prod_{j \neq i} (C_{1j})^{a_j}} \quad \text{رابطه (۲)}$$

د) ارسال پرس و جو و بازیابی اطلاعات: هر عضو گروه (U_i) ، m_i را

به موتور جستجو ارسال می کند. زمانی که پاسخ a_i را دریافت کرد، به کل گروه به صورت همه پختی ارسال می کند. هر کاربر با توجه به پرس و جوی ارسالی پاسخ a_i خود را برمی دارد.

روش دوم از رویکردهای سمت سرور

موتور جستجو، پرس و جو را به عنوان ورودی می گیرد و گزارش پرس و جوی ناشناس شده و پروفایل کاربر را به عنوان خروجی تولید می کند. زیر سیستم ها در این بخش توصیف می شوند [۳]:

- الف) ایجاد پروفایل و ناشناس سازی: این فاز اصلی ترین فاز است.
- از توصیه کننده ها^{۱۶} به منظور تصحیح غلط املائی استفاده می شود.
- آنتروپی ها^{۱۷}: تعداد نتایج این کلمه در مقایسه با بقیه عبارات های موجود در پرس و جو محاسبه می شود تا کلمه اصلی در کل عبارت پرس و جو پیدا شود.
- طبقه بندی کننده^{۱۸}: طبقه یا دسته متناظر با پرس و جو پیدا می شود.
- ناشناس کننده^{۱۹}: گزارش های دسته بندی شده به عنوان ورودی این مرحله است. برای هر دسته دو مجموعه پرس و جو و کاربر داریم. هنگامی که یک مجموعه از هر دسته به ماکزیم مقدار k می رسد به طور تصادفی یک کاربر و یک پرس و جو انتخاب و به یکدیگر متصل می شوند. با این کار کلیت حفظ می شود ولی ویژگی مخصوص او حذف می شود. یک مورد خاص زمانی است که در یک دسته در مجموعه کاربران، همه کاربران مشابه باشند که با توجه به تعداد پرس و جوی دریافتی در ثانیه توسط موتورهای جستجو احتمال خیلی پایینی دارد. در

ایده اصلی این روش به این صورت است: هر کاربر که می خواهد پرس و جویی را ارسال کند، پرس و جوی خود را ارسال نمی کند بلکه پرس و جوی کاربر دیگر را ارسال می کند و پرس و جوی او نیز توسط فرد دیگری ارسال می شود. کاربران از پرس و جوهای یکدیگر اطلاعی ندارند زیرا پرس و جوها رمزنگاری می شوند. بنابراین موتورهای جستجو نمی توانند برای کاربران، پروفایل های واقعی ایجاد کنند. نود مرکزی، n کاربر را گروه بندی می کند. برای هر کاربر U_i به طوری که $i \in \{1, \dots, n\}$ یکی از پرس و جوهای $n-1$ کاربر دیگر را ارسال می کند. U_i از مبدا پرس و جو دریافتی، اطلاعی ندارد. برای رسیدن به این هدف ابتدا همه پرس و جوها بین کاربران به طور تصادفی توزیع می شود. این کار با استفاده از عملیات رمزنگاری، بازسازی و جایگشت انجام می شود. سپس پرس و جوی دریافتی توسط کاربران ارسال می شود. موتور جستجو، نتیجه پرس و جوها را به کاربران به صورت همه پختی^{۱۵} ارسال می کند. هر کاربر نتایج مربوط به خود را برمی دارد و دیگر نتایج را دور می ریزد. در این پروتکل فرض شده است که همه کاربران از پروتکل پیروی می کنند و هیچ تبانی بین دو کاربر وجود ندارد. این پیش فرض ها به دلیل رسیدن به این هدف است که موتور جستجو نتواند برای کاربران پروفایل ایجاد کند [۱]. این پروتکل شامل چهار زیر پروتکل است: **الف) ایجاد گروه**: زمانی که کاربر U_i می خواهد به موتور جستجوی وب یک پرس و جو ارسال کند یک درخواستی را به گروه مرکزی C ارسال می کند. گروه C درخواست های n کاربر را دریافت می کند و یک گروه n عضوی از کاربران $\{U_1, \dots, U_n\}$ ایجاد می کند. به کاربران اعلام می کند که عضو یک گروه هستند. کاربران، یک کانال ارتباطی بین خودشان دارند. هر کاربر U_i می تواند به بقیه گروه پیام ارسال کند. از این به بعد دیگر نیازی به حضور نود مرکزی C نیست. **ب) تولید کلید برای گروه**: کاربران $\{U_1, \dots, U_n\}$ بر سر یک کلید p که یک عدد اول است، توافق می کنند. به طوری که $p=2q+1$. q نیز در اینجا یک عدد اول است. کاربران $\{U_1, \dots, U_n\}$ با استفاده از الگوریتم ElGamal، کلید عمومی y را تولید می کنند. هر کاربر U_i کلید خصوصی x_i را برای خود نگه می دارد. **ج) ناشناس کردن پرس و جوها**: برای هر کاربر گروه $\{U_1, \dots, U_n\}$ که پرس و جوی m_i را ارسال و پرس و جوی m_i را از یک عضو دیگر از گروه دریافت می کند، مراحل زیر را انجام می دهیم: برای هر کاربر U_i به طوری که $i \in \{1, \dots, n\}$ است، این کارها را انجام می دهیم: هر کاربر U_i ، یک مقدار تصادفی r_i را تولید می کند و پرس و جوی m_i خود را با استفاده از کلید عمومی گروه، y رمزنگاری می کند (با استفاده از الگوریتم ElGamal). نتیجه به صورت رابطه (۱) می شود:

$$E_y(m_i, r_i) = (C_{1i}, C_{2i}) = C_i^D \quad \text{رابطه (۱)}$$

کاربر U_i پرس و جوی رمز شده C_i^D را به کاربر دیگر گروه U_j ارسال می کند به طوری که $j \neq i$. یک ترکیب خاص برای کاربران مثلاً از 1 تا n فرض

الگوریتم ۲

ورودی: گزارش‌های دسته‌بندی شده، k, δ
خروجی: گزارش‌های ناشناس شده

۱. برای هر گزارش از مجموعه گزارش‌های دسته‌بندی شده این کارها انجام شود:
۲. $user, query, category \leftarrow \log$
۳. $users [category] \leftarrow user;$
۴. $query [category] \leftarrow query;$
۵. اگر اندازه $users [category]=k$ بود، سپس
۶. اگر $u \in users [category]$ و $u \neq user$ سپس
۷. یک کاربر را به صورت تصادفی از مجموعه کاربران انتخاب کن
۸. یک پرس‌وجو را به صورت تصادفی از مجموعه پرس‌وجوها انتخاب کن
۹. ارسال u و q
۱۰. در غیراینصورت $\delta = k * k$
۱۱. پایان

الگوریتم ۲- الگوریتم ناشناس‌سازی

الگوریتم ۳

ورودی: گزارش‌های دسته‌بندی شده، k, δ
خروجی: گزارش پرس‌وجوها

۱. برای هر گزارش از مجموعه گزارش‌های دسته‌بندی شده این کارها انجام شود:
۲. $user, query, category \leftarrow \log$
۳. $users [category] \leftarrow user;$
۴. $query [category] \leftarrow query;$
۵. اگر اندازه $users [category]=k$ بود، سپس
۶. اگر $u \in users [category]$ و $u \neq user$ سپس
۷. یک کاربر را از مجموعه کاربران انتخاب کن
۸. یک پرس‌وجو را از مجموعه پرس‌وجوها انتخاب کن
۹. ارسال u و q
۱۰. در غیراینصورت $\delta = k * k$
۱۱. پایان

الگوریتم ۳- کاهش ناشناس‌سازی

این شرایط یک ضریب δ در نظر می‌گیریم که در k ضرب می‌کنیم تا سائز مجموعه افزایش یابد. الگوریتم ۲ این فرآیند را توصیف می‌کند. (ب) کاهش ناشناس‌سازی^{۲۰}: برعکس فاز ناشناس‌سازی عمل می‌کند. با این کار تست می‌کند که آیا حمله کننده می‌تواند موفق شود یا خیر. در الگوریتم ۳ حمله کننده سعی می‌کند تا هر پرس‌وجو را به کاربر خودش ربط دهد.

۳-۳- مقایسه روش‌های بررسی شده

در جدول ۱ به مقایسه روش‌های بررسی شده می‌پردازیم.

۴- نتیجه گیری و کارهای آینده

در این مقاله به بررسی تهدیدها و راه حل‌های ارائه شده جهت برقراری حفظ حریم خصوصی کاربران پرداختیم. در اکثر الگوریتم‌های ارائه شده، کاربران و سرورهای واسط غیرمهاجم در نظر گرفته شده‌اند. در صورتی که در محیط واقعی ممکن است مهاجمی خود را به عنوان کاربر عادی معرفی کند. همچنین ممکن است مهاجم به سرورهای واسط نیز حمله کند. در صورتی که در روش‌های بررسی شده به نحوه برقراری امنیت فیزیکی و نرم افزاری این سرورها اشاره‌ای نشده است. الگوی جستجوی کاربران ممکن است تغییر کند. در پژوهش‌های آینده باید این شرایط در نظر گرفته شود. تحلیل رفتار کاربران در کنار مشخصات آنها می‌تواند نتایج جستجوی بهتری را ارائه کند. تمام URLهایی که یک کاربر برای یافتن نتایج دنبال می‌کند، دلیل بر مرتبط بودن با موضوع پرس‌وجو را ندارد. باید زمان توقف در یک صفحه، میزان کلمات آن صفحه و مواردی از این قبیل نیز در نظر گرفته شوند.

منابع

- [1] J. Roca, A. Viejo and J. Joancomarti, "Preserving User's Privacy in Web Search Engines", Elsevier, Computer Communications, (2009), 1541-1551.
- [2] H. Aljifri and D. Navarro, "Search Engines and Privacy", Elsevier, Computers & Security, (2004), 379-388.
- [3] D. Estrems, J. Roca and A. Viejo, "Working at the Web Search Engine Side to Generate Privacy-Preserving User Profiles", Elsevier, Expert Systems With Applications, (2016), 523-535.
- [4] S. Malthankar and S. Kolte, "Client Side Privacy Protection Using Personalized Web Search", 7th International Conference on Communication, Computing and Virtualization, Mumbai, India, (2016), 1029 - 1035.
- [5] J. M. Saji, K. BhongleJ. and et al., "Advancement in Personalized Web Search Engine With Customized Privacy Protection", Springer, Nature Singapore Pte Ltd, (2018), 405-413.

جدول ۱ - مقایسه راه حل های ارائه شده جهت حفظ حریم خصوصی کاربران موتورهای جستجوی وب

نام نویسنده	نام الگوریتم	سمت سرور یا کلاینت	پیش فرض ها	مزایا	معایب	عملکرد
J. [۱] Roca, A. Viejo	ندارد	سرور	سرور میانی و کلاینتها تبانی نمی کنند	عملکرد بهتر نسبت به روش های مشابه که از رمزنگاری استفاده می کنند.	تاخیر در پاسخ به پرس و جوها.	هر کاربر پرسجوی خود را رمز می کند و رندم توسط کاربر دیگری ارسال می شود. کاربر پاسخ پرسجوهای دریافتی را رمزگشایی می کند تا پاسخ پرس و جو خود را بیابد.
D. [۳] Estrems, J. Roca	UUP	سرور	داده های اصلی در یک پایگاه داده امن نگهداری می شوند.	احتمال موفقیت حمله کننده ۱/۸٪ است، حفظ کلیت علاقه گروهی از کاربران به موضوعی خاص	مصرف زیاد حافظه	جهت ناشناس سازی پرسجوها و کاربران در هر دسته به صورت تصادفی جایگشت داده می شوند.
S. [۴] Malthankar S. Kolte و	UPS	کلاینت	سرورها غیرقابل اعتماد هستند و هر کلاینت تنها به خود اعتماد دارد.	هزینه کم، ایجاد پروفایل آنلاین، داشتن حق انتخاب کاربر برای هر پرس و جو	کاربر باید ثبت نام کند و اطلاعات هویتی ارائه دهد	کاربر برای هر پرس و جو می تواند خصوصی و عمومی بودن را برای گزارش آن پرس و جو در پروفایل خود تعیین کند.
Y. Xu [۱۴] B. Zhang و	ندارد	کلاینت	اگر عبارتی در داکویمنت کاربر تکرار شود، جز علائق کاربر است.	پنهان کردن هر بخش از پروفایل توسط کاربر	احتمال دسته بندی عبارات به عنوان والد و فرزندی به جای دسته بندی در یک گروه کلی تر	ایجاد پروفایل سلسله مراتبی برای کاربران و دسته بندی پرس و جوها در گروه های کلی تر، کاربر امکان مخفی کردن هر بخش از درخت سلسله مراتبی را دارد.

[6] J. D. Velásquez, "Web Mining and Privacy Concerns: Some Important Legal Issues to be Consider Before Applying any Data and Information Extraction Technique in Web-Based Environments", Elsevier, Expert Systems with Applications, (2013), 5228–5239.

[7] J. Wang, and H. Raghav, "An Exploration of Risk Information Search via a Search Engine: Queries and Clicks in Healthcare and Information Security", Elsevier, Decision Support Systems, (2012), 395–405.

[8] A. Cooper, "A Survey of Query Log Privacy-Enhancing Techniques from a Policy Perspective", ACM Transactions on the Web, New York, (2008), 19–27.

[9] S. M. Beitzel, E. Jensen and A. Chowdhury, "Hourly Analysis of a Very Large Topically Categorized Web Query Log", 27th annual international ACM SIGIR conference on research and development in information retrieval, Sheffield, United Kingdom, (2004).

[10] A. Korolova, K. Kenthapadi, N. Mishra, and A. Ntoulas, "Releasing Search Queries and Clicks Privately", In Proceedings of the 18th international world wide web conference, Madrid, Spain, (2009).

[11] Y. He and J. F. Naughton, "Anonymization of Set-Valued Data via Topdown, Local Generalization", In Proceedings of the proceedings of the VLDB endowment, (2009).

[12] J. Soria and J. Domingo, "Big Data Privacy: Challenges to Privacy Principles and Models", Data Science and Engineering, (2015), 1–8.

[13] L. Rutkowski, M. Jaworski, and L. Pietruczuk, "Decision Trees for Mining Data Streams Based on the Gaussian Approximation", IEEE Transactions on Knowledge and Data Engineering, (2014), 108 - 119.

[14] Y. Xu and B. Zhang and K, "Privacy-Enhancing Personalized Web Search", International World Wide Web Conference Committee (IW3C2), Lyon (France), (2007), 8–12.

زیرنویس

- ¹ Log
- ² Query
- ³ Personalize
- ⁴ Batch Processing
- ⁵ Query Removal

-
- 6 Statistical Disclosure Control (SDC)
 - 7 Real Time Processing
 - 8 Static
 - 9 Data Stream
 - 10 Private
 - 11 Public
 - 12 Documents
 - 13 Overlap
 - 14 Chide & Parent
 - 15 Broadcast
 - 16 Recommender
 - 17 Entropies
 - 18 Classifier
 - 19 Anonymizer
 - 20 De-anonymizer

Archive of SID