



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

دانشکده مهندسی کامپیوتر

پیشنهاد رساله دکتری مهندسی کامپیوتر (گرایش نرم افزار)

حریم خصوصی تفاضلی در جستجوی وب شخصی سازی شده

Differential Privacy in Personalized Web Search

دانشجو: احسان عدالت

استاد راهنما: مهران سلیمان فلاح

آذر ماه ۱۴۰۰

چکیده

موتورهای جستجو نتایج پرس‌وجوهای کاربران را شخصی‌سازی نموده و نتایجی دقیق‌تر، متناسب با پسندهای کاربران، فراهم می‌آورند. برای این منظور، موتورهای جستجو همواره به دنبال استخراج اطلاعات زمینه و یادگیری الگوی رفتار جستجوی کاربران هستند. از سوی دیگر، کاربران نسبت به پسندهای خود حساسند و آن‌ها را خصوصی می‌دانند. بنابراین، حفظ حریم خصوصی کاربران در تعامل با موتورهای جستجو موضوعی مهم است. در این پیشنهاد رساله، حفظ حریم خصوصی کاربران در تعامل با موتور جستجو را بر اساس مدل حریم خصوصی تفاضلی مطرح می‌کنیم و به بیان جزئیات این مسئله و ارائه پیشنهادی برای حل آن می‌پردازیم.

موتور جستجو اطلاعات زمینه کاربران، مانند موقعیت جغرافیایی، و رفتار جستجوی آنان را، شامل پرس‌وجوهای انجام شده و سندهایی که در پی جستجوهای خود مراجعه نموده‌اند، در یک پایگاه‌داده ذخیره می‌نماید. آنگاه، با استفاده از داده‌های موجود در این پایگاه‌داده، پروفایل کاربران را ایجاد می‌کند. موضوع مهم آن است که بین داده‌های ثبت شده در این پایگاه‌داده همبستگی وجود دارد، به گونه‌ای که از سایر رکوردهای پایگاه‌داده می‌توان درباره یک رکورد حذف شده یا تغییر یافته اطلاعاتی به دست آورد. علاوه بر آن، این پایگاه‌داده رشدیابنده است و در طول زمان رکوردهای تازه‌ای به آن اضافه می‌شود. بنابراین، در عملی کردن خط‌مشی حریم خصوصی تفاضلی باید به چالش‌های همبستگی داده‌ها و رشدیابنده بودن پایگاه‌داده توجه کنیم.

موضوع مهم دیگر آن است که موتور جستجو یک موجودیت جعبه‌سیاه است و روش آن در تشکیل پروفایل کاربران و نیز بازیابی سندهای جستجو شده مشخص نیست. این موضوع چالشی دیگر در عملی کردن خط‌مشی حریم خصوصی تفاضلی کاربران است. با توجه به این چالش، حل مسئله را در سه گام انجام می‌دهیم. در گام نخست، خط‌مشی حریم خصوصی را با فرض مشخص بودن روش یادگیری موتور جستجو عملی می‌کنیم. در گام دوم، خط‌مشی را در حالتی عملی می‌کنیم که روش یادگیری موتور جستجو به شکل یک توزیع احتمال بر روی مجموعه‌ای از روش‌ها داده شده است. در گام سوم، از فنون یادگیری ماشین برای تخمین روش موتور جستجو در تشکیل پروفایل کاربران و نیز بازیابی سندهای جستجو شده بهره می‌جوییم.

موتور جستجو پایگاه‌داده فعالیت‌های کاربران را به طور کامل در دست دارد و پروفایل کاربران را از آن ایجاد می‌کند. بنابراین، برای حفظ حریم خصوصی کاربران باید پایگاه‌داده فعالیت‌های آنان را به شکل مغشوش شده سنتز کنیم. پایگاه‌داده مغشوش شده با اضافه کردن تعدادی رکورد جدید یا تغییر رکوردهای موجود ایجاد می‌شود. در جستجوی وب شخصی‌سازی شده، با ارسال پرس‌وجوهای پوششی، علاوه بر پرس‌وجوی اصلی، و نیز گسترش پرس‌وجوی اصلی پایگاه‌داده مغشوش شده ایجاد شده و موتور جستجو به طور دقیق از پسندهای کاربران آگاه نمی‌شود. متناسب با خط‌مشی حریم خصوصی کاربران، روش موتور جستجو در یادگیری پروفایل کاربران، و نیز همبستگی موجود میان رکوردهای پایگاه‌داده تعداد و نوع پرس‌وجوهای پوششی و شیوه گسترش پرس‌وجوی اصلی را تعیین می‌کنیم. تغییرات ایجاد شده در پرس‌وجوها به گونه‌ای است که نتایج شخصی‌سازی شده، تا آنجا که ممکن است، برای کاربران مفید باشند.

کلمه‌های کلیدی: جستجوی وب شخصی‌سازی شده، حریم خصوصی تفاضلی، سنتز پایگاه‌داده مغشوش شده، بازیابی اطلاعات شخصی‌سازی شده

فهرست مطالب

۱ فصل اول: مقدمه	۱
۱-۱ جستجوی وب	۲
۱-۱-۱ نمایه سازی	۲
۱-۱-۲ پرس و جو و پاسخ	۴
۱-۱-۳ شخصی سازی جستجوی وب	۵
۱-۱-۴ حریم خصوصی در جستجوی وب شخصی سازی شده	۶
۲-۱ اهداف رساله	۱۰
۳-۱ ساختار مطالب پیشنهاد رساله	۱۳
۲ فصل دوم: پیش زمینه	۱۵
۱-۲ حریم خصوصی	۱۶
۲-۲ حفظ حریم خصوصی مبتنی بر افراز کردن	۱۷
۳-۲ حریم خصوصی مبتنی بر تصادفی کردن	۲۲
۴-۲ روش های ایجاد پرو فایل کاربر	۳۲
۵-۲ آنتولوژی محاسباتی	۳۳
۳ فصل سوم: کارهای پژوهشی پیشین	۳۵
۱-۳ کارهای پیشین حفظ حریم خصوصی در جستجوی وب شخصی سازی شده	۳۶
۳-۲ نقد و بررسی پژوهش های پیشین	۴۳
۳-۳ پژوهش های مرتبط با چالش های مطرح شده مسئله	۴۹
۴ فصل چهارم: پیشنهاد رساله	۵۳
۱-۴ بیان کلی مسئله و گام های پیشنهادی برای حل آن	۵۴
۲-۴ شرح مسئله	۵۷

۴-۲-۱	مدل سامانه شخصی سازی نتایج جستجوی کاربران	۵۷
۴-۲-۲	مدل مهاجم	۵۹
۴-۲-۳	حریم خصوصی تفاضلی و سنتز تاریخچه جستجوها	۶۰
۴-۲-۴	چالش های حل مسئله	۶۶
۵-۲-۴	گام نخست در حل مسئله با فرض جعبه سفید بودن موتور جستجو	۶۷
۴-۲-۶	گام دوم در حل مسئله با فرض جعبه خاکستری بودن موتور جستجو	۶۹
۴-۲-۷	گام سوم در حل مسئله با فرض جعبه سیاه بودن موتور جستجو	۷۱
۴-۳	ارزیابی	۷۱
۴-۴	هدف های رساله و زمان بندی	۷۳
۵	ضمیمه الف: نحوه انجام پژوهش جهت انتخاب موضوع رساله	۷۵
۶	ضمیمه ب: پیاده سازی الگوریتم ایجاد پروفایل کاربر	۸۱
۷	ضمیمه پ: واژه نامه انگلیسی به فارسی	۸۵
۸	ضمیمه ت: واژه نامه فارسی به انگلیسی	۸۹
۹	مراجع	۹۳

فهرست اشکال

- شکل ۱-۱: فرایند نمایه‌سازی در موتور جستجو. ۳
- شکل ۱-۲: فرایند پرس‌وجو در موتور جستجو. ۴
- شکل ۲-۱: حفظ حریم خصوصی در نگهداری، تحلیل و انتشار اطلاعات. ۱۶
- شکل ۲-۲: حریم خصوصی تفاضلی در حضور یا نبود شخص X ۲۴
- شکل ۲-۳: سازوکار لاپلاس. ۲۵
- شکل ۲-۴: نمونه پروفایل شبکه نحوی. ۳۲
- شکل ۴-۱: شخصی‌سازی نتایج پرس‌وجوی کاربر با روش مرتب‌سازی. ۵۴
- شکل ۴-۲: شخصی‌سازی نتایج پرس‌وجوی کاربر با روش گسترش پرس‌وجو. ۵۵
- شکل ۴-۳: مدل مفهومی حفظ حریم خصوصی در جستجوی وب شخصی‌سازی شده بر اساس سنتز پایگاه‌داده مغشوش شده مبتنی بر حریم خصوصی تفاضلی. ۶۳
- شکل ۴-۴: بخشی از آنتولوژی محاسباتی موضوع‌ها. ۶۴
- شکل ۴-۵: وضعیت خوشه‌ها (دایره‌های بنفش) و موضوع‌ها (دایره‌های آبی و قرمز) در سه حالت $k=1$ (a)، $k=2$ (b) و $k=3$ (c). ۶۸
- شکل ۴-۶: وضعیت موضوع‌ها و خوشه‌ها در حالت (c) بعد از مغشوش کردن پایگاه‌داده. ۶۹
- شکل ۴-۷: نمونه‌ای از رکوردهای ذخیره شده در پایگاه‌داده AOL. ۷۲
- شکل ۵-۱: متدولوژی پژوهشی استفاده شده. ۷۶
- شکل ۵-۲: روند انتخاب موضوع پژوهشی. ۷۷
- شکل ۶-۱: پیاده‌سازی الگوریتم ایجاد پروفایل کاربر. ۸۴

فهرست جداول

- جدول ۴-۱: زمان‌بندی اجرای فعالیت‌های رساله پیشنهادی ۷۴
- جدول ۵-۱: خلاصه‌ای از متودولوژی استفاده شده برای جستجو و پژوهش ۷۷
- جدول ۵-۲: پایگاه داده‌های علمی مورد جستجو ۷۹
- جدول ۵-۳: نام کنفرانس‌های مرتبط با حریم خصوصی تفاضلی و جستجوی وب شخصی‌سازی شده ۷۹
- جدول ۵-۴: نام مجله‌های مرتبط با حریم خصوصی تفاضلی و جستجوی وب شخصی‌سازی شده ۸۰

فصل اول

مقدمه

نیاز افراد به موتورهای جستجوی وب برای یافتن اطلاعات در اینترنت غیرقابل انکار است. موتورهای جستجو نتایج جستجو را شخصی‌سازی می‌کنند تا افراد بهتر و سریع‌تر مستندهای حاوی اطلاعات مورد نظر خود را پیدا کنند. برای این منظور، اطلاعات زمینه^۱ کاربران و رفتار جستجوی^۲ آنان را جمع‌آوری می‌کنند و بر اساس آن‌ها پاسخ بهتری به پرس‌وجوی کاربران می‌دهند. با وجود این، ممکن است کاربران به جمع‌آوری اطلاعات مربوط به خود در بعضی موضوع‌ها علاقمند نباشند. به عبارت دقیق‌تر، کاربران می‌خواهند، تا آنجا که ممکن است، موتور جستجو اطلاعات کمتری درباره پسندهای^۳ آنان به دست آورد.

در ادامه، ابتدا شیوه عملکرد موتورهای جستجو را شرح می‌دهیم. آنگاه، مسئله حفظ حریم خصوصی^۴ در جستجوی وب شخصی‌سازی^۵ شده، که موضوع این پیشنهاد رساله است، مطرح می‌گردد. ساختار این پیشنهاد رساله نیز در قسمت پایانی می‌آید.

۱-۱ جستجوی وب

موتورهای جستجوی وب با نمایه‌سازی^۶ صفحه‌های وب امکان پردازش و پاسخ به پرس‌وجوی کاربران را فراهم می‌کنند.

۱-۱-۱ نمایه‌سازی

در شکل ۱-۱ فرایند نمایه‌سازی آمده است. در زیرفرایند استخراج متن، متن‌ها از سندها (صفحه‌های وب) جمع‌آوری می‌شوند. این کار توسط خزشگرهای^۷ وب، نرم‌افزارهای خودکاری که با بررسی صفحه‌ها و دنبال کردن پیوندها متن‌ها را استخراج می‌کنند، انجام می‌شود. همچنین، خزشگرهای وب از سرنخ‌های خبری^۸، که در قالب XML هستند، استفاده کرده و صفحه‌های جدید، فیلم‌ها، عکس‌ها و مانند آن را در اختیار موتور جستجو قرار می‌دهند. موتور جستجو تمام سندهای پیدا شده را در یک مخزن نگهداری می‌کند. در دنیای وب، قالب‌های مختلفی، مانند HTML، XML، و PDF، برای سندها وجود دارند. موتور جستجو همه این سندها را به سندی در قالب XML تبدیل می‌کند تا پردازش آن‌ها ساده‌تر شود.

^۱ context information

^۲ search behavior

^۳ preference

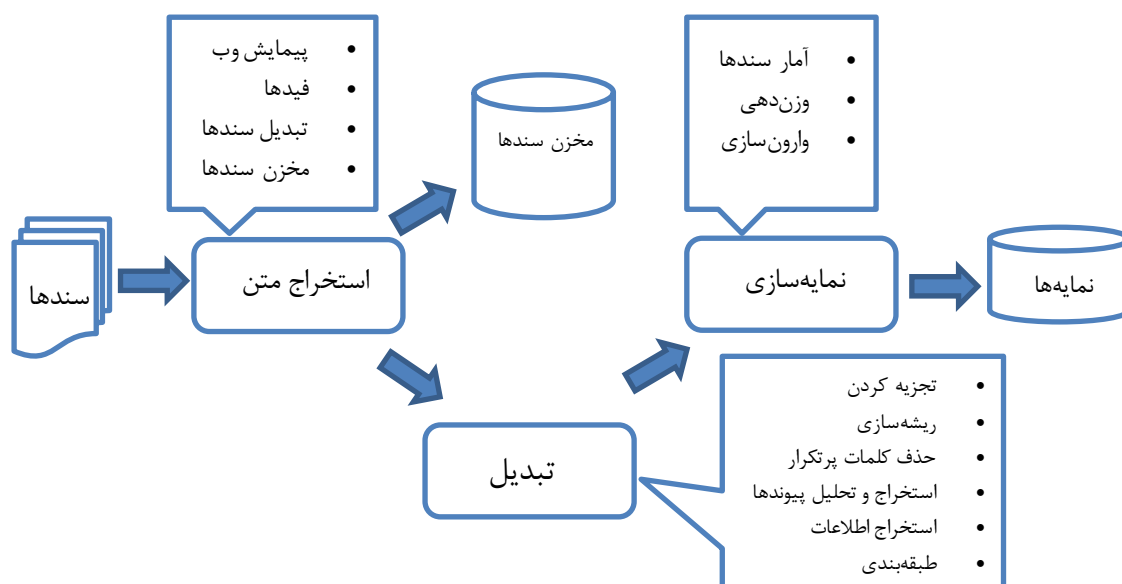
^۴ privacy

^۵ personalized

^۶ indexing

^۷ crawler

^۸ feed



شکل ۱-۱: فرایند نمایه سازی در موتور جستجو [۱].

زیرفرایند تبدیل متن ابتدا متن استخراج شده از یک سند را تجزیه^۱ کرده و عنوان سند، پیوندها، و کلمه‌ها را جدا می‌کند. آنگاه، کلمه‌های موسوم به ایست^۲، کلمه‌هایی که در اغلب متن‌ها به فراوانی یافت می‌شوند، کنار گذاشته می‌شوند. همچنین، هر کلمه با ریشه^۳ خود جایگزین می‌شود. سپس، اطلاعات مکانی، تاریخ‌ها، افراد معروف، و مانند آن استخراج می‌شوند. همچنین، بر اساس اطلاعات استخراج شده و فراداده‌های موجود در سند، به آن سند برچسب زده می‌شود. به این ترتیب، سند طبقه‌بندی می‌شود.

در زیر فرایند نمایه سازی، اطلاعات آماری در مورد سندها مانند بسامد کلمه‌ها، مکان رخداد کلمه‌ها در سند، و تعداد کل کلمه‌های سند استخراج می‌شوند. سپس الگوریتم‌های به کار گرفته شده در بازیابی اطلاعات^۴، مانند الگوریتم‌های مبتنی بر آماره TF-IDF، با استفاده از اطلاعات آماری به دست آمده و برچسب سند، به کلمه‌های موجود در سند وزن می‌دهد. در آخر، وارون سازی^۵ صورت می‌گیرد که با آن به هر کلمه تعدادی سند مرتبط خواهد شد و اطلاعات به شکل کلمه-سند ذخیره می‌شوند. به ساختمان داده حاصل از وارون سازی نمایه می‌گویند. در نمایه، هر کلمه همراه با سندهای مرتبط با آن، مکان رخداد آن کلمه در هر سند، و وزن آن کلمه در هر سند ذخیره می‌شود.

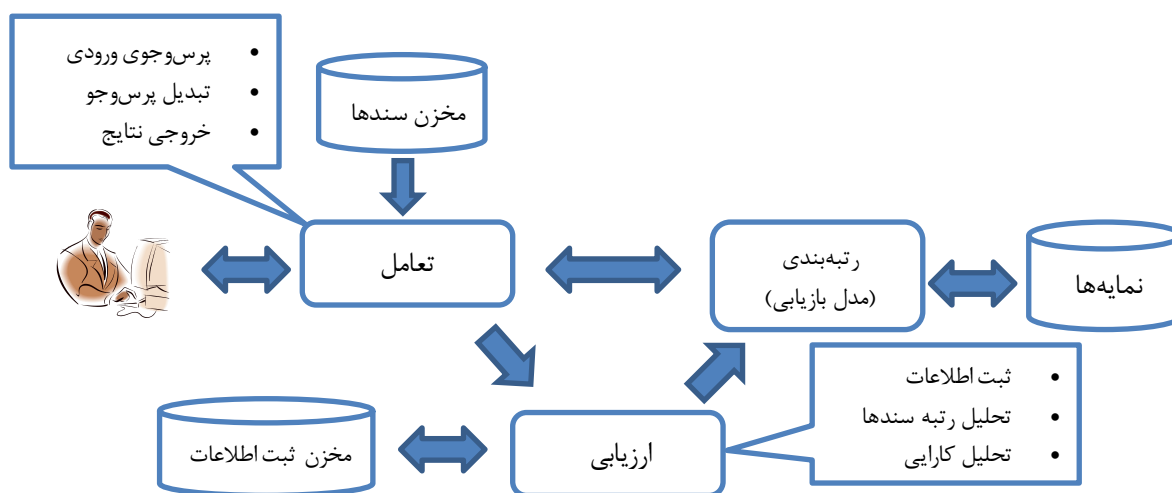
^۱ parse

^۲ stop word

^۳ stem

^۴ information retrieval

^۵ inversion



شکل ۱-۲: فرایند پرس‌وجو در موتور جستجو [۱].

۱-۱-۲ پرس‌وجو و پاسخ

در شکل ۱-۲، فرایند پرس‌وجو و پاسخ آمده است. در زیرفرایند تعامل کاربر، موتور جستجو یک رابط کاربری، برای ارسال رشته پرس‌وجو، به کاربر می‌دهد. بیشتر پرس‌وجوهای کاربران به شکل تعدادی کلمه کلیدی هستند. با وجود این، ممکن است زبان پرس‌وجو امکاناتی برای پرس‌وجوهای دقیق‌تر داشته باشد، مانند الزام وجود همه کلمه‌های جستجو شده در نتایج یافت شده یا امکان جستجو در یک سایت خاص. موتور جستجو با تجزیه پرس‌وجو، ریشه‌سازی، و حذف کلمه‌های ایست، رشته پرس‌وجو را به رشته‌ای دیگر تبدیل می‌کند. در این قسمت، اشتباه‌های نوشتاری نیز بررسی می‌شوند. علاوه بر آن، با اضافه کردن تعدادی کلمه به رشته جستجوی کاربر، تعدادی پرس‌وجوی دقیق‌تر به کاربر پیشنهاد داده می‌شود. در پایان، نتایج برگردانده شده از موتور جستجو، شامل ترتیبی از پیوندها و خلاصه‌ای از اطلاعات موجود در سندهای یافت شده، به کاربر نمایش داده می‌شود.

الگوریتم‌های رتبه‌بندی^۱، با استفاده از نمایه حاصل از فرایند نمایه‌سازی، سندها را بر اساس میزان ارتباط با پرس‌وجوی کاربر (بر مبنای یک مدل بازیابی) رتبه‌بندی می‌کنند. ممکن است کاربران از جزئیات الگوریتم‌های رتبه‌بندی و مدل بازیابی بی‌اطلاع باشند.

^۱ ranking

در زیرفرایند ارزیابی، عملکرد موتور جستجو از نظر زمان^۱ و کیفیت پاسخ به جستجوی کاربران ارزیابی می‌شود. برای این کار، تمام رفتار جستجوی کاربران، شامل پرس‌وجوهای انجام شده، پیوندهای کلیک شده پس از دریافت پاسخ جستجو، و مدت زمان صرف شده برای هر سند، ثبت می‌شوند. از نتایج ارزیابی می‌توان در بهبود الگوریتم رتبه‌بندی و مدل بازیابی اطلاعات بهره برد. همچنین، موتور جستجو از اطلاعات به کار گرفته شده در ارزیابی در موارد زیر استفاده می‌کند:

- پیشنهاد رشته‌های پرس‌وجوی جدید به کاربران
- بررسی اشتباه نوشتاری در رشته جستجوی کاربران
- ارائه نتایج جستجو به صورت شخصی‌سازی شده به کاربران

۳-۱-۱ شخصی‌سازی جستجوی وب

موتورهای جستجو برای ارائه نتایج دقیق‌تر به کاربران، پاسخ به پرس‌وجوها را بر اساس پسندهای کاربران شخصی‌سازی می‌کنند. بنابراین، از آنجایی که پسندهای کاربران متفاوت است، پرس‌وجوهای یکسان توسط کاربران متفاوت نتایجی متفاوت خواهد داشت. برای شخصی‌سازی، موتور جستجو بر اساس یک مدل کاربر^۲، پروفایل کاربران را تشکیل می‌دهد. پروفایل کاربر از روی اطلاعات زمینه^۳، مانند مکان کاربر، زمان، خصوصیات شخصیتی، دانش، سن، و جنسیت کاربر، و نیز اطلاعات مربوط به رفتار کاربر در جستجوی سندها (رفتار جستجو^۴)، شامل پرس‌وجوهای کاربر، سندهای مراجعه شده، و مدت زمان و ترتیب پرداختن به سندها، ساخته می‌شود.

با توجه به طول دوره زمانی در نظر گرفته شده برای جمع‌آوری اطلاعات، پروفایل کاربر به یکی از دو گونه کوتاه‌مدت یا بلندمدت خواهد بود [۲]. در پروفایل کوتاه‌مدت، اطلاعات مربوط به کاربر در طول یک جلسه جستجو استخراج می‌شود. این نوع پروفایل زمانی مفید است که کاربران بخواهند پرس‌وجوهایی بر خلاف پرس‌وجوهای مرسوم خود داشته باشند. برای مثال، یک پزشک بخواهد در مورد یک موضوع حقوقی جستجو کند. در مقابل، پروفایل بلندمدت بر پسندهای بلندمدت کاربر دلالت دارد.

^۱ response time

^۲ user model

^۳ context information

^۴ search behavior

موتورهای جستجو از دو روش برای شخصی سازی نتایج جستجوها استفاده می کنند. در روش اول، پیوندهای برگردانده شده با توجه به پروفایل کاربر مرتب می شوند. در روش دوم، بر اساس پروفایل کاربر، تعدادی کلمه به رشته پرس و جوی کاربر اضافه شده، و آنگاه، پرس و جوی گسترش یافته^۱ پردازش می شود.

۴-۱-۱ حریم خصوصی در جستجوی وب شخصی سازی شده

ساختن پروفایل کاربر برای شخصی کردن نتایج جستجو لازم است. این پروفایل از روی اطلاعات زمینه و رفتار جستجوی کاربر ساخته می شود. از آنجایی که این اطلاعات اغلب خصوصی^۲ است، حفظ حریم خصوصی کاربران در جستجوی وب شخصی سازی شده اهمیت پیدا می کند. منظور از حریم خصوصی، حقوق افراد در کنترل اطلاعات خصوصی آنان و جلوگیری از دسترسی غیر مجاز به آنها است [۳]. کاربران ممکن است خط مشی های حریم خصوصی متفاوتی داشته باشند. ممکن است کاربری نسبت به ذخیره شدن اطلاعات خود توسط موتور جستجو بی تفاوت باشد، و در مقابل، کاربر دیگری انتظار داشته باشد که موتور جستجو هیچ گونه اطلاعات شخصی او را ذخیره نکند. البته کاربرانی هم هستند که نسبت به شناسایی برخی پسندهای خود حساس ترند. برای مثال، کاربری درباره شناسایی شدن گرایش سیاسی خود حساس است ولی در مورد علاقه ورزشی خود حساسیتی ندارد.

در بیشتر موتورهای جستجو، امکانی وجود دارد که کاربران می توانند ذخیره سازی اطلاعات خود را غیر فعال کنند. با این کار، موتور جستجو تعهد داده است که هیچ گونه داده ای درباره کاربران ذخیره نمی کند. بنابراین، موتور جستجو هیچ اطلاع قبلی از پسندهای کاربران نخواهد داشت. با انتخاب این امکان، کاربران از ویژگی مفید شخصی سازی نتایج جستجو محروم می شوند. بنابراین، باید طرحی ارائه شود که توازن بین حفظ حریم خصوصی کاربران و شخصی سازی نتایج جستجوی وب فراهم آورد. به عبارت دیگر، کاربران باید بتوانند میزان مجاز افشای پسندهای خود را تعیین کنند.

سوالاتی که ممکن است در اینجا مطرح شود این است که چگونه می توان هم پسندهای کاربر را برای حفظ حریم خصوصی او پنهان نمود و هم از خدمت شخصی سازی نتایج استفاده کرد. برای پاسخ به این سوال از یک مثال استفاده می کنیم. فرض کنیم، پسند حساس کاربر گرایش سیاسی او است. سازوکاری که برای پنهان سازی گرایش سیاسی کاربر پیاده سازی می شود، به شکلی است که موتور جستجو گرایش سیاسی دقیق کاربر را

^۱ expanded query

^۲ private

شناسایی نمی‌کند. با وجود این، موتور جستجو از علاقمندی کاربر به سیاست آگاهی پیدا می‌کند. آنگاه، موتور جستجو می‌تواند با توجه به علاقمندی کاربر به سیاست نتایج را شخصی‌سازی کند، بدون اینکه به طور دقیق بداند که به کدام گرایش سیاسی علاقمند است.

یکی از مشکلات در عملی کردن^۱ خط‌مشی حریم خصوصی در جستجوی وب شخصی‌سازی شده آن است که موتور جستجو یک سامانه جعبه‌سیاه است. به عبارت دیگر، الگوریتم‌های به کار گرفته شده در موتور جستجو نامعلوم فرض می‌شوند. در این حالت، فقط می‌توان پرس‌وجو ارسال نمود و پاسخ آن را دریافت کرد. همچنین، فرض بر آن است که موتور جستجو یک موجودیت درستکار ولی کنجکاو^۲ است. چنین موجودیتی پاسخ پرس‌وجوها را همیشه به درستی می‌دهد ولی اطلاعات حساس کاربران را نیز نگهداری می‌کند.

پژوهش‌هایی در رابطه با مسئله حفظ حریم خصوصی در بازیابی اطلاعات شخصی‌سازی شده^۳، با تأکید بر جعبه‌سیاه بودن موتور جستجو، انجام شده است [۴]–[۲۰]. در برخی از این پژوهش‌ها، پروتکلی جدید و در برخی دیگر، طرحی بر مبنای مدل‌های شناخته شده حریم خصوصی پیشنهاد شده است. برخی از مدل‌های حریم خصوصی به کار گرفته شده k -بی‌نامی^۴ [۵]، [۷]، [۱۶]، [۱۷]، [۱۸]، [۱۳]، [۱۴]، [۲۰]، l -تنوع^۵ (آنتروپی [۱۵])، و انکارپذیری قابل قبول^۶ [۹] هستند. در روش‌های ارائه شده، تعدادی پرس‌وجوی پوششی به همراه پرس‌وجوی اصلی ارسال شده و سربرار زیادی بر موتور جستجو تحمیل می‌شود. گفتنی است که مدل‌های حریم خصوصی شناخته شده‌ای مانند l -تنوع (متمایز^۷، احتمالاتی و بازگشتی)، (t, n) -نزدیکی^۸، δ -حضور^۹، حریم خصوصی تفاضلی^{۱۰} و پاسخ تصادفی شده^{۱۱} مورد استفاده قرار نگرفته‌اند. نگاشت مدل‌های حریم خصوصی به بازیابی اطلاعات شخصی‌سازی شده، به گونه‌ای که بازتاب‌دهنده مفهوم درستی از حریم خصوصی باشد، خود یک مسئله مهم است.

حریم خصوصی مانند بسیاری از مفاهیم امنیت ماهیتی کیفی دارد و تعریف دقیقی از آن ارائه نشده است. پژوهش‌گران با توجه به دیدگاه و برداشت خود از حریم خصوصی برای آن مدلی ارائه کرده‌اند. برای مثال،

^۱ enforcement

^۲ honest-but-curious

^۳ personalized information retrieval (PIR)

^۴ k -anonymity

^۵ l -diversity

^۶ plausible deniability

^۷ distinct

^۸ (t, n) -closeness

^۹ δ -presense

^{۱۰} differential privacy

^{۱۱} randomized response

مدل‌های k -بی‌نامی یا مدل حریم خصوصی تفاضلی ارائه شده‌اند. در مدل حریم خصوصی تفاضلی به پایگاه‌داده آماری^۱ توجه می‌شود. منظور از پایگاه‌داده آماری، پایگاه‌داده‌ای است که پرس‌وجوهای آماری مانند میانگین، واریانس و غیره روی آن اجرا می‌شود. همچنین، در مدل حریم خصوصی تفاضلی پاسخ به پرس‌وجوهای آماری تحلیل‌گر مغشوش شده و در اختیار او قرار می‌گیرد. بنابراین، منظور از حریم خصوصی در این مدل، افشا نشدن داده‌های حساس کاربران در پاسخ به پرس‌وجوهای آماری تحلیل‌گر است. در مدل‌های مبتنی بر افراز کردن (مانند k -بی‌نامی)، کل پایگاه‌داده در اختیار تحلیل‌گر قرار گرفته و منتشر می‌شود. به این دلیل، رکوردهای پایگاه‌داده به شکلی تغییر پیدا می‌کنند که خط‌مشی حریم خصوصی هیچ کاربری نقض نشود. بنابراین، منظور از حریم خصوصی در این مدل، افشا نشدن داده‌های حساس موجود در تک‌تک رکوردهای پایگاه‌داده است. همان‌طور که دیده می‌شود، ماهیت مدل حریم خصوصی تفاضلی با مدل‌های مبتنی بر افراز کردن متفاوت است و به نظر می‌رسد مقایسه آن‌ها درست نباشد. با وجود این، در ادبیات حوزه حریم خصوصی موارد زیر به عنوان اشکال‌های مدل‌های مبتنی بر افراز کردن اشاره شده است.

مدل‌های مبتنی بر افراز کردن^۲، بر اساس مدلی شکل می‌گیرند که از دانش پس‌زمینه^۳ مهاجم وجود دارد. بر اساس این مدل، مشخص می‌شود که چه داده‌هایی از پایگاه‌داده باید تغییر پیدا کنند. اشکال اساسی که به این مدل‌ها وارد می‌شود، این موضوع است که مدلی که از دانش پس‌زمینه ساخته می‌شود، ممکن است همیشه کامل نباشد. در نتیجه، مدل حریم خصوصی به اطلاعات پس‌زمینه‌ای که در نظر گرفته نشده است، آسیب‌پذیر است. همچنین، همان‌طور که در قسمت ۲-۲ از فصل ۲ پیشنهاد رساله توضیح داده شده است، حمله‌های کمینه کردن^۴، ترکیب^۵ و پیش‌زمینه^۶ نیز بر روی مدل‌های مبتنی بر افراز کردن قابل اجرا است. گفتنی است، حمله‌های کمینه کردن و پیش‌زمینه بر مدل حریم خصوصی تفاضلی قابل اجرا نیستند، زیرا در این مدل بر خلاف مدل‌های مبتنی بر افراز کردن، کل پایگاه‌داده منتشر نشده و تنها پاسخ به پرس‌وجوهای آماری مغشوش می‌شوند. همچنین، مدل حریم خصوصی تفاضلی نسبت به اطلاعات پیش‌زمینه مقاوم است و مستقل از این اطلاعات حریم خصوصی کاربران را حفظ می‌کند (ماهیت مدل حریم خصوصی تفاضلی متفاوت است). گفتنی

^۱ statistical database

^۲ partition-based model

^۳ background knowledge

^۴ minimality attack

^۵ compositional attack

^۶ foreground attack

است، در مدل حریم خصوصی تفاضلی قضیه‌های ترکیب^۱ وجود دارند. این قضیه‌ها، باعث می‌شوند که مدل حریم خصوصی تفاضلی نسبت به حمله ترکیب نیز مقاوم باشد.

در مسئله حفظ حریم خصوصی کاربران در جستجوی وب شخصی‌سازی شده، دیدگاه ما نسبت به مسئله بدین شرح است. موتور جستجو (تحلیل‌گر) تاریخچه جستجوهای کاربر (پایگاه‌داده) را به طور کامل در اختیار دارد. همچنین، موتور جستجو پرس‌وجوی خود را روی پایگاه‌داده اجرا می‌کند. پرس‌وجوی موتور جستجو که الگوریتم ایجاد پروفایل کاربر است، ماهیتی آماری دارد. بنابراین، پایگاه‌داده در این مسئله آماری است. از آنجایی که موتور جستجو تمام پایگاه‌داده را در اختیار دارد، حفظ اطلاعات حساس در تک‌تک رکوردهای آن معنایی ندارد و هدف ما از حفظ حریم خصوصی جلوگیری از افشا شدن اطلاعات حساس کاربر در پاسخ به پرس‌وجوی موتور جستجو (ایجاد پروفایل کاربر) است. بنابراین، مناسب‌ترین مدل برای حفظ حریم خصوصی در این دیدگاه مدل حریم خصوصی تفاضلی است.

در بعضی از پژوهش‌های پیشین، از مدل‌های مبتنی بر افراز کردن نیز به دلیل برداشت متفاوت از حریم خصوصی استفاده شده است. در پژوهش مرجع [۵]، $k-1$ پرس‌وجوی پوششی به منظور پنهان کردن پرس‌وجوی اصلی کاربر ارسال می‌شود. در پژوهش مرجع [۱۵]، از مدل l -تنوع آنتروپی استفاده می‌شود. در این پژوهش، $k-1$ کلمه کلیدی به پرس‌وجوی جستجوی کاربر اضافه می‌شود تا کلمه کلیدی اصلی کاربر پنهان بماند. در پژوهش مرجع [۱۲]، به ازای تمام پرس‌وجوهای جستجوی کاربر $k-1$ پرس‌وجوی پوششی ایجاد و ارسال می‌شود. این سازوکار، به شدت سودمندی طرح را تحت تاثیر قرار می‌دهد. در پژوهش‌های مرجع [۹] و [۱۰] از مدل انکارپذیری قابل قبول استفاده شده است. در این پژوهش‌ها، هدف این است که کاربر بتواند علاقه‌مندی خود به یک موضوع حساس را با وجود $m-1$ موضوع پوششی دیگر در پروفایل خود به طور قابل قبول انکار کند. علاقه‌مندی کاربر به این m موضوع باید تقریباً به اندازه یکدیگر باشد. در نتیجه، کاربر می‌تواند به طور قابل قبول علاقه‌مندی خود به موضوع حساس را انکار کند. نقد و بررسی پژوهش‌های پیشین در فصل ۳ قسمت ۳-۲ (نقد و بررسی پژوهش‌های پیشین) آمده است. طرح‌های پیشین حفظ حریم خصوصی در جستجوی وب شخصی‌سازی شده، هیچ فرضی در رابطه با روش به دست آوردن پروفایل کاربران، در موتور جستجو، وجود ندارد. با وجود این، پرس‌وجوهای پوششی بر اساس آماره‌هایی مانند TF-IDF و آنتروپی کلمه‌های موجود در پرس‌وجوی اصلی تولید می‌شوند. این موضوع، به دلیل اینکه به روش به کار

^۱ composition theorem

گرفته شده در موتور جستجو برای تشکیل پروفایل کاربران توجه نمی‌شود، منجر به عدم کارایی این طرح‌ها می‌شود.

۱-۲ اهداف رساله

در این پیشنهاد رساله، راه‌کار پیشنهادی خود را برای حفظ حریم خصوصی کاربران در جستجوی وب شخصی-سازی شده مطرح خواهیم کرد. پیشنهاد ما بر اساس مدل حریم خصوصی تفاضلی است. موتور جستجو تمامی فعالیت‌های کاربران، دربرگیرنده اطلاعات زمینه و رفتار جستجوی کاربران، را جمع‌آوری می‌کند. این اطلاعات در پایگاه‌داده‌ای ذخیره می‌شوند که هر رکورد آن شامل پرس‌وجوی جستجوی کاربر و نیز شناسه سندهایی است که کاربر در پی اعلام نتایج جستجو به سند مربوط به آن‌ها مراجعه نموده است. کاربران باید با مغشوش نمودن اطلاعات خود (اضافه کردن نویز) پایگاه‌داده را به شکل مغشوش شده سنتز کنند.^۱ به بیان دیگر، کاربران بر اساس دانش خود از پرس‌وجوهایی که موتور جستجو از پایگاه‌داده خواهد داشت و همچنین، خط‌مشی حریم خصوصی خود، با اضافه کردن تعدادی رکورد علاوه بر هر رکورد اصلی یا تغییر هر رکورد اصلی در زمان جستجو، پایگاه‌داده مغشوش شده را ایجاد می‌کنند.

برای سنتز پایگاه‌داده مغشوش شده بر اساس مدل حریم خصوصی تفاضلی در جستجوی وب شخصی‌سازی شده، باید پرس‌وجوی موتور جستجو بر روی پایگاه‌داده را شناسایی کنیم. پس از شناسایی پرس‌وجوی موتور جستجو، باید مفاهیم همسایگی^۲، حساسیت^۳، سازوکار تصادفی کردن^۴، و شیوه تحلیل و اندازه‌گیری سودمندی^۵ و اتلاف حریم خصوصی^۶ را به شکل صوری بیان کنیم. در سنتز پایگاه‌داده مغشوش شده، متناسب با حساسیت پرس‌وجوهای تحلیل‌گر (موتور جستجو)، به پایگاه‌داده نویز اضافه می‌شود. ارائه طرحی با کمترین مقدار نویز اضافه شده، به صورتی که سطح قابل قبولی از مفید بودن حفظ شده و خط‌مشی حریم خصوصی نیز برآورده شود، یک چالش در حل این مسئله است. برای اضافه کردن نویز، کاربران می‌توانند تعدادی رکورد پوششی علاوه بر رکوردهای اصلی خود به پایگاه‌داده اضافه کنند. همچنین می‌توانند، رکوردهای اصلی خود را، در جهت حفظ حریم خصوصی، با گسترش پرس‌وجوی جستجوی خود یا تغییر سندهای مراجعه شده در پی جستجو مغشوش کنند.

^۱ synthesize

^۲ neighborhood

^۳ sensitivity

^۴ randomization mechanism

^۵ utility

^۶ privacy loss

رکوردهای مختلف پایگاه داده، در مسئله حفظ حریم خصوصی در جستجوی وب شخصی سازی شده، همبستگی^۱ دارند. به عبارت دیگر، به دلیل ارتباط رکوردهای پایگاه داده با هم، در صورت حذف یا تغییر یک رکورد از پایگاه داده، از سایر رکوردهای موجود در پایگاه داده می توان درباره آن رکورد تغییر یافته یا حذف شده اطلاعاتی به دست آورد. بنابراین، برخلاف طرح اصلی حریم خصوصی تفاضلی که رکوردها را مستقل از هم در نظر می گیرد، باید ارتباط بین رکوردها را نیز در حفظ حریم خصوصی در نظر بگیریم.

پایگاه داده اطلاعات مربوط به جستجوهای کاربران اندازه ثابتی ندارد. به عبارت دیگر، تعداد رکوردهای این پایگاه داده با ارسال هر پرس و جوی جستجوی کاربران افزایش پیدا می کند. همچنین، موتور جستجو پرس و جوی خود را روی آخرین نسخه پایگاه داده، قبل از دریافت یک پرس و جوی جستجوی جدید، اجرا می کند تا پروفایل کاربر را به دست آورد. در حریم خصوصی تفاضلی، چنانکه در تعریف ابتدایی در نظر گرفته شده است، تعداد رکوردهای پایگاه داده ثابت در نظر گرفته می شود. در مرجع [۲۱]، راه کاری برای عملی کردن حریم خصوصی تفاضلی در پایگاه داده های رشدیابنده^۲ ارائه شده است که در آن با بزرگ تر شدن پایگاه داده، میزان نویز اضافه شده به پاسخ ها بر اساس آخرین وضعیت پایگاه داده و نیز نویزهای اضافه شده به پاسخ های قبلی محاسبه می شود. استفاده از ایده مطرح در این مرجع را باید در سنتز پایگاه داده مغشوش شده بررسی کنیم.

در حریم خصوصی تفاضلی، میزان اتلاف حریم خصوصی متناظر با فاش شدن هر یک از رکوردهای پایگاه داده یکسان در نظر گرفته می شود. به عبارت دیگر، سازوکار حفظ حریم خصوصی برای تمام رکوردها به صورت یکسان عمل می کند. در مسئله حفظ حریم خصوصی در جستجوی وب شخصی سازی شده، اغلب کاربران پسندهای حساس محدود و مشخصی دارند. بنابراین، میزان اتلاف حریم خصوصی برای رکوردهای مختلف پایگاه داده یکسان نیست. سازوکار تصادفی کردن را باید با در نظر گرفتن تنوع اتلاف حریم خصوصی^۳ [۲۲] طراحی کنیم. بر این اساس، مغشوش کردن اطلاعات را می توان فقط برای حفاظت از پسندهای حساس کاربر انجام داد. طراحی چنین سازوکاری، چالش دیگری در حفظ حریم خصوصی در جستجوی وب شخصی سازی شده است. تأثیر رشدیابنده بودن پایگاه داده بر عملی کردن حریم خصوصی تفاضلی با وجود تنوع در اتلاف حریم خصوصی بررسی نشده است و در این رساله به آن پرداخته می شود.

تولید و ارسال پرس و جوی پوششی و گسترش یافته باید بر اساس روش یادگیری پروفایل در موتور جستجو انجام شود. با وجود این، موتور جستجو یک سامانه جعبه خاکستری یا جعبه سیاه است و روشی را که

^۱ correlation^۲ growing database^۳ variation of privacy loss

برای یادگیری پروفایل کاربر استفاده می‌کند و مدل بازیابی استفاده شده در آن، (به صورت کامل) مشخص نیست. همچنین، شیوه و میزان تاثیر تاریخچه جستجوی کاربر در شخصی‌سازی نتایج جستجو که در پروفایل کوتاه‌مدت و بلندمدت نمود پیدا می‌کند، مشخص نیست. همان‌طور که در بالا گفته شد، علاوه بر مشخص نبودن الگوریتم یادگیری پروفایل کاربر، چالش‌های دیگری شامل همبستگی رکوردها، رشدیابنده بودن پایگاه‌داده، و یکسان نبودن ائتلاف حریم خصوصی برای موضوع‌های مختلف وجود دارند. بنابراین، برای سادگی روند حل مسئله، آن را در سه گام تعریف می‌کنیم. در گام اول حل مسئله، موتور جستجو جعبه‌سفید فرض می‌شود. به عبارت دیگر، الگوریتم یادگیری پروفایل کاربر و الگوریتم بازیابی در موتور جستجو به طور کامل مشخص هستند. در گام دوم، موتور جستجو جعبه‌خاکستری فرض می‌شود. به بیان دیگر، فرض می‌کنیم یک توزیع احتمال روی مجموعه‌ای از روش‌های یادگیری پروفایل کاربر داده شده است. ممکن است این توزیع بر روی مقادیر مختلف پارامتری باشد که در روش یادگیری موثر است. همچنین، در این گام الگوریتم بازیابی دانسته فرض می‌شود.

در گام سوم، موتور جستجو را جعبه‌سیاه فرض می‌کنیم. با استفاده از فنون یادگیری ماشین می‌توان روش یادگیری پروفایل کاربر را به دست آورد. گفتنی است، در به دست آوردن روش یادگیری پروفایل کاربر، نیاز است که به پروفایل کاربر دسترسی داشته باشیم. همان‌طور که گفته شد، موتور جستجو یک سامانه جعبه‌سیاه است و امکان محاسبه مستقیم پروفایل کاربر وجود ندارد. با وجود این، با مقایسه نتایج شخصی‌سازی شده و شخصی‌سازی نشده در پاسخ به پرس‌وجوهای کاربر، می‌توانیم پروفایل کاربر را پیدا کنیم [۷]. نکته مهم این است که به دست آوردن پروفایل کاربر و روش یادگیری آن، دارای خطا است. بنابراین، مدل حریم خصوصی تفاضلی را باید با در نظر گرفتن وجود خطا طراحی کنیم. این موضوع در میزان تصادفی کردن پرس‌وجوها و، در نتیجه، در مفید بودن نتایج جستجوی شخصی‌سازی شده تأثیرگذار است. گفتنی است، روش مبتنی بر فنون یادگیری ماشین که در گام سوم استفاده می‌شود، در گام دوم نیز قابل استفاده است. با وجود این، رویکرد ما در گام دوم مبتنی بر نظریه احتمالات است. مقایسه نتایج گام دوم و سوم از جهت میزان سودمندی و حفظ حریم خصوصی بخش دیگری از رساله خواهد بود.

بنابراین اهداف ما در تعریف رساله پیش‌رو به صورت زیر تعریف می‌شود:

- ارائه طرحی الهام گرفته از روش سنتز پایگاه‌داده مغشوش شده مبتنی بر مدل حریم خصوصی تفاضلی در جستجوی وب شخصی‌سازی شده با فرض مشخص بودن روش یادگیری پروفایل کاربر. این قسمت، گام نخست در این رساله است و در انجام آن دو چالش زیر وجود دارد.

○ همبستگی رکوردهای پایگاه داده

○ پایگاه داده رشدیابنده

- ارائه طرحی الهام گرفته از روش سنتز پایگاه داده مغشوش شده مبتنی بر مدل حریم خصوصی تفاضلی با فرض مشخص بودن پرس و جوی مسئول پایگاه داده به صورت تخمینی.
- ارائه طرحی الهام گرفته از روش سنتز پایگاه داده مغشوش شده مبتنی بر مدل حریم خصوصی تفاضلی در جستجوی وب شخصی سازی شده با در نظر گرفتن موتور جستجو به صورت یک سامانه جعبه سیاه.

بیانیه رساله نیز با توجه به اهداف گفته شده به صورت زیر تعریف می شود.

« در جستجوی وب شخصی سازی شده، موتور جستجو اطلاعات رفتار جستجوی هر کاربر را در یک پایگاه داده (تاریخچه جستجوها) نگهداری می کند. موتور جستجو با اجرای الگوریتم یادگیری پروفایل کاربر بر روی پایگاه داده، پروفایل کاربر را استخراج کرده و بر اساس آن نتایج جستجو را شخصی سازی می کند. رکوردهای پایگاه داده شامل اطلاعات رفتار جستجوی کاربر هستند و به مرور افزایش پیدا می کنند. همچنین، رکوردها با هم همبستگی دارند. از آنجایی که، الگوریتم یادگیری پروفایل کاربر ماهیت آماری دارد و هدف از حفظ حریم خصوصی، نمایان نشدن موضوعهای حساس مورد علاقه کاربر در نتایج شخصی سازی شده او توسط موتور جستجو است، مدل حریم خصوصی تفاضلی را به کار می بریم. با ارائه تعریف جدیدی از دو تاریخچه جستجوی همسایه، می توان همبستگی رکوردها و رشدیابنده بودن تاریخچه جستجوها را در نظر گرفت. همچنین، با روشی الهام گرفته شده از سنتز پایگاه داده، حتی بدون داشتن اطلاعات کامل از الگوریتم یادگیری پروفایل کاربر توسط موتور جستجو، می توان هم خط مشی حریم خصوصی کاربران را برآورده نمود و هم به خدمت شخصی سازی دارای سودمندی برای کاربر دست یافت.»

۱-۳ ساختار مطالب پیشنهاد رساله

ساختار مطالب پیشنهاد رساله به این شرح است. در فصل ۲، مطالب پیش زمینه لازم توضیح داده خواهند شد. در قسمت ۱-۲، مفهوم حریم خصوصی بیان می شود. در قسمت ۲-۲، مدل های مبتنی بر افراز کردن مانند k -بی نامی بررسی می شوند. در قسمت ۲-۳، مدل های مبتنی بر تصادفی کردن مانند حریم خصوصی تفاضلی شرح داده می شوند. در قسمت ۲-۴، روش های ایجاد پروفایل کاربر توضیح داده می شوند. در پایان در قسمت ۲-۵ آنتولوژی محاسباتی بیان خواهد شد. در فصل ۳، کارهای پیشین این پیشنهاد رساله بررسی شده اند. در

قسمت ۱-۳، پژوهش‌های پیشین حفظ حریم خصوصی در جستجوی وب شخصی‌سازی شده بر اساس مدل‌های حریم خصوصی به‌کار گرفته شده، نقد و بررسی شده‌اند. در قسمت ۰، پژوهش‌های مرتبط با چالش‌های حل مسئله مطرح شده آمده‌اند. در فصل ۴، راه‌کار پیشنهادی برای حل مسئله مطرح شده در این پیشنهاد رساله توضیح داده می‌شود. در قسمت ۱-۴، به بیان کلی مسئله و گام‌های پیشنهادی حل آن می‌پردازیم. در قسمت ۲-۴، مسئله شرح داده شده و چالش‌های موجود در پیش‌برد رساله بیان می‌شوند. در قسمت ۴-۲-۴، هدف‌های رساله و زمان‌بندی انجام آن آورده شده‌اند.

فصل دوم

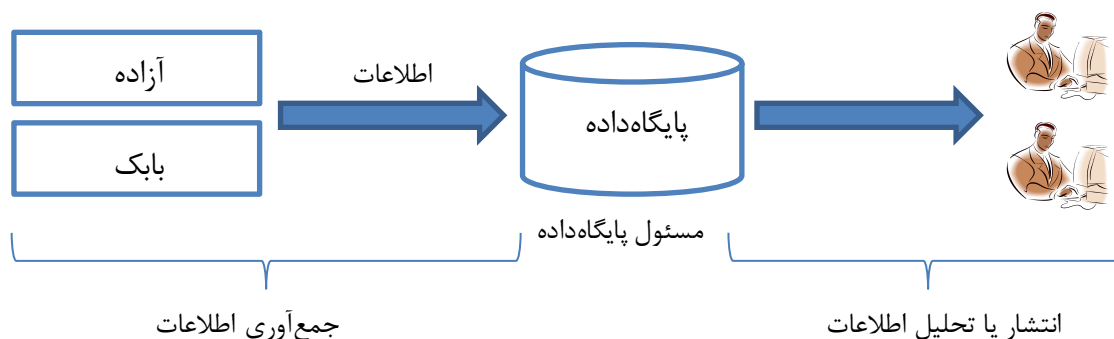
پیش زمینه

در این فصل، مطالب پیش‌زمینه لازم برای تعریف پیشنهاد رساله توضیح داده خواهند شد. در این فصل روش‌های مختلف حفظ حریم خصوصی افراد در نگهداری و انتشار یا تحلیل اطلاعات خصوصی افراد توضیح داده می‌شود.

۲-۱ حریم خصوصی

در چند دهه گذشته سرعت تولید اطلاعات افزایش یافته است. دولت‌ها، سازمان‌ها و تجارت‌های مختلف اطلاعات مختلف و زیادی در رابطه با افراد جمع‌آوری کرده و آن‌ها را در پایگاه‌داده‌ها نگهداری می‌کنند. مسئول پایگاه‌داده، وظیفه حفظ و انتشار آن‌ها را برای تحلیل‌های بعدی بر عهده دارد. افراد انتظار دارند به میزان حساسیت خود نسبت به افشا شدن اطلاعات خود حریم خصوصی داشته باشند. در لغت‌نامه وبستر، لغت حریم خصوصی همراه با معانی خلوت و محفوظ از نفوذ غیر مجاز آمده است. همچنین، در مرجع [۳]، تعریف حریم خصوصی به شکل حقوق افراد در کنترل اطلاعات خصوصی آنان و جلوگیری از دسترسی غیر مجاز به آن‌ها آمده است.

همان‌طور که در شکل ۲-۱ مشخص است، فرآیند جمع‌آوری و تحلیل اطلاعات در دو مرحله اتفاق می‌افتد. مرحله اول، مرحله جمع‌آوری اطلاعات است که مسئول پایگاه‌داده، داده‌های افراد را دریافت کرده و نگهداری می‌کند. اگر مسئول پایگاه‌داده مورد اعتماد نباشد، در همین مرحله حریم خصوصی افراد نقض می‌شود. با وجود این، در بسیاری از سناریوها مسئول پایگاه‌داده مورد اعتماد است. مرحله دوم، مرحله انتشار یا تحلیل اطلاعات است. در این مرحله، مسئول پایگاه‌داده به شکل خصوصی به پرس‌وجوی افراد عادی پاسخ می‌دهد یا اینکه کل پایگاه‌داده را برای تحلیل و پژوهش خصوصی کرده و آنگاه منتشر می‌کند. به طور پیش‌فرض، افراد عادی مورد اعتماد نیستند. در نتیجه فرایند تحلیل و انتشار باید حریم خصوصی افراد را حفظ کند.



شکل ۲-۱: حفظ حریم خصوصی در نگهداری، تحلیل و انتشار اطلاعات [۷۱]

در ادامه این فصل، روش‌های مختلف حفظ حریم خصوصی افراد شرح داده می‌شود. در روش اول، حفظ حریم خصوصی مبتنی بر افراز کردن^۱ توضیح داده می‌شود. در این روش، رکوردهای پایگاه داده بر اساس یک ملاک خاص، تعدادی از ویژگی‌های هر رکورد بی‌نام می‌شوند و رکوردها به دسته‌های متمایز افراز می‌شود. آنگاه، کل پایگاه داده منتشر شده و در اختیار پژوهشگران قرار می‌گیرد. در روش دوم، حفظ حریم خصوصی مبتنی بر تصادفی کردن^۲ شرح داده خواهد شد. در این روش مالکان به اطلاعات ارسالی خود در هنگام جمع‌آوری آن‌ها نویز اضافه می‌کنند یا مسئول پایگاه داده به نتیجه پرس‌وجوی ارسالی توسط پژوهشگران نویز اضافه می‌کند. این تقسیم‌بندی از مقاله مرجع [۲۳] گرفته شده است.

۲-۲ حفظ حریم خصوصی مبتنی بر افراز کردن

در مدل‌های مبتنی بر افراز کردن، به طور کلی ستون‌های جدول پایگاه داده به سه دسته تقسیم می‌شوند.

- ستون‌هایی که به طور مستقیم موجب ارتباط رکورد با یک فرد خاص می‌شوند مانند ستون نام و کد ملی.
- ستون‌هایی که به طور غیرمستقیم و از طریق جدولی کمکی، موجب ارتباط رکورد با یک فرد خاص می‌شوند. به این ستون‌ها مانند جنسیت، کدپستی، و تاریخ تولد شبه‌شناسه^۳ گفته می‌شود.
- ستون داده‌های حساس مانند ستون نوع بیماری، مقدار حقوق و مانند آن.

در این مدل‌ها ستون‌های دسته اول حذف می‌شوند و اطلاعات ستون‌های دسته دوم، عمومی می‌شوند^۴. عمومی شدن داده‌های ستون‌های دسته دوم (شبه‌شناسه‌ها) باعث می‌شود، تعدادی رکورد با مقادیر یکسان برای ستون‌های شبه‌شناسه به وجود بیاید. مجموعه رکوردها با مقدار ستون‌های شبه‌شناسه یکسان، تشکیل کلاس هم‌ارزی^۵ می‌دهند.

❖ k -بی‌نامی^۶

^۱ partition-based
^۲ randomization-based
^۳ quasi-Identifiers
^۴ generalized
^۵ equivalence class
^۶ k-anonymity

در مدل k -بی‌نامی [۲۳] و [۲۴] تعداد رکوردها در هر کلاس هم‌ارزی، حداقل باید k تا باشد. عمومی کردن مقادیر ستون‌های شبه‌شناسه هم باید به صورت کمینه^۱ صورت بگیرد تا کمترین میزان از دست دادن اطلاعات اتفاق بیفتد.

❖ l -تنوع^۲

در مدل k -بی‌نامی دو نوع حمله امکان‌پذیر است. حمله اول بر اساس همگن بودن مقدارهای ستون حساس در حداقل یکی از کلاس‌های هم‌ارزی تشکیل شده در پایگاه‌داده است. به عبارت دیگر، در یک کلاس هم‌ارزی تمام مقدارهای ستون حساس یکسان باشند. حمله دوم بر اساس اطلاعات پیشین مهاجم است. مهاجم با اطلاعات پیشین خود می‌تواند تعدادی از سطرهای نامرتب با قربانی را از کلاس هم‌ارزی مربوط به او نادیده بگیرد و در بهترین حالت قربانی را به یک سطر خاص از پایگاه‌داده ربط دهد. آنگاه، مهاجم می‌تواند مقدار حساس مربوط به قربانی را دقیقاً بدست بیاورد. برای جلوگیری از این دو حمله، مدل l -تنوع [۲۶] ارائه شده است.

تعریف: مجموعه رکوردهای پایگاه‌داده با مقدار ستون‌های شبه‌شناسه یکسان، تشکیل کلاس هم‌ارزی می‌دهند. یک کلاس هم‌ارزی گفته می‌شود که l -تنوع دارد اگر حداقل l متغیر «خوش-نمایان‌شده»^۳ در ستون حساس آن وجود داشته باشد. یک جدول گفته می‌شود که l -تنوع دارد، اگر همه کلاس‌های هم‌ارزی در آن l -تنوع داشته باشند.

«خوش-نمایان‌شده» می‌تواند تعاریف متفاوتی داشته باشد. از جمله:

۱. متمایز^۴: وجود l مقدار حساس متمایز در ستون حساس در هر کلاس هم‌ارزی.
۲. احتمالاتی: تکرار دفعات حضور یک مقدار حساس در یک کلاس هم‌ارزی حداکثر $1/l$ باشد.
۳. آنتروپی: اگر رابطه زیر برای هر کلاس هم‌ارزی برقرار باشد:

$$Entropy(E) = - \sum_{s \in S} Pr[E, s] \log Pr[E, s] \geq \log l$$

$Pr[E, s]$: احتمال حضور مقدار حساس s در کلاس هم‌ارزی E .

^۱ minimal

^۲ l-diversity

^۳ well-represented

^۴ distinct

برای اینکه کلاس هم‌ارزی دارای l -تنوع باشد، آنتروپی کل جدول باید بیشتر از $\log(l)$ باشد. گفتنی است، این تعریف سخت‌گیرانه است. به عبارت دیگر، اگر در یک جدول یک یا چند مقدار حساس، پرتکرار باشد، مقدار آنتروپی جدول کم خواهد بود و عملی کردن l -تنوع آنتروپی ممکن نخواهد بود.

۴. (c, l) -تنوع بازگشتی^۱: این تعریف تضمین می‌کند در یک کلاس هم‌ارزی مقدار پرتکرار، بیش از حد تکرار نشود و مقادیر کم‌تکرار هم نادر نباشند. اگر m مقدار حساس، در یک کلاس هم‌ارزی وجود داشته باشد و r_i تعداد تکرار i -امین $(1 \leq i \leq m)$ مقدار حساس باشد که به صورت نزولی مرتب شده‌اند، آنگاه رابطه زیر باید برقرار باشد تا کلاس هم‌ارزی، (c, l) -تنوع بازگشتی شود:

$$r_1 \leq c(r_l + r_{l+1} + \dots + r_m)$$

c در رابطه بالا، یک عدد ثابت و حقیقی بزرگتر از صفر است. با حذف کردن یک مقدار حساس از کلاس هم‌ارزی، بقیه رکوردها باید $(c, l-1)$ -تنوع بازگشتی باشند. فرض می‌شود، ۱-تنوع همیشه برقرار است.

❖ t -نزدیکی^۲

در مدل t -نزدیکی از مرجع [۲۶] و [۲۷] گفته می‌شود که توزیع احتمال مقادیر حساس در کلاس هم‌ارزی باید حداکثر به اندازه t از توزیع احتمال مقادیر حساس در کل پایگاه‌داده فاصله داشته باشد. این تعریف، نسبت به ضعف‌های مدل l -تنوع در حفظ حریم خصوصی قوی‌تر است. در مدل‌های قبلی، با توجه به پراکندگی مقادیر حساس در کلاس هم‌ارزی می‌توان استنتاج‌های آماری درباره اطلاعات افراد داشت. برای مثال، اگر سه نوع مختلف از بیماری‌های گوارشی در یک کلاس هم‌ارزی وجود داشته باشد، آنگاه، این کلاس ۳-تنوع است، با وجود این، اطلاعات فردی که در این کلاس قرار بگیرد، نشان دهنده بیماری گوارشی آن فرد است. اگر فاصله توزیع احتمال در کلاس هم‌ارزی را از توزیع احتمال یک مجموعه n تایی از رکوردهای یک پایگاه‌داده در نظر بگیریم، تعریف (n, t) -نزدیکی را خواهیم داشت. اگر n به اندازه کل پایگاه‌داده باشد، این تعریف همان تعریف t -نزدیکی خواهد شد.

❖ δ -حضور^۳

^۱ recursive diversity

^۲ closeness

^۳ presence

در مدل δ -حضور [۲۹] گفته می‌شود که احتمال حضور افراد در جدول عمومی شده، باید در بازه δ قرار گیرد. اگر جدول عمومی T^* و جدول خصوصی T (که T زیرمجموعه T^* است) داده شود، مدل δ -حضور برای جدول عمومی شده T^* برای تمام رکوردهای t برقرار است اگر:

$$\delta_{min} \leq \Pr[t \in T|T^*] \leq \delta_{max}$$

❖ انکارپذیری قابل قبول^۱

در این قسمت انکارپذیری قابل قبول [۳۰] برای حفظ حریم خصوصی در تولید پایگاه داده ساختگی^۲ توضیح داده می‌شود. در تولید پایگاه داده ساختگی، هدف این است که پایگاه داده اصلی بر اساس یک خط‌مشی حریم خصوصی به پایگاه داده‌ای جدید تبدیل شود. پایگاه داده ساختگی از تغییر رکوردهای پایگاه داده اصلی به دست می‌آید. در انکارپذیری قابل قبول با بیان غیرصوری می‌توان گفت که مهاجم با هر اطلاعات پیشین نمی‌تواند استنتاج کند که یک رکورد خاص در پایگاه داده ساختگی، از کدام رکورد خاص در پایگاه داده اصلی به دست آمده است. یک سازوکار، انکارپذیری قابل قبول را ارضا می‌کند، اگر حداقل $k > 0$ رکورد در پایگاه داده اصلی وجود داشته باشد که همگی با احتمال یکسان بتوانند، یک رکورد خاص در پایگاه داده ساختگی را تولید کنند. در این روش نیاز به اضافه کردن نویز به داده‌ها نیست. این سازوکار از دو قسمت تشکیل شده است، قسمت اول مدل‌های تولیدکننده و قسمت دوم، آزمون حریم خصوصی است. در قست اول، باید مدلی طراحی شود که با استفاده از آن بتوان رکوردهایی تولید کرد که بیشترین سودمندی را داشته باشند. قسمت دوم تضمین می‌کند که حریم خصوصی افراد حفظ می‌شود. این الگوریتم طوری طراحی می‌شود که به راحتی بتواند ارتباط هر رکورد تولیدی در پایگاه داده ساختگی با رکوردهای پایگاه داده اصلی را به طور قابل قبول انکار کند.

اگر M یک مدل تولید داده باشد که رکوردی مانند d را دریافت کند آنگاه، رکورد ساختگی y را با احتمال $\Pr[M(d) = y]$ می‌سازد. $k > 0$ یک عدد طبیعی و $\gamma > 0$ یک عدد حقیقی، پارامترهای حریم خصوصی در این مدل هستند. با استفاده از پارامترهای گفته شده انکارپذیری قابل قبول به صورت زیر تعریف می‌شود.

تعریف انکارپذیری قابل قبول: برای هر پایگاه داده D با $|D| \geq k$ و هر رکورد y که توسط مدل احتمالاتی تولیدکننده داده M تولید شده است ($d_1 \in D$ و $M(d_1) = y$)، می‌گوییم y با (k, γ) -انکارپذیری قابل قبول انتشارپذیر است، اگر حداقل $k - 1$ رکورد منحصر به فرد $d_2, d_3, \dots, d_k \in D \setminus \{d_1\}$ وجود داشته باشند که برای هر $i, j \in 1, \dots, k$ رابطه زیر برای آن‌ها برقرار باشد.

^۱ plausible deniability

^۲ synthetic database

$$\gamma^{-1} \leq \frac{\Pr[M(d_i) = y]}{\Pr[M(d_j) = y]} \leq \gamma$$

هر چه k بیشتر باشد، مجموعه رکوردهای مشابه غیرقابل تمایز بیشتر می‌شود. هر چه γ نیز به یک نزدیکتر باشد، میزان تمایزپذیری رکوردهای مشابه کمتر می‌شود.

سازوکار انکارپذیری قابل قبول: برای مدل M ، پایگاه‌داده D ، پارامترهای حریم خصوصی k و γ ، رکورد y منتشر می‌شود یا چیزی منتشر نمی‌شود. مراحل تولید و انتشار y به شکل زیر است:

۱. انتخاب تصادفی d از D .

۲. تولید رکورد y ($M(d) = y$).

۳. اجرای آزمون حریم خصوصی.

آزمون حریم خصوصی: در مرحله آزمون حریم خصوصی برای مدل M ، پایگاه‌داده D ، رکوردهای d و y ، پارامترهای حریم خصوصی k و γ ، با توجه به شرایط زیر y به عنوان خروجی منتشر می‌شود.

۱. $i \geq 0$ تنها مقدار طبیعی باشد که در رابطه زیر صدق کند:

$$\gamma^{-i-1} \leq \Pr[M(d) = y] \leq \gamma^{-i}$$

۲. $k' \geq k$ تعداد رکوردهای $d_a \in D$ باشد که در رابطه زیر صدق کند:

$$\gamma^{-i-1} \leq \Pr[M(d_a) = y] \leq \gamma^{-i}$$

۳. اگر $k' \geq k$ باشد آنگاه آزمون موفق بوده و y منتشر می‌شود.

❖ حمله به مدل‌های مبتنی بر افراز کردن

در ادبیات این حوزه تعدادی حمله بر روی مدل‌های مبتنی بر افراز کردن ارائه شده‌اند که در ادامه به طور مختصر توضیح داده می‌شوند.

- **حمله کمینه‌کردن^۱ [۳۱]:** در مدل‌های مختلف مبتنی بر افراز کردن، از الگوریتم‌های بر پایه کمینه‌کردن استفاده می‌شود تا کمترین میزان از دست‌دادن اطلاعات اتفاق بیفتد. در حمله کمینه‌کردن، فرض می‌شود که مهاجم از الگوریتم عمومی کردن ستون‌ها اطلاع دارد. این اطلاع، در کنار جدول منتشر شده و اطلاعات عمومی خارجی موجود درباره اشخاص موجب نقض حریم خصوصی می‌شود.

^۱ minimality attack

- **حمله ترکیب^۱ [۲۳]:** در این حمله از انتشار پایگاه‌داده‌های عمومی شده مختلف که شامل اطلاعات اشخاص مشترک است، سو استفاده می‌شود. برای مثال پایگاه‌داده اول، با مدل ۴-بی‌نامی از بیمارستان الف و پایگاه‌داده دوم، با مدل ۶-بی‌نامی از بیمارستان ب منتشر شود. اگر ما بدانیم آزاده به هر دو بیمارستان مراجعه کرده است، با استفاده از اطلاعات پس‌زمینه‌ای که از او داریم، می‌توانیم کلاس هم‌ارزی مربوط به آزاده در هر دو پایگاه‌داده را شناسایی کرده و سپس اطلاعات حساس موجود در دو کلاس را اشتراک بگیریم. اطلاعات اشتراکی، اعتقاد ما نسبت به بیماری آزاده را بالاتر می‌برد.
- **حمله پیش‌زمینه^۲ [۳۲]:** پایگاه‌داده بی‌نام شده که منتشر می‌شود، خود دارای اطلاعاتی درباره اشخاص است. این الگوهای اطلاعاتی از طریق الگوریتم‌های یادگیری ماشین قابل دست‌یابی است. برای مثال، از اطلاعات بیماری افراد می‌توان به این الگو رسید که ژاپنی‌ها با احتمال کمتری بیماری قلبی دارند. با استفاده از این اطلاعات که اطلاعات پیش‌زمینه گفته می‌شوند، می‌توان در مورد اطلاعات حساس افراد تصمیم‌گیری کرد. تفاوت این حمله با حمله پس‌زمینه در این است که اطلاعات کمکی در اینجا از خود داده‌ها بدست می‌آید و لازم نیست در پایگاه‌داده‌های خارجی آن‌ها را جست‌وجو کرد.

۲-۳ حریم خصوصی مبتنی بر تصادفی کردن

در مدل‌های مبتنی بر تصادفی کردن، در هنگام جمع‌آوری اطلاعات خصوصی افراد یا در هنگام انتشار و تحلیل، اطلاعات به صورت تصادفی تغییر می‌کند. در این قسمت، دو روش پاسخ تصادفی شده^۳ و حریم خصوصی تفاضلی^۴ به عنوان روش‌های حفظ حریم خصوصی مبتنی بر تصادفی کردن توضیح داده خواهند شد.

❖ پاسخ تصادفی شده

در بسیاری از موقعیت‌ها افراد علاقه ندارند که عقیده خود را نسبت به یک موضوع ابراز کنند. برای مثال، سوال «آیا شما به حزب الف علاقمند هستید؟». در چنین موقعیت‌هایی مصاحبه‌کننده با استفاده از پاسخ تصادفی شده می‌تواند امکانی را فراهم کند که فرد بتواند پاسخ خود را به راحتی انکار کند. در مرجع [۳۳]، سازوکار مصاحبه‌کننده به شکل شرح داده شده در ادامه آمده است. در پاسخ به سوال، فرد با احتمال p حقیقت را می‌گوید و با احتمال $(1 - p)$ ، دوباره به صورت تصادفی و با احتمال p می‌گوید «بله» و با احتمال $(1 - p)$

^۱ compositional attack

^۲ foreground attack

^۳ randomized response

^۴ differential privacy

می‌گوید «خیر». با این روش در صورتی که حقیقت پاسخ «بله» باشد، فرد با احتمال $p + (1 - p)p$ می‌گوید «بله» و با احتمال $(1 - p)(1 - p)$ می‌گوید «خیر». با این روش، فرد همیشه می‌تواند پاسخ خود به سوال را انکار کند.

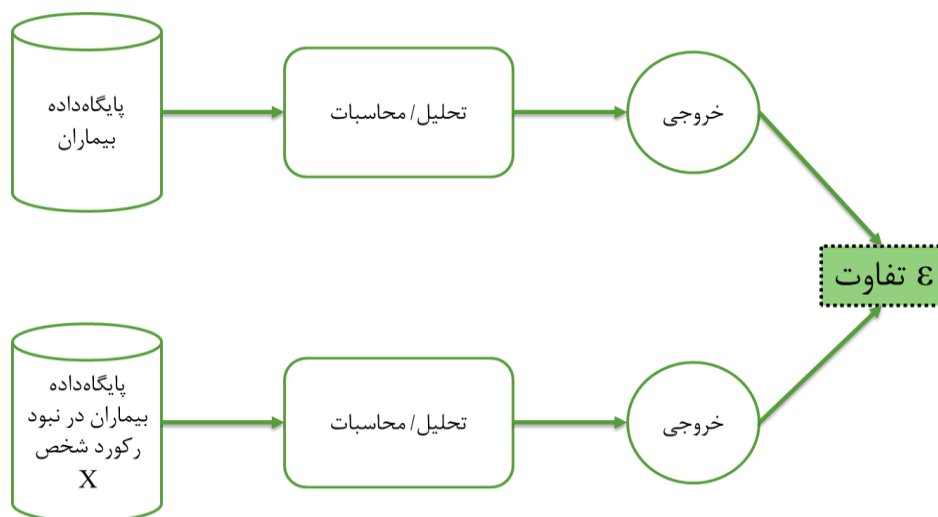
❖ حریم خصوصی تفاضلی

در مدل حریم خصوصی تفاضلی این فرض وجود دارد که اطلاعات خصوصی افراد مانند، نوع بیماری، حقوق ماهیانه، علائم بیماری، و مانند آن در یک پایگاه داده ذخیره شده است. قرار است این اطلاعات برای انجام پژوهش، در اختیار پژوهشگران قرار بگیرد. البته پژوهشگران اجازه دارند، پرس‌وجوهای آماری از قبیل میانگین، میانه، واریانس، و مانند آن را به پایگاه داده ارسال کرده و نتیجه را مشاهده کنند. در اصطلاح به این گونه از پایگاه داده‌ها، پایگاه داده آماری^۱ گفته می‌شود. هدف، ارائه مدلی برای این سناریو است که افراد مختلف بدون نگرانی از افشای اطلاعات خصوصی خود در این پژوهش شرکت کنند. به عبارت دیگر، شرکت یا عدم شرکت در این پژوهش نباید در میزان اطلاعات پژوهشگر و یا مهاجم در مورد یک شخص خاص موثر باشد. برای این هدف حریم خصوصی تفاضلی ارائه شده است که در زیر تعریف آن آمده است.

تعریف [۳۴][۳۵]: سازوکار تصادفی M ، ϵ -خصوصی تفاضلی است اگر برای هر دو پایگاه داده همسایه D_1 و D_2 که در یک رکورد با هم تفاوت دارند (با هم همسایه هستند) و به ازای هر $S \subseteq \text{Range}(M)$ رابطه زیر برقرار باشد:

$$\Pr[M(D_1) \in S] \leq e^\epsilon \Pr[M(D_2) \in S]$$

^۱ statistical database



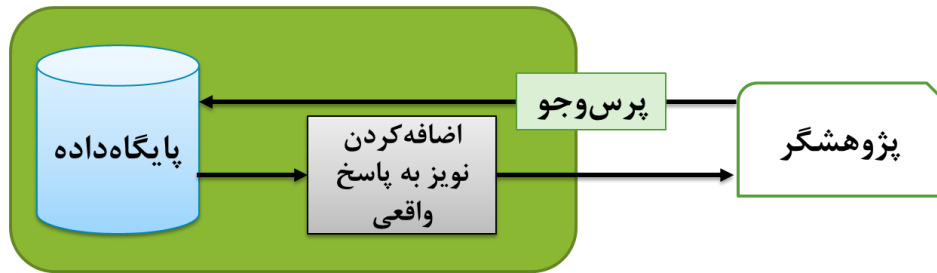
شکل ۲-۲: حریم خصوصی تفاضلی در حضور یا نبود شخص X .

سازوکار M که به ازای آن، رابطه بالا برقرار باشد، نگرانی‌های افراد شرکت‌کننده در پژوهش در مورد افشای اطلاعات خصوصی آنان را بر طرف می‌کند. یعنی شرکت یا عدم شرکت فردی خاص در پژوهش، در اعتقاد مهاجم نسبت به پایگاه‌داده آماری تغییر چندانی ایجاد نمی‌کند (شکل ۲-۲).

در تعریف، قید به ازای تمام پایگاه‌داده‌های همسایه‌ی ممکن آمده است، دلیل این نوشته این است که حریم خصوصی برای تمام افراد حاضر در پایگاه‌داده (تمام رکوردها) تضمین شود. همچنین، حریم خصوصی یک فرد بدون در نظر گرفتن سایر افراد موجود در پایگاه‌داده بررسی شود. در این تعریف، پایگاه‌داده همسایه به صورت وجود یا عدم وجود یک رکورد تعریف شده است، دلیل این موضوع این است که حضور یا نبود یک فرد در پایگاه‌داده بررسی شود. پارامتر ϵ ، تنظیم‌کننده میزان حریم خصوصی افراد است. هرچه ϵ به صفر نزدیک‌تر باشد، حریم خصوصی بیشتر و کامل‌تر عملی خواهد شد. در تعریف، گفته شده به ازای تمام S های ممکن، به این دلیل که تمام خروجی‌های ممکن و محتمل برای پرس‌وجوی پژوهشگر در نظر گرفته شود.

در ادامه دو سازوکار مختلف برای تصادفی‌کردن معرفی خواهد شد. آنگاه، قضیه‌های ترکیب و پساپردازش برای حفظ حریم خصوصی پرس‌وجوهای پیچیده توضیح داده می‌شوند. در ادامه این قسمت، تعابیر مختلف از همسایگی و اتلاف حریم خصوصی^۱ توضیح داده می‌شوند. همچنین، روش‌های مختلف صوری‌سازی برای حریم خصوصی تفاضلی ارائه می‌شوند. در پایان هم حریم خصوصی تفاضلی محلی بررسی خواهد شد.

^۱ privacy loss



شکل ۲-۳: سازوکار لاپلاس

(۱) سازوکار تصادفی کردن^۱

برای اینکه رابطه حریم خصوصی تفاضلی معرفی شده در تعریف برقرار باشد، باید سازوکار M تعریف شود. از میان سازوکارهای مختلف معرفی شده، در اینجا دو سازوکار لاپلاس^۲ و نمایی^۳ معرفی می‌شوند.

- سازوکار لاپلاس^۴

سازوکار لاپلاس برای پرس‌و‌جوهای مناسب است که خروجی آن‌ها یک عدد حقیقی است، برای مثال «تعداد افرادی که بیماری دیابت دارند چند نفر است؟». اگر f تابع پرس‌و‌جو و X پایگاه‌داده باشد، مقداری نویز متناسب با f به خروجی $f(X)$ به صورت تصادفی اضافه می‌شود. به عبارت دیگر، $f(X) + \eta$ به عنوان خروجی ارسال می‌شود که $\eta \sim \text{Lap}\left(\frac{c}{\epsilon}\right)$ به دست می‌آید. $\text{Lap}(\sigma)$ توزیع احتمال لاپلاس با انحراف معیار σ است. c میزان حساسیت^۵ تابع f است (شکل ۲-۳).

تعریف [۳۴][۳۵]: برای تمام پایگاه‌داده‌های همسایه D_1 و D_2 ، که در یک رکورد با هم تفاوت دارند،

حساسیت $f: \mathcal{D} \rightarrow \mathbb{R}^d$ به صورت زیر تعریف می‌شود:

$$c = \max_{D_1, D_2} \|f(D_1) - f(D_2)\|_1$$

قضیه [۳۴][۳۵]: برای $f: \mathcal{D} \rightarrow \mathbb{R}^d$ ، سازوکار لاپلاس، ϵ -خصوصی تفاضلی است.

- سازوکار نمایی^۶

^۱ randomization mechanism^۲ laplace^۳ exponential^۴ laplace mechanism^۵ sensitivity^۶ exponential mechanism

سازوکار نمایی، سازوکاری عمومی است که سازوکار لاپلاس را نیز شامل می‌شود و برای حفظ حریم خصوصی تفاضلی در پاسخ‌گویی به پرس‌وجوهای مختلف (به طور خاص پرس‌وجوهای با خروجی غیر عدد حقیقی) استفاده می‌شود. برای مثال، در پاسخ به پرس‌وجوی «کدام ملیت در دانشگاه بیشترین فراوانی را دارد؟» می‌توان از این سازوکار استفاده کرد. اگر دامنه پاسخ‌ها، R و تابع سودمندی^۱، $u: \mathbb{R}^d \times R \rightarrow \mathbb{R}$ (با ثابت بودن پایگاه داده، کاربر علاقه دارد، عضوی از دامنه با خروجی بیشینه تابع سودمندی را دریافت کند). را در نظر بگیریم و حساسیت تابع سودمندی u را به صورت زیر با ثابت بودن دامنه R تعریف کنیم:

$$\Delta u \equiv \max_{r \in R} \max_{D_1, D_2} |u(D_1, r) - u(D_2, r)|$$

تعریف زیر را برای سازوکار نمایی خواهیم داشت.

تعریف [۳۶]: سازوکار نمایی $M_E(D, u, R)$ ، عضو $r \in R$ را با احتمال متناسب با $\exp\left(\frac{\epsilon u(D, r)}{2\Delta u}\right)$ انتخاب

کرده و به عنوان خروجی ارسال می‌کند.

قضیه [۳۶]: سازوکار نمایی M_E ، ϵ -خصوصی تفاضلی است.

(۲) ترکیب^۲

تا اینجا سازوکارهای مختلفی برای حریم خصوصی تفاضلی معرفی شدند. در این قسمت، دو قضیه مطرح می‌شوند که به وسیله آن‌ها می‌توان سازوکارهایی برای حفظ حریم خصوصی برای پرس‌وجوهای پیچیده‌تر را ساخت.

قضیه ترکیب ترتیبی^۳ [۳۶]: اگر سازوکارهای M_1 تا M_n با پارامترهای ϵ_1 تا ϵ_n به صورت تفاضلی خصوصی باشند، آنگاه اجرای ترتیبی آن‌ها به اندازه $\epsilon = \sum_{i=1}^n \epsilon_i$ به صورت تفاضلی خصوصی است.

قضیه ترکیب موازی^۴ [۳۶]: اگر سازوکارهای M_1 تا M_n با پارامترهای ϵ_1 تا ϵ_n به صورت تفاضلی خصوصی باشند، آنگاه اجرای موازی آن‌ها روی افرازهای مشخصی از پایگاه داده، به اندازه $\epsilon = \max_{i=1 \dots n} \epsilon_i$ به صورت تفاضلی خصوصی است.

(۳) پس‌پردازش^۵

^۱ utility function

^۲ composition

^۳ sequential composition

^۴ parallel composition

^۵ post-processing

همواره خروجی سازوکارهای تصادفی کردن، پردازش می‌شوند و سپس، نتیجه به پژوهشگر داده می‌شود. قضیه پس‌پردازش نشان می‌دهد که پس‌پردازش روی یک سازوکار به صورت تفاضلی خصوصی، هم حریم خصوصی تفاضلی را حفظ می‌کند.

قضیه پس‌پردازش [۳۶]: اگر M یک سازوکار ϵ -خصوصی تفاضلی و D پایگاه‌داده مورد نظر برای تحلیل باشد، همچنین، $f: R \rightarrow R'$ یک نگاشت تصادفی دلخواه باشد، آنگاه، $f(M(D))$ هم یک سازوکار ϵ -خصوصی تفاضلی است.

با استفاده از این دو قضیه ترکیب و پس‌پردازش، می‌توان بلاک‌هایی حافظ حریم خصوصی تفاضلی تعریف کرد. آنگاه می‌توان، آن‌ها را به صورت موازی و یا ترتیبی ترکیب کرد و بر روی خروجی آن‌ها پس‌پردازش نیز داشت تا در مجموع الگوریتمی پیچیده یا برنامه‌ای بزرگ و حافظ حریم خصوصی تفاضلی [۳۸] ایجاد کرد.

(۴) همسایگی

در تعریف و مقاله اصلی حریم خصوصی تفاضلی [۳۵]، پایگاه‌داده به شکل آرایه‌ای n در m در نظر گرفته می‌شود. آنگاه، همسایگی دو پایگاه‌داده D_1 و D_2 ، فاصله همینگ^۱ میان آن‌ها تعریف می‌شود. با وجود این، در مسائل و کاربردهای مختلف نیاز است که با توجه به ساختار پایگاه‌داده و نوع رکوردهای آن، تعبیر جدیدی از همسایگی نیز ارائه شود. کارهای مختلفی در مورد داده‌های با ساختار گراف، داده‌های جریانی^۲، عملیات روی مجموعه‌ها، تصاویر، داده‌های ژن افراد، اطلاعات مکانی، و بازیابی اطلاعات شخصی‌سازی شده انجام شده‌اند. در فصل ۳، به طور کامل به بازیابی اطلاعات شخصی‌سازی شده خواهیم پرداخت.

برای نمونه، در شبکه‌های اجتماعی رابطه دوستی میان افراد در قالب گراف نمایش داده می‌شود که گره‌ها افراد و یال‌ها رابطه دوستی هستند. در اینجا، پایگاه‌داده‌های همسایه را برای مثال می‌توان در حضور یا عدم حضور یک فرد (گره) و یا وجود یا عدم وجود رابطه دوستی (یال) تعریف کرد. در پایگاه‌داده‌ای که مسیرهای روزانه یک فرد در آن نگهداری می‌شود، دو پایگاه‌داده با هم همسایه هستند اگر، در یک مسیر با هم تفاوت داشته باشند. تفاوت این دو مسیر می‌تواند در یک نقطه مکانی باشد یا در یک بازه زمانی، مکان‌های متفاوتی داشته باشند یا به طور کامل دو مسیر با هم متفاوت باشند.

(۵) تنوع در اتلاف حریم خصوصی

^۱ hamming distance

^۲ streaming data

میزان حساسیت افراد مختلف نسبت به رکوردهای خود در پایگاه‌داده می‌تواند متفاوت باشد [۲۲]. با وجود این، در تعریف اصلی حریم خصوصی تفاضلی، حساسیت همه افراد به یک اندازه فرض شده است و همگی به اندازه ϵ اجازه می‌دهند که اطلاعات آنان فاش شود. برای اینکه متغیر بودن حساسیت افراد مختلف نسبت به اطلاعات آنان را در تعریف حریم خصوصی تفاضلی بیاوریم، به جای ϵ از تابع $\Psi: d_i \rightarrow \mathbb{R}$ استفاده می‌کنیم که یک رکورد از پایگاه‌داده می‌گیرد و مقدار ϵ متناسب با آن را بر می‌گرداند. فرض شده است که هر رکورد مربوط به یک شخص خاص است. همچنین، با توجه به مفهوم همسایگی، مفهوم رکورد هم تغییر می‌کند. بنابراین، تعریف حریم خصوصی تفاضلی به شکل زیر در می‌آید. علاوه بر آن، یک آستانه

$$Pr[M(D_1) \in S] \leq e^{\Psi(d_i)} Pr[M(D_2) \in S]$$

۶ روش‌های صوری‌سازی^۱

با توجه به تعاریف مختلف ارائه شده برای حریم خصوصی تفاضلی می‌توان آن را به شکل‌های مختلف صوری کرد [۲۲]:

- $Pr[M(D_1) \in S] \leq e^\epsilon Pr[M(D_2) \in S]$

به ازای هر دو پایگاه‌داده همسایه، D_1 و D_2 که $S \subseteq O$ که مجموعه تمام خروجی‌های سازوکار تصادفی M است. M سازوکار ما برای حفظ حریم خصوصی است. برای مثال، اضافه کردن نویز از توزیع لاپلاس به خروجی تابع شمارش. برای اولین بار در مرجع [۳۴] و [۳۵] حریم خصوصی تفاضلی به شکل بالا صوری شده است.

- $\frac{Pr[t \in D | M(D) = o]}{Pr[t \notin D | M(D) = o]} \leq e^{\epsilon \frac{Pr[t \in D]}{Pr[t \notin D]}}$

در این صوری‌سازی، توانایی مهاجم به صورت بیزین در نظر گرفته می‌شود. در این تعریف، اعتقاد مهاجم قبل و بعد از اجرای سازوکار M بر روی پایگاه‌داده در حضور و عدم حضور رکورد t ، مقایسه می‌شود. گفته می‌شود که اعتقاد مهاجم قبل از اجرای M و بعد از آن در حضور و نبود رکورد t در نهایت باید به اندازه ϵ تفاوت داشته باشد.

- $SD(M(D), M(D_{-i})) \leq \epsilon$
 $SD(X, Y) = \max_{S \subseteq D} |Pr[X \in S] - Pr[Y \in S]|$

^۱ formalism method

منظور از SD در تعریف بالا، تابع تفاضل آماری^۱ است. در این تعریف، بودن یا نبودن رکورد با نشانه i ، نباید بیشتر از ϵ روی تفاضل دو توزیع احتمال موخر^۲، موثر باشد. در اینجا، بر خلاف تعریف‌های قبلی، تفاضل توزیع‌های احتمالی در نظر گرفته می‌شود که بازتاب دهنده اعتقاد مهاجم نسبت به پایگاه‌داده هستند.

(۷) حریم خصوصی تفاضلی محلی^۳

در مدل حریم خصوصی تفاضلی، مسئول پایگاه‌داده مورد اعتماد است. با وجود این، به دلایل مختلف این موجودیت می‌تواند مورد اعتماد نباشد. تعدادی از دلایل عدم اعتماد به مسئول پایگاه‌داده می‌تواند شامل فروش اطلاعات کاربران، تغییر مدیریت و عوض شدن سیاست‌های سازمان، اجبار مقام قضایی برای افشای اطلاعات و مانند آن باشد. توضیحات این قسمت با مطالعه مراجع [۳۹]–[۴۶] به دست آمده است.

در LDP، به جای اینکه به پاسخ پرس‌وجوهای پایگاه‌داده نویز اضافه شود، مالکان به داده‌های خود نویز اضافه می‌کنند و در عمل، داده‌های مغشوش در پایگاه‌داده ذخیره می‌شوند. برای بیان صوری، فرض می‌شود اطلاعات کاربر، یک پایگاه‌داده با اندازه یک است. به عبارت دیگر، v و v' دو پایگاه‌داده با اندازه یک و همسایه هستند. همچنین، v^* داده مغشوشی است که کاربر به مسئول پایگاه‌داده ارسال می‌کند. علاوه بر آن، \mathcal{A} سازوکار تصادفی‌کردنی است که کاربر روی داده‌های خود در جهت مغشوش کردن اجرا می‌کند. به ازای هر مقدار برای v ، v' و $v^* \in \text{Range}(\mathcal{A})$ رابطه زیر برقرار است:

$$\Pr[\mathcal{A}(v) = v^*] \leq e^\epsilon \Pr[\mathcal{A}(v') = v^*]$$

گفته می‌شود که \mathcal{A} ، ϵ -خصوصی تفاضلی محلی است. به عبارت دیگر، مسئول پایگاه‌داده با دیدن v^* نمی‌تواند با قطعیت بگوید که داده کاربر v یا v' است. از سازوکارهای تصادفی‌کردن \mathcal{A} می‌توان به سازوکار لاپلاس یا استفاده از پاسخ تصادفی شده^۴ اشاره کرد. سازوکار \mathcal{A} ، متناسب با الگوریتم تحلیل‌گر به داده‌ها نویز اضافه می‌کند. این الگوریتم می‌تواند پردازش‌های مختلفی از جمله محاسبه میانگین، شمارش، و الگوریتم‌های پیچیده مانند به دست آوردن مقدار پارامتر یک توزیع احتمال یا الگوریتم‌های یادگیری ماشین باشد. نویز اضافه شده باید توازن درستی میان حریم خصوصی و سودمندی^۵ ایجاد کند.

^۱ statistical difference

^۲ posterior

^۳ local differential privacy (LDP)

^۴ randomized response (RR)

^۵ utility

برای مثال، در صورتی که از سازوکار لاپلاس استفاده شود، لازم است تا حساسیت الگوریتم محاسبه شود. در حالت محافظه‌کارانه، هر کاربر می‌تواند از حساسیت عمومی استفاده کند. استفاده از حساسیت عمومی موجب اضافه شدن بیشترین نویز به داده‌ها می‌شود. اگر پایگاه‌داده‌های D_1 و D_2 همسایه و دلخواه باشند، حساسیت عمومی از رابطه زیر به دست می‌آید.

$$GS = \max_{D_1, D_2} \|f(D_1) - f(D_2)\|_1$$

در مقابل، هر کاربر می‌تواند متناسب با داده خود حساسیت محلی الگوریتم را محاسبه کرده و سپس با توجه به آن داده خود را مغشوش کند. اگر پایگاه‌داده D_1 شامل داده کاربر باشد و پایگاه‌داده D_2 همسایه با D_1 و دلخواه باشد، حساسیت محلی به صورت زیر محاسبه می‌شود.

$$LS = \max_{D_2} \|f(D_1) - f(D_2)\|_1$$

گفتنی است، نویز اضافه شده متناسب با حساسیت عمومی حد بالای نویز اضافه شده است و استفاده از حساسیت محلی نویز کمتری به داده‌ها اضافه می‌کند.

۸) حریم خصوصی تفاضلی در پایگاه‌داده دارای همبستگی^۱

در مرجع [۴۷]، طرحی برای پایگاه‌داده دارای همبستگی ارائه شده است. گفتنی است، طرح ارائه شده، برای حریم خصوصی تفاضلی محلی نیست. در پایگاه‌داده D ، اگر رکورد r_i با $k-1$ رکورد دیگر دارای همبستگی باشد، به مجموعه این k رکورد رکوردهای همبسته^۲ گفته می‌شود. این مجموعه با CoR_{r_i} نمایش داده می‌شود. $\{r_i, r_j \in D | \text{all } r_j \text{ are correlated to } r_i\}$ آنگاه به D پایگاه‌داده همبسته گفته می‌شود. اگر $|k| = 1$ باشد، پایگاه‌داده دارای رکوردهای مستقل از هم است. حذف هر کدام از رکوردهای همبسته ممکن است اثر متفاوتی بر دیگر رکوردها داشته باشد. بنابراین، مفهوم درجه همبستگی^۳ تعریف می‌شود. اگر r_i و r_j با هم همبسته باشند، درجه همبستگی آن‌ها با $\delta_{ij} \in [-1, 1]$ و $|\delta_{ij}| \geq \delta_0$ نمایش داده می‌شود که δ_0 آستانه همبستگی است. بنابراین، $|\delta_{ij}| = 0$ به معنی استقلال و $|\delta_{ij}| = 1$ به معنای همبستگی کامل دو رکورد r_i و r_j است. در نتیجه، همبستگی رکوردها را می‌توان در قالب ماتریس Δ نمایش داد. ویژگی‌های این ماتریس عبارتند از: (۱) قطر ماتریس همگی یک هستند. (۲) ماتریس نسبت به قطر متقارن است. به عبارت دیگر، $\delta_{ij} = \delta_{ji}$. (۳) اگر $|\delta_{ij}| < \delta_0$ خانه متناظر با آن صفر می‌شود. (۴) تنها تعدادی از رکوردها با دیگران همبستگی دارند.

^۱ correlation

^۲ correlated records

^۳ correlation degree

$$\Delta = \begin{pmatrix} \delta_{11} & \cdots & \delta_{1n} \\ \vdots & \ddots & \vdots \\ \delta_{n1} & \cdots & \delta_{nn} \end{pmatrix}$$

با توجه به نوع اطلاعات ذخیره شده در پایگاه داده، باید رکوردهای همبسته را شناسایی کرده و ماتریس Δ را ایجاد کرد. با استفاده از ماتریس Δ ، مفهوم حساسیت همبسته^۱ تعریف می‌شود. برای این تعریف، نیاز است حساسیت رکورد^۲ تعریف شود. برای ماتریس Δ و پرس‌وجوی f حساسیت رکورد r_i به شکل زیر تعریف می‌شود:

$$CS_i = \sum_{j=0}^n |\delta_{ij}| (\|f(D_j) - f(D_{-j})\|_1)$$

حساسیت رکورد r_i ، میزان تاثیر حذف رکورد r_i بر همه رکوردهای موجود در پایگاه داده را اندازه‌گیری می‌کند. اگر پایگاه داده همبستگی نداشته باشد، حساسیت رکورد r_i ، مثل تعریف اصلی حساسیت خواهد شد. بر این اساس حساسیت همبسته برای پرس‌وجوی f به شکل زیر تعریف می‌شود:

$$CS_q = \max_{i \in q} (CS_i)$$

q مجموعه رکوردهایی است که در پاسخ به پرس‌وجوی f وجود دارند. حساسیت همبسته مرتبط با تابع پرس‌وجو است. از میان رکوردهای موجود در پاسخ به پرس‌وجو، حساسیت رکوردی که بیشترین حساسیت همبسته دارد، به عنوان حساسیت همبسته پرس‌وجو انتخاب می‌شود. در حالت ساده‌انگارانه، حساسیت پرس‌وجو از ضرب تعداد رکوردهای همبسته با هم در حساسیت پرس‌وجو در حالت نبود همبستگی به دست می‌آید. در مرجع [۴۷]، اثبات شده است که حساسیت همبسته از حساسیت حالت ساده‌انگارانه، بهینه و کمتر است. پس از تعریف حساسیت همبسته، می‌توان سازوکار لاپلاس را به شکل

$$\hat{f}(D) = f(D) + \text{Lap}\left(\frac{CS_q}{\epsilon}\right)$$

نوشت. بنابراین، رابطه حریم خصوصی تفاضلی به ازای هر دو پایگاه داده همسایه D و D' و به ازای تمام

مقادیر $S \subseteq \text{Renge}(M)$ به شکل

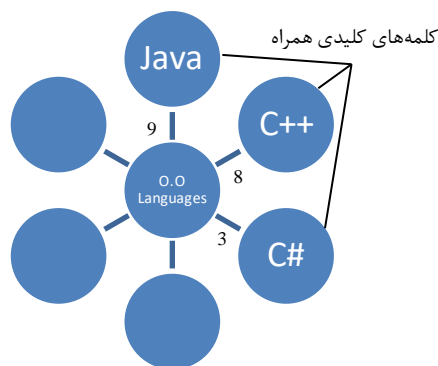
$$\Pr[M(D) = S] \leq e^\epsilon \Pr[M(D') = S]$$

تعریف می‌شود.

^۱ correlated sensitivity

^۲ record sensitivity

۴-۲ روش‌های ایجاد پروفایل کاربر



شکل ۴-۲: نمونه پروفایل شبکه نحوی

همان طور که گفته شد، دو روش گسترش پرس‌وجوی کاربر و مرتب‌سازی نتایج جستجو برای شخصی‌سازی نتایج جستجو استفاده می‌شود. در هر دوی این روش‌ها از پروفایل کاربر برای شخصی‌سازی استفاده می‌شود. در این بخش، روش‌های مختلف ایجاد و نمایش پروفایل کاربر توضیح داده می‌شود. مطالب این

بخش از مراجع [۴۸] و [۴۹]، آمده‌اند. روش‌های نمایش پروفایل کاربر شامل، پروفایل کلمه کلیدی^۱، پروفایل شبکه نحوی^۲ و پروفایل مفهومی^۳ است. در پروفایل کلمه کلیدی، هر کدام از کلمه‌های کلیدی یا گروهی از آن‌ها بازتاب‌دهنده پسندهای کاربر هستند. کلمه‌های کلیدی از سندهای مراجعه شده کاربر استخراج می‌شوند. به هر کدام از کلمه‌های کلیدی مقداری به عنوان وزن آن کلمه بر اساس آمارهای مانند TF-IDF اختصاص پیدا می‌کند. وزن هر کلمه نمایانگر ارزش آن کلمه در پروفایل کاربر است.

اشکالی که روش پروفایل کلمه‌های کلیدی دارد، عدم تمایز میان کلمه‌ها با معانی مختلف است. برای مقابله با این اشکال پروفایل شبکه‌های نحوی ارائه شده است. این پروفایل، با شبکه نحوی وزن‌دار نمایش داده می‌شود که هر گره از آن یک مفهوم خاص را بازتاب می‌دهد. به هر گره، تعدادی کلمه کلیدی مرتبط با آن مفهوم وصل می‌شوند. این کلمه‌های کلیدی کله‌هایی هستند که در سندهای مراجعه شده کاربر همراه هم ظاهر می‌شوند. یال‌های موجود در شبکه، نشان‌دهنده ارتباط کلمه‌های کلیدی با مفاهیم هستند. روی هر یال مقداری به عنوان وزن قرار می‌گیرد. وزن هر یال بازتاب‌دهنده شدت علاقه کاربر به آن مفهوم است.

پروفایل مفهومی مانند پروفایل شبکه نحوی، ساختار گراف دارد. با این تفاوت که به جای کلمه‌های کلیدی، مفاهیم انتزاعی مورد پسند کاربر در گره‌ها قرار می‌گیرند. همچنین، پروفایل مفهومی شبیه به پروفایل کلمه‌های کلیدی است. از این جهت که هر گره (مفهوم)، همراه با یک مقدار به عنوان وزن آن مفهوم است. در این روش روی یال‌ها وزنی وجود ندارد. برای ایجاد هر کدام از پروفایل‌های مختلف گفته شده در بالا، روش‌های گوناگون مختلفی از جمله فنون یادگیری ماشین (با سرپرست^۴ یا بدون سرپرست^۱) مطرح شده‌اند. در این بخش به عنوان

^۱ keyword profile

^۲ semantic network profile

^۳ concept profile

^۴ supervised

نمونه روش ایجاد پروفایل کلمه‌های کلیدی بر اساس آماره TF-IDF و الگوریتم k-means شرح داده می‌شود. رابطه TF-IDF در زیر آمده است. ارزش هر کلمه t در سند d با $w_{t,d}$ نمایش داده می‌شود.

$$w_{t,d} = tf_{t,d} \log \left(\frac{N}{df_t} \right)$$

در رابطه بالا، $tf_{t,d}$ بیان‌گر تعداد دفعاتی است که کلمه t در سند d تکرار شده است. همچنین، df_t نماینده تعداد سندهایی است که کلمه t در آن‌ها آمده است و N نشان دهنده تعداد کل سندها است. بنابراین، $\log \left(\frac{N}{df_t} \right)$ ارزش کلمه t ، را صوری می‌کند. به عبارت دیگر، هرچه یک کلمه نادرتر باشد، ارزشمندتر است. در مرجع [۴۹]، موتور جستجو از روش مرتب‌سازی نتایج پرس‌وجو برای شخصی‌سازی استفاده می‌کند.

بر روی پرس‌وجوی جستجو و سندهای کلیک شده کاربر، پیش‌پردازش‌های حذف کلمه‌های ایست و استخراج ریشه کلمه‌ها انجام می‌شود. از پرس‌وجوی جستجو و سندهای کلیک شده، کلمه‌ها به همراه تعداد تکرار آن‌ها استخراج می‌شوند. تکرار کلمه‌ها در IDF هر کلمه ضرب می‌شود تا بردار وزن کلمه‌ها در پرس‌وجوی جستجو و سندها به دست بیاید. با استفاده از الگوریتم k-means، هر کدام از این بردارها در یکی از k خوشه^۲ با توجه به کمترین فاصله کسینوسی هر بردار و مرکز خوشه، قرار می‌گیرد. آنگاه، بردار مرکز هر خوشه دوباره محاسبه و به‌روز می‌شود. بردار مرکز خوشه، نماینده تمام بردارهای موجود در آن خوشه است و از میانگین مقادیر بردارهای آن خوشه ساخته می‌شود. k حداکثر تعداد پسندهایی است که موتور جستجو در پروفایل کاربر به شکل برداری نگهداری می‌کند. بنابراین، خروجی فرایند ایجاد پروفایل کاربر بردار مرکز خوشه‌ها است. برای شخصی‌سازی نتیجه جستجوی کاربر، بردار وزن کلمه‌های سندهای نتیجه از جستجوی کاربر، نیز به دست می‌آیند. آنگاه، با استفاده از محاسبه فاصله کسینوسی بردار سندهای نتیجه و بردارهای مرکز خوشه‌ها، سندهای نتیجه جستجوی کاربر بر اساس کمترین فاصله مجدداً مرتب می‌شوند.

۵-۲ آنتولوژی محاسباتی^۳

همان طور که در مرجع [۵۰] آمده است، Ontology یک اسم غیرقابل شمارش است و منظور هستی‌شناسی به عنوان شاخه‌ای در فلسفه است که به «طبیعت وجود» می‌پردازد. در مقابل آنتولوژی

^۱ unsupervised

^۲ cluster

^۳ ontology

محاسباتی که یک اسم قابل شمارش است، در لغت‌نامه در حوزه معنایی «رایانش» به معنای «مجموعه‌ای از مفاهیم در یک موضوع یا دامنه خاص و ارتباط بین آن‌ها» آمده است. همان طور که در مرجع [۵۱] آمده است، رده‌بندی^۱ (یا ردگان) تعریف‌کننده کلاس‌های اشیا و ارتباط بین آن‌ها است. کلاس‌ها، زیرکلاس‌ها، و ارتباط بین موجودیت‌ها ابزارهای مهم در رده‌بندی هستند. آنتولوژی محاسباتی از رده‌بندی به همراه مجموعه‌ای از قواعد استنتاج استفاده می‌کند. بنابراین، با استفاده از قواعد استنتاج موجود در آنتولوژی محاسباتی می‌توان دانش و اطلاعات جدید به دست آورد. برای مثال، در یک آنتولوژی محاسباتی این قاعده گفته شده است که «اگر فردی جمهوری‌خواه باشد، آن فرد جنگ‌طلب است». با این قاعده می‌توان استنتاج کرد که «جورج بوش جمهوری‌خواه است، بنابراین، جورج بوش جنگ‌طلب است». همچنین، در مرجع [۵۱] آمده است که پایگاه‌دانش^۲ شامل آنتولوژی محاسباتی و نمونه‌های مفاهیم موجود در آن است. در مقابل، در یک پایگاه‌داده، داده‌ها در قالب یک مدل داده^۴ مانند مدل داده رابطه‌ای^۵ ذخیره می‌شوند. برای مثال، داده‌های مربوط به مقاله‌های یک ناشر (شامل عنوان مقاله، نویسندگان، سال انتشار و غیره) در پایگاه‌داده ذخیره می‌شوند. با وجود این، در پایگاه‌دانش آنتولوژی محاسباتی مربوط به موضوع‌های پژوهشی در علوم کامپیوتر به همراه مقاله‌های مرتبط با هر موضوع وجود دارد.

همان طور که در مرجع [۵۱] آمده است، برای نشان دادن آنتولوژی‌های محاسباتی از زبان‌های آنتولوژی محاسباتی مانند OWL استفاده می‌شود. گفتنی است، کتابخانه‌هایی برای نگهداری آنتولوژی‌های محاسباتی از پیش تعریف شده وجود دارد. تعدادی از این کتابخانه‌ها در زیر آمده‌اند:

- [ONKI^۶](https://onki.fi/en/browser/)
- [oeGov^۷](http://oegov.us/)
- [OLS^۸](https://www.ebi.ac.uk/ols/index)

^۱ taxonomy

^۲ knowledge base

^۳ instance

^۴ data model

^۵ relational

^۶ Link to: <https://onki.fi/en/browser/>

^۷ Link to: <http://oegov.us/>

^۸ Link to: <https://www.ebi.ac.uk/ols/index>

فصل سوم

کارهای پژوهشی پیشین

۳-۱ کارهای پیشین حفظ حریم خصوصی در جستجوی وب شخصی سازی شده

در این قسمت، پژوهش های مرتبط با حریم خصوصی در جستجوی وب شخصی سازی شده بررسی می شوند. این پژوهش ها در سه دسته مدل های مبتنی بر افراز کردن، مدل های مبتنی بر تصادفی کردن، و سایر روش ها و مقاله ها دسته بندی شده اند. در پایان، خلاصه این پژوهش ها در **Error! Reference source not found.** آورده شده اند.

❖ مدل های مبتنی بر افراز کردن

در مرجع [۱۲]، در سامانه سمت کلاینت برنامه ای اجرا می شود. در این برنامه، همراه با هر پرس و جوی کاربر، k پرس و جوی دیگر نیز ارسال می شود. به این پرس و جوها، پرس و جوهای پوششی گفته می شود. هدف از ارسال آن ها، پنهان کردن پرس و جوی اصلی کاربر است. پرس و جوهای پوششی از موضوع های مختلفی ایجاد می شوند، که آنتروپی آن موضوع ها، در بازه ϵ از موضوع پرس و جوی اصلی است. گفتنی است، پروفایلی از کاربر که بازتاب دهنده پسندهای دقیق او است، در برنامه موجود در کلاینت نیز ایجاد و ذخیره می شود. پاسخ به پرس و جوی ارسال شده به موتور جستجو، بر اساس پروفایل موجود در کلاینت نیز مجددا مرتب می شود.

در مرجع [۵]، از مدل حریم خصوصی شبیه به k -بی نامی استفاده می شود. در طرح ارائه شده، حداقل $k - 1$ پرس و جو همراه با پرس و جوی اصلی کاربر ارسال می شود. پرس و جوهای پوششی بر اساس فاصله مفهومی از پرس و جوی اصلی انتخاب می شوند. هر چه این فاصله بیشتر باشد، حریم خصوصی کاربر نیز بیشتر حفظ می شود. ارسال پرس و جوهای پوششی به معنای ارسال پرس و جوهایی بی ارتباط با پرس و جوی اصلی است. هر چه فاصله مفهومی کمتر باشد، موتور جستجو اطلاعات بیشتری درباره مفهوم مورد نظر کاربر به دست می آورد. برای مثال، اگر «ورزش آبی» موضوع پرس و جوی کاربر باشد، موضوع های «کوهنوردی» و «شنا و شیرجه» می توانند موضوع های پرس و جوهای پوششی باشند. پرس و جوی پوششی درباره موضوع «کوهنوردی»، مفهوم مورد نظر کاربر را بیشتر پنهان می کند. هر چه k بیشتر باشد، حریم خصوصی نیز بیشتر حفظ می شود.

برای ایجاد پرس و جوهای پوششی، پرس و جوی اصلی از نظر نحوی و ساختار جمله بررسی می شود و قسمت های اصلی آن استخراج می شود. سپس، با محاسبه آنتروپی قسمت ها، قسمتی که اطلاعات بیشتری (در تئوری اطلاعات و مفهوم آنتروپی) دارد، به این معنی که احتمال وقوع آن کمتر است، به عنوان مفهوم اصلی

پرسوجوی کاربر انتخاب می‌شود. سپس، با استفاده از هستی‌شناسی^۱ موضوعها و مقدار پارامتر c ، که فاصله مفهومی را تعیین می‌کند، پرسوجوهای پوششی تولید می‌شوند. هر چه c بیشتر باشد، حریم خصوصی بیشتر حفظ می‌شود.

در مرجع [۱۵]، از مدل حریم خصوصی l -تنوع آنتروپی استفاده شده است. در این مرجع، همراه با هر کلمه کلیدی، $k - 1$ کلمه کلیدی دیگر نیز ارسال می‌شود. به عبارت دقیق‌تر، کلمه‌های کلیدی با هم OR می‌شوند تا موتور جستجو نتواند با اطمینان مشخص کند که کلمه کلیدی اصلی کاربر کدام بوده است. کلمه‌ها از یک انبار کلمه‌های کلیدی استخراج می‌شوند. این کلمه‌ها، شامل کلمه‌های رایج در زبان انگلیسی هستند. به طور مثال، برای بدست آوردن کلمه‌ها و تعداد تکرار آن‌ها، می‌توان از بررسی کردن متن اخبار NBC استفاده کرد. این $k - 1$ کلمه، باید طوری انتخاب شود که $H(Q_0) \geq H(k)$ باشد. تابع H ، میزان آنتروپی متغیر تصادفی ورودی را محاسبه می‌کند. Q_0 ، متغیر تصادفی مربوط به پرسوجوهای مختلف است. به عبارت دیگر، در اینجا از مدل l -تنوع آنتروپی استفاده می‌شود. حداقل l پرسوجو با هم ارسال می‌شوند که رابطه $H(Q_0) \geq H(k)$ برای آن‌ها برقرار است. اگر توزیع احتمال پرسوجوها یکنواخت باشد، مقدار آنتروپی، بیشینه و برابر با $\log k$ است. اگر موتور جستجو به طور قطع بداند، پرسوجو مربوط به کاربر چیست، مقدار آنتروپی برابر با صفر می‌شود. طبق پروتکلی که در مقاله ارائه شده است، تضمین می‌شود که کلمه‌ها طوری اضافه می‌شوند که رابطه بالا برقرار باشد.

مرجع‌های [۸]، [۹]، و [۱۰] مربوط به یک رساله دکتری هستند و همگی بر اساس فرضها و محدودیت‌های زیر نوشته شده‌اند:

- حریم خصوصی به عنوان قابلیت عدم تمایز در پسندهای کاربر در نظر گرفته می‌شود، نه عدم شناسایی هویت کاربر.
- تنها در مورد موتورهای جستجو بحث شده است. سامانه‌های توصیه‌گر^۲ دیگر مثل Amazon هم می‌توانند مورد بحث قرار گیرند.
- در نظر گرفتن موتور جستجو به صورت جعبه سیاه. الگوریتم و نحوه عملکرد موتور جستجو در نظر گرفته نمی‌شود. فقط ارسال پرسوجو و دریافت پاسخ آن مورد نظر است.

^۱ ontology

^۲ recommender system

- تنها یادگیری پروفایل کاربر در موتور جستجو بر اساس تنوع تبلیغ‌های ارائه شده از آن در نظر گرفته می‌شود. رتبه‌بندی پیوندها در پاسخ به جستجو کاربر در نظر گرفته نمی‌شود. همچنین، جستجوهای مرتبط با اخبار یا آب و هوا نیز در نظر گرفته نمی‌شوند. محتوا و مفهوم پرس‌وجوها و الگوی کلیک شدن پیوندها نیز در نظر گرفته نمی‌شود.
- متن پرس‌وجو در یادگیری موتور جستجو بسیار مهم است. روش‌هایی برای دور زدن این یادگیری باید ارائه شود.
- ورودی‌های متنی در نظر گرفته می‌شود. ورودی‌ها می‌توانند، متنوع باشند. برای مثال عکس، موقعیت مکانی، و مانند آن.
- حفظ حریم خصوصی در جلسه‌های کوتاه‌مدت (۲۰ دقیقه تعامل) در نظر گرفته شده است. یادگیری پروفایل بلندمدت در موتور جستجو در نظر گرفته نشده است.
- تعامل کاربر با موتور جستجو تنها در ارسال پرس‌وجو، کلیک کردن پیوند تبلیغ‌ها و دوباره بازگشت به صفحه نتایج خلاصه می‌شود. همان‌طور که گفته شد، این فرایند می‌تواند پیچیده‌تر باشد.
- پلتفرم مورد بررسی، مرورگرهای نصب شده بر روی کامپیوترهای شخصی است. پلتفرم موبایل و اطلاعات ارسالی در آن پلتفرم در نظر گرفته نمی‌شود.
- در محاسبات، برای پیش‌پردازش از مدل ساده فضای برداری و آماره TF-IDF استفاده می‌شود. همان‌طور که پیش از گفته شد، امروزه روش‌های پیچیده‌تری مثل مدل زبانی، یادگیری ماشین و شبکه عصبی نیز در موتورهای جستجو استفاده می‌شود.

هدف مقاله مرجع [۸]، این است که نقض حریم خصوصی و موضوع‌های یادگرفته شده توسط موتور جستجو را بیابد و به کاربر گزارش دهد. مدل حریم خصوصی استفاده شده در این مرجع ϵ -عدم تمایزپذیری است. این مدل بر این اساس است که اعتقاد موتور جستجو درباره علاقه کاربر به موضوع c ، بعد از k تعامل ارسال پرس‌وجو و دریافت پاسخ آن، نباید بیشتر از e^ϵ باشد. این مفهوم به شکل زیر صوری می‌شود.

$$\frac{\Pr[X_c = c | \Omega_k, \mathcal{E}_k]}{\Pr[X_c = c | \mathcal{E}_1]} \leq e^\epsilon$$

X_c : متغیر تصادفی مربوط به پسند کاربر به موضوع c .

Ω_k : سه‌تایی پرس‌وجوی جستجو، پاسخ و پیوند تبلیغ کلیک شده.

\mathcal{E}_n : دانش پیشین موتور جستجو در مرحله n ام.

برای اندازه‌گیری مقادیر در رابطه بالا برای موضوع‌های مختلف، به جای استفاده از تک‌تک تعامل‌ها، از پرس‌وجوهای کاوشگر^۱ استفاده می‌شود. پرس‌وجوهای کاوشگر، در بازه‌های زمانی مشخص به موتور جستجو ارسال می‌شوند. این نوع از پرس‌وجو، بر اساس تعدادی فرض ساخته می‌شود. رشته کلمه‌ها در این پرس‌وجو، به شکلی انتخاب می‌شود که در راستای موضوع‌های جستجو شده اخیر کاربر باشد. همچنین، این پرس‌وجو مستقل از پرس‌وجوهای قبلی است. در نتیجه لازم نیست سابقه جستجوی کاربر در نظر گرفته شود، چون تمام آن‌ها در پرس‌وجوی کاوشگر و نتیجه آن نهفته است. علاوه بر آن، در پاسخ به این پرس‌وجو هیچ پیوند تبلیغی کلیک نمی‌شود و تنها متن تبلیغ‌های ارائه شده توسط موتور جستجو برای تحلیل‌های بعدی در نظر گرفته می‌شوند. فرض آخر این است که، اگر بر اساس هر پرس‌وجوی کاوشگر، ϵ -عدم تمایزپذیری برقرار باشد، آنگاه، کل جلسه ϵ -عدم تمایزپذیری را برآورده می‌کند. البته این مقاله به دنبال یافتن نقض حریم خصوصی است. بنابراین، در هر مرحله که رابطه ϵ -عدم تمایزپذیری برقرار نباشد، هشدار لازم به کاربر داده می‌شود. شیوه پردازش متن تبلیغ‌ها و ساخت پرس‌وجوی کاوشگر از پیش‌پردازش رشته پرس‌وجو و پاسخ به تعامل‌های پیشین کاربر با موتور جستجو استفاده می‌شود. این پیش‌پردازش بر اساس روش فضای برداری و با به‌کارگیری آماره TF-IDF انجام می‌شود.

هدف مقاله مرجع [۹]، این است که روشی ارائه کند که نقض حریم خصوصی را تشخیص دهد و برای دفاع در برابر آن روشی را ارائه کند. مدل حریم خصوصی استفاده شده در این مرجع (m, ϵ) -انکارپذیری قابل قبول است. در این مدل، برای هر موضوع حساس، $m - 1$ موضوع غیرحساس وجود دارد که در یک جلسه با k تعامل ارسال پرس‌وجو و دریافت پاسخ آن رابطه زیر برقرار باشد:

$$e^{-\epsilon} \leq \prod_{j=0}^{k-1} \frac{\Pr[\Omega_{k-j} | X = x, \mathcal{E}_{k-j}]}{\Pr[\Omega_{k-j} | X = x_i, \mathcal{E}_{k-j}]} \leq e^{\epsilon}$$

در رابطه بالا، x موضوع حساس و x_i موضوع غیرحساس است. در این رابطه، m -بی‌نامی دیده می‌شود. به ازای هر موضوع حساس باید $m - 1$ موضوع غیرحساس دیگر وجود داشته باشد که با پارامتر ϵ قابل تمایز نیستند. همچنین، شکل ضعیف حریم خصوصی تفاضلی نیز دیده می‌شود. به جای تمام موضوع‌های موجود، موضوع حساس، در نزدیکی m موضوع غیرحساس قرار می‌گیرد که میزان این نزدیکی با مقدار ϵ تعیین می‌شود.

^۱ prob query

گفتنی است، در این مقاله اثبات شده است که اگر رابطه ϵ -عدم تمایزپذیری مرجع [۸] میان موضوع‌های جستجو شده کاربر وجود داشته باشد، آنگاه، رابطه $(m, 4\epsilon)$ -انکارپذیری قابل قبول نیز برقرار خواهد بود.

در این مقاله، برای حفظ حریم خصوصی با کمک مدل (m, ϵ) -انکارپذیری قابل قبول روش موضوع نایب^۱ ارائه شده است. در این روش، تعدادی موضوع غیر حساس انتخاب شده و تعدادی پرس‌وجو در رابطه با آن‌ها تولید می‌شود. سپس، پرس‌وجوهای کاربر با پرس‌وجوهای موضوع‌های نایب درهم‌ریخته می‌شوند و همگی به موتور جستجو ارسال می‌شود. نتایج نشان می‌دهد که با این روش می‌توان حداکثر میزان (m, ϵ) -انکارپذیری قابل قبول را به ازای $\epsilon = 0$ و $m = 2$ عملی کرد. اشکال این روش این است که روش ارائه شده، به صورت برخط و در عمل که پرس‌وجوهای کاربر از پیش مشخص نیستند، کارایی ندارد. علاوه بر آن، پرس‌وجوهای ارسالی کاربر و رفتار جستجوی او در جستجوهای مختلف با هم همبستگی دارند. به عبارت دیگر، نمی‌توان پرس‌وجویی در مورد یک موضوع ارسال کرد که با موضوع‌های دیگر هیچ ارتباطی نداشته باشد. همچنین، صرف ارسال پرس‌وجو کفایت نمی‌کند. رفتار جستجوی کاربر در یادگیری پروفایل او نیز موثر است.

هدف مقاله مرجع [۱۰]، این است که روشی ارائه کند که اطلاعات شخصی به اندازه محدود و متناسب با خدمت مورد نظر در اختیار کاربران قرار گیرد. این مقاله بر اساس تعدادی فرض علاوه بر فرض‌های مطرح شده کلی، روش خود را مطرح می‌کند. فرض اول این است که تعامل کاربر با موتور جستجو در یک جلسه و در قالب دنباله‌ای از ورودی‌خروجی‌ها (مجموعه Z) است. فرض دیگر این است که کاربران موضوع‌های حساس (مورد پسند) خود را تعیین می‌کنند. با استفاده از یک تابع می‌توان، هر جفت ورودی‌خروجی را به یک مجموعه موضوع نگاشت کرد. در این مقاله، موازنه سودمندی و حریم خصوصی تعیین کننده موضوعات حساس کاربر است.

مدل حریم خصوصی استفاده شده در این مرجع δ -انکارپذیری قابل قبول است. کاربر u می‌تواند پسند خود به موضوع c (ارتباط یک جفت ورودی‌خروجی با موضوع حساس c) را به طور قابل قبول انکار کند، اگر رابطه زیر برقرار باشد.

$$\Pr[z \in Z_k^{u,c} | z \in Z_{att,k}] \leq \delta$$

در رابطه بالا، z یک جفت ورودی‌خروجی، $Z_k^{u,c}$ دنباله‌ای از ورودی‌خروجی‌ها در گام k ام است که کاربر u آن را در دسته موضوع c قرار می‌دهد، و $Z_{att,k}$ دنباله‌ای از ورودی‌خروجی‌ها در گام k ام است که مهاجم آن‌ها را در اختیار دارد. این مدل با (m, ϵ) -انکارپذیری قابل قبول معرفی شده در مقاله مرجع [۹] متفاوت است. در مقاله مرجع [۹]، ورودی‌خروجی مشاهده شده را می‌توان به تعدادی موضوع متفاوت نگاشت کرد. در این مقاله،

^۱ proxy topic

δ مشخص می‌کند که کاربر به چه میزان انتظار دارد که مهاجم بتواند مشاهده خود را به یک موضوع حساس نگاشت کند. اگر حاصل احتمال بیشتر از δ باشد، آنگاه، کاربر نمی‌تواند پسند خود را انکار نماید. رابطه بالا را می‌توان به شکل زیر نوشت:

$$\Pr[z \in Z_k^{u,c} | z \in Z_k] \leq \delta \Pr[z \in Z_{att,k} | z \in Z_k]$$

در رابطه بالا، Z_k دنباله ورودی خروجی‌ها در گام k ام است. $\Pr[z \in Z_{att,k} | z \in Z_k]$ نیز نشان دهنده توان مهاجم است. برای یک مهاجم عمومی، این مقدار برابر با یک است. با وجود این، برای یک مهاجم محلی، مقدار این احتمال مقداری کمتر از یک است. مهاجم محلی قدرت بیشتری دارد و برای مثال می‌تواند مشخص کند که کدام جفت ورودی خروجی‌ها، مربوط به P_c (پراکسی کاربر c) است.

در این مقاله برای بررسی حریم خصوصی کاربر ابزار 3PS ارائه شده است. معماری ابزار 3PS از سه قسمت مجموعه کاربران، مجموعه سامانه‌های نایب^۱ (به منظور ایجاد هویت‌های گروهی^۲)، و موتور جستجو تشکیل شده است. اشکالی که این روش دارد این است که مجموعه پراکسی‌ها توسط موتور جستجو مدیریت می‌شوند. این موضوع به معنای عدم تطابق ابزار ارائه شده با موتورهای جستجوی فعلی است. به همین دلیل، برای نشان دادن درستی روش ارائه شده، یک سامانه توصیه‌گر ساده با توجه به روش مطرح شده پیاده‌سازی شده است. نحوه انتخاب سامانه‌های نایب به این شکل است که کاربر با توجه به رابطه‌های زیر نزدیک‌ترین سامانه نایب به موضوع مورد جستجو خود را انتخاب می‌کند تا بیشترین استفاده از خدمت شخصی‌سازی را ببرد.

$$\min_{p \in P} \sum_{c \in C} |\Pr[z \in Z_{u,k}^{u,c} | z \in Z_{u,k}] - \Pr[z \in Z_{u,k}^{u,c} | z \in Z_{p,k}]|$$

$$\Pr[z \in Z_k^{u,c} | z \in Z_{p,k}] \leq \delta$$

❖ مدل‌های مبتنی بر تصادفی کردن

با بررسی ادبیات این حوزه به این نتیجه می‌رسیم که هنوز از مدل‌های مبتنی بر تصادفی کردن برای حفظ حریم خصوصی کاربران در تعامل با موتورهای جستجو استفاده نشده است. مهم‌ترین این مدل‌ها حریم خصوصی تفاضلی است که این پیشنهاد رساله بر اساس این مدل تهیه شده است.

❖ سایر روش‌ها و مقاله‌ها

^۱ proxy system

^۲ group identity

در مقاله مرجع [۱۹]، چهار سطح برای حفظ حریم خصوصی در ارتباط با موتورهای جستجو تعریف شده است. در سطح اول (شبه هویت)، اطلاعات مربوط به هویت کاربر مثل موقعیت مکانی یا آدرس IP محافظت می‌شود (این اطلاعات از پرس‌وجوها حذف می‌شود). موتور جستجو بر اساس اطلاعات حفاظت نشده می‌تواند پروفایل ایجاد کند. در سطح دوم (هویت گروهی)، گروهی از کاربران با یک هویت واحد با موتور جستجو ارتباط دارند. در این حالت، نمی‌توان از بازیابی اطلاعات شخصی‌سازی شده به طور کلی برای هر شخص استفاده کرد. با وجود این، اگر کاربران به درستی گروه‌بندی شوند، کاربران با پسندهای مشترک در یک گروه خواهند بود و می‌توان از بازیابی اطلاعات شخصی‌سازی شده برای کل گروه بهتر استفاده کرد. موتور جستجو بر اساس اطلاعات گروهی می‌تواند پروفایل ایجاد کند. یک راه پیاده‌سازی استفاده از سامانه نایب است. پرس‌وجوها از طریق سامانه نایب به موتور جستجو ارسال می‌شوند. به جای یک سامانه نایب، ممکن است گروهی از آن‌ها نیز وجود داشته باشند. در این طرح، هویت کاربر پشت سامانه نایب پنهان می‌شود.

در سطح سوم (بدون هویت)، هیچ اطلاعات هویتی درباره کاربر، حتی اطلاعات گروهی، در موتور جستجو وجود نخواهد داشت. در این حالت عملیات رتبه‌بندی نتایج باید در سامانه سمت کاربر اتفاق بیفتد. یکی از راه‌های پیاده‌سازی، استفاده از شبکه TOR است. با استفاده از TOR، موتور جستجو هیچ اطلاعاتی درباره کاربر نمی‌تواند بدست بیاورد. بنابراین، نمی‌تواند هیچ پروفایلی برای کاربر ایجاد کند. در سطح چهارم (بدون اطلاعات شخصی)، موتور جستجو از هویت و اطلاعات شخصی فرد هیچ اطلاعی ندارد. برای پیاده‌سازی می‌توان از یک شخص ثالث مورد اعتماد استفاده کرد که این شخص به جای کاربر پرس‌جوهای او را ارسال می‌کند. در این حالت موتور جستجو باید به طور قانونی از نگهداری اطلاعات افراد منع شود. این سطح بالاترین حالت حریم خصوصی است که رسیدن به آن آسان نیست.

در مقاله مرجع [۴]، برای حفظ حریم خصوصی کاربران در ارتباط با موتورهای جستجو، یک پروتکل ارائه شده است. در این پروتکل، b کاربر به صورت نظیر به نظیر^۱ با هم در ارتباط هستند. یک حافظه مشترک میان آن‌ها وجود دارد. هر کاربر پرس‌وجوی خود را با کلید یک سیستم رمز متقارن، رمز کرده و در حافظه مشترک قرار می‌دهد و منتظر می‌ماند تا یکی از b-1 کاربر دیگر، پرس‌وجوی او را برداشته و با کلید مشترک گروه ترجمه کرده و به جای او به موتور جستجو ارسال کند. کاربر ارسال کننده، پاسخ دریافتی را دوباره با کلید مشترک رمز کرده در حافظه مشترک قرار می‌دهد تا کاربر اصلی پاسخ را بردارد. با این پروتکل، موتور جستجو نمی‌تواند برای هر کاربر یک پروفایل ایجاد کند. چون پرس‌وجوهای کاربران بین b کاربر پخش می‌شود. همچنین، به دلیل وجود رمزنگاری، یک نفوذگر نمی‌تواند پرس‌وجوها و پاسخ‌ها را ببیند. علاوه بر آن، چون b کاربر در یک گروه

^۱ peer to peer

وجود دارند، کاربران یک گروه می‌توانند پروتکل را به هر کدام از افراد موجود در گروه نسبت دهند و نمی‌توانند به طور مشخص، یک فرد را به اجرای آن نسبت دهند.

در این مقاله، یک پروتکل دیگر با همین فرض‌ها و البته با طراحی متفاوت در نحوه توزیع کلید و همچنین، انتخاب کاربر ارسال‌کننده پرس‌وجو، ارائه شده است. ضعف‌ها و نقاط قوت هر کدام از این دو پروتکل بیان شده‌اند و در نهایت یک پروتکل بر اساس نقاط قوت دو پروتکل قبلی طراحی شده است.

در مقاله مرجع [۶]، آمده است که موتورهای جستجو با دریافت پرس‌وجوهای کاربر، با اطلاعاتی بیشتر از آنچه که کاربر انتظار دارد، پروفایل کاربر را ایجاد می‌کند. در این مقاله طرحی ارائه می‌شود که بر اساس آن اطلاعات کاربر عمومی می‌شود تا موتور جستجو بیش از حد مورد نیاز در مورد کاربر اطلاعات نداشته باشد. طرح ارائه شده نیاز به هیچ فرد ثالث مورد اعتمادی ندارد. در این مقاله، یک برنامه روی سامانه کاربر اجرا می‌شود که پروفایل دقیق کاربر و همچنین، پروفایل سفارشی‌شده او بر اساس هستی‌شناسی موضوع‌ها را نگهداری می‌کند. کاربر می‌تواند موضوع‌های حساس خود را مشخص کند. بر اساس موضوع‌های حساس، پروفایل سفارشی‌شده تهیه می‌شود.

وقتی کاربر یک پرس‌وجو ارسال می‌کند، پروفایلی عمومی‌شده بر اساس محتوای پرس‌وجو و همچنین، پروفایل سفارشی‌شده ایجاد می‌شود. پروفایل عمومی‌شده به همراه پرس‌وجو، برای موتور جستجو ارسال می‌شود. موتور جستجو بر اساس پروفایل ارسالی به پرس‌وجو پاسخ می‌دهد. فهرست پیوندهای نتیجه، به برنامه اجراشده در سامانه کاربر داده می‌شود. برنامه یا همان نتایج را به کاربر نمایش می‌دهد یا اینکه بر اساس پروفایل دقیق او، پیوندها را مجدداً مرتب می‌کند.

در مقاله مرجع [۱۴]، همراه با پرس‌وجوهای کاربر، تعدادی پرس‌وجوی پوششی به صورت تصادفی ارسال می‌شود تا موتور جستجو را گمراه کند. برای این کار یک افزونه برای مرورگر پیاده‌سازی شده است تا این عملیات اتفاق بیفتد.

۲-۳ نقد و بررسی پژوهش‌های پیشین

در این قسمت، پژوهش‌هایی که در قسمت قبل توضیح داده شده‌اند، نقد و بررسی می‌شوند و اشکال‌های هر کدام از طرح‌های ارائه شده بیان می‌شوند.

اشکالی که طرح ارائه شده در پژوهش مرجع [۱۲] دارد این است که در این طرح نمی‌توان از امکان شخصی‌سازی پاسخ‌ها که در موتور جستجو وجود دارد، استفاده کرد. برای مثال، نمی‌توان از امکان موتور جستجوی گوگل که اخبار محلی را نمایش می‌دهد، استفاده کرد. در این طرح، همراه با هر پرس‌وجو، k

پرس‌وجوی پوششی دیگر نیز ارسال می‌شود. در نتیجه، این طرح از نظر سودمندی برای کاربر بسیار ضعیف است و سربار زیادی هم بر موتور جستجو ایجاد می‌کند. علاوه بر آن، پرس‌وجوهای پوششی با توجه به موضوع پرس‌وجوی اصلی ایجاد می‌شوند. با وجود این، به تارخچه جستجوهای پیشین کاربر توجه نمی‌شود. بنابراین، در ایجاد پرس‌وجوی پوششی به آنچه که موتور جستجو از کاربر می‌داند توجه نمی‌شود. همچنین، در این طرح، در ارسال پرس‌وجوهای پوششی به جای توجه به تک‌تک رکوردها، به پروفایل ایجاد شده در موتور جستجو توجه دارد. بنابراین، از مدل‌های مبتنی بر افراز کردن (مانند k -بی‌نامی) بهره نمی‌برد.

هدف طرح ارائه شده در مرجع [۵]، پنهان کردن پرس‌وجوهای اصلی با ارسال پرس‌وجوهای پوششی است. در این طرح، موضوع پرس‌وجوی اصلی کاربر پیدا می‌شود. تعدادی موضوع در فاصله مفهومی c از موضوع پرس‌وجوی اصلی از آنالوژی محاسباتی موضوع‌ها پیدا می‌شود. بر اساس این موضوع‌ها، $k - 1$ پرس‌وجو پوششی ارسال می‌شود. بنابراین، طرح ارائه شده بر اساس مدل k -بی‌نامی است. موضوع پرس‌وجوی اصلی بر اساس آماره TF-IDF در قالب یک آرایه استخراج می‌شود. اشکالی که این طرح دارد این است که ارسال پرس‌وجوهای پوششی می‌تواند در حجم زیاد و با تعداد کاربران بالا، باعث ایجاد سربار زیاد در شبکه شود. این سربار به این دلیل می‌تواند زیاد محسوب شود که پرس‌وجوهای پوششی بدون توجه به الگوریتم یادگیری پروفایل در موتور جستجو ایجاد می‌شوند.

هدف طرح ارائه شده در مرجع [۱۵]، پنهان کردن کلمه کلیدی اصلی کاربر با اضافه کردن $k - 1$ کلمه کلیدی به آن است. کلمه‌های کلیدی پوششی با توجه به آنالوژی آن‌ها انتخاب می‌شود. بنابراین، مدل حریم خصوصی این پژوهش، l -تنوع آنالوژی است. اشکالی که طرح ارائه شده دارد این است که در آن کلمه‌های کلیدی با هم OR می‌شوند و نمی‌توان جمله یا عبارت داشت. پس در عمل طرح ارائه شده نمی‌تواند پرس‌وجوهای واقعی را پشتیبانی کند. طرح ارائه شده، به ازای تمام پرس‌وجوها اجرا می‌شود. بنابراین، سودمندی طرح برای موضوع‌هایی که برای کاربر حساس نیستند، کاهش پیدا می‌کند. همچنین، در اضافه کردن کلمه‌های کلیدی به پرس‌وجوی اصلی به پرس‌وجوهای پیشین کاربر توجه نمی‌شود. طرح ارائه شده پرس‌وجوها را مستقل از هم در نظر می‌گیرد.

به روش ارائه شده در مرجع [۸]، تعدادی اشکال و نقد وارد است. در تعریف ϵ -عدم تمایزپذیری، می‌توان به جای نسبت دو احتمال شرطی از فاصله میان توزیع احتمال اولیه و نهایی استفاده کرد. برای رابطه فاصله دو توزیع احتمال، روش‌های متعددی وجود دارد. در پیش‌پردازش بدون توجه به عملکرد موتور جستجو، از روش فضای برداری و آماره TF-IDF استفاده می‌شود. همان طور که گفته شد، موتورهای جستجو از روش‌های پیچیده‌تری مثل یادگیری ماشین، شبکه عصبی و یا مدل زبانی برای یادگیری پروفایل کاربر استفاده می‌کنند.

ارسال پرس‌وجوی کاوشگر در یادگیری پروفایل کاربر توسط موتور جستجو موثر است. در صورتی که، در انتخاب آن دقت نشود، ممکن است، اعتقاد موتور جستجو نسبت به یک موضوع تغییر کند. همچنین، این پرس‌وجو بر اساس آماره TF-IDF تولید می‌شود. در این روش مکان کلمه‌ها در رشته در نظر گرفته نمی‌شود، با وجود این، در موتورهای جستجو، کلمه‌هایی که زودتر ظاهر می‌شوند، در نتایج بیشتر تاثیر دارند. بنابراین، بدون در نظر گرفتن ویژگی‌های مهم برای موتور جستجو در یادگیری پروفایل کاربر، نمی‌توان پرس‌وجوهای کاوشگر مفید تولید و ارسال کرد.

در پژوهش‌های مرجع [۹] و [۱۰] از مدل انکارپذیری قابل قبول استفاده شده است. در این پژوهش‌ها، هدف این است که کاربر بتواند علاقه‌مندی خود به یک موضوع حساس را با وجود $m - 1$ موضوع پوششی دیگر در پروفایل خود به طور قابل قبول انکار کند. علاقه‌مندی کاربر به این m موضوع باید تقریباً به اندازه یک‌دیگر باشد. در نتیجه، کاربر می‌تواند به طور قابل قبول علاقه‌مندی خود به موضوع حساس را انکار کند. در مرجع [۹]، رابطه حریم خصوصی به شکلی نوشته شده است که بتوان از ابزار ارائه شده در مقاله مرجع [۸]، در این مقاله نیز استفاده کرد. گفتنی است، این مقاله، تمام ضعف‌های مقاله مرجع [۸] را، به دلیل استفاده از ابزار ارائه شده در آن مرجع، به ارث می‌برد. به عنوان جمع‌بندی ضعف‌های مطرح در این مقاله در زیر آمده‌اند.

- حفظ حریم خصوصی تنها در پروفایل کوتاه‌مدت در نظر گرفته شده است.
- شخصی‌سازی تنها در تنوع تبلیغ‌های توصیه شده موتور جستجو در نظر گرفته می‌شود.
- رتبه‌بندی پیوندهای نتیجه شده در پاسخ به جستجوی کاربر در نظر گرفته نمی‌شود.
- شخصی‌سازی در جستجوهای مرتبط با اخبار یا آب و هوا در نظر گرفته نمی‌شوند.
- محتوا و مفهوم رشته پرس‌وجوها شامل کلمه‌های ظاهرشده، تعداد تکرار کلمه‌ها، ترتیب کلمه‌ها و مانند آن در نظر گرفته نمی‌شود.
- الگوی کلیک پیوندهای نتیجه جستجو کاربر و شیوه گردش میان پیوندها در نظر گرفته نمی‌شود. تعامل کاربر با موتور جستجو تنها در ارسال پرس‌وجو، کلیک کردن پیوند تبلیغ، و بازگشت به صفحه نتایج خلاصه می‌شود.
- تنها ورودی‌های متنی به عنوان رشته جستجو در نظر گرفته می‌شود. ورودی‌ها می‌توانند متنوع باشند. به عنوان نمونه می‌توان به ورودی‌های از نوع عکس و فیلم اشاره کرد.
- پلتفرم مورد بررسی مرورگرهای نصب‌شده بر روی کامپیوترهای شخصی است. پلتفرم تلفن همراه و اطلاعات ارسالی از آن پلتفرم مانند موقعیت جغرافیایی و آمار استفاده از برنامه‌های کاربردی در نظر گرفته نمی‌شود.

• در محاسبات، برای پیش‌پردازش و تولید پرس‌وجوی پوششی از مدل بازایی فضای برداری، مبتنی بر آماره TF-IDF، استفاده می‌شود. امروزه، روش‌های پیچیده‌تری مثل مدل زبانی، یادگیری ماشین و شبکه عصبی نیز در موتورهای جستجو استفاده می‌شوند. استفاده از مدل مبتنی بر آماره TF-IDF، به این معنی است که فرض شده است که موتور جستجو از این روش برای یادگیری پروفایل کاربر استفاده می‌کند که فرض درستی نیست، چون موتور جستجو یک سامانه جعبه‌سیاه است و اطلاعی از روش یادگیری پروفایل کاربر در دسترس نیست.

اشکالی که روش ارائه شده در مرجع [۱۰] دارد این است که رابطه انتخاب موضوع نایب، الزاما بهترین سامانه نایب مربوط به موضوع c را انتخاب نمی‌کند، چون عملیات جمع بر روی قدر مطلق تفاضل‌ها اتفاق می‌افتد. ممکن است یک سامانه نایب شباهت کمتری به موضوع c نسبت به سامانه نایب دیگر داشته باشد، ولی حاصل جمع برای آن در رابطه بالا کمتر شود. همچنین، برای عملی کردن این طرح لازم است که موتور جستجو پیاده‌سازی خود را تغییر دهد.

اشکالی که طرح ارائه شده در مرجع [۴] دارد این است که برای استفاده از پروتکل مطرح شده در مقاله نیاز است که تعداد زیادی از کاربران در آن شرکت کنند و موفقیت‌آمیز بودن آن به شرکت افراد در پروتکل بستگی دارد. همچنین، سربر اجرایی پروتکل بالاست [۵]. علاوه بر آن ممکن است بعضی از کاربرانی که در پروتکل شرکت می‌کنند، موذی باشند و پرس‌وجوهای نامناسب و غیرقانونی ارسال کنند. در این مقاله سازوکاری برای جلوگیری از این دسته از کاربران وجود ندارد [۵]. اشکال اصلی روش مرجع [۶] نیز این است که موتورهای جستجوی موجود باید پیاده‌سازی خود را مطابق با طرح ارائه شده عوض کنند.

اشکالی که روش ارائه شده در مرجع [۱۴] دارد این است که ارسال پرس‌وجوهای تصادفی باعث می‌شود که QoS^1 موتور جستجو پایین بیاید و کاربر نتواند از خدمت شخصی‌سازی پرس‌وجوها استفاده کند. همچنین، در مرجع [۵۲] نشان داده شده است که با استفاده از الگوریتم‌های طبقه‌بندی در یادگیری ماشین، می‌توان با احتمال بالاتر از ۰٫۸ و برای بعضی از کاربرها با احتمال نزدیک به یک، پرس‌وجوهای اصلی کاربر را از پرس‌وجوهای پوششی، تشخیص داد. به عبارت دیگر، این ابزار در حفظ حریم خصوصی کارایی ندارد.

در پژوهش مرجع [۵۳]، موتور جستجو درست‌کارولی کنج‌کاو در نظر گرفته شده است. در این پژوهش، هدف از حفظ حریم خصوصی، افشا نشدن موضوع‌های مورد علاقه کاربر در پروفایل او است. برای این کار، از روش

^۱ quality of service

رمزنگاری هومومورفیک کامل^۱ استفاده شده است. مدل سامانه در نظر گرفته شده در این پژوهش به این شکل است که پروفایل کاربر و فرایند شخصی سازی در سمت دستگاه کاربر (کلاینت) با توجه به تاریخچه جستجوهای ذخیره شده او در کلاینت ایجاد می شود. در مدل تهدید ارائه شده نیز، تهدیدهای به مخاطره افتادن کلاینت، موتور جستجو، و کانال ارتباطی میان این دو موجودیت در نظر گرفته شده است. با توجه به مدل سامانه و روش حفظ حریم خصوصی، موتور جستجو باید الگوریتم بازیابی خود را به شکلی تغییر دهد که بتواند سندهای مرتبط با پرس و جوی رمز شده را با ویژگی موجود در رمزنگاری هومومورفیک کامل بازیابی کند. گفتنی است، در رمزنگاری هومومورفیک کامل، می توان عملیات جمع، ضرب، مقایسه و غیره را بر روی مقادیر رمز شده انجام داد. در برخی از پژوهش های پیشین، موتور جستجو به شکل درستکارولی کنجکاو در نظر گرفته نشده است. در این پژوهش ها، مانند مراجع [۵۴-۶۰]، موتور جستجو مورد اعتماد است. در برخی از این دسته از پژوهش ها، موتور جستجو باید تاریخچه جستجوی کاربران را برای پژوهش گران یا شرکت های ثالث بدون افشا شدن اطلاعات حساس کاربران منتشر کند. در این پژوهش ها، موتور جستجو از حریم خصوصی تفاضلی برای ایجاد یک پایگاه داده ساختگی از تاریخچه جستجوی کاربران بهره می برد. در بعضی دیگر از پژوهش ها مانند مرجع [۵۳]، شرکت (های) ثالث می توانند به صورت تعاملی پرس و جو ارسال کنند و موتور جستجو به صورت مغشوش شده به پرس و جوی آن (ها) بر اساس مدل حریم خصوصی تفاضلی پاسخ می دهد.

در پژوهش مرجع [۵۷]، موتور جستجو مورد اعتماد است. پایگاه داده تاریخچه جستجوی همه کاربران موتور جستجو به شرکت های ثالث داده می شود. از این اطلاعات در جهت ایجاد تبلیغات بهتر و دقیق تر استفاده خواهد شد. در این پژوهش، موتور جستجو بر اساس مدل حریم خصوصی تفاضلی پایگاه داده ساختگی تاریخچه جستجوهای کاربران را منتشر خواهد کرد. برای حفظ سودمندی، به جای انتشار اطلاعات آماری مغشوش شده از تاریخچه جستجوها (تعداد پرس و جوها درباره هر موضوع، هیستوگرام، و غیره)، هر کدام از پرس و جوی جستجوی کاربران مغشوش می شود. از آنجایی که پرس و جوی جستجو یک عبارت متنی است (مانند عبارت «نسخه جدید تلفن سامسونگ»)، نمی توان از سازوکار لاپلاس استفاده کرد. علاوه بر آن، برای حفظ سودمندی پرس و جوی مغشوش شده باید از نظر مفهومی به پرس و جوی اصلی تا حد امکان نزدیک باشد. برای مثال، عبارت «نسخه جدید تبلت اپل» به جای عبارت «نسخه جدید تلفن سامسونگ» منتشر شود. برای این کار، از طریق آنتولوژی محاسباتی موضوع (های) پرس و جو استخراج می شود. با توجه به موضوع (های) استخراج شده، موضوعی از آنتولوژی را در نظر می گیرد که موضوع کلی پرس و جو است و شامل همه موضوع های آن می شود. با استفاده از

^۱ fully homomorphic encryption (FHE)

سازوکار نمایی از میان موضوع‌های موجود در درختواره موضوع کلی پرس‌وجوی اصلی، یک موضوع انتخاب شده و بر اساس آن یک پرس‌وجوی جدید ایجاد می‌شود. موضوعی احتمال بالاتری برای انتخاب دارد که به موضوع پرس‌وجوی اصلی شبیه‌تر باشد. برای اندازه‌گیری شباهت دو موضوع معیاری در مقاله ارائه شده است.

در پژوهش مرجع [۵۸]، رکوردهای پایگاه‌داده به شکل $\langle uid, qc, c \rangle$ است. uid شناسه کاربر، qc جفت پرس‌وجو و سند کلیک شده، و c تعداد دفعات این جستجو است. با استفاده از مدل حریم خصوصی تفاضلی تقریبی مقادیر c مغشوش می‌شود و پایگاه‌داده ساختگی ایجاد می‌شود. در صورت صفر شدن c ، رکورد مرتبط با آن از پایگاه‌داده ساختگی حذف می‌شود. در نتیجه، پایگاه‌داده ساختگی به شرکت ثالث داده می‌شود. برای اندازه‌گیری سودمندی، فاصله توزیع احتمال رکوردها در پایگاه‌داده اصلی و ساختگی با روش Kullback-Leibler محاسبه می‌شود. برای افزایش کارایی، سازوکار حفظ حریم خصوصی نیز به صورت توزیع شده تعریف شده است.

در پژوهش مرجع [۵۹]، n کلمه از هر پرس‌وجو با هم در نظر گرفته می‌شود تا مفهوم پرس‌وجو از کلمه‌های کنار هم استخراج شود. از هر کاربر به اندازه d_0 ، که به صورت تصادفی انتخاب می‌شود، عبارت جستجو انتخاب می‌شود. تعداد سندهای کلیک شده به ازای هر عبارت نیز محاسبه و در $c(ph_i, u_j)$ نگهداری می‌شود. بر اساس مفهوم و شباهت عبارت‌ها به هم، آن‌ها خوشه‌بندی می‌شوند. به طور تصادفی عبارتی از هر خوشه انتخاب می‌شود. مجموع تعداد دفعات کلیک شدن سند u_j به ازای هر عبارت موجود در آن خوشه انتخابی، به عنوان $c'(ph_i, u_j)$ در نظر گرفته می‌شود. این مقدار با مقداری نویز از توزیع $Lap\left(\frac{d_0}{\epsilon_0}\right)$ مغشوش می‌شود. مقدار مغشوش شده بعد از نرمال شدن به عنوان امتیاز عبارت‌های موجود در این خوشه در نظر گرفته می‌شود. بر اساس بیشینه امتیازها، برای هر کاربر d_1 رکورد انتخاب می‌شود. تعداد $c(ph_i, u_j)$ با مقداری نویز از $Lap\left(\frac{d_1}{\epsilon_1}\right)$ جمع و منتشر می‌شود.

در پژوهش مرجع [۶۰]، دنباله سندهای مشاهده شده کاربر در جلسه‌های مختلف بررسی شده است. در این مقاله، پایگاه‌داده شامل مجموعه‌ای از این جلسه‌ها است. مدل حفظ حریم خصوصی در اینجا، مدل حریم خصوصی تفاضلی است. در این کار قرار است تعداد جلسه‌هایی که در زمان t صفحه i را مشاهده می‌کنند، به صورت مغشوش شده منتشر شود. برای این منظور، تعداد جلسه‌های مشاهده کننده صفحه i در یک زمان مشخص با مقداری نویز از توزیع لاپلاس جمع می‌شود. به کمک داده‌های عمومی موجود، مدل وضعیت فضای

فعالیت‌های مختلف (مانند جستجو درباره بیماری، مطالعه اخبار اقتصادی، و غیره) یاد گرفته می‌شود. از این مدل، برای افزایش سودمندی مقادیر مغشوش شده استفاده می‌شود.

۳-۳ پژوهش‌های مرتبط با چالش‌های مطرح شده مسئله

در این قسمت، کارهای پژوهشی مرتبط با چالش‌های بیان شده درباره مسئله مطرح شده در این پیشنهاد رساله آمده‌اند. رویکرد این قسمت، مقایسه و تطابق کارهای پژوهشی موجود با مسئله مطرح در این پیشنهاد رساله است. یکی از چالش‌های مطرح شده، پایگاه‌داده رشدیابنده است. در مقاله مرجع [۶۱]، طرحی بر اساس مدل حریم خصوصی تفاضلی محلی ارائه شده است. در این طرح، میانگین مقادیر ارسال شده کاربران محاسبه می‌شود. بر خلاف مسئله مطرح در این پیشنهاد رساله، که تمام اطلاعات پایگاه‌داده مربوط به یک کاربر است، در این طرح، n کاربر حضور دارند. اطلاعات کاربران تنها یک بیت داده در نظر گرفته شده است که به صورت مستمر برای مسئول پایگاه‌داده ارسال می‌شود. سازوکار تصادفی کردن بر مبنای تغییر داده کاربران است. با وجود این، در مسئله مطرح شده در این پیشنهاد رساله، اطلاعات کاربر شامل تاریخچه جستجوهای او است. همچنین، به دلیل از دست رفتن سودمندی، امکان تغییر اطلاعات کاربر وجود ندارد و تنها با اضافه کردن رکورد جدید می‌توانیم، پایگاه‌داده را مغشوش کنیم. به عبارت دیگر، علاوه بر اضافه شدن اطلاعات جدید در طول زمان، به هنگام مغشوش کردن نیز تعدادی رکورد به پایگاه‌داده اضافه می‌شوند. به منظور کاهش سرعت رشد مقدار ϵ ، به دلیل فراخوانی‌های مکرر سازوکار تصادفی کردن (قضیه ترکیب ترتیبی)، در هر بازه زمانی که به اندازه مناسب اطلاعات کاربران تغییر کند، سازوکار تصادفی کردن اجرا می‌شود. در این پژوهش، فرض شده است که اطلاعات کاربر در فاصله‌های زمانی طولانی تغییر می‌کند. برای مثال، اطلاعات مربوط به تنظیمات مرورگر یا فراوانی شکلهای استفاده شده در صفحه‌کلید. با وجود این، در مسئله مطرح شده ما، به ازای هر جستجوی کاربر محتوای پایگاه‌داده تغییر می‌کند و پروفایل کاربر به موجب آن، به‌روز می‌شود.

در پژوهش مرجع [۶۱]، در بازه‌های زمانی مشخص، با استفاده از پروتکل رای‌گیری، از کاربران درباره ارسال مغشوش اطلاعات نظرسنجی می‌شود. در صورتی که نتیجه رای‌گیری ارسال اطلاعات شد، با استفاده از پاسخ تصادفی شده اطلاعات کاربران مغشوش و ارسال می‌شود. در مقاله مرجع [۲۱]، طرحی مبتنی بر مدل حریم خصوصی تفاضلی برای پایگاه‌داده رشدیابنده ارائه شده است. در این مقاله، روشی ارائه شده است که طرح‌های حریم خصوصی برای پایگاه‌داده با اندازه ثابت را برای پایگاه‌داده رشدیابنده هماهنگ می‌کند. ایده مطرح در این پژوهش، باید در جهت استفاده در سنتز پایگاه‌داده مغشوش شده بررسی شود.

در مقاله‌های مرجع‌های [۶۲]–[۶۴]، از مدل حریم خصوصی تفاضلی برای خصوصی کردن الگوریتم k -means استفاده شده است. ایده‌های ارائه شده در این مراجع مناسب برای پاسخ‌گویی به مسئله مطرح شده در این پیشنهاد رساله نیستند. در این مراجع، هر رکورد پایگاه داده مربوط به یک کاربر است. با وجود این، در مسئله ما، تمام رکوردها مربوط به یک کاربر هستند. در سازوکار تصادفی کردن، اطلاعات مربوط به کاربران با تغییر دادن آن‌ها مغشوش می‌شوند. اما در مسئله ما، به دلیل از دست دادن سودمندی، امکان تغییر دادن رکوردها نیست و باید تعدادی رکورد جدید برای مغشوش کردن ارسال شود. در این مراجع، اندازه پایگاه داده ثابت است. ولی در مسئله ما، پایگاه داده رشدیابنده است. همچنین، در مسئله ما، بین رکوردها همبستگی نیز وجود دارد. هدف این مراجع، خصوصی کردن عضویت رکورد یک کاربر در یک خوشه است. در مقابل هدف طرح ما، خصوصی کردن نزدیکی مرکز یک خوشه به هر کدام از مرکزهای زیرخوشه‌های آن با توجه به خط‌مشی حریم خصوصی است. به عبارت دیگر، در مسئله ما، خوشه‌ها به شکل سلسله مراتبی هستند. چالش اصلی که در این مراجع به آن پرداخته شده، کاهش تعداد دورهای اجرای الگوریتم است تا مقدار ϵ کندتر رشد کند. در مسئله ما، به دلیل رشدیابنده بودن پایگاه داده، به ازای اضافه شدن هر رکورد جدید، حداقل یک دور الگوریتم باید اجرا شود. این موضوع، چالشی مهم در اندازه ϵ است که ما در مسئله خود باید به آن توجه کنیم.

چالش دیگر در حل مسئله، همبستگی رکوردهای پایگاه داده است. در مرجع [۴۷]، که در قسمت ۳–۲، جزئیات آن توضیح داده شد، طرحی برای حفظ حریم خصوصی تفاضلی در پایگاه داده دارای همبستگی ارائه شده است. برای اینکه بتوانیم نگاشت این طرح به مسئله خود را انجام دهیم، باید به سه سوال زیر پاسخ دهیم:

- چطور می‌توان رکوردهای همبسته را شناسایی کرد؟
- چطور می‌توان حساسیت را برای رکوردهای همبسته تعریف کرد؟
- چطور می‌توان سازوکار حریم خصوصی تفاضلی در سنتز پایگاه داده را برای پایگاه داده دارای همبستگی بازطراحی کرد؟

در مرجع [۶۵]، مدلی برای حریم خصوصی تفاضلی با وجود همبستگی میان رکوردهای پایگاه داده ارائه شده است. در این مقاله، دو پایگاه داده وابسته $D(L, R)$ و $D'(L, R)$ همسایه هستند، اگر تغییر یک رکورد D_i در $D(L, R)$ به D'_i در $D'(L, R)$ موجب تغییر $L - 1$ رکورد در $D'(L, R)$ به دلیل وجود وابستگی احتمالاتی R میان آن‌ها شود. بر این اساس رابطه حریم خصوصی تفاضلی به صورت

$$\max_{D, D'} \frac{\Pr[M(D(L, R) = S)]}{\Pr[M(D'(L, R) = S)]} \leq e^\epsilon$$

تعریف می‌شود. در مدل ارائه شده حساسیت پرس‌وجوها مانند مقاله مرجع [۴۷]، بر اساس ماتریس Δ تعریف می‌شود. تفاوتی که وجود دارد این است که در این مقاله، ایده‌های محاسبه همبستگی میان رکوردها صورتی شده است. البته ماتریس تولید شده نسبت به مقاله مرجع [۴۷]، دقیق‌تر است و موجب می‌شود نویز کمتری به اطلاعات اضافه شود. علاوه بر آن، برخلاف مقاله قبل که همبستگی به صورت خطی محاسبه می‌شود، در این مقاله، همبستگی میان رکوردها به صورت احتمالاتی محاسبه می‌شود. گفتنی است، سازوکار تصادفی کردن لاپلاس مانند مقاله مرجع [۴۷] محاسبه می‌شود. برای حل مسئله تعریف شده در این پیشنهاد رساله، باید از ایده‌های مطرح برای تولید ماتریس Δ در این مقاله استفاده کرد.

در مرجع [۶۶]، طرحی بر اساس مدل حریم خصوصی تفاضلی محلی و با در نظر گرفتن همبستگی در میان اطلاعات ارائه شده است. در این طرح قرار است، اطلاعات مغشوش و ارسال شده کاربران به سرور، در قالب یک پایگاه‌داده در اختیار پژوهشگران قرار گیرد. به عبارت دیگر، قرار است بر اساس اطلاعات کاربران، پایگاه‌داده‌ای ساختگی ایجاد شود. برای اینکه پایگاه‌داده ساختگی برای پژوهشگران مفید باشد، علاوه بر اطلاعات مغشوش شده، کاربران اطلاعاتی درباره ارتباط و همبستگی اطلاعات نیز به سرور ارسال می‌کنند. تفاوت‌هایی که این طرح با مسئله ما دارد این است که در مسئله ما، تمام اطلاعات مربوط به یک کاربر است. همچنین، پایگاه‌داده رشدیابنده است و علاوه بر آن، با اضافه کردن رکوردهای پوششی، اطلاعات کاربر مغشوش می‌شود. همچنین، به طور مشخص یک پرس‌وجو بر روی اطلاعات کاربر انجام شده و پروفایل او ایجاد می‌شود. همبستگی میان اطلاعات باید در زمان مغشوش کردن اطلاعات توسط خود کاربر در نظر گرفته شود تا پروفایل ایجاد شده، افشا کننده پسندهای حساس او نباشد.

در مرجع [۶۷]، طرحی برای حریم خصوصی تفاضلی ارائه شده است که در آن علاوه بر همبستگی میان رکوردها، اطلاعات پیشین مهاجم در رابطه با پایگاه‌داده نیز در نظر گرفته شده است. همان طور که در فصل ۲ نیز توضیح داده شد، در مدل حریم خصوصی تفاضلی اگر مهاجم به تمام رکوردهای پایگاه‌داده به غیر از یکی دسترسی داشته باشد، با مشاهده نتیجه پرس‌وجوی مغشوش شده، نمی‌توان مقدار رکورد مخفی را به دست آورد. این توضیح تا زمانی که رکوردهای پایگاه‌داده با هم همبستگی نداشته باشند، درست است. در این مقاله، طرحی برای حفظ حریم خصوصی در وجود همبستگی میان رکوردها و اطلاعات پیشین مهاجم ارائه شده است. در مسئله مطرح در این پیشنهاد رساله، موتور جستجو (مهاجم) تمام رکوردهای پایگاه‌داده را در اختیار دارد و

تحلیل‌گر است. در این حالت، کاربران باید با مغشوش کردن اطلاعات خود، پایگاه‌داده سنتز شده را در اختیار موتور جستجو قرار دهند تا بتوانند استنتاج موتور جستجو از پایگاه‌داده را مغشوش کنند.

در پژوهش مرجع [۶۸]، آمده است که برای حفظ حریم خصوصی هر رکورد از پایگاه‌داده می‌توان مقدار ستون شبه‌شناسه را به طور تصادفی تغییر داد. روش تغییر این مقدار بستگی به نوع داده دارد. برای داده‌های عددی، از توزیع لاپلاس و برای داده‌های دارای دسته‌بندی (مانند اندازه پیراهن) از توزیع نمایی می‌توان کمک گرفت. برای بقیه داده‌هایی که عددی و دارای دسته‌بندی نیستند، (مانند شغل فرد) باید از آنتولوژی محاسباتی کمک گرفت. در این حالت، بر اساس مقادیر موجود در ستون شبه‌شناسه مفهوم مشترک آن‌ها در آنتولوژی محاسباتی استخراج می‌شود. در این پژوهش، تعریفی برای محاسبه فاصله دو مفهوم از آنتولوژی ارائه شده است. بر اساس این تعریف، میانگین و واریانس مفاهیم موجود در ستون مورد نظر محاسبه می‌شود. از توزیع نرمال $\mathcal{N}(0, \sigma)$ که انحراف معیار آن ضریبی از انحراف معیار مقادیر موجود در ستون است، مقداری نویز به تصادف انتخاب می‌شود. با توجه به مقدار و علامت نویز، مفهومی از آنتولوژی در نزدیکی مفهوم اصلی موجود در ستون جایگزین می‌شود. به نویز اضافه شده در این سازوکار، نویز محتوایی^۱ گفته می‌شود.

در پژوهش مرجع [۶۹]، بر اساس مدل l -تنوع آنتروپی راهکاری برای حفظ حریم خصوصی در جستجوی وب شخصی‌سازی شده ارائه شده است. در این پژوهش، بر اساس موضوع پرس‌وجوی اصلی کاربر و پرس‌وجوهای مرتبط پیشین او در تاریخچه جستجوها، $l - 1$ موضوع دیگر از آنتولوژی محاسباتی موضوع‌ها انتخاب می‌شود. موضوع‌های پوششی به شکلی انتخاب می‌شوند که از نظر آنتروپی به مفهوم پرس‌وجوی اصلی نزدیک باشند. همچنین، بعد از ارسال پرس‌وجوهای اصلی و پوششی، با استفاده از نظریه بیز، اعتقاد موتور جستجو نسبت به موضوع‌های مورد علاقه کاربر به‌روز می‌شود. این اعتقاد نباید از یک مقدار آستانه بیشتر شود. اشکال‌هایی در این پژوهش وجود دارد. در شناسایی موضوع پرس‌وجوی اصلی کاربر، سندهای کلیک شده او در نظر گرفته نشده است. علاوه بر آن، همه پرس‌وجوها حساس در نظر گرفته شده‌اند. به عبارت دیگر، خط‌مشی کاربر در نظر گرفته نشده است. در نتیجه، به ازای هر پرس‌وجوی کاربر، تعدادی پرس‌وجوی پوششی وجود دارد. بنابراین، سربار طرح بالا است. همچنین، در ایجاد پرس‌وجوهای پوششی، الگوریتم یادگیری پروفایل کاربر نیز در نظر گرفته نمی‌شود. بنابراین، سازوکار حفظ حریم خصوصی به شکل دقیق و با کمترین میزان ارسال پرس‌وجوی پوششی نمی‌تواند کار کند.

^۱ semantic noise

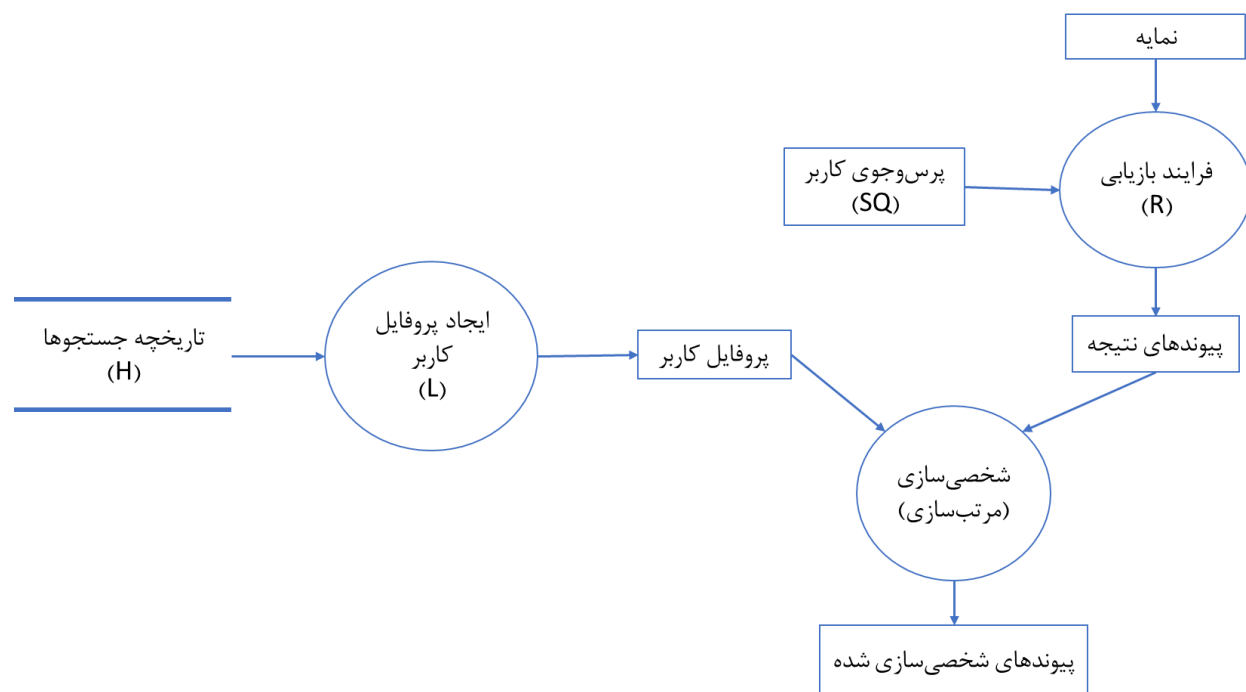
فصل چهارم

پیشنهاد رساله

در این فصل، ابتدا مسئله را به طور کلی بیان کرده و گام‌های پیشنهادی حل آن را مطرح می‌کنیم. آنگاه، مسئله را با توجه به گام‌های حل آن شرح داده و چالش‌ها و هدف‌های رساله را ارائه می‌کنیم. در پایان این فصل، زمان‌بندی انجام رساله پیشنهاد می‌شود.

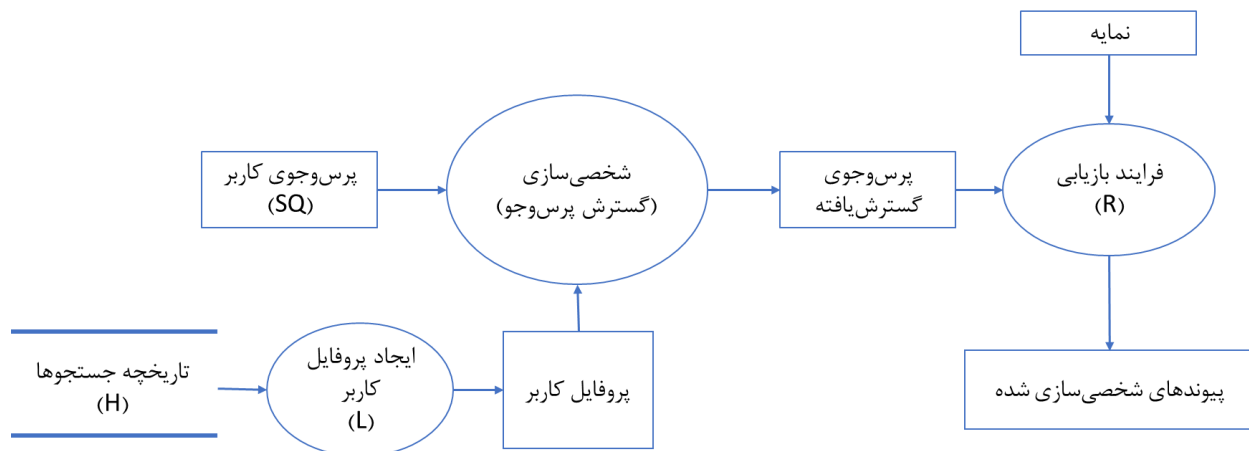
۴-۱ بیان کلی مسئله و گام‌های پیشنهادی برای حل آن

موتورهای جستجو در صفحه‌های وب می‌خزند. اطلاعات به دست آمده از این خزش نمایه‌سازی شده و یک نمایه، که نشان می‌دهد هر کلمه در کدام صفحه‌های وب موجود است، استخراج می‌شود. بر اساس پرس‌وجوی جستجوی دریافت شده از کاربر و نمایه استخراج شده، فرایند بازیابی لیستی مرتب شده از پیوندها به صفحه‌های وب را به کاربر باز می‌گرداند (فرایند R در شکل ۴-۱ و شکل ۴-۲). موتور جستجو، همچنین، پرس‌وجوهای جستجو و سندهایی را که کاربر در پی هر جستجو به آن‌ها مراجعه می‌نماید، که روی هم تاریخچه جستجوهای آن کاربر نامیده می‌شود، در یک پایگاه‌داده ذخیره می‌کند (پایگاه‌داده H در شکل ۴-۱ و شکل ۴-۲). با استفاده از اطلاعات ذخیره شده در این پایگاه‌داده، موتور جستجو پروفایل کاربر را ایجاد می‌نماید (فرایند L در شکل ۴-۱ و شکل ۴-۲). در پاسخ به هر پرس‌وجوی جستجوی جدید کاربر، موتور جستجو نتایج را برای آن کاربر شخصی‌سازی می‌کند. شخصی‌سازی نتایج بر اساس پروفایل استخراج شده کاربر و به یکی از دو روش



شکل ۴-۱: شخصی‌سازی نتایج پرس‌وجوی کاربر با روش مرتب‌سازی.

مرتب‌سازی^۱ یا گسترش پرس‌وجوی کاربر^۲ انجام می‌شود. در روش مرتب‌سازی، پیوندهای حاصل از فرایند بازیابی مجدداً بر اساس پروفایل کاربر مرتب می‌شوند (شکل ۴-۱). در روش گسترش پرس‌وجوی کاربر، رشته‌ای از کلمه‌ها، بر اساس پروفایل کاربر، به پرس‌وجوی کاربر اضافه شده و پرس‌وجوی حاصل به جای پرس‌وجوی اصلی پردازش می‌شود (شکل ۴-۲).



شکل ۴-۲: شخصی‌سازی نتایج پرس‌وجوی کاربر با روش گسترش پرس‌وجو.

با استخراج پروفایل کاربران، موتور جستجو از پسندهای کاربران مطلع می‌شود. با وجود این، بیشتر کاربران نسبت به افشای پسندهای خود حساسیت دارند و آن را در خط‌مشی حریم خصوصی خود اعلام می‌کنند. سوال اصلی آن است که چطور می‌توان حریم خصوصی کاربران را در جستجوی وب شخصی‌سازی شده حفظ کرد. منظور از حفظ حریم خصوصی، کنترل میزان افشا شدن پسندهای حساس کاربران است. در این پیشنهاد رساله، از مدل حریم خصوصی تفاضلی برای بیان خط‌مشی حریم خصوصی کاربران و عملی کردن آن استفاده می‌کنیم. از آنجایی که موتور جستجو تمامی پایگاه‌داده را در اختیار دارد و پرس‌و‌جوهای خود را به منظور ایجاد پروفایل کاربر بر روی آن اجرا می‌کند، از روشی الهام گرفته شده از سنتز^۳ پایگاه‌داده مغشوش شده بر اساس مدل حریم خصوصی تفاضلی استفاده می‌کنیم. در این روش، کاربران با توجه به خط‌مشی حریم خصوصی خود و الگوریتم به کار گرفته شده توسط موتور جستجو برای یادگیری پروفایل کاربر با اضافه کردن تعدادی پرس‌وجوی جستجو یا تغییر پرس‌وجوی جستجوی ارسالی خود پایگاه‌داده مغشوش شده را ایجاد می‌نمایند. بنابراین، هدف ما در این پیشنهاد رساله حفظ حریم خصوصی کاربران در جستجوی وب شخصی‌سازی شده بر اساس روش سنتز پایگاه‌داده مغشوش شده مبتنی بر مدل حریم خصوصی تفاضلی است.

^۱ sort

^۲ query expansion

^۳ synthesize

در بسیاری از موارد، الگوریتم یادگیری پروفایل کاربر و نیز فرایند بازیابی در موتور جستجو (فرایندهای R و L در شکل ۴-۱ و شکل ۴-۲) دانسته نیستند. در این رساله، از آنجایی که چالش‌های دیگری نیز در نگاشت مدل حریم خصوصی تفاضلی به جستجوی وب شخصی‌سازی شده وجود دارد، حل مسئله را در چند گام پیشنهاد می‌کنیم. در هر یک از این گام‌ها، راه‌حلی را برای مسئله با در نظر گرفتن برخی از فرض‌ها در رابطه با دانسته بودن یا نادانسته بودن الگوریتم‌های به کار گرفته شده توسط موتور جستجو برای یادگیری پروفایل کاربران و نیز برای بازیابی اطلاعات ارائه می‌کنیم. وابسته به میزان اطلاعات موجود درباره این الگوریتم‌ها، موتور جستجو جعبه سیاه، جعبه خاکستری، یا جعبه سفید می‌نامیم.

در گام نخست، فرض می‌کنیم چگونگی یادگیری پروفایل کاربر و نیز بازیابی اطلاعات توسط موتور جستجو به صورت کامل مشخص است. به عبارت دیگر، موتور جستجو یک سامانه جعبه سفید است. در این حالت، نگاشت درستی از مدل حریم خصوصی تفاضلی به مسئله حفظ حریم خصوصی در جستجوی وب شخصی‌سازی شده لازم است. برای این کار، باید مفاهیم به کار رفته در مدل حریم خصوصی تفاضلی، مانند همسایگی، حساسیت، سازوکار تصادفی کردن، و تحلیل و اندازه‌گیری سودمندی و اتلاف حریم خصوصی، را به شکل صوری بیان کنیم. همچنین، در جستجوی وب شخصی‌سازی شده، پایگاه داده رشدیابنده است و میان رکوردهای مختلف آن همبستگی وجود دارد. در نگاشت حریم خصوصی تفاضلی به این مسئله، باید به این موضوع نیز توجه داشت.

در گام دوم، فرض می‌کنیم روش بازیابی اطلاعات توسط موتور جستجو به صورت کامل مشخص است. همچنین، فرض می‌کنیم تقریبی از روش یادگیری پروفایل کاربر توسط موتور جستجو در دست است. به عبارت دیگر، موتور جستجو یک سامانه جعبه خاکستری است. در این گام، فرض بر آن است که یک توزیع احتمال روی مقادیر مختلف پارامترهای یک روش یادگیری خاص، مانند k در k -means، داده شده است. آنگاه، با در نظر گرفتن چالش‌های مطرح شده در گام نخست، مدل حریم خصوصی تفاضلی را برای حفظ حریم خصوصی در جستجوی وب شخصی‌سازی به کار می‌بریم.

در گام سوم، موتور جستجو را جعبه سیاه فرض می‌کنیم. بنابراین، روش یادگیری پروفایل کاربر و فرایند بازیابی هر دو نادانسته فرض می‌شوند. در این گام، با بهره‌گیری از فنون یادگیری ماشین، تخمینی از روش به کار گرفته شده توسط موتور جستجو برای یادگیری پروفایل کاربر به دست می‌آوریم. آنگاه، مدل حریم خصوصی تفاضلی را مانند دو گام قبل به مسئله حفظ حریم خصوصی در جستجوی وب شخصی‌سازی شده می‌نگاریم. آنچه در گام سوم انجام خواهد شد، در گام دوم نیز انجام شدنی است. با وجود این، راه‌حل ارائه شده در گام دوم

مبتنی بر فنون یادگیری ماشین نیست، بلکه به صورت مستقیم از نظریه احتمالات بهره می‌جوید. به نظر می‌رسد، این موضوع منجر به محافظه‌کارانه بودن راه‌حل شود. مقایسه این دو راه‌حل، به خصوص به دلیل وجود خطا در به‌کارگیری فنون یادگیری ماشین، بخشی از پژوهش مربوط به موضوع این پیشنهاد رساله خواهد بود.

ماهیت رساله پیشنهادی در گام‌های نخست و دوم حل مسئله نظری است. سازوکارهای حفظ حریم خصوصی ارائه شده در این گام‌ها و میزان سودمندی آن‌ها، به صورت صوری اثبات و راستی‌آزمایی می‌شوند. با وجود این، با توجه به فرض‌های موجود در قسمت ۴-۲-۱ (مدل سامانه شخصی‌سازی نتایج جستجوی کاربران) یک موتور جستجو پیاده‌سازی و سازوکار حریم خصوصی ارائه شده برای آن عملی می‌شود. گام سوم حل مسئله، ماهیتی نظری و پیاده‌سازی دارد. برای این گام، لازم است الگوریتم استخراج پروفایل کاربر از موتور جستجو تعریف و پیاده‌سازی شود. با توجه به پروفایل به‌دست آمده، باید سازوکار حریم خصوصی تفاضلی برای آن تعریف و عملی شود. در این گام نیز، از موتور جستجوی استفاده شده در دو گام پیشین، البته با فرض جعبه‌سیاه بودن الگوریتم‌های آن، استفاده خواهد شد.

۴-۲ شرح مسئله

در این قسمت، مدل سامانه در نظر گرفته شده برای شخصی‌سازی نتایج جستجوی کاربران، مدل مهاجم، و نیز مدل حریم خصوصی تفاضلی شرح داده می‌شوند. به عبارت دیگر، جزئیات مربوط به کارکرد اجزای نشان داده شده در شکل ۴-۱ و شکل ۴-۲، که بازتاب دهنده فرض‌ها و انتخاب‌های ما در مورد این اجزا هستند، توضیح داده می‌شوند. آنگاه، چالش‌های حل مسئله را بیان کرده و گام‌های پیشنهادی مطرح شده در قسمت ۴-۱ را با توجه به مدل سامانه شرح می‌دهیم.

۴-۲-۱ مدل سامانه شخصی‌سازی نتایج جستجوی کاربران

پیش‌تر گفته شد که تاریخچه جستجوهای هر کاربر در یک پایگاه‌داده ذخیره می‌شود. این پایگاه‌داده را با H نشان می‌دهیم. هر رکورد H را یک سه‌تایی به شکل $\langle id, sq, cd \rangle$ در نظر می‌گیریم که در آن id یک شناسه یکتا، sq یک پرس‌وجوی جستجو، و cd مجموعه‌ای از اسندها است که کاربر در پی جستجوی sq به آن‌ها مراجعه نموده است. sq می‌تواند متن، تصویر، یا فیلم باشد. در این پیشنهاد رساله، ما تنها پرس‌وجوهای جستجو

از نوع متن را در نظر خواهیم گرفت. سندهای موجود در *cd* سندهای کلیک شده^۱ نیز نامیده می‌شوند. گفتنی است که رفتار جستجوی کاربر، افزون بر پرس‌وجوهای جستجوی کاربر، شامل مجموعه سندهای کلیک شده، مدت زمان پرداختن به هر سند، الگو و ترتیب مراجعه به سندها، و مانند آن است. با وجود این، در این پیشنهاد رساله، رفتار جستجوی کاربر فقط شامل پرس‌وجوهای جستجو و سندهای کلیک شده در نظر گرفته می‌شود.

فرایند بازیابی، که آن را با R نشان می‌دهیم، پرس‌وجوی جستجو و نمایه را دریافت کرده و لیستی از پیوندها به سندهای مربوط به آن پرس‌وجو را بازمی‌گرداند. این فرایند در گذر زمان بر اساس مدل‌های احتمالاتی و زبانی، با استفاده از فنون یادگیری ماشین، بهبود می‌یابد. در این پیشنهاد رساله، فرض می‌کنیم فرایند R ثابت است. با وجود این، R نادانسته فرض می‌شود. همان‌طور که در قسمت ۴-۲ آمده است، پروفایل کاربر می‌تواند از نوع کلمه کلیدی، شبکه نحوی، یا شبکه مفهومی باشد. همچنین، فرایند ایجاد پروفایل کاربر، که آن را با L نشان می‌دهیم، روش‌های یادگیری ماشین را بر روی تاریخچه جستجوی مربوط به یک کاربر به کار می‌گیرد [۴۹][۴۸].

در گام‌های یکم و دوم حل مسئله، که در قسمت ۴-۱ بیان شده است، فرض می‌شود که پروفایل کاربر از نوع کلمه کلیدی بوده و فرایند L ، فرایند یادگیری پروفایل کاربر، مبتنی بر آماره TF-IDF و الگوریتم *meansk*- است. با این انتخاب، نتیجه فرایند L بردار مرکز خوشه‌ها در الگوریتم *meansk*- است. شرح جزئیات مربوط به الگوریتم یادگیری پروفایل کاربر مبتنی بر آماره TF-IDF و الگوریتم *meansk*- در قسمت ۴-۲ آمده است. موضوع دیگر آن است که ممکن است دو نوع پروفایل یکی بر اساس رفتار جستجوی کاربر در بلندمدت و دیگری مبتنی بر رفتار جستجوی او در کوتاه‌مدت ایجاد شود. در این پیشنهاد رساله، فقط پروفایل مربوط به رفتار جستجوی بلندمدت کاربر در نظر گرفته می‌شود. به عبارت دیگر، بر خلاف برخی از موارد که در آن‌ها فعالیت‌های اخیر کاربر تأثیر بیشتری در پروفایل او دارند، اثر فعالیت‌های اخیر و قدیم کاربر در پروفایل ساخته شده برای کاربر یکسان در نظر گرفته می‌شود. در گام سوم، که در آن موتور جستجو جعبه‌سیاه فرض می‌شود، نوع پروفایل کاربر و فرایند L نادانسته فرض می‌شوند.

شخصی‌سازی نتایج جستجوی کاربر به روش گسترش پرس‌وجوی کاربر یا مرتب‌سازی نتایج جستجو انجام می‌شود. در هر دو روش، شخصی‌سازی بر اساس پروفایل کاربر و پرس‌وجوی جستجوی کاربر انجام می‌شود. بنابراین، می‌توان هر یک از روش‌ها را در نظر گرفت. در این پیشنهاد رساله، روش مرتب‌سازی نتایج (شکل ۴-۱) انتخاب شده است. موضوع دیگر آن است که موتورهای جستجو از اطلاعات زمینه‌ای کاربر، مانند موقعیت

^۱ clicked documents

جغرافیایی و جنسیت، نیز در شخصی سازی نتایج جستجو استفاده می کنند. با وجود این، در این پیشنهاد رساله، بر روی اطلاعات به دست آمده از رفتار جستجوی کاربران تمرکز می نماییم.

۲-۲-۴ مدل مهاجم

در این قسمت، فرض های در نظر گرفته شده درباره توانایی های مهاجم، که در اینجا همان موتور جستجو است، بررسی می شوند. مجموعه این فرض ها مدل مهاجم نامیده می شود. در مسئله مطرح شده در این پیشنهاد رساله، موتور جستجو تمامی پایگاه داده را در اختیار دارد و در نقش تحلیل گر، الگوریتم یادگیری پروفایل کاربر را بر روی آن اجرا می کند. همچنین، موتور جستجو درستکار ولی کنجکاو در نظر گرفته شده است. به عبارت دیگر، اگرچه موتور جستجو ممکن است از پروفایل کاربر اطلاعاتی به دست آورد که حریم خصوصی کاربر را نقض کند، اما هیچگاه بر خلاف پروتکل توافق شده با کاربران عمل نکرده و اقدامات مخرب انجام نخواهد داد.

در جستجوی وب شخصی سازی شده، موتور جستجو اطلاعات مربوط به رفتار جستجوی هر کاربر را به صورت جداگانه ضبط می نماید. بنابراین، پایگاه داده دربرگیرنده اطلاعات مربوط به یک کاربر خاص فرض می شود. در نتیجه، میان رکوردهای پایگاه داده همبستگی وجود دارد. این موضوع به تحلیل گر در استخراج اطلاعات بیشتر درباره هر کاربر کمک می کند. گفتنی است که در طول تعامل هر کاربر با موتور جستجو، رکوردهای جدیدی به پایگاه داده اضافه می شوند. به عبارت دیگر، پایگاه داده رشدیابنده است. با اضافه شدن رکوردهای جدید، موتور جستجو پروفایل دقیق تری از کاربر ایجاد کرده و به طور دقیق تر می تواند از پسندهای حساس کاربران آگاهی پیدا کند. در مسئله مطرح شده در این پیشنهاد رساله، در صورتی که کاربران به درستی و متناسب با الگوریتم یادگیری پروفایل کاربر اطلاعات خود را مغشوش نکنند، مهاجمان می توانند، برخلاف خط مشی کاربران، اطلاعات حساس آن ها را به دست آورند. با وجود این، موتور جستجو یک سامانه جعبه سیاه بوده و کاربران روش به کار گرفته شده برای بازیابی و یادگیری پروفایل کاربر را نمی دانند.

به طور خلاصه، توانایی موتور جستجو (مهاجم) در به دست آوردن پسندهای حساس کاربران به ویژگی هایی زیر از جستجوی وب شخصی سازی شده وابسته است.

- موتور جستجو پایگاه داده (تاریخچه جستجوهای یک کاربر) را به طور کامل در دست دارد.
- رکوردهای موجود در یک پایگاه داده همبسته هستند.
- پایگاه داده رشدیابنده است و در طول زمان رکوردهای جدید به آن اضافه می شود.

- کاربر از فرایند ایجاد پروفایل کاربر و فرایند بازیابی اطلاعاتی ندارد و از دید او موتور جستجو یک سامانه جعبه سیاه است.

۳-۲-۴ حریم خصوصی تفاضلی و سنتز تاریخچه جستجوها

در این قسمت، ابتدا با یک مثال نشان می‌دهیم چگونه با سنتز پایگاه داده مغشوش شده می‌توان حریم خصوصی کاربران را مبتنی بر مدل حریم خصوصی تفاضلی حفظ نمود. در این مثال، پرس‌وجوی تحلیل‌گر بر روی پایگاه داده میانگین مقدارهای عددی است که در پایگاه داده ذخیره شده‌اند. همچنین، پایگاه داده ثابت است و به ازای هر کاربر فقط یک مقدار عددی (رکورد) در این پایگاه داده قرار دارد. فرض دیگر آن است که رکوردها همبسته نیستند. آنگاه، دشواری‌های به‌کارگیری این رویکرد را در حفظ حریم خصوصی در جستجوی وب شخصی‌سازی شده مطرح می‌کنیم. این دشواری‌ها ناشی از عدم وجود تعریف خط‌مشی حریم خصوصی بر اساس مدل حریم خصوصی تفاضلی در جستجوی وب شخصی‌سازی شده و نیز درست نبودن فرض‌هایی است که رویکرد ارائه شده وابسته به آن‌ها است.

کاربران، که تعداد آن‌ها n فرض می‌شود، اطلاعات حساس خود را، که یک عدد طبیعی است، به مسئول پایگاه داده اعلام می‌کنند. بنابراین، پایگاه داده $H \in \mathbb{N}^n$ ایجاد می‌شود. مسئول پایگاه داده، پایگاه داده مغشوش شده \hat{H} را متناسب با پرس‌جوی تحلیل‌گر (تابع میانگین) و با اضافه کردن تعدادی رکورد پوششی، که با $CR \in \mathbb{N}^k$ نمایش داده می‌شود، سنتز و منتشر می‌کند. بنابراین، $\hat{H} = H \cdot CR$ که در آن « \cdot » عملگر الحاق است، منتشر خواهد شد. همان طور که گفته شد، تحلیل‌گر تابع میانگین $L: \mathbb{N}^{n+k} \rightarrow \mathbb{R}$ را روی \hat{H} اجرا خواهد کرد. سازوکار سنتز پایگاه داده مغشوش شده را می‌توان با تابع $M: \mathbb{N}^n \rightarrow \Delta(\mathbb{N}^{n+k})$ بازنمود. این تابع پایگاه داده H متشکل از n عدد طبیعی را به توزیع احتمال $M(H)$ روی پایگاه داده‌های متشکل از $n+k$ عدد طبیعی می‌نگارد. به عبارت دیگر، M به تصادف یک پایگاه داده بازمی‌گرداند. احتمال اینکه M پایگاه داده سنتز شده \hat{H} را بازگرداند $M(H)(\hat{H})$ است. سازوکار M ، ϵ -خصوصی تفاضلی است اگر برای هر دو پایگاه داده همسایه $H, H' \in \mathbb{N}^n$ و هر $r \in \mathbb{R}$ رابطه زیر برقرار باشد:

$$\left(\sum_{\{\hat{H} \in \mathbb{N}^{n+k} \mid L(\hat{H})=r\}} M(H)(\hat{H}) \right) \leq e^\epsilon \left(\sum_{\{\hat{H} \in \mathbb{N}^{n+k} \mid L(\hat{H})=r\}} M(H')(\hat{H}) \right).$$

مسئول پایگاه داده می‌تواند از سازوکاری مبتنی بر سازوکار لاپلاس، در جهت سنتز پایگاه داده مغشوش شده استفاده کند. در سازوکار لاپلاس، میانگین رکوردها با مقداری نویز از توزیع لاپلاس جمع شده و مقدار میانگین

مغشوش شده به دست می‌آید. روش محاسبه میانگین مغشوش شده در زیر آمده است. R متغیری تصادفی نمایان گر مقادیری مختلف برای میانگین پایگاه‌داده‌های $\hat{H} \in \mathbb{N}^{n+k}$ است. میزان حساسیت تابع L با S_L نشان داده شده است. گفتنی است، روش محاسبه حساسیت یک تابع در قسمت ۲-۳ توضیح داده شده است.

$$\eta \sim \text{Lap}\left(0, \frac{S_L}{\epsilon}\right)$$

$$R = L(H) + \eta$$

به عبارت دیگر، R از توزیع لاپلاس زیر به دست می‌آید:

$$R \sim \text{Lap}\left(L(H), \frac{S_L}{\epsilon}\right)$$

در سازوکار M ، مسئول پایگاه‌داده به صورت تصادفی از توزیع احتمال $\text{Lap}\left(L(H), \frac{S_L}{\epsilon}\right)$ مقدار r را انتخاب می‌کند. احتمال انتخاب r ، $f_R(r)$ است که به صورت زیر محاسبه می‌شود:

$$f_R(r) = \frac{\epsilon}{2S_L} e^{\left(-\frac{|r-L(H)|}{\frac{S_L}{\epsilon}}\right)}$$

بعد از آن، از فضای نمونه پایگاه‌داده‌های $\hat{H} \in \mathbb{N}^{n+k}$ ، پایگاه‌داده‌هایی با میانگین r را جدا کرده و بر اساس توزیع احتمال $\Gamma(\mathbb{N}^{n+k})$ ، یکی از آن‌ها را به تصادف انتخاب می‌کند. متغیر تصادفی \mathbb{H} نمایان گر پایگاه‌داده‌های $\hat{H} \in \mathbb{N}^{n+k}$ با میانگین r است که احتمال رخداد آن‌ها از توزیع Γ به دست می‌آید. توزیع Γ ممکن است توزیعی یکنواخت، توزیعی نورمال، یا هر توزیع دیگری باشد. بنابراین، احتمال انتخاب پایگاه‌داده \hat{H} از این توزیع، $f_{\mathbb{H}}(\hat{H})$ است. در نتیجه، احتمال انتخاب پایگاه‌داده \hat{H} با میانگین r به صورت $f_R(r) \times f_{\mathbb{H}}(\hat{H})$ است. بنابراین، احتمال سنتز پایگاه‌داده‌های \hat{H} با میانگین r به شکل زیر محاسبه می‌شود:

$$\begin{aligned} & \sum_{\{\hat{H} \in \mathbb{N}^{n+k} \mid L(\hat{H})=r\}} M(H)(\hat{H}) \\ &= \sum_{\{\hat{H} \in \mathbb{N}^{n+k} \mid L(\hat{H})=r\}} f_R(r) \times f_{\mathbb{H}}(\hat{H}) \\ &= f_R(r) \sum_{\{\hat{H} \in \mathbb{N}^{n+k} \mid L(\hat{H})=r\}} f_{\mathbb{H}}(\hat{H}) \\ &= f_R(r) \times 1 = f_R(r) \end{aligned}$$

نسبت احتمال سنتز پایگاه‌داده‌های \hat{H} به ازای هر مقدار میانگین $r \in \mathbb{R}$ از پایگاه‌داده‌های همسایه $H, H' \in \mathbb{N}^n$ به صورت زیر محاسبه می‌شود:

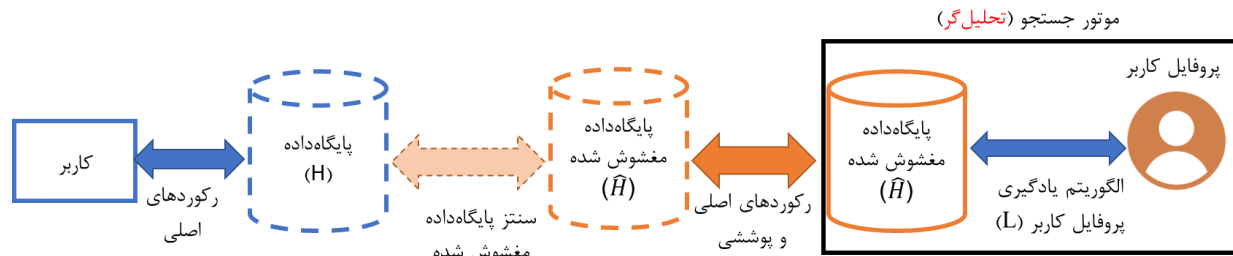
$$\frac{\sum_{\{\hat{H} \in \mathbb{N}^{n+k} \mid L(\hat{H})=r\}} M(H)(\hat{H})}{\sum_{\{\hat{H} \in \mathbb{N}^{n+k} \mid L(\hat{H})=r\}} M(H')(\hat{H})} = \frac{f_R(r)}{f_{R'}(r)} = e^{\left(-\frac{|r-L(H)|}{\frac{S_L}{\epsilon}} + \frac{|r-L(H')|}{\frac{S_L}{\epsilon}}\right)} \leq e^{\left(\frac{|L(H)-L(H')|}{\frac{S_L}{\epsilon}}\right)} \leq e^{\epsilon}$$

بنابراین، سازوکار M ، ϵ -خصوصی تفاضلی است. برای محاسبه سودمندی سازوکار ارائه شده، امید ریاضی تفاضل میانگین پایگاه داده‌های \hat{H} و H را به صورت زیر محاسبه می‌کنیم:

$$E[L(\hat{H}) - L(H)] = E[L(\hat{H})] - E[L(H)] = E[R] - L(H) = L(H) - L(H) = 0$$

همان طور که دیده می‌شود، مقادیری که تحلیل گر به عنوان $L(\hat{H})$ محاسبه می‌کند، به صورت میانگین با $L(H)$ برابر است. به عبارت دیگر، سازوکار ارائه شده بیشترین سودمندی را دارد. در مثال ارائه شده، سازوکاری برای سنتز پایگاه داده \hat{H} بر اساس مدل حریم خصوصی تفاضلی بیان شد. حال با توجه به این مثال، جزئیات مربوط به مسئله مطرح در این پیشنهاد رساله را شرح می‌دهیم.

در شکل ۳-۴، مدل مفهومی حفظ حریم خصوصی در جستجوی وب شخصی سازی شده بر اساس سنتز پایگاه داده مغشوش شده مبتنی بر حریم خصوصی تفاضلی آمده است. در این شکل، قسمت‌های مغشوش شده با رنگ نارنجی و موجودیت‌های فرضی با خط چین نشان داده شده‌اند. در این مدل مفهومی، پایگاه داده H شامل تاریخچه جستجوهای یک کاربر است. الگوریتم اجرا شده تحلیل گر (موتور جستجو) بر روی H ، که با L در شکل ۴-۱ و شکل ۲-۴ نشان داده شده است، الگوریتم یادگیری پروفایل کاربر است. از آنجایی که این پروفایل خود مبتنی بر TF-IDF و آماره TF-IDF فرض شده است، H پایگاه داده‌ای آماری خواهد بود. جزئیات الگوریتم یادگیری پروفایل کاربر مبتنی بر TF-IDF و آماره TF-IDF در قسمت ۴-۲ شرح داده شده است. برای عملی کردن خط‌مشی حریم خصوصی کاربر، دو راهکار وجود دارد. در راهکار اول، متناسب با الگوریتم یادگیری پروفایل کاربر، خط‌مشی حریم خصوصی او، و آنتولوژی محاسباتی موضوع‌ها تعدادی رکورد پوششی شامل پرس‌وجوی جستجو و سندهای کلیک شده در پی آن ایجاد می‌شوند. در راهکار دوم، با توجه به الگوریتم یادگیری پروفایل کاربر، خط‌مشی حریم خصوصی او، و آنتولوژی محاسباتی با اضافه شدن رشته‌ای از کلمه‌ها پرس‌وجوهای جستجوی کاربر گسترش پیدا می‌کنند. همچنین، تعدادی از سندهای نتیجه جستجو نیز کلیک می‌شوند. به عبارت دیگر، از روش گسترش پرس‌وجو استفاده می‌شود. در رساله، ما راهکار اول را انتخاب کرده‌ایم. در این راهکار، برخلاف مراجع [۱۲] و [۱۴]، بر اساس خط‌مشی حریم خصوصی کاربر و تنها برای پرس‌وجوهای مرتبط با موضوع‌های حساس او، رکوردهای پوششی ایجاد می‌شوند. همچنین، برخلاف مراجع [۹]، [۱۰]، و [۵] الگوریتم یادگیری پروفایل کاربر نیز در نظر گرفته می‌شوند. این کار، باعث افزایش سودمندی و حفظ حریم خصوصی می‌شود. در نتیجه، به جای پایگاه داده H ، پایگاه داده مغشوش شده \hat{H} در اختیار موتور



شکل ۴-۳: مدل مفهومی حفظ حریم خصوصی در جستجوی وب شخصی سازی شده بر اساس سنتز پایگاه داده مغشوش شده مبتنی بر حریم خصوصی تفاضلی

جستجو قرار می گیرد. برای سنتز پایگاه داده \bar{H} ، کاربر با ایجاد رکوردهای پوششی (CR)، شامل پرس و جویهای پوششی و سندهای کلیک شده در پی آنها (راهکار اول)، پایگاه داده تاریخچه جستجوها را مغشوش می نماید. در مدل حریم خصوصی تفاضلی، خط مشی حریم خصوصی با استفاده از پارامتر ϵ تعیین می شود. هرچه مقدار ϵ کمتر باشد، حریم خصوصی کاربر بیشتر حفظ می شود. در جستجوی وب شخصی سازی شده، کاربران موضوع های حساس خود را همراه با میزان حفظ حریم خصوصی برای هر کدام اعلام می کنند. موضوع مهم در این رابطه این است که کاربران انتظار دارند، از خدمت شخصی سازی نتایج جستجو در کنار حفظ حریم خصوصی استفاده کنند. به عبارت دیگر، پرسش مهمی که باید پاسخ داده شود این است که چگونه می توان هم موضوع های کاربر را برای حفظ حریم خصوصی او پنهان نمود و هم از خدمت شخصی سازی نتایج استفاده کرد. برای پاسخ به این پرسش از یک مثال استفاده می کنیم. فرض کنیم، موضوع حساس کاربر تیم فوتبال مورد علاقه او است. سازوکاری که برای پنهان سازی تیم فوتبال کاربر پیاده سازی می شود به شکلی است که موتور جستجو تیم فوتبال مورد علاقه دقیق کاربر را شناسایی نمی کند. با وجود این، موتور جستجو از علاقمندی کاربر به فوتبال آگاهی پیدا می کند. آنگاه، موتور جستجو می تواند با توجه به علاقمندی کاربر به فوتبال، بدون اینکه به طور دقیق بداند که به کدام تیم فوتبال علاقمند است، نتایج را شخصی سازی کند.

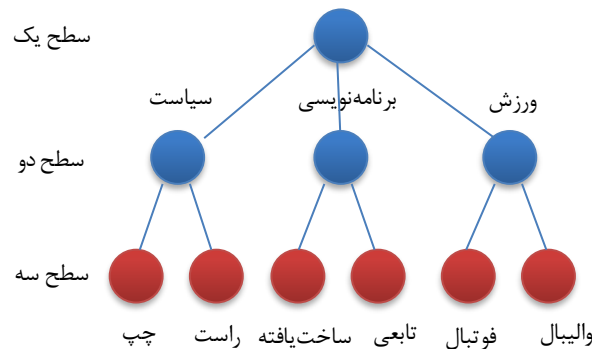
موضوع های مختلف، مانند فوتبال و تیم های فوتبال، با هم ارتباط دارند و ممکن است خط مشی حریم خصوصی کاربران نسبت به هر کدام از آنها متفاوت باشد. بنابراین، در عملی کردن خط مشی حریم خصوصی کاربران از ارتباط موضوع های مختلف استفاده می شود. به عنوان مثال، پرس و جویهای پوششی در رابطه با تیم های فوتبال دیگر ایجاد می شود. برای این کار، از آنتولوژی محاسباتی^۱ مربوط به موضوع های مختلف استفاده خواهیم کرد.

^۱ computational ontology: به فصل ۲ قسمت ۲-۵ مراجعه شود.

کاربران در خط‌مشی حریم خصوصی، موضوع‌های حساس خود را همراه با پارامتر ϵ برای آن موضوع اعلام می‌کنند. آنتولوژی محاسباتی موضوع‌ها به شکل یک درخت با ساختار $T(t, \mathbb{T})$ تعریف می‌شود. در این ساختار، t یک موضوع و \mathbb{T} مجموعه زیردرخت‌های T است. هر عضو از \mathbb{T} نیز با ساختار $T'(t', \mathbb{T}')$ تعریف می‌شود که T' هم یک درخت از موضوع‌ها است. در این پیشنهاد رساله، مجموعه تمام موضوع‌های موجود در T را با $topic_T$ نشان می‌دهیم. علاوه بر آن، موضوع‌های موجود در سطح i از T را با مجموعه T_i و موضوع‌های فرزند موضوع t با مجموعه $child(t)$ نمایش داده می‌شوند. در این پیشنهاد رساله، خط‌مشی حریم خصوصی کاربر با توجه به آنتولوژی محاسباتی موضوع‌های T ، در قالب یک تابع بیان می‌شود. دامنه این تابع، موضوع‌های موجود در T ($t \in topic_T$) و برد آن، $\mathbb{R}^{\geq 0}$ است. بنابراین، تابع خط‌مشی حریم خصوصی کاربر است. گفتنی است، برای هر موضوع t موجود در هر سطح i از درخت آنتولوژی محاسباتی T ، رابطه $P(t_i) \geq P(t')$ به ازای تمام t' های عضو $child(t_i)$ باید برقرار باشد. در شکل ۴-۴، به عنوان نمونه بخشی از آنتولوژی محاسباتی موضوع‌ها آمده است. گفتنی است، آنتولوژی محاسباتی T که بازتاب درستی از موضوع‌های ممکن و موجود است را باید انتخاب کنیم.

برای مثال یک خط‌مشی P به شکل زیر تعریف می‌شود. در این خط‌مشی، برای تمام موضوع‌های t_3 موجود در مجموعه T_3 ($t_3 \in T_3$)، رابطه $0 \leq P(t_3) = \epsilon_1 < \infty$ برقرار است. همچنین، برای هر $t_2 \in T_2$ و هر $t_1 \in T_1$ ، $P(t_2) = P(t_1) = \infty$ و برای موضوع‌های با سطح بزرگتر از دو ($t_{>2} \in T_{>2}$)، $P(t_{>2}) = 0$ فرض می‌شود. بنابراین در شکل ۴-۴، با توجه به خط‌مشی گفته شده، موضوع‌های حساس او با رنگ قرمز نشان داده شده‌اند. برای مثال، رابطه‌های زیر برقرار است.

$$\begin{aligned} P(\text{فوتبال}) &= P(\text{والیبال}) = \epsilon_1 \\ P(\text{سیاست}) &= P(\text{ورزش}) = \infty \\ P(\text{استقلال}) &= P(\text{پرسپولیس}) = 0 \end{aligned}$$



شکل ۴-۴: بخشی از آنتولوژی محاسباتی موضوع‌ها.

با توجه به خط‌مشی حریم خصوصی P ، کاربر نمی‌خواهد موتور جستجو تیم فوتبال مورد علاقه او (پرسپولیس یا استقلال) را تشخیص دهد. با وجود این، اشکالی ندارد که علاقه او به فوتبال را بفهمد. این تحلیل، برای موضوع‌های «سیاست» و «برنامه‌نویسی» نیز درست است.

پیش‌تر گفته شد که تمام رکوردهای موجود در پایگاه‌داده H مربوط به یک کاربر خاص هستند. بنابراین، رکوردهای موجود در پایگاه‌داده همبستگی دارند. ارتباط میان رکوردهای پایگاه‌داده را با دو ایده متفاوت می‌توان بررسی کرد. در ایده اول، تعریف جدیدی برای همسایگی دو پایگاه‌داده باید ارائه شود. برای مثال، دو پایگاه‌داده همسایه هستند اگر در حضور یا عدم حضور رکوردهای موضوع t با هم متفاوت باشند. متناسب با تعریف همسایگی باید حساسیت الگوریتم یادگیری پروفایل کاربر نیز محاسبه و سازوکار تصادفی کردن برای آن تعریف شود. در ایده دوم، از ایده‌های مطرح در مرجع‌های [۴۷] و [۶۵]، می‌توان سازوکار تصادفی کردن را با وجود همبستگی در پایگاه‌داده طراحی کرد. در این ایده، همسایگی دو تاریخچه جستجو در متفاوت بودن یک رکورد در دو پایگاه‌داده تعریف می‌شود. در طراحی این سازوکار، باید مفهوم همبستگی رکوردها در پایگاه‌داده تاریخچه جستجوهای کاربر به درستی صوری شود. به عبارت دیگر، رابطه رکوردها با هم باید تعریف شده و همبستگی آن‌ها باید بر اساس معیار مشخصی کمی شود. علاوه بر آن، حساسیت الگوریتم یادگیری پروفایل کاربر با توجه به وجود همبستگی (حساسیت همبسته) باید محاسبه شود. در قسمت ۳-۲، ایده‌های مطرح در مراجع [۴۷] و [۶۵] به همراه نحوه محاسبه حساسیت همبسته شرح داده شده است. انتخاب هر کدام از این ایده‌ها، با چالش‌هایی روبه‌رو است که بهترین راه‌حل برای مسئله مطرح در این پیشنهاد رساله باید انتخاب شود.

همان‌طور که گفته شد، در مسئله مطرح در این پیشنهاد رساله پایگاه‌داده رشدیابنده است و موتور جستجو بعد از دریافت هر رکورد جدید، پروفایل کاربر را مجدداً محاسبه می‌کند. با رشد پایگاه‌داده، پروفایل ایجاد شده دقیق‌تر می‌شود. بنابراین، با توجه به آنچه که موتور جستجو از کاربر می‌داند، در بازه‌های زمانی مشخص از رشد پایگاه‌داده باید سازوکار سنتز پایگاه‌داده مغشوش شده را اجرا کرد. پیش از این گفته شد، کاربر در خط‌مشی حریم خصوصی خود برای هر موضوع مقداری برای ϵ مشخص می‌کند. با توجه به ایده‌های مطرح در مرجع‌های [۲۱] و [۶۱]، همسایگی در قالب همسایگی دو جریان پایگاه‌داده^۱ تعریف می‌شود. دو جریان پایگاه‌داده X و X' با هم همسایه هستند، اگر این دو پایگاه‌داده رشدیابنده، فقط در یک رکورد با هم تفاوت داشته باشند. بنابراین، اگر سازوکار حفظ حریم خصوصی در گام i ، که با M_i نشان داده می‌شود، به اندازه ϵ_i خصوصی تفضلی باشد، رابطه $\sum_0^i \epsilon_i < \epsilon$ طبق قضیه ترکیب ترتیبی باید به ازای هر i برقرار باشد. چالش اصلی تعیین مقدار مناسب ϵ_i و زمان اجرا برای هر گام است، به شکلی که بیشترین سودمندی را برای کاربر داشته باشد.

^۱ Database Stream

گفتنی است، در مسئله مطرح در این پیشنهاد رساله کاربران به طور مستقیم از نتیجه اجرای فرایند ایجاد پروفایل کاربر (فرایند L) روی پایگاه داده استفاده می کنند. به عبارت دیگر، نتیجه فرایند L منجر به تولید پروفایل کاربر می شود که موتور جستجو از آن در جهت شخصی سازی نتایج استفاده می کند. بنابراین، عملی کردن سازوکار M به طور مستقیم بر سودمندی نتایج جستجو برای کاربر موثر است. در نتیجه، معیاری برای اندازه گیری سودمندی M باید ارائه شود. این معیار، نشان دهنده میزان مفید بودن پیوندها به سندهایی است که موتور جستجو به عنوان نتیجه پردازش یک پرس و جوی جستجو به کاربر بر می گرداند. بر اساس این معیار، سازوکاری بهینه در حفظ حریم خصوصی کاربر باید ارائه شود.

با توجه به ویژگی های مطرح شده مسئله، بعد از حل چالش های مربوط به همبستگی و رشدیابنده بودن پایگاه داده، می توان پروفایل مغشوش شده را متناسب با خط مشی کاربر محاسبه کرد. پروفایل مغشوش شده با اضافه کردن نویز به مرکز خوشه های به دست آمده از الگوریتم $meansk$ حاصل می شود. با در نظر گرفتن پروفایل مغشوش شده، می توان پایگاه داده مغشوش شده H را سنتز کرد. گفتنی است، الگوریتمی برای ایجاد رکوردهای پوششی (CR) با در نظر گرفتن پروفایل مغشوش شده باید ارائه شود. همان طور که دیده می شود، برای عملی کردن این سازوکار نیاز است تا کاربر الگوریتم یادگیری پروفایل را بداند. پیش تر گفته شد که موتور جستجو از دید کاربر یک سامانه جعبه سیاه است. بنابراین، با توجه به وجود چالش های گوناگون در حل مسئله، سازوکار سنتز پایگاه داده مغشوش شده در سه گام و با در نظر گرفتن جعبه سفید، جعبه خاکستری، و جعبه سیاه بودن موتور جستجو بررسی خواهد شد.

۴-۲-۴ چالش های حل مسئله

با توجه به مدل حریم خصوصی تفاضلی محلی مطرح شده در قسمت ۴-۲-۳، چالش های زیر برای حل مسئله وجود خواهند داشت.

- کاربر بر اساس آنتولوژی محاسباتی T ، خط مشی حریم خصوصی P را اعلام می کند. آنتولوژی محاسباتی T مناسب برای مسئله مطرح در این پیشنهاد رساله را باید انتخاب کنیم.
- پیش تر گفته شد که همبستگی میان رکوردهای پایگاه داده باید در نظر گرفته شوند. در یک ایده، گفته شد که باید تعریف جدیدی از همسایگی ارائه شود. در ایده دیگر، با توجه به تئوری های موجود درباره همبستگی، الگوریتم محاسبه همبستگی میان رکوردهای پایگاه داده H باید ارائه شود. در صورت انتخاب

هرکدام از ایده‌های گفته شده، حساسیت الگوریتم یادگیری پروفایل کاربر و سازوکار تصادفی کردن متناظر با هر ایده باید بررسی و صوری شود.

- برای انتخاب سازوکاری برای سنتز پایگاه داده مغشوش شده \hat{H} ، باید معیاری برای اندازه‌گیری سودمندی کاربر ارائه شود. همچنین، سازوکار نهایی باید برای همه حالت‌های ممکن برای خط‌مشی حریم خصوصی کاربر درست باشد. به عبارت دیگر، حالت‌های مختلف برای موضوع‌های حساس کاربر در سطح‌های گوناگونی از آنتولوژی محاسباتی موضوع‌ها و همچنین، مقدارهای متفاوت برای پارامتر ϵ مربوط به هر موضوع باید بررسی و صوری شوند.
- در سنتز پایگاه داده مغشوش شده \hat{H} ، فرض کردیم الگوریتم تولید رکوردهای پوششی از پروفایل مغشوش شده را در دست داریم. این الگوریتم باید بررسی و ارائه شود.
- درباره چالش رشدیابنده بودن پایگاه داده، باید اندازه مناسب برای ϵ_i در هر گام i و زمان اجرای هر گام در رشد پایگاه داده تعیین شود. همچنین، اثر وجود موضوع‌ها با اندازه‌های ϵ مختلف، نیز باید بررسی شود.
- سازوکار ارائه شده، مرتبط با نوع الگوریتم ایجاد پروفایل کاربر است. با وجود این، موتور جستجو از دید کاربر یک سامانه جعبه سیاه است. راه حلی برای چالش جعبه سیاه بودن موتور جستجو باید ارائه شود.

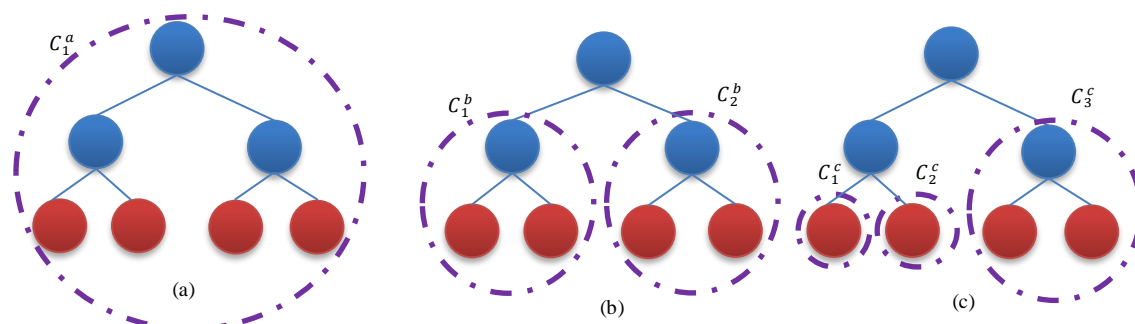
۴-۲-۵ گام نخست در حل مسئله با فرض جعبه سفید بودن موتور جستجو

در این قسمت، با فرض مشخص بودن فرایند ایجاد پروفایل کاربر (جعبه سفید بودن موتور جستجو)، شهودی از سازوکار سنتز پایگاه داده مغشوش شده \hat{H} ، را با در نظر گرفتن مثالی از خط‌مشی P (قسمت ۳-۲-۴) شرح می‌دهیم. با توجه به مدل سامانه تعریف شده در قسمت ۱-۲-۴، ما از الگوریتم $meansk$ - با تعداد خوشه مشخص k در جهت ایجاد پروفایل کاربر استفاده می‌کنیم. برای توضیح سازوکار M ، موضوع‌های جستجو شده کاربر در لحظه عملی کردن آن با توجه به آنتولوژی محاسباتی T و خط‌مشی P ، در شکل ۵-۴ و شکل ۶-۴ در قالب درختی دارای رنگ نمایش داده شده‌اند. با توجه به خط‌مشی P ، موضوع‌های حساس ($P(t) < \infty$) با دایره‌های قرمز رنگ نشان داده شده‌اند. با توجه به اندازه k در الگوریتم $meansk$ -، ممکن است الگوریتم یادگیری پروفایل کاربر یکی از حالت‌های سه گانه در شکل ۵-۴ را در جهت ایجاد پروفایل به کار بگیرد. حالت‌های دیگر که شامل تعداد متفاوتی برای k و موضوع‌های جستجو شده است، باید بررسی شوند. در شکل ۵-۴، وضعیت خوشه‌ها با دایره‌های خط‌چین و بنفش در این سه حالت نشان داده شده است. در سازوکار سنتز پایگاه داده

مغشوش شده \bar{H} ، هدف ما تولید پایگاه داده‌ای است که اجرای الگوریتم meansk - روی آن منجر به تولید پروفایل مغشوش شده کاربر می‌شود. گفتنی است، پروفایل مغشوش شده کاربر در اثر اضافه کردن نویز به مرکز خوشه‌های به دست آمده از الگوریتم meansk - حاصل می‌شود. بنابراین، توزیع بردارها در هر خوشه باید متناسب با خط‌مشی حریم خصوصی کاربر مغشوش گردد.

در حالت (a)، با توجه به آنتولوژی محاسباتی موضوع‌ها (شکل ۴-۴)، خوشه C_1^a شامل بیش از یک موضوع از سطح دو (موضوع‌های غیر حساس) است. در حالتی که فراوانی بردارهای مربوط به یک موضوع حساس بسیار بیشتر از سایر موضوع‌ها باشد، مرکز خوشه C_1^a ، نمایان‌گر علاقه کاربر به آن موضوع حساس است. برای سنتز پایگاه داده مغشوش شده \bar{H} ، از ایده سازوکار مطرح شده در قسمت ۳-۲-۴ استفاده می‌کنیم. در این حالت، با در نظر گرفتن الگوریتم meansk - با اضافه کردن مقداری نویز، متناسب با خط‌مشی P ، مرکز خوشه C_1^a را مغشوش می‌کنیم. به عبارت دیگر، پروفایل مغشوش شده را محاسبه می‌کنیم. سپس، با استفاده از مرکز خوشه مغشوش شده تعدادی رکورد پوششی به CR اضافه می‌کنیم و پایگاه داده \bar{H} را سنتز می‌کنیم. برای مثال، فرض کنید خوشه C_1^a شامل موضوع‌های سیاست و ورزش (از سطح دو در شکل ۴-۴) باشد. همچنین، فراوانی رکوردهای جستجو شده کاربر درباره موضوع سیاست و گرایش چپ بسیار بیشتر از دیگر موضوع‌ها باشد. در این شرایط مرکز خوشه به موضوع گرایش چپ سیاسی بسیار نزدیک است. یکی از حالت‌هایی که مرکز این خوشه بعد از مغشوش شدن ممکن است به آن متمایل شود، موضوع‌های ورزش یا گرایش راست است. کاربر با اضافه کردن تعدادی رکورد در رابطه با این موضوع‌ها می‌تواند پایگاه داده را مغشوش کند. به طور کلی، تعیین تعداد و موضوع رکوردهای پوششی موجود در CR با توجه به پروفایل مغشوش شده چالشی مهم است که الگوریتمی برای آن باید طراحی شود.

در حالت (b)، با توجه به آنتولوژی محاسباتی موضوع‌ها (شکل ۴-۴) و خط‌مشی حریم خصوصی کاربر، در



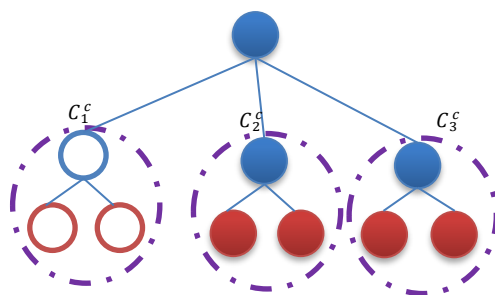
شکل ۴-۵: وضعیت خوشه‌ها (دایره‌های بنفش) و موضوع‌ها (دایره‌های آبی و قرمز) در سه حالت (a) $k=1$ ، (b) $k=2$ و (c) $k=3$.

هر خوشه تنها یک موضوع از سطح دو (موضوع‌های غیر حساس) وجود دارد. بنابراین، در خوشه‌های c_1^b و c_2^b ، توزیع بردارها، با توجه به خط‌مشی حریم خصوصی کاربر، باید به گونه‌ای باشند که بردار مرکز هر کدام از خوشه‌ها، افشا کننده پسند کاربر نباشند. مانند آنچه که در حالت (a) توضیح داده شد، بردار مرکز هر کدام از خوشه‌ها را باید مغشوش کرد.

در حالت (c)، با توجه به آنتولوژی محاسباتی موضوع‌ها (شکل ۴-۴) و خط‌مشی P ، در خوشه‌های c_1^c و c_2^c موضوعی از سطح دو (موضوع‌های غیر حساس) وجود ندارد و این خوشه‌ها، شامل موضوع‌های سطح سه (موضوع‌های حساس) است. بنابراین، برای این خوشه‌ها، نمی‌توان از ایده حالت (a) استفاده کرد. زیرا، مرکز این خوشه‌ها به روشنی نمایان‌گر علاقه کاربر به موضوع‌های حساس هستند. برای این حالت، مرکز خوشه‌های c_1^c و c_2^c در شکل ۴-۵، باید به شکلی مغشوش شوند که شکل خوشه‌ها و موضوع‌ها مانند خوشه‌های موجود در شکل ۴-۶ شوند. در شکل ۴-۶، موضوع‌های حساس جدید اضافه شده با دایره‌های سفید و حاشیه قرمز و موضوع‌های غیر حساس با دایره‌های سفید و حاشیه آبی نمایش داده شده است. چالش مهم برای حالت (c)، ارائه الگوریتمی است که در آن مرکز خوشه c_1^c در شکل ۴-۵، به شکلی مغشوش می‌شود که تعدادی رکورد درباره موضوع‌های جدید اضافه خواهند شد و خوشه c_1^c در شکل ۴-۶ ایجاد می‌شود. همچنین، رکوردهای موجود در خوشه‌های c_2^c و c_1^c در شکل ۴-۵، در خوشه c_2^c در شکل ۴-۶ ادغام می‌شوند.

۴-۲-۶ گام دوم در حل مسئله با فرض جعبه‌خاکستری بودن موتور جستجو

در گام دوم، موتور جستجو به صورت جعبه‌خاکستری در نظر گرفته می‌شود. همان طور که گفته شد، در این حالت یک توزیع احتمال روی مقادیر مختلف از یک یا چند پارامتر از یک روش خاص ایجاد پروفایل کاربر مشخص است. با توجه به مدل سامانه تعریف شده در قسمت ۴-۲-۱، ما از الگوریتم $meansk$ - در جهت ایجاد پروفایل کاربر استفاده می‌کنیم. بنابراین در این قسمت، یک توزیع احتمال روی تعداد خوشه‌ها در الگوریتم k -



شکل ۴-۶: وضعیت موضوع‌ها و خوشه‌ها در حالت (c) بعد از مغشوش کردن پایگاه داده.

means مشخص است. پیش تر تعداد خوشه‌ها در این الگوریتم با k نمایش داده شد. فرض می‌کنیم، در الگوریتم k -means با احتمال برابر $\frac{1}{3}$ ، k ممکن است مقدارهای یک، دو، و سه داشته باشد. با توجه به خط‌مشی P که در قسمت ۳-۲-۴، توضیح داده شد و موضوع‌های جستجو شده کاربر که پیش از این با ساختار درخت نمایش داده شده‌اند، در شکل ۵-۴، وضعیت خوشه‌ها (دایره‌های بنفش) در این سه حالت آمده‌اند. در این قسمت نیز مانند گام نخست شهودی از سازوکار سنتز پایگاه‌داده مغشوش شده \hat{H} را بیان خواهیم کرد.

سازوکار تصادفی کردن برای حالت‌های (a)، (b)، و (c) به صورت مستقل در گام نخست شرح داده شده‌اند. در این گام، یک توزیع احتمال روی k نیز وجود دارد. سازوکار سنتز پایگاه‌داده مغشوش شده \hat{H} ، باید به شکلی ارائه شود، که با توجه به توزیع احتمال داده شده و با در نظر گرفتن بیشینه سودمندی برای کاربر، پایگاه‌داده را سنتز کند. این سازوکار، باید تعداد، نسبت، و نوع رکوردهای پوششی را با توجه به سازوکار سنتز در حالت‌های (a)، (b)، و (c)، و همچنین، توزیع احتمال داده شده روی k تعیین کند. آنچه مشخص است، سازوکار ارائه شده، باید به صورت محافظه‌کارانه عمل کند. البته این دیدگاه موجب کاهش سودمندی می‌شود. روشی که در این گام استفاده خواهد شد، بر اساس نظریه احتمالات خواهد بود. برای مثال، یکی از سازوکارهای تصادفی کردن ممکن است بر اساس اضافه کردن رکوردهای پوششی به نسبت $\frac{1}{3}$ از سازوکارهای گفته شده در حالت‌های (a)، (b)، یا (c) باشد. سازوکار دیگر ممکن است این باشد که محافظه‌کارانه‌ترین حالت تصادفی کردن از میان حالت‌های (a)، (b)، یا (c) را انتخاب کنیم و بر اساس آن رکوردهای پوششی را ارسال نماییم. کارایی هر کدام از این سازوکارها باید بر اساس میزان سودمندی روش برای کاربر و همچنین، میزان اتلاف حریم خصوصی محاسبه شده و بهترین آن‌ها انتخاب شود.

همان طور که مشخص است، با توجه به توزیع احتمال گفته شده روی k ، اضافه کردن نویز بر اساس حالت (c)، محافظه‌کارانه است و بیشترین میزان حفظ حریم خصوصی در برابر کمترین مقدار سودمندی را به همراه دارد. برای مثال، در صورتی که تعداد خوشه‌ها در الگوریتم دو باشد و ما بر اساس حالت (c) نویز ایجاد کنیم، در موتور جستجو سندهای مربوط به دو تا از خوشه‌های موجود در شکل ۶-۴ برای مثال C_1^c و C_2^c ، با هم ادغام می‌شوند. در این حالت، بردار مرکز خوشه ادغام شده، از مقدار اصلی خود بسیار منحرف می‌شود. در نتیجه، سودمندی بر خلاف حریم خصوصی بسیار کاهش پیدا می‌کند. برای مثال، اگر خوشه‌های ورزش و سیاست موجود باشند و سازوکار موضوع حقوق را اضافه کند، آنگاه، در صورتی که تعداد واقعی خوشه‌ها در الگوریتم دو باشد، موضوع حقوق با یکی از خوشه‌ها، برای مثال، موضوع سیاست ادغام می‌شود. بنابراین، بردار مرکز این خوشه به موضوع حقوق متمایل می‌شود. در صورتی که فقط باید بازتاب‌دهنده علاقه کاربر به سیاست باشد.

۷-۲-۴ گام سوم در حل مسئله با فرض جعبه‌سیاه بودن موتور جستجو

در گام سوم موتور جستجو جعبه‌سیاه در نظر گرفته می‌شود. به عبارت دیگر، فرایند ایجاد پروفایل کاربر و بازیابی (فرایند L و R در شکل ۴-۱) نادانسته فرض می‌شوند. برای به کارگیری مدل حریم خصوصی تفاضلی در سنتز پایگاه داده مغشوش شده \hat{H} ، باید فرایند L مشخص باشد. در این گام، ابتدا به کمک فنون یادگیری ماشین فرایند L را به دست می‌آوریم. همانطور که در شکل ۴-۱ مشخص است، برای به دست آوردن فرایند L ، به پروفایل کاربر نیاز است. با وجود این، به دلیل جعبه‌سیاه بودن موتور جستجو پروفایل کاربر در دست نیست. بنابراین، ابتدا باید پروفایل کاربر را به دست آورد. با استفاده از ایده مقاله مرجع [۷] و مقایسه نتایج جستجو در حالت شخصی سازی شده و شخصی سازی نشده می‌توان پروفایل کاربر را به شکل آرایه‌ای از موضوع‌ها و میزان علاقمندی کاربر به آن‌ها به دست آورد. به عبارت دیگر، در این ایده معکوس فرایند R یاد گرفته می‌شود.

می‌توان با استفاده از پروفایل به دست آمده، به کمک فنون یادگیری ماشین، فرایند L را نیز به دست آورد. با مشخص شدن فرایند L ، می‌توانیم مدل حریم خصوصی تفاضلی را به مسئله بنگاریم و مانند گام نخست، سازوکار سنتز پایگاه داده مغشوش شده \hat{H} را صوری کنیم. برای این کار، دو ایده مطرح می‌شود. در ایده اول، در یک پیش‌پردازش فرایند L ، شناخته می‌شود. آنگاه، سازوکار تصادفی کردن عملی می‌شود. در ایده دوم، به صورت افقی و در هنگام تعامل با موتور جستجو، فرایند L به صورت تدریجی شناخته می‌شود و متناسب با شناخت فرایند L ، سازوکار تصادفی کردن عملی می‌شود.

به دست آوردن معکوس فرایند R و نیز فرایند L با کمک فنون یادگیری ماشین به طور ذاتی دارای خطا هستند. اثر این خطا در سازوکار سنتز پایگاه داده مغشوش شده \hat{H} و همچنین، میزان سودمندی و اتلاف حریم خصوصی باید تحلیل شود. علاوه بر آن، به کمک فنون یادگیری ماشین می‌توان مقدار دقیق پارامتر نامعلوم در حالت جعبه‌خاکستری (گام دوم) را نیز به دست آورد و سپس، مانند گام سوم مسئله را حل کرد. گفتنی است، رویکرد گام دوم استفاده از نظریه احتمالات است. لازم است، هر دو رویکرد مطرح شده در گام دوم و سوم از نظر میزان سودمندی و اتلاف حریم خصوصی مقایسه شوند.

۳-۴ ارزیابی

سازوکاری که برای حفظ حریم خصوصی کاربر در جستجوی وب شخصی سازی شده ارائه خواهد شد، بر اساس مدل حریم خصوصی تفاضلی است. در مرحله اول ارزیابی، باید صحت سازوکار ارائه شده بررسی شود. به عبارت دیگر، باید نشان داده شود که سازوکار ارائه شده طبق خط‌مشی حریم خصوصی کاربر به درستی حریم خصوصی

او را حفظ می‌کند. برای این کار، ما از روش اثبات صوری و ریاضی استفاده خواهیم کرد. در مرحله دوم ارزیابی، باید سودمندی سازوکار ارائه شده بررسی شود. به عبارت دیگر، باید بررسی کنیم که بعد از عملی کردن سازوکار حفظ حریم خصوصی، پیوندهای پیشنهادی موتور جستجو بعد از هر جستجوی کاربر چقدر برای او مفید است. برای این کار، باید تفاوت پیوندهای پیشنهادی موتور جستجو با حضور سازوکار حفظ حریم خصوصی و بدون آن را بررسی کنیم. در صورتی که، تفاوتی وجود نداشته باشد، بیشترین میزان سودمندی حاصل شده است. برای اندازه‌گیری این تفاوت باید تعدادی سنجه انتخاب و مورد ارزیابی قرار بگیرند. برای این کار، لازم است مجموعه‌ای از پرس‌وجوها و سندهای محک^۱ انتخاب شوند [۱]. با استفاده از این محک‌ها، تاریخچه جستجوی کاربر ایجاد شده و بر اساس آن در موتور جستجو پروفایل کاربر ساخته می‌شود.

AnonID	Query	QueryTime	ItemRank	ClickURL
142	rentdirect.com	2006-03-01 07:17:12		
142	www.prescriptionfortime.com	2006-03-12 12:31:06		
142	merit release appearance	2006-04-22 23:51:18		
142	www.bonsai.wbff.org	2006-05-06 08:49:34		
142	loislaw.com	2006-05-12 22:43:36		
142	rapny.com	2006-05-18 09:21:57		
142	whitepages.com	2006-05-19 19:36:31		
217	lottery	2006-03-01 11:58:51	1	http://www.calottery.com
217	lottery	2006-03-01 11:58:51	1	http://www.calottery.com
217	ameriprise.com	2006-03-01 14:06:23	1	http://www.ameriprise.com
217	susheme	2006-03-02 12:31:08		
217	united.com	2006-03-03 14:54:13		
217	mizuno.com	2006-03-07 22:41:17	1	http://www.mizuno.com

شکل ۴-۷: نمونه‌ای از رکوردهای ذخیره شده در پایگاه داده AOL

سنجه‌های بررسی سودمندی میزان ارتباط سندهای نتیجه شده با پرس‌وجو و همچنین، میزان علاقه‌مندی کاربر به آن سندها را اندازه‌گیری می‌کنند. سنجه‌های بررسی سودمندی باید برای هر دو حالت حضور و عدم حضور سازوکار حفظ حریم خصوصی اندازه‌گیری و مقایسه شوند. یکی از این سنجه‌ها ممکن است MAP^۲ باشد. در این سنجه، با استفاده از معیارهای دقت^۳ و فراخوانی^۴ میزان ارتباط سندهای نتیجه شده به پرس‌وجو و ترتیب نمایش آن‌ها به کاربر اندازه‌گیری می‌شود. علاوه بر آن سنجه‌ای برای تناسب سندهای نتیجه شده و پروفایل کاربر برای اندازه‌گیری میزان شخصی‌سازی باید انتخاب شود. برای ارزیابی می‌توان از پایگاه داده AOL استفاده کرد. این پایگاه داده دارای حدود ۱۷ میلیون پرس‌وجو و سندهای کلیک شده در پی ارسال هر پرس‌وجو

^۱ benchmark

^۲ mean average precision (MAP)

^۳ precision

^۴ recall

است که توسط حدود ۶۵۷ هزار کاربر در بازه ۱ مارس تا ۳۱ مه سال ۲۰۰۶ ایجاد شده است. با وجود قدیمی بودن، این پایگاه داده بزرگترین پایگاه داده انگلیسی از تاریخچه جستجوی کاربران است. نمونه‌ای از رکوردهای این پایگاه داده در شکل ۴-۷ دیده می‌شود.

۴-۴ هدف‌های رساله و زمان‌بندی

با توجه به مسئله مطرح شده و چالش‌های عنوان شده آن، اهداف ما در رساله به صورت زیر تعریف می‌شود.

- ارائه طرحی برای سنتز پایگاه داده مغشوش شده بر اساس مدل حریم خصوصی تفاضلی برای جستجوی وب شخصی‌سازی شده با فرض جعبه سفید بودن موتور جستجو.
- ارائه طرحی برای سنتز پایگاه داده مغشوش شده بر اساس مدل حریم خصوصی تفاضلی برای جستجوی وب شخصی‌سازی شده با فرض جعبه خاکستری بودن موتور جستجو.
- ارائه طرحی برای سنتز پایگاه داده مغشوش شده بر اساس مدل حریم خصوصی تفاضلی برای جستجوی وب شخصی‌سازی شده با فرض جعبه سیاه بودن موتور جستجو.

در جدول ۴-۱ زمانبندی پیشنهادی انجام رساله آمده است. همان طور که مشخص است تعدادی از این فعالیت‌ها با هم هم‌پوشانی دارند.

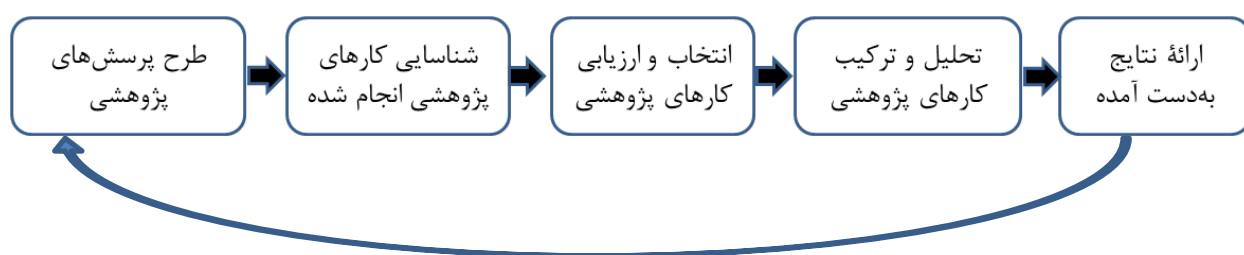
جدول ۴-۱: زمان‌بندی اجرای فعالیت‌های رساله پیشنهادی

شماره فعالیت	فعالیت‌ها	مدت زمان (برحسب ماه)
۱	حل مسئله با فرض جعبه‌سفید بودن موتور جستجو	۸
۱-۱	ارائه روشی برای بیان خط‌مشی حریم خصوصی تفاضلی (همسایگی و پارامتر ϵ) و راهکاری برای چالش همبستگی رکوردها	۴
۱-۲	ارائه راهکاری برای مسئله پایگاه‌داده رشدیابنده	۲
۱-۳	ارائه معیاری برای محاسبه سودمندی و ارائه الگوریتم سنتز پایگاه‌داده مغشوش شده	۲
۲	حل مسئله با فرض جعبه‌خاکستری بودن موتور جستجو	۶
۱-۲	بررسی الگوریتم یادگیری پروفایل کاربر همراه با تخمین به شکل یک توزیع احتمال روی مقادیر پارامترهای آن و تاثیر این تخمین در صوری کردن طرح و چالش‌های مطرح (همبستگی رکوردها و پایگاه‌داده رشدیابنده و غیره)	-
۳	حل مسئله با فرض جعبه‌سیاه بودن موتور جستجو	۶
۳-۱	شناسایی الگوریتم ایجاد پروفایل کاربر به وسیله فنون یادگیری ماشین	۳
۳-۲	ارائه طرحی برای سنتز پایگاه‌داده مغشوش شده بر اساس مدل حریم خصوصی تفاضلی برای روش به دست آمده در ۱-۳، با در نظر گرفتن وجود خطا و همچنین، چالش‌های همبستگی رکوردها و پایگاه‌داده رشدیابنده و غیره	۳
۴	مقایسه روش مطرح در ۳ و ۲ از نظر میزان سودمندی و اتلاف حریم خصوصی	۲
۵	تدوین رساله	۲
۶	مطالعه کارهای پژوهشی مرتبط جدید	۲۴

ضمیمه الف

نحوه انجام پژوهش جهت انتخاب موضوع رساله

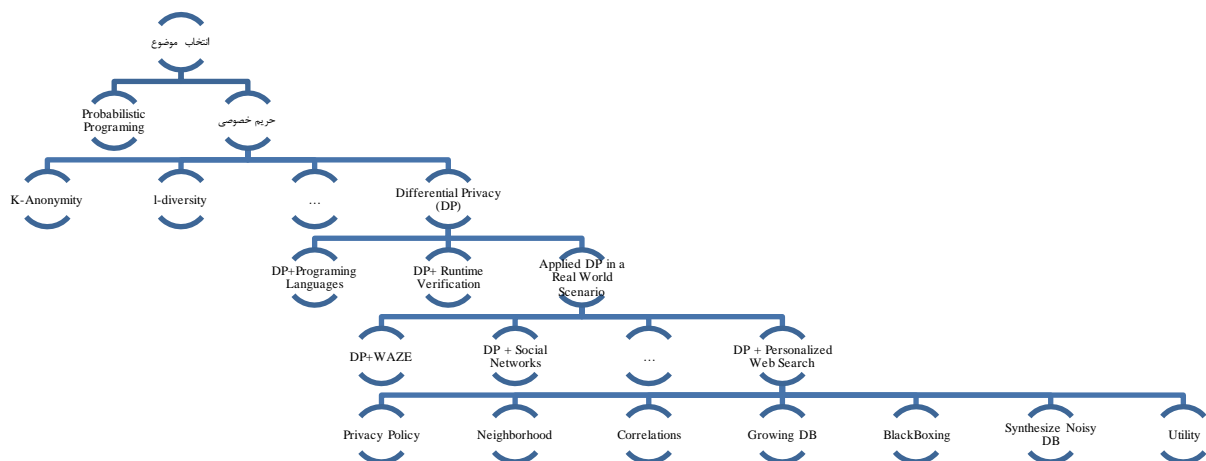
با انجام مشورت با استاد راهنما، زمینه پژوهشی اینجانب برای انجام رساله دکترا در حوزه حریم خصوصی *تفاضلی* تعیین گردید. در این راستا، متدولوژی استفاده شده به منظور جمع‌آوری پژوهش‌ها و مطالعه آن‌ها روش SLR^۱ انتخاب شد. با استفاده از این روش می‌توان یک دسته‌بندی سازماندهی شده از پژوهش‌های انجام شده در یک زمینه را ارائه داد. خروجی حاصل از این روش قابل تکرار و راستی‌آزمایی توسط سایر پژوهش‌گران است که این امر باعث شده است تا روش SLR از مقبولیت بالایی در میان پژوهش‌گران برخوردار باشد. متدولوژی SLR استفاده شده در پیشنهاد رساله پیش‌رو دارای پنج گام است که در شکل ۵-۱ نمایش داده شده است [۵۸].



شکل ۵-۱: متدولوژی پژوهشی استفاده شده

توضیحات نوشته شده در هر گام گویای فعالیت‌هایی است که در آن گام باید انجام شود. در شکل ۵-۲، روند انتخاب موضوع آمده است. در هر سطح از موضوع‌های موجود در این شکل، به ازای هر موضوع فرایند موجود در شکل ۵-۱ اجرا شده است. گفتنی است، کارهای انجام شده در این پیشنهاد رساله در رابطه با هر یک از این گام‌ها در جدول ۵-۱ نشان داده شده است.

^۱ Systematic Literature Review



شکل ۲-۵: روند انتخاب موضوع پژوهشی

جدول ۱-۵: خلاصه‌ای از متودولوژی استفاده شده برای جستجو و پژوهش

فعالیت انجام شده	گام
<p>در این گام، به ترتیب پرسش‌های زیر مطرح شده است:</p> <ul style="list-style-type: none"> • حوزه‌های تحقیقاتی در رابطه با حریم خصوصی کدام است؟ • حوزه‌های تحقیقاتی در رابطه با حریم خصوصی تفاضلی کدام است؟ • حوزه‌های تحقیقاتی در رابطه با حریم خصوصی تفاضلی و طراحی زبان‌های برنامه‌نویسی کدام است؟ • حوزه‌های تحقیقاتی در رابطه با حریم خصوصی تفاضلی و راستی‌آزمایی زمان‌اجرا کدام است؟ • حوزه‌های تحقیقاتی در رابطه با حریم خصوصی تفاضلی در فضای کاربردی کدام است؟ • حوزه‌های تحقیقاتی در رابطه با حریم خصوصی تفاضلی در جستجوی وب شخصی‌سازی شده کدام است؟ 	طرح پرسش‌های پژوهشی

<ul style="list-style-type: none"> • در این گام، پایگاه‌داده‌های علمی نشان داده شده در جدول ۲-۵ مورد جستجو قرار گرفت. علاوه بر پایگاه‌داده‌های این جدول از پایگاه‌داده scholar.google.com نیز استفاده شده است. کلمه‌های کلیدی مرتبط با موضوع‌های زیر همراه با «Differential Privacy» در این پایگاه‌داده‌ها جستجو شد. این موضوع‌ها عبارتند از: <ul style="list-style-type: none"> ❖ Personalized Web Search ❖ Neighborhood ❖ Correlations ❖ Growing DB ❖ BlackBoxing ❖ Synthetize Noisy DB <p>لازم به ذکر است که برای انجام جستجو از امکانات موجود در بخش جستجوی پیشرفته پایگاه‌داده‌ها استفاده گردیده است. البته از امکانات ویژه‌ای که برخی از این پایگاه‌داده‌ها ارائه کرده‌اند نیز استفاده شده است. به عنوان نمونه، انجام جستجو در کلیدواژه‌ها، چکیده و عنوان‌های مربوط به مقاله‌ها.</p> <ul style="list-style-type: none"> • همچنین، عمل جستجو از اوایل سال ۲۰۱۹ میلادی تا انتهای ماه آگوست سال ۲۰۲۱ در جریان بوده است. اما، در خصوص کارهای پژوهشی به‌دست آمده، هیچ نوع محدودیتی از این نظر که این کارها در چه زمانی انجام شده‌اند، در نظر گرفته نشده است. 	<p>شناسایی کارهای پژوهشی انجام شده</p>
<ul style="list-style-type: none"> • معیارهای گنجاندن کارهای پژوهشی جهت مطالعه دقیق‌تر: مهم‌ترین معیار در انتخاب کارهای پژوهشی وجود کلیدواژه‌های استفاده شده در چکیده و مخصوصاً عنوان مقاله‌ها بوده است. • معیارهای حذف کارهای پژوهشی: این معیارها عبارتند از حذف مقاله‌های تکراری و مقاله‌هایی که به عنوان مقدمه یا پیش‌گفتار در کنفرانس‌ها و مجله‌ها چاپ شده‌اند. • روش‌های انتخاب کارهای پژوهشی: پس از مطالعه چکیده، مقدمه، نتیجه‌گیری و نگاهی اجمالی به متن هر مقاله کارهای پژوهشی مرتبط انتخاب گردیده است. • در این گام تعدادی مقاله انتخاب گردید. این مقاله‌ها از کنفرانس‌ها 	<p>انتخاب و ارزیابی کارهای پژوهشی</p>

و ژورنال‌های جدول ۳-۵ و • جدول ۴-۵ گردآوری شده است.	
در این مرحله، کارهای پژوهشی انجام شده در دسته‌های زیر طبقه‌بندی شدند:	تحلیل و ترکیب کارهای پژوهشی
<ul style="list-style-type: none"> • مفهومی‌های بنیادی حریم خصوصی و حریم خصوصی تفاضلی • حریم خصوصی تفاضلی برای ایجاد زبان‌های برنامه‌نویسی • حریم خصوصی تفاضلی در راستی‌آزمایی زمان‌اجرا • حریم خصوصی تفاضلی در شبکه‌های اجتماعی • حریم خصوصی تفاضلی در موقعیت مکانی کاربران • حریم خصوصی تفاضلی در جستجوی وب شخصی‌سازی شده • سایر موارد مرتبط 	
<p>نتیجه به‌دست آمده پس از مطالعه مقاله‌های انتخاب شده این است که غالب کارهای پژوهشی فعلی، از حریم خصوصی تفاضلی برای حفظ حریم خصوصی کاربران در جستجوی وب شخصی‌سازی شده استفاده نکرده‌اند. بنابراین پرسش مطرح شده در این پیشنهاد رساله به‌صورت زیر خواهد بود:</p> <p>حفظ حریم خصوصی کاربران در جستجوی وب شخصی‌سازی شده بر اساس مدل حریم خصوصی تفاضلی چگونه خواهد بود؟</p>	ارائه نتایج به‌دست آمده

جدول ۲-۵: پایگاه‌داده‌های علمی مورد جستجو

ناشر	آدرس اینترنتی
Springer	https://link.springer.com
ACM	https://dl.acm.org
Elsevier	https://www.sciencedirect.com
IEEE	https://ieeexplore.ieee.org

جدول ۳-۵: نام کنفرانس‌های مرتبط با حریم خصوصی تفاضلی و جستجوی وب شخصی‌سازی شده

ردیف	نام کنفرانس
۱	ACM Transactions on Privacy and Security
۲	International Conference on World Wide Web - WWW
۳	International ACM SIGIR Conference on Research and Development in

Information – SIGIR	
International Conference on Trust, Security and Privacy in Computing and Communications-TRUSTCOM	۴
IEEE Symposium on Security and Privacy-SP	۵
ACM conference on Information and knowledge management-CIKM	۶
Knowledge discovery and data mining-KDD	۷
ACM SIGMOD-SIGACT-SIGART Conference on Principles of Database Systems-PODS	۸
International Conference on Data Engineering-ICDE	۹
Very Large Data Bases-VLDB	۱۰
International Colloquium on Automata Languages and Programming-ICALP	۱۱
International Conference on Web Information Systems Engineering-WISE	۱۲
ACM Conference on Computer and Communications Security-CCS	۱۳
Extending Database Technology-EDBT	۱۴
International Conference on Database Theory-ICDT	۱۵

جدول ۴-۵: نام مجله‌های مرتبط با حریم خصوصی تفاضلی و جستجوی وب شخصی‌سازی شده

نام مجله	ردیف
ACM SIGIR Forum	۱
Information Sciences	۲
Data & Knowledge Engineering	۳
IEEE Transactions on Information Forensics and Security	۴
Online Information Review	۵
ACM Transactions on Knowledge Discovery from Data	۶
IEEE Transactions on Knowledge and Data Engineering	۷
ACM Transactions on Knowledge Discovery from Data	۸

ضمیمه ب

پیاده‌سازی الگوریتم ایجاد پروفایل کاربر

پروفایل ساخته شده بر اساس مدل سامانه گفته شده در فصل ۴ قسمت ۴-۲-۱ است. پروفایل ساخته شده، از نوع پروفایل کلمه کلیدی^۱ است. در این الگوریتم، سندها با جدا کردن متن اصلی از سند، حذف کلمه‌های ایست، و ریشه‌سازی پردازش شده و لیستی از کلمه‌های هر سند به دست می‌آید. سپس، با استفاده از آماره TF-IDF بردار هر سند را محاسبه کردیم. بردارهای به دست آمده، با استفاده از الگوریتم k-means با اندازه $k = 10$ پردازش شدند و پروفایل کاربر را ایجاد کردیم. پروفایل کاربر ساخته شده شامل ۱۰ بردار است. این بردارها، بردارهای مرکز خوشه‌ها در الگوریتم k-means هستند. کد شکل ۱-۶، پروفایل کاربری را می‌سازد که ۱۰۰۰۰ سند را در اینترنت جستجو کرده است:

```
from selectolax.parser import HTMLParser
import pickle
from domain_utils import *
import re
import string
from nltk.stem import PorterStemmer

class CustomTokenizer:

    def __init__(self, n_pages=10000, path_stopwords='/content/ProfileBuilder/stopwords.txt'):
        self.path_stopwords = path_stopwords
        self.FOLDER = '/content/ProfileBuilder/clickedDocs/'
        self.n_pages = n_pages
        self.stemmer = PorterStemmer()
        self.docs_tokens = {}
        with open(self.path_stopwords, 'r') as stop_file:
            self.stop_words = stop_file.readlines()
        self.stop_words = list(map(lambda x: x[:-1], self.stop_words))

    @staticmethod
    def get_text(html):
        tree = HTMLParser(html)
        if tree.body is None:
            return None

        for tag in tree.css('script'):
            tag.decompose()
        for tag in tree.css('style'):
            tag.decompose()

        text = tree.body.text(separator='')
        return text
```



```
def get_docs_tokens(self):
    return self.docs_tokens

def process_page(self, code, doc_text):
    doc_text = self.get_text(doc_text)
    if doc_text is None:
        doc_text = 'none'
    tokens = self.tokenize(doc_text)
    self.docs_tokens[code] = tokens

def tokenize(self, doc_text):
    tokens = doc_text.split()
    tokens = [''.join(c for c in t if c not in string.punctuation) for t in tokens]
    tokens = [t.lower() for t in tokens]
    tokens = map(replace_digits, tokens)
    tokens = map(self.stemmer.stem, tokens)
    tokens = [t for t in tokens if not lesseq_two_letters(t)]
    tokens = [t for t in tokens if t not in self.stop_words]
    tokens = [t for t in tokens if t]
    return tokens

def preprocess_documents(self):
    for filename in range(self.n_pages):
        with open(self.FOLDER + str(filename)) as f:
            doc_text = f.read()

            if doc_text is None:
                print('empty doc')
                print(filename)
                continue

            self.process_page(int(filename), doc_text)

# removing digits and returning the word
def replace_digits(st):
    return re.sub('\d', '', st)

# returns true if the word has less or equal 2 letters
def lesseq_two_letters(word):
    return len(word) <= 2

CT = CustomTokenizer()
CT.preprocess_documents()

from sklearn.feature_extraction.text import TfidfVectorizer

docsTokens = CT.get_docs_tokens()
corpus = []
for key in docsTokens.keys():
    corpus.append(''.join(docsTokens[key]))
vectorizer = TfidfVectorizer()
X = vectorizer.fit_transform(corpus)
```

```
print(X.shape)

### OUTPUT→ (10000, 63089)

from sklearn.cluster import KMeans

kmeans = KMeans(n_clusters=10, random_state=0).fit(X)
print(kmeans.cluster_centers_)

### OUTPUT→
array([[0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
        0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
       [9.51832544e-06, 0.00000000e+00, 0.00000000e+00, ...,
        0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
       [1.23299192e-04, 0.00000000e+00, 3.82269553e-06, ...,
        0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
       ...,
       [0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
        0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
       [1.99277923e-04, 0.00000000e+00, 0.00000000e+00, ...,
        0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
       [0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
        0.00000000e+00, 0.00000000e+00, 0.00000000e+00]])
```

شکل ۶-۱: پیاده‌سازی الگوریتم ایجاد پروفایل کاربر

ضمیمه پ

واژه‌نامه انگلیسی به فارسی

Adaptive	وفقی
Anonymity	بی‌نامی
Bias	متمایل
Closeness	نزدیکی
Composition	ترکیب
context information	اطلاعات زمینه
Count Query	پرس‌وجوی شمارش
Data Type	نوع داده
Database Curator	مسئول پایگاه‌داده
Differential Privacy	حریم خصوصی تفاضلی
Distinct	متمايز
Diversity	تنوع
Equivalence Class	کلاس هم‌ارزی
Explainer	توضیح‌دهنده
Exponential	نمایی
Feed	سرنخ‌های خبری
Formalism Method	روش صوری‌سازی
Generalized	عمومی شده
Global	سراسری
Group Identity	هویت‌های گروهی
Honest but curious	درستکار ولی کنج‌کاو
Indexing	نمایه‌سازی
Indistinguishability	قابلیت عدم تمایز
Interpretable Machine Learning	یادگیری ماشین تفسیرپذیر
inversion	وارون‌سازی
Laplace	لاپلاس
Local Differential Privacy	حریم خصوصی تفاضلی محلی
Mechanism	سازوکار
Minimal	کمینه

Ontology	آنتولوژی محاسباتی
page rank	رتبه‌بندی صفحه
Parallel Composition	ترکیب موازی
Parse	پارس
Partition-based	مبتنی بر افراز کردن
peer to peer	نظیر به نظیر
Personalized Information Retrieval	داده‌کاوی شخصی‌شده
Plausible Deniability	انکارپذیری قابل قبول
Posterior	موخر
Post-Processing	پسپردازش
Preference	پسند
Presense	حضور
privacy	حریم خصوصی
Privacy Loss	اتلاف حریم خصوصی
proxy system	سامانه نایب
Proxy Topic	موضوع نایب
Query Expansion	گسترش پرس‌وجو
Randomization Mechanism	سازوکار تصادفی کردن
Randomization-based	مبتنی بر تصادفی کردن
Randomized Response	پاسخ تصادفی شده
ranking	رتبه‌بندی
recommender system	سامانه توصیه‌گر
Recursive Diversity	تنوع بازگشتی
response time	زمان پاسخ
retrieval	بازیابی
Score Function	تابع امتیاز
search behavior	رفتار جستجو
Sensitivity	حساسیت

Sequential Composition	ترکیب ترتیبی
Statistical Database	پایگاه داده آماری
Stem	ریشه
Stop Words	کلمه‌های ایست
Streaming Data	داده‌های جریانی
Synthetic Database	پایگاه داده ساختگی
Utility	سودمندی
Utility Function	تابع سودمندی
Web Crawler	خزش‌گر وب
Well-represented	خوش-نمایان شده

ضمیمه ت

واژه‌نامه فارسی به انگلیسی

Adaptive	وفقی
Privacy Loss	اتلاف حریم خصوصی
context information	اطلاعات زمینه
Plausible Deniability	انکارپذیری قابل قبول
retrieval	بازیابی
Anonymity	بی نامی
Parse	پارس
Randomized Response	پاسخ تصادفی شده
Statistical Database	پایگاه داده آماری
Synthetic Database	پایگاه داده ساختگی
Count Query	پرس و جوی شمارش
Post-Processing	پساپردازش
Preference	پسند
Score Function	تابع امتیاز
Utility Function	تابع سودمندی
Composition	ترکیب
Sequential Composition	ترکیب ترتیبی
Parallel Composition	ترکیب موازی
Diversity	تنوع
Recursive Diversity	تنوع بازگشتی
Explainer	توضیح دهنده
privacy	حریم خصوصی
Differential Privacy	حریم خصوصی تفاضلی
Local Differential Privacy	حریم خصوصی تفاضلی محلی
Sensitivity	حساسیت

Presense	حضور
Web Crawler	خزش گر وب
Well-represented	خوش-نمایان شده
Personalized Information Retrieval	داده کاوی شخصی شده
Streaming Data	داده های جریانی
Honest but curious	درستکار ولی کنجکاو
ranking	رتبه بندی
page rank	رتبه بندی صفحه
search behavior	رفتار جستجو
Formalism Method	روش صوری سازی
Stem	ریشه
response time	زمان پاسخ
Mechanism	سازوکار
Randomization Mechanism	سازوکار تصادفی کردن
recommender system	سامانه توصیه گر
proxy system	سامانه نایب
Global	سراسری
Feed	سرنخ های خبری
Generalized	عمومی شده
Utility	سودمندی
Indistinguishability	قابلیت عدم تمایز
Equivalence Class	کلاس هم ارزی
Stop Words	کلمه های ایست
Minimal	کمینه
Query Expansion	گسترش پرس و جو

Laplace	لاپلاس
Partition-based	مبتنی بر افراز کردن
Randomization-based	مبتنی بر تصادفی کردن
Distinct	متمایز
Bias	متمایل
Database Curator	مسئول پایگاه داده
Posterior	موخر
Proxy Topic	موضوع نایب
Closeness	نزدیکی
peer to peer	نظیر به نظیر
Indexing	نمایه سازی
Exponential	نمایی
Data Type	نوع داده
Ontology	آنتولوژی محاسباتی
Group Identity	هویت های گروهی
inversion	وارون سازی
Interpretable Machine Learning	یادگیری ماشین تفسیرپذیر

مراجع

- [1] B. Croft, D. Metzler, and T. Strohman, *Search Engines: Information Retrieval in Practice*, 1st ed. Addison-Wesley Publishing Company, 2015.
- [2] J. Liu, C. Liu, and N. J. Belkin, "Personalization in Text Information Retrieval: A Survey," *Journal of the Association for Information Science and Technology*, vol. 71, no. 3, pp. 349–369, Mar. 2020.
- [3] S. L. Garfinkel, "De-Identification of Personal Information," *National Institute of Standards and Technology*, 2015. <http://dx.doi.org/10.6028/NIST.IR.8053> (accessed Jan. 10, 2021).
- [4] J. Domingo-Ferrer, M. Bras-Amorós, Q. Wu, and J. Manjón, "User-Private Information Retrieval Based on a Peer-to-Peer Community," *Data & Knowledge Engineering*, vol. 68, no. 11, pp. 1237–1252, Nov. 2009.
- [5] D. Sánchez, J. Castellà-Roca, and A. Viejo, "Knowledge-Based Scheme to Create Privacy-Preserving but Semantically-Related Queries for Web Search Engines," *Information Sciences*, vol. 218, pp. 17–30, 2013.
- [6] G. Chen, H. Bai, L. Shou, K. Chen, and Y. Gao, "UPS: Efficient Privacy Protection in Personalized Web Search," in *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information - SIGIR '11*, 2011, p. 615.
- [7] A. Majumder and N. Shrivastava, "Know Your Personalization: Learning Topic Level Personalization in Online Services," in *Proceedings of the 22nd international conference on World Wide Web - WWW '13*, 2013, pp. 873–884.
- [8] P. Mac Aonghusa and D. J. Leith, "Don't Let Google Know I'm Lonely," *ACM Transactions on Privacy and Security*, vol. 19, no. 1, pp. 1–25, Aug. 2016.
- [9] P. Mac Aonghusa and D. J. Leith, "Plausible Deniability in Web Search - From Detection to Assessment," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 4, pp. 874–887, 2018.
- [10] P. Mac Aonghusa and D. Leith, "3PS - Online Privacy through Group Identities," *arXiv preprint:1811.11039*, 2018.
- [11] A. Gervais, R. Shokriy, A. Singlay, S. Capkun, and V. Lendersz, "Quantifying Web-Search Privacy," in *Proceedings of the ACM Conference on Computer and Communications Security*, 2014, pp. 966–977.
- [12] W. U. Ahmad, M. M. Rahman, and H. Wang, "Topic Model based Privacy Protection in Personalized Web Search," in *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval - SIGIR '16*, 2016, pp. 1025–1028.
- [13] A. Petit, T. Cerqueus, S. Ben Mokhtar, L. Brunie, and H. Kosch, "PEAS: Private, Efficient and Accurate Web Search," in *2015 IEEE Trustcom/BigDataSE/ISPA-TRUSTCOM'15*, 2015, pp. 571–580.
- [14] D. C. Howe and H. Nissenbaum, "TrackMeNot: Resisting Surveillance in Web Search," in *Lessons from the Identity Trail: Anonymity, Privacy and Identity in a Networked Society*, L. Kerr, C. Lucock, and V. Steeves, Eds. Oxford University Press, 2009, pp. 417–436.
- [15] J. Domingo-Ferrer, A. Solanas, and J. Castellà-Roca, "H(κ)-Private Information Retrieval from Privacy-Uncooperative Queryable Databases," *Online Information Review*, vol. 33, no. 4, pp. 720–744, 2009.
- [16] M. Radhika and V. Vijaya Chamundeeswari, "Privacy Protection in Personalized Web

- Search using Obfuscation,” *ARNP Journal of Engineering and Applied Sciences*, vol. 10, no. 7, pp. 3225–3227, 2020.
- [17] E. Balsa, C. Troncoso, and C. Diaz, “OB-PWS: Obfuscation-Based Private Web Search,” in *2012 IEEE Symposium on Security and Privacy*, 2012, pp. 491–505.
 - [18] Y. Xu, K. Wang, G. Yang, and A. W. C. Fu, “Online Anonymity for Personalized Web Services,” in *Proceeding of the 18th ACM conference on Information and knowledge Management - CIKM '09*, 2009, pp. 1497–1500.
 - [19] X. Shen, B. Tan, and C. Zhai, “Privacy Protection in Personalized Search,” *ACM SIGIR Forum*, vol. 41, no. 1, pp. 4–17, Jun. 2007.
 - [20] Y. Zhu, L. Xiong, and C. Verdery, “Anonymizing User Profiles for Personalized Web Search,” in *Proceedings of the 19th international conference on World wide web - WWW '10*, 2010, p. 1225.
 - [21] R. Cummings, S. Krehbiel, K. A. Lai, and U. Tantipongpipat, “Differential Privacy for Growing Databases,” in *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018-NeurIPS 2018*, 2018, pp. 8878–8887.
 - [22] D. Desfontaines and B. Pejó, “SoK: Differential Privacies,” *Proceedings on Privacy Enhancing Technologies*, vol. 2020, no. 2, pp. 288–313, Apr. 2020.
 - [23] S. R. Ganta, S. P. Kasiviswanathan, and A. Smith, “Composition Attacks and Auxiliary Information in Data Privacy,” in *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2008, pp. 265–273.
 - [24] L. Sweeney., “k-Anonymity: A Model for Protecting Privacy,” *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 10, no. 05, pp. 557–570, Oct. 2002.
 - [25] P. Samarati and L. Sweeney, “Generalizing Data to Provide Anonymity when Disclosing Information,” in *Proceedings of the seventeenth ACM SIGACT-SIGMOD-SIGART symposium on Principles of database systems - PODS '98*, 1998, vol. 23, no. 3, p. 188.
 - [26] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkatasubramanian, “l-Diversity: Privacy Beyond k-Anonymity,” *ACM Transactions on Knowledge Discovery from Data*, vol. 1, no. 1, pp. 3-es, Mar. 2007.
 - [27] N. Li, T. Li, and S. Venkatasubramanian, “t-Closeness: Privacy Beyond k-Anonymity and l-Diversity,” in *2007 IEEE 23rd International Conference on Data Engineering*, 2007, no. 2, pp. 106–115.
 - [28] Ninghui Li, Tiancheng Li, and S. Venkatasubramanian, “Closeness: A New Privacy Measure for Data Publishing,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 7, pp. 943–956, Jul. 2010.
 - [29] M. E. Nergiz, M. Atzori, and C. Clifton, “Hiding the Presence of Individuals from Shared Databases,” in *Proceedings of the ACM SIGMOD International Conference on Management of Data*, 2007, pp. 665–676.
 - [30] V. Bindshaedler, R. Shokri, and C. A. Gunter, “Plausible Deniability for Privacy-Preserving Data Synthesis,” *Proceedings of the VLDB Endowment*, vol. 10, no. 5, pp. 481–492, Jan. 2017.
 - [31] R. C. W. Wong, A. W. C. Fu, K. Wang, and J. Pei, “Minimality Attack in Privacy Preserving Data Publishing,” in *Proceedings of the 33rd International Conference on Very Large Data Bases*, 2007, pp. 543–554.

- [32] R. C.-W. Wong, A. W.-C. Fu, K. Wang, P. S. Yu, and J. Pei, "Can the Utility of Anonymized Data be Used for Privacy Breaches?," *ACM Transactions on Knowledge Discovery from Data*, vol. 5, no. 3, pp. 1–24, Aug. 2011.
- [33] S. L. Warner, "Randomized Response: A Survey Technique for Eliminating Evasive Answer Bias," *Journal of the American Statistical Association*, vol. 60, no. 309, pp. 63–69, Mar. 1965.
- [34] C. Dwork, "Differential Privacy," in *33rd International Colloquium Automata, Languages and Programming (ICALP)*, 2006, vol. 4052, pp. 1–12.
- [35] C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating Noise to Sensitivity in Private Data Analysis," in *Theory of Cryptography Conference*, 2006, pp. 265–284.
- [36] C. Dwork and A. Roth, "The Algorithmic Foundations of Differential Privacy," *Foundations and Trends® in Theoretical Computer Science*, vol. 9, no. 3–4, pp. 211–407, 2013.
- [37] C. Dwork and A. Roth, "The Algorithmic Foundations of Differential Privacy," *Foundations and Trends® in Theoretical Computer Science*, vol. 9, no. 3–4, pp. 211–407, 2013.
- [38] D. Winograd-Cort, A. Haeberlen, A. Roth, and B. C. Pierce, "A Framework for Adaptive Differential Privacy," in *Proceedings of the ACM on Programming Languages*, 2017, vol. 1, pp. 1–29.
- [39] Q. Ye and H. Hu, "Local Differential Privacy: Tools, Challenges, and Opportunities," in *International Conference on Web Information Systems Engineering*, 2019, pp. 13–23.
- [40] Z. Qin, Y. Yang, T. Yu, I. Khalil, X. Xiao, and K. Ren, "Heavy Hitter Estimation over set-valued data with local differential privacy," in *Proceedings of the ACM Conference on Computer and Communications Security*, 2016, pp. 192–203.
- [41] R. Chen, H. Li, A. K. Qin, S. P. Kasiviswanathan, and H. Jin, "Private Spatial Data Aggregation in the Local Setting," in *2016 IEEE 32nd International Conference on Data Engineering (ICDE 2016)*, 2016, pp. 289–300.
- [42] G. Fanti, V. Pihur, and Ú. Erlingsson, "Building a RAPPOR with the Unknown: Privacy-Preserving Learning of Associations and Data Dictionaries," *Proceedings on Privacy Enhancing Technologies*, vol. 2016, no. 3, pp. 41–61, 2016.
- [43] Y. Wang, X. Wu, and D. Hu, "Using Randomized Response for Differential Privacy Preserving Data Collection," in *Proceedings of the Workshops of the (EDBT/ICDT) 2016 Joint Conference*, 2016, pp. 1–12.
- [44] Ú. Erlingsson, V. Pihur, and A. Korolova, "RAPPOR: Randomized Aggregatable Privacy-Preserving Ordinal Response," in *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security - CCS '14*, 2014, pp. 1054–1067.
- [45] B. Avent, A. Korolova, D. Zeber, T. Hovden, and B. Livshits, "Blender: Enabling local search with a hybrid differential privacy model," in *26th USENIX Security Symposium*, 2017, pp. 747–764.
- [46] X. Xiong, S. Liu, D. Li, Z. Cai, and X. Niu, "A Comprehensive Survey on Local Differential Privacy," *Security and Communication Networks*, vol. 2020, pp. 1–29, Oct. 2020.
- [47] Tianqing Zhu, Ping Xiong, Gang Li, and Wanlei Zhou, "Correlated Differential Privacy: Hiding Information in Non-IID Data Set," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 2, pp. 229–242, Feb. 2015.

- [48] A. Kacem, “Personalized Information Retrieval based on Time-Sensitive User Profile,” Université Paul Sabatier (Toulouse 3), 2017.
- [49] S. Gauch, M. Speretta, A. Chandramouli, and A. Micarelli, “User profiles for personalized information access,” *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 4321 LNCS, pp. 54–89, 2007.
- [50] N. Guarino, D. Oberle, and S. Staab, “What Is an Ontology?,” in *Handbook on Ontologies*, Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 1–17.
- [51] H. P. Alesso and C. F. Smith, *Thinking on the Web: Berners-Lee, Gödel and Turing*. 605 Third Avenue New York, NY, United States: Wiley-Interscience, 2008.
- [52] S. T. Peddinti and N. Saxena, “On the Privacy of Web Search Based on Query Obfuscation: A Case Study of TrackMeNot,” in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 6205 LNCS, 2010, pp. 19–37.
- [53] P. Kodeswaran and E. Viegas, “Applying Differential Privacy to Search Queries in a Policy Based Interactive Framework,” in *Proceeding of the ACM first international workshop on Privacy and anonymity for very large databases - PAVLAD '09*, 2009, p. 25.
- [54] A. El-Ansari, A. Beni-Hssane, M. Saadi, and M. El Fissaoui, “PAPIR: Privacy-Aware Personalized Information Retrieval,” *Journal of Ambient Intelligence and Humanized Computing*, vol. 12, no. 10, pp. 9891–9907, 2021.
- [55] H. A. Feild, J. Allan, and J. Glatt, “CrowdLogging: Distributed, Private, and Anonymous Search Logging,” in *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information - SIGIR '11*, 2011, no. 2, p. 375.
- [56] S. Zhang, H. Yang, and L. Singh, “Anonymizing Query Logs by Differential Privacy,” in *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, Jul. 2016, pp. 753–756.
- [57] D. Sánchez, M. Batet, A. Viejo, M. Rodríguez-García, and J. Castellà-Roca, “A Semantic-Preserving Differentially Private Method for Releasing Query Logs,” *Information Sciences*, vol. 460–461, pp. 223–237, Sep. 2018.
- [58] Y. Hong, J. Vaidya, H. Lu, P. Karras, and S. Goel, “Collaborative Search Log Sanitization: Toward Differential Privacy and Boosted Utility,” *IEEE Transactions on Dependable and Secure Computing*, vol. 12, no. 5, pp. 504–518, Sep. 2015.
- [59] X. Meng, Z. Xu, B. Chen, and Y. Zhang, “Privacy-Preserving Query Log Sharing Based on Prior N-Word Aggregation,” in *2016 IEEE Trustcom/BigDataSE/ISPA*, Aug. 2016, pp. 722–729.
- [60] L. Fan, L. Bonomi, L. Xiong, and V. Sunderam, “Monitoring Web Browsing Behavior with Differential Privacy,” in *WWW 2014 - Proceedings of the 23rd International Conference on World Wide Web*, 2014, pp. 177–187.
- [61] M. Joseph, A. Roth, J. Ullman, and B. Waggoner, “Local Differential Privacy for Evolving Data,” in *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018 (NeurIPS 2018)*, 2018, pp. 2381–2390.
- [62] U. Stemmer, “Locally Private k-Means Clustering,” in *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, 2020, pp. 548–559.
- [63] K. Nissim and U. Stemmer, “Clustering Algorithms for the Centralized and Local

- Models,” in *Proceedings of Algorithmic Learning Theory*, volume 83 of *Proceedings of Machine Learning Research*, 2018, pp. 619–653.
- [64] U. Stemmer and H. Kaplan, “Differentially Private k-Means with Constant Multiplicative Error,” in *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018 (NeurIPS 2018)*, 2018, pp. 5436–5446.
 - [65] C. Liu, S. Chakraborty, and P. Mittal, “Dependence Makes You Vulnerable: Differential Privacy Under Dependent Tuples,” 2016.
 - [66] X. Ren *et al.*, “LoPub: High-Dimensional Crowdsourced Data Publication With Local Differential Privacy,” *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 9, pp. 2151–2166, Sep. 2018.
 - [67] Y. Li, X. Ren, S. Yang, and X. Yang, “Impact of Prior Knowledge and Data Correlation on Privacy Leakage: A Unified Analysis,” *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 9, pp. 2342–2357, Sep. 2019.
 - [68] M. Rodriguez-Garcia, M. Batet, and D. Sanchez, “Semantic Noise: Privacy-Protection of Nominal Microdata through Uncorrelated Noise Addition,” in *2015 IEEE 27th International Conference on Tools with Artificial Intelligence (ICTAI)*, Nov. 2015, vol. 2016-Janua, pp. 1106–1113.
 - [69] W. U. Ahmad, K.-W. Chang, and H. Wang, “Intent-aware Query Obfuscation for Privacy Protection in Personalized Web Search,” in *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, Jun. 2018, pp. 285–294.
 - [70] D. Denyer and D. Tranfeld, “Producing a Systematic Review,” in *The SAGE Handbook of Organizational Research Methods*, London: Sage Publications, 2009, pp. 671–689.
 - [71] T. Zhu, G. Li, W. Zhou, and P. S. Yu, *Differential Privacy and Applications*. Springer International Publishing, 2017.