

بازبینی

User k-anonymity for privacy preserving
data mining of query logs

داریوش حسن پور آده

۹۳۰۸۱۶۴

این مقاله به مشکل حفظ حریم شخصی افراد در اطلاعات بدست آمده از موتورهای جستجو پرداخته است. برای نشان دادن اهمیت مساله نمونه مثال از شرکت AOL آورده است که برای منظور کمک به جامعه پژوهشی بازایی اطلاعات باعث گردید اطلاعات شخصی برخی از کاربران افشا شود برای همین مساله برای حفظ اطلاعات شخصی افراد قبل از اینکه گزارشات جستجوها^۱ برای عموم ارائه شود نیاز است که یک عملیات پنهان سازی^۲ بروی داده ها اعمال شود؛ ولی تعیین اینکه تا چه حد داده ها نیاز به پنهان سازی^۳ دارند، سخت می باشد. بنابراین سعی دارد که روشی ارائه دهد که توازن بین حفظ اطلاعات شخصی افراد و عدم از دست رفت اطلاعات واقعا مفید موجود در داده های جمع آوری شده از موتورهای جستجو ارائه دهد.

این مقاله با نام بردن چندین کار قبلی مبنی بر اینکه آنها نیز از روش k -anonymity استفاده کرده اند ولی بطور کامل تضمین نمی کردند که معیار k در k -anonymity رعایت شده است - زیرا با حذف داده ها به این منظور می رسیدند. نوآوری این مقاله در ارائه ی روشی برای تضمین معیار k -anonymity در داده های جستجو با استفاده از میکرو-دسته^۴ بدون اینکه هیچ یک از داده ها صراحتا از \log search حذف شده باشند.

این مقاله سپس به شرح مفهوم micro-aggregates پرداخته است. که شرح داده که شامل دو بخش می شود پارتیشن کردن و بعد متراکم کردن داده ها است. و همچنین آورده که برای این منظور یک جنگلی از درختان باید تشکیل شود که هر درخت متعلق به یکی از کاربران است و شاخه های این درختان شامل جستجوهای آن کاربر می باشد. و سپس معیار فاصله را برای هر جستجو معرفی کرده است.

بطور خلاصه اگر بخواهیم روش مطرح شده با استفاده از micro-aggregates را بگوییم می توان گفت که ابتدا داده ها با توجه به معیار فاصله مطرح شده ریز-دسته بندی^۵ می شوند سپس مرکز دسته ها به عنوان نماینده آن دسته می شود و یک کوئری^۶ به عنوان نماینده آن کوئری ها ارائه می شود. که علاوه بر اینکه نمایانگر کوئری های هم دسته خود می باشد به علت نوع انتزاعی ای این نماینده دارد به خوبی می تواند اطلاعات شخصی افراد را مخفی نگه دارد - در اینجا k تعداد کاربران موجود در هر دسته می باشد.

از مزایای روش ارائه شده در این مقاله علاوه بر سادگی روش می توان گفت که تضمین می کند معیار k در k -anonymity رعایت شده است و سپس اینکه اطلاعات شخصی افراد را می تواند با نسبت خوبی حفظ کند بدین نحو که بعد از پنهان سازی داده ها به ازای هر کوئری k کاربر مرتبط با آن کوئری می توان نسبت داد؛ همچنین مقاله آورده است که روش پیشنهادی از نظر محاسباتی حداکثر کارایی را دارد.

از آنجایی که روش مطرح شده در حالت پایه در واقع دسته بندی می باشد بنابراین تمامی معایب مطرح شده در دسته بندی را به طور ضمنی دارد مانند بررسی داده ها با ابعاد بالا نیاز به زمان بیشتری می خواهد، کارایی روش به فاصله تعریف شده وابسته می باشد. و همچنین مقاله در مورد راه حل هایی در مقابله به کوئری های پرت^۷ و تاثیر داده های پرت بر میزان پنهان سازی داده ها حرفی نزده و همان طور که می دانیم خوشه بندی به داده های پرت حساس می باشد که مقاله از کنار این موضوع بدون بررسی گذشته است. همچنین یکی از معایب دیگری که می توان به روش ارائه شده نسبت داد این است که میزان تعیین k مناسب نیز سخت می باشد که در اینجا کل عملیات

^۱ Search log

^۲ Anonymize

^۳ Degree of privacy

^۴ micro-aggregates

^۵ Clustering

^۶ Query

^۷ Outlier

بر مبنای مقدار k بنا شده است. از طرف دیگر روش ارائه شده همه‌ی ویژگی‌های موجود در رکوردها با یک دید یکسان نگاه می‌کند در حالی که اهمیت مخفی‌سازی برخی ویژگی‌ها از برخی ارجحیت دارد و این ارجحیت در خوشه‌بندی‌های اعمال شده تأثیر داده نمی‌شود مثلاً اهمیت پنهان‌سازی زمان جستجو از لینک کلیک شده بعد از جستجو کمتر است و این ارجحیت می‌تواند در نحوه‌ی شکل یافتن میکرو-دسته‌ها موثر باشد که مقاله به این موضوع نیز پرداخته است.