

مروری بر روش‌های حفظ حریم خصوصی در انتشار

داده‌ها

رضا ابراهیمی آتانی^{۱*} و مهدی صادق پور^۲

^۱دانشیار گروه مهندسی کامپیوتر دانشگاه گیلان، گیلان، رشت، ایران
rebrahimi@guilan.ac.ir

^۲دانش آموخته کارشناسی ارشد مهندسی کامپیوتر، نرم‌افزار دانشگاه گیلان، گیلان، رشت، ایران
Mehdi.sadeghpour@live.com

چکیده

با گسترش خدمات الکترونیکی و سهولت در جمع‌آوری و ذخیره‌سازی داده‌ها، حجم وسیعی از اطلاعات افراد در قالب رکوردهایی درون پایگاه‌داده سازمان‌های خصوصی و دولتی ذخیره شده است. این رکوردها به‌طور معمول شامل اطلاعات حساس شخصی از جمله میزان درآمد و نوع بیماری هستند که باید از خارج از دسترس سایرین نگه داشته شود. اگر مخفی کردن اطلاعات حساس رکوردها به‌درستی انجام گیرد، تأثیر مخربی در سودمندی داده‌ها ندارد؛ زیرا در اکثر مواقع دریافت‌کنندگان داده به‌دنبال مشخصات افراد و داده‌های سطح رکورد نیستند و می‌خواهند قوانین و الگوهای کلی و در سطح مجموعه‌داده به‌دست بیاورند. "حفظ حریم خصوصی در انتشار داده‌ها" راه‌کارهایی برای تضمین محرمانه‌ماندن اطلاعات حساس یک مجموعه‌داده در هنگام انتشار آن در محیطی غیر قابل اعتماد است. در این فرایند سعی بر آن است که ضمن مخفی‌نگه‌داشتن اطلاعات حساس، داده‌های منتشرشده همچنان برای عملیات کشف دانش مفید باقی بماند. این نوشتار با هدف مرور روش‌های حفظ حریم خصوصی داده‌ها به نگارش درآمده است.

واژگان کلیدی: حفظ حریم خصوصی، اشتراک‌گذاری داده‌ها، داده‌کاوی، گمنام‌سازی

۱- مقدمه

طبق ماده ۱۲ اعلامیه جهانی حقوق بشر^۲ (UDHR) قلمرو خصوصی افراد نباید مورد مداخله قرار گیرد و تمام اشخاص حق دارند به‌وسیله قانون در برابر چنین مداخلات و تعرضاتی محافظت شوند [۲]. جدا از قوانینی که برای محافظت از حریم خصوصی وجود دارد و خاطیان را مجازات می‌کند، جلب اعتماد مشتریان عاملی برای احتیاط دوچندان در منتشر کردن داده‌ها است؛ زیرا افراد حاضر نیستند به شرکتی مراجعه کنند که ممکن است اطلاعات شخصی‌شان را در اختیار سایرین قرار دهد.

رکوردهای یک مجموعه‌داده را می‌توان به‌عنوان نمونه‌هایی از یک جامعه آماری در نظر گرفت. این جامعه آماری می‌تواند بیماران یک درمانگاه، کاربران یک سایت اینترنتی، کارمندان یک اداره و یا مثال‌هایی از این دست

توسعه و نفوذ خدمات الکترونیکی در زندگی بشر موجب شده است، سازمان‌ها و شرکت‌های مختلف، حجم گسترده‌ای از اطلاعات مشتریان خود را به‌طور روزمره ذخیره کنند. این داده‌ها منبع با ارزشی برای کشف الگوها و قوانین سودمند به‌وسیله داده‌کاوی هستند؛ به همین دلیل ممکن است پژوهش‌گران، شرکت‌ها و سازمان‌های مرتبط، خواهان دریافت و استفاده از این داده‌ها باشند. در چنین شرایطی چگونه می‌توان بدون نقض کردن حریم شخصی، مجموعه‌داده‌ای را که حاوی اطلاعات مشتریان است، در اختیار نهادهای دیگر قرار داد؟ حفظ حریم خصوصی در انتشار داده‌ها^۱ (PPDP) به‌عنوان راه‌کاری برای حل این مشکل معرفی شده است [۱].

^۲ Universal Declaration of Human Rights (UDHR)

^۱ Privacy-Preserving Data Publishing (PPDP)

- **مالکان رکورد^۱:** افرادی هستند که اطلاعات مربوط به آن‌ها در مجموعه داده ثبت شده است.
 - **منتشرکننده داده‌ها^۲:** سازمانی است که مجموعه داده را در اختیار دارد و اقدام به انتشار آن می‌کند.
 - **مهاجم:** فرض بر این است که بعد از انتشار داده‌ها، مهاجم قادر خواهد بود نسخه‌ای از آن را به دست آورد. وی تلاش می‌کند، اطلاعات حساس یک یا چند مالک رکورد (که قربانی^۳ نامیده می‌شوند) را از مجموعه داده استخراج کند. اطلاعات در دسترس مهاجم فقط به داده‌های منتشرشده محدود نیست و امکان دارد از دانش پیش‌زمینه^۴ خود نیز برای انجام حملات استفاده کند.
 - **دریافت‌کننده داده‌ها^۵:** پژوهش‌گران علمی یا سازمان‌های تجاری که به دنبال کشف دانش از مجموعه داده منتشر شده هستند.
 - **گمنام‌سازی^۶ داده‌ها:** فرآیندی است که داده‌های خام^۷ را به داده‌های گمنام‌شده تبدیل می‌کند. داده‌های گمنام‌شده در مقابل نوع خاصی از حملات که در یک مدل حریم خصوصی^۸ از حملات شده است، از افشای اطلاعات مصون می‌مانند. مدل‌های حفظ حریم خصوصی داده در بخش ۵ تشریح شده‌اند.
- همانطور که در شکل ۱۰ (۱) نمایش داده شده است، دو مرحله اصلی فرایند PPDP جمع‌آوری و انتشار داده‌ها هستند [۵]. در مرحله نخست منتشرکننده داده‌ها اطلاعات

¹ Record Owners

² Data Publisher

³ Victim

⁴ Background Knowledge

⁵ Data Recipient

⁶ Anonymization

⁷ Raw Data

⁸ Privacy Model

باشد. مجموعه داده منتشرشده باید به گونه‌ای باشد که انجام بررسی‌های دلخواه بر روی آن، منجر به آشکارشدن مشخصات یک فرد خاص نشود و فقط اطلاعاتی از کل جامعه قابل یادگیری باشد. این هدف از طریق گمنام‌سازی داده‌ها محقق می‌شود.

با این توضیحات ممکن است، پرسیده شود که اگر هدف نهایی، استخراج یک مدل از کل داده‌ها است، پس چرا به جای انتشار مجموعه داده، آن مدل استخراج و منتشر نگردد؟ در پاسخ به این سؤال باید مدنظر داشت که انتشار مجموعه داده امکان به کارگیری روش‌های مختلف برای کاویدن داده‌ها را فراهم می‌آورد. انتشار داده‌های شرکت Netflix [۳] مثال خوبی در این زمینه است. این شرکت برای بهبود سامانه توصیه‌گر خود، رقابتی عمومی برای معرفی الگوریتم‌های مختلف برپا کرد. در این شرایط انتشار داده‌ها تنها روشی بود که انعطاف لازم برای آزمودن الگوریتم‌های گوناگون را فراهم می‌کرد.

هدف اصلی این نوشتار ارائه یک مقاله مروری برای معرفی روش‌های حفظ حریم خصوصی در انتشار داده‌ها است که بدین ترتیب ادامه خواهد یافت: در بخش ۲، به بررسی PPDP پرداخته و تعاریف مورد نیاز ارائه می‌شود. در بخش ۳ کاربردهای PPDP با ذکر چند مثال بازگو می‌شود. بخش ۴ به معرفی توابع پرکاربرد در گمنام‌سازی می‌پردازد. در بخش ۵ مدل‌های مختلف حفظ حریم خصوصی مرور شده و نقاط قوت و ضعف آن‌ها بررسی می‌شود. و بخش انتهایی به نتیجه‌گیری و جمع‌بندی این نوشتار اختصاص دارد.

۲- حفظ حریم خصوصی در انتشار داده‌ها

فرآیند PPDP راه‌کاری برای تضمین حریم خصوصی داده‌ها در محیط‌های غیر قابل اعتماد است. در این فرایند سعی بر آن است که ضمن مخفی نگه‌داشتن اطلاعات محرمانه، داده‌های منتشرشده همچنان برای عملیات کشف دانش مفید باقی بماند [۴]. در زیر، مفاهیم پایه‌ای PPDP که در ادامه مورد استفاده قرار می‌گیرند، تعریف شده‌اند:

به‌تنهایی برای شناسایی مالک رکورد کافی نیست؛ اما با استفاده از ترکیب آن‌ها، در اغلب موارد می‌توان این کار را انجام داد. مانند ترکیب کدپستی و تاریخ تولد و جنسیت.

- **صفات حساس^۳:** صفاتی که صاحب رکورد مایل نیست عموم از آن‌ها اطلاع داشته باشند؛ مانند: نوع بیماری، میزان حقوق و ...
- **صفات غیر حساس^۴:** صفاتی که در هیچ کدام از دسته‌های قبل جای نمی‌گیرند.

شکل (۲) مثالی را نمایش می‌دهد که در آن صفات یک مجموعه داده، به چهار نوع ذکر شده در بالا تقسیم‌بندی شده‌اند.

کد ملی	نام	تحصیلات	جنسیت	وضعیت تاهل	تاریخ تولد	شبه‌شناسها	
						صفت غیرحساس	صفت حساس
۲۲۲۵۵۵۹	رضا محبی	دیپلم	مرد	متاهل	۱۳۶۴/۵/۲۳	ساعات کار در هفته	میزان حقوق
۲۲۲۷۷۷۱	تیما راهنما	کارشناسی	مرد	مجرد	۱۳۵۹/۲/۹	۴۰	۹۱۰
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

(شکل-۲): مثالی از انواع صفات در PPDP

نخستین قدم برای شکستن رابطه بین رکوردها و مالکان رکورد، حذف شناسه‌های صریح از مجموعه داده است. اما همانطور که در شکل (۳) نمایش داده شده است، حذف این شناسه‌ها لزوماً باعث حفظ حریم خصوصی نمی‌شود [۶]. در این شکل داده‌های پزشکی که فاقد شناسه صریح هستند با استفاده از شبه‌شناسه‌ها (جنسیت، تاریخ تولد و کدپستی) به داده‌هایی که از قبل در اختیار فرد مهاجم بوده است (فهرست رأی‌دهندگان) پیوندهایی^۵ شده است؛ یعنی مهاجم از شبه‌شناسه‌ها برای تشخیص ارتباط بین شناسه صریح مالک رکورد و صفت حساس وی استفاده کرده است.

در مقابله با چنین حملاتی می‌توان با استفاده از عملیات گمنام‌سازی موجب شد تا شبه‌شناسه‌های هر

را از کاربران خود جمع‌آوری می‌کند. فرض بر این است که مالکان رکورد به منتشرکننده داده‌ها اطمینان دارند و حاضرند اطلاعات حساس خود را با او به اشتراک بگذارند، اما دریافت‌کنندگان داده، مورد اعتماد نیستند و اطلاعات حساس باید از آن‌ها مخفی نگاه داشته شود. بنابراین در مرحله بعد روش‌های حفظ حریم خصوصی روی داده‌ها اعمال می‌شود و خروجی در اختیار دریافت‌کنندگان داده قرار می‌گیرد. هر دریافت‌کننده ممکن است داده‌های مورد نیاز خود را از چندین منتشرکننده جمع‌آوری کند. به‌عنوان مثال، یک شرکت بیمه برای انجام داده‌کاوی در زمینه رابطه بین شغل افراد و یک بیماری خاص، نیاز به داده‌های ذخیره‌شده در بیمارستان‌های شهر دارد. در این مثال بیماران مالکان رکورد، بیمارستان‌ها منتشرکنندگان داده و شرکت بیمه دریافت‌کننده داده‌ها است.



(شکل-۱): جمع‌آوری و انتشار داده‌ها

همان‌طور که ذکر شد در فرآیند گمنام‌سازی، پردازش‌هایی بر روی صفات ذخیره‌شده در مجموعه داده انجام می‌گیرد تا در هنگام انتشار عمومی داده‌ها، اطلاعات حساس افشا نشود. پردازشی که روی هر صفت انجام می‌گیرد به نوع آن بستگی دارد. انواع صفات در PPDP عبارتند از [۵]:

- **شناسه صریح^۱:** مجموعه صفاتی که به‌طور مستقیم برای شناسایی مالک رکورد مورد استفاده قرار می‌گیرند، مانند: کد ملی، شماره دانشجویی و ...
- **شبه‌شناسه^۲:** مجموعه‌ای از صفات که هر یک

¹ Explicit Identifier
² Quasi Identifier

³ Sensitive Attributes
⁴ Non Sensitive Attributes
⁵ Linking

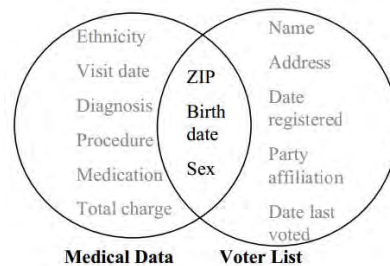
درمانی از مراجعه‌کنندگان ثبت می‌کنند، بهره می‌گیرند. به‌عنوان نمونه ممکن است پژوهش‌گران به دنبال رابطه بین یک بیماری خاص با وزن فرد باشند. داده‌هایی که قرار است در اختیار آنان قرار بگیرد باید به‌گونه‌ای باشد که اطلاعات شخصی بیماران افشا نشود. در ایالات متحده از سال ۱۹۹۶ مجموعه قوانینی با عنوان "قانون انتقال و پاسخ‌گویی بیمه سلامت"^۳ تدوین شده است که بر اساس یکی از بندهای آن مراکز درمانی اجازه ندارند اطلاعات خام بیماران را با سایر سازمان‌ها به اشتراک بگذارند، اما هیچ منعی برای انتشار داده‌های پزشکی گمنام‌شده وجود ندارد [۹].

یکی از سناریوهایی که کاربرد PPDP در مراکز درمانی را نشان می‌دهد، سرویس انتقال خون صلیب سرخ هنگ‌کنگ^۴ است [۱۰]. روند کاری این مرکز در شکل (۴) نمایش داده شده است. ابتدا خون مورد نیاز از اهداکنندگان دریافت و پس از تست سلامت آن، در اختیار بیمارستان‌های مختلف قرار داده می‌شود. بیمارستان‌ها از این خون برای درمان بیماران استفاده می‌کنند. تمام جزئیات نحوه استفاده از خون، مانند مشخصات دریافت‌کننده، نوع عمل انجام‌شده، دلایل تزریق خون و ... در پایگاه داده بیمارستان ثبت می‌شود. مرکز انتقال خون به‌طور دوره‌ای این داده‌ها را از بیمارستان‌ها جمع‌آوری می‌کند و با تحلیل آن اطلاعاتی از قبیل تخمین میزان خون مورد نیاز هر بیمارستان در آینده به‌دست می‌آورد. داده‌هایی که از طریق بیمارستان‌ها در اختیار مرکز انتقال خون قرار می‌گیرد، نباید شامل اطلاعات شخصی بیماران باشد. این امر از طریق گمنام‌سازی داده‌ها محقق می‌شود.



(شکل-۴): روند کار سرویس انتقال خون صلیب سرخ هنگ‌کنگ

رکورد در چندین رکورد دیگر نیز تکرار شود و مهاجم در پیوندهای داده‌ها دچار ابهام شود. پس از اعمال گمنام‌سازی بر روی داده‌های خام، میزانی از اطلاعات دیگر در داده‌های خروجی در دسترس نیست. هرچه میزان حریم خصوصی بیشتر مد نظر باشد، ممکن است، سودمندی داده‌ها پایین‌تر بیاید و هرچه سودمندی بیشتری مد نظر باشد، ممکن است، به کاهش سطح حریم خصوصی بیانجامد [۷]. بنابراین یک روش موفق در حفظ حریم خصوصی داده باید تعادل مناسبی بین این دو پارامتر برقرار کند.



(شکل-۳): پیوندهای به منظور بازشناسی داده‌ها [۶]

۳- کاربردهای حفظ حریم خصوصی در انتشار داده‌ها

استخراج دانش از داده‌های جمع‌آوری‌شده توسط سازمان‌های مختلف همواره مورد توجه پژوهش‌گران بوده است. در همین حال، حفظ حریم خصوصی مالکین رکورد، اشتراک‌گذاری داده‌ها^۱ را با چالش مواجه می‌کند. PPDP به‌عنوان یکی از روش‌های حفظ حریم خصوصی داده‌ها، می‌تواند در بسترهای مختلفی از جمله مراکز درمانی و تجاری، سامانه‌های آگاه از موقعیت^۲، موتورهای جستجو، رسانه‌های اجتماعی و ... مورد استفاده قرار بگیرد [۸، ۱]. در ادامه دو مورد از کاربردهای ذکرشده مورد بحث قرار می‌گیرد.

۳-۱- مراکز درمانی

پژوهش‌گران علم پزشکی، شرکت‌های بیمه و سایر نهادهایی که در حوزه سلامت فعالیت می‌کنند، از داده‌هایی که مراکز

³ Health Insurance Portability and Accountability Act (HIPPA)

⁴ Hong Kong Red Cross Blood Transfusion Service

¹ Data Sharing

² Location-aware Systems

۲-۳- رسانه‌های اجتماعی

در سال‌های اخیر شبکه‌های اجتماعی مانند Facebook و Myspace، برنامه‌های پیام‌رسان مانند Skype و Telegram و سایر رسانه‌های اجتماعی کاربران زیادی را به خود جذب کرده‌اند. همان‌طور که در شکل (۵) به تصویر کشیده شده است، آمارها نشان می‌دهد، حجم اطلاعاتی که هر دقیقه در رسانه‌های اجتماعی تولید می‌شود بسیار درخور توجه است. وجود روابط خاص مانند دوستی^۱، دنبال کردن^۲ و پسندیدن^۳ در این داده‌ها، موقعیت مناسبی برای درک الگوهای رفتاری انسان‌ها ایجاد کرده است. کند و کاو رسانه‌های اجتماعی^۴ (SMM) ترکیبی از علوم مختلف از جمله: داده‌کاوی، جامعه‌شناسی، علوم شبکه^۵ و آمار است، که برای استخراج دانش از داده‌های رسانه‌های اجتماعی^۶ به کار می‌رود [۱۱].

اعتماد در رسانه‌های اجتماعی [۱۲]، نقش رسانه‌های اجتماعی در بحران‌ها [۱۳] و جستجوی اجتماعی [۱۴] از جمله حوزه‌هایی هستند که در SMM مورد بررسی قرار می‌گیرد.

علی‌رغم این ویژگی‌های مفید تعداد بسیار کمی از داده‌های اجتماعی به‌صورت عمومی منتشر شده است که مهم‌ترین دلیل آن نگرانی از افشای اطلاعات خصوصی کاربران این برنامه‌ها است [۱]. PDP با حذف این مانع می‌تواند امکان دسترسی پژوهش‌گران را به این منبع باارزش اطلاعاتی فراهم آورد.



(شکل-۵): فعالیت کاربران در برخی از رسانه‌های اجتماعی طی یک دقیقه (برگرفته از [۱۵])

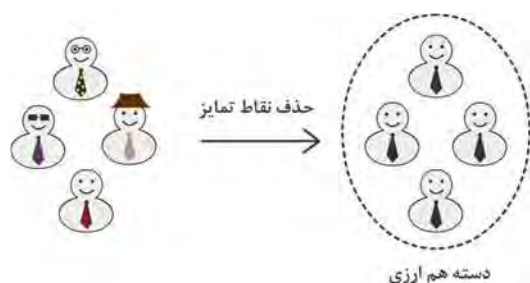
۴- عملگرهای گمنام سازی

در طی فرآیند گمنام‌سازی توابعی بر روی رکوردها اعمال می‌شود تا بتوان مجموعه‌داده را در دسترس عموم قرار داد. این توابع به‌عنوان عملگرهای گمنام‌سازی شناخته می‌شوند. درواقع عملگرهای گمنام‌سازی ابزار لازم برای دستیابی به مدل‌های حفظ حریم خصوصی (بخش ۵) را فراهم می‌کنند. این عملگرها براساس روش کارکردشان به سه گروه تقسیم شده‌اند که در ادامه معرفی می‌شود.

۴-۱- گمنامی از طریق دسته‌سازی

در این روش هر رکورد به همراه چندین رکورد دیگر در گروهی قرار داده می‌شود که مقدار شبه‌شناسه‌های تمام اعضای آن یکسان است. چنین گروهی، گروه هم‌ارزی^۷ نامیده می‌شود. در این حالت مهاجم با جستجو کردن شبه‌شناسه‌های قربانی در مجموعه‌داده، به جای رکوردی یکتا، به گروه هم‌ارزی که رکورد قربانی عضو آن است، دست می‌یابد و درنتیجه با اطمینان کامل صفت حساس وی را نمی‌تواند تشخیص دهد.

هرچه دامنه مقادیر و تعداد شبه‌شناسه‌ها بیشتر باشد، احتمال پیدا کردن دسته‌ای از رکوردها که بتوانند در یک گروه هم‌ارزی قرار بگیرند کمتر می‌شود. به‌همین دلیل لازم است توابعی بر روی رکوردهای اولیه اعمال شود که نقاط تمایز آن‌ها را کمتر کند (شکل ۶). مهم‌ترین توابعی که برای دسته‌سازی رکوردها مورد استفاده قرار می‌گیرند عمومی‌سازی^۸ و فرونشانی^۹ [۱۶] هستند.



(شکل-۶): گمنامی از طریق دسته‌سازی

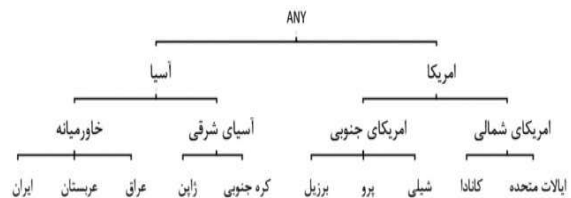
- 1 Friendship
- 2 Follow
- 3 Like
- 4 Social Media Mining (SMM)
- 5 Network Science
- 6 Social Media Data

- 7 Equivalence Group
- 8 Generalization
- 9 Suppression

۱-۱-۴- عمومی سازی

در عمومی سازی مقدار یک صفت به بازه‌های بزرگتر تعمیم داده می‌شود، برای مثال نشانی دقیق محل سکونت یک فرد را می‌توان با شهری که در آن زندگی می‌کند، جایگزین کرد و در مرحله بعد شهر را به استان تعمیم داد. عمومی سازی صفات به کمک درخت طبقه‌بندی^۱ انجام می‌گیرد. در شکل (۷) نمونه‌ای از درخت طبقه‌بندی برای صفت "ملیت" نمایش داده شده است. به منظور عمومی سازی مقدار یک صفت، باید گره متناظر با آن در درخت طبقه‌بندی پیدا و سپس از مقدار گره والد برای آن صفت استفاده شود.

برگ‌های درخت طبقه‌بندی شامل کلیه مقادیری هستند که صفت مورد نظر می‌تواند داشته باشد. گره ریشه نیز نشان‌دهنده عمومی‌ترین مقدار است که به‌طور معمول با "ANY" نمایش داده می‌شود. شکل (۸) مثالی را نمایش می‌دهد که در آن عملگر عمومی سازی روی صفت "ملیت" اعمال شده است و در نتیجه مقدار این شبه‌شناسه در هر دو رکورد یکسان شده است.



(شکل-۷): درخت طبقه‌بندی صفت ملیت

ایجاد درخت طبقه‌بندی برای هر صفت نیازمند داشتن اطلاعات پیش‌زمینه در مورد آن است؛ به عنوان نمونه، برای تشکیل درخت شکل (۷) لازم است بدانیم کشورهای شیلی، پرو و برزیل در آمریکای جنوبی واقع شده‌اند. به همین دلیل به‌طور معمول این کار به صورت غیر خودکار و توسط یک انسان خبره (آشنا با موجودیت‌های سیستم) انجام می‌گیرد، اما برای صفاتی که مقادیر عددی می‌پذیرند ساخت درخت می‌تواند به صورت خودکار و از طریق بازه‌بندی انجام گیرد [۱۷].

^۱ Taxonomy Tree

جنسیت	ملیت	سن	بیماری
مرد	ایران	۲۸	آنفولانزا
مرد	عربستان	۲۳	دیابت

جنسیت	ملیت	سن	بیماری
مرد	خاورمیانه	۲۸	آنفولانزا
مرد	خاورمیانه	۲۳	دیابت

(شکل-۸): اعمال یک مرحله عمومی سازی روی صفت ملیت

۴-۱-۲- فرونشانی

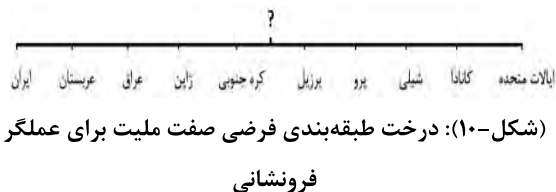
عملگر فرونشانی برای از بین بردن تمایز شبه‌شناسه‌ها، مقدار صفت را حذف می‌کند و به جای آن مقدار گمشده^۲ قرار می‌دهد. شکل (۹) نمونه‌ای از نحوه کار این عملگر را نشان می‌دهد. همان‌طور که دیده می‌شود، برخلاف عمومی سازی که مقادیر قبل و بعد از گننام سازی را از لحاظ معنایی منسجم نگه می‌دارد، در این حالت نمی‌توان هیچ برداشتی از ملیت رکوردها داشت. بنابراین استفاده از این عملگر باید با احتیاط صورت گیرد. در غیر این صورت تعداد مقادیر گمشده زیاد می‌شود و کیفیت داده‌ها به شدت کاهش می‌یابد.

جنسیت	ملیت	سن	بیماری
مرد	ایران	۲۸	آنفولانزا
مرد	عربستان	۲۳	دیابت

جنسیت	ملیت	سن	بیماری
مرد	?	۲۸	آنفولانزا
مرد	?	۲۳	دیابت

(شکل-۹): اعمال فرونشانی روی صفت ملیت

فرونشانی را می‌توان حالت خاصی از عمومی سازی در نظر گرفت که در آن از درخت طبقه‌بندی به عمق یک استفاده می‌شود. به عنوان مثال عمومی سازی شکل (۸) به وسیله درخت نمایش داده در شکل (۱۰) نتیجه‌ای مشابه با فرونشانی در پی خواهد داشت.



(شکل-۱۰): درخت طبقه‌بندی فرضی صفت ملیت برای عملگر فرونشانی

فرونشانی و عمومی سازی به دو روش سراسری^۳ و محلی^۴ مورد استفاده قرار می‌گیرند [۱۸]. در روش

^۲ Missing Value

^۳ Global

^۴ Local

موجب می‌شود استفاده از روش تجزیه با مشکل مواجه شود [۲۰].

بیماری	سن	جنسیت
دیابت	۳۳	مرد
دیابت	۳۵	مرد
هیپاتیت	۳۳	مرد
هیپاتیت	۳۶	زن
دیابت	۴۰	زن

گروه	سن	جنسیت
۱	۳۳	مرد
۱	۳۵	مرد
۱	۳۳	مرد
۲	۳۶	زن
۲	۴۰	زن

بیماری	سن	جنسیت
دیابت	۳۳	مرد
دیابت	۳۵	مرد
هیپاتیت	۳۳	مرد
هیپاتیت	۳۶	زن
دیابت	۴۰	زن

شکل-۱۱): تجزیه مجموعه داده به دو جدول ST و QIT

۴-۲-۲-۲- برش‌دهی

روش برش‌دهی [۲۱] مجموعه داده را به دو صورت افقی و عمودی تقسیم‌بندی می‌کند. در برش‌دهی عمودی، صفت‌هایی که مقدارشان وابستگی زیادی به هم دارند شناسایی شده و در ستونی مجزا قرار می‌گیرند. در برش‌دهی افقی، مجموعه داده به زیرمجموعه‌هایی بدون اشتراک از رکوردها تقسیم می‌شود؛ سپس مقادیر موجود در هر ستون از هر زیر مجموعه به صورت تصادفی جابه‌جا می‌شوند. شکل (۱۲) نمونه‌ای از عملکرد برش‌دهی را نمایش می‌دهد. "تجزیه" را می‌توان حالت خاصی از برش‌دهی دانست که در آن شبه‌شناسه‌ها در یک ستون و صفات حساس در ستون دیگر قرار می‌گیرند.

با استفاده از این روش سودمندی داده تا حد زیادی حفظ می‌شود، زیرا الگویی که صفات وابسته بر طبق آن در کنار هم قرار می‌گیرند، دست نخورده باقی می‌ماند. این روش برای حفظ حریم خصوصی داده نیز موثر است، زیرا صفاتی که قرار گرفتن آن‌ها کنار هم کم تکرار است و برای شناسایی مالک رکورد مفید می‌تواند باشد، از هم جدا می‌شوند.

بیماری	سن	جنسیت	کدپستی، بیماری
سوه هاضمه	۴۹۷۰۶	مرد	(۲۲، مرد) (۴۷۹۰۵، آنفولانزا)
آنفولانزا	۴۷۹۰۶	زن	(۲۲، زن) (۴۷۹۰۶، سوه هاضمه)
آنفولانزا	۴۷۹۰۵	زن	(۳۳، زن) (۴۷۹۰۵، پرونشیت)
پرونشیت	۴۷۹۰۵	زن	(۵۲، زن) (۴۷۹۰۶، آنفولانزا)
آنفولانزا	۴۷۳۰۲	مرد	(۵۴، مرد) (۴۷۳۰۴، ورم معده)
سوه هاضمه	۴۷۳۰۲	مرد	(۶۰، مرد) (۴۷۳۰۲، آنفولانزا)
سوه هاضمه	۴۷۳۰۴	مرد	(۶۰، مرد) (۴۷۳۰۲، سوه هاضمه)
ورم معده	۴۷۳۰۴	زن	(۶۴، زن) (۴۷۳۰۴، سوه هاضمه)

شکل-۱۲): نمونه‌ای از نحوه عملکرد روش برش‌دهی [۲۱]

سراسری وقتی عملگر روی صفت a از رکورد r اعمال می‌شود و مقدار آن را از v به v' تغییر می‌دهد، باعث می‌شود صفت a در تمام رکوردهایی که در آن‌ها مقدار v دارد، به v' تغییر مقدار دهد. تأثیر روش محلی در سطح رکورد است و فقط مقدار یک صفت را تغییر می‌دهد.

۴-۲-۴- شکستن رابطه بین صفت‌ها

برخلاف عملگرهای گمنام‌سازی که در بخش قبل معرفی شد، در این روش مقدار شبه‌شناسه‌ها تغییر نمی‌یابد. این عملگرها از طریق تفکیک صفات مجموعه داده به جداول مختلف موجب می‌شوند، مهاجم از روی شبه‌شناسه‌های یک رکورد نتواند به صفت حساس آن دست پیدا کند. در ادامه، تجزیه^۱ و برش‌دهی^۲ که دو روش شناخته شده در این زمینه هستند، معرفی می‌شوند.

۴-۲-۴-۱- تجزیه

تجزیه [۱۹] مجموعه داده اولیه را به دو جدول تقسیم می‌کند که یکی از آنها شامل صفات حساس (جدول ST) و دیگری دربردارنده شبه‌شناسه‌ها (جدول QIT) است. همانطور که در شکل (۱۱) قابل مشاهده است، برای مرتبط‌سازی ST و QIT صفت جدیدی به نام "گروه" به هر دو جدول اضافه می‌شود. در واقع پس از تجزیه، به جای اینکه صفات حساس به رکوردها تعلق داشته باشند به گروه‌ها نسبت داده می‌شوند. بدین ترتیب ارتباط مستقیم بین شبه‌شناسه‌ها و صفت حساس یک رکورد شکسته می‌شود.

تجزیه مقدار شبه‌شناسه‌ها را تغییر نمی‌دهد. این ویژگی در کاربردهایی مانند پاسخ‌دهی به پرس‌وجو بسیار با ارزش است. در مقابل، انتشار داده‌ها در دو جدول جداگانه استفاده از روش‌های استاندارد داده‌کاوی را دشوار می‌سازد. همچنین ممکن است در برخی سناریوها صفات حساس به سادگی از شبه‌شناسه‌ها قابل تفکیک نباشد که این امر

¹ Anatomization

² Slicing

۴-۳- تغییر مقدار داده‌ها

در این روش‌ها مقادیر داده با مقدار غیر واقعی دیگری جایگزین می‌شوند. تغییر مقادیر باید به گونه‌ای انجام گیرد که اطلاعات آماری محاسبه‌شده از داده‌های خام و گمنام‌شده تفاوت چندانی نداشته باشد. در این حالت در سطح رکورد نمی‌توان به داده‌ها اطمینان کرد، یعنی رکوردها ممکن است مربوط به افراد حقیقی نباشند.

اضافه کردن نوفه به داده‌ها یکی از روش‌های رایج در این زمینه است که برای مخفی کردن صفات عددی مورد استفاده قرار می‌گیرد [۲۲]. در این روش، مقدار v با $v + r$ تعویض می‌شود که r یک مقدار تصادفی است. r باید طوری انتخاب شود که هم مقادیر تولیدشده به اندازه کافی با مقادیر اصلی فاصله داشته باشند و هم اطلاعات آماری مانند میانگین صفت دست‌نخورده باقی بماند. تولید داده‌های مصنوعی^۱، روشی دیگر برای مخفی‌نگه‌داشتن مقدار واقعی داده‌ها است. بدین منظور ابتدا مدلی آماری از روی داده‌ها ساخته می‌شود؛ سپس از روی مدل بدست‌آمده تعدادی رکورد تولید شده و به جای داده‌های اصلی منتشر می‌شود [۲۳].

۵- مدل‌های حفظ حریم خصوصی

پیش از اینکه بخواهیم برای حفظ حریم خصوصی تلاش کنیم، باید مشخص کنیم چه زمانی این هدف برآورده می‌شود. تعریف ایده‌آل از حفظ حریم خصوصی هنگام انتشار داده‌ها، در [۲۴] ارائه شده است. بر طبق این تعریف، بعد از آنکه داده‌ها در اختیار مهاجم قرار می‌گیرد، نسبت به زمانی که داده‌ها را در اختیار نداشته است، نباید هیچ دانشی درمورد مالکان رکورد به او اضافه شود. با در نظر گرفتن اطلاعاتی که ممکن است از قبل در اختیار مهاجم باشد، این سطح از حریم خصوصی قابل دستیابی نیست [۲۵] به همین دلیل در بیشتر کارهایی که در حوزه حریم خصوصی داده انجام می‌گیرد از مدل‌های تعدیل

^۱ Synthetic Data Generation

یافته‌تری استفاده می‌شود که در ادامه به شرح آن‌ها پرداخته می‌شود.

۵-۱- مدل پیونددهی رکورد

در حملات پیونددهی رکورد مهاجم شبه‌شناسه‌های قربانی را با رکوردهای انتشاریافته تطبیق می‌دهد. بعد از انجام این کار دسته‌ای کوچک از رکوردها باقی می‌ماند که ممکن است، متعلق به قربانی باشد. پیونددهی رکورد زمانی رخ می‌دهد که مهاجم موفق شود رکورد مربوط به قربانی را به‌طور یکتا شناسایی کند [۱].

در شکل (۱۳) صفات تحصیلات، جنسیت و سن به عنوان شبه‌شناسه‌ها و حقوق به‌عنوان صفت حساس در نظر گرفته شده است. پس از انتشار این جدول اگر مهاجم به دنبال رکوردی با شبه‌شناسه‌های <کارشناسی، مرد، ۲۹> باشد، به راحتی می‌تواند آن را پیدا کند، زیرا تنها یک رکورد متناظر وجود دارد و در نتیجه حمله پیونددهی رکورد با موفقیت انجام می‌گیرد.

روش k -anonymity [۶] برای جلوگیری از چنین حملاتی، مجموعه داده را به دسته‌های هم‌ارزی با اندازه k تقسیم می‌کند. در این حالت به‌ازای هر مقدار ممکن از شبه‌شناسه‌ها، یا هیچ رکوردی در مجموعه داده وجود ندارد و یا k رکورد (و بیشتر) موجود است، در نتیجه مهاجم حداکثر با درجه اطمینان $1/k$ می‌تواند رکورد قربانی را شناسایی کند. این روش در قسمت "ب" از شکل (۱۳) برای گمنام‌سازی مورد استفاده قرار گرفته است. همانطور که دیده می‌شود هر یک از رکوردهای این جدول دست‌کم با یک رکورد دیگر شبه‌شناسه‌های یکسانی دارد. به چنین جدولی 2-anonymous می‌گویند. لازم به ذکر است که در کاربردهای عملی مقادیر بزرگتری برای k در نظر گرفته می‌شود.

در برخی پایگاه داده‌ها مانند فهرست بیماران یک درمانگاه، ممکن است n رکورد متعلق به یک فرد باشد. در چنین مواردی پس از اعمال k -anonymity این رکوردها در

برای مقابله با حمله پیونددهی صفت است. برای تأمین I-diversity در یک جدول، در هر دسته هم‌ارزی باید دست کم I مقدار مختلف برای صفت حساس موجود باشد. این خاصیت سبب می‌شود تا مهاجم بعد از پیدا کردن دسته‌ای که قربانی به آن تعلق دارد، با مقادیر یکسانی از صفت حساس مواجه نشود.

اگر در یک دسته هم‌ارزی I مقدار متفاوت، ولی نزدیک به هم وجود داشته باشد باز هم مشکل عدم تنوع بروز پیدا می‌کند. برای مثال قسمت "ب" از شکل (۱۴) مدل 3-diversity را ارضا می‌کند؛ اما مقادیر صفت حساس در نخستین گروه هم‌ارزی بسیار نزدیک به هم هستند. برای رفع این نقص، در [۲۸] روشی به نام t-Closeness ارائه شد که بر طبق آن تفاوت توزیع صفت حساس در هر دسته هم‌ارزی نسبت به توزیع آن در کل جدول باید از حد آستانه t کمتر باشد. این روش از فاصله Earth Mover [۲۹] برای سنجش تفاوت بین توزیع‌ها استفاده می‌کند. در شکل (۱۴) نتیجه گمنام‌سازی یک مجموعه داده به روش‌های I-diversity و t-Closeness مقایسه شده است.

برقرار کردن t-Closeness در یک جدول، نیازمند استفاده زیاد از عملگرهای گمنام‌سازی است که باعث کاهش قابل توجه در کیفیت داده‌ها می‌شود. برای برطرف کردن این ضعف روشی انعطاف‌پذیرتر به نام (n,t)-Closeness [۳۰] ارائه شده است. اگر به‌ازای هر گروه هم‌ارزی E1 مجموعه‌ای از رکوردها با نام E2 موجود باشد که:

- بیش از n رکورد داشته باشد.
- E1 زیرمجموعه آن باشد.
- فاصله توزیع صفت حساس در E1 و E2 کمتر از t باشد.

آنگاه آن مجموعه داده (n,t)-Closeness را ارضا می‌کند. در این روش این هر چه n کوچکتر در نظر گرفته شود کیفیت داده‌ها کمتر لطمه می‌بیند و در حالتی که n برابر با تعداد کل رکوردها باشد معادل روش t-Closeness است.

یک دسته هم‌ارزی قرار می‌گیرند. در چنین شرایطی مهاجم با درجه اطمینان n/k رکورد قربانی را به دست می‌آورد که با نیاز تعریف شده در k-anonymity مطابقت ندارد. برای حل این مشکل می‌توان از مفهومی عمومی‌تر به نام (X,Y)-anonymity [۲۶] استفاده کرد.

در این روش X و Y دو زیر مجموعه از صفات مجموعه داده هستند که اشتراکی با هم ندارند. برای برآورده شدن حریم خصوصی در (X,Y)-anonymity، هر مقدار در X باید دست کم با k مقدار متمایز از Y ارتباط داشته باشد. اگر Y را شناسه و X را مجموعه شبه‌شناسه‌ها در نظر بگیریم، با انتخاب مقدار مناسب برای k می‌توان حداکثر تعداد رکوردهای متعلق به یک فرد در یک دسته هم‌ارزی را کنترل کرد.

ب) 2-anonymous

حقوق	سن	جنسیت	تحصیلات
۱۰۰۰	[۲۵-۳۰]	مرد	ارشد و بالاتر
۱۰۵۰	[۳۰-۳۵]	مرد	پایین تر از ارشد
۹۰۰	[۳۰-۳۵]	مرد	پایین تر از ارشد
۱۱۰۰	[۲۵-۳۰]	مرد	ارشد و بالاتر
۹۵۰	[۲۵-۳۰]	زن	پایین تر از ارشد
۹۵۰	[۲۵-۳۰]	زن	پایین تر از ارشد
۹۵۰	[۲۵-۳۰]	زن	پایین تر از ارشد

الف) مجموعه داده خام

حقوق	سن	جنسیت	تحصیلات
۱۰۰۰	۲۶	مرد	کارشناسی ارشد
۱۰۵۰	۳۰	مرد	کارشناسی
۹۰۰	۳۴	مرد	کاردانی
۱۱۰۰	۲۸	مرد	دکترا
۹۵۰	۲۵	زن	کارشناسی
۹۵۰	۲۷	زن	کارشناسی
۹۵۰	۲۶	زن	کاردانی

شکل-۱۳: جدول 2-anonymous

۵-۲- مدل پیونددهی صفت

همانطور که گفته شد در حمله پیونددهی رکورد مهاجم در صدد تشخیص رکورد مربوط به قربانی است. اما در بیش‌تر موارد به دست آوردن صفت حساس قربانی، مهاجم را به هدف خود می‌رساند و نیازی به شناسایی رکورد نیست [۱]. برای مثال در شکل (۱۳) مقدار صفت حساس در سه رکورد آخر برابر است و مهاجم با در اختیار داشتن شبه‌شناسه‌های <کاردانی، زن، ۲۶> بدون آنکه رکورد قربانی را شناسایی کند، از میزان حقوق وی مطلع می‌شود. این مشکل ناشی از عدم تنوع در مقدار صفت حساس در یک دسته هم‌ارزی است.

I-diversity [۲۷] یکی از روش‌های مورد استفاده

T*		
کدبستی	سن	ملیت
۴۷*	*	امریکا
۴۷*	*	امریکا
۴۷*	*	امریکا
۴۸*	*	اروپا
۴۸*	*	اروپا

P			
نام	کدبستی	سن	ملیت
Alice	۴۷۹۰۶	۳۵	ایالات متحده
Bob	۴۷۹۰۳	۵۹	کانادا
Chris	۴۷۹۰۶	۴۲	ایالات متحده
Dirk	۴۷۶۲۰	۱۸	برزیل
Eunice	۴۷۶۲۰	۲۲	برزیل
Frank	۴۷۶۳۳	۶۳	یرو
Gail	۴۸۹۷۳	۳۳	اسپانیا
Harry	۴۸۹۷۲	۴۷	بلغارستان
Iris	۴۸۹۷۰	۵۳	فرانسه

0.167-closeness (ب)				3-diverse (ب)				الف) مجموعه داده خام			
حقوق	سن	کدبستی	حقوق	حقوق	سن	کدبستی	حقوق	کدبستی	سن	حقوق	
۳۰۰۰	≤ ۴۰	۴۷۶۷*	۱	۳۰۰۰	۲*	۴۷۶**	۱	۳۰۰۰	۲۹	۴۷۶۷۷	
۵۰۰۰	≤ ۴۰	۴۷۶۷*	۲	۴۰۰۰	۲*	۴۷۶**	۲	۴۰۰۰	۲۲	۴۷۶۰۲	
۹۰۰۰	≤ ۴۰	۴۷۶۷*	۳	۵۰۰۰	۲*	۴۷۶**	۳	۵۰۰۰	۲۷	۴۷۶۷۸	
۶۰۰۰	≥ ۴۰	۴۷۹۰*	۴	۶۰۰۰	≥ ۴۰	۴۷۹**	۴	۶۰۰۰	۴۳	۴۷۹۰۵	
۱۱۰۰۰	≥ ۴۰	۴۷۹۰*	۵	۱۱۰۰۰	≥ ۴۰	۴۷۹**	۵	۱۱۰۰۰	۵۳	۴۷۹۰۹	
۸۰۰۰	≥ ۴۰	۴۷۹۰*	۶	۸۰۰۰	≥ ۴۰	۴۷۹**	۶	۸۰۰۰	۴۷	۴۷۹۰۶	
۴۰۰۰	≤ ۴۰	۴۷۶۰*	۷	۷۰۰۰	۳*	۴۷۶**	۷	۷۰۰۰	۳۰	۴۷۶۰۵	
۷۰۰۰	≤ ۴۰	۴۷۶۰*	۸	۹۰۰۰	۳*	۴۷۶**	۸	۹۰۰۰	۳۶	۴۷۶۲۳	
۱۰۰۰۰	≤ ۴۰	۴۷۶۰*	۹	۱۰۰۰۰	۳*	۴۷۶**	۹	۱۰۰۰۰	۳۲	۴۷۶۰۷	

(شکل-۱۴): مقایسه روش‌های t-Closeness و I-diversity

(شکل-۱۵): روش δ-Presence

همان‌طور که گفته شد، در روش δ-Presence مهاجم حداکثر δ٪ از وجود رکورد قربانی در جدول منتشر شده اطمینان دارد، بنابراین احتمال این که بتواند با موفقیت حمله پیونددهی رکورد یا صفت را انجام دهد نیز حداکثر δ٪ است؛ لذا این روش برای مقابله با حملات پیونددهی رکورد و صفت نیز کارایی دارد. عیب عمده روش δ-Presence این است که در محاسبه احتمالات فرض می‌شود منتشرکننده داده و مهاجم به یک جدول خارجی مشترک (مجموعه داده P) دسترسی دارند. در حالی که چنین فرضی در دنیای واقعی زیاد محتمل نیست.

۵-۴- مدل احتمالاتی

در مدل‌های احتمالاتی هدف جلوگیری از پیونددهی رکورد قربانی به یک رکورد، صفت و یا جدول خاص نیست، بلکه اکثر این مدل‌ها تلاش می‌کنند تغییر دانش مهاجم، قبل و بعد از مشاهده داده منتشر شده، به مقدار کوچکی محدود شود.

Differential Privacy [۲۵] یکی از سخت‌گیرانه‌ترین

مدل‌های حفظ حریم خصوصی است. بر طبق این مدل، انتشار یا عدم انتشار رکورد قربانی نباید در نتیجه بررسی‌هایی که مهاجم روی مجموعه داده انجام می‌دهد تأثیر قابل توجهی بگذارد. به بیان ریاضی، اگر بررسی‌های مهاجم توسط تابع تصادفی^۱ F مدل شود، به‌ازای هر مجموعه داده D و D' که در وجود یک رکورد اختلاف دارند باید رابطه زیر برقرار باشد:

^۱ Randomized Function

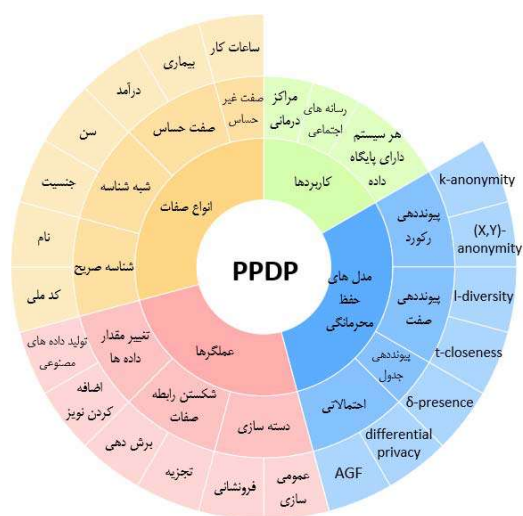
۵-۳- مدل پیونددهی جداول

در برخی شرایط، صرف تشخیص این که رکورد قربانی در مجموعه داده انتشار یافته وجود دارد یا خیر، می‌تواند نقض حریم خصوصی تلقی شود. به‌عنوان نمونه زمانی که یک بیمارستان داده‌های مربوط به اشخاصی که دچار بیماری خاصی هستند را منتشر می‌کند، پی بردن به هویت هر یک از افراد، نوع بیماری وی را آشکار می‌سازد. در زمینه مدل پیونددهی جداول نسبت به مدل پیونددهی رکورد و صفت، کارهای کمتری انجام شده است که معروف‌ترین آن‌ها δ-Presence است [۳۱].

طبق این مدل، احتمال وجود رکورد قربانی در مجموعه داده انتشار یافته باید به سقف δ٪ محدود شود. این مفهوم با مثالی که در شکل (۱۵) آورده شده است بیشتر شرح داده می‌شود. فرض می‌شود جدول P از قبل در دسترس عموم قرار دارد. جدول T زیرمجموعه‌ای از رکوردهای P است که مقدار صفت حساس در آن‌ها یکسان است. این جدول پس از گمنام‌سازی با نام T* منتشر می‌شود. T* نسخه گمنام‌سازی شده T است که پس از انتشار آن مهاجم نباید با احتمال بیشتر از d قادر باشد که بگوید یکی از رکوردهای P در T* وجود دارد. احتمال این که هریک از رکوردهای a-f در T* حضور داشته باشند برابر با $\frac{1}{|T^*|}$ است. این احتمال برای رکوردهای g-i برابر با $\frac{1}{|T^*|}$ است. بنابراین مهاجم حداکثر با احتمال - می‌تواند از وجود یک رکورد دلخواه از P در جدول T* اطمینان داشته باشد.

ارزشی برای استخراج دانش به وسیله داده‌کاوی هستند. این داده‌ها به‌طور معمول شامل اطلاعات حساسی از موجودیت‌های سامانه هستند. به‌همین دلیل پیش از این‌که داده‌ها منتشر شوند، باید هرگونه ارتباط صریح بین مالکان رکورد و صفات حساس از آن‌ها حذف شود. PPDP شاخه‌ای از حفظ حریم خصوصی داده است که به بررسی روش‌های انتشار یک مجموعه داده گمنام می‌پردازد. در این مقاله فرایند PPDP معرفی شد و برخی از کاربردهای آن ذکر شدند. شکل (۱۶) خلاصه‌ای از مطالب مطرح‌شده را نمایش می‌دهد.

بخشی از شکل به دسته‌بندی عملگرهای گمنام‌سازی و ذکر مثال‌هایی از هر کدام مربوط است. تقسیم رکوردها در گروه‌های هم‌ارزی، شکستن رابطه بین صفات رکوردها و تعویض مقدار داده‌ها با مقادیر غیر واقعی، سه روش عمده عملگرهای گمنام‌سازی هستند. در جدول (۱) مقایسه‌ای بین عملگرهای ذکرشده انجام گرفته است.



(شکل-۱۶): نمای کلی از حفظ حریم خصوصی در انتشار داده‌ها

همان‌طور که گفته شد دستیابی به حریم خصوصی مطلق غیر ممکن است، یعنی نمی‌توان کاری کرد که پس از انتشار مجموعه داده، مهاجم هیچ‌گونه اطلاعاتی از مالکین رکورد به دست نیاورد. به‌همین دلیل با توجه به کاربرد، مدل‌های مختلفی برای سطوح متفاوت از حفظ حریم خصوصی

$\forall S \in Range(F) (\Pr[F(D) = S] \leq e^\epsilon \times \Pr[F(D') = S])$ که s صفت حساس و ϵ پارامتری است که هرچه کوچک‌تر انتخاب شود، میزان اطلاعات تأثیر یک رکورد بر نتیجه بررسی‌ها کمتر می‌شود. Differential Privacy تعریف محکمی از حریم خصوصی ارائه می‌دهد که تفاوتی میان صفات حساس و سایر صفات قائل نیست. دستیابی به چنین هدفی مستلزم اضافه کردن نوفه فراوان به داده‌ها است که موجب کاهش کیفیت مجموعه داده می‌شود. در مقاله [۳۲] با در نظر داشتن این مشکل، مدلی تعدیل یافته‌تر معرفی شده است: اطمینانی که مهاجم پس از مشاهده مجموعه داده گمنام‌شده نسبت به صفت حساس قربانی بدست می‌آورد نباید بیشتر از یک حدس تصادفی باشد.

مدل حفظ حریم خصوصی معرفی شده در [۳۲] (AGF) بر اساس دو مفهوم اطمینان مورد انتظار^۱ و اطمینان مشاهده شده^۲ شکل گرفته است. اطمینان مشاهده شده یک رکورد، بیانگر این است که مهاجم با چه احتمالی می‌تواند صفت حساس آن را به قربانی نسبت دهد. هرچه اطمینان مشاهده شده رکوردها کمتر باشد مهاجم با اطمینان کمتری می‌تواند بین صفات حساس و شناسه‌های صریح رابطه برقرار کند، بنابراین باید مقدار آن را به سقف معینی محدود کرد.

فرض می‌شود مجموعه داده شامل n رکورد باشد و D_0 نیز مجموعه داده‌ای با اندازه n است که به صورت تصادفی از فضای نمونه رکوردها ایجاد شده است. برطبق این مدل، اطمینان مشاهده شده یک رکورد در مجموعه داده گمنام شده، نباید از احتمال دیده شدن آن رکورد در D_0 بیشتر باشد. این احتمال اطمینان مورد انتظار نامیده می‌شود. برای تأمین حریم خصوصی طبق قاعده ذکرشده، به‌ازای هر رکورد موجود در مجموعه داده، باید رابطه زیر برقرار باشد:

$$\text{Observed Confidence}(r) \leq \text{Expected Confidence}(r)$$

۶- جمع‌بندی و نتیجه‌گیری

پایگاه داده سازمان‌ها، درمانگاه‌ها و شرکت‌های بزرگ منابع با

¹ Expected Confidence

² Observed Confidence

گرافی و ... اشاره کرد. یکی دیگر از مسائل چالش برانگیز که کم اهمیت در نظر گرفته شده‌اند، به اطلاعات حساس دست خواهد یافت. اگر تعداد شبه‌شناسه‌ها زیاد انتخاب شود، کیفیت داده‌های خروجی کمتر می‌شود. ارائه راه‌حل در PPDP نحوه انتخاب شبه‌شناسه‌ها است. اگر صفات کمی به‌عنوان شبه‌شناسه انتخاب شوند مهاجم از طریق صفاتی برای مسائلی از این دست، لزوم پژوهش بیشتر در این زمینه را ایجاد می‌کند.

پیشنهاد شده است. در جدول (۲) مدل‌هایی که در این نوشتار معرفی شدند، بر اساس ویژگی‌ها و نقاط ضعف هر کدام مقایسه شده‌اند.

از جمله مباحث باز پژوهشی در حوزه PPDP می‌توان به ارائه مدل‌های حفظ حریم خصوصی جدید، روش‌های جدید گمنام‌سازی که مجموعه داده خروجی آن‌ها برای عمل داده‌کاوی خاصی بهینه شده باشد و بررسی چگونگی گمنام‌سازی مجموعه داده‌هایی با ساختار متفاوت (داده‌های توزیع شده بین چند عضو، پایگاه داده

(جدول-۱): مقایسه عملگرهای گمنام‌سازی

نام روش	نحوه عملکرد	رکوردها معرف یک نمونه واقعی هستند؟	قابلیت استفاده از روش‌های استاندارد داده‌کاوی	میزان تغییر مقدار صفات
عمومی‌سازی	تعمیم مقدار صفت به بازه‌ای بزرگتر	بله	بله	نسبتاً کم
فرونشانی	حذف مقدار صفت	بله	بله	نسبتاً زیاد
تجزیه	جدا کردن صفات حساس از شبه‌شناسه‌ها	تا حدودی	خیر	بدون تغییر
برش‌دهی	جدا کردن صفات مرتبط به هم تعویض مقدار آن‌ها در رکوردهای مختلف	خیر	بله	بدون تغییر
اضافه کردن نوفه	تولید داده‌های جدید به وسیله افزودن یک مقدار تصادفی به داده‌های اولیه	خیر	بله	نسبتاً کم
تولید داده‌های مصنوعی	تولید داده‌های جدید با توجه به خصوصیات آماری داده‌های اولیه	خیر	بله	کاملاً عوض می‌شود

(جدول-۲): مقایسه مدل‌های حفظ حریم خصوصی

نام مدل	رابطه‌ای که در مجموعه داده باید برقرار شود	ویژگی‌ها	نقطه ضعف
k-anonymity	گروه‌بندی رکوردها در دسته‌های هم‌ارزی به اندازه k	محافظت در برابر پیونددهی رکورد، تخریب کم	چند رکورد از یک فرد در یک دسته هم‌ارزی قرار می‌گیرد
(X,Y)-anonymity	هر مقدار در زیرمجموعه X حداقل با k مقدار متمایز از Y ارتباط داشته باشد	رفع ضعف ذکر شده برای k-anonymity	مقدار صفت حساس در اعضای یک دسته هم‌ارزی می‌تواند یکسان باشد
l-diversity	وجود l مقدار متنوع از صفت حساس در هر دسته هم‌ارزی	محافظت در برابر پیونددهی صفت، تخریب کم	l مقدار متنوع می‌توانند نزدیک به هم باشند
t-closeness	تفاوت توزیع صفت حساس در هر دسته هم‌ارزی نسبت به توزیع آن در کل جدول از t کمتر باشد	رفع ضعف ذکر شده برای l-diversity	تخریب زیاد
(n,t)-closeness	هر گروه هم‌ارزی ابرمجموعه‌ای با بیش از n رکورد داشته باشد که فاصله توزیع صفت حساس در آن دو کمتر از t است.	رفع ضعف ذکر شده برای t-closeness	امکان وجود چند صفت حساس در مجموعه داده در نظر گرفته نشده است
δ -presence	احتمال وجود رکورد قربانی در مجموعه داده انتشار یافته به $\delta\%$ محدود شود	محافظت در برابر پیونددهی رکورد، صفت، جدول	فرض دسترسی منتشرکننده داده و مهاجم به یک جدول خارجی مشترک (P)
differential privacy	انتشار یا عدم انتشار یک رکورد، در نتیجه بررسی‌هایی که روی مجموعه داده انجام می‌گیرد تأثیر قابل توجهی نداشته باشد	نزدیک‌ترین مدل به تعریف ایده‌آل از حریم خصوصی	تخریب زیاد
AGF	اطمینان مشاهده‌شده رکوردها کوچکتر از اطمینان مورد انتظار باشد	رفع ضعف ذکر شده برای differential privacy	فرض مستقل بودن مقدار صفات

- Connected World," Proceedings of the fifth ACM international conference on Web search and data mining, pp. 93-102, 2012.
- [13] S. Kumar, X. Hu and H. Liu, "A Behavior Analytics Approach to Identifying Tweets from Crisis Regions," Proceedings of the 25th ACM conference on Hypertext and social media, pp. 255-260, 2014.
- [14] D. Horowitz and S. Kamvar, "The Anatomy of a Large-Scale Social Search Engine," In Proceedings of the 19th international conference on World wide web, pp. 431-440, 2010.
- [15] "What happens in an Internet minute?," [Online]. Available: <http://www.intel.com/content/www/us/en/communications/internet-minute-infographic.html>. [Accessed 24 8 2015].
- [16] P. Samarati, "Protecting Respondents' Identities in Microdata Release," IEEE Transactions on Knowledge and Data Engineering, vol. 13, no. 6, pp. 1010-1027, 2001.
- [17] B. Fung, K. Wang and P. Yu, "Anonymizing Classification Data for Privacy Preservation," IEEE Transactions on Knowledge and Data Engineering, vol. 19, no. 5, pp. 711-725, 2007.
- [18] K. LeFevre, D. J. DeWitt and R. Ramakrishnan, "Incognito: Efficient Full-Domain K-Anonymity," Proceedings of the 2005 ACM SIGMOD international conference on Management of data, pp. 49-60, 2005.
- [19] X. Xiao and Y. Tao, "Anatomy: Simple and effective privacy preservation," In Proc. of the 32nd International Conference on Very Large Data Bases, pp. 139-150, 2006.
- [20] V. S. Susan and T. Christopher, "A Survey on Privacy Preservation in," International Journal of Computer Science and Mobile Computing, vol. 3, no. 3, pp. 188-193, 2014.
- [21] T. Li, N. Li, J. Zhang and I. Molloy, "Slicing: A New Approach for Privacy Preserving Data Publishing," IEEE Transactions on Knowledge and Data Engineering, vol. 24, no. 3, pp. 561 - 574, 2012.
- [22] R. Brand, "Microdata Protection through Noise Addition," In Inference Control in
- [1] B. C. M. Fung, K. Wang, A. Wai-Chee Fu and P. S. Yu, Introduction to Privacy-Preserving Data Publishing: Concepts and Techniques, Chapman and Hall/CRC, 2010.
- [2] "Universal Declaration of Human Rights," [Online]. Available: <http://www.un.org/en/documents/udhr/>.
- [3] J. Bennett and S. Lanning, "The Netflix Prize," Proceedings of the KDD Cup Workshop, pp. 3-6, August 2007.
- [4] D. Nettleton, "Data Privacy and Privacy-Preserving Data Publishing," in Commercial Data Mining: Processing, Analysis and Modeling for Predictive Analytics Projects, Morgan Kaufmann, 2014, pp. 266-277.
- [5] B. Fung, K. Wang and P. Yu, "Privacy-Preserving Data Publishing: A Survey of Recent Developments," ACM Computing Surveys, vol. 42, no. 4, 2010.
- [6] L. Sweeney, "k-Anonymity: A Model for Protecting Privacy," International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, vol. 10, no. 5, pp. 557-570, 2002.
- [7] A. Korolova, "Protecting Privacy When Mining and Sharing User Data," PhD thesis, 2012.
- [8] K. S. Babu, "Utility-Based Privacy Preserving Data Publishing," PhD thesis, 2013.
- [9] "Summary of the HIPAA Privacy Rule," [Online]. Available: <http://www.hhs.gov/ocr/privacy/hipaa/understanding/summary/index.html>. [Accessed 22 8 2015].
- [10] N. Mohammed, B. C. M. Fung, P. C. K. Hung and C. K. Lee, "Healthcare Data: A Case Study on the Blood Transfusion Service," Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 1285-1294, June 2009.
- [11] R. Zafarani, M. A. Abbasi and H. Liu, Social Media Mining An Introduction, Cambridge University Press, 2014.
- [12] J. Tang, H. Gao and H. Liu, "mTrust: Discerning Multi-Faceted Trust in a



رضا ابراهیمی آتانی دانشیار گروه مهندسی رایانه دانشکده فنی و مهندسی دانشگاه گیلان است. نامبرده مدرک دکترای خود را در سال ۱۳۸۹ در رشته مهندسی الکترونیک از دانشگاه علم و

صنعت ایران دریافت کرد. ایشان عضو پیوسته انجمن رمز ایران و انجمن های بین المللی IACR و IEEE هستند. از ایشان تاکنون چهار عنوان کتاب و بیش از یکصد مقاله در مجلات و کنفرانس های داخلی و بین المللی به چاپ رسیده است. زمینه پژوهشی مورد علاقه وی، طراحی و پیاده سازی الگوریتم های رمزنگاری، امنیت شبکه و امنیت نرم افزار است.

مهدی صادق پور متولد ۱۳۶۹، مدرک کارشناسی و کارشناسی ارشد خود را به ترتیب در شهریور ماه ۱۳۹۱ و آذر ماه ۱۳۹۳ در رشته مهندسی کامپیوتر نرم افزار از دانشگاه گیلان دریافت کرده است. زمینه پژوهشی مورد علاقه ایشان مهندسی نرم افزار، داده کاوی امنیت شبکه است.

ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 414-423, August 2006.

- [27] A. Machanavajjhala, D. Kifer, J. Gehrke and M. Venkatasubramanian, "L-diversity: Privacy beyond k-anonymity," ACM Transactions on Knowledge Discovery from Data (TKDD), vol. 1, no. 1, 2007 .
- [28] N. Li, T. Li and S. Venkatasubramanian, "t-Closeness: Privacy Beyond k-Anonymity and l-Diversity," IEEE 23rd International Conference on Data Engineering, pp. 106 - 115, April 2007.
- [29] Y. Rubner, C. Tomasi and L. J. Guibas, "The Earth Mover's Distance as a Metric for Image Retrieval," International Journal of Computer Vision, vol. 40, no. 2, pp. 99 - 121, 2000.
- [30] N. Li, T. 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, 2010.
- [31] M. E. Nergiz, M. Atzori and C. Clifton, "Hiding the Presence of Individuals from Shared Databases," InProc. of ACM International Conference on Management of Data, pp. 665-676, June 2007 .
- [32] A. S. Sattar, J. Li, X. Ding, J. Liu and M. Vincent, "A general framework for privacy preserving data publishing," Knowledge-Based Systems, vol. 54, pp. 276-287, 2013. Statistical Databases, From Theory to Practice, pp. 97-116, 2002.
- [23] D. B. Rubin, "Discussion Statistical Disclosure Limitation," Journal of Official Statistics, vol. 9, no. 2, p. 461-468, 1993.
- [24] T. Dalenius, "Towards a Methodology for Statistical Disclosure Control," Statistik Tidskrift, vol. 15, p. 429-222, 1977.
- [25] C. Dwork, "Differential Privacy," in Automata, Languages and Programming, Springer Berlin Heidelberg, 2006, pp. 1-12.
- [26] K. Wang and B. C. M. Fung, "Anonymizing sequential releases," Proceedings of the 12th