



ارایهٔ یک روش جدید انتشار داده‌ها با حفظ محرمانگی با هدف بهبود دقت طبقه‌بندی روی داده‌های گمنام

رضا ابراهیمی‌آتانی* و مهدی صادق‌پور

گروه مهندسی کامپیوتر، دانشکدهٔ فنی، دانشگاه گیلان، رشت، ایران

چکیده

با توسعهٔ روزافزون خدمات دولت الکترونیکی، اطلاعات شخصی افراد در قالب پایگاه‌های داده در دستگاه‌ها و ارگان‌های دولتی و خصوصی ذخیره شده است. در بسیاری از موارد برای پردازش و استخراج دانش از این منابع دادهٔ بزرگ و بارز، نیاز به انتشار منابع داده و در اختیار گذاشتن اطلاعات به سایر نهادها و شرکت‌ها پدید می‌آید که این امر موجب ایجاد چالش‌های امنیتی در نقض حریم خصوصی افراد می‌شود. در این مقاله ضمن بررسی کامل پیشینهٔ پژوهش، حفظ محرمانگی در انتشار داده‌ها، یک روش کارآمد برای گمنام‌سازی ارائه می‌شود که هدف آن حفظ دقت طبقه‌بندی روی داده‌های گمنام است. این روش با بهره‌گیری از درخت تصمیم از انتشار اطلاعاتی که تأثیر کمی بر سودمندی داده‌های خروجی دارد و حذف آن‌ها موجب تأمین محرمانگی می‌شود، جلوگیری می‌کند. یکی از چالش‌های طرح‌هایی که از عمل‌گر گمنام‌سازی عمومی‌سازی استفاده می‌کنند، نیازمندی به ساخت درخت طبقه‌بندی برای هر شبه‌شناسه است که بیش‌تر به صورت خودکار صورت می‌گرفت. در طرح پیشنهادی نیازی به ساخت درخت طبقه‌بندی نیست. نتایج شبیه‌سازی و ارزیابی‌های انجام‌شده نشان می‌دهد، میان دقت الگوریتم‌های طبقه‌بندی که روی مجموعه داده استاندارد گمنام‌شده توسط این روش و مجموعه دادهٔ اولیه آموزش دیده‌اند، تفاوت اندکی وجود دارد.

واژگان کلیدی: حفظ محرمانگی، طبقه‌بندی، گمنام‌سازی، درخت تصمیم، عمل‌گر فرونشانی.

A New Privacy Preserving Data Publishing Technique Conserving Accuracy of Classification on Anonymized Data

Reza Ebrahimi Atani* & Mehdi Sadeghpour

Department of Computer Engineering, University of Guilan, Rasht, Iran

Abstract

Data collection and storage has been facilitated by the growth in electronic services, and has led to recording vast amounts of personal information in public and private organizations databases. These records often include sensitive personal information (such as income and diseases) and must be covered from others access. But in some cases, mining the data and extraction of knowledge from these valuable sources, creates the need for sharing them with other organizations. This would bring security challenges in user's privacy. The concept of privacy is described as sharing of information in a controlled way. In other words, it decides what type of personal information should be shared and which group or person can access and use it. "Privacy preserving data publishing" is a solution to ensure secrecy of sensitive information in a data set, after publishing it in a hostile environment. This process aimed to hide sensitive information and keep published data suitable for knowledge discovery techniques. Grouping data set records is a broad approach to data anonymization. This technique prevents access to sensitive attributes of a specific record by eliminating the distinction between a

* Corresponding author

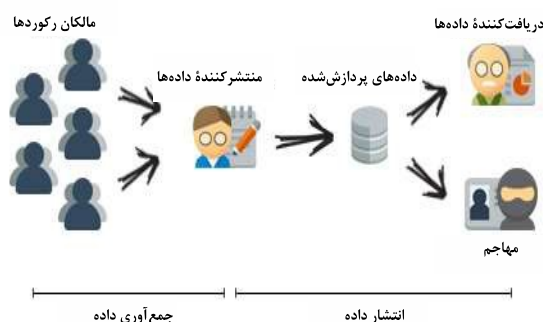
* نویسندهٔ عهده‌دار مکاتبات

number of data set records. So far a large number of data publishing models and techniques have been proposed but their utility is of concern when a high privacy requirement is needed. The main goal of this paper to present a technique to improve the privacy and performance data publishing techniques. In this work first we review previous techniques of privacy preserving data publishing and then we present an efficient anonymization method which its goal is to conserve accuracy of classification on anonymized data. The attack model of this work is based on an adversary inferring a sensitive value in a published data set to as high as that of an inference based on public knowledge. Our privacy model and technique uses a decision tree to prevent publishing of information that removing them provides privacy and has little effect on utility of output data. The presented idea of this paper is an extension of the work presented in [20]. Experimental results show that classifiers trained on the transformed data set achieving similar accuracy as the ones trained on the original data set.

Keywords: Privacy preservation, Data sharing, Anonymization, Classification, Decision tree, Suppression.

• **دریافت کننده داده‌ها:** پژوهش‌گران علمی یا سازمان‌های تجاری که به دنبال کشف دانش از مجموعه داده منتشر شده هستند.

• **گمنام‌سازی داده‌ها:** فرآیندی است که داده‌های خام را به داده‌های گمنام‌شده تبدیل می‌کند. داده‌های گمنام‌شده در مقابل نوع خاصی از حملات که در یک مدل محرمانگی تعریف شده است، از افشای اطلاعات مصون می‌مانند. در PPDP سعی بر این است که با پردازش مجموعه داده اولیه هویت مالکان رکوردها، پنهان نگه‌داشته شوند تا به حریم شخصی آنها لطمه‌ای وارد نشود. همان‌طور که در شکل (۱) نمایش داده شده است، دو مرحله اصلی فرآیند PPDP جمع‌آوری و انتشار داده‌ها هستند [4]. در مرحله نخست منتشرکننده داده‌ها اطلاعات 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 محرمانگی روی داده‌ها اعمال می‌شود و خروجی در اختیار دریافت‌کنندگان داده قرار می‌گیرد.



(شکل-۱): جمع‌آوری و انتشار داده‌ها
(Figure-1): Data collection and data publishing

۱- مقدمه

با توسعه و نفوذ هر چه بیشتر خدمات الکترونیکی در زندگی روزمره افراد و سازمان‌ها، سامانه‌های رایانه‌ای و پایگاه‌های داده، محل پردازش و جمع‌آوری اطلاعات بسیار مهمی از زندگی افراد شده است. نظر به حجم بسیار زیاد اطلاعات در این منابع داده، کشف الگوها و قوانین سودمند، روش‌های داده‌کاوی را به ابزاری کلیدی برای استفاده از این منابع باارزش داده تبدیل کرده است. مالکان این مجموعه داده‌ها ممکن است، به دلایل مختلفی تصمیم به انتشار آنها بگیرند؛ اما به دلیل وجود اطلاعات حساس و شخصی مربوط به افراد، نمی‌توانند این داده‌ها را به صورت خام انتشار دهند. فرآیند انتشار با حفظ محرمانگی داده^۱ (PPDP) راه‌کاری برای حل این مشکل در محیط‌های غیرقابل اعتماد است [1,4].

در این فرآیند سعی بر آن است که ضمن مخفی نگه‌داشتن اطلاعات حساس، داده‌های منتشرشده، همچنان برای عملیات کشف دانش مفید باقی بماند [3]. در ادامه، مفاهیم پایه‌ای PPDP که در مقاله مورد استفاده قرار می‌گیرند، تعریف شده‌اند:

- **مالکان رکورد:** افرادی هستند که اطلاعات مربوط به آنها در مجموعه داده ثبت شده است.
- **منتشرکننده داده‌ها:** 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 نامیده می‌شوند) را از مجموعه داده استخراج کند. اطلاعات در دسترس مهاجم فقط به داده‌های منتشرشده محدود نیست و امکان دارد از دانش پیش‌زمینه خود نیز برای انجام حملات استفاده کند.

^۱ Privacy-Preserving Data Publishing (PPDP)

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.

صفات رکوردها در یک مجموعه‌داده از لحاظ کاربرد در PPDP به چند دسته تقسیم می‌شوند. به صفاتی که مالک رکورد مایل نیست عموم از آنها اطلاع داشته باشند، صفت حساس گفته می‌شود (مانند نوع بیماری و یا میزان حقوق یک فرد). از سوی دیگر صفاتی که به‌صراحت مالک رکورد را مشخص می‌کنند، شناسه صریح نامیده می‌شوند. هدف PPDP پنهان‌ساختن تناظر بین شناسه صریح و صفت حساس رکوردها است. نخستین قدم در گمنام‌سازی مجموعه داده، حذف شناسه‌های صریح از رکوردها است؛ اما همان‌طور که در [5] نشان داده شده است، حتی در صورت حذف شناسه صریح به‌کمک دسته‌ای دیگر از صفات که شبه‌شناسه^۱ نامیده می‌شوند، می‌توان صفت حساس مالکان رکوردها را به‌دست آورد. کدپستی، تاریخ تولد و جنسیت را به‌عنوان مثالی از شبه‌شناسه‌ها می‌توان نام برد که هر یک به‌تنهایی برای تعیین مالک رکورد کافی نیستند؛ اما با ترکیب آنها می‌توان به هویت یک فرد خاص رسید.

برای جلوگیری از این مشکل می‌توان با استفاده از عمل‌گرهای گمنام‌سازی، موجب شد تا شبه‌شناسه‌های هر رکورد در چندین رکورد دیگر نیز تکرار و مهاجم در شناسایی مالک رکورد دچار ابهام شود. از میان عمل‌گرهای گمنام‌سازی پرکاربرد، می‌توان به عمل‌گرهای عمومی‌سازی^۲ و فرونشانی^۳ اشاره کرد [2]. در عمل‌گر عمومی‌سازی، مقدار یک صفت به‌کمک درخت طبقه‌بندی^۴ به بازه‌ای بزرگ‌تر تعمیم داده می‌شود؛ برای مثال در صفت نشانی، مقدار ایران با خاورمیانه جایگزین می‌شود و در مرحله بعد می‌توان خاورمیانه را با آسیا جایگزین کرد.

مشکل این روش، نیاز به تشکیل درخت طبقه‌بندی برای هر شبه‌شناسه است که در بیش‌تر موارد این کار به‌صورت دستی و توسط یک فرد خبره انجام می‌گیرد [3]. علاوه‌براین ممکن است، بین درخت‌های طبقه‌بندی تشکیل‌شده توسط

افراد مختلف، تفاوت وجود داشته باشد و موجب به‌دست‌آمدن نتایج گوناگون شود. همچنین زمانی که از عمل‌گر عمومی‌سازی سلولی^۵ در یک مجموعه‌داده استفاده شود، مجموعه یادشده دیگر برای عمل طبقه‌بندی^۶ مناسب نیست و منجر به مشکل پدیده اکتشاف داده^۷ می‌شود [4]. در مقابل، عمل‌گر فرونشانی، مقدار شبه‌شناسه موردنظر را با مقدار گمشده^۸ جایگزین می‌کند. مقدار گمشده که به‌طورمعمول با علامت “?” نمایش داده می‌شود، بیان‌گر ناشناخته‌بودن مقدار صفت موردنظر است. استفاده از عمل‌گر فرونشانی باید با دقت انجام گیرد؛ زیرا تعداد زیاد شبه‌شناسه‌هایی که مقدار گمشده دارند، موجب کاهش کیفیت مجموعه داده می‌شود [2].

در این مقاله الگوریتم جدیدی برای حفظ محرمانگی طبق چهارچوب معرفی‌شده در [20] معرفی می‌شود. در روش پیشنهادی از عمل‌گر فرونشانی استفاده می‌شود، در نتیجه مشکلات یادشده را برای عمومی‌سازی به همراه نخواهد داشت. این روش از درخت تصمیمی که از روی داده‌های خام ایجاد می‌شود، به‌منظور انتخاب شبه‌شناسه‌ها برای فرونشانی استفاده می‌کند و مجموعه‌داده گمنام‌شده توسط آن برای عمل طبقه‌بندی مناسب است. نتایج به‌دست‌آمده از آزمایش‌ها نشان می‌دهد که دقت طبقه‌بندی در این حالت، تفاوت اندکی با طبقه‌بندی روی مجموعه داده اولیه دارد.

در ادامه ساختار مقاله بر این اساس ارائه می‌شود: در بخش دوم به‌اختصار پیشینه پژوهش درخصوص انتشار داده با حفظ محرمانگی ارائه می‌شود. در بخش سوم مدل‌های حفظ محرمانگی و روابط ریاضی مورد نیاز برای معرفی الگوریتم آورده شده است. در بخش چهارم مقاله، الگوریتم پیشنهادی تشریح می‌شود. بخش پنجم نتایج پیاده‌سازی و ارزیابی طرح گمنام‌سازی پیشنهادی ارائه و مورد نقد و بررسی واقع و سرانجام در بخش ششم، مقاله جمع‌بندی و نتیجه‌گیری می‌شود.

۲- پیشینه پژوهش

دسته‌ای از پژوهش‌های انجام‌شده در زمینه PPDP به‌دنبال حفظ کیفیت داده‌های گمنام‌شده برای انجام یک عمل داده‌کاوی معین هستند. به بیان دیگر، تلاش می‌شود پردازش رکوردها به‌گونه‌ای انجام گیرد که مجموعه‌داده گمنام‌شده

^۵ Cell Generalization

^۶ Classification

^۷ Data Exploration

^۸ Missing Value

^۱ Quasi Identifier

^۲ Generalization

^۳ Suppression

^۴ Taxonomy Tree

علاوه بر این که یک مدل حفظ محرمانگی را ارضا می کند، برای استفاده در عمل داده کاوی خاصی بهینه شده باشد.

طبقه بندی یکی از عملیات اصلی داده کاوی محسوب می شود. به همین دلیل پژوهش های بسیاری برای کاهش اختلاف نتیجه طبقه بندی کردن روی داده های خام و گمنام، انجام شده است. پژوهش ارایه شده در [21] را می توان از نخستین طرح های کاربردی در این زمینه به شمار آورد. این روش برای هر عمل عمومی سازی G ، میزان محرمانگی که بر اثر G به دست می آید (P) و مقدار اطلاعاتی را که از دست می رود (I) محاسبه می کند. تا زمانی که مدل k -anonymity در کل مجموعه داده برقرار نشده است، عملی که معیار $\frac{I(G)}{P(G)}$ را کمینه کند، روی داده ها اعمال می شود.

روش معرفی شده در [22] روند گمنام سازی را به صورت بالابه پایین انجام می دهد؛ به عبارت دیگر کار را از حالتی که تمام شبه شناسه ها در عمومی ترین حالت خود قرار دارند، آغاز می کند و با استفاده از عمل گر خاص سازی در هر مرحله، مقدار یک شبه شناسه را یک گره در درخت طبقه بندی پایین می آورد (عکس عمل عمومی سازی). در این روش نیز انتخاب شبه شناسه ها برای خاص سازی، با توجه به میزان گمنامی حاصل شده و اطلاعات از دست رفته انجام می شود. اجرای الگوریتم تا زمانی که خاص سازی بیشتر موجب نقض مدل محرمانگی نشود، ادامه می یابد. این روش نیازی به درخت طبقه بندی برای شبه شناسه های عددی ندارد و درخت مورد نیاز این شبه شناسه ها طی اجرای الگوریتم ساخته می شود. این ایده در [10] بهبود بیشتری پیدا کرده است و مطابق آن اگر درخت طبقه بندی برای شبه شناسه ای موجود نباشد از عمل گر فرونشانی استفاده می شود. این روش از فرونشانی سراسری استفاده می کند؛ یعنی وقتی که فرونشانی بر شبه شناسه q با مقدار x اعمال می شود، q در تمام رکوردهایی که در آن ها مقدار x دارد، فرونشاندن می شود.

در روش KACTUS [23] درخت طبقه بندی از فرآیند گمنام سازی حذف شده است و فقط از عمل گر فرونشانی استفاده می شود. این روش از فرونشانی محلی استفاده می کند که باعث می شود، فرونشاندن شبه شناسه ای از یک رکورد، تأثیری بر سایر رکوردها نداشته باشد و در نتیجه کیفیت مجموعه داده کمتر کاهش می یابد. شکل (۲) نحوه کارکرد این روش را برای برقراری k -anonymity با پارامتر $k=100$ نشان می دهد. در قسمت الف، صفت D مورد پردازش قرار می گیرد. از میان فرزندان این گره، فقط یکی بیش از صد رکورد دارد و دو گره دیگر در مجموع ۶۵ رکورد دارند؛ بنابراین ۳۵ رکورد از

طبقه A برای ایجاد دسته صدتایی نگه داشته می شود و ۱۴۵ رکورد دیگر با شبه شناسه های $\{A='a3', B='b1', C='c1', D='d2'\}$ به مجموعه رکوردهای گمنام شده اضافه می شوند. در قسمت ب، صفت E مورد بررسی قرار می گیرد. به دلیل این که هیچ کدام از فرزندان آن بیش از صد رکورد ندارند، پس از فرونشانی E با هم ادغام می شوند (قسمت ج) و هیچ رکوردی به مجموعه رکوردهای گمنام شده اضافه نمی شود. این روند تا پیمایش کل درخت ادامه می یابد.

هرچند دقت طبقه بندی روی داده های گمنام شده توسط این روش در حدی مطلوب است، اما به دلیل این که مجموعه داده را طبق مدل محرمانگی k -anonymity گمنام می کند، رکوردهای منتشر شده در مقابل حملاتی از قبیل حمله هم نوعی و حمله دانش پیش زمینه آسیب پذیر هستند [23] و [24].

به طور اصولی دست یابی به محرمانگی مطلق غیر ممکن است؛ یعنی نمی توان کاری کرد که پس از انتشار مجموعه داده، مهاجم هیچ گونه اطلاعاتی از مالکان رکورد به دست نیاورد. به همین دلیل با توجه به کاربرد، مدل های مختلفی برای سطوح متفاوت از حفظ محرمانگی پیشنهاد شده است که هر کدام در مقابله با حملات معینی کارآمد هستند. در جدول (۱) مقایسه ای بین عملگرهای گمنامی انجام گرفته است.

اگر در زمان گمنام سازی از کاربردی که داده ها در آن استفاده خواهند شد، آگاه باشیم، می توانیم داده هایی متناسب با آن کاربرد ایجاد کنیم. طبقه بندی یکی از عملیات مهم داده کاوی است. به همین دلیل انتشار داده هایی که کیفیت مناسب برای طبقه بندی داشته باشند، بسیار با اهمیت است. در این شرایط ملاک باکیفیت بودن داده ها، تشابه نتایج طبقه بندی روی داده های خام و گمنام است.

در این مقاله ضمن استفاده از ایده فرونشانی شبه شناسه های با بهره اطلاعاتی کمتر، به جای گمنام سازی طبق k -anonymity از مدل احتمالاتی ارائه شده در [4] استفاده می شود که در مقایسه با مدل شناخته شده Differential privacy [7] به عملکرد خوبی در حفظ سودمندی و محرمانگی داده ها دست یافته است.

۳- مدل های حفظ محرمانگی در انتشار داده

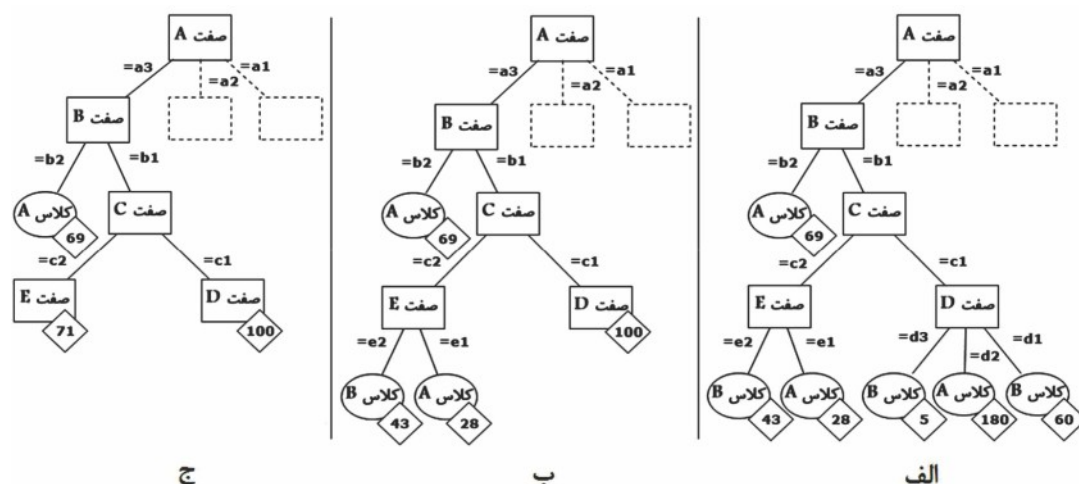
پیش از تلاش برای حفظ محرمانگی باید مشخص کنیم چه زمانی این هدف برآورده می شود. تعریف ایده آل از حفظ

در ادامه به شرح آن‌ها پرداخته می‌شود.

۱-۳- مدل پیونددهی رکورد

در حملات پیونددهی رکورد، مهاجم شبه‌شناسه‌های قربانی را با رکوردهای انتشار یافته تطبیق می‌دهد؛ بعد از انجام این کار دسته‌ای کوچک از رکوردها باقی می‌ماند، که ممکن است، متعلق به قربانی باشد.

محرمانگی هنگام انتشار داده‌ها، در [12] ارائه شده است. بر طبق این تعریف، بعد از آن که داده‌ها در اختیار مهاجم قرار می‌گیرد، نسبت به زمانی که داده‌ها را در اختیار نداشته است، نباید هیچ دانشی در مورد مالکان رکورد به او اضافه شود. با در نظر گرفتن اطلاعاتی که ممکن است از قبل در اختیار مهاجم باشد، این سطح از محرمانگی قابل دستیابی نیست [13]. به همین دلیل در بیش‌تر کارهایی که در حوزه محرمانگی داده انجام می‌گیرد از مدل‌های تعدیل یافته‌تری استفاده می‌شود که



(شکل-۲): روند اجرای الگوریتم kACTUS
(Figure-2): kACTUS Algorithm operation flow

(جدول-۱): مقایسه عملگرهای گمنام‌سازی

(Table-1): A comparison on anonymity operators

نام روش	نحوه عملکرد	رکوردها معرف یک نمونه واقعی هستند؟	قابلیت استفاده از روش‌های استاندارد داده‌کاوی	میزان تغییر مقدار صفات
عمومی سازی	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
برش دهی	جدا کردن صفات مرتبط به هم و 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	به‌طور کامل عوض می‌شود

به راحتی می‌تواند آن را پیدا کند؛ زیرا تنها یک رکورد متناظر وجود دارد و در نتیجه حمله پیونددهی رکورد با موفقیت انجام می‌شود.

روش k -anonymity [5] برای جلوگیری از چنین حملاتی، مجموعه داده را به دسته‌های هم‌ارزی با اندازه k ، تقسیم می‌کند. در این حالت به ازای هر مقدار ممکن از

پیونددهی رکورد، زمانی رخ می‌دهد که مهاجم موفق شود رکورد مربوط به قربانی را به‌طور یکتا شناسایی کند [1]. به عنوان مثال در جدول (۳) صفات تحصیلات، جنسیت و سن به عنوان شبه‌شناسه‌ها و حقوق به عنوان صفت حساس در نظر گرفته شده است. پس از انتشار این جدول اگر مهاجم به دنبال رکوردی با شبه‌شناسه‌های <کارشناسی، مرد، ۲۹> باشد،

۲-۳- مدل پیونددهی صفت

همان‌طور که گفته شد، در حمله پیونددهی رکورد مهاجم در صدد تشخیص رکورد مربوط به قربانی است؛ اما در بیش‌تر موارد به‌دست‌آوردن صفت حساس قربانی، مهاجم را به هدف خود می‌رساند و نیازی به شناسایی رکورد نیست [1]. برای مثال در جدول (۲) مقدار صفت حساس در سه رکورد آخر برابر است و مهاجم با دراختیارداشتن 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. l -diversity [15] یکی از روش‌های مورد استفاده برای مقابله با حمله پیونددهی صفت است. برای تأمین l -diversity در یک جدول، در هر دسته هم‌ارزی باید دست‌کم l 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 l مقدار متفاوت ولی نزدیک به هم وجود داشته باشد، باز هم مشکل عدم تنوع بروز پیدا می‌کند.

برای مثال قسمت «ب» از جدول (۳) مدل 3 -diversity را ارضا می‌کند؛ اما مقادیر صفت حساس در نخستین گروه هم‌ارزی بسیار نزدیک به هم هستند. برای رفع این نقص، در [16] 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 [17] برای سنجش تفاوت بین توزیع‌ها استفاده می‌کند. در جدول (۴) نتیجه گمنام‌سازی یک مجموعه داده به روش‌های l -diversity و t -Closeness مقایسه شده است.

برقرارکردن 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 تکنیکی انعطاف‌پذیرتر به نام (n,t) -Closeness ارائه شده است [18]. 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 هرچه 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.

شبه‌شناسه‌ها، یا هیچ رکوردی در مجموعه داده وجود ندارد و یا k رکورد (و بیشتر) موجود است، در نتیجه مهاجم بیش‌تر با درجه اطمینان $1/k$ می‌تواند رکورد قربانی را شناسایی کند. این روش در قسمت «ب» از جدول (۳) 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.

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 یک دسته هم‌ارزی قرار می‌گیرند. در این حالت مهاجم با درجه اطمینان n/k رکورد قربانی را به‌دست می‌آورد که با نیاز تعریف‌شده در k -anonymity مطابقت ندارد. برای حل این مشکل می‌توان از مفهومی عمومی‌تر به نام (X,Y) -anonymity [14] استفاده کرد. در این روش X و Y دو زیرمجموعه از صفات مجموعه داده هستند که اشتراکی با هم ندارند. برای برآورده شدن محرمانگی در (X,Y) -anonymity، هر مقدار در X باید دست‌کم 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 بیشینه تعداد رکوردهای متعلق به یک فرد در یک دسته هم‌ارزی را کنترل کرد.

(جدول ۲-): مثالی از گمنام‌سازی k -anonymity با $k=2$

(Table-2): A example of k -anonymity with $k=2$

الف) مجموعه داده خام

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

ب) جدول 2 -anonymous

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

(جدول-۳): مقایسه روش‌های t-Closeness و l-diversity

(Table-3): A comparison between l-diversity and t-closeness techniques

الف) مجموعه داده خام				ب) 3-diverse				پ) 0.167-closeness			
کدپستی	سن	حقوق		کدپستی	سن	حقوق		کدپستی	سن	حقوق	
۴۷۶۷۷	۲۹	۳۰۰۰	۱	۴۷۶**	۲*	۳۰۰۰	۱	۴۷۶۷*	≤ ۴۰	۳۰۰۰	۱
۴۷۶۰۲	۲۲	۴۰۰۰	۲	۴۷۶**	۲*	۴۰۰۰	۲	۴۷۶۷*	≤ ۴۰	۵۰۰۰	۳
۴۷۶۷۸	۲۷	۵۰۰۰	۳	۴۷۶**	۲*	۵۰۰۰	۳	۴۷۶۷*	≤ ۴۰	۹۰۰۰	۸
۴۷۹۰۵	۴۳	۶۰۰۰	۴	۴۷۹**	≥ ۴۰	۶۰۰۰	۴	۴۷۹۰*	≥ ۴۰	۶۰۰۰	۴
۴۷۹۰۹	۵۲	۱۱۰۰۰	۵	۴۷۹**	≥ ۴۰	۱۱۰۰۰	۵	۴۷۹۰*	≥ ۴۰	۱۱۰۰۰	۵
۴۷۹۰۶	۴۷	۸۰۰۰	۶	۴۷۹**	≥ ۴۰	۸۰۰۰	۶	۴۷۹۰*	≥ ۴۰	۸۰۰۰	۶
۴۷۶۰۵	۳۰	۷۰۰۰	۷	۴۷۶**	۳*	۷۰۰۰	۷	۴۷۶۰*	≤ ۴۰	۴۰۰۰	۲
۴۷۶۷۳	۳۶	۹۰۰۰	۸	۴۷۶**	۳*	۹۰۰۰	۸	۴۷۶۰*	≤ ۴۰	۷۰۰۰	۷
۴۷۶۰۷	۳۲	۱۰۰۰۰	۹	۴۷۶**	۳*	۱۰۰۰۰	۹	۴۷۶۰*	≤ ۴۰	۱۰۰۰۰	۹

(جدول-۴): روش δ-Presence

(Table-4): δ-Presence technique
P

نام	کدپستی	سن	ملیت	
Alice	۴۷۹۰۶	۳۵	ایالات متحده	a
Bob	۴۷۹۰۳	۵۹	کانادا	b
Chris	۴۷۹۰۶	۴۲	ایالات متحده	c
Dirk	۴۷۶۳۰	۱۸	برزیل	d
Eunice	۴۷۶۳۰	۲۲	برزیل	e
Frank	۴۷۶۳۳	۶۳	پرو	f
Gail	۴۸۹۷۳	۳۳	اسپانیا	g
Harry	۴۸۹۷۲	۴۷	بلغارستان	h
Iris	۴۸۹۷۰	۵۲	فرانسه	i

T*

کدپستی	سن	ملیت	
۴۷*	*	امریکا	b
۴۷*	*	امریکا	c
۴۷*	*	امریکا	f
۴۸*	*	اروپا	h
۴۸*	*	اروپا	i

جداول نسبت به مدل پیونددهی رکورد و صفت، کارهای کمتری انجام شده که معروف‌ترین آن‌ها δ-Presence است [19].

طبق این مدل، احتمال شناسایی رکورد قربانی در مجموعه داده انتشار یافته باید به سقف $\delta\%$ محدود شود. این مفهوم با مثالی که در جدول (۴) 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. T^* نسخه گمنام‌سازی شده T است که پس از انتشار آن مهاجم نباید با احتمال بیشتر از δ 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 $\frac{3}{6}$ است. این احتمال برای رکوردهای g-i برابر با $\frac{2}{3}$ است. بنابراین مهاجم حداکثر با احتمال $\frac{2}{3}$ می‌تواند از وجود یک رکورد دلخواه از P در جدول T^* اطمینان داشته باشد.

همان‌طور که گفته شد، در روش δ-Presence مهاجم حداکثر $\delta\%$ از وجود رکورد قربانی در جدول منتشرشده اطمینان دارد؛ بنابراین احتمال این که بتواند با موفقیت حمله پیونددهی رکورد یا صفت را انجام دهد، نیز حداکثر $\delta\%$ 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. درحالی که چنین فرضی در دنیای واقعی زیاد محتمل نیست.

۳-۲-۱- مدل پیونددهی جداول

در برخی شرایط، صرف تشخیص این که رکورد قربانی در مجموعه داده انتشار یافته وجود دارد یا خیر، می‌تواند نقض محرمانگی تلقی شود. به عنوان نمونه زمانی که یک بیمارستان داده‌های مربوط به اشخاصی را که دچار بیماری خاصی هستند، منتشر می‌کند، پی‌بردن به هویت هر یک از افراد، نوع بیماری وی را آشکار می‌سازد. در زمینه مدل پیونددهی

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

طرح ارائه‌شده در [13] یکی از سخت‌گیرانه‌ترین مدل‌های حفظ محرمانگی 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 F مدل شود، به‌ازای هر مجموعه داده D و D' که در وجود یک رکورد اختلاف دارند، باید رابطه زیر برقرار باشد:

$$\forall S \in \text{Range}(F) (\Pr[F(D) = S] \leq e^\epsilon \times \Pr[F(D') = S])$$

که S صفت حساس و ϵ پارامتری است که هرچه کوچک‌تر انتخاب شود، میزان تأثیر یک رکورد بر نتیجه بررسی‌ها کمتر می‌شود. Differential Privacy تعریف محکمی از محرمانگی 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. طرح ارائه‌شده در [20] با در نظر داشتن این مشکل مدلی تعدیل‌یافته‌تر معرفی کرده است: اطمینانی که مهاجم پس از مشاهده مجموعه داده گمنام‌شده نسبت به صفت حساس قربانی به‌دست می‌آورد، نباید بیش‌تر از یک حد تصادفی باشد. از آنجاکه طرح ارائه‌شده در این مقاله بر پایه همین مدل بنا شده است، در ادامه با جزئیات کامل‌تری تشریح می‌شود.

مدل حفظ محرمانگی معرفی‌شده در [20] (AGF) بر اساس دو مفهوم اطمینان مورد انتظار (EC) و اطمینان مشاهده‌شده (OC) شکل گرفته است. اطمینان مشاهده‌شده یک رکورد، بیان‌گر این است که مهاجم با چه احتمالی می‌تواند صفت حساس آن را به قربانی نسبت دهد. اگر رکورد قربانی (v) عضو دسته هم‌ارزی به‌اندازه $|E(v)|$ باشد و تعداد تکرار صفت حساس v در این دسته هم‌ارزی را $\#S(v)$ بنامیم، از دید مهاجم احتمال آن که صفت حساس v متعلق به قربانی باشد، برابر $\frac{\#S(v)}{|E(v)|}$ است. به‌دلیل این که در این روش در ابتدا از رکوردها با نرخ β نمونه‌برداری می‌شود، باید احتمال گزینش رکورد قربانی در مرحله نمونه‌برداری نیز مدنظر قرار بگیرد. همان‌طور که رابطه (۱) نشان می‌دهد از ضرب این دو احتمال،

اطمینان مشاهده‌شده رکورد r به‌دست می‌آید.

$$\text{ObservedConfidence}(r) = \beta \times \frac{\#S(r)}{|E(r)|} \quad (1)$$

هرچه اطمینان مشاهده‌شده 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 به‌صورت تصادفی از فضای نمونه رکوردها ایجاد شده است. برطبق این مدل، اطمینان مشاهده‌شده یک رکورد در مجموعه داده گمنام‌شده، نباید از احتمال دیده‌شدن آن رکورد در D_0 بیش‌تر باشد. این احتمال اطمینان مورد انتظار نامیده می‌شود. با فرض این که مجموعه داده شامل n رکورد باشد، احتمال وجود رکورد r در D_0 برابر با احتمال مکمل "۰" پیروزی در n آزمون با احتمال پیروزی $\Pr(r)$ است و با استفاده از توزیع دوجمله‌ای طبق رابطه (۲) محاسبه می‌شود:

$$\begin{aligned} \text{ExpectedConfidence}(r) &= 1 - f(0; n, \Pr(r)) \\ &= 1 - (1 - \Pr(r))^n \end{aligned} \quad (2)$$

که $\Pr(r)$ با فرض مستقل بودن صفات از یکدیگر به‌وسیله ضرب احتمال شبه‌شناسه‌ها در احتمال صفت حساس به‌دست می‌آید:

$$\Pr(r) = (\prod_i \Pr(q_i)) \times \Pr(s) \quad (3)$$

با توجه به رابطه (۱) و (۲)، برای تأمین محرمانگی طبق قاعده ذکرشده، به‌ازای هر رکورد موجود در مجموعه داده، باید رابطه (۴) برقرار باشد:

$$\text{ObservedConfidence}(r) \leq \text{ExpectedConfidence}(r) \quad (4)$$

وقتی رابطه (۴) برای رکوردی برقرار نبود، یکی از شبه‌شناسه‌ها باید برای عمومی‌سازی انتخاب شود. این کار بر اساس معیاری به نام Left Domain Size (LDS) انجام می‌گیرد. LDS شبه‌شناسه‌ای با مقدار v ، تعداد گره‌هایی است که در درخت طبقه‌بندی در سطح برابر یا بالاتر از v قرار دارد. شبه‌شناسه‌ای که مقدارش کمترین LDS را دارد برای عمومی‌سازی برگزیده خواهد شد.

در حین عملیات گمنام‌سازی، ممکن است، عمل‌گر عمومی‌سازی روی تعداد زیادی از شبه‌شناسه‌های یک رکورد اعمال شود. چنین رکوردی کیفیت داده‌های خروجی را کاهش می‌دهد. به همین منظور از معیار تخریب برای سنجش میزان کیفیت رکوردها استفاده می‌شود و رکوردهایی که میزان تخریب آن‌ها از حدّ از پیش تعریف‌شده δ' بیش‌تر باشد از مجموعه خروجی حذف می‌شوند. میزان تخریب رکورد r با

شبه‌شناسه‌های q_0 تا q_d از رابطه (۵) به دست می‌آید.

$$\delta = \frac{1}{d} \sum_{i=1}^d \left(1 - \frac{\text{Current level of } q_i}{\text{Maximal level in } q_i}\right) \quad (5)$$

اگر تخریب "۰" باشد به معنی این است که تمام شبه‌شناسه‌های رکورد، مقدار اصلی خود را حفظ کرده‌اند و

تخریب "۱" به معنی عمومی‌سازی تمام شبه‌شناسه‌ها تا گروه ریشه است.

در جدول (۵) مدل‌های مطرح حرمانگی بر اساس ویژگی‌ها و نقاط ضعف هر کدام مقایسه شده‌اند.

(جدول-۵): مقایسه مدل‌های حفظ حرمانگی
(Table-5): A Comparison between Privacy models

نام مدل	رابطه‌ای که در مجموعه داده باید برقرار شود	ویژگی‌ها	نقطه ضعف
k-anonymity	گروه‌بندی رکوردها در دسته‌های هم ارزی 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 هم ارزی 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 می‌تواند 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 هم ارزی 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 حرمانگی	تخریب زیاد
AGF	اطمینان مشاهده شده رکوردها کوچکتر از اطمینان مورد انتظار باشد	رفع ضعف ذکر شده برای differential privacy	فرض مستقل بودن مقدار صفات

۴- طرح الگوریتم پیشنهادی

در این بخش از مقاله، الگوریتم جدیدی برای حفظ حرمانگی در انتشار داده‌ها ارائه می‌شود که طبق چهارچوب معرفی شده در [20] طراحی و خروجی آن برای استفاده در طبقه‌بندی، بهینه شده است. یکی از مشکلات موجود در روش [20]، نیاز به حضور فردی خبره به منظور تشکیل درخت طبقه‌بندی برای هر شبه‌شناسه است. ممکن است افراد مختلف، درخت‌های طبقه‌بندی مختلفی برای صفات در نظر بگیرند که این امر موجب به دست آمدن نتایج گوناگون می‌شود. علاوه بر این مجموعه داده حاصل از اجرای این روش عام‌منظوره است، یعنی در کاربردهای مختلف داده کاوی می‌توان از آن استفاده کرد؛ ولی برای استفاده در کاربرد خاصی بهینه نیست. هنگامی که داده‌ها با هدف استفاده در طبقه‌بندی منتشر می‌شوند، می‌توان با تغییر در نحوه گمنام‌سازی، به نتایج بهتری دست یافت.

به دلیل این که [20] از عمل‌گر عمومی‌سازی سلولی

برای گمنام‌سازی استفاده می‌کند، در نحوه طبقه‌بندی رکوردهای خروجی ابهام به وجود می‌آید [1]. به عنوان نمونه اگر مقدار صفت «ملیت» در برخی رکوردها از ایران به آسیا تبدیل شود، الگوریتم طبقه‌بندی که به وسیله داده‌های گمنام‌شده آموزش دیده و قرار است رکورد جدیدی را طبقه‌بندی کند، باید از رابطه آسیا و ایران اطلاع داشته باشد؛ زیرا در داده‌های واقعی صفت ملیت هیچ‌گاه مقدار آسیا ندارد. همچنین ممکن است، در بعضی موارد قوانین طبقه‌بندی متناقضی به وجود بیاید. به عنوان مثال اگر صفت «ملیت» مقدار ایران داشت، آن رکورد متعلق به طبقه یک، و اگر مقدار آسیا داشت، متعلق به طبقه دو است.

در روش پیشنهادی از عمل‌گر فرونشانی استفاده می‌شود؛ در نتیجه مشکلات ذکر شده برای عمومی‌سازی را به همراه نخواهد داشت. در این روش، انتخاب شبه‌شناسه‌ها برای فرونشانی به وسیله درخت تصمیمی که از روی داده‌های خام ایجاد می‌شود، صورت می‌پذیرد.

شبه‌شناسه A در تعیین طبقه این رکورد نقشی ندارد، بنابراین فرونشاندن و به‌دلیل مشابه در قسمت (۳) نیز مقدار شبه‌شناسه C حذف می‌شود.

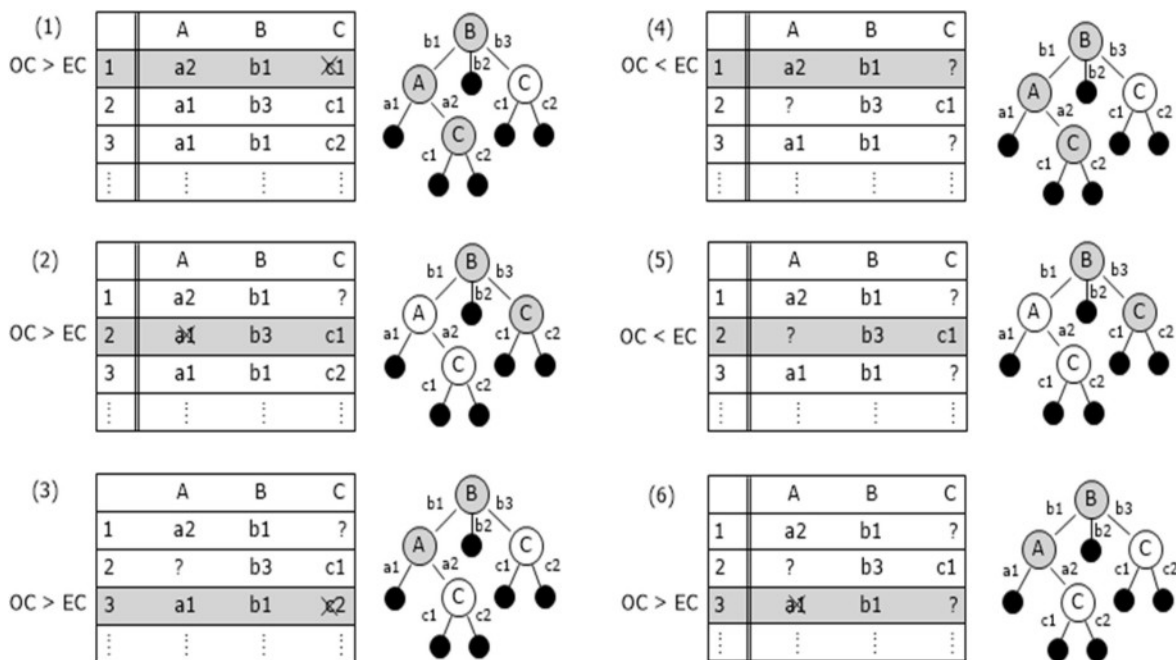
بعد از این‌که تمام رکوردها یک دور پردازش شدند، در قسمت (۴) دوباره نوبت رکورد نخست فرا می‌رسد. به‌دلیل فرونشانی‌های انجام‌شده در دور نخست، اطمینان مشاهده‌شده (OC) رکوردهای نخست و دوم کاهش یافته و قاعده محرمانگی برای این دو رکورد برقرار شده است؛ بنابراین در قسمت‌های (۴) و (۵) هیچ‌گونه تغییری در مقدار شبه‌شناسه‌ها صورت نمی‌پذیرد. در قسمت (۶)، رابطه $OC < EC$ برای رکورد سوم برقرار نیست. با توجه به این‌که مقدار C در دور قبل حذف شده بود، در این دور پایین‌ترین گره، یعنی A برای فرونشانی انتخاب می‌شود. این روال تا زمانی که رابطه (۴) برای تمام رکوردها برقرار شود، ادامه می‌یابد.

به‌دلیل این‌که روش پیشنهادی از فرونشانی استفاده می‌کند، تنها دو حالت برای شبه‌شناسه q_i از رکورد r متصور است: مقدار اصلی خود را حفظ کرده و یا مقدار گمشده گرفته است. میزان تخریب رکورد r با شبه‌شناسه‌های q_0 تا q_d از به‌دست می‌آید.

$$\delta = \frac{1}{d} \sum_{i=1}^d g(r, q_i), \quad g(r, q_i) = \begin{cases} 0: q_i \text{ isn't missed from } r \\ 1: q_i \text{ is missed from } r \end{cases} \quad (6)$$

در روش پیشنهادی، انتخاب شبه‌شناسه‌ها برای فرونشانی، به‌وسیله درخت تصمیم انجام می‌شود و در این کار اولویت با شبه‌شناسه‌هایی است که در پایین درخت قرار دارند. بدین منظور ابتدا درخت تصمیم از روی داده‌ها ساخته می‌شود؛ سپس اطمینان مشاهده‌شده و اطمینان مورد انتظار برای هر یک از رکوردها محاسبه شده و برقراربودن رابطه (۴) بررسی می‌شود. در صورتی‌که این رابطه برقرار نباشد، مسیر رکورد در درخت تصمیم مشخص شده و شبه‌شناسه‌ای که در انتهای مسیر قرار دارد، فرونشاندن می‌شود. با فرونشانی شبه‌شناسه‌ها، نقاط تمایز رکوردها کمتر می‌شود و اندازه دسته‌های هم‌ارزی افزایش می‌یابد. با توجه به رابطه (۱) این امر منجر به کاهش اطمینان مشاهده‌شده رکوردها می‌شود و به‌تدریج قاعده محرمانگی برای تمام رکوردهای مجموعه داده برقرار می‌شود.

در شکل (۳) مثالی ساده از نحوه کارکرد روش پیشنهادی نمایش داده شده است. برای ساده‌سازی فقط سه رکورد نخست آورده شده است. مجموعه داده شامل شبه‌شناسه‌های A، B و C است. کار از پردازش رکورد نخست در قسمت (۱) آغاز می‌شود؛ فرض می‌شود که این رکورد رابطه (۴) را ارضا نمی‌کند؛ همان‌طور که در درخت متناظر دیده می‌شود، گره C آخرین گرهی است که در مسیر رکورد نخست، روی آن تصمیم‌گیری می‌شود؛ لذا مقدار این شبه‌شناسه باید حذف شود. رابطه $OC < EC$ برای رکورد دوم هم برقرار نیست.



(شکل-۳): نحوه عملکرد روش پیشنهادی
(Figure-3): Operation of the proposed technique

اگر رکورد r هنوز رابطه (۴) را ارضا نمی‌کند، باید عمل‌گر فرونشانی روی یکی از شبه‌شناسه‌های آن اعمال شود. مرحله نه با استفاده از درخت تصمیم DT از میان شبه‌شناسه‌ها، صفت q که کمترین بهره اطلاعاتی را دارد، انتخاب می‌کند. نحوه این کار در الگوریتم ارائه‌شده در شکل (۵) شرح داده می‌شود. مرحله ده فرونشانی را روی شبه‌شناسه q از رکورد r اعمال می‌کند.

مرحله دوازده میزان تخریب رکورد را از طریق δ محاسبه می‌کند و در صورتی که از حد آستانه δ' بیشتر باشد، رکورد مربوطه از مجموعه D^* حذف می‌شود. مرحله نوزده مجموعه D^* را که شامل رکوردهای گمنام شده است، در خروجی قرار می‌دهد.

نحوه انتخاب صفتی که باید فرونشاندن شود در شکل (۵) بیان شده است. برای این که اعمال فرونشانی کمترین تأثیر را در نتیجه طبقه‌بندی بگذارد، در هر مرحله شبه‌شناسه‌هایی که دورتر از ریشه درخت طبقه‌بندی هستند، فرونشاندن می‌شوند. ورودی‌های الگوریتم رکورد r مجموعه شبه‌شناسه‌ها Q و درخت تصمیم DT هستند و پس از اتمام کار شبه‌شناسه q_i برگردانده می‌شود.

QIWithLowestInfoGain(r, Q, DT)

Output: q_i

```

1:  $q_i \leftarrow Q[0]$ 
2:  $CQ \leftarrow \emptyset$ 
3: FOR each  $a \in Q$  DO
4:   IF  $a$  is not missed from  $r$  THEN
5:      $CQ \leftarrow CQ \cup \{a\}$ 
6:   END IF
7: END FOR
8:  $n \leftarrow \text{Root}(DT)$ 
9: WHILE  $n \notin \text{Leaves}(DT)$  DO
10:  IF  $n \in CQ$  THEN
11:     $q_i \leftarrow n$ 
12:     $CQ \leftarrow CQ - \{n\}$ 
13:  END IF
14:   $n \leftarrow \text{NextNode}(DT, n, r)$ 
15: END WHILE
16: IF  $|CQ| > 0$  THEN
17:   $q_i \leftarrow \text{RandomSelect}(CQ)$ 
18: END IF
19: RETURN  $q_i$ 
```

(شکل-۵): انتخاب یک شبه‌شناسه برای فرونشانی
(Figure-5): Selection of a quasi identifier for Suppression

الگوریتم ارائه‌شده در شکل (۴) روند گمنام‌سازی داده‌ها را نمایش می‌دهد. مجموعه داده (D)، مجموعه شبه‌شناسه‌ها (Q)، احتمال مقادیر مختلف برای هر شبه‌شناسه (P)، نرخ نمونه‌گیری (β) و حد تخریب مورد قبول (δ') به عنوان ورودی به الگوریتم داده می‌شود و با اجرای الگوریتم مجموعه داده گمنام شده (D^*) ایجاد می‌شود. مرحله یک شبه‌شناسه‌های عددی را بازه‌بندی و مقادیر این صفات را با بازه‌ای که در آن قرار دارند، تعویض می‌کند. طول بازه‌ها می‌تواند از حد پایین و بالای اعداد در هر صفت به دست آید. این کار به منظور کم‌تر ساختن پراکندگی مقادیر در صفات عددی، انجام می‌شود. مرحله دو درخت تصمیم را از روی مجموعه داده تشکیل می‌دهد. این درخت در مراحل بعد برای انتخاب صفتی که کمترین بهره اطلاعاتی را دارد، مورد استفاده قرار می‌گیرد. مرحله سه با نرخ β و به صورت تصادفی از مجموعه داده نمونه‌برداری می‌کند.

در مراحل ۴-۵ بررسی می‌شود که اگر اصل حرمانگی تعریف‌شده در رابطه (۴) برای تمامی رکوردها صادق باشد، اجرای الگوریتم خاتمه یابد؛ در غیر این صورت رکوردها برای پردازش بیشتر به مرحله شش فرستاده می‌شوند. مراحل ۶-۷ اطمینان مورد انتظار و اطمینان مشاهده‌شده رکورد r را از رابطه (۱) و رابطه (۲) محاسبه می‌کنند.

Anonymize (D, Q, P, β, δ')

Output: D^*

```

1:  $D^* \leftarrow \text{DiscretizeNumericQIs}(D, Q)$ 
2:  $DT \leftarrow \text{GenerateDecisionTree}(D^*)$ 
3:  $D^* \leftarrow \text{RandomSample}(D^*, \beta)$ 
4: WHILE all records in  $D^*$  does not satisfy Eq. 3-4 DO
5:   FOR each  $r \in D^*$  DO
6:      $ec \leftarrow \text{ExpectedConfidence}(r, D^*, Q, P)$ 
7:      $oc \leftarrow \text{ObservedConfidence}(r, D^*, Q)$ 
8:     IF  $ec < oc$  THEN
9:        $q \leftarrow \text{QIWithLowestInfoGain}(r, Q, DT)$ 
10:       $r \leftarrow \text{Suppress}(r, q)$ 
11:     ELSE
12:        $\delta = \text{Distortion}(r, Q)$ 
13:       IF  $\delta > \delta'$  THEN
14:          $D^* \leftarrow D^* - r$ 
15:       END IF
16:     END IF
17:   END FOR
18: END WHILE
19: RETURN  $D^*$ 
```

(شکل-۴): روند گمنام‌سازی داده‌ها
(Figure-4): Anonymization technique

گرفته می‌شود. لازم به ذکر است که خواندن داده‌ها از پایگاه داده و فهرست‌گذاری روی شبه‌شناسه‌ها می‌تواند باعث بهبود سرعت اجرای الگوریتم شود. فایل گمنام‌سازی شده به صورت arff. {نام الگوریتم گمنام‌سازی}. {نام مجموعه داده} در پوشه out ذخیره می‌شود. فرمت arff علاوه بر داده‌ها شامل سرآیندی است که فهرستی از نام و نوع صفات در آن نگهداری می‌شود. نتایج آماری حاصل از اجرای برنامه (زمان اجرا و دقت طبقه‌بندی توسط الگوریتم‌های مختلف) در فایل متنی با نام results.txt ذخیره می‌شود.

۵-۱- ارزیابی و نتایج به دست آمده

در این بخش به بررسی روش پیشنهادی (PM) و مقایسه عملکرد آن با روش ارائه شده در [20] پرداخته می‌شود. همان‌طور که گفته شد، نیاز روش‌های مبتنی بر عمومی‌سازی به درخت طبقه‌بندی، به عنوان نقطه ضعف آن‌ها تلقی می‌شود. برای بررسی این موضوع عملکرد روش AGF درحالی‌که درخت طبقه‌بندی در ورودی به آن داده نشده نیز مورد ارزیابی قرار گرفته است (در ادامه AGFIL نامیده می‌شود). در این حالت درختی به عمق یک که تمام مقادیر دامنه صفت در برگ‌ها و مقدار ANY در ریشه قرار دارد به صورت خودکار ایجاد می‌شود.

درخت تصمیم مورد نیاز روش پیشنهادی توسط الگوریتم J48 ایجاد و هنگام بازه‌بندی شبه‌شناسه‌های عددی، طول بازه‌ها ۵٪ دامنه مقادیر صفت در نظر گرفته شده است. در آزمایش‌های انجام شده، برای همه روش‌ها نرخ نمونه‌گیری برابر با ۹۰٪ و حد تخریب مورد قبول ۰,۶ است. درخت‌های طبقه‌بندی مورد استفاده روش AGF مطابق [26] ایجاد شده‌اند.

روند ارزیابی دقت طبقه‌بندی، در شکل (۶) نمایش داده شده است. در ابتدا مجموعه داده به ده قسمت مساوی تقسیم می‌شود. یک قسمت به عنوان مجموعه آزمون و نه قسمت دیگر به عنوان مجموعه آموزش انتخاب می‌شود؛ سپس الگوریتم گمنام‌سازی بر روی مجموعه آموزش اعمال و الگوریتم طبقه‌بندی با استفاده از آن آموزش داده می‌شود. در ادامه رکوردهای مجموعه آزمون برای اندازه‌گیری دقت طبقه‌بند ساخته شده، مورد استفاده قرار می‌گیرد. این فرآیند تا زمانی که تمام ده قسمت به عنوان مجموعه آزمون انتخاب شوند، ادامه می‌یابد و در انتها میانگین دقت طبقه‌بندی در تمام دورها محاسبه می‌شود.

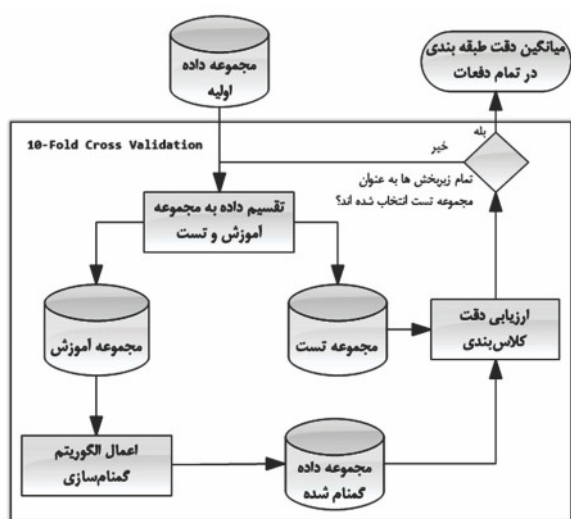
مراحل ۵-۲ مجموعه‌ای از شبه‌شناسه‌ها ایجاد می‌کند که مقدارشان در رکورد r هنوز فرونشاندن نشده است. شبه‌شناسه‌ای که در انتها بازگردانده می‌شود از این مجموعه (CQ) انتخاب خواهد شد.

مراحل ۸-۹، بررسی صفت‌ها را از ریشه درخت شروع می‌کند و تا زمانی که گره فعلی یکی از برگ‌های درخت نیست، به این کار ادامه می‌دهد. مراحل ۱۰-۱۲ در صورتی که گره فعلی عضوی از مجموعه شبه‌شناسه‌های CQ باشد، آن را به عنوان کم‌اهمیت‌ترین شبه‌شناسه‌ای که تاکنون بررسی شده قرار می‌دهد. همچنین حذف این شبه‌شناسه از CQ موجب می‌شود شبه‌شناسه‌ای که بیش از یک بار در مسیر ریشه تا برگ ظاهر می‌شود، فقط در نخستین دفعه مورد بررسی قرار گیرد.

مراحل ۱۶-۱۷ بررسی می‌کنند تمام شبه‌شناسه‌ها در مسیر ریشه تا درخت ملاقات شده‌اند یا خیر. در صورتی که هنوز چند شبه‌شناسه در مجموعه CQ وجود داشته باشد، یکی از آن‌ها به صورت تصادفی به عنوان شبه‌شناسه خروجی انتخاب می‌شود. مرحله نوزده صفت qi را که انتظار می‌رود فرونشاندن آن تأثیر کمی در نتیجه طبقه‌بندی داشته باشد در خروجی قرار می‌دهد.

۵- پیاده‌سازی طرح پیشنهادی و نتایج شبیه‌سازی و ارزیابی

پیاده‌سازی الگوریتم با استفاده از زبان جاوا و به کمک کتابخانه وکا [25] انجام گرفته است. وکا یکی از مشهورترین سکوها یادگیری ماشین است که به زبان جاوا و توسط دانشگاه وایکاتو نیوزلند توسعه پیدا کرده است. کاربر می‌تواند از طریق واسط گرافیکی و واسط خط فرمان از وکا استفاده کند و یا آن را به عنوان یک کتابخانه به پروژه برنامه‌نویسی خود اضافه کند. با توجه به امکانات فراوانی که API وکا برای برنامه‌نویس فراهم می‌کند، نیازی به کدنویسی برای عملیات سطح پایین نیست و بسیاری از قسمت‌های لازم برای پیاده‌سازی و ارزیابی الگوریتم مانند طبقه‌ها و توابع مفید برای کار با مجموعه داده‌ها و الگوریتم‌های طبقه‌بندی، توسط این کتابخانه انجام می‌گیرد. مجموعه داده‌های ورودی باید با فرمت csv ذخیره شده باشند. نخستین خط از فایل ورودی شامل نام صفات مجموعه داده است. طبق قرارداد نام شبه‌شناسه‌ها با #QI خاتمه می‌یابد و آخرین صفت به عنوان صفت حساس در نظر



(شکل-۶): روند ارزیابی الگوریتم‌های گمنام‌سازی
(Figure-6): Evaluation of anonymization algorithms

(جدول-۶): صفات‌های مجموعه‌داده Adult
(Table-6): Attributes of Adult Data set

تعداد/دامنه مقادیر	نوع	صفت	تعداد/دامنه مقادیر	نوع	صفت
14	گسسته	Occupation	17 - 90	پیوسته	Age
8	گسسته	Work-class	0 - 99999	پیوسته	Capital-gain
7	گسسته	Marital-status	0 - 4356	پیوسته	Capital-loss
6	گسسته	Relationship	1 - 16	پیوسته	Education-num
2	گسسته	Sex	13492 - 1490400	پیوسته	Final-weight
40	گسسته	Native-country	1 - 99	پیوسته	Hours-per-week
5	گسسته	Race	16	گسسته	Education

(جدول-۷): صفات‌های انتخاب شده به‌عنوان شبه‌شناسه
(Table-7): Selected attributes for quasi identifier

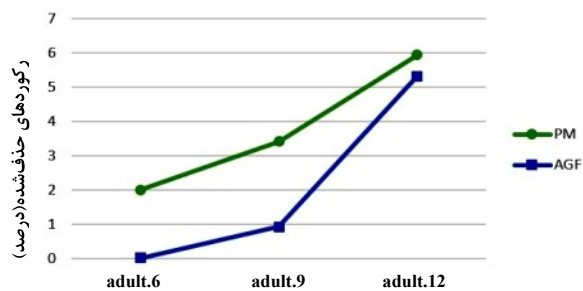
عنوان	شبه‌شناسه‌ها
Adult.6	age, education, sex, race, occupation, native-country
Adult.9	age, education, sex, race, occupation, native-country, workclass, marital-status, capital-loss
Adult.12	age, education, sex, race, occupation, native-country, workclass, marital-status, capital-loss, relationship, hours-per-week, capital-gain

۵-۲- دقت طبقه‌بندی

برای ارزیابی تأثیر گمنام‌سازی بر دقت طبقه‌بندی از الگوریتم‌های J48، NaiveBayes، PART و Logistic استفاده شده است. در نخستین آزمایش نه صفت از مجموعه‌داده به‌عنوان شبه‌شناسه انتخاب شده‌اند. همان‌طور که در شکل

(۷) نمایش داده شده است نتایج به‌دست‌آمده نشان می‌دهد دقت طبقه‌بندی‌هایی که با مجموعه‌داده گمنام‌شده توسط روش پیشنهادی آموزش دیده‌اند، بیشتر است. برای مقایسه بهتر، دقت طبقه‌بندی بر روی داده‌های گمنام‌شده نیز اندازه‌گیری شده است. نتایج به‌دست‌آمده برای این حالت، در نمودارها با

بیشتر است. با این حال وقتی تعداد زیادی از صفت‌ها به‌عنوان شبه‌شناسه در نظر گرفته شوند، در روش AGF عمومی‌سازی تا سطح بالایی از درخت طبقه‌بندی انجام خواهد شد و نتیجه مشابه استفاده از عملگر فرونشانی خواهد بود.

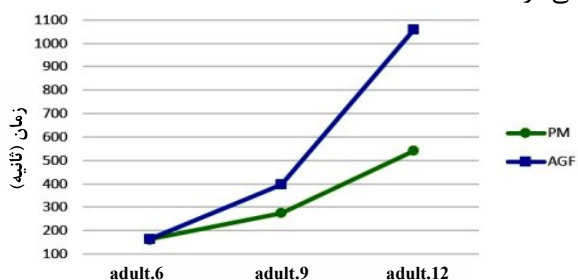


(شکل-۹): درصد رکوردهای غیرقابل انتشار
(Figure-9): A comparison on Unpublishable records

۴-۵- زمان اجرا

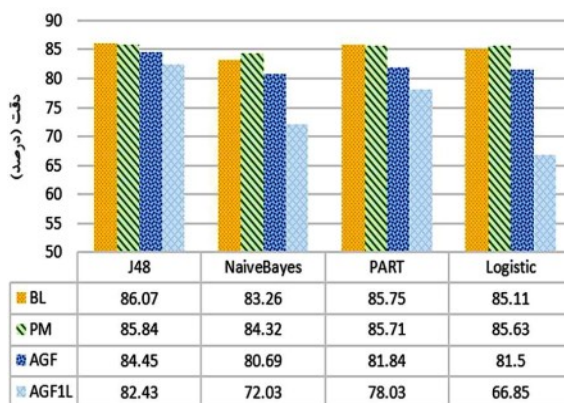
تعداد شبه‌شناسه‌ها و اندازه مجموعه‌داده عواملی هستند که بر زمان اجرای الگوریتم تأثیر دارند. برای آزمودن تأثیر این دو عامل، آزمایش‌های جداگانه‌ای انجام شده است. زمان اجرای کدها، بر روی سیستمی با پردازنده AMD Phenom II 1090T (3.2 GHz) و 4GB حافظه اصلی اندازه‌گیری شده است.

با افزایش تعداد شبه‌شناسه‌ها نقاط تمایز رکوردها بیشتر شده و انتظار می‌رود پردازش بیشتری برای برقرار کردن رابطه (۴) نیاز باشد. به‌منظور بررسی تأثیر تعداد شبه‌شناسه‌ها بر زمان مورد نیاز برای گمنام‌سازی، از مجموعه‌داده Adult با ۶، ۹ و ۱۲ شبه‌شناسه استفاده شد. همان‌طور که در شکل (۱۰) نمایش داده شده است، زمان مورد نیاز برای گمنام‌سازی توسط الگوریتم پیشنهادشده در مقایسه با AGF پایین‌تر است و همچنین تأثیرپذیری کمتری از تعداد شبه‌شناسه‌ها دارد. علت بهبود سرعت اجرا این است که AGF مقدار شبه‌شناسه‌ها را سطح به سطح در درخت طبقه‌بندی بالا می‌برد؛ اما در الگوریتم پیشنهادی مقدار شبه‌شناسه‌ها در یک مرحله حذف می‌شود.

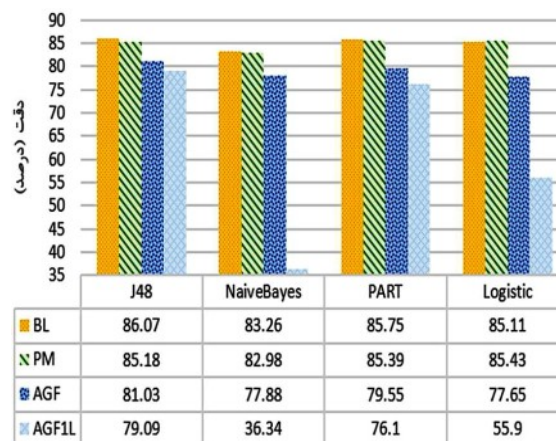


(شکل-۱۰): تأثیر تعداد شبه‌شناسه‌ها بر زمان اجرا
(Figure-10): The impact of Quasi Identifiers on simulation run time

نام BL نمایش داده شده است. به‌منظور بررسی تأثیر تعداد شبه‌شناسه‌ها بر کیفیت خروجی هر الگوریتم، در آزمایش دوم از دوازده شبه‌شناسه استفاده شده است. با توجه به نمودار شکل (۸) روش پیشنهادی در این حالت توانسته است، دقت طبقه‌بندی را نزدیک به حالت قبل نگاه دارد؛ درحالی‌که روش‌های AGF و AGF1L با افت دقت بیشتری مواجه شده‌اند.



(شکل-۷): دقت طبقه‌بندی - مجموعه‌داده Adult.9
(Figure-7): Classification accuracy in Adult.9 Data set



(شکل-۸): دقت طبقه‌بندی - مجموعه‌داده Adult.12
(Figure-8): Classification accuracy in Adult.12 Data set

۳-۵- رکوردهای غیرقابل انتشار

گفته شد که در طول اجرای الگوریتم عملگرهای دسته‌سازی مقدار تعدادی از صفات هر رکورد را تغییر می‌دهند و به‌وسیله معیار تخریب مشخص می‌شود که رکورد موردنظر کیفیت لازم را برای انتشار دارد یا خیر. شکل (۹) نسبت رکوردهای غیرقابل انتشار به کل رکوردهای مجموعه‌داده را نشان می‌دهد.

همان‌طور که انتظار می‌رفت، به‌دلیل استفاده از فرونشانی در روش پیشنهادی تعداد رکوردهای غیرقابل انتشار

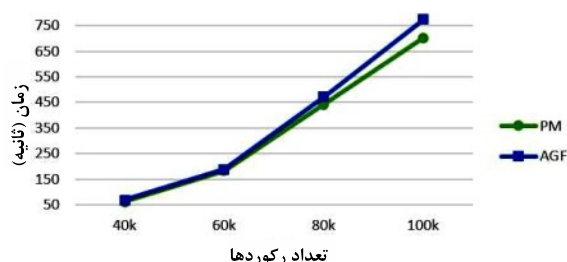
در طبقه‌بندی گمنام می‌شوند و یا در حالتی که رکوردهای مجموعه داده کم تعداد است، استفاده از روش پیشنهادی مناسب نخواهد بود.

7-References

۷- مراجع

- [1] B. C. M. Fung, K. Wang, A. Wai-Chee Fu and P. S. Yu, (2010), Introduction to Privacy-Preserving Data Publishing: Concepts and Techniques, Chapman and Hall/CRC.
- [2] J. Bennett and S. Lanning, (2007), "The Netflix Prize", Proceedings of the KDD Cup Workshop, pp. 3-6.
- [3] D. Nettleton, (2014), "Data Privacy and Privacy-Preserving Data Publishing," in Commercial Data Mining: Processing, Analysis and Modeling for Predictive Analytics Projects, Morgan Kaufmann, pp. 266-277.
- [4] B. Fung, K. Wang and P. Yu, (2010), "Privacy-Preserving Data Publishing: A Survey of Recent Developments", ACM Computing Surveys, vol. 42, no. 4,
- [5] L. Sweeney, (2002), "k-Anonymity: A Model for Protecting Privacy", International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, vol. 10, no. 5, pp. 557-570.
- [6] K. S. Babu, (2013), Utility-Based Privacy Preserving Data Publishing, PhD thesis, National Institute of Technology Rourkela.
- [7] N. Mohammed, B. C. M. Fung, P. C. K. Hung and C. K. Lee, (2009), "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.
- [8] D. Dou and S. Coulondre, (2012), "Detecting Privacy Violations in Multiple Views Publishing," in Database and Expert Systems Applications, Springer-Verlag Berlin Heidelberg, 506-513.
- [9] A. Anjum and G. Raschia, (2013), "Anonymizing Sequential Releases under Arbitrary Updates", Proceedings of the Joint EDBT/ICDT 2013 Workshops, pp. 145-154.
- [10] B. Fung, K. Wang and P. Yu, (2007), "Anonymizing Classification Data for Privacy Preservation", IEEE Transactions on Knowledge and Data Engineering, vol. 19, no. 5, pp. 711-725.
- [11] V. S. Susan and T. Christopher, (2014), "A Survey on Privacy Preservation in Data Publishing", International Journal of Computer Science and Mobile Computing, vol. 3, no. 3, pp. 188-193.

به منظور بررسی میزان مقیاس پذیری، در آزمایش‌های مختلف ۴۰، ۶۰، ۸۰ و ۱۰۰ هزار رکورد به ورودی الگوریتم داده شدند. همان‌طور که گفته شد، مجموعه داده Adult شامل ۴۸۸۴۲ رکورد است. برای افزایش تعداد رکوردها، برای هر صفت، احتمال داشتن مقادیر مختلف حساب شد؛ سپس بر طبق احتمالات به دست آمده، مقادیر مختلف صفات در کنار هم قرار داده شد و رکوردهای مورد نیاز تولید شدند. در این آزمایش شش صفت به عنوان شبه‌شناسه انتخاب شدند. نتایج به دست آمده در شکل (۱۱) نمایش داده شده است.



(شکل-۱۱): تأثیر تعداد رکوردها بر زمان اجرا

(Figure-11): The impact of records on simulation run time

(جدول-۸): مقایسه روش پیشنهادی و AGF

(Table-8): Comparison of proposed technique with AGF

روش	کاربرد	نیاز به درخت طبقه‌بندی	زمان	رکوردهای غیرقابل انتشار	دقت طبقه‌بندی
PM	طبقه‌بندی	خیر*	کمتر*	بیشتر	بیشتر*
AGF	عام‌منظوره*	بله	بیشتر	کمتر*	کمتر

۶- جمع‌بندی و نتیجه‌گیری

در این مقاله استفاده از عمل‌گر فرونشانی برای دستیابی به چهارچوب محرمانگی تعریف‌شده در [20] مورد بررسی قرار گرفت و به جای به کارگیری عمومی‌سازی سلولی از فرونشاندن صفاتی که بهره اطلاعاتی کمتری دارند، استفاده شد. اگرچه عمل‌گر فرونشانی نسبت به عمومی‌سازی، مقدار صفات را کم‌مفهوم‌تر می‌کند، ولی همان‌طور که نتایج به دست آمده نشان می‌دهد، استفاده درست از این عمل‌گر، علاوه بر این که می‌تواند نتیجه طبقه‌بندی را نزدیک به طبقه‌بندی روی داده‌های خام نگه دارد، مشکلات ساخت درخت طبقه‌بندی برای هر شبه‌شناسه را به همراه ندارد. در جدول (۸) روش پیشنهادی و روش ارائه‌شده در [20] با یکدیگر مقایسه شده‌اند. در هر ویژگی، روشی که عملکرد بهتری داشته با علامت ستاره (*) مشخص شده است. با توجه به این جدول، در شرایطی که داده‌ها با هدفی غیر از استفاده

Publishing", Journal of Information Security, vol. 4, no. 2, pp. 101-112.

[25] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann and I. H. Witten, (2009), "The WEKA Data Mining Software: An Update", ACM SIGKDD Explorations Newsletter, vol. 11, no. 1, pp. 10-18.

[26] "Taxonomy trees of the Adult data set", [Online]. Available: <http://ddm.cs.sfu.ca/dmsoft/Privacy/products/adultHierarchy.txt>. [Accessed 8 May 2016].

[27] "UCI Machine Learning Repository: Adult Data Set", [Online]. Available: <http://archive.ics.uci.edu/ml/datasets/Adult>. [Accessed 8 May 2016].

[28] M. Nergiz, C. Clifton and A. Nergiz, (2009), "Multirelational k-Anonymity", IEEE Transactions on Knowledge and Data Engineering, vol. 21, no. 8, pp. 1104-1117.

[۲۹] مهدی صادق پور، "حفظ محرمانگی در انتشار داده‌ها به وسیله گمنام‌سازی دسته‌ای"، پایان نامه کارشناسی ارشد مهندسی کامپیوتر – نرم‌افزار، دانشگاه گیلان، ۱۳۹۴.

[29] Mehdi Sadeghpour, "Privacy Preserving Data Publishing using Group Based Anonymization", MSc thesis in Software engineering, University of Guilan, 2015.



رضا ابراهیمی آتانی مدرک کارشناسی خود را از دانشگاه گیلان در سال ۱۳۸۱ دریافت و همچنین درجه کارشناسی ارشد و دکترای خود را از دانشگاه علم و صنعت در

سال ۱۳۸۳ و ۱۳۸۹ اخذ کرده است. در حال حاضر، ایشان عضو هیئت علمی و دانشیار در دانشکده فنی و مهندسی دانشگاه گیلان است. زمینه‌های پژوهشی مورد علاقه ایشان عبارتند از: طراحی و تحلیل الگوریتم‌ها و پروتکل‌های رمزنگاری و امنیتی جهت کاربرد در شبکه‌های رایانه‌ای و بی‌سیم.

نشانی رایانامه ایشان عبارت است از:

rebrahimi@guilan.ac.ir



مهدی صادق پور مدرک کارشناسی و کارشناسی ارشد خود را در سال‌های ۱۳۹۲ و ۱۳۹۴ از گروه مهندسی کامپیوتر دانشگاه گیلان اخذ کرد. موضوعات مورد علاقه ایشان مهندسی نرم‌افزار و طراحی و

توسعه سامانه‌های امن و با حفظ حریم خصوصی سازمانی است.

نشانی رایانامه ایشان عبارت است از:

mehdi.sadeghpour@live.com

[12] T. Dalenius, (1977), "Towards a Methodology for Statistical Disclosure Control", Statistik Tidskrift, vol. 15, 429-222.

[13] C. Dwork, (2006), "Differential Privacy," in Automata, Languages and Programming, Springer Berlin Heidelberg, pp. 1-12.

[14] K. Wang and B. C. M. Fung, (2006), "Anonymizing sequential releases", Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 414-423.

[15] A. Machanavajjhala, D. Kifer, J. Gehrke and M. Venkitasubramanian, (2007), "l-diversity: Privacy beyond k-anonymity", ACM Transactions on Knowledge Discovery from Data, vol. 1, no. 1, article 3.

[16] N. Li, T. Li and S. Venkatasubramanian, (2007), "t-Closeness: Privacy Beyond k-Anonymity and l-Diversity", IEEE 23rd International Conference on Data Engineering, pp. 106 - 115.

[17] Y. Rubner, C. Tomasi and L. J. Guibas, (2000), "The Earth Mover's Distance as a Metric for Image Retrieval", International Journal of Computer Vision, vol. 40, no. 2, pp. 99 - 121.

[18] N. Li, T. Li and S. Venkatasubramanian, (2010), "Closeness A New Privacy Measure for Data Publishing", IEEE Transactions on Knowledge and Data Engineering, vol. 22, no. 7, pp. 943-956.

[19] M. E. Nergiz, M. Atzori and C. Clifton, (2007), "Hiding the Presence of Individuals from Shared Databases", InProc. of ACM International Conference on Management of Data, pp. 665-676.

[20] A. S. Sattar, J. Li, X. Ding, J. Liu and M. Vincent, (2013), "A general framework for privacy preserving data publishing", Knowledge-Based Systems, vol. 54, 276-287.

[21] K. Wang, P. Yu and S. Chakraborty, (2004), "Bottom-Up Generalization: A Data Mining Solution to Privacy Protection", Fourth IEEE International Conference on Data Mining, pp. 249 - 256.

[22] B. Fung, K. Wang and Y. P.S, (2005), "Top-Down Specialization for Information and Privacy Preservation", Proc. 21st International Conference on Data Engineering, pp. 205-216.

[23] S. Kisilevich, L. Rokach, Y. Elovici and B. Shapira, (2010), "Efficient Multidimensional Suppression for K-Anonymity", IEEE Transactions on Knowledge and Data Engineering, vol. 22, no. 3, pp. 334 - 347.

[24] A. Hussien, N. Hamza and A. Hefny, (2013), "Attacks on Anonymization-Based Privacy-Preserving: A Survey for Data Mining and Data