

روشی جدید برای حفظ حریم خصوصی داده‌ها بر اساس تجزیه نامنفی ماتریس

لعیا علی احمدی پور و عزت ولی پور

بخش علوم کامپیوتر، دانشکده ریاضی و کامپیوتر، دانشگاه شهید باهنر کرمان، کرمان، ایران
بخش ریاضی کاربردی، دانشکده ریاضی و کامپیوتر، دانشگاه شهید باهنر کرمان، کرمان، ایران

تاریخ پذیرش: ۱۴۰۱/۷/۷

تاریخ دریافت: ۱۴۰۱/۵/۱۲

نوع مقاله: علمی-پژوهشی

چکیده

نگرانی شرکت‌ها از افشا و نقض حریم خصوصی کاربران در حال افزایش است و این امر منجر به تمرکز بسیاری از محققان بر روی توسعه روش‌های حفظ حریم خصوصی داده‌ها شده است. این روش‌ها روی داده‌های اصلی تأثیر می‌گذارند و با حفظ ویژگی‌ها و ارتباط بین آن‌ها، داده را بصورت دیگری منتشر می‌کنند. در این مقاله با استفاده از روش تجزیه نامنفی ماتریس‌ها، الگوریتمی برای تولید داده‌هایی با سطوح امنیتی متفاوت، جهت انتشار برای یک مجموعه داده اصلی پیشنهاد می‌شود. پیاده‌سازی‌ها با دو رویکرد متفاوت روی مجموعه داده‌های متنوع استاندارد نشان می‌دهند که الگوریتم پیشنهادی علاوه بر تولید داده‌هایی با سطوح امنیتی مختلف، ساختار محدودیت‌های داده‌ی اصلی را نیز حفظ می‌کند. روش مطرح شده در این مقاله روی مجموعه داده‌های با ابعاد بالا پیاده‌سازی شده، در حالیکه برخی روش‌های ریزتجمیع‌کننده مبتنی بر فازی قابلیت پیاده‌سازی روی آن‌ها را ندارند. همچنین، نتایج تجربی با استفاده از روش c -میانگین فازی نشان می‌دهد که میزان اطلاعات از دست رفته نیز بسیار ناچیز است. لذا، الگوریتم پیشنهادی می‌تواند داده‌هایی را منتشر کند که در عین حفظ حریم خصوصی، برای پردازشگران قانونی قابل استناد باشند.

عبارات و کلمات کلیدی: حفظ حریم خصوصی داده، روش خوشه‌بندی تجزیه نامنفی ماتریس، داده ماسک شده، داده منتشر شده، اطلاعات از دست رفته.

۱ مقدمه

در سال‌های اخیر رویکردها و روش‌های متنوعی جهت حفظ حریم خصوصی^۱ گسترش یافته‌اند. به طور کلی، مکانیزم‌های حافظ حریم خصوصی داده‌ها را می‌توان در دو رده‌ی روش‌های داده-محور (اهداف کلی)^۲ و روش‌های محاسبات-محور (اهداف خاص)^۳ قرار داد [۱۳]. در واقع روش‌هایی که برای انتشار داده مستقل از کاربرد گذشته آن استفاده می‌شوند، را داده-محور گویند [۱۲]. از طرفی روش‌های گسترش یافته برای مجموعه داده‌های خاص را محاسبات-محور می‌نامند [۱۵، ۱۶]. یکی از مهمترین زیررده‌های روش‌های داده محور، روش‌های آشفته‌ساز^۴ [۹] مانند ریزتجمیع‌کننده‌ها^۵ [۳، ۴] و روش‌های تعویض رتبه^۶ [۱۱] هستند که بر پایه آشفته کردن داده‌ها عمل می‌کنند. روش‌های ریزتجمیع‌کننده از یک الگوریتم خوشه‌بندی^۷ جهت افراز داده‌های اصلی به مجموعه‌هایی با حداقل k و حداکثر $2k$ عضو استفاده می‌کنند. سپس نماینده هر خوشه را تعیین و نهایتاً سایر اعضا خوشه را با نماینده جانشین می‌کنند [۱۴]. این جانشانی جزئیات اطلاعات اعضا را پوشانده و حریم خصوصی را حفظ می‌کند. علاوه بر این، یک فایل ریزتجمیع شده، به نوعی در خاصیت k -ناشناسی^۸ صدق می‌کند زیرا هر کدام از حداقل k عضو خوشه از سایرین قابل تمایز نیست. لازم به ذکر است که به مجموعه داده‌ای که پس از انجام عملیات محافظتی و حفظ حریم خصوصی به عنوان خروجی در اختیار سایرین قرار می‌گیرد، به اصطلاح داده منتشر شده^۹ می‌گویند. مسئله ریزتجمیع کننده روی داده‌های تک متغیره در زمان چند جمله‌ای^{۱۰} قابل پیاده سازی است [۸] در حالیکه ثابت شده است این مسئله برای داده‌های چند متغیره NP -سخت^{۱۱} است [۱۰].

تورا [۱۴] یک روش ریزتجمیع‌کننده بر پایه خوشه‌بندی C -میانگین فازی ارائه می‌دهد که

¹Privacy preserving

²Data-driven (general purpose)

³Computation-driven (specific purpose)

⁴Perturbative methods

⁵Microaggregation

⁶Swap ranking

⁷Clustering algorithm

⁸ k -anonymity

⁹Published data

¹⁰Polynomial time

¹¹NP- Hard

این روش توانایی بررسی برقراری یا عدم برقراری برخی از محدودیت‌های^{۱۲} اصلی روی مجموعه داده ماسک شده را نیز دارد. هر چند روش مذکور روی مجموعه داده‌های با ابعاد و اندازه‌ی بزرگ قابل پیاده‌سازی نیست.

در این مقاله روش جدیدی برای حفظ حریم خصوصی داده‌ها معرفی می‌شود که از تجزیه نامنفی ماتریس‌ها به عنوان روش خوشه‌بندی بهره می‌برد. ایده استفاده از خوشه‌بندی در الگوریتم پیشنهادی مشابه با روش‌های ریزتجمیع‌کننده است. هر چند خوشه‌بندی مورد استفاده در اینجا بر عکس ریزتجمیع‌کننده‌ها نامقید است. علاوه بر این هدف ما تولید چندین مجموعه داده ماسک شده^{۱۳} با سطوح امنیتی متفاوت است که همگی بتوانند کاندید داده منتشر شده باشند. در این صورت تصمیم‌گیرنده از آزادی بیشتری برای انتخاب داده منتشر شده برخوردار است. همچنین امروزه، بسیاری از تصمیمات کاربران، مبتنی بر پردازش داده‌های منتشر شده هستند. از طرفی مکانیزم‌های حفظ حریم خصوصی از کیفیت داده‌ها می‌کاهند، لذا معیارهایی برای سنجش کیفیت داده‌های منتشر شده پیشنهاد شده‌اند. یکی از مهمترین آن‌ها، معیار اطلاعات از دست رفته (IL)^{۱۴} است [۱۲]. بنابراین در پیاده‌سازی الگوریتم پیشنهادی روی مجموعه داده‌های واقعی علاوه بر تولید کاندیداهای متفاوت برای انتشار، به بررسی برقراری ساختار محدودیت‌های^{۱۵} اصلی روی داده‌های ماسک شده نیز پرداخته شده است. همچنین با استفاده از خوشه‌بندی c -میانگین فازی مقدار IL را نیز برای تمام داده‌های تولید شده محاسبه شده است. نهایتاً، این امکان را برای تصمیم‌گیرنده فراهم شده که با در نظر گرفتن تعامل بین حفظ حریم خصوصی و اطلاعات از دست رفته، یکی از داده‌های ماسک شده را برای انتشار انتخاب کند.

ادامه مقاله به صورت زیر ارائه شده است. در بخش ۲ مفاهیم و پیش‌نیازهای اولیه‌ی مورد استفاده در مقاله مطرح شده‌اند. بخش ۳ به جزئیات روش پیشنهادی و الگوریتم آن اختصاص یافته است. در بخش ۴ نتایج تجربی پیاده‌سازی الگوریتم و در پایان نتیجه‌گیری و پیشنهادهایی برای تحقیقات آتی مطرح شده‌اند.

¹² Constraints

¹³ Masked Data

¹⁴ Information loss

¹⁵ Constraints

۲ مفاهیم اولیه و پیش‌نیازها

در این بخش مفاهیم خوشه‌بندی تجزیه نامنفی ماتریس و اطلاعات از دست رفته به اختصار معرفی می‌شوند.

۱.۲ روش خوشه‌بندی تجزیه نامنفی ماتریس

در این مقاله مجموعه‌ی همه ماتریس‌های نامنفی $m \times n$ را با نماد $\mathbb{R}_{\geq 0}^{m \times n}$ نمایش داده شده‌اند.

تعریف ۱.۲. ([۱۷]) فرض کنید $D \in \mathbb{R}_{\geq 0}^{m \times n}$ و $k \leq \text{rank}(D)$ داده شده باشند. تجزیه D به ماتریس‌های پایه^{۱۶}

$W \in \mathbb{R}_{\geq 0}^{m \times k}$ و ضرایب^{۱۷} $H \in \mathbb{R}_{\geq 0}^{k \times n}$ بطوریکه $D \approx WH$ را NMF گویند.

یک دیدگاه متداول برای یافتن ماتریس‌های H و W حل مساله بهینه‌سازی زیر است [۷]:

$$\min_{W \geq 0, H \geq 0} \|D - WH\|_F, \quad (۱)$$

که در آن $\|\cdot\|_F$ محاسبه نرم فروبینوس^{۱۸} هست. بریا و همکاران [۲] برخی از الگوریتم‌های حل مساله (۱) را مرور کرده‌اند و به تعدادی از کاربردهای آن پرداخته‌اند.

NMF بطور گسترده‌ای در حوزه‌های مختلف علوم داده استفاده می‌شود. به عنوان

مثال در خوشه‌بندی، یک مجموعه از بردارهای داده $d_1, \dots, d_n \in \mathbb{R}^m$ وجود دارند که آن‌ها را به عنوان ستون‌های ماتریس D قرار می‌دهیم:

$$D = [d_1, \dots, d_n].$$

^{۱۶}Basis matrix

^{۱۷}Coefficient matrix

^{۱۸}Frobenius norm

منظور از یک تجزیه k -رتبه از ماتریس D یافتن تجزیه NMF i -ی است که در آن ماتریس‌های W و H دارای رتبه k هستند که مقدار k از قبل داده شده است. از منظر علوم داده این تجزیه معادل خوشه‌بندی مشاهدات d_i در k خوشه است. در این کاربرد از NMF ، $W_{m \times k}$ به عنوان ماتریس مراکز خوشه تعبیر می‌شود که ستون w_k بیانگر مرکز خوشه k -ام است. همچنین $H_{k \times n}$ ماتریس درجات عضویت خوشه‌ها است بطوریکه درایه h_{ij} نشان دهنده درجه عضویت d_j در خوشه i -ام است [۷].

۲.۲ اطلاعات از دست رفته

هر چند داده‌های ماسک شده برای حفظ حریم خصوصی و امنیت کاربران ضروری هستند، اما از طرف دیگر بخش مهمی از بقا، کارایی و توسعه شرکت‌ها به این داده‌ها و نتایج حاصل از پردازش آن‌ها وابسته است. در صورتی که داده ناصحیح یا مخدوش در اختیار آن‌ها قرار گیرد، می‌تواند روی تصمیمات کلیدی برای توسعه، تاثیر نامطلوب بگذارد. این تاثیر منفی بنا به کاربرد داده اثر مستقیمی روی جامعه هدف دارد؛ به عنوان مثال در مصارف پزشکی این امر ممکن است سلامت افراد را به خطر اندازد. همچنین در بسیاری از موارد تصمیم برپایه اطلاعات نادرست براحتی منجر به کاهش درآمدهای مالی شرکت و یا حتی ورشکستگی آن می‌گردد. بنابراین بررسی و علم به اینکه ماسک کردن داده برای امنیت چه میزان منجر به تحریف و از دست رفتن اطلاعات حساس داده‌ها می‌شود، اهمیت زیادی دارد. دانشمندان علوم داده معیارهایی^{۱۹} را برای اندازه‌گیری میزان اطلاعات از دست رفته گسترش داده‌اند.

تعریف ۲.۲. ([۱۲]) فرض کنید D مجموعه داده اصلی روی دامنه X و D' مجموعه داده ماسک شده آن روی دامنه X' باشند. برای یک تحلیل داده^{۲۰} مفروض (یعنی $f: X \rightarrow X'$) میزان اطلاعات از دست رفته برای مجموعه داده‌های D و D' به صورت

¹⁹Measures

²⁰Data analysis

زیر تعریف می شود:

$$IL_f(D, D') = \text{divergence}(f(D), f(D')), \quad (۲)$$

که دیورژانس راهی برای مقایسه دو عضو X' می باشد. همچنین برای هر دو عضو $D_1, D_2 \in X'$ روابط زیر باید برقرار باشند:

- $\text{divergence}(D_1, D_1) = 0$,
- $\text{divergence}(D_1, D_2) \geq 0$,
- $\text{divergence}(D_1, D_2) = \text{divergence}(D_2, D_1)$.

ملاحظه ۳.۲. روشن است که هر تابع نرم یا شبه نرم روی X' کاندیدایی برای محاسبه دیورژانس است. همچنین با توجه به نوع و کاربری داده ها، معیارهای متنوعی برای محاسبه دیورژانس ارائه می شوند (برای مشاهده مثال ها و توضیحات بیشتر به [۱۲] و منابع آن رجوع کنید).

۳ روش پیشنهادی

در این مقاله جهت حفظ امنیت و حریم خصوصی داده های مورد مطالعه، از دیدگاه خوشه بندی NMF استفاده می کنیم. در دنیای واقعی لازم است که داده ها برای استفاده در کاربردهای مختلف سطوح امنیتی متفاوتی داشته باشند. لذا، ما در این مقاله داده های ماسک شده متنوعی از نظر میزان شباهت به داده های اصلی تولید می کنیم تا قدرت اختیار تصمیم گیرنده را برای استفاده از این داده ها بالاتر ببریم. علاوه بر این به مساله حفظ محدودیت های لازم روی ویژگی های داده ها نیز می پردازیم و برقراری یا عدم برقراری آن ها را روی داده های ماسک شده بررسی می کنیم.

الگوریتم ۳ فرایند اجرای روش پیشنهادی را نشان می دهد.

[!h]

ورودی:

- ماتریس داده اصلی $D = [v_1, \dots, v_N] \in \mathbb{R}^{F \times N}$
- محدودیت‌های روی ویژگی‌ها (C_i) برای $i = 1, \dots, I$.

خروجی:

- یک خانواده از ماتریس‌های داده ماسک شده $A = \{A^k | k = 3, \dots, \text{rank}(D)\}$ است.
- خطاهای $error^k$ و e_i^k .

گام ۱. رتبه ماتریس D را بیابید و قرار دهید $r = \text{rank}(D)$.

گام ۲. برای رتبه‌های $k = 3, \dots, r$

خوشه بندی $NMF(D, k)$ را محاسبه کنید.

ماتریس‌های H^k و W^k را بیابید.

برای $j = 1, \dots, N$

بردارهای داده ماسک شده

$$a_j^k \approx \sum_{l=1}^k H_{lj}^k W_l^k \in \mathbb{R}^F \quad (3)$$

را محاسبه کنید.

قرار دهید $A^k = [a_j^k]_{j=1, \dots, N}$.

میزان خطای نسبی را حساب کنید:

$$error^k = \frac{1}{FN} \left(\sum_{j=1}^N \|v_j - a_j^k\|_1 \right). \quad (4)$$

گام ۳. برای ماتریس‌های ماسک شده A^k ، $k = 3, \dots, r$

برای $i = 1, \dots, I$

اگر محدودیت C_i روی A^k برقرار است قرار دهید $e_i^k = 0$ ،

روشی جدید برای حفظ حریم خصوصی داده‌ها بر اساس تجزیه نامنفی ماتریس — ۸۰

در غیر اینصورت با توجه به نوع محدودیت میزان خطای نقض آن،
 e_i^k ، را محاسبه کنید.

توضیح الگوریتم روش پیشنهادی. الگوریتم ۳ تعداد N داده را به عنوان ورودی می‌پذیرد که دارای F ویژگی کمی هستند. همچنین با توجه به تاثیرات بین داده‌ای در دنیای واقعی ممکن است محدودیت‌هایی روی آنها وجود داشته باشد. به عنوان مثال فرض کنید که یک کارخانه قصد تولید N محصول را داشته باشد و ویژگی‌های اول، دوم و سوم به ترتیب نشان‌دهنده هزینه نیروی انسانی، هزینه مواد اولیه و هزینه کل برای تولید هر واحد از محصولات باشند. بنابراین ویژگی سوم (سطر ۳-ام ماتریس D) ترکیب خطی از دو ویژگی دیگر است. لذا داده‌های ماسک شده نیز باید در این محدودیت صدق کنند. گام ۱ جهت ساخت داده‌های ماسک شده با سطوح امنیتی متفاوت (A^k) رتبه ماتریس داده (D) را محاسبه می‌کند.

گام ۲ تجزیه‌های NMF با رتبه‌های $k = 3, \dots, \text{rank}(D)$ را بدست می‌آورد. لازم به ذکر است تجزیه NMF برای رتبه‌های $k = 1, 2$ نتایج قابل استنادی تولید نمی‌کند، بنابراین این رتبه‌ها بررسی نشده‌اند. همچنین در این گام برای هر ماتریس ماسک شده A^k میزان خطای نسبی $error^k$ یعنی میزان تفاوت نسبی A^k با داده اصلی D محاسبه می‌شود.

گام ۳ خطاهای حاصل از نقض محدودیت e_i^k را برای هر مجموعه داده‌ی ماسک شده A^k و هر محدودیت C_i محاسبه می‌کند.

ملاحظه ۱.۳. الگوریتم ۳ متناظر با هر ماتریس ماسک شده A^k ، $k = 3, \dots, r$ یک دنباله از خطاهای $error^k$ و e_i^k ، $i = 1, \dots, I$ نیز تولید می‌کند. تصمیم‌گیرنده می‌تواند بر اساس خطاها، شرایط موجود و بویژه نوع داده مورد مطالعه؛ داده مناسب را جهت انتشار انتخاب نماید. به عنوان نمونه در مسائل داده‌کاوی بهتر است A^* -یی انتخاب شود که برای آن

$$error^* = \min_{k=3, \dots, r} error^k.$$

| هزینه ۱۶٪ | هزینه ۷٪ | کل هزینه | ساعات پرداخت شده | نرخ دستمزد | دستمزد | مجموع ساعت کل |
|-----------|----------|----------|------------------|------------|--------|---------------|
| f_1 | f_2 | f_3 | f_4 | f_5 | f_6 | f_7 |
| ۱۵ | ۲۳ | ۴۲۰۱ | ۲۳ | ۵۰ | ۱۱۵۰ | ۳۷ |
| ۱۲ | ۴۳ | ۵۹۹۳ | ۲۸ | ۷۰ | ۱۹۶۰ | ۳۷ |
| ۶۴ | ۲۲۹ | ۳۱۹,۲۷ | ۱۲ | ۸۴ | ۱۰۰۸ | ۲۵ |
| ۱۲ | ۴۵ | ۶۲,۰۷ | ۲۹ | ۷۳ | ۲۱۱۷ | ۳۰ |
| ۲۸ | ۳۹ | ۷۴,۲۱ | ۹ | ۳۰ | ۲۷۰ | ۴۰ |
| ۷۱ | ۱۰۲ | ۱۹۱,۵ | ۱۰ | ۶۳ | ۶۳۰ | ۲۰ |
| ۲۳ | ۶۴ | ۹۵,۱۶ | ۹ | ۷۴ | ۶۶۶ | ۱۰ |
| ۲۵ | ۱۰۲ | ۱۳۸,۱۴ | ۷۲ | ۳۰ | ۲۱۶۰ | ۸۰ |
| ۴۸ | ۲۳۰ | ۳۰۱,۷۸ | ۲۶ | ۳۰ | ۷۸۰ | ۳۵ |
| ۳۲ | ۵۰ | ۹۰,۶۲ | ۶ | ۴۵ | ۲۷۰ | ۱۵ |
| ۹۰ | ۲۰۰ | ۳۱۸,۴ | ۸ | ۴۵ | ۳۶۰ | ۱۵ |
| ۱۶ | ۱۰۰ | ۱۲۵,۵۶ | ۳۴ | ۵۵ | ۱۸۷۰ | ۴۵ |

جدول ۱: مجموعه داده‌های اصلی

از طرفی در مسائل آماری که ارتباطات میان داده‌ای از اهمیت زیادی برخوردارند، بهتر است A^* -ی جهت انتشار انتخاب شود که در آن میزان خطاهای نقض محدودیتها، e_i^k ، به نسبت کمتر باشد.

تورا [۱۳، ۱۴] مثال کوچکی از یک مجموعه داده همراه با تعدادی محدودیت ارایه داده است. در اینجا جهت روشنتر شدن روش پیشنهادی، الگوریتم ۳ را روی این مثال پیاده‌سازی می‌کنیم. لازم به ذکر است که این مجموعه داده در حالت کلی و گسترش یافته آن [۵] در بخش ۴ مطالعه خواهد شد.

مثال ۲.۳. جدول ۱ شامل داده‌های دوازده نفر است. این جدول اطلاعات افراد روی هفت ویژگی هزینه ۱۶٪، هزینه ۷٪، هزینه کل، ساعات پرداخت شده، نرخ دستمزد، مجموع دستمزد و ساعت کل را شامل می‌شود. همانطور که در جدول ملاحظه می‌کنید مقادیر این ویژگی‌ها با $f_j, j = 1, \dots, 7$ نشان داده شده‌اند. همچنین داده‌ها در سه محدودیت زیر صدق می‌کنند:

$$\begin{aligned}
 C_1: & f_3 = \alpha_1 f_1 + \alpha_2 f_2, \\
 C_2: & f_6 = f_4 * f_5, \\
 C_3: & f_4 \leq f_7.
 \end{aligned}
 \tag{5}$$

روشی جدید برای حفظ حریم خصوصی داده ها بر اساس تجزیه نامنفی ماتریس — ۸۲

ورودی: روشن است که D برابر با ترانواده ماتریس داده در جدول ۱ است. همچنین

$$.F = \gamma \text{ و } N = 12$$

$$.r = \text{rank}(D) = 6 \quad \text{گام ۱.}$$

گام ۲. برای تولید خانواده داده‌های ماسک شده $A = \{A^3, \dots, A^6\}$ ، گام ۲ چهار بار تکرار میشود. همچنین در هر تکرار میزان خطای نسبی ($error^k$) از (۴) محاسبه شده و که این مقادیر در جدول ۲ ارائه شده‌اند.

گام ۳. این گام به بررسی برقراری محدودیت‌ها روی ماتریس‌های ماسک شده اختصاص دارد.

بررسی محدودیت C_1

ابتدا روی ماتریس D با حل مساله کمترین مربعات^{۲۱}:

$$\min_{(\alpha_1, \alpha_2)} \|f_2 - (\alpha_1 f_1 + \alpha_2 f_2)\|_2 \quad (6)$$

مقدار $(\alpha_1^*, \alpha_2^*) = (1.16, 1.07)$ را محاسبه می‌کنیم. لازم به ذکر است که جهت حل مساله (۶) از روش تجزیه QR [۷] استفاده می‌کنیم. سپس مساله (۶) را به طور مشابه برای هر کدام از A^k ‌ها بازنویسی و (α_1^k, α_2^k) را محاسبه می‌کنیم. در اینجا برای محاسبه مقدار خطای نقض محدودیت C_1 ، فرمول زیر را پیشنهاد می‌کنیم:

$$e_1^k = |\alpha_1^k - \alpha_1^*| + |\alpha_2^k - \alpha_2^*| \quad (7)$$

بررسی محدودیت C_2

$f_6 = f_4 * f_5$ یک محدودیت غیرخطی است. ما فرمول زیر را جهت خطای حاصل از نقض محدودیت C_2 برای هر A^k ، تعریف می‌کنیم:

²¹Least square

$$e_{\Psi}^k = \frac{1}{12} \left(\sum_{j=1}^{12} \left| A_{\Psi j}^k - (A_{\Psi j}^k * A_{\Delta j}^k) \right| \right). \quad (8)$$

بررسی محدودیت C_3

در این حالت صدق یا عدم صدق نامساوی $f_{\Psi} \leq f_{\Delta}$ باید تحقیق شود. لذا

برای هر A_k ،

(۹)

$$e_{\Psi}^k = \begin{cases} \frac{1}{12} \sum_{A_{\Psi j}^k > A_{\Delta j}^k} (A_{\Psi j}^k - A_{\Delta j}^k) & \text{if } \exists j \in \{1, \dots, 12\} : A_{\Psi j}^k > A_{\Delta j}^k \\ 0 & \text{O.W.} \end{cases}$$

خلاصه نتایج را در جدول ۲ ملاحظه می‌کنید. علامت \checkmark نشان‌دهنده برقراری محدودیت در ماتریس داده ماسک شده است. البته در صورت عدم برقراری محدودیت، مقدار خطای نقض نیز ارائه شده است. در پایان تصمیم گیرنده آزادی زیادی جهت انتخاب یک داده ماسک شده برای انتشار دارد. به عنوان نمونه A^3 و A^4 هر دو تقریباً در تمام محدودیتها صدق می‌کنند هر چند میزان خطای نسبی متفاوتی دارند. در شرایطی که ترجیح می‌دهیم داده با شباهت بیشتر به داده اصلی (درعین وجود تفاوت) را منتشر کنیم، بهتر است A^3 را به عنوان داده منتشر شده انتخاب کنیم. اما اگر خطر افشا خیلی زیاد باشد، بهتر است داده A^4 را منتشر کنیم.

| A^k | A^3 | A^4 | A^5 | A^6 |
|--|------------------------|------------------------|------------------------|------------------------|
| $error^k$ | ۰/۱۱۱ | ۰/۱۸۷ | ۰/۱۵۳ | ۰/۱۵۴ |
| $(\alpha_{\Psi}^k, \alpha_{\Delta}^k)$ | (۱/۱۶, ۱/۰۷) | (۱/۱۶, ۱/۰۷) | (۱/۱۵, ۱/۰۵۹) | (۱/۱۲, ۱/۰۵) |
| C_{Ψ}^k | \checkmark | \checkmark | $e_{\Psi}^k = 0.017$ | $e_{\Psi}^k = 0.053$ |
| C_{Δ}^k | $e_{\Delta}^k = 0.002$ | $e_{\Delta}^k = 0.003$ | $e_{\Delta}^k = 0.001$ | $e_{\Delta}^k = 0.002$ |
| C_{Ψ}^k | \checkmark | \checkmark | \checkmark | \checkmark |

جدول ۲: نتایج اجرا الگوریتم ۳ روی مثال ۲.۳

۴ نتایج تجربی

در این بخش عملکرد روش پیشنهادی از دو دیدگاه مورد مطالعه قرار می‌گیرد. در دیدگاه اول، الگوریتم ۳ روی مجموعه داده استاندارد *Census data set* [۵] که یک محک ۲۲ در روش‌های ریزتجمع‌کننده و حافظ حریم خصوصی است، پیاده‌سازی می‌شود. رویکرد دوم به بررسی کارایی و نحوه عملکرد داده‌های ماسک شده در کاربردهای عملی به عنوان نمونه: داده‌کاوی، می‌پردازد.

پیاده‌سازی‌ها در نرم‌افزار مبتلب ۲۰۲۱ (a) در یک سیستم با پردازنده مرکزی Core i7 با حافظه اصلی ۸ Gb انجام شده‌اند. از آنجا که NMF روشی است که در اجراهای متفاوت نتایج یکسان تولید نمی‌کند و الگوریتم ۳ از NMF بهره می‌برد، لذا نتایج اعلام شده مقدار میانگین ده بار اجرا هستند.

۱.۴ رویکرد اول پیاده‌سازی

در رویکرد اول شکل اصلی مجموعه داده مطرح شده در مثال ۲.۳ یعنی داده استاندارد *Census data set* [۵] بررسی شده است. در این مجموعه، اطلاعات هزاروهشتاد نفر ($N = 1080$) روی سیزده ویژگی مختلف ($F = 13$) گردآوری شده که سه محدودیت (۵) نیز در آن برقرار است. خلاصه نتایج پیاده‌سازی الگوریتم ۳ روی مجموعه داده *Census data set* در جدول ۳ آمده است. توجه کنید که در این مثال $rank(D) = 12$ محاسبه شده است.

در توضیح جدول ۳ موارد زیر قابل ذکر هستند:

- با توجه به اینکه $rank(D) = 12$ به کمک الگوریتم ۳، نه مجموعه داده ماسک شده $(A^3 - A^{12})$ تولید می‌شود.

²²Benchmark

| A^k | A^7 | A^2 | A^5 | A^6 | A^7 | A^8 | A^9 | A^{10} | A^{11} | A^{12} |
|--------------|---|---|--|---|---|---|---|---|---|---|
| $error^k$ | ۰٫۰۰۰۵ | ۰٫۰۰۰۵ | ۰٫۰۰۱۱ | ۰٫۰۰۱ | ۰٫۰۰۰۱ | ۰٫۰۰۰۱ | ۰٫۰۰۱۲ | ۰٫۰۰۰۳ | ۰ | ۰ |
| (α_1) | $(\begin{smallmatrix} ۰٫۰۰۳۶ \\ ۰٫۰۰۴۱۰ \end{smallmatrix})$ | $(\begin{smallmatrix} ۰٫۰۰۳۶ \\ ۰٫۰۰۴۱۱ \end{smallmatrix})$ | $(\begin{smallmatrix} ۰٫۰۰۳۷۵ \\ ۰٫۰۰۴۰۸ \end{smallmatrix})$ | $(\begin{smallmatrix} ۰٫۰۰۳۶ \\ ۰٫۰۰۲۳۶ \end{smallmatrix})$ | $(\begin{smallmatrix} ۰٫۰۰۱۷۶ \\ ۰٫۰۰۰۶۴۹ \end{smallmatrix})$ | $(\begin{smallmatrix} ۰٫۰۰۲۲ \\ ۰٫۰۰۳۵۷ \end{smallmatrix})$ | $(\begin{smallmatrix} ۰٫۰۰۲۳ \\ ۰٫۰۰۴۲۸ \end{smallmatrix})$ | $(\begin{smallmatrix} ۰٫۰۰۲۳ \\ ۰٫۰۰۲۳۴ \end{smallmatrix})$ | $(\begin{smallmatrix} ۰٫۰۰۸۲ \\ ۰٫۰۰۰۳۳ \end{smallmatrix})$ | $(\begin{smallmatrix} ۰٫۰۰۴۵ \\ ۰٫۰۰۴۹۵ \end{smallmatrix})$ |
| e_1^k | ۰٫۰۰۷۸ | ۰٫۰۰۸۷ | ۰٫۰۰۴۳ | ۰٫۰۰۲۵ | ۰٫۰۰۱۶۲ | ۰٫۰۰۰۳۰ | ۰٫۰۰۲۰۴ | ۰٫۰۰۱۱۱ | ۰٫۰۰۰۸۰ | ۰٫۰۰۱۲۷ |
| e_1^k | ۷۲×10^{-7} | ۸۲×10^{-7} | ۲×10^{-7} | ۲×10^{-7} | ۵×10^{-7} | ۵×10^{-7} | ۱×10^{-7} | ۲×10^{-7} | ۳۸×10^{-6} | ۷×10^{-7} |
| e_1^k | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

جدول ۳: نتایج پیاده‌سازی الگوریتم ۳ روی مجموعه داده *Census data set*

• در سطر دوم برابر بودن میزان خطاهای نسبی برای دو نوع داده ماسک شده به معنای یکی بودن آن‌ها نیست. به عنوان نمونه اگر چه $error^2 = error^3$ اما مسلماً $A^2 \neq A^3$. در واقع این‌ها دو مجموعه داده ماسک شده متفاوت هستند که از نظر تصمیم‌گیرنده با توجه به شرایط دیگر ممکن است یکی از آن‌ها بر دیگری ارجح باشد.

• سطرهای سوم و چهارم به بررسی محدودیت اول (۵) می‌پردازند. با حل مساله (۶) روی مجموعه داده اصلی *Census data set*، $\alpha_1 = ۰٫۰۰۰۳۵$ و $\alpha_2 = ۰٫۰۰۴۱۷$ محاسبه شده‌اند. سطر سوم مقادیر α_1 و α_2 حاصل از حل مساله (۶) برای داده‌های ماسک شده متفاوت و سطر چهارم میزان خطای نقض محدودیت اول (۷) را نمایش می‌دهند.

• در سطر پنجم، میزان خطای نقض محدودیت دوم (۵) محاسبه شده از رابطه (۹) آمده است. با توجه به مقادیر حاصل، از نظر یک تصمیم‌گیرنده خوشبین محدودیت C_2 در مجموعه‌های ماسک شده نیز برقرار است. هر چند در این جدول ما برای تصمیم‌گیرنده بدبین خطاها را اعلام کرده‌ایم.

• نماد ✓ در سطر آخر نشانگر برقراری محدودیت سوم (۵) در تمام مجموعه‌های ماسک شده است.

در مقایسه با کارهای تورا [۱۳، ۱۴] با توجه به NP-سخت بودن روش‌های ریزتجمیع‌کننده، روش فازی پیشنهادی خود را روی زیر مجموعه بسیار کوچکی از مجموعه داده [۵] (به ترتیب با ۱۲ داده و سه ویژگی؛ ۱۲ داده و هفت ویژگی) پیاده‌سازی کرده است. روش

روشی جدید برای حفظ حریم خصوصی داده‌ها بر اساس تجزیه نامنفی ماتریس — ۸۶

های ایشان بر پایه یافتن مراکز خوشه متفاوت فقط با اعمال نویز با توزیع های نرمال متفاوت است و همچنین در [۱۴] تنها برقراری محدودیت اول و در [۱۳] سه محدودیت (۵) را بررسی کرده است.

۲.۴ رویکرد دوم پیاده‌سازی

در این بخش میزان اطلاعات از دست رفته را در کاربرد عملی (خوشه‌بندی) زمانیکه داده‌ها با روش پیشنهادی جهت تامین امنیت ماسک شده‌اند، می‌سنجیم. بدین منظور از شش مجموعه داده استاندارد مورد استفاده در روش‌های داده‌کاوی استفاده می‌کنیم. ابتدا برای هر مجموعه داده با روش c - میانگین فازی (FCM)^{۲۳} [۱] مجموعه داده اصلی را خوشه‌بندی می‌کنیم، سپس عملکرد روش FCM را با استفاده از معیار دقت نرخ خوشه‌بندی^{۲۴}:

$$CR_{FCM} = \frac{\text{تعداد نمونه‌هایی که بدرستی خوشه‌بندی شده‌اند}}{\text{تعداد کل نمونه‌ها}}, \quad (10)$$

ارزیابی می‌کنیم. در ادامه، روش FCM را روی مجموعه داده‌های ماسک شده پیاده‌سازی و CR_{FCM} را برای آن‌ها محاسبه می‌کنیم. در پایان، برای بررسی میزان اطلاعات از دست رفته از تعریف ۲ استفاده می‌کنیم. در این راستا، تابع f مساله خوشه‌بندی است و دیورژانس را مقایسه CR_{FCM} ها در نظر می‌گیریم و معیار زیر را پیشنهاد می‌دهیم:

$$IL(D, A^k) = |CR_{FCM}(D) - CR_{FCM}(A^k)|. \quad (11)$$

جزئیات مجموعه‌های داده و نتایج پیاده‌سازی در ادامه توضیح داده شده است.

- مجموعه داده Haberman's Survival شامل سیصد و شش داده ($N = 306$)
- روی سه ویژگی ($F = 3$) است که در دو کلاس دسته‌بندی می‌شوند [۲۱]. نرخ

²³Fuzzy c-means (FCM)

²⁴Clustering rate (CR)

خوشه‌بندی مجموعه داده اصلی برابر با $CR_{FCM} = ۰/۶۵۰۳$ و $rank(D) = ۳$ هستند. خلاصه نتایج پیاده‌سازی را در جدول ۴ مشاهده می‌کنید.

| | |
|------------|-----------------------|
| | $A^۳$ |
| $error^k$ | $۴/۰۹ \times ۱۰^{۱۴}$ |
| CR_{FCM} | ۰/۵۲۶۱ |
| IL | ۰/۱۲۴۲ |

جدول ۴: مجموعه داده *Haberman's Survival*

• مجموعه داده Iris شامل صد و پنجاه داده ($N = ۱۵۰$) روی چهار ویژگی ($F = ۴$) است که در سه کلاس دسته‌بندی می‌شوند [۲۰]. نرخ خوشه‌بندی مجموعه داده اصلی برابر با $CR_{FCM} = ۰/۸۹۳۳$ و $rank(D) = ۴$ هستند. خلاصه نتایج پیاده‌سازی را در جدول ۵ مشاهده می‌کنید.

| | | |
|------------|--------|------------------------|
| | $A^۳$ | $A^۴$ |
| $error^k$ | ۰ | $۲/۰۳۷ \times ۱۰^{۱۵}$ |
| CR_{FCM} | ۰/۹۰۰۰ | ۰/۸۹۳۳ |
| IL | ۰/۰۱ | ۰ |

جدول ۵: مجموعه داده *Iris*

• مجموعه داده Glass شامل دویست و چهارده داده ($N = ۲۱۴$) روی نه ویژگی ($F = ۹$) است که در شش کلاس دسته‌بندی می‌شوند [۲۲]. نرخ خوشه‌بندی مجموعه داده اصلی برابر با $CR_{FCM} = ۰/۵۳۷۴$ و $rank(D) = ۹$ هستند. خلاصه نتایج پیاده‌سازی را در جدول ۶ مشاهده می‌کنید.

| | | | | | | | |
|------------|-------|-------|-------|-------|-------|-------|-----------------------|
| | $A^۳$ | $A^۴$ | $A^۵$ | $A^۶$ | $A^۷$ | $A^۸$ | $A^۹$ |
| $error^k$ | ۴/۰۲ | ۳/۷۴ | ۲/۲۴ | ۳/۲۳ | ۴/۱۶ | ۲/۷۱ | $۶/۲۴ \times ۱۰^{۱۲}$ |
| CR_{FCM} | ۰/۵۷ | ۰/۵۷ | ۰/۴۰ | ۰/۵۶ | ۰/۵۷ | ۰/۵۶ | ۰/۵۳ |
| IL | ۰/۰۴ | ۰/۰۳ | ۰/۰۱ | ۰/۰۲ | ۰/۰۵ | ۰/۱ | ۰ |

جدول ۶: مجموعه داده *Glass*

- مجموعه داده Breast cancer شامل ششصد و نود و نه داده ($N = 699$) روی ده ویژگی ($F = 10$) است که در دو کلاس دسته‌بندی می‌شوند [۱۹]. نرخ خوشه‌بندی روی مجموعه داده‌ی اصلی برابر با $CR_{FCM} = 0.9836$ و $rank(D) = 9$ هستند. خلاصه نتایج پیاده‌سازی را در جدول ۷ مشاهده می‌کنید.

| | A^2 | A^4 | A^6 | A^8 | A^{10} | A^{12} | A^{14} |
|------------|--------|--------|--------|--------|----------|----------|--------------------|
| $error^k$ | ۱٫۲۷ | ۱٫۴۸ | ۱٫۸۳ | ۲٫۳۵ | ۳٫۹۱ | ۱۰٫۲۵ | 1×10^{15} |
| CR_{FCM} | ۰٫۹۵۵۷ | ۰٫۹۶۵۷ | ۰٫۹۶۱۴ | ۰٫۹۵۵۷ | ۰٫۹۵۷۱ | ۰٫۹۵۵۷ | ۰٫۹۵۴۲ |
| IL | ۰٫۰۲۷ | ۰٫۰۱۷ | ۰٫۰۲۲ | ۰٫۰۲۷ | ۰٫۰۲۶ | ۰٫۰۲۷ | ۰٫۰۲۹ |

جدول ۷: مجموعه داده Breast Cancer

- مجموعه داده Wine شامل صد و هفتاد و هشت داده ($N = 178$) روی سیزده ویژگی ($F = 13$) است که در سه کلاس دسته‌بندی می‌شوند [۱۸]. نرخ خوشه‌بندی روش c - میانگین فازی روی داده‌های اصلی برابر با $CR_{FCM} = 0.6854$ است. از آنجا که $rank(D) = 13$ پس الگوریتم ۳ یازده مجموعه داده متفاوت ($A^3 - A^{13}$) ایجاد می‌کند. خلاصه نتایج پیاده‌سازی الگوریتم ۳، روش FCM و اطلاعات ازدست رفته در جدول ۸ نشان داده شده‌اند.

| | A^3 | A^4 | A^5 | A^6 | A^7 | A^8 | A^9 | A^{10} | A^{11} | A^{12} | A^{13} |
|------------|--------|--------|--------|--------|--------|--------|--------|----------|----------|----------|------------------------|
| $error^k$ | ۰٫۲۴ | ۱٫۲۰ | ۰٫۰۷ | ۱٫۲۴ | ۱٫۲۳ | ۱٫۲۴ | ۲٫۲۵ | ۱٫۲۵ | ۱٫۹۸ | ۱٫۹۹۸ | 8.225×10^{10} |
| CR_{FCM} | ۰٫۶۸۵۴ | ۰٫۶۸۵۴ | ۰٫۷۴۷۲ | ۰٫۶۸۵۴ | ۰٫۶۸۵۴ | ۰٫۶۸۵۴ | ۰٫۶۸۵۴ | ۰٫۶۸۵۴ | ۰٫۶۸۵۴ | ۰٫۶۸۵۴ | ۰٫۶۸۵۴ |
| IL | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ |

جدول ۸: مجموعه داده Wine

- مجموعه داده Heart SPECT شامل دویست و شصت و هفت داده ($N = 267$) روی بیست و دو ویژگی ($F = 22$) است که در دو کلاس دسته‌بندی می‌شوند [۲۳]. نرخ خوشه‌بندی روی مجموعه داده‌ی اصلی برابر با $CR_{FCM} = 0.6292$ و $rank(D) = 22$ هستند. خلاصه نتایج پیاده‌سازی را در جدول‌های ۹ و ۱۰ مشاهده می‌کنید.

| | | | | | | | | | |
|------------|-------|-------|-------|-------|-------|-------|-------|----------|----------|
| | A^3 | A^4 | A^5 | A^6 | A^7 | A^8 | A^9 | A^{10} | A^{11} |
| $error^k$ | ۴/۱۵ | ۴/۹۱ | ۵/۴۱ | ۵/۷۳ | ۶/۴۲ | ۶/۸۲ | ۶/۹۸ | ۷/۶۶ | ۷/۸۳ |
| CR_{FCM} | ۰/۵۹ | ۰/۶۰ | ۰/۶۱ | ۰/۶۱ | ۰/۶۱ | ۰/۶۱ | ۰/۶۲ | ۰/۶۲ | ۰/۶۲ |
| IL | ۰/۰۳ | ۰/۰۲۶ | ۰/۰۱ | ۰/۰۱ | ۰/۰۱ | ۰/۰۱ | ۰/۰۰۷ | ۰/۰۰۳ | ۰/۰۰۷ |

جدول ۹: مجموعه داده SPECT Heart (قسمت اول)

| | | | | | | | | | | | |
|------------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| | A^{12} | A^{13} | A^{14} | A^{15} | A^{16} | A^{17} | A^{18} | A^{19} | A^{20} | A^{21} | A^{22} |
| $error^k$ | ۸/۸۷ | ۹/۸۹۱ | ۰/۰۲ | ۱۰/۸۹ | ۱۳/۴۰ | ۱۴/۷۴ | ۱۸/۲۷ | ۱۸/۲۰ | ۲۱/۳۳ | ۳۶/۴۲ | ۴۶/۷۱۵ |
| CR_{FCM} | ۰/۶۱ | ۰/۶۲ | ۰/۶۳ | ۰/۶۱ | ۰/۶۲ | ۰/۶۱ | ۰/۶۱ | ۰/۶۲ | ۰/۶۱ | ۰/۶۱ | ۰/۶۰ |
| IL | ۰/۰۱ | ۰/۰۰۷ | ۰/۰۰۳ | ۰/۰۱ | ۰ | ۰/۰۱ | ۰/۰۱ | ۰/۰۰۳ | ۰/۰۱ | ۰/۰۱ | ۰/۰۲ |

جدول ۱۰: مجموعه داده SPECT Heart (قسمت دوم)

در جدول‌های ۱۰-۴ میزان خطای نسبی $error^k$ برای $k = rank(D)$ از سایر k ها بیشتر و عدد خیلی بزرگی است. از طرف دیگر میزان اطلاعات از دست رفته برای تمام k ها تقریباً یکسان و مقدار ناچیزی است. لذا انتخاب $A^{rank(D)}$ به عنوان داده منتشر شده، ضمن حفظ حریم خصوصی به طور موثر، برای کاربران قانونی که قصد پردازش داده‌ها را دارند؛ نیز بسیار مناسب است.

۵ نتیجه‌گیری و پیشنهاد برای کارهای آتی

این مقاله به مساله چالش برانگیز انتشار داده‌ها با توجه به حفظ حریم خصوصی، حفظ ویژگی‌ها و محدودیت‌های داده‌های اصلی و کاهش اطلاعات از دست رفته، پرداخته است. در این راستا یک الگوریتم بر اساس تجزیه نامنفی ماتریس‌ها طراحی شده که داده‌هایی با سطوح امنیت متفاوت را برای انتشار تولید می‌کند. رفتار الگوریتم روی یک مجموعه داده استاندارد پیچیده در مسائل امنیت که همراه با محدودیت‌های ساختاری است، بررسی شده است. نتایج تجربی نشان می‌دهند که روش پیشنهادی در مقایسه با برخی از روش‌های فازی علاوه بر حفظ امنیت، ساختار محدودیت‌های داده‌های اصلی با اندازه و ابعاد بزرگ را همچنان روی داده‌های منتشر شده نیز حفظ می‌کند. از طرف دیگر نتایج پیاده‌سازی الگوریتم روی سه مجموعه داده متنوع نشان دهنده کارایی الگوریتم پیشنهادی

روشی جدید برای حفظ حریم خصوصی داده ها بر اساس تجزیه نامنفی ماتریس — ۹۰

از منظر افزایش امنیت و کاهش اطلاعات از دست رفته در مساله خوشه‌بندی، روی این مجموعه‌ها است. در پایان، با توجه به کارایی روش مطرح شده در این مقاله، توسعه روش تجزیه نامنفی ماتریس‌ها بصورت فازی به عنوان یک ریزتجمیع‌کننده را برای تحقیقات آتی پیشنهاد می‌کنیم.

مراجع

- [1] Bezdek, J.C., Ehrlich, R, Full, W. (1984) FCM: the fuzzy c-means clustering algorithm, *Computers and Geo-sciences* , 10 , 191-203.
- [2] Berrya, M.W, Browne, M, Langville, A.N, Pauca, V.P and R.J. Plemmons. (2007) Algorithms and applications for approximate nonnegative matrix factorization, *Computational Statistics and Data Analysis*, 52, 155–173.
- [3] Castro, O., Gentile, C., Spagnolo-Arrizabalaga, E. (2022) An algorithm for the microaggregation problem using column generation, *Computers and Operations Research*, 68, 105817.
- [4] César Fadel, A., Satoru Ochi, L., André de Moura Brito, J., Silva Semaan, G. (2021) Microaggregation heuristic applied to statistical disclosure control, *Information Sciences*, 548, 37-55.
- [5] Domingo-Ferrer, J., and Mateo-Sanz J.M. (2002) Practical data-oriented microaggregation for statistical disclosure control, *IEEE Transactions on Knowledge and Data Engineering*, 14, 189-201.
- [6] Edgar, B. Antoni, M. Agusti, A. (2022) Privacy-preserving process mining: A microaggregation-based approach, *Journal of Information Security and Applications*, 68, 103235.
- [7] Elden, L. (2007) *Matrix Methods in Data Mining and Pattern Recognition*, Society for Industrial and Applied Mathematics 106-110.

- [8] Hansen, S.L., Mukherjee, S. (2003) A polynomial algorithm for optimal univariate microaggregation, in IEEE Transactions on Knowledge and Data Engineering, 4, 1043-1044.
- [9] Kiran, A., and Shirisha, N. (2022) K-Anonymization approach for privacy preservation using data perturbation techniques in data mining. Materials Today: Proceedings.
- [10] Oganian, A, Domingo-Ferrer, J. (2000) On the Complexity of Optimal Microaggregation for Statistical Disclosure Control, Statistical Journal of the United Nations Economic Commission for Europe, 4, 345–354.
- [11] Rodriguez-Garcia, M., Batet, M., Sánchez, D. (2019) Utility-preserving privacy protection of nominal data sets via semantic rank swapping, Information Fusion, 45, 282-295.
- [12] Torra, V. (2017) Masking methods. In: Torra, V. (ed.) *Data Privacy: Foundations, New Developments and the Big Data Challenge*. Studies in Big Data, 28, 191–238.
- [13] Torra, V. (2008) Constrained microaggregation: adding constraints for data editing, Transactions on data privacy, 1, 86–104.
- [14] Torra, V. (2020) Fuzzy Clustering-based Microaggregation to Achieve Probabilistic K-anonymity for Data with Constraints, Journal of Intelligent and Fuzzy Systems, 39, 5999–6008.
- [15] Vaidya, J, Zhu, Y. and C. Clifton, *Privacy Preserving Data Mining*, in Advances in Information Security, Springer, 19 2006, 1-121.
- [16] Yao, A.C. (1982) Protocols for secure computations, 23rd Annual Symposium on Foundations of Computer Science, 160-164.

- [17] Wang, Y. X, and Zhang, Y. J. (2013) Nonnegative Matrix Factorization: A Comprehensive Review, in IEEE Transactions on Knowledge and Data Engineering, 6, 1336-1353.
- [18] <https://archive.ics.uci.edu/ml/datasets/wine>.
- [19] [https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+\(original\)](https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+(original)).
- [20] <https://archive.ics.uci.edu/ml/datasets/iris>.
- [21] <https://archive.ics.uci.edu/ml/datasets/haberman's+survival>.
- [22] <https://archive.ics.uci.edu/ml/datasets/glass+identification>.
- [23] <https://archive.ics.uci.edu/ml/datasets/spect+heart>.