

تاریخ دریافت مقاله: ۹۶/۰۴/۲۹  
تاریخ پذیرش مقاله: ۹۶/۰۹/۱۲

## روش بی‌نام‌سازی داده‌های حساس با حفظ سودمندی مبتنی بر مرتب‌سازی داده‌ها

محمد طه عسکری\*

کارشناسی ارشد نرم‌افزار - دانشکده فنی و مهندسی - دانشگاه آزاد دامغان - ایران  
پست الکترونیکی: taha.askaree@gmail.com

رضا مرتضوی

استادیار - دانشکده فنی و مهندسی - دانشگاه دامغان - دامغان - ایران  
پست الکترونیکی: r\_mortazavi@du.ac.ir

### چکیده:

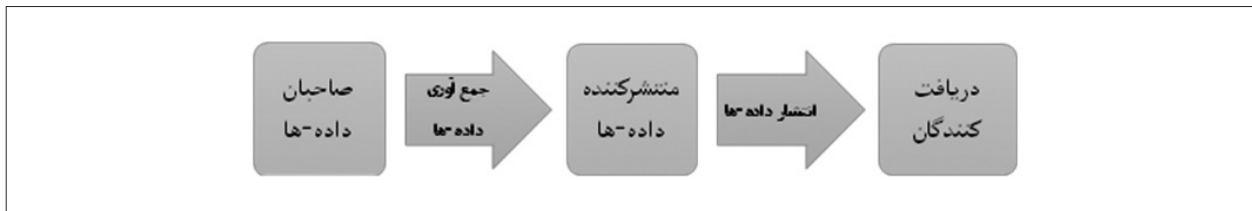
داده‌ها و حریم خصوصی را میسر ساخته است. واژه‌های کلیدی: حریم خصوصی، کنترل افشای آماری، ریزتجمیع داده‌گرا، k-بی‌نامی، ائتلاف اطلاعات

### ۱- مقدمه

پیشرفت‌های مداوم در فناوری‌های رایانه‌ای شرکت‌ها را قادر می‌سازد تا حجم عظیمی از داده‌های شخصی را جمع‌آوری نمایند، که منابع ارزشمند برای پژوهشگران و شرکت‌های دولتی می‌باشند [۱]. اسکن حریم خصوصی را به عنوان «حق فرد برای تعیین این‌که چه نوع اطلاعات (شخصی) در مورد وی در حال تبادل است» یا «کنترل یک فرد روی اطلاعات شخصی‌اش» تعریف می‌کند [۲]. PPDP<sup>۱</sup> یک ناحیه پژوهشی است که تلاش می‌کند داده‌ها را قبل از انتشار برای حفاظت اطلاعات حساس تغییر دهد. بین حریم خصوصی و سودمندی داده‌ها در PPDP توازن وجود دارد

ریزتجمیع خانوادگی از روش‌های کنترل افشای آماری (SDC) از ریزداده‌ها (رکوردهای اشخاص و یا شرکت‌ها) برای پوشش ریزداده‌ها است. داده‌های منتشر شده باید ضمن حفظ سودمندی، حریم خصوصی صاحبان داده‌ها را حفظ کند. بنابراین یک توازن بین سودمندی داده‌ها و حریم خصوصی وجود دارد. در این مقاله دو الگوریتم ریزتجمیع پیشنهاد شده است که ENFPN و NFPN++ نامیده می‌شود که ابتدا داده‌ها را با توجه به رکورد آخر و رکورد ماقبل آخر و پنج رکورد آخر مرتب می‌سازد، سپس یک بخش‌بندی با کمترین ائتلاف سودمندی نسبت به داده‌های مرتب شده جستجو می‌کند. نتایج تجربی نشان می‌دهد که روش پیشنهادی به کمترین ائتلاف اطلاعات نسبت به روش‌های ریزتجمیع گذشته دست یافته است و توازن بهتر بین سودمندی

\* نویسنده مسئول



شکل ۱: فرایند انتشار داده‌ها، بر گرفته از [۱۱]

[۱۳] و سرویس‌هایی بر اساس مکان (LBS) [۱۴] است که برای خوشه‌بندی یک مجموعه از رکوردها در گروه‌هایی با حداقل  $k$  رکورد می‌باشد [۱۵]. ریزتجمیع یکی از بیشترین روش‌هایی است که برای رضایتمندی  $k$ -بی‌نامی<sup>۸</sup> مورد استفاده قرار می‌گیرد که شامل دو مرحله است: ۱- بخش‌بندی داده‌های داخل گروه‌ها با حداقل  $k$  عنصر ۲- جابجایی هر رکورد با مرکز ثقل گروهش [۳]. چندین مهارت از  $k$ -بی‌نامی و متغیرهایش در مقالات امنیت پایگاه داده پیشنهاد شده است مانند  $(\alpha, k)$ -بی‌نامی [۱۶]،  $(L, \alpha)$ -تنوع<sup>۹</sup> [۱۷] و  $k-(\rho, \alpha)$ -بی‌نامی حساس<sup>۱۱</sup> [۱۸].

<sup>۱۲</sup>NFPN-MHM [۳] یک روش ریزتجمیع داده‌گرا است. این مقاله توسعه‌هایی را بر روی NFPN-MHM پیشنهاد می‌کند که NFPN++ و Enhanced-NFPN با بیشترین صحت و کمترین ائتلاف سودمندی نامیده می‌شوند.

بخش‌های این مقاله مطابق ذیل سازماندهی شده است: بخش ۲ یک نظر اجمالی روی روش‌های ریزتجمیع دارد. بخش ۳ شامل تعریف‌ها و مقدمه‌های مورد نیاز برای شرح روش‌های پیشنهادی است. بخش ۴ ارائه روش‌های پیشنهادی و بخش ۵ نمایش مقایسه‌ها بین روش‌های پیشنهادی و دیگر روش‌ها و نهایتاً بخش ۶ نتیجه می‌باشد.

## ۲- پیشینه

این بخش چند روش ریزتجمیع را به طور خلاصه شرح می‌دهد که بر روی ریزداده‌ها اعمال می‌شوند. پیشنهاد همه روش‌ها پیدا کردن یک گروه‌بندی از داده‌ها با کمترین ائتلاف

[۳]. تغییر یک پایگاه داده با کاهش کیفیت باید به گونه‌ای باشد که اطلاعات حساس افشا نشوند؛ تحلیل داده‌های تغییر یافته شبیه به تحلیل داده‌های اصلی است [۴]. داده‌های منتشر شده باید بی‌نام<sup>۲</sup> باشند درحالی‌که طبق حریم خصوصی داده‌ها حفظ شده هستند. بنابراین قبل از انتشار چنین داده‌هایی باید چند دستورالعمل برای ایجاد داده‌های بی‌نام اعمال شود. چندین مدل و سازوکار برای نیازمندی‌های این بی‌نامی<sup>۳</sup> در پایگاه داده‌ها، داده‌کاوی<sup>۴</sup> [۵، ۶] و مقالات آماری پیشنهاد شده است [۷-۹]. این حفاظت که در ازای کاهش سودمندی داده‌های محافظت شده برای کاربران مجاز نظیر پژوهشگران به دست آمده است، آن‌ها را مجبور به استفاده از روش‌های خاص برای استخراج دانش از داده‌های منتشر شده می‌نماید [۱۰]. شکل ۱ فرایند انتشار داده‌ها را نشان می‌دهد.

از آنجا که حریم خصوصی یک مفهوم چند وجهی می‌باشد، مسئله تضمین حفاظت مناسب از حریم خصوصی اشخاص آسان نیست: اطلاعات معین (حساس) راجع به اشخاص باید پوشیده بماند، هویت اشخاص باید محافظت شده باشد، یا اعمال اشخاص نباید قابل ردگیری<sup>۵</sup> باشد. عامل دیگری که بر پیچیدگی مسئله می‌افزاید وجود منابع متعدد داده‌هاست که تحلیل آن‌ها و مرتبط ساختن آن‌ها با هم، منجر به نشت نابجای اطلاعاتی می‌شود که نباید افشا شوند [۱۲].

ریزتجمیع<sup>۶</sup> یک تکنیک در SDC<sup>۷</sup> با کاربردهایی در علم کامپیوتر نظیر انتشار داده‌ها با حفظ حریم خصوصی

8- Location-Based Services  
9- k-anonymity  
10- diversity  
11- sensitive k-anonymity  
12- Nearest Far Point Next MHM

2- anonymous  
3- anonymity  
4- data mining  
5- traceable  
6- microaggregation  
7- Statistical Disclosure Control

اطلاعات و خطر افشا است که سودمندی بیشتر داده‌ها و حریم خصوصی را به دنبال دارد. این مسئله می‌تواند در زمان چندجمله‌ای برای مجموعه داده‌های تک‌متغیره حل شود، اما یک مسئله NP-Hard برای مجموعه داده‌های چند متغیره می‌باشد [۱۹].

یک مجموعه (فایل) ریزداده، شامل مجموعه‌ای از رکوردهاست که هر یک حاوی اطلاعات یک موجودیت مشخص (مانند فرد یا سازمان) هستند، و معمولاً به صورت یک ماتریس با  $n$  سطر ( $n$  رکورد) و  $v$  ستون (مقادیر  $v$  خصیصه برای هر رکورد) در نظر گرفته می‌شود [۲۰]. در ریزداده شناسه‌ها<sup>۱۳</sup> مجموعه‌ای از خصیصه‌ها هستند که حاوی اطلاعاتی است که صریحاً موجودیت متعلق به رکورد (صاحب رکورد) را مشخص می‌کنند [۲۱]؛ مانند نام یا شماره حساب بانکی، و یا کد پرسنلی کارمندان یک سازمان. و همچنین نیمه‌شناسه‌ها<sup>۱۴</sup> مجموعه خصیصه‌هایی هستند که به صورت بالقوه صاحب رکورد را می‌شناسانند [۲۱]؛ برخلاف خصیصه‌های شناسه، به تنهایی نمی‌توانند یک موجودیت را شناسایی کنند، اما ممکن است بتوان با کنار هم قرار دادن چند خصیصه نیمه‌شناسه، دقیقاً به یک موجودیت مشخص رسید [۲۰]. به عنوان مثال، ممکن است سن، کد پستی و درآمد ماهیانه افراد در کنار هم، آن‌ها را از یکدیگر به خوبی متمایز نمایند.

روش‌های ریزتجمیع به دو دسته اندازه‌ثابت و داده‌گرا تقسیم می‌شوند. روش‌های اندازه‌ثابت داده‌ها را درون گروه‌هایی با اندازه  $k$  تفکیک می‌کنند.

یکی از معروف‌ترین روش‌های ریزتجمیع اندازه‌ثابت MDAV<sup>۱۵</sup> [۲۲] است. این روش مکرراً خوشه‌ها را تولید می‌کند. در هر تکرار، ابتدا مرکز ثقل از رکوردهای خوشه‌بندی نشده را دوباره محاسبه و دو خوشه جدید تشکیل می‌دهد. اولین خوشه شامل  $r$  دورترین رکورد از مرکز ثقل سراسری و  $k-1$  نزدیک‌ترین همسایه‌اش و دومین آن شامل  $s$  دورترین رکورد از  $r$  و  $k-1$  نزدیک‌ترین

همسایه‌ها از  $s$  می‌باشد. نهایتاً عضوهای خوشه‌های جدید از لیست را حذف می‌کند. این تکرارها تا کمتر از  $m < 2k$  رکوردهای خوشه‌بندی نشده باقیمانده ادامه دارد. اگر  $m > k-1$  باشد همه رکوردهای باقیمانده یک خوشه جدید را تشکیل می‌دهند یا در غیر این صورت هر رکورد باقیمانده به نزدیک‌ترین خوشه‌اش متصل می‌شود.

روش دیگر اندازه‌ثابت شناخته شده CBFS<sup>۱۶</sup> [۲۳] است. CBFS شبیه MDAV است اما یک‌بار مرکز ثقل را محاسبه می‌کند و یک خوشه جدید در هر تکرار تشکیل می‌دهد. این روش سریع‌تر از MDAV است چون مرکز ثقل فقط یک‌بار محاسبه می‌شود. ساختمان داده‌های درخت KD در KD-CBFS [۲۴] با شتاب دادن به رویه جستجوی  $k$  نزدیک‌ترین همسایه‌ها مورد استفاده است. این پیچیدگی CBFS را از  $O(n^2)$  به  $O(n \log n)$  بهبود می‌دهد و قابلیت اجرا برای حجم داده‌های بزرگ را ایجاد می‌کند.

یک نسخه از MDAV پیشنهاد شده است که MDAV-generic [۲۵] نامیده می‌شود. آستانه پایانی تکرارهای MDAV عمومی با  $3k$  تنظیم شده است.

در مقابل، روش‌های داده‌گرا خوشه‌های با اندازه متغیر را تولید می‌کنند. نسخه داده‌گرای MDAV، V-MDAV [۱۵] است. بعد از ساخت هر خوشه با اندازه  $V$ -MDAV،  $k$  خوشه را تا  $2k-1$  رکورد گسترش می‌دهد.  $e_{min}$  یک رکورد خوشه‌بندی نشده را مشخص می‌کند که حداقل فاصله را از نزدیک‌ترین عضو خوشه‌بندی شده‌اش دارد. این فاصله،  $d_{in}$  و فاصله بین  $e_{min}$  و نزدیک‌ترین رکورد خوشه‌بندی نشده‌اش،  $d_{out}$  می‌باشد. رکورد  $e_{min}$  به خوشه اضافه می‌شود، اگر  $d_{in} < \gamma d_{out}$  باشد که  $\gamma$  یک آستانه تعریف شده توسط کاربر است.

روش دیگر GSMS<sup>۱۷</sup> [۲۶] است. این روش دارای دو مرحله می‌باشد، در مرحله اول  $\left\lfloor \frac{N}{k} \right\rfloor$  تکرار انجام می‌شود که در هر تکرار خوشه نامزدی (خوشه‌ای شامل رکورد نامزد برای مرکز خوشه بودن و  $k-1$  نزدیک‌ترین همسایه‌اش) که

13- identifiers  
14- quasi-identifiers  
15- Maximum Distance to Average Vector

16- Centroid-Based Fixed-Size  
17- Group Selection Based on Sequential Minimization of SSE

باعث کمینه شدن  $SSE^{18}$  سجاری از رکوردهای باقیمانده شود را رد می‌کند. در پایان این مرحله، خوشه  $k$  عضوی تشکیل می‌گردد. در مرحله دوم تعداد اندک رکوردهای باقیمانده را به نزدیک‌ترین خوشه‌ها منتسب می‌نماید.  $MHM^{19}$  [۲۷] یک راه حل بهینه را در زمان چند جمله‌ای برای مجموعه داده‌های تک‌متغیره میسر می‌سازد. این روش همه داده‌ها را در یک ترتیب صعودی مرتب می‌کند، سپس خوشه‌های با مقدار  $SSE$  کمینه برای داده‌های مرتب شده را جستجو می‌کند. جزئیات بیشتر این الگوریتم در بخش ۴ شرح داده شده است.  $MHM$  فقط برای مجموعه داده‌های تک‌متغیره بهینه است، اگرچه روش‌های اکتشافی گوناگونی وجود دارند که مرتب‌سازی داده‌های چندمتغیره را پیشنهاد می‌کنند که با اعمال  $MHM$  یک قسمت‌بندی بهینه با رابطه مرتب‌سازی داده شده ارائه می‌نماید (لزوماً این معادل با بخش‌بندی بهینه مجموعه داده‌ها نیست). به عنوان مثال ژینیتا و همکاران [۲۸] یک روش بر اساس منحنی‌های فضا-پرکن<sup>۲۰</sup> هیلبرت پیشنهاد کرده‌اند تا همه رکوردهای داده‌ها در یک توالی مرتب شود. این روش خیلی سریع است اما مجموعه داده‌های محافظت شده با اتلاف اطلاعات بالا را تولید می‌کند.

$IMHM^{21}$  [۱] یک پیاده‌سازی کارا از  $MHM$  است که وزن‌های لبه را به صورت تکراری محاسبه می‌کند، بنابراین  $O(nk)$  عملیات حسابی نیاز دارد.

روش دیگری که  $MHM$  استفاده می‌کند مرتب‌سازی داده‌ها با روش نزدیک‌ترین نقطه بعدی (NPN) [۲۹] است.  $NPN$  یک توالی از رکوردها، در شروع با انتخاب دورترین رکورد از مرکز ثقل را تولید می‌کند. در هر مرحله اگر  $r$  انتهای توالی باشد، نزدیک‌ترین رکورد به  $r$  به آن اضافه می‌شود، تا همه رکوردها به توالی اضافه شوند. ساختمان توالی با پیروی از نزدیک‌ترین رکوردها می‌تواند یک الگوریتم  $NPN$  سردرگم را تولید کند و یک رکورد دور را

18- Sum of Square Error

19- Multivariate Hansen Mukherjee algorithm

20- space-filling

21- Iterative MHM

22- Nearest Point Next MHM

به انتهای توالی اضافه می‌کند.

NFPN-MHM [۳] این مسئله را نشانی‌دهی می‌کند. این

روش با پیچیدگی زمانی یکسان با  $NPN$  اجرا می‌شود. این مقاله گسترش‌هایی را روی  $NFPN-MHM$  پیشنهاد می‌کند. روش‌های پیشنهادی جدید معمولاً سودمندی بیشتر و مجموعه داده‌های محافظت شده را تولید می‌نمایند.

### ۳- تعریف‌ها

یک مجموعه داده عددی  $T$  نرمال شده شامل  $n$  رکورد موجود است. هر رکورد  $d$  خصوصیت عددی دارد و می‌تواند به عنوان یک نقطه در فضای  $d$ -بعدی ارائه شود. برای هر رکورد داده  $x_i$ ،  $1 < i < n$ ،  $x_{ij}$  زامین خصوصیت از  $x_i$  را مشخص می‌کند. یک عدد صحیح  $k > 1$  در اولین قدم از الگوریتم ریزتجمیع داده شده است. این روش، داده‌ها را درون  $g$  گروه مجزا با حداقل  $k$  عضو تقسیم می‌کند به طوری که  $\sum_{i=1}^g n_i = n$  که  $n_i$  اندازه گروه نام است.  $x_{ij}^p$  زامین خصوصیت از آمین رکورد از  $p$  امین گروه می‌باشد،  $c^p$  مرکز ثقل  $p$  امین گروه و میانگین برداری از همه اعضای گروه است.  $c^p$  مانند رابطه ۱ محاسبه می‌شود:

$$c_j^p = \frac{1}{n_p} \sum_{i=1}^{n_p} x_{ij}^p \quad (1)$$

که  $c_j^p$  زامین خصوصیت از  $c^p$  و  $n_p$  اندازه گروه است. مرکز ثقل سراسری که به وسیله  $c$  مشخص شده، مرکز ثقل همه رکوردها در  $T$  است و می‌تواند با یک شیوه مشابه محاسبه شود. در انتشار داده‌ها با حفظ حریم خصوصی یک توازن بین سودمندی داده‌ها و محرمانگی وجود دارد.  $IL^{22}$  یک معیار محبوب برای تعیین کیفیت سودمندی داده‌های بی‌نام است و مانند رابطه ۲ محاسبه می‌شود:

$$IL = 100 \times \frac{SSE}{SST} \quad (2)$$

که  $SSE$  (رابطه ۳) جمع مربع خطا برای هر رکورد از مرکز ثقل خوشه‌اش است و  $SST$  (رابطه ۴) جمع مربع خطا برای هر رکورد از مرکز ثقل سراسری و رابطه ۵ فاصله

23- Information Loss

بین دو رکورد است :

$$SSE = \sum_{p=1}^g \sum_{i=1}^{n_p} distance(x_i^p, c^p) \quad (3)$$

$$SST = \sum_{i=1}^n distance(x_i, c) \quad (4)$$

$$distance(x_i, x_j) = \sum_{l=1}^d (x_{il} - x_{jl})^2 \quad (5)$$

مرتب شده همانند NFPN-MHM استفاده می‌کنند، اما تفاوت در شیوه مرتب‌سازی داده‌ها است. این دو مرحله از این دو الگوریتم عبارتند از: ۱- مرتب‌سازی داده‌ها و ۲- خوشه‌بندی با استفاده از MHM. این دو مرحله مطابق ذیل توصیف می‌شوند: مراحل‌نهایی تجمیع و غیرنرمال‌سازی مجموعه داده‌های محافظت شده همانند کارهای قبلی در مقالات هستند و در اینجا بررسی نمی‌شود.

#### ۴-۱- گام مرتب‌سازی داده‌ها

اولین قدم مرتب کردن داده‌ها در یک توالی از رکوردها است، به گونه‌ای که مجاورت رکوردهای نزدیک در فضای دامنه اصلی به قدر ممکن در توالی به‌دست آمده حفظ شوند. به این ترتیب داده‌های محافظت شده سودمندتر تولید می‌شود زیرا در گام بعدی، به طور پی‌درپی فقط رکوردهای مجاور با یکدیگر خوشه‌بندی می‌شوند، بنابراین خوشه‌های متراکم‌تر تولید می‌شوند. شکل ۲ رویه مرتب‌سازی الگوریتم NFPN++ و شکل ۳ الگوریتم ENFPN را نشان می‌دهد.

دورترین نقطه از مرکز ثقل سراسری سر توالی است. نقاط بعدی به انتهای توالی بر اساس رابطه ۹ اضافه می‌شوند.  $x_i$  را به‌عنوان انتهای توالی در نظر بگیرید،  $x_j$  به  $x_i$  اضافه خواهد شد اگر و فقط اگر فاصله رابطه ۹ برای  $x_j$  کمینه باشد:

$$D_{ij} = \begin{cases} d_{ij} & \text{if } (d_{ij} < 1) \\ distance(x_j, x_i) & \text{otherwise} \end{cases} \quad (9)$$

که  $d_{ij}$  طبق رابطه ۱۰ تعریف می‌شود:

$$d_{ij} = \frac{distance(x_j, x_i)}{distance(x_j, c)} \quad (10)$$

فاصله بین  $x_j$  و  $x_i$  مانند  $d_{ij}$  محاسبه می‌شود فقط وقتی که  $x_j$  نزدیک‌تر به  $x_i$  از مرکز ثقل است (خط ۶ و ۷ از شکل ۴). در شکل ۲، رکورد آخر و رکورد rLast ماقبل آخر می‌باشد. همچنین ۲ وزن یا اهمیت رکورد آخر و (1-2) وزن رکورد ماقبل آخر می‌باشد. همچنین می‌توان با توجه به

وقتی  $IL=0$  است یعنی این‌که هر رکورد با مرکز ثقل خوشه‌اش معادل است، بنابراین اطلاعاتی از دست رفته نیست و داده‌ای تغییر پیدا نکرده است.

از طرف دیگر توازن به محرمانگی اشاره دارد. خطر افشا با معیار خطر افشای اطلاعات حساس بیان می‌شود.  $DLD^{24}$  (رابطه ۶) یک روش محبوب برای تعیین کیفیت خطر افشا است:

$$DLD = 100 \times \frac{m}{n} \quad (6)$$

که  $m$  تعداد کلی رکوردهای منطبق است. فرض کنید  $T'$  نسخه بی‌نام شده از  $T$  است و رکورد  $x_i \in T$  با یک رکورد بی‌نام شده  $x'_i \in T'$  تعویض شده است. رکورد  $x_i$  منطبق است اگر و فقط اگر رابطه ۷ برقرار باشد:

$$x_i = \arg \min_{x_j} (distance(x_j, x'_i)) \quad (7)$$

معمولاً افزایش اتلاف اطلاعات، کاهش خطر افشا و بالعکس را نتیجه می‌دهد. بنابراین چند معیار وجود دارد که توازن بین  $IL$  و  $DLD$  سنجیده می‌شود.  $SI^{25}$  یک معیار است که به وسیله رابطه ۸ محاسبه می‌شود.

$$SI = a \cdot DLD + (1-a) \cdot IL \quad (8)$$

که  $0 \leq a \leq 1$  اهمیت  $DLD$  در  $SI$  را کنترل می‌کند.

#### ۴- روش‌های پیشنهادی

در الگوریتم NFPN-MHM [۳] تشکیل توالی بر اساس آخرین نقطه صورت می‌گرفت. در روش اول به رکورد آخر و ماقبل آخر و در روش دوم به پنج رکورد آخر نیز توجه شده که نتایج بهتری را به دنبال دارد. NFPN++ و Enhanced-NFPN الگوریتم MHM را روی داده‌های

24- Distance base Linkage Disclosure  
25- Score Index

**Input:** dataset  $T$   
**Output:** ordered sequence  $o$  on records of  $T$

```

1  $l \leftarrow T$ 
2  $o \leftarrow []$ 
3  $c \leftarrow \text{global\_centroid}$ 
4  $rNew \leftarrow$  furthest record from  $c$ 
5 add  $rNew$  to the end of  $o$ 
6 remove  $rNew$  from  $l$ 
7 while  $l$  is not empty
8    $n \leftarrow \text{NearestFarSearch}(rNew, c, l)$ 
9   add  $n$  to the end of  $o$ 
10   $rLast \leftarrow rNew$ 
11   $r \leftarrow n$ 
12   $\gamma \leftarrow \text{weight}$ 
13   $rNew \leftarrow \text{Mean}(r * \gamma, rLast * (1-\gamma))$ 
14 return  $o$ 

```

شکل ۲: الگوریتم مرتب‌سازی NFPN++

**Input:** dataset  $T$   
**Output:** ordered sequence  $o$  on records of  $T$

```

1  $l \leftarrow T$ 
2  $o \leftarrow []$ 
3  $c \leftarrow \text{global\_centroid}$ 
4  $rNew \leftarrow$  furthest record from  $c$ 
5 add  $rNew$  to the end of  $o$ 
6 remove  $rNew$  from  $l$ 
7 while  $l$  is not empty
8   if  $\text{size}(o) \geq 5$ 
9      $rNew \leftarrow \text{Mean}(\text{last five records in } o)$ 
10     $n \leftarrow \text{NearestFarSearch}(rNew, c, l)$ 
11    add  $n$  to the end of  $o$ 
12  else
13     $rNew \leftarrow \text{Mean}(\text{all record in } o)$ 
14     $n \leftarrow \text{NearestFarSearch}(rNew, c, l)$ 
15    add  $n$  to the end of  $o$ 
16 return  $o$ 

```

شکل ۳: الگوریتم مرتب‌سازی ENFPN

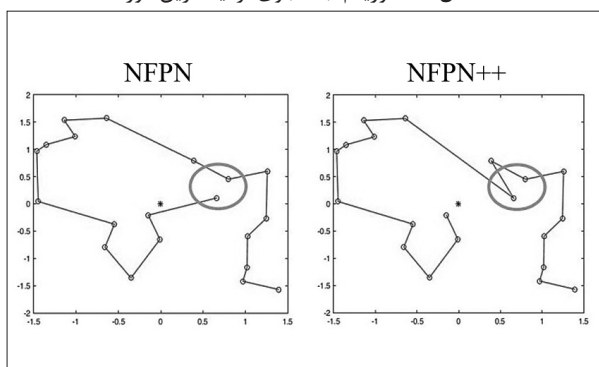
**Input:** record  $r$ , record  $c$ , list  $l$   
**Output:** record  $n$  in  $l$ , furthest record from  $c$  between nears to  $r$

```

1  $mind \leftarrow +\infty$ 
2  $n \leftarrow \text{NULL}$ 
3 For each record  $x$  in  $l$ 
4    $d1 \leftarrow \text{EuclideanDistance}(x, r)$ 
5    $d2 \leftarrow \text{EuclideanDistance}(x, c)$ 
6   if ( $d2 > d1$ )
7      $d \leftarrow d1 / d2$ 
8   else
9      $d \leftarrow d1$ 
10  if ( $d < mind$ )
11     $mind \leftarrow d$ 
12     $n \leftarrow x$ 
13 return  $n$ 

```

شکل ۴: الگوریتم جستجوی نزدیک‌ترین دور



شکل ۵: مقایسه بین مرتب‌سازی‌های NFPN++ و NFPN روی یک نمونه داده. هر محور یک صفت را نمایش می‌دهد.

اهمیت DLD، IL و SI مقادیر ۲‌های گوناگون قرار داد که در بخش نتایج شرح داده خواهد شد.

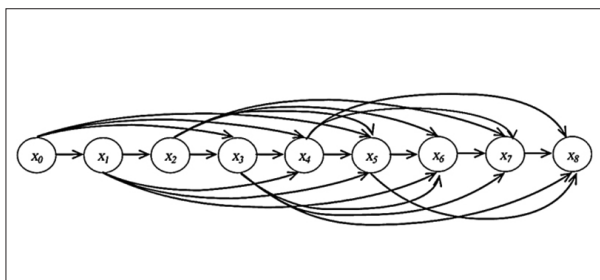
در شکل ۳، خط ۸ الگوریتم، اگر اندازه دنباله ۵ بیشتر از پنج باشد میانگین پنج رکورد انتهای توالی به‌عنوان معیار تصمیم‌گیری برای پارامتر  $d1$  در شکل ۴ می‌باشد. در غیر این صورت میانگین کل رکوردهای دنباله ۵ که کمتر از ۵ می‌باشد در نظر گرفته می‌شود. با انجام آزمایش‌ها بر روی میانگین چند رکورد آخر (به‌عنوان مثال از میانگین ۲ رکورد تا ۱۰ رکورد آخر) بهترین نتیجه برای میانگین ۵ رکورد آخر به‌دست آمد. شکل ۵ مقایسه بین مرتب‌سازی NFPN و NFPN++ و شکل ۶ مقایسه بین مرتب‌سازی NFPN و ENFPN برای مجموعه داده نمونه یکسان را نشان می‌دهد.

#### ۴-۲- خوشه‌بندی داده‌های مرتب شده

دومین قدم از ENFPN و NFPN++ اعمال الگوریتم ریزتجمیع بهینه برای داده‌های تک‌متغیره روی مجموعه داده  $T$  چندمتغیره با رابطه ترتیب تهیه شده از داده‌ها در گام قبل است. الگوریتم MHM [۲۷] مجموعه داده  $T$  با  $n$  رکورد را می‌پذیرد، یک  $k$  صحیح و ترتیبی از رکوردها به‌عنوان ورودی و بخش‌بندی بهینه با نسبت به صفر، رضایت‌مندی  $k$ -بی‌نامی با حداقل مقدار اتلاف اطلاعات را بر می‌گرداند. MHM رکوردهای مرتب شده در یک گراف جهت‌دار بدون دایره  $G$  را در نظر می‌گیرد که هر رکورد به‌عنوان یک گره در نظر گرفته می‌شود.

همان‌طور که در شکل ۷ نشان داده شده، MHM ابتدا رکورد موقت  $x_0$  را قبل از گره  $x_1$  اضافه می‌کند و سپس لبه‌های شروع  $k$  از هر  $x_i, 0 \leq i \leq n-k, x_{i+j}$  به  $x_{i+j-1}, 1 \leq j \leq k$  و  $i+j \leq n$  را اضافه می‌کند. هر لبه  $e_{ij}$  به‌عنوان یک خوشه شامل رکوردها  $\{x_{0[i]}, \dots, x_{0[i-1]}\}$  فرض می‌شود، و وزن لبه با مجموع مربعات خطا از این خوشه معادل است. بعد از اضافه شدن لبه‌ها به  $G$ ، مسیرهای چندگانه از  $x_0$  به  $x_n$  در  $G$  وجود دارند. هر مسیر یک  $k$ -بخش‌بندی معتبر از داده‌ها را نشان می‌دهد و طول آن SSE نهایی از این خوشه‌بندی را تعریف می‌کند. کوتاه‌ترین مسیر از  $x_0$  به  $x_n$  خوشه‌بندی





شکل ۷: گراف ایجاد شده به وسیله MHM برای یک مجموعه داده تک متغیره با ۸ رکورد،  $k=3$

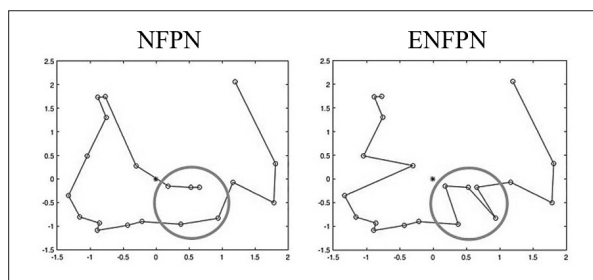
جدول ۱: تعیین مقدار  $\gamma$  برای IL کمینه

	EIA	Census	Tarra gona
$k = 3$	0.1	0.2	0.3
$k = 4$	0.3	0.5	0.6
$k = 5$	0.5	0.5	0.6
$k = 6$	0.5	0.5	0.6
$k = 7$	0.5	0.5	0.6
$k = 8$	0.5	0.5	0.6
$k = 9$	0.5	0.5	0.6
$k = 10$	0.4	0.5	0.6

سپس چون آزمون آماری معنی دار است با استفاده از آزمون‌های تعقیبی آماری (آزمون LSD) به مقایسه روش پیشنهادی با دیگر روش‌ها پرداخته شد. نتایج ضمن تایید برتری روش پیشنهادی نسبت به روش‌های پیشین، به برتری نسبی روش NFPN++ اشاره دارد. همچنین اختلاف به شدت معنی دار روش پیشنهادی با روش NPN مشاهده شد.

مقادیر خاکستری در جداول ۴ و ۵ عدم موفقیت روش پیشنهادی را نشان می‌دهد. اما نسبت به دو روش دیگر که موفق نبوده‌اند، موفق بوده است.

نتایج تایید می‌کند که هیچ یک از IL و یا DLD برای مقایسه به صورت منحصر به فرد کافی نیستند، به این خاطر در کاربردهای متفاوت، اهمیت متفاوتی به هر یک از آن‌ها می‌دهند. SI برای در نظر گرفتن هر دوی IL و DLD با وزن‌های مختلف استفاده می‌شود.  $\alpha$  سطح اهمیت خطر افشا است، بنابراین اهمیت اتلاف اطلاعات به وسیله  $(1-\alpha)$



شکل ۶: مقایسه بین مرتب‌سازی‌های NFPN و ENFPN روی یک نمونه داده. هر محور یک صفت را نمایش می‌دهد.

بهتر با رابطه ترتیب داده شده از داده‌ها را مشخص می‌کند که حداقل اتلاف اطلاعات را به وسیله تقسیم SSE بر SST را متحمل می‌شود.

### ۵- نتایج تجربی

نتایج NFPN++ با توجه به اهمیت ۳ معیار DLD، IL و SI با مقادیر مختلف  $\gamma$  ارائه شده است. تعیین مقدار دقیق  $\gamma$  با توجه به جداول ۱ و ۲ و ۳ می‌باشد. این مقادیر با آزمایش‌های مختلف از ۰٫۱ تا ۰٫۹ به دست آمده و بهترین آن‌ها طبق این جدول‌ها می‌باشد. جداول ۴ و ۵ و شکل‌های ۸ تا ۱۰ نتایج آزمایش‌ها را در مقایسه با سه روش NPN و MDAV نشان می‌دهند. روش پیشنهاد شده مقادیر IL سراسری کمینه را در ۱۸ از ۲۴ آزمایش نشان می‌دهد. همچنین مقادیر DLD کمینه را برای تمام آزمایش‌ها نشان می‌دهد.

برای مقایسه و اثربخشی روش‌ها، استفاده از روش آنالیز واریانس توصیه شده است. بدین منظور با ورود داده‌ها در نرم‌افزار SPSS و با استفاده از جداول آنالیز واریانس به بررسی و مقایسه روش‌ها با یکدیگر و همچنین مقایسه مؤلفه‌های دیگر از قبیل مجموعه داده‌ها و پارامتر  $k$  پرداخته شد.

شکل‌های ۸ و ۹ به ترتیب نمودارهای حاصل از جداول ۴ و ۵ را نشان می‌دهد.

نتایج به دست آمده از آنالیز واریانس، اختلاف معنی دار آماری در استفاده از روش‌ها را نشان داد ( $P < 0.05$ )، مقدار احتمال باید کمتر از ۰٫۰۵ باشد تا اختلاف معنی دار باشد).

جدول ۲: تعیین مقدار  $\gamma$  برای DLD کمینه

	EIA	Census	Tarra gona
k = 3	0.9	0.9	0.9
k = 4	0.9	0.9	0.9
k = 5	0.9	0.9	0.9
k = 6	0.9	0.9	0.9
k = 7	0.9	0.9	0.9
k = 8	0.9	0.9	0.9
k = 9	0.9	0.9	0.9
k = 10	0.9	0.9	0.9

جدول ۳: تعیین مقدار  $\gamma$  برای SI کمینه

EIA	$\alpha = 0.3$	$\gamma = 0.7$
	$\alpha = 0.5$	$\gamma = 0.7$
	$\alpha = 0.7$	$\gamma = 0.7, 0.9$
Census	$\alpha = 0.3$	$\gamma = 0.3, 0.5$
	$\alpha = 0.5$	$\gamma = 0.7$
	$\alpha = 0.7$	$\gamma = 0.8, 0.9$
Tarra gona	$\alpha = 0.3$	$\gamma = 0.6, 0.9$
	$\alpha = 0.5$	$\gamma = 0.9$
	$\alpha = 0.7$	$\gamma = 0.9$

جدول ۴: نتایج مقایسه‌های IL برای الگوریتم NFPN++

	Method	k = 3	k = 4	k = 5	k = 6	k = 7	k = 8	k = 9	k = 10
EIA	MDAV	0.48	0.67	1.67	1.31	2.17	2.87	3.18	3.83
	NPN	0.5	0.69	1	1.31	1.92	2.14	2.36	2.49
	NFPN	<b>0.41</b>	0.61	<b>0.86</b>	<b>1.1</b>	1.83	2.08	2.29	2.45
	<b>NFPN++</b>	<b>0.43</b>	<b>0.6</b>	<b>0.86</b>	<b>1.11</b>	<b>1.74</b>	<b>1.92</b>	<b>2.11</b>	<b>2.18</b>
Census	MDAV	5.69	7.49	9.09	10.38	<b>11.58</b>	<b>12.39</b>	<b>13.29</b>	<b>14.16</b>
	NPN	6.21	8.96	11.05	13.28	15.31	17.36	18.65	20.23
	NFPN	5.63	7.6	9.3	11.03	12.74	14.35	15.43	16.68
	<b>NFPN++</b>	<b>5.47</b>	<b>7.35</b>	<b>8.92</b>	<b>10.28</b>	<b>11.76</b>	<b>12.95</b>	<b>14.13</b>	<b>15.42</b>
Tarragona	MDAV	16.93	19.55	22.46	26.33	27.52	29.69	31.21	33.19
	NPN	17.53	21.69	28.17	30.63	32.21	34.01	35.96	38.72
	NFPN	15.47	18.79	21.81	25.87	28.81	30.28	31.07	33.23
	<b>NFPN++</b>	<b>15.23</b>	<b>18.21</b>	<b>21.55</b>	<b>24.8</b>	<b>27.32</b>	<b>28.65</b>	<b>30.37</b>	<b>32.4</b>

می‌دهند. روش پیشنهاد شده مقادیر IL سراسری کمینه را در ۱۸ از ۲۷ آزمایش نشان می‌دهد. همچنین مقادیر DLD کمینه را در ۱۹ از ۳۳ آزمایش نشان می‌دهد. شکل‌های ۱۱ و ۱۲ به ترتیب نمودارهای حاصل از جداول ۶ و ۷ را نشان می‌دهد.

شکل ۱۳ نتایج مقایسه‌های SI با مقادیر مختلف  $\alpha$  برای الگوریتم ENFPN را نشان می‌دهد.

با توجه به جدول ۶ نتایج نشان می‌دهد روش پیشنهادی برای مجموعه داده EIA موفق نبوده است اما

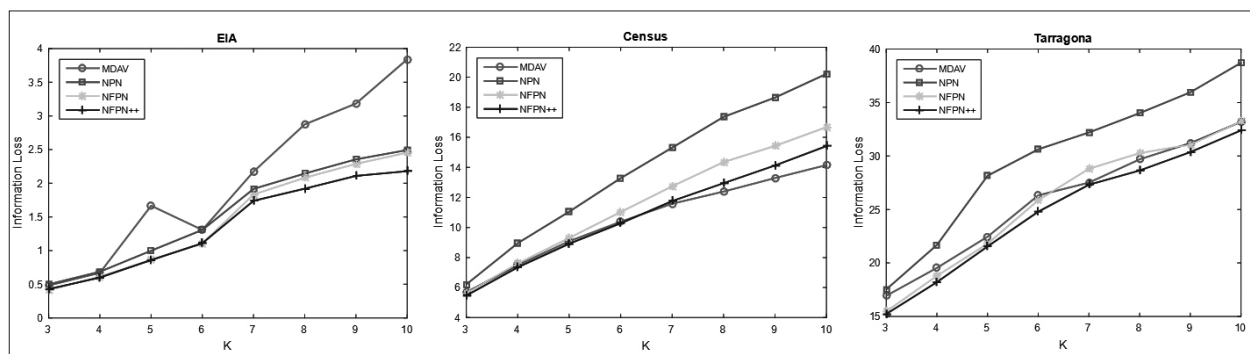
مشخص می‌شود. مقادیر پایین‌تر اتلاف اطلاعات صحت بیشتری را نشان می‌دهد که معمولاً هزینه بالاتر خطر افشا را به دنبال دارد. اگرچه یک کاهش در SI بهبود سراسری در هر دوی IL و DLD را نشان می‌دهد و یک روش ساده برای مقایسه کردن روش‌های بی‌نام‌سازی متفاوت را میسر می‌سازد. شکل ۱۰ نتایج مقایسه‌های SI با مقادیر مختلف  $\alpha$  با توجه به ۲های جدول ۳ برای الگوریتم NFPN++ را نشان می‌دهد.

جداول ۶ و ۷ و شکل‌های ۱۱ تا ۱۳ نتایج آزمایش‌های ENFPN را در مقایسه با دو روش NPN و NFPN نشان

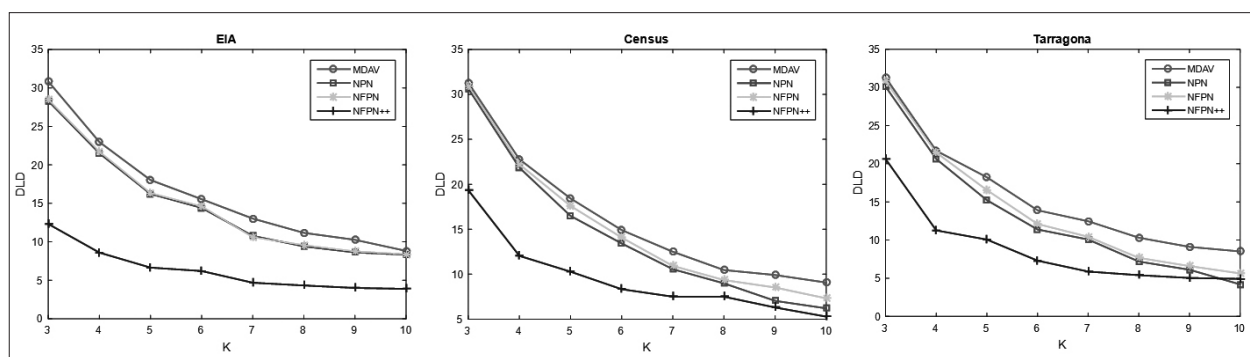


جدول ۵: نتایج مقایسه‌های DLD برای الگوریتم NFPN++

	Method	k = 3	k = 4	k = 5	k = 6	k = 7	k = 8	k = 9	k = 10
EIA	MDAV	31.23	23.39	18.35	15.88	13.39	11.39	10.51	9.07
	NPN	28.67	21.85	16.54	14.71	11.05	9.68	8.9	8.6
	NFPN	28.86	22.07	16.64	14.88	10.90	9.82	9.02	8.65
	<b>NFPN++</b>	<b>12.32</b>	<b>8.60</b>	<b>6.65</b>	<b>6.21</b>	<b>4.69</b>	<b>4.33</b>	<b>4.01</b>	<b>3.89</b>
Census	MDAV	31.3	22.78	18.43	14.91	12.5	10.46	9.91	9.07
	NPN	30.65	21.85	16.48	13.43	10.56	8.98	7.04	6.2
	NFPN	30.93	22.22	17.59	14.07	10.93	9.35	8.52	7.31
	<b>NFPN++</b>	<b>19.35</b>	<b>12.04</b>	<b>10.28</b>	<b>8.33</b>	<b>7.50</b>	<b>7.50</b>	<b>6.30</b>	<b>5.28</b>
Tarragona	MDAV	31.41	21.82	18.47	14.03	12.47	10.31	9.11	8.51
	NPN	30.34	20.86	15.23	11.39	10.19	7.19	6.12	4.2
	NFPN	31.06	21.58	16.67	12.23	10.43	7.79	6.71	5.64
	<b>NFPN++</b>	<b>20.62</b>	<b>11.27</b>	<b>10.07</b>	<b>7.31</b>	<b>5.88</b>	<b>5.40</b>	<b>5.04</b>	<b>4.92</b>



شکل ۸: نمودارهای حاصل از جدول ۴

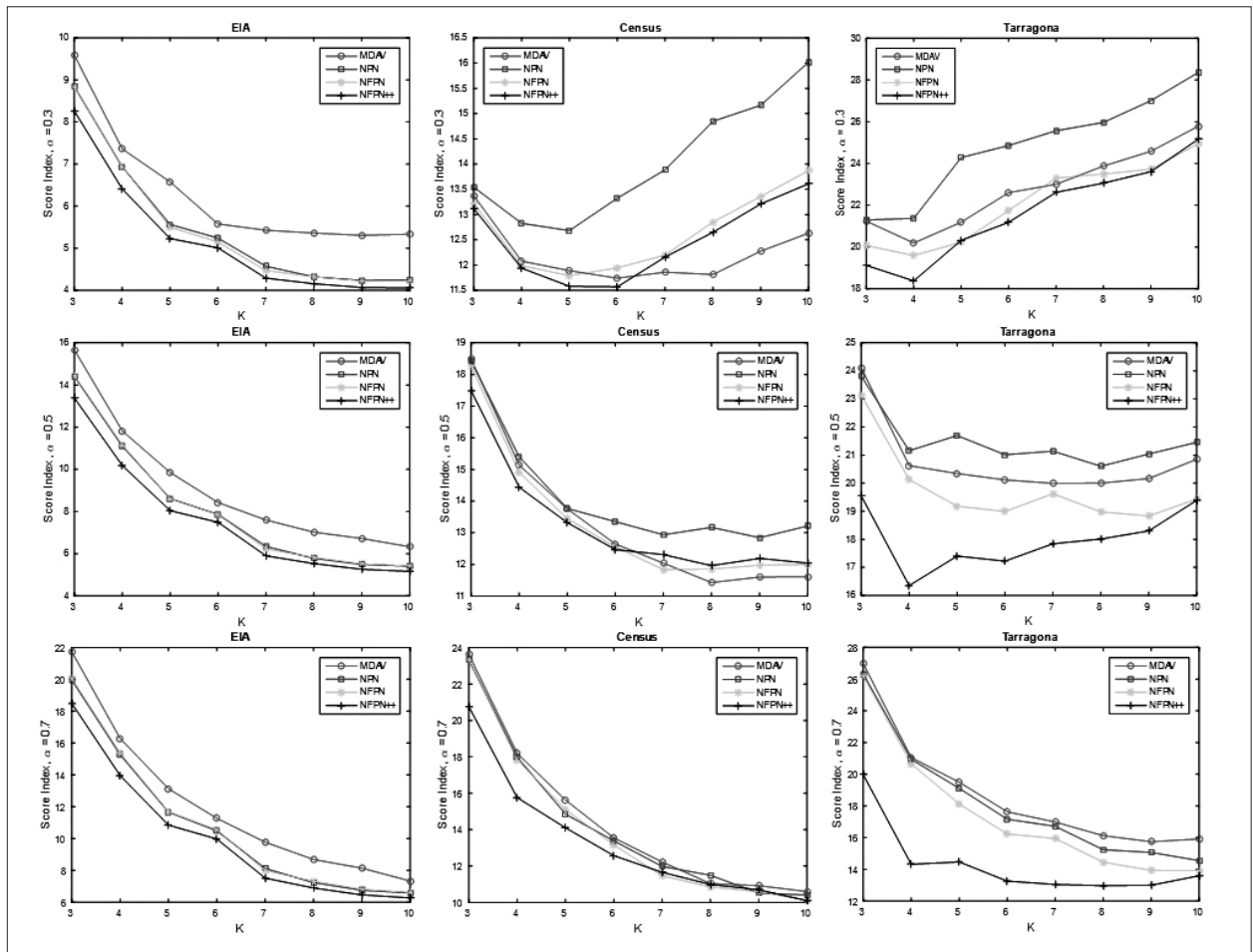


شکل ۹: نمودارهای حاصل از جدول ۵

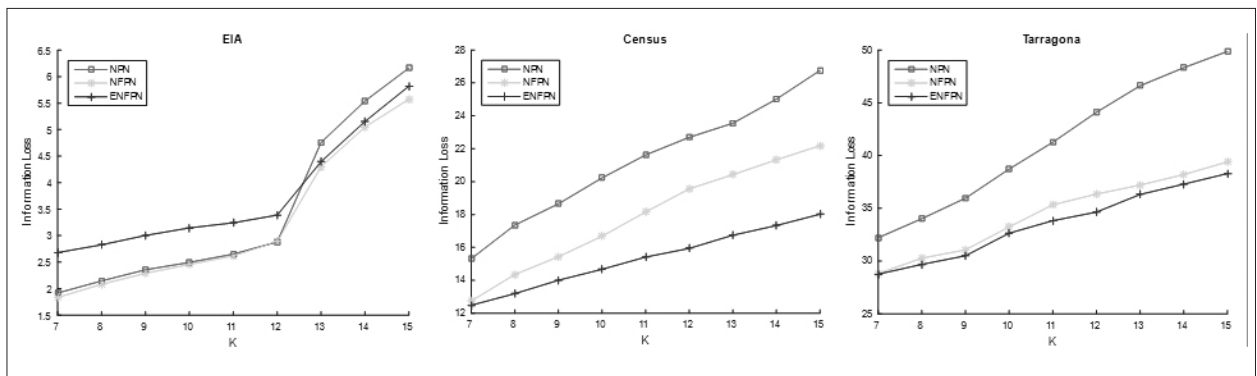
### نتیجه

یک توازن بین حریم خصوصی و سودمندی داده‌ها در انتشار ریزداده‌ها وجود دارد. ریزتجمیع یکی از رویکردهای مورد استفاده در بسیاری از موارد برای حل

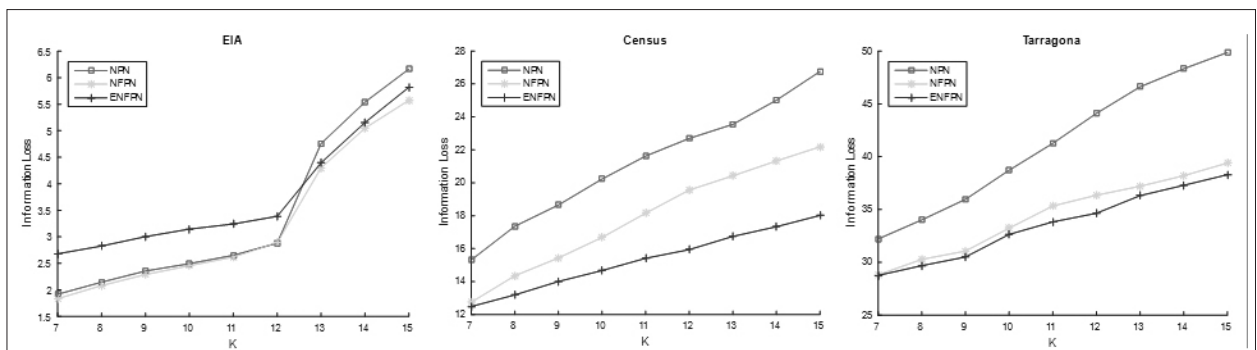
از طرفی طبق جدول ۷ برای تمام آزمایش‌ها موفق عمل کرده است. از طرفی دیگر با آزمایش‌های انجام گرفته مشاهده گردید که معیار  $ll$  برای  $k > 6$  نتایج بهتری را به دنبال دارد.



شکل ۱۰: نتایج مقایسه‌های SI با مقادیر مختلف  $\alpha$  با توجه به  $\gamma$  های جدول ۳ برای الگوریتم NFPN++



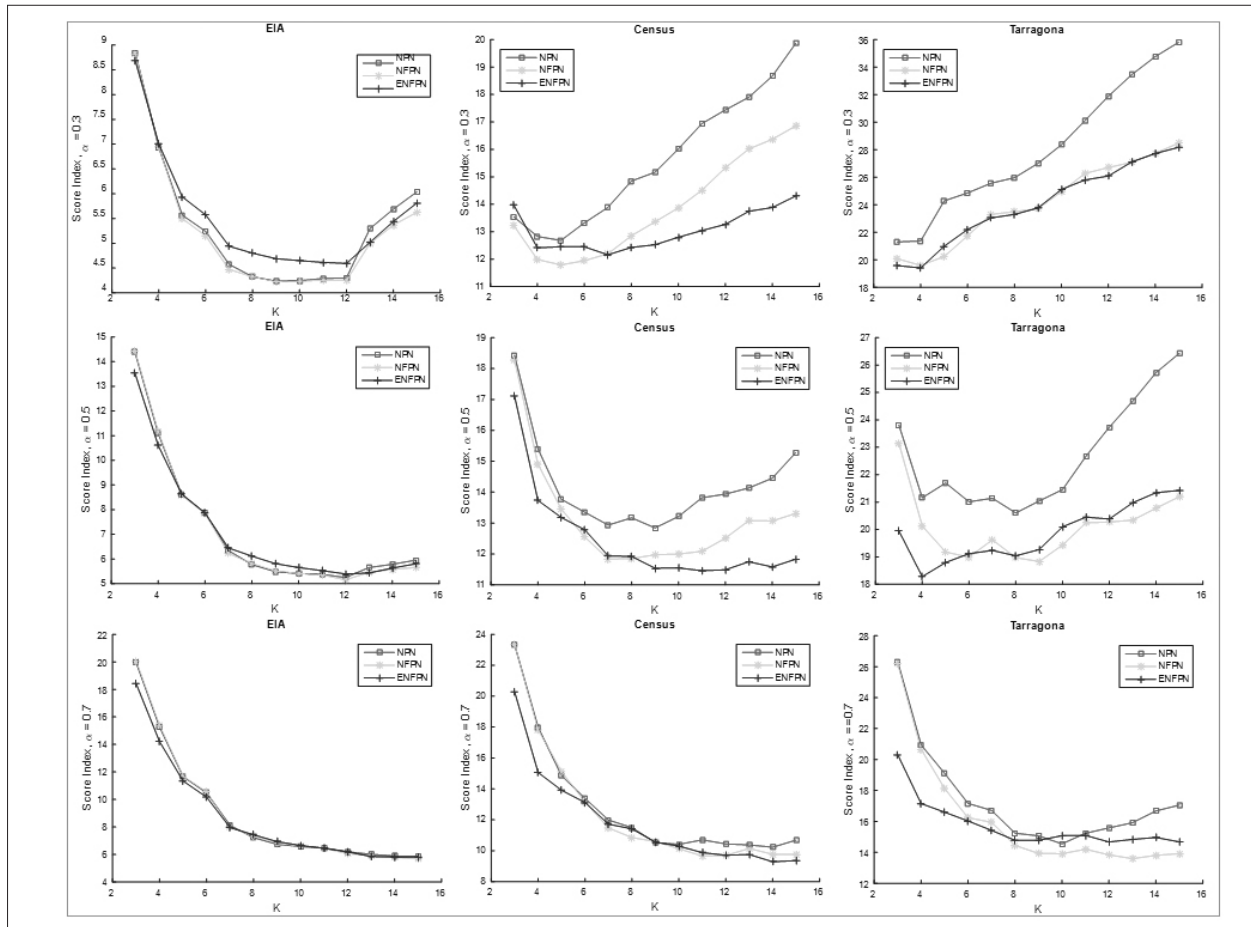
شکل ۱۱: نمودارهای حاصل از جدول ۶



شکل ۱۲: نمودارهای حاصل از جدول ۷

جدول ۶: نتایج مقایسه‌های IL برای الگوریتم ENFPN

	Method	k = 7	k = 8	k = 9	k = 10	k = 11	k = 12	k = 13	k = 14	k = 15
EIA	NPN	1.92	2.14	2.36	2.49	2.65	<b>2.88</b>	4.76	5.54	6.16
	NFPN	<b>1.83</b>	<b>2.08</b>	<b>2.29</b>	<b>2.45</b>	<b>2.62</b>	2.89	<b>4.30</b>	<b>5.05</b>	<b>5.57</b>
	ENFPN	2.68	2.83	3.00	3.14	3.24	3.39	4.40	5.16	5.82
Census	NPN	15.31	17.36	18.65	20.23	21.62	22.69	23.55	25.01	26.76
	NFPN	12.74	14.35	15.43	16.68	18.15	19.56	20.42	21.32	22.18
	ENFPN	<b>12.49</b>	<b>13.19</b>	<b>14.01</b>	<b>14.67</b>	<b>15.42</b>	<b>15.94</b>	<b>16.74</b>	<b>17.33</b>	<b>18.02</b>
Tarragona	NPN	32.21	34.01	35.96	38.72	41.27	44.11	46.62	48.33	49.87
	NFPN	28.81	30.28	31.07	33.23	35.33	36.35	37.21	38.20	39.40
	ENFPN	<b>28.76</b>	<b>29.69</b>	<b>30.51</b>	<b>32.64</b>	<b>33.83</b>	<b>34.65</b>	<b>36.32</b>	<b>37.29</b>	<b>38.30</b>



شکل ۱۳: نتایج مقایسه‌های SI با مقادیر مختلف  $\alpha$  برای الگوریتم ENFPN

جدول ۷: نتایج مقایسه‌های DLD برای الگوریتم ENFPN

	Method	$k = 3$	$k = 4$	$k = 5$	$k = 6$	$k = 7$	$k = 8$	$k = 9$	$k = 10$	$k = 11$	$k = 12$	$k = 13$
EIA	NPN	28.67	21.85	16.54	14.71	11.05	9.68	8.9	8.6	8.11	7.62	6.55
	NFPN	28.86	22.07	16.64	14.88	10.90	9.82	9.02	8.65	8.09	7.43	6.62
	<b>ENFPN</b>	<b>25.73</b>	<b>19.60</b>	<b>15.40</b>	<b>13.61</b>	<b>10.22</b>	<b>9.41</b>	<b>8.63</b>	<b>8.16</b>	<b>7.82</b>	<b>7.40</b>	<b>6.48</b>
Census	NPN	30.65	21.85	16.48	<b>13.43</b>	<b>10.56</b>	<b>8.98</b>	<b>7.04</b>	<b>6.2</b>	<b>6.02</b>	<b>5.19</b>	<b>4.72</b>
	NFPN	30.93	22.22	17.59	14.07	10.93	9.35	8.52	7.31	6.02	5.46	5.74
	<b>ENFPN</b>	<b>25.00</b>	<b>17.04</b>	<b>15.00</b>	<b>13.61</b>	11.39	10.65	9.07	8.43	7.50	7.04	6.76
Tarragona	NPN	30.34	20.86	15.23	11.39	10.19	<b>7.19</b>	<b>6.12</b>	<b>4.2</b>	<b>4.08</b>	<b>3.36</b>	<b>2.76</b>
	NFPN	31.06	21.58	16.67	12.23	10.43	7.79	6.71	5.64	5.16	4.20	3.48
	<b>ENFPN</b>	<b>20.86</b>	<b>15.47</b>	<b>13.31</b>	<b>11.39</b>	<b>9.71</b>	8.39	8.03	7.55	7.07	6.12	5.64

IDF to hide sensitive itemsets. Applied Intelligence, 2013. 38(4): p. 502-510.

[6] Yin, Y., Kaku, I., Tang, J., & Zhu, J., Privacy-preserving data mining. Data mining. Springer, Berlin, 2011.

[7] Domingo-Ferrer, J., & Torra, V., Disclosure control methods and information loss for microdata. Confidentiality, disclosure, and data access: theory and practical applications for statistical agencies, 2001: p. 91-110.

[8] Oganian, A., & Karr, A. F., Masking methods that preserve positivity constraints in microdata. Journal of Statistical Planning and Inference, 2011. 141(1): p. 31-41.

[9] Domingo-Ferrer, J., Solanas, A., & Martinez-Balleste, A., Privacy in Statistical Databases: k-Anonymity Through Microaggregation. in Granular Computing (GrC). 2006.

[10] Xu, C., Wang, Y., Gu, Y., Lin, S., & Yu, G., Efficient fuzzy ranking queries in uncertain databases. Applied Intelligence, 2012. 37(1): p. 47-59.

[11] Kiran, P., & Kavya, N., A Survey on Methods, Attacks and Metric for Privacy Preserving Data Publishing. International Journal of Computer Applications, 2012. 53(18): p. 20-28

[12] De Capitani Di Vimercati, S., Foresti, S., Livraga, G., & Samarati, P., Data privacy: definitions and techniques. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, 2012. 20(06): p. 793-817.

[13] Domingo-Ferrer, J., & Sramka, M., Disclosure Control by Computer Scientists: An Overview and an Application of Microaggregation to Mobility Data Anonymization. 2011: p. 1082-1086.

[14] Rebollo-Monedero, D., Forné, J., & Soriano, M., An algorithm for k-anonymous microaggregation and clustering inspired by the design of distortion-optimized quantizers. Data & Knowledge Engineering, 2011. 70(10): p. 892-921.

[15] Solanas, A., Martinez-Balleste, A., & Domingo-Ferrer, J., V-MDAV: a multivariate microaggregation with variable group size. in 17th COMPSTAT Symposium of the IASC, Rome. 2006.

[16] Wong, R. C.-W., Li, J., Fu, A. W.-C., & Wang, K., ( $\alpha$ , k)-anonymity: an enhanced k-anonymity model for privacy preserving data publishing. in Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data

این مسئله است. روش‌های ریزتجمیع گوناگون در مقالات پیشنهاد شده است. این مقاله روش ریزتجمیع NFPN++ و ENFPN، را به‌عنوان توسعه روش پیش‌تر ارائه شده NFPN-MHM با یک صحت بهبود یافته پیشنهاد کرد. این روش‌ها معمولاً توازن بهتر بین اتلاف اطلاعات و خطر افشا نسبت به روش‌های دیگر در سایر مقالات را نشان می‌دهند. همچنین انعطاف‌پذیری لازم را برای اهمیت دادن به هر یک از اتلاف اطلاعات، کنترل افشا و شاخص امتیاز را دارند.

### سپاس‌گزاری

در اینجا لازم هست از آقای دکتر مصطفی سالاری برای به اشتراک‌گذاری کد NFPN و کمک‌های ایشان سپاس‌گزاری نمایم.

### مراجع

- [1] Mortazavi, R., Jalili, S., & Gohargazi, H., Multivariate microaggregation by iterative optimization. Applied intelligence, 2013. 39(3): p. 529-544.
- [2] Schoeman, F.D., Philosophical dimensions of privacy: An anthology. 1984: Cambridge University Press.
- [3] Salari, M., Jalili, S., & Mortazavi, R., A utility preserving data-oriented anonymization method based on data ordering. in Telecommunications (IST), 2014 7th International Symposium on. 2014. IEEE.
- [4] Torra, V., Navarro-Arribas, G., & Stokes, K., An Overview of the Use of Clustering for Data Privacy, in Unsupervised Learning Algorithms. 2016, Springer. p. 237-251.
- [5] Hong, T.-P., Lin, C.-W., Yang, K.-T., & Wang, S.-L., Using TF-

Engineering, IEEE Transactions on, 2005. 17(7): p. 902-911.

[23] Solé, M., Muntés-Mulero, V., & Nin, J., Efficient microaggregation techniques for large numerical data volumes. International Journal of Information Security, 2012. 11(4): p. 253-267.

[24] Domingo-Ferrer, J., & Torra, V., Ordinal, continuous and heterogeneous k-anonymity through microaggregation. Data Mining and Knowledge Discovery, 2005. 11(2): p. 195-212.

[25] Panagiotakis, C., & Tziritas, G., Successive group selection for microaggregation. Knowledge and Data Engineering, IEEE Transactions on, 2013. 25(5): p. 1191-1195

[26] Hansen, S. L., & Mukherjee, S., A polynomial algorithm for optimal univariate microaggregation. IEEE Transactions on Knowledge & Data Engineering, 2003(4): p. 1043-1044.

[27] Ghinita, G., Karras, P., Kalnis, P., & Mamoulis, N., Fast data anonymization with low information loss. in Proceedings of the 33rd international conference on Very large data bases. 2007. VLDB Endowment.

[28] Domingo-Ferrer, J., Martínez-Ballesté, A., Mateo-Sanz, J. M., & Sebé, F., Efficient multivariate data-oriented microaggregation. The VLDB Journal, 2006. 15(4): p. 355-369.

mining. 2006. ACM.

[17] Sun, X., Li, M., & Wang, H., A family of enhanced-diversity models for privacy preserving data publishing. Future Gener ComputSyst, 2011. 27(3): p.348-356.

[18] Sun, X., Sun, L., & Wang, H., Extended k-anonymity models against sensitive attribute disclosure. Computer Communications, 2011. 34(4): p. 526-535.

[19] Oganian, A., & Domingo-Ferrer, J., On the complexity of optimal microaggregation for statistical disclosure control. Statistical Journal of the United Nations Economic Commission for Europe, 2001. 18(4): p. 345-353.

[20] Navarro-Arribas, G., & Torra, V., Information fusion in data privacy: A survey. Information Fusion, 2012. 13(4): p. 235-244.

Fung, B., Wang, K., Chen, R., & Yu, P. S., Privacy-preserving data publishing: A survey of recent developments. ACM Computing Surveys (CSUR), 2010. 42(4): p. 14.

[21] Domingo-Ferrer, J., & Mateo-Sanz, J. M., Practical data-oriented microaggregation for statistical disclosure control. Knowledge and Data Engineering, IEEE Transactions on, 2002. 14(1): p. 189-201.

[22] Laszlo, M., & Mukherjee, S., Minimum spanning tree partitioning algorithm for microaggregation. Knowledge and Data

## چاپ دوم منتشر شد !

**برای کسب اطلاعات بیشتر و تهیه کتاب  
با شماره تلفن زیر تماس حاصل فرمایید  
۶۶۴۱۲۸۶۱ (انجمن انفورماتیک ایران)**

جیسون فرید و دیوید هاین مایرهنسون، بنیانگذاران شرکت نرم افزاری 37signals هستند. محصولات تولید شده توسط شرکت آن‌ها میلیون‌ها کاربر در سراسر جهان دارد. آن‌ها در این کتاب راه‌های موفقیت شرکتشان را با شما در میان می‌گذارند. این کتاب در فهرست پرفروش‌ترین کتاب‌های روزنامه نیورکتایمز قرار داشته است.

