

## خوشه‌بندی داده‌های جریان داده دارای برچسب

زهرا بیات

دانشجوی کارشناسی ارشد دانشکده ریاضی، آمار و علوم کامپیوتر، پردیس علوم، دانشگاه تهران  
پست الکترونیکی: bayat.zahra@ut.ac.ir

هدیه ساجدی\*

استادیار دانشکده ریاضی، آمار و علوم کامپیوتر، پردیس علوم، دانشگاه تهران  
پست الکترونیکی: hhsajedi@ut.ac.ir

### چکیده

واژه‌های کلیدی: توده جریان داده<sup>۲</sup>، جریان داده<sup>۴</sup>، خوشه‌بندی<sup>۵</sup>، خوشه‌بندی نظارت شده<sup>۶</sup>.

امروزه مسئله پیداکردن الگوریتم‌های خوشه‌بندی نظارت‌شده برای داده‌های پویا و جریان داده اهمیت زیادی دارد. محققان سعی می‌کنند برای حل کردن این مسئله الگوریتم‌های جدید ارائه نمایند و یا الگوریتم‌های موجود را بهبود دهند. در میان این الگوریتم‌ها، روش SAIC<sup>۱</sup> برای خوشه‌بندی داده‌های پویا با خوشه‌های با اندازه و شکل دلخواه ارائه شده‌است. در این روش، تعداد خوشه‌ها به طور خودکار توسط الگوریتم مشخص می‌شود اما این الگوریتم قادر به تشخیص صحیح خوشه‌های مسائل تک دسته‌ای نیست. این امر موجب می‌شود که بعضاً در خوشه‌بندی جریان داده اختلال ایجاد شود. در این مقاله ضمن توضیح علت ایجاد این مشکل، الگوریتم ISAIC<sup>۲</sup> برای بهبود الگوریتم SAIC پیشنهاد می‌شود. همچنین عملکرد الگوریتم ISAIC با الگوریتم SAIC روی چند مجموعه داده مورد مقایسه قرار گرفته و نتایج آرایه شده‌است. میزان بهبود دقت دسته‌بندی بر روی مجموعه داده‌های مورد آزمایش حداقل صفر و حداکثر حدود ۶۵٪ است.

### ۱-مقدمه

واژه جریان داده اشاره بر توالی سریع اطلاعاتی که بالقوه حجیم و پیوسته هستند دارد [۱۲]. کاربردهای گوناگونی وجود دارند که این‌گونه داده‌ها را تولید می‌کنند که از جمله آن‌ها می‌توان به سیستم‌های نظارت بر شبکه، سیستم‌های تحلیل فروش، سیستم‌های ارتباطی و جریان کلیک کاربران اشاره کرد [۵].

تفاوت جریان داده با داده ایستا را می‌توان در موارد زیر برشمرد [۶]:

- داده‌ها در جریان داده برخلاف داده‌های ایستا به صورت برخط و به مرور دریافت می‌شوند.
- هیچ کنترلی بر روی ترتیب ورود عناصر جریان داده وجود ندارد.
- در حالت کلی اندازه یک جریان داده نامحدود است.

3- Data stream of chunks

4- Data stream

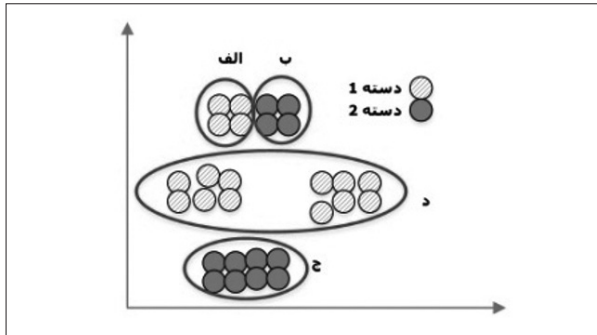
5- Clustering

6- Supervised clustering

1- Supervised Adaptive Incremental Clustering

2- Improved Supervised Adaptive Incremental Clustering

\* نویسنده مسئول



شکل ۱: نمونه‌ای از خوشه‌بندی نظارت شده.

شده‌اند در فاصله زیادی از این توده‌های جریان داده باشند. در این حالت بدون توجه به شکل داده‌های این توده، تنها یک خوشه در مرکز توده جدید توسط الگوریتم SAIC معرفی می‌شود.

در این مقاله یک نسخه بهبود یافته از الگوریتم SAIC برای حل مسئله خوشه‌بندی ارائه می‌شود که نواقص الگوریتم SAIC را برطرف می‌نماید. الگوریتم SAIC یک الگوریتم تصادفی است که بر مبنای مقایسه خوشه‌ها عمل می‌کند و این امر بخصوص وقتی با توده جریان داده تک دسته‌ای روبرو می‌شویم ممکن است مشکل ایجاد کند. لازم به ذکر است پیچیدگی محاسباتی و زمان اجرای الگوریتم بهبود یافته ISAIC تفاوت قابل ملاحظه‌ای نسبت به الگوریتم SAIC ندارد.

ساختار مقاله در ادامه شرح داده می‌شود. در بخش ۲ کارهای مرتبط مطرح می‌شود. در بخش ۳ الگوریتم بهبود یافته ISAIC معرفی می‌شود. بخش ۴ آزمایش‌ها را شرح می‌دهد و بخش ۵ شامل جمع‌بندی مقاله است.

## ۲- کارهای مرتبط

در این بخش ابتدا تعدادی از الگوریتم‌های خوشه‌بندی مطرح برای مجموعه داده‌های ایستا به صورت خلاصه شرح داده می‌شوند و سپس به خوشه‌بندی جریان داده پرداخته می‌شود. الگوریتم‌های خوشه‌بندی نظارت شده، همان‌طور که شکل ۱ نشان می‌دهد، یک مجموعه داده را به چندین خوشه طوری تقسیم می‌نمایند که داده‌های عضو یک خوشه دارای برچسب یکسان باشند.

با توجه به خصوصیات جریان داده روش‌های خوشه‌بندی باید به صورت افزایشی باشند زیرا تمامی داده‌ها به یکباره در دسترس نیستند [۷].

فرآیند خوشه‌بندی نظارت شده سعی دارد که یک مجموعه داده را به چندین خوشه تقسیم نماید به طوری که داده‌های قرار گرفته در یک خوشه دارای یک برچسب باشند و تعداد خوشه‌ها کمینه باشد. در حال حاضر روش‌های متعددی برای خوشه‌بندی داده‌ها وجود دارد که بر اساس نوع داده‌ها، شکل خوشه‌ها، فاصله داده‌ها و غیره عمل خوشه‌بندی را انجام می‌دهند [۱۳، ۱۴].

در میان الگوریتم‌های خوشه‌بندی جریان داده روش SAIC، یک الگوریتم تقریبی است که در سال ۲۰۱۷ توسط Zheng و دیگران ارائه شده است. الگوریتم SAIC به خوشه‌بندی توده‌های جریان داده می‌پردازد. در این الگوریتم فرض می‌شود که در هر لحظه یک توده از داده‌ها وارد می‌شود. تفاوت توده‌های جریان داده با جریان داده این است که وقتی با توده‌های جریان داده کار می‌کنیم در واقع با ترکیبی از جریان داده و داده‌های ایستا روبرو هستیم و در واقع با هر توده می‌توان مشابه داده‌های ایستا برخورد کرد. بخصوص وقتی با داده‌های حجیم روبرو هستیم گاهی برای پردازش این داده‌ها آن‌ها را به توده‌های جریان داده تبدیل می‌کنیم. الگوریتم SAIC مبتنی بر یک معیار فاصله است و هر یک از فاصله‌های همیلتونی، اقلیدسی و یا فاصله ماهالانوبیس قابل استفاده هستند. این الگوریتم تصادفی بر مبنای مقایسه فاصله بین خوشه‌ها است و نه تنها فاصله یک داده از خوشه‌های با برچسب مشابه در تعیین موقعیت آن داده موثر است بلکه فاصله داده با خوشه‌های با برچسب مخالف نیز تاثیرگذار است.

عدم کارایی این الگوریتم زمانی است که تمام داده‌های یک توده از جریان داده یک برچسب داشته باشند و در واقع به یک دسته تعلق داشته باشند و مجموعه خوشه‌هایی که تا آن زمان در این جریان داده ایجاد

## – الگوریتم‌های چگالی محور

ویژگی‌های نزدیک به یکدیگر با هم ادغام شده و این کار ادامه می‌یابد تا چند خوشه مجزا حاصل شود. مشکل این روش‌ها حساس بودن به نوفه و پیچیدگی فضایی بالا است. در روش‌های خوشه‌بندی بالا به پایین<sup>۹</sup> ابتدا تمام داده‌ها به‌عنوان یک خوشه در نظر گرفته شده و با به‌کارگیری یک الگوریتم تکرار شونده هر بار این خوشه با در نظر گرفتن یک معیار شباهت داده‌ها، به خوشه‌های مجزا تقسیم می‌شوند. این کار ادامه می‌یابد تا یک یا چند خوشه یک عضوی ایجاد شود [۲].

## – الگوریتم‌های خوشه‌بندی جریان داده

الگوریتم‌های خوشه‌بندی جریان داده توسط [۴] به دو دسته یک فازی و دو فازی تقسیم می‌شوند. الگوریتم‌های خوشه‌بندی یک فازی، خوشه‌بندی جریان داده را یک نسخه پیوسته از خوشه‌بندی داده‌های ایستا در نظر می‌گیرند. محدودیتی که این روش‌ها دارند این است که وزن‌های یکسانی را برای داده‌های قدیمی و جدید در نظر می‌گیرند. الگوریتم STREAM<sup>۱۰</sup> و الگوریتم STREAM-E<sup>۱۱</sup> [۱۰] نمونه‌هایی از الگوریتم‌های یک فازی هستند. در ادامه الگوریتم STREAM به‌طور خلاصه شرح داده می‌شود. در الگوریتم STREAM خوشه‌ها بر اساس تمام داده‌هایی که تا زمان جاری دریافت شده‌اند، محاسبه می‌شوند. به عبارت دیگر در این الگوریتم با استفاده از روش تقسیم و حل<sup>۱۲</sup> جریان‌های داده به جزءهایی تقسیم و خوشه‌های موجود در آن بر اساس الگوریتم k-میانگین محاسبه می‌شوند [۸].

الگوریتم‌های دوفازی از فاز برخط و فاز برون خط تشکیل می‌شوند. در فاز برخط که به صورت پیوسته در حال انجام است، اطلاعات آماری از داده‌ها گردآوری می‌شود. در واقع در این مرحله نوعی پیش‌پردازش بر روی داده‌ها انجام می‌شود. در فاز برون خط با توجه به این اطلاعات پیش‌پردازش شده، نتیجه نهایی خوشه‌بندی به دست

در این الگوریتم‌ها فرض می‌شود که خوشه‌ها مناطقی با چگالی بیشتر هستند که توسط مناطق با چگالی کمتر از هم جدا شده‌اند. یکی از مطرح‌ترین الگوریتم‌ها در این زمینه الگوریتم DBSCAN<sup>۷</sup> است [۲]. روش این الگوریتم به این صورت است که هر داده متعلق به یک خوشه، در دسترس چگالی برای سایر داده‌های همان خوشه است، ولی در دسترس چگالی سایر داده‌های خوشه‌های دیگر نیست (چگالی داده در همسایگی به مرکز داده و شعاع همسایگی دلخواه  $\epsilon$  در نظر گرفته می‌شود). مزیت این روش این است که تعداد خوشه‌ها به صورت خودکار مشخص می‌شود. از نقاط قوت این الگوریتم، تشخیص نوفه و داده‌های پرت است. مشکل این روش پیچیدگی زمانی بالای آن است.

## – الگوریتم k-میانگین

این روش در عین سادگی یک روش بسیار کاربردی و پایه چند روش دیگر مثل خوشه‌بندی فازی می‌باشد. در این الگوریتم ابتدا به تعداد دلخواه نقطه به‌عنوان مرکز خوشه در نظر گرفته می‌شود. سپس با بررسی هر داده، آن داده به نزدیک‌ترین مرکز خوشه نسبت داده می‌شود. پس از اتمام این کار با گرفتن میانگین در هر خوشه، مراکز خوشه‌ها و به دنبال آن خوشه‌ها اصلاح می‌شود. از جمله مشکلات این روش این است که بهینگی آن وابسته به انتخاب اولیه مراکز است. مشکل دیگر الگوریتم k-میانگین این است که تعداد خوشه‌ها به‌عنوان ورودی توسط کاربر باید تنظیم شود [۲].

## – خوشه‌بندی سلسله مراتبی

خوشه‌بندی سلسله مراتبی خود شامل دو نوع روش خوشه‌بندی است:

در روش‌های خوشه‌بندی پایین به بالا<sup>۸</sup> ابتدا هر داده به‌عنوان یک خوشه در نظر گرفته می‌شود. سپس در ادامه با به‌کارگیری یک الگوریتم هر بار خوشه‌های دارای

9- Top-Down

10- Stream clustering

11- Evolution-based technique for stream clustering

12- Divide and conquer

7- Density-based spatial clustering of applications with noise

8- Bottom-Up

می‌آید [۹]. الگوریتم CluStream<sup>۱۳</sup>، الگوریتم HPStream<sup>۱۴</sup>، الگوریتم DenStream<sup>۱۵</sup> و الگوریتم D-Stream<sup>۱۶</sup> نمونه‌هایی از الگوریتم‌های دو فاز می‌باشند.

### – الگوریتم SAIC

الگوریتم SAIC، الگوریتمی یک فاز است که برای خوشه‌بندی توده‌های جریان داده به کار می‌رود. الگوریتم SAIC داده‌ها را در هر توده به طور جداگانه خوشه‌بندی می‌کند ولی در زمان خوشه‌بندی هر توده، خوشه‌های ایجاد شده در توده‌های قبلی را نادیده نمی‌گیرد، گاهی این خوشه‌ها بروز می‌شوند و گاهی خوشه جدید ایجاد می‌شود. اطلاعات هر خوشه علاوه بر برچسب آن خوشه، شامل تعداد اعضای آن خوشه نیز هست. الگوریتم SAIC خوشه‌بندی را بر اساس مقایسه برچسب داده‌ها و فاصله بین خوشه‌ها انجام می‌دهد.

تعداد اعضای خوشه‌ها دو کاربرد دارد. اولین کاربرد آن، حذف نقاط دورافتاده و نوفه است. اگر در پایان خوشه‌بندی یک توده از جریان داده، خوشه‌ای تعداد اعضای کمی داشته باشد (از حد آستانه کمتر باشد) آن خوشه حذف می‌شود.

دومین کاربرد آن، بروزسانی مرکز خوشه‌ها است. در زمان بروزسانی یک خوشه، وقتی داده‌ای وارد یک خوشه می‌شود، مرکز خوشه تغییر می‌کند و از طرفی در این الگوریتم تنها مرکز خوشه نگهداری می‌شود.

در الگوریتم SAIC داده‌ها یکی یکی بررسی می‌شوند و در نهایت یا یک خوشه جدید ایجاد می‌شود یا مرکز یک خوشه بروزسانی می‌شود. اگر برچسب داده جدید باشد یا اگر کمترین فاصله این داده با خوشه‌های با برچسب مشابه بیشتر از کمترین فاصله این داده با خوشه‌های با برچسب مخالف خود باشد، خوشه جدید ایجاد می‌شود.

دلیل این‌که این الگوریتم قادر به تشخیص خوشه‌ها با اشکال خاص می‌باشد به نوع بررسی آن برمی‌گردد که

13- a method of clustering data streams, based on the concept of microclusters

14- High dimensional projected data stream clustering

15- Density based stream clustering

16- clustering stream data using density

هم فاصله با خوشه‌های با برچسب مشابه و هم فاصله با خوشه‌های با برچسب مخالف را در نظر می‌گیرد.

به علت آن‌که داده‌ها تک تک وارد و بررسی می‌شوند و مبنای این الگوریتم مقایسه است، ترتیب نقاط ورودی در نتیجه خوشه‌بندی تاثیرگذار است. به منظور کاهش این تاثیر، الگوریتم داده‌های توده را به هم ریخته و کار خوشه‌بندی را آنقدر تکرار می‌کند که در دو مرحله متوالی تعداد خوشه‌ها تغییر نکند [۸].

### ۳- روش پیشنهادی

مشکل الگوریتم SAIC زمانی است که در یک توده از جریان داده، داده‌های با برچسب یکسان به هم نزدیک و از داده‌های با برچسب مخالف دور باشند و مجموعه خوشه‌هایی که تا آن زمان در این جریان داده ایجاد شده‌اند در فاصله زیادی از این داده‌ها قرار گرفته باشند. در واقع این الگوریتم بر مبنای مقایسه فاصله هر داده ورودی با خوشه‌های موجود کار می‌کند. خوشه‌بندی داده‌های درون یک توده مشابه روش k-میانگین است با این تفاوت که به تعیین تعداد خوشه‌های اولیه نیاز ندارد؛ چون داده‌ها دارای برچسب می‌باشند الگوریتم بر اساس مقایسه داده جاری با خوشه‌های قبلی، خوشه جدید را در صورت لزوم تولید می‌کند. در روش SAIC فرض شده است که آنقدر تنوع برچسب (و در نتیجه دسته) در هر توده داده موجود است که بتوانیم خوشه‌بندی را بدون نیاز به فرض خوشه‌های اولیه ایجاد کنیم. اما ممکن است در عمل یک توده داده شامل داده‌هایی با تنوع کافی برچسب نباشد یا پراکندگی داده‌ها به گونه‌ای باشد که تنوع برچسب تاثیری در خوشه‌بندی توده داده نداشته باشد. در این شرایط تعداد خوشه‌های حاصل در هر توده برابر با تعداد برچسب‌های موجود در آن توده خواهد بود و خوشه‌بندی قادر به تشخیص شکل داده نیست.

در این مقاله جهت حل این مشکل در الگوریتم پیشنهادی ISAIC برای خوشه‌بندی هر توده، تعدادی خوشه اولیه

در نظر می‌گیریم. به این منظور پارامتر  $k$  در نظر گرفته می‌شود که تعداد خوشه‌های اولیه را مشخص می‌کند. الگوریتم پیشنهادی، هر توده را چندین بار خوشه‌بندی می‌کند تا این‌که از بی‌اثر بودن ترتیب داده‌ها اطمینان حاصل شود. بنابراین با انتخاب  $k$  داده اول به عنوان مراکز خوشه‌های اولیه، با تنوع خوبی از خوشه‌های اولیه روبرو می‌شویم. از طرفی این الگوریتم زمان اجرا را افزایش نمی‌دهد. در الگوریتم SAIC تا وقتی که در دو دور متوالی خوشه‌بندی، تغییر ایجاد نشود، الگوریتم ادامه می‌یابد و خوشه‌بندی مجدد انجام می‌شود. در هر دور این الگوریتم تعداد خوشه‌ها افزایش می‌یابد.

در الگوریتم ISAIC با در نظر گرفتن یک تعداد خوشه اولیه، تعداد دورهای اجرای الگوریتم کاهش می‌یابد و الگوریتم سریعتر همگرا می‌شود.

انتخاب مقدار متغیر  $k$  به عوامل ذیل وابسته است:

- تعداد برجسب‌های مشاهده شده در یک توده: به ازای هر برجسب حداقل یک خوشه تشکیل می‌شود. هر چه تعداد برجسب‌های درون یک توده بیشتر باشد، حداقل تعداد خوشه‌هایی که آن توده می‌تواند داشته باشد افزایش می‌یابد. بنابراین اگر تنوع برجسب کمتر باشد الگوریتم SAIC تعداد خوشه‌های کمتری را شناسایی می‌کند و ممکن است در تشخیص شکل خوشه و در نهایت خوشه‌بندی با مشکل مواجه شود.

- شماره توده: چون با ورود هر توده تعدادی خوشه تولید می‌شود که در روند خوشه‌بندی توده بعدی استفاده می‌شود، انتظار داریم با افزایش شماره توده، تعداد خوشه‌ها روندی افزایش داشته باشد و در نتیجه به مقدار کمتری برای متغیر  $k$  نیاز باشد.

- پراکندگی داده‌های درون یک توده: با افزایش پراکندگی داده‌های درون یک توده نیاز به  $k$  با مقدار بیشتر است. به عنوان مثال برای خوشه‌بندی توده‌ای به شکل دایره به شعاع ۱۰ در مقایسه با توده‌ای به شکل دایره و شعاع ۱ به تعداد خوشه‌های بیشتری نیاز است تا بتوان شکل توده را

به درستی تشخیص داد.

- تعداد داده‌های درون هر توده: هر چه تعداد داده‌های درون یک توده کمتر باشد به  $k$  کوچک‌تری نیاز است.

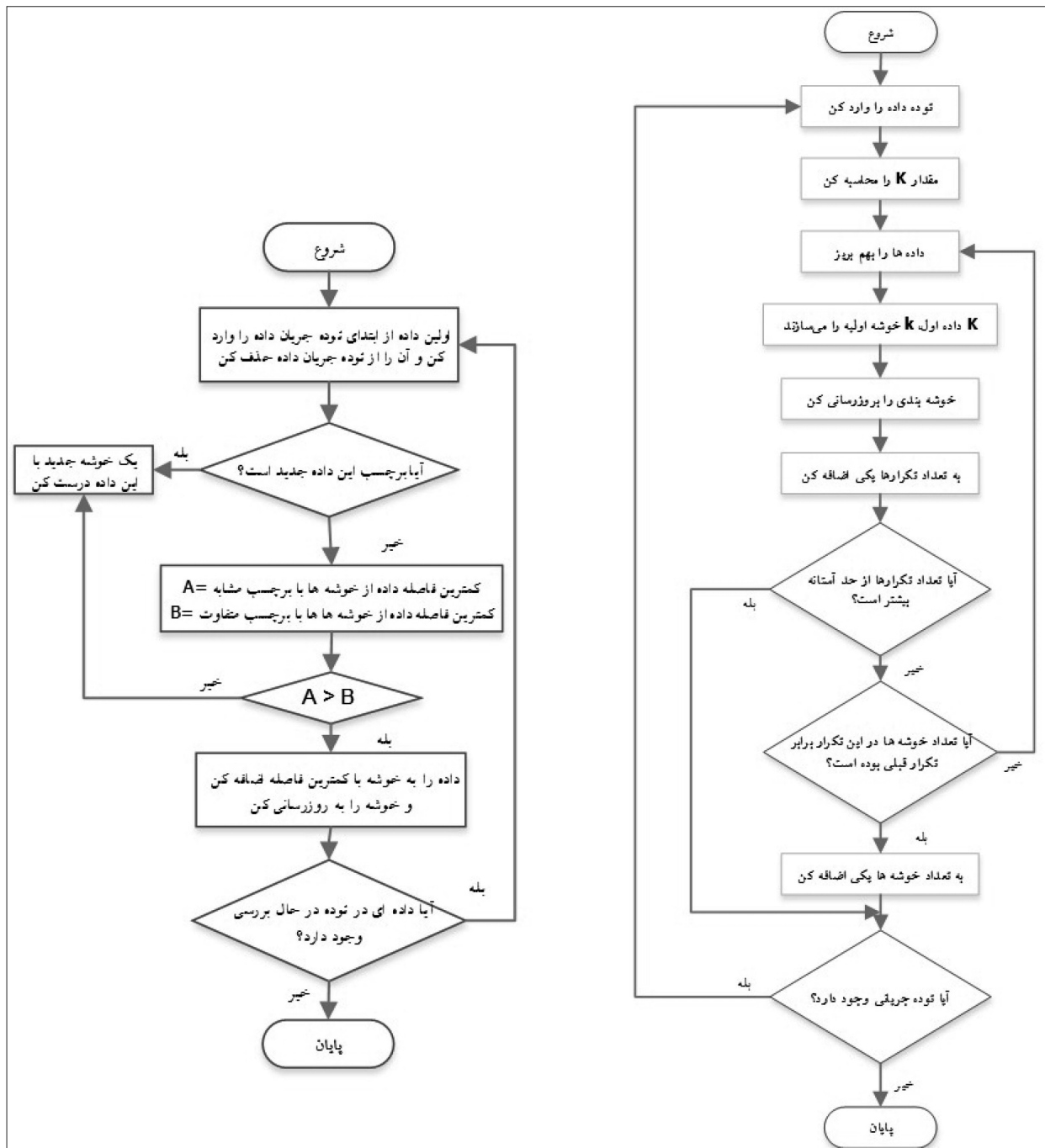
همه عوامل مذکور در تعیین مقدار  $k$  موثر هستند ولی برای تعیین  $k$  تنها به آخرین عامل توجه کرده‌ایم؛ زیرا عوامل اول و دوم معیارهای قطعی برای تعیین  $k$  نیستند و محاسبه سومین عامل متحمل هزینه زیادی در زمان اجرای الگوریتم است.

هر چه مقدار  $k$  بیشتر باشد باعث افزایش دقت می‌شود ولی از سوی دیگر با افزایش  $k$  امکان خوشه‌بندی بیش از اندازه لازم افزایش می‌یابد. به گونه‌ای که در بدترین حالت هر یک از داده‌ها معرف یک خوشه می‌شوند.

در الگوریتم ISAIC مقدار  $k$  طبق رابطه (۱) در نظر گرفته شده است که در آن تعداد داده‌های هر توده است. متغیر  $C$  با توجه به پراکندگی داده مقادری می‌شود. اگر  $C$  عددی بزرگ باشد، باعث کاهش تعداد خوشه‌های اولیه خواهد شد و ممکن است خوشه‌بندی به درست انجام نشود و الگوریتم پیشنهادی مشابه الگوریتم SAIC عمل می‌کند. اگر  $C$  مقداری کوچک داشته باشد باعث افزایش تعداد خوشه‌های اولیه می‌شود و اگر  $C$  برابر با یک باشد، برای هر داده یک خوشه در نظر گرفته می‌شود.

$$k = NS/C \quad (1)$$

شکل ۲ روند کلی الگوریتم ISAIC را نشان می‌دهد. همان‌طور که در شکل ۲(الف) ملاحظه می‌شود الگوریتم ISAIC توده‌های جریان داده را به ترتیب دریافت می‌کند. به دلیل این‌که در گام خوشه‌بندی داده‌ها ترتیب ورود داده‌ها مهم است، این الگوریتم داده‌های یک توده را به هم می‌ریزد سپس خوشه‌بندی می‌کند و تعداد خوشه‌ها را می‌شمارد. در صورتی که تعداد خوشه‌ها برابر با تعداد خوشه‌های مرحله قبل باشد، عمل خوشه‌بندی توده جاری خاتمه می‌یابد و پس از آن الگوریتم منتظر ورود توده جدید می‌شود. در صورتی که تعداد خوشه‌ها برابر با تعداد خوشه‌های مرحله قبل نباشد الگوریتم بار دیگر



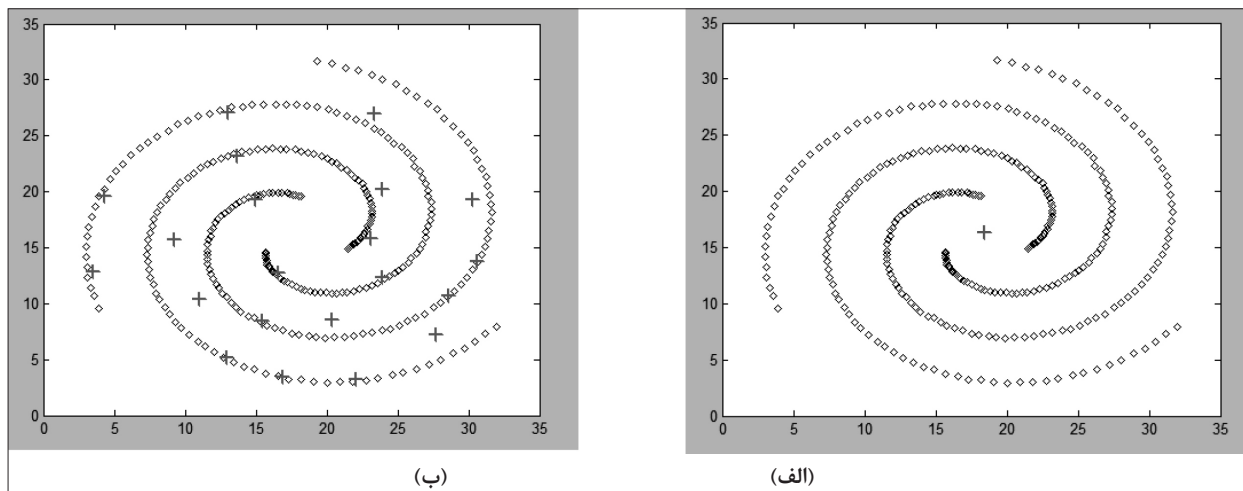
شکل ۲: الف- روند کلی الگوریتم خوشه‌بندی پیشنهادی ISAIC. ب- خوشه‌بندی هر توده جریان داده (الف) (ب)

خوشه‌بندی شکل داده‌ها به درستی تشخیص داده شود. در شکل ۲(ب) روال خوشه‌بندی نشان داده شده است.

#### ۴- آزمایش‌ها

برای ارزیابی روش ISAIC و مقایسه با روش SAIC مجموعه داده‌های متفاوتی به کار گرفته شده است.

داده‌ها را به هم می‌ریزد و خوشه‌بندی را تکرار می‌کند. این عملیات تکرار می‌شود تا جایی که یا تعداد تکرارها برابر حد آستانه شود یا تعداد خوشه‌ها در دو تکرار متوالی تغییر نکند. قبل از ورود به گام خوشه‌بندی در الگوریتم ISAIC در صورتی که تعداد خوشه‌های موجود از یک حد آستانه‌ای کمتر باشد  $k$  خوشه اولیه در نظر گرفته می‌شود، تا در گام



شکل ۳: نتیجه خوشه‌بندی برای مجموعه داده Spiral۱-الف- روش SAIC-ب- روش ISAIC.

این دو الگوریتم را روی این جریان داده بررسی می‌کنیم. جریان داده مورد آزمایش دارای سه برچسب است که توده ۱ دارای برچسب ۱ و توده‌های ۲ و ۳ به ترتیب دارای برچسب ۲ و ۳ هستند. زمان اجرای الگوریتم SAIC روی این مجموعه داده ۰,۲۰۲۷ ثانیه است. با به‌کارگیری الگوریتم ISAIC روی همین مجموعه داده زمان ۰,۲۶۰۸ ثانیه است.

در شکل ۴، داده‌های سیاه رنگ، برچسب ۱ دارند و داده‌های سبز و آبی رنگ به ترتیب برچسب‌های ۲ و ۳ دارند. از طرفی خوشه‌های قرمز دارای برچسب ۱ هستند و خوشه‌های آبی و بنفش رنگ به ترتیب دارای برچسب‌های ۲ و ۳ هستند. عمده اختلاف روش‌های SAIC و ISAIC در ستون دوم متناظر با توده اول شکل ۴ قابل مشاهده است. هر دو الگوریتم در این مرحله داده‌ها را به درستی خوشه‌بندی کرده‌اند ولی روش ISAIC شکل توده را نیز تشخیص داده است. با ورود توده جدید، نتیجه خوشه‌بندی روش SAIC درست نمی‌باشد زیرا انتظار داریم خوشه‌های قرمز رنگ (که معرف داده‌های با برچسب ۱ هستند) به داده‌های سیاه رنگ (که دارای برچسب ۱ هستند) نسبت به خوشه‌های دیگر نزدیک‌تر باشند. همان‌طور که در شکل ۴-ج نشان داده شده است روش ISAIC جریان داده را به درستی خوشه‌بندی کرده است.

ترتیب توده‌ها در جریان خوشه‌بندی با روش SAIC مهم

برخی از این مجموعه داده‌ها با روش SAIC به خوبی خوشه‌بندی می‌شوند و برخی با این روش به درستی خوشه‌بندی نمی‌شوند. از سوی دیگر برخی از این مجموعه داده‌ها ایستا هستند و برخی دیگر جریان داده هستند. در مجموع شش مجموعه داده برای ارزیابی به‌کار گرفته شده است که دو مورد از آن‌ها ایستا و بقیه جریان داده هستند. به‌عنوان اولین مثال مجموعه داده ایستای Spiral۱ را مورد ارزیابی قرار می‌دهیم که از روی مجموعه داده Spiral۱ [۱۳]، با تغییر برچسب تمام داده‌ها به یک به‌دست آمده است. متغیر C طبق آزمایش‌های متعدد برابر با ۳۰ در نظر گرفته شده است.

نتایج آزمایش‌های در شکل ۳ نشان می‌دهد که برای مجموعه داده Spiral1 تک دسته‌ای روش SAIC تنها یک خوشه در مرکز شکل می‌یابد (شکل ۳-الف) در حالی که روش پیشنهادی شکل را تشخیص می‌دهد (شکل ۳-ب). اگرچه این امر در مسئله تک‌هدفه مشکلی ایجاد نمی‌کند ولی اگر یک توده جریان داده از یک دسته باشد و از خوشه‌های شکل گرفته تا آن لحظه (در توده‌های قبل)، به اندازه کافی دور باشد ممکن است در نهایت منجر به خطای قابل ملاحظه‌ای گردد.

به‌عنوان دومین مثال، دو الگوریتم SAIC و ISAIC را روی جریان داده pathbased با در نظر گرفتن سه توده مطابق شکل ۴ آزمایش کرده و زمان اجرا و نحوه عملکرد

الف: داده خام	توده ۱	توده ۲	توده ۳	کل داده
ب: نتیجه خوشه‌بندی SAIC				
ج: نتیجه خوشه‌بندی ISAIC				

شکل ۴: الف) توده جریان داده pathbased، نتیجه خوشه‌بندی الگوریتم‌های (ب) SAIC و (ج) ISAIC روی جریان داده pathbased.

در حالی که این مشکلات همان‌طور که در شکل ۵ (ج) نشان داده شده است در الگوریتم ISAIC رخ نمی‌دهد.

شکل ۶ نمای کلی دو روش را روی جریان داده Ring-Ball2 با ترتیب (۱) نشان می‌دهد مجموعه داده‌ها و خوشه‌ها همزمان در این شکل نشان داده شده است. در شکل ۶-الف انباشتگی کلاس ۲ در نقطه‌ای نزدیک به مرکز کاملاً مشهود است در این نقطه خوشه‌ای با برچسب ۱ به نادرستی قرار داده شده است. در حالی که وقتی این توده‌های جریان داده را به‌عنوان ورودی الگوریتم ISAIC قرار دادیم همان‌طور که در شکل ۶-ب مشهود است دیگر این انباشتگی نابجا رویت نشد.

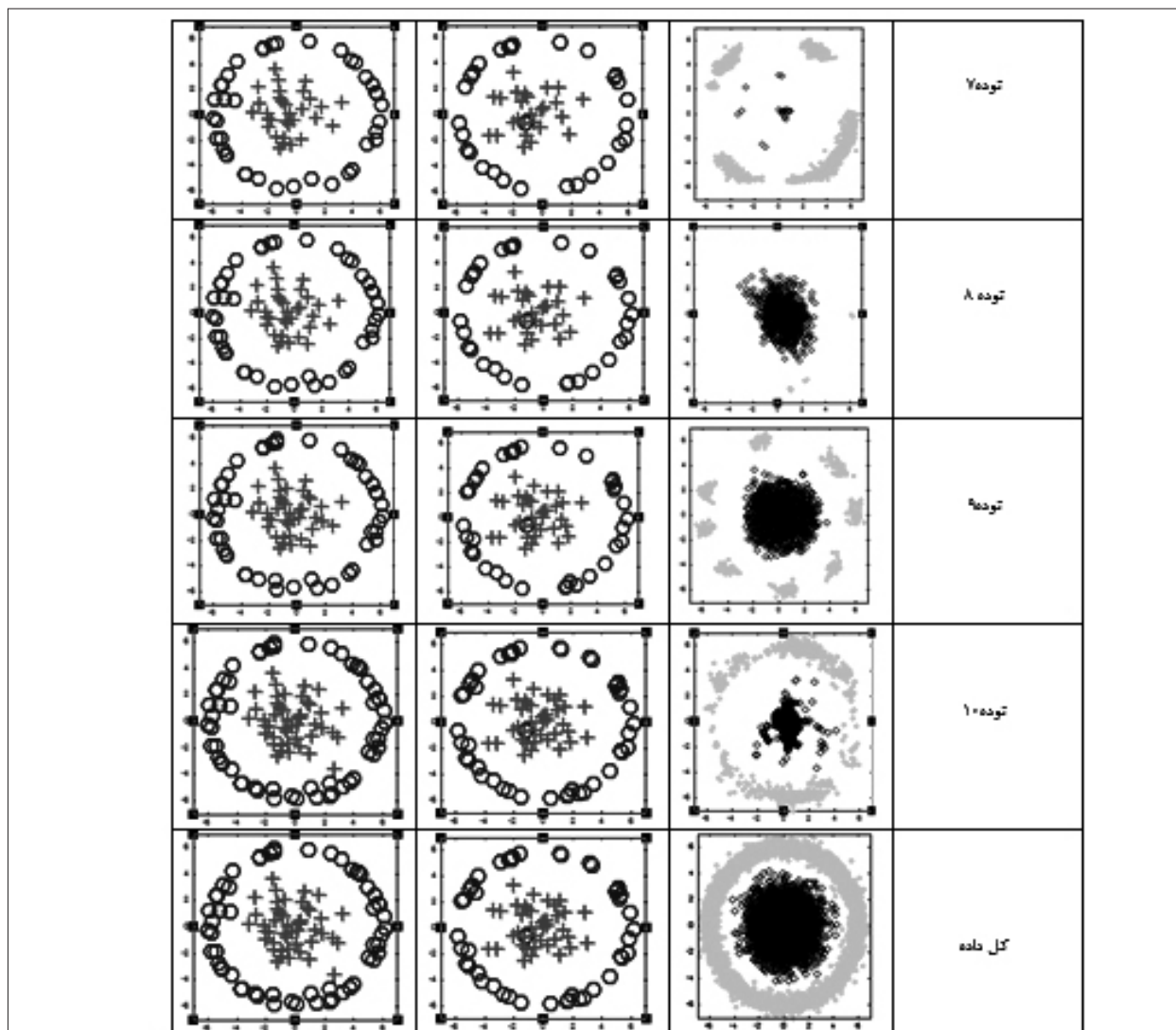
با جابجایی توده اول و هفتم روی جریان داده Ring-Ball2 ترتیب (۲) را تعریف می‌کنیم. همان‌طور که گفته شد ترتیب ورود توده‌ها در الگوریتم SAIC بسیار تاثیرگذار است. این دفعه همان‌طور که در شکل ۷ می‌بینید با تغییر ترتیب توده‌ها هر دو الگوریتم نتیجه خوبی می‌دهند، بنابراین این امر نشان می‌دهد که الگوریتم SAIC تحت تاثیر ترتیب ورود توده‌هاست در حالی که به دلیل این‌که جریان داده برخط است هیچ کنترلی بر ترتیب ورود این توده‌ها نداریم. همان‌طور که در شکل ۷ مشاهده می‌شود روی جریان داده Ring-Ball2 با ترتیب (۲) هر دو الگوریتم خوب عمل کرده و هیچ انباشتگی یا خوشه با برچسب نادرست دیده نمی‌شود. حال این آزمایش را روی دو مجموعه داده دیگر که

است ولی این امر هیچ خللی در روش ISAIC ایجاد نمی‌کند. در واقع در تمام مسائلی که روش SAIC عملکرد خوبی دارد، روش ISAIC نیز عملکرد خوبی دارد. به این منظور این دو روش را روی جریان داده Ring-Ball2 [۸] یک مرتبه با همان ترتیب ذکر شده توسط Zheng و دیگران [۹] و مرتبه بعد این ترتیب را تغییر می‌دهیم و مورد آزمایش قرار می‌دهیم

در شکل ۵-الف داده‌های اولیه جریان داده نشان داده شده‌اند. توده‌های نشان داده شده در این قسمت به ترتیب (۱) وارد می‌شوند بعد از پردازش به وسیله الگوریتم تنها خوشه‌های ایجاد شده توسط الگوریتم باقی می‌ماند و این توده داده از بین می‌رود. در نهایت همان‌طور که مقایسه شکل‌های ۵ (ب) و ۵ (ج) نشان می‌دهد بعد از ورود توده اول، الگوریتم SAIC یک خوشه با برچسب ۱ در مرکز شکل ایجاد می‌کند که با توجه به این‌که این خوشه نماینده تعداد زیادی داده است با ورود توده‌های بعدی نیز حذف نمی‌شود (در سطر توده اول و ستون (ب) در شکل ۵). این در حالی است که هیچ داده‌ای با برچسب ۱ در حوالی مرکز شکل هیچگاه وارد نشده است و در این نواحی فقط داده با برچسب ۲ وارد شده است. علاوه بر این خوشه با برچسب نادرست که در الگوریتم SAIC رخ می‌دهد، چگالی خوشه‌های با برچسب ۲ نیز در حوالی این خوشه نادرست زیاد است؛ در شکل ۶ این مسئله مشهود است.



(ج) نتیجه خوشبجدي ISAIC	(ب) نتیجه خوشبجدي SAIC	(الف) داده اولیه	
			گروه ۱
			گروه ۲
			گروه ۳
			گروه ۴
			گروه ۵
			گروه ۶



شکل ۵: توده‌های جریان داده Ring-Ball2 با ترتیب (۱)

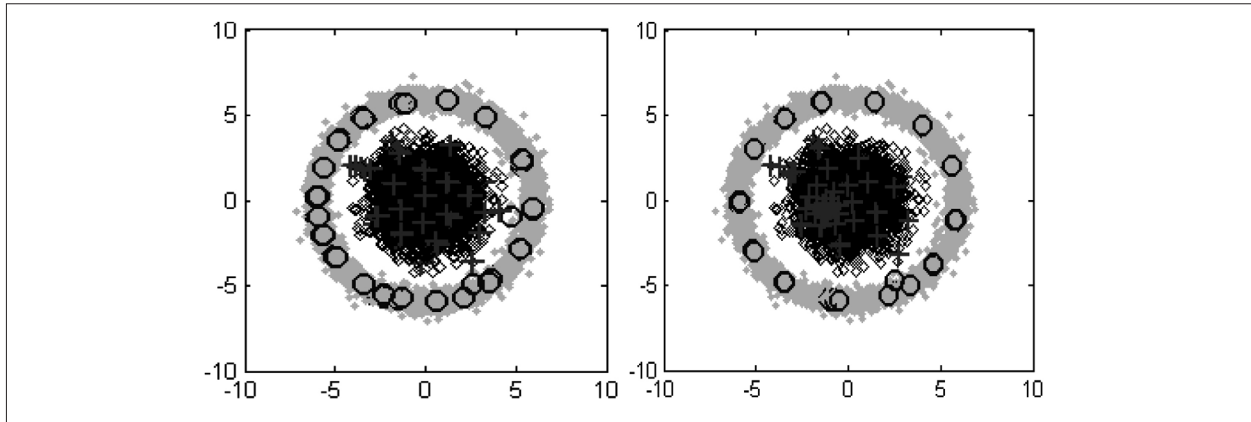
درستی نشان داده‌اند. زمان اجرای دو الگوریتم نزدیک به هم بوده است. در این الگوریتم برای بی‌تاثیر کردن ترتیب نقاط توده، داده‌های توده را به هم می‌ریزیم و این کار را آنقدر تکرار می‌کنیم که تعداد خوشه‌ها در ۲ تکرار متوالی برابر یا حداقل به اندازه حد آستانه تکرار صورت گرفته باشد. این امر سبب می‌شود اولاً در هر بار اجرا زمان صرف شده کمی تغییر کند دوم این‌که زمان اجرای دو الگوریتم نزدیک به هم باشند. آنچه که مشهود است زمان اجرای هر دو الگوریتم کسری از دقیقه است که نشان از کارا بودن الگوریتم دارد.

نتیجه تکرار این آزمایش روی جریان داده Spril در

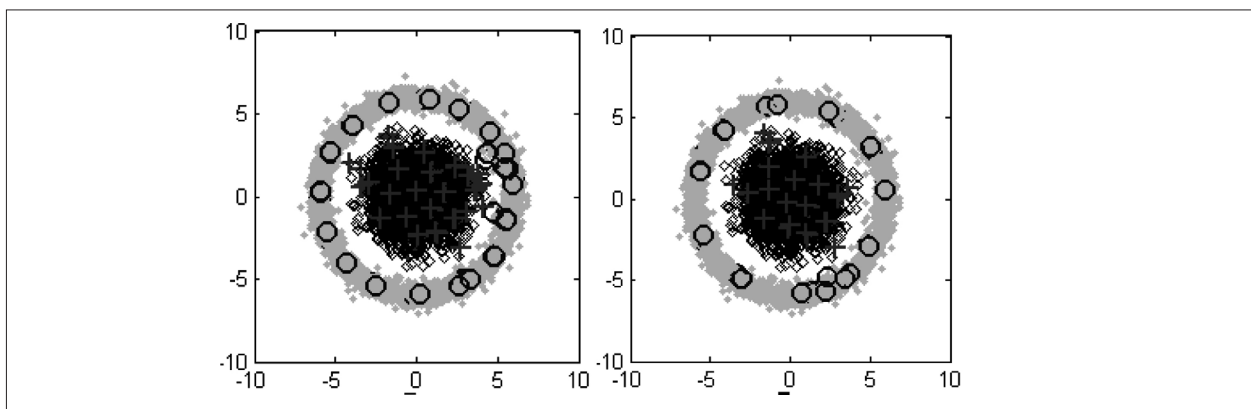
یکی ایستا و دیگری جریان داده است امتحان می‌کنیم. برای این مجموعه داده‌ها روش SAIC خوب عمل می‌کند و نشان می‌دهیم روش ISAIC نیز خوب عمل می‌کند. مجموعه داده ایستای مورد بررسی Ringball1 [۱۱] است.

نتیجه آزمایش روی مجموعه داده ایستا Ringball1 با الگوریتم SAIC، ۶,۳۹۴۳ ثانیه طول کشیده است و الگوریتم SAIC، ۶,۵۴۵۱ ثانیه به طول انجامیده است. در شکل ۸ نتیجه آزمایش نشان داده شده است.

همان‌طور که در شکل ۸ مشاهده می‌شود هر دو الگوریتم نه تنها توانسته‌اند خوشه‌هایی با برچسب صحیح را ایجاد کنند بلکه شکل داده ایستای Ringball1 را نیز به



شکل ۶: الف- اجرای الگوریتم SAIC روی Ring-Ball2 با ترتیب (۱) - ب- اجرای الگوریتم ISAIC روی Ring-Ball2 با ترتیب (۱)



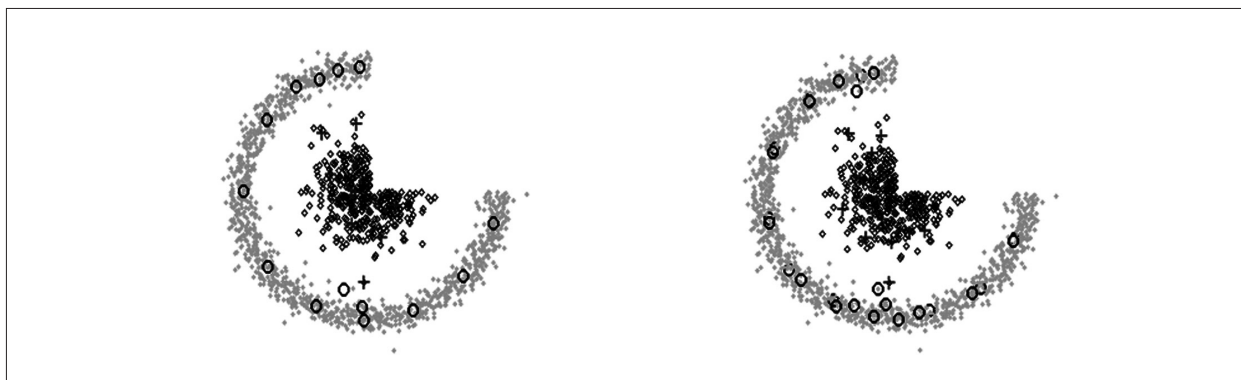
شکل ۷: الف- اجرای الگوریتم SAIC روی جریان داده Ring-Ball2 با ترتیب (۲) - ب- اجرای الگوریتم ISAIC روی جریان داده Ring-Ball2 با ترتیب (۲).

داده بررسی می‌کنیم. زمان اجرای الگوریتم SAIC روی این مجموعه داده ۰,۲۹۹۵ ثانیه است. با به‌کارگیری الگوریتم ISAIC روی همین مجموعه داده زمان ۰,۳۰۷۹ ثانیه است. شکل ۱۰ نشان می‌دهد هر دو الگوریتم SAIC و ISAIC خوشه‌ها برچسب درستی دارند و شکل را تشخیص می‌دهند.

در آزمایشی دیگر، یک جریان داده جدید با نام circles تعریف می‌کنیم. این جریان داده شامل دو دایره به مرکز (۲ و ۲) و به شعاع‌های ۰.۸ و ۲، یک دایره به مرکز (۲ و ۲) و به شعاع ۱ و دو دایره به مرکز (۱۶ و ۲) و به شعاع‌های ۰.۵ و ۲ است. توده‌هایی که منجر به شکل‌گیری این جریان داده شده‌اند در شکل ۱۱ و در ستون دوم آن به ترتیب نشان داده شده‌اند. نکته قابل توجهی که در این توده‌ها وجود دارد تنوع برچسب می‌باشد و برخلاف مثال‌های قبلی که

شکل ۹ نشان داده شده است. در قسمت ۹ (الف) داده‌های اولیه نشان داده شده‌اند که به ترتیب نمایانگر توده اول، توده دوم، توده سوم و نهایتاً کل داده‌ها هستند. در ۹ (ب) نتیجه خوشه‌بندی الگوریتم SAIC بعد از ورود هر یک از توده‌ها است. در شکل ۹ (ج) نتیجه خوشه‌بندی الگوریتم ISAIC بعد از ورود توده‌ها به ترتیب ذکر شده است.

همان‌طور که در شکل ۹ مشهود است هر دو الگوریتم توانسته‌اند به خوبی این توده جریان داده را خوشه‌بندی کنند. نتیجه خوشه‌بندی توسط این دو الگوریتم در شکل تکمیلی ۱۰ نشان داده می‌شود. در شکل-الف داده‌های برچسب دار و خوشه‌های معرفی شده توسط دو الگوریتم SAIC و ISAIC هم‌زمان نشان داده شده‌اند که کمک می‌کنند بتوان نتیجه خوشه‌بندی را بهتر ملاحظه کرد. زمان اجرا و نحوه عملکرد این دو الگوریتم را با هم روی این جریان



(الف) (ب)

شکل ۸: اجرای الگوریتم‌های (الف) SAIC و (ب) ISAIC روی مجموعه داده Ringball1.

کل داده	توده ۳	توده ۲	توده ۱	الف: داده‌های اولیه
				ب: نتیجه خوشه‌بندی SAIC
				ج: نتیجه خوشه‌بندی ISAIC

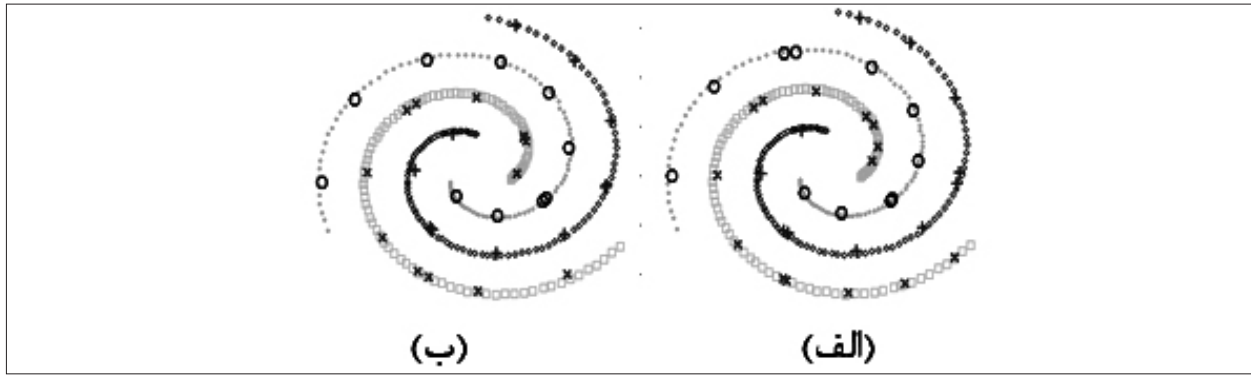
شکل ۹: توده‌های جریان داده Spiral به ترتیب از راست به چپ

مناسب برای داده‌های اولیه سبز و مشکی رنگ با شعاع بیشتر را به درستی پیدا کند. بنابراین اشتباها به ازای این مقادیر خوشه‌هایی در مرکز داده‌های سبز و سیاه رنگ با شعاع کوچکتر تشخیص داده در حالی که این خوشه‌ها به جای نزدیکی به داده‌های با شعاع بزرگتر به داده‌های با شعاع کوچکتر نزدیکتر هستند. شکل ۱۲ (ب) اجرای الگوریتم ISAIC را نشان می‌دهد. این الگوریتم علاوه بر این که خوشه‌بندی صحیحی انجام داده است شکل را هم به درستی تشخیص داده است. نکته قابل ملاحظه‌ای که در این روش‌ها وجود دارد زمان اجرای آن‌هاست که بسیار مناسب است و در تمام این الگوریتم‌ها مشابه است.

جدول ۱ دقت دسته‌بندی روی چند مجموعه داده را با به‌کارگیری الگوریتم‌های SAIC و ISAIC نشان می‌دهد. نتایج نشان می‌دهد که کارایی الگوریتم ISAIC مشابه یا بالاتر

مشکل خوشه‌بندی تنها در توده اول روی می‌داد در این مثال حتی در توده ۴ هم با این مشکل مواجه هستیم. در ستون دوم شکل ۱۱ این جریان داده نشان داده شده است که شامل ۵ توده و در کل شامل ۳ دسته است. در ستون سوم شکل ۱۱ نتیجه خوشه‌بندی این جریان داده به ترتیب با ورود این توده‌ها نشان داده شده است. آخرین سطر از این شکل در ابتدا تمام داده‌های اولیه حاصل از این جریان داده را در ستون دوم خود نشان می‌دهد. سپس در ستون سوم خوشه‌بندی این جریان داده با ورود هر یک از این توده‌ها نشان داده شده است.

در شکل ۱۲ داده‌های اولیه و نتیجه خوشه‌بندی همزمان نشان داده شده‌اند. شکل ۱۲ (الف) نتیجه خوشه‌بندی حاصل از الگوریتم SAIC را نشان می‌دهد. همان‌طور که در این شکل مشهود است. الگوریتم نتوانسته خوشه‌های



شکل ۱۰: اجرای الگوریتم‌های (الف) SAIC و (ب) ISAIC روی مجموعه داده Spiral.

توده اول	داده اولیه	نتیجه روش ISAIC
توده دوم		
توده سوم		
توده چهارم		
توده پنجم		
کل توده‌ها		

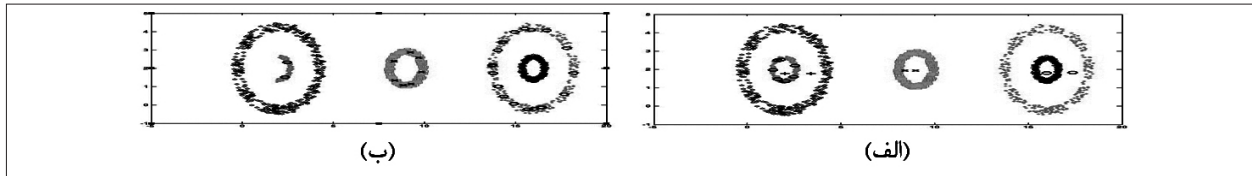
شکل ۱۱: توده‌های اولیه جریان داده circles در ستون دوم و حاصل خوشه‌بندی ISAIC در ستون سوم نشان داده شده است.

چالش‌های پیش روی این مسائل را ذکر کرده سپس به طور خلاصه الگوریتم SAIC را که به‌تازگی معرفی شده است بیان نمودیم. بعد از ذکر محاسن این روش، مشکل پیش روی این الگوریتم بیان شده و راهکاری برای حل این مشکل ارائه شد. در ادامه با ذکر چند نمونه مسئله جریان داده، برتری الگوریتم بهبود یافته نسبت به الگوریتم اولیه نشان داده شد. این الگوریتم عملکرد مناسبی دارد و نسبت به الگوریتم اولیه افزایش قابل توجهی در زمان اجرا

از الگوریتم SAIC است. الگوریتم ISAIC روز مجموعه داده‌های Spiral، Ring-Ball، Pathbased، و Circles نتایج بهتری به‌دست آورده است. میزان بهبود دقت دسته‌بندی بر روی مجموعه داده‌های مورد آزمایش حداقل صفر و حداکثر حدود ۶۵٪ است.

##### ۵- جمع بندی

در این مقاله ابتدا اهمیت الگوریتم‌های جریان داده و



شکل ۱۲: الف خوشه‌بندی بر اساس SAIC ب: خوشه‌بندی بر اساس ISAIC

جدول ۱: دقت دسته‌بندی روی چند مجموعه داده با به‌کارگیری الگوریتم‌های SAIC و ISAIC.

مجموعه داده	دقت دسته‌بندی (%) الگوریتم SAIC	دقت دسته‌بندی (%) الگوریتم ISAIC
Spiral1	100	100
Pathbased	61.391	97.055
Ring-Ball <sub>2</sub> (ترتیب ۱)	83.020	97.055
Ring-Ball <sub>2</sub> (ترتیب ۲)	97.360	97.701
Ring-Ball1	89.226	89.309
Spiral	79.010	81.010
Circles	23.090	98.289

Clustering – Algorithms and Benefits”, In proceedings of the 16th IEEE International Conference on Tools with Artificial Intelligenc (ICTAI04), Boca Raton, Florida, November (2004) 774-776.

ندارد. در کارهای آینده عملکرد روش پیشنهادی بر روی مجموعه داده‌های جریان داده واقعی مورد ارزیابی قرار خواهد گرفت.

### منابع

- [1] Laiwen Zheng, Hong Huo, Yiyou Guo, Tao Fang, Supervised Adaptive Incremental Clustering for data stream of chunks, Neurocomputing 219 (2017) 502–517.
- [2] Kevin. Machine learning a probabilistic perspective. MIT Press, 2012. 875. ISBN 0262018020. Murphy
- [3] Maryam Mousavi, Azuraliza Abu Bakar, and Mohammadmahdi Vakilian, Data Stream Clustering Algorithms: A Review, Int. J. Advance Soft Compu. Appl, Vol. 7, No. 3, November 2015 ISSN 2074-8523.
- [4] Gaber, M. M 2012. Advances in data stream mining. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery. 2(1) 79-85.
- [۵] فاتحی، کاوان؛ هادی مقدر و محمدرضا زارع میرک آباد، ۱۳۹۳، مروری بر الگوریتم‌ها و چالش‌های موجود در خوشه‌بندی جریان داده، دومین کنفرانس بین‌المللی دستاوردهای نوین در علوم مهندسی و پایه، به‌صورت الکترونیکی، مرکز علمی کاوشگر علم،
- [6]Wan, L., Ng, W. K., Dang, X. H., Yu, P. S., Zhang, K 2009. Density-based clustering of data streams at multiple resolutions. ACM Transactions on Knowledge Discovery from Data, 3(3) 1-28.
- [7] Ziemiński, R. Z 2013. DeltaDens–Incremental Algorithm for On–Line Density–Based Clustering. New Trends in Databases and Information Systems. 163-172.
- [8] O’callaghan, L., Mishra, N., Meyerson, A., Guha, S., Motwani, R 2002. Streaming-data algorithms for high-quality clustering. Proceedings of the 18th IEEE International Conference on Data Engineering. 685- 694.
- [9] Gama, J 2012. A survey on learning from data streams: current and future trends. Progress in Artificial Intelligence. 45-55.
- [10] Udommanetanakit, K., Rakthanmanon, T., Waiyamai, K 2007. E-stream: Evolution-based technique for stream clustering. Advanced Data Mining and Applications. 605-615.
- [11] C.-D. Wang, J.-H. Lai, D. Huang, W.-S. Zheng, Svsstream: a support vector-based algorithm for clustering data streams, IEEE Trans. Knowl. Data Eng. 25 (6) (2013) 1410–1424.
- [12] Ericsson, “5G for the networked society beyond 2020,” Mobile World Congress 2013, February 2013.
- [۱۳] اسماعیلی، مهدی. (۱۳۹۲). مفاهیم و تکنیک‌های داده کاوی. کاشان: دانشگاه آزاد اسلامی.
- [14] Eick, C.F., Zeidat, N., and Zhenghong, Z., “Supervised