

time series and combination clustering. In the first step, a time series data set is segmented using a fixed window size, and each segment is clustered by applying a hierarchical clustering algorithm and Euclidean distance. Also, we use a logarithmic relation based on the length of the time series data set to determine the number of components, selecting the best outcomes using various internal criteria including intergroup variance, Kalinsky-Harbaz, and Dunn index. In the second step, the results of the first stage are processed using ensemble clustering and the final clustering label is obtained. We develop two novel algorithms based on different internal criteria for selecting the best segmentations: the first one in which we consider only one internal criterion and the second one in which we consider three internal criteria simultaneously. Moreover, we run various settings on 28 datasets with 10 replications for the two presented algorithms, checking the final precision using an external RAND index. Then, in order to identify the best settings for the proposed algorithms we applied Wilkison statistical test. Statistical comparison of the results of the two new algorithms on 28 data sets with some algorithms in the related literature indicates significant improvement in terms of error rate and execution time. Finally, the findings acquired based on the best settings of the proposed algorithms indicate that the suggested method has the best RAND index among the previous algorithms in the literature for 32% of the dataset tiers.

Key Words: Time series, Clustering, Data mining, Sub-series.

ali.ghorbanian@mail.um.ac.ir
h-razavi@um.ac.ir

یک رویکرد جدید به منظور خوشه‌بندی سری‌های زمانی با استفاده از ترکیب زیرسری‌های زمانی

علی قربانیان (دانشجوی دکتری)

حمیده رضوی* (دانشیار)

گروه مهندسی صنایع، دانشکده مهندسی، دانشگاه فردوسی مشهد

چکیده

خوشه‌بندی سری‌های زمانی فرایندی است که سری‌های زمانی را با توجه به خصوصیات آن‌ها گروه‌بندی می‌کند. در پژوهش‌های پیشین به شباهت موجود بین قطعات یک سری زمانی به منظور خوشه‌بندی توجه کمتری شده است. در این مقاله یک رویکرد جدید دومرحله‌ای بر اساس قطعه‌بندی سری زمانی و خوشه‌بندی ترکیبی ارائه شده است. در مرحله اول یک مجموعه داده سری زمانی با استفاده از اندازه پنجره ثابت قطعه‌بندی شده و هر قطعه به طور جداگانه خوشه‌بندی شده است. سپس با استفاده از معیارهای درونی، بهترین نتایج حاصله انتخاب شده است. در مرحله دوم نتایج حاصل از مرحله اول با استفاده از خوشه‌بندی ترکیبی، پردازش شده و برچسب نهایی خوشه‌بندی حاصل شده است. نتایج الگوریتم ارائه شده نشان‌دهنده افزایش کارایی خوشه‌بندی به میزان ۲,۹۲ درصد و رسیدن به عدد ۶۷,۲۵ می‌باشد. همچنین بررسی عملکرد الگوریتم با بهترین نتایج ادبیات نیز نشان‌دهنده بهترین کارایی با حداقل هزینه زمانی می‌باشد.

کلمات کلیدی: سری‌های زمانی، خوشه‌بندی، قطعه‌بندی، خوشه‌بندی ترکیبی

A new approach to time series clustering by combination of sub-series

Ali Ghorbanian (Phd Student)

Hamideh Razavi* (Associate Professor)

Department Of Industrial Engineering, Faculty Of Engineering, Ferdowsi University Of Mashhad (FUM)

Abstract

Time series-clustering, defined as deriving trends and archetypes from sequential data, divides time series into groups considering their characteristics. Previous work mainly focused on distance criterion and clustering algorithm to cluster the time series so few researchers have investigated the similarities between the segments of a time series. To address this research gap, we propose a new two-step approach based on sub-

* نویسنده مسئول. آدرس پستی: مشهد- دانشگاه فردوسی مشهد- دانشکده مهندسی- گروه مهندسی صنایع

یادگیری عمیق اشاره نمود [۷, ۸]. در این تحقیق از روش خوشه‌بندی سری زمانی کل استفاده می‌شود که در آن یک مجموعه از سری‌های زمانی با توجه به معیارهای مشابهت خوشه‌بندی می‌گردد [۹].

یکی از راهکارهایی که به منظور افزایش دقت خوشه‌بندی مورد توجه قرار گرفته است، استفاده از معیارهای فاصله‌ای مخصوص سری‌های زمانی می‌باشد. معمولاً معیارهای فاصله‌ای متفاوتی برای سری‌های زمانی با توجه به ماهیت این نوع داده‌ها استفاده می‌گردد. تعدادی از رایج‌ترین معیارهای فاصله عبارت است از هاسدورف، همینگ^۱، DTW^۱، اقلیدسی و LCSS^۲ [۱۰].

در مطالعات قبلی دو رویکرد شامل تعریف، بهبود و ترکیب معیارهای فاصله مورد توجه قرار گرفته است [۱۱-۱۳]. در رویکرد اول با معرفی معیار فاصله جدید و یا بهبود معیارهای موجود سعی در بهبود خوشه‌بندی داشته‌اند. اگرچه توسعه معیارهای فاصله‌ای جدید، بهبود و یا ترکیب آن‌ها می‌تواند تا حدودی دقت خوشه‌بندی این نوع داده‌ها را افزایش دهد اما هزینه زمانی (زمان اجرای الگوریتم) زیادی را نیز در بردارد.

منظور از دقت خوشه‌بندی شاخص‌های بیرونی خوشه‌بندی مانند خلوص، شاخص رند^۳ و آنتروپی می‌باشد [۱۴]. به عنوان مثال در پژوهشی که توسط رحیم خان و زکریا با ارائه یک الگوریتم مناسب به منظور معیار LCSS انجام شد تا حدود ۵۰ درصد زمان محاسبات کاهش یافت ولی هم‌زمان باعث کاهش دقت گردید [۱۵]. به همین دلیل، سلیمانی و عابسی با تغییر این معیار به صورت فازی دقت خوشه‌بندی را بهبود بخشیده‌اند ولی در روش ایشان زمان افزایش یافته‌است [۱۲]. همچنین کمال زاده و همکاران یک معیار فاصله جدید به منظور خوشه‌بندی سری‌های زمانی با طول بلند با استفاده از روابط هندسی تعریف کرده‌اند [۱۱]. برای سری‌های زمانی با طول متفاوت، وانگ و همکاران از یک معیار فاصله بر مبنای اختلاف سطح زیر منحنی دو سری زمانی استفاده نموده‌اند [۱۶]. در مطالعات دیگر در همین راستا با معرفی معیار^۴ WDTW،^۴ MP^۵ و^۶ MSCD دقت خوشه‌بندی افزایش یافته‌است [۱۷-۲۰].

در رویکرد دوم به منظور استفاده از مزایای هر یک از معیارهای فاصله نظیر DTW،^۷ DDTW و LCSS از ترکیب این معیارها به منظور افزایش دقت خوشه‌بندی استفاده شده‌است [۵, ۱۳, ۲۱]. منظور از این مزایا شناسایی نقاط پرت و همچنین شناسایی انتقال داده‌ها در طول زمان در یک سری زمانی می‌باشد. با وجود این که

خوشه‌بندی سری‌های زمانی در علوم مختلفی مانند ستاره‌شناسی، بیولوژی و آب‌وهوا استفاده می‌گردد. این نوع خوشه‌بندی مانند خوشه‌بندی سایر داده‌ها یک نوع یادگیری بدون نظارت می‌باشد که در آن اطلاعاتی از برچسب اشیا موجود نمی‌باشد. برچسب‌ها اعدادی هستند که از فرایند خوشه‌بندی حاصل گشته و نشان‌دهنده سری‌های زمانی قرار گرفته در یک خوشه می‌باشند. تعریف معیارهای فاصله‌ای مناسب برای سری‌های زمانی، موضوع پژوهش‌های متعددی بوده است. علت آن است که با توجه به ماهیت داده‌های سری زمانی، معیارهای فاصله‌ای ویژه‌ای مورد نیاز است. دقت خوشه‌بندی و زمان اجرای الگوریتم، دو چالش اصلی در خوشه‌بندی سری‌های زمانی می‌باشد. خوشه‌بندی سری‌های زمانی می‌تواند به صورت مستقیم و یا غیرمستقیم در صنایع و خدمات مورد استفاده قرار گیرد. یکی از کاربردهای اصلی خوشه‌بندی سری‌های زمانی کشف الگوهای رفتاری در مورد تقاضا مصرف به‌عنوان نمونه در حوزه انرژی الکتریکی، مصرف گاز طبیعی، مصرف آب و غیره می‌باشد. هدف از این کار شناخت الگوهای رفتاری مصرف‌کنندگان به منظور مدیریت مصرف متقاضیان در زمان‌های مختلف می‌باشد. همچنین به منظور کشف الگوهای مشابه در حوزه سلامت مانند همه‌گیری کوید-۱۹ نیز از این تکنیک استفاده می‌گردد [۱]. به صورت غیرمستقیم نیز خوشه‌بندی سری‌های زمانی به منظور پیش‌بینی این نوع داده مورد استفاده قرار می‌گیرد. در واقع تحقیقات نشان داده است که پیش‌بینی سری‌های زمانی در خوشه‌های یکسان می‌تواند نتایج بهتری را داشته باشد.

خوشه‌بندی سری زمانی شامل سه نوع خوشه‌بندی، خوشه‌بندی سری زمانی کل، خوشه‌بندی توالی و خوشه‌بندی نقطه‌ای می‌باشد. در خوشه‌بندی سری زمانی کل برخلاف دو نوع دیگر مجموعه از سری‌های زمانی با توجه به معیارهای مشابهت مانند داشتن کمترین فاصله از یکدیگر در گروه‌های مختلف قرار می‌گیرند. به عنوان مثال در شکل ۲ چندین سری زمانی در دو گروه با رنگ آبی و سبز خوشه‌بندی شده‌اند که نشان‌دهنده خوشه‌بندی کل می‌باشد. از روش‌های خوشه‌بندی سری‌های زمانی می‌توان روش‌های نیمه نظارت‌شده [۲, ۳]، رویکردهای ترکیب و جنگل تصادفی [۴]، ترکیب روش‌های خوشه‌بندی بر مبنای فاصله و چگالی [۵]، استفاده از روش‌های خوشه‌بندی وزنی [۶] و استفاده از شبکه‌های عصبی و رویکردهای

^۵ Matrix profile

^۶ Maximum shifting correlation distance

^۷ Derivative dynamic time warping

^۱ Dynamic time warping

^۲ Longest common subsequence

^۳ Rand Index

^۴ Weighted dynamic time warping

ترکیب معیارها دقت خوشه‌بندی را افزایش می‌دهد به دلیل افزایش حجم محاسبات، زمان بیشتری صرف می‌کند. به‌عنوان نمونه استفاده از معیار DDTW به‌منظور خوشه‌بندی ۸۴ مجموعه داده، زمانی نزدیک به ۸۰ ساعت برای هر مجموعه داده در برداشته است در حالی که متوسط شاخص رند برای هر مجموعه داده نیز برابر با ۶۰ می‌باشد [۲۲]. علیرغم زمان اجرای طولانی، مشاهده می‌گردد که دقت خوشه‌بندی نیز افزایش چشم‌گیری نداشته است به همین سبب در این پژوهش از فاصله اقلیدسی استفاده می‌گردد.

به‌منظور برخورد با چالش هزینه زمانی نیز چندین روش مورد توجه قرار گرفته‌است. روش اول، استفاده از الگوریتم‌های چندمرحله‌ای به‌منظور کاهش حجم مجموعه داده و روش دوم استفاده از مشخصه‌های سری‌های زمانی به‌عنوان متغیرها برای خوشه‌بندی می‌باشد. به‌عنوان نمونه در روش اول آقابزرگی و همکاران، با استفاده از یک الگوریتم خوشه‌بندی در فاز اول توانستند زمان خوشه‌بندی را کاهش دهند [۲۳]. همچنین ژانگ و همکاران نیز با استفاده از یک شبکه وزن‌دار جهت در یک الگوریتم دوم‌مرحله‌ای، پیچیدگی مسئله خوشه‌بندی را کاهش داده‌اند [۲۴]. ماناکوا و تاچنکو با به‌کارگیری معیار هایپکینز و سیستم نمونه‌گیری، یک الگوریتم دوم‌مرحله‌ای را توسعه داده‌اند [۲۵]. در روش دوم، از مشخصه‌های سری زمانی به‌منظور خوشه‌بندی استفاده می‌شود. این مشخصه‌ها می‌تواند به‌صورت مستقیم از سری زمانی استخراج گردد. ازجمله مشخصه‌های پرکاربرد می‌توان به معیارهای واریانس، همبستگی مرتبه اول، خطی بودن، انحنای فصلی بودن، نقطه اوج و فرورفتگی اشاره نمود [۲۶]. به‌عنوان مثال زو و همکاران به‌منظور استخراج مشخصه‌های سری زمانی از تبدیل سه نوع گراف بازگشتی، پدیداری و شبکه‌های انتقال استفاده نمودند [۲۷]. داسیلوا با استخراج مشخصه‌های متوسط درجه، متوسط طول مسیر، تعداد اجتماعات، ضریب خوشه‌بندی و چگالی، از تبدیل سری زمانی به سه نوع گراف NVG، HVG و QG جهت خوشه‌بندی استفاده کرده است [۲۸]. همچنین فریرا و ژاویه به‌منظور خوشه‌بندی یک مجموعه داده از دو رویکرد ϵ -NN و k -NN به‌منظور تبدیل مجموعه داده به یک شبکه پیچیده بدون وزن استفاده نمودند [۲۹]. استفاده از این رویکرد اگرچه موجب کاهش زمان محاسبات می‌گردد اما دقت خوشه‌بندی را کاهش می‌دهد.

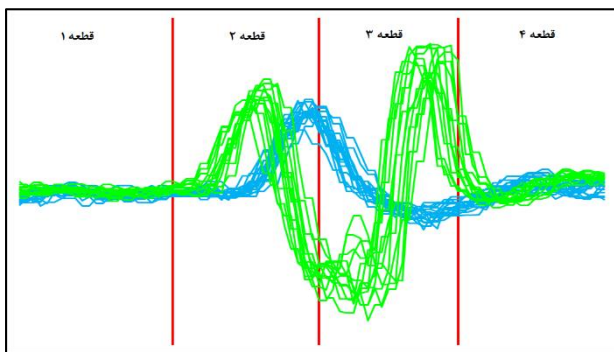
مطالعات اخیر نشان داده‌است که قطع‌بندی یک سری زمانی می‌تواند نقش مهمی در افزایش دقت خوشه‌بندی سری‌های زمانی داشته باشد. به‌عنوان مثال، گویجو و همکاران در مرحله اول یک سری زمانی را قطع‌بندی نموده و در مرحله بعد با توجه به خصوصیات آماری هر یک از قطع‌های ایجادشده، مشخصه‌های آن‌ها را استخراج نموده‌اند. سپس از مشخصات استخراج‌شده به‌منظور خوشه‌بندی

نهایی استفاده کرده‌اند [۲۲]. همچنین بوناسینا و همکاران به‌منظور خوشه‌بندی از ترکیب رویکردهای تبدیل سری زمانی به شبکه و قطع‌بندی، استفاده نموده‌اند [۳۰]. عملکرد مناسب الگوریتم‌های ارائه‌شده بر مبنای قطع‌بندی جایگاه آن را در مباحث خوشه‌بندی تثبیت نموده است.

در پژوهش‌های پیشین در بیشتر موارد به‌منظور بهبود کارایی خوشه‌بندی یک مجموعه داده سعی بر استفاده از معیارهای فاصله جدید و یا الگوریتم‌های جدید در حوزه یادگیری ماشین بوده‌است. اما نتایج ارائه‌شده نشان داده است که معمولاً معیارهای فاصله‌ای جدید مورد استفاده می‌توانند بسیار زمان‌بر باشند و استفاده عملی از این نوع معیارهای فاصله کاربردی نمی‌باشد. همچنین نتایج حاصل از این پژوهش‌ها نشان داده است که استفاده از الگوریتم‌های خاص شاید بتواند در بعضی از مجموعه داده‌ها نتایج ملموسی داشته باشد، اما این مورد همه مجموعه داده‌ها جامعیت ندارد. با توجه به پژوهش‌های انجام‌گرفته در حوزه داده‌کاوی می‌توان به نقش مهم خوشه‌بندی ترکیبی در افزایش کارایی خوشه‌بندی در داده‌های مختلف اشاره نمود، اما در هیچ‌کدام از پژوهش‌های انجام‌گرفته در حوزه سری‌های زمانی این تکنیک موردتوجه قرار نگرفته است. همچنین پژوهش‌های اخیر نشان داده است که زیر سری‌های ایجادشده با استفاده از قطع‌بندی یک سری زمانی می‌تواند دارای اهمیت فراوانی به‌منظور خوشه‌بندی یک مجموعه داده باشد. با توجه به موارد اشاره شده هدف از این پژوهش ارائه یک الگوریتم با کارایی بالا و سریع به‌منظور خوشه‌بندی یک مجموعه داده سری زمانی می‌باشد. در همین راستا در تحقیق فوق نشان داده‌ایم که لزوماً همه قطعات ایجادشده در یک مجموعه داده سری زمانی نمی‌توانند نماینده خوبی به‌منظور استفاده در خوشه‌بندی نهایی باشد، برای رفع این مسئله نشان داده‌ایم که می‌توان با ابتدا با استفاده از معیارهای درونی بهترین قطعات ایجادشده را انتخاب نمود و با استفاده از بهترین قطعات ایجادشده و خوشه‌بندی ترکیبی این قطعات، یک الگوریتم سریع و دقیق به‌منظور خوشه‌بندی داشته باشیم.

در این مقاله، با توجه به نقش مهم قطع‌بندی سری‌های زمانی به‌منظور خوشه‌بندی، دو الگوریتم خوشه‌بندی بر مبنای قطع‌بندی و خوشه‌بندی ترکیبی ارائه شده است. الگوریتم‌های ارائه‌شده دارای سه گام اصلی قطع‌بندی، خوشه‌بندی و خوشه‌بندی ترکیبی می‌باشند. در الگوریتم اول به‌منظور خوشه‌بندی ترکیبی نهایی از یک معیار درونی به‌منظور انتخاب تعداد مشخص از نتایج به دست آمده استفاده می‌گردد، اما در الگوریتم دوم به‌طور هم‌زمان از سه معیار درونی برای این انتخاب استفاده می‌گردد.

شکل ۱ فلوجارت کلی روش پژوهش را نمایش می‌دهد. بعد از مقدمه و مرور ادبیات طرحی بانام خوشه‌بندی ترکیبی بر مبنای

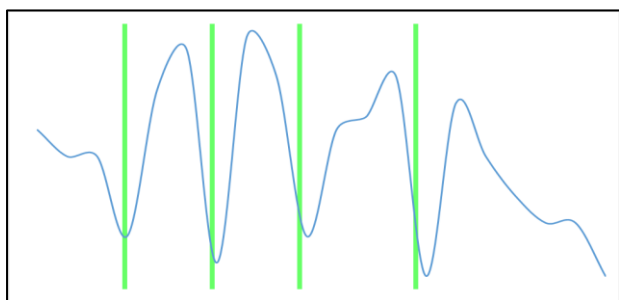


شکل ۲: قطعه‌بندی یک سری زمانی

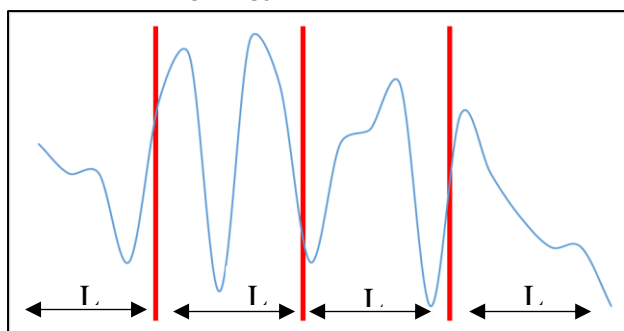
۲-۱ قطعه‌بندی

قطعه‌بندی یکی از روش‌های تحلیل سری‌های زمانی می‌باشد که در آن به دو روش متفاوت می‌توان توالی از زیرسری‌ها را ایجاد نمود. منظور از توالی از زیر سری‌های، سری‌های زمانی با ابعاد کوچک‌تر از سری زمانی اصلی می‌باشند که در کنار یکدیگر سری زمانی اصلی را تشکیل می‌دهند. به‌عنوان مثال در شکل ۲ چهار زیر سری تشکیل شده با استفاده از قطعه‌بندی در کنار یکدیگر سری زمانی اولیه را شکل می‌دهند.

در روش اول یک سری زمانی با استفاده از الگوریتم‌های مشخصی می‌تواند به زیرمجموعه‌هایی با طول‌های غیر یکسان تبدیل گردند [۳۱, ۳۲]. شکل ۳ (الف) این نوع قطعه‌بندی را نمایش می‌دهد. در روش دوم که در این پژوهش نیز مورد استفاده قرار گرفته است ابتدا پنجره‌ای به طول L تعریف گردیده و سپس سری زمانی به زیرمجموعه‌هایی با طول‌های برابر با L تقسیم می‌گردد [۳۳]. شکل ۳ (ب) این نوع قطعه‌بندی را نمایش می‌دهد.



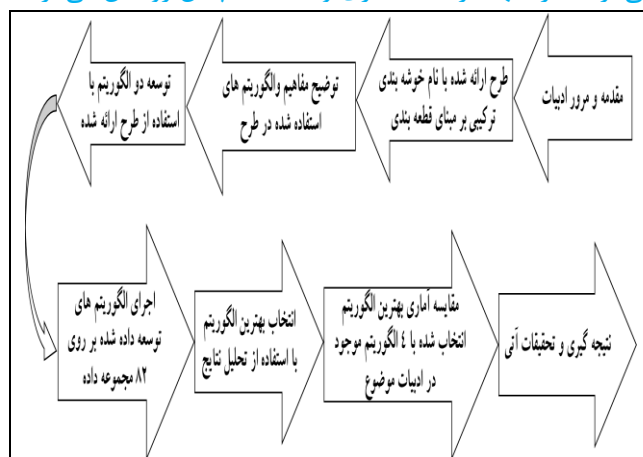
الف- قطعه‌بندی با طول متغیر



ب- قطعه‌بندی با طول ثابت

شکل ۳: انواع مختلف قطعه‌بندی سری زمانی

قطعه‌بندی ارائه می‌شود، سپس مفاهیم و الگوریتم‌های استفاده‌شده در طرح به‌صورت کامل توضیح داده می‌شود. در گام بعدی با توجه به طرح ارائه‌شده دو الگوریتم توسعه می‌یابد. در مرحله بعد الگوریتم‌های موردنظر روی چندین مجموعه داده اجرا می‌گردد و با توجه به معیارهای مختلف مانند کارایی و زمان بهترین الگوریتم انتخاب می‌گردد. به‌منظور بررسی کارایی الگوریتم انتخاب‌شده نتایج آن با ۴ الگوریتم موجود در ادبیات با استفاده از تست آماری مقایسه می‌گردد. در انتها نیز نتیجه‌گیری و مطالعات پیش رو بیان می‌گردد.



شکل ۱: فلوچارت روش پژوهش

۲- خوشه‌بندی ترکیبی بر مبنای قطعه‌بندی

در مدل این پژوهش به‌جای این که یک مجموعه‌داده سری زمانی به‌صورت مستقیم خوشه‌بندی گردد، ابتدا این مجموعه‌داده به قطعات مساوی تبدیل شده و سپس هر قطعه ایجادشده به‌صورت جداگانه خوشه‌بندی می‌گردد. در پایان نیز با استفاده از الگوریتم‌های خوشه‌بندی ترکیبی نتایج به‌دست‌آمده از قطعات مختلف با یکدیگر ترکیب می‌گردند. الگوریتم دارای سه گام اصلی شامل قطعه‌بندی با اندازه ثابت، خوشه‌بندی هر یک از قطعات ایجادشده و ترکیب نتایج حاصل شده می‌باشد.

شکل ۲ یک مجموعه‌داده سری‌زمانی با دو خوشه را نمایش می‌دهد که با رنگ سبز و آبی مشخص شده است، مشاهده می‌شود که قطعات شماره ۱ و ۴ خوشه‌بندی را به‌خوبی نمایش نمی‌دهند درحالی که قطعات شماره ۲ و ۳ خوشه‌ها به‌خوبی متمایز شده‌اند. با توجه به شکل اگر بتوان در خوشه‌بندی مجموعه‌داده از قطعاتی که خوشه‌بندی را به درستی نمایش داده‌اند استفاده نمود نتایج خوشه‌بندی کل مجموعه‌داده بهبود خواهد داشت، به همین سبب استفاده از خوشه‌بندی قطعات و ترکیب آن‌ها می‌تواند نتایج دقیق‌تری داشته باشد.

۲-۲ خوشه‌بندی قطعات ایجادشده

الگوریتمی که به منظور خوشه‌بندی قطعات ایجادشده در گام اول مورد استفاده قرار گرفته است از نوع سلسله مراتبی تجمیعی می‌باشد در این الگوریتم از دو نوع فاصله دورترین (d_{max}) و میانگین (d_{mean}) برای فاصله بین خوشه‌ها به منظور پیوند استفاده شده است. روابط شماره (۱) و (۲) به ترتیب این دو نوع فاصله را نمایش می‌دهند. A و B نمایش‌دهنده خوشه‌ها و a و b اعضا هر خوشه می‌باشند همچنین d نشان‌دهنده فاصله بین دو عضو می‌باشد [۱۳].

$$(۱) \quad d_{max} = \max \{d(a,b) : a \in A, b \in B\}$$

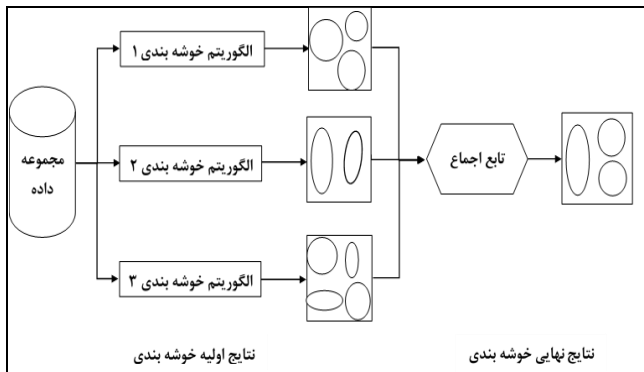
$$(۲) \quad d_{mean} = \frac{1}{|A||B|} \sum_{a \in A} \sum_{b \in B} d(a,b)$$

به منظور محاسبه فاصله بین دو سری زمانی می‌توان از معیارهای فاصله‌ای متفاوتی از قبیل فاصله اقلیدسی، DTW و بزرگ‌ترین زیر دنباله مشترک (LCSS) استفاده نمود [۳۶-۳۴]. با توجه به بالا بودن هزینه زمانی معیارهای DTW و LCSS در این پژوهش فاصله اقلیدسی استفاده شده است [۳۷-۳۹]. اگر دو سری زمانی $X = \{x_t\}_{t=1}^n$ و $Y = \{y_t\}_{t=1}^n$ با طول N موجود باشد در این صورت فاصله اقلیدسی (ED) با رابطه (۳) قابل محاسبه می‌باشد. البته باید توجه داشت که این فاصله صرفاً هنگامی می‌تواند مورد استفاده قرار بگیرد که دو سری زمانی دارای طول یکسان باشند [۳۴].

$$(۳) \quad ED(X, Y) = \sqrt{\sum_{t=1}^n (x_t - y_t)^2}$$

۳-۲ خوشه‌بندی ترکیبی

خوشه‌بندی ترکیبی یکی از موضوعات مطرح‌شده در حوزه خوشه‌بندی می‌باشد که به ندرت در حوزه سری‌های زمانی مورد توجه قرار گرفته است. از آنجایی که همه روش‌های خوشه‌بندی لزوماً برای یک مجموعه داده دارای نتایج مطلوب نمی‌باشد، لذا می‌توان به طور هم‌زمان از چند روش خوشه‌بندی استفاده نمود. سپس با استفاده از تکنیک‌های مشخصی خوشه‌بندی نهایی را انتخاب نمود. شکل ۴ یک نمایش کلی از این رویکرد را نشان می‌دهد [۴۰].



شکل ۴: خوشه‌بندی ترکیبی

خوشه‌بندی ترکیبی شامل دو قسمت تولید و ترکیب می‌باشد. در قسمت تولید می‌توان از دو رویکرد همگن و ناهمگن جهت تولید برچسب‌های متفاوت استفاده نمود. روش‌های ترکیب را می‌توان به چهار دسته اصلی شامل مستقیم، مبتنی بر خصوصیات، تشابهات زوجی و رویکردهای مبتنی بر گراف تقسیم‌بندی نمود. در این پژوهش به منظور تولید برچسب‌های مختلف از قطعه‌بندی و در بخش ترکیب نیز از یک الگوریتم مبتنی بر گراف با نام LWGP استفاده شده است. در این روش ابتدا فاصله بین خوشه‌های ایجادشده برای هر شیء محاسبه می‌گردد و سپس از این فواصل و استفاده از یک الگوریتم خوشه‌بندی به منظور برچسب نهایی هر یک از اشیا استفاده می‌گردد [۴۱].

۴-۲ معیارهای درونی

با توجه به اینکه در خوشه‌بندی از برچسب هر یک از اشیا اطلاعی در دست نیست، معیارهای درونی می‌توانند اهمیت ویژه‌ای داشته باشند. هدف از این معیارها، ارزیابی ساختار خوشه‌های ایجادشده بر اساس شباهت درون خوشه‌ای و تمایز بین خوشه‌های می‌باشد. جدول ۱ سه معیار مورد استفاده در این پژوهش را نمایش می‌دهد. در این روابط n و k به ترتیب نشان‌دهنده تعداد اشیا و خوشه‌ها می‌باشند، همچنین x و y نیز اشیا خاص داخل خوشه c را نمایش می‌دهند [۴۲].

جدول ۱: معیارهای درونی ارزیابی خوشه‌بندی

ردیف	معیار	نماد	رابطه ریاضی	مقدار مطلوب
۱	واریانس بین گروهی ^۱ [۱۳]	V	$\frac{1}{n-k} \sum_{i=1}^k \sum_{x \in c_i} dist(x, c_i)$	کمینه

^۱ Inter-group Variance

بیشینه	$\min_i \left\{ \min_j \left(\frac{\min_{x \in C_i, y \in C_j} (dist(x, y))}{\max_k \left\{ \max_{x, y \in C_k} (dist(x, y)) \right\}} \right) \right\}$	D	دان [۴۳] ^۱	۲
بیشینه	$\frac{n-k}{n-1} \frac{\sum_{i=1}^k n_i dist^2(x, c_i)}{\sum_{i=1}^k \sum_{x \in C_i} dist^2(x, c_i)}$	CH	کالینسکی- هارباز [۴۴] ^۲	۳

۵-۲ معیارهای بیرونی

k قطعه تبدیل شده است، به ترتیب تعداد ۲ تا k برچسب خوشه‌بندی نیز وجود خواهد داشت.

گام سوم: در این گام نتایج خوشه‌بندی موجود با استفاده از یک الگوریتم خوشه‌بندی ترکیبی با یکدیگر ترکیب می‌گردد تا یک نتیجه واحد برای هر مرحله قطعه‌بندی از ۲ تا k قطعه ایجاد گردد. در پایان این گام، تعداد $k-1$ برچسب خوشه‌بندی برای مجموعه داده سری زمانی وجود خواهد داشت. در انتهای این گام نتایج خوشه‌بندی مجموعه داده بدون قطعه‌بندی نیز به این تعداد اضافه می‌گردد و تعداد برچسب‌های موجود برابر با k خواهد بود.

گام چهارم: در ادامه با استفاده از یک معیار داخلی مشخص تعداد m برچسب از میان k برچسب ایجاد شده، انتخاب می‌شوند ($m \leq k$).

گام پنجم: در گام نهایی m برچسب حاصل شده مجدداً با استفاده از یک الگوریتم خوشه‌بندی ترکیبی با یکدیگر ترکیب می‌گردند و یک برچسب خوشه‌بندی نهایی حاصل می‌گردد.

انتخاب مقدار k و m می‌تواند تا حدودی پیچیده باشد اما به نظر می‌رسد می‌توان مقدار k را طوری انتخاب نمود که برای سری‌های زمانی طولانی، بزرگ نباشد و همچنین برای سری‌های زمانی با طول کوتاه نیز بتواند ساختار سری زمانی را برای قطعه‌ها حفظ نماید. با توجه به این موارد به نظر می‌رسد استفاده از یک رابطه لگاریتمی بتواند یک مقدار k مناسب را برای استفاده در الگوریتم ایجاد نماید. بنابراین رابطه (۵) برای نحوه محاسبه مقدار k در این پژوهش به کار گرفته شده است. در این رابطه، L برابر با طول سری زمانی می‌باشد. با توجه به نتایج عددی به دست آمده برای پارامتر m ، مقادیر عددی ۲ یا ۳ مناسب به نظر می‌رسد.

$$k = \log_{10} L \quad (5)$$

معیارهای بیرونی با فرض مشخص بودن برچسب اشیاء به منظور بررسی دقت الگوریتم مورد استفاده قرار می‌گیرند. در این مطالعه به منظور مقایسه با پژوهش‌های انجام گرفته از معیار بیرونی رند (RI) استفاده شده است که رابطه (۴) آن را نمایش می‌دهد.

$$RI = \frac{TP + TN}{TP + TN + FN + FP} \quad (4)$$

در رابطه (۴) TP^3 نشان‌دهنده تعداد اشیاء می‌باشد که در کلاس و خوشه یکسانی قرار دارند. TN^4 تعداد اشیاء هستند که در خوشه‌های متفاوت قرار دارند و کلاس آن‌ها نیز متفاوت می‌باشد. FP^5 نشان‌دهنده اشیایی می‌باشد که دارای خوشه‌بندی متفاوت می‌باشند در حالی که کلاس آن‌ها یکسان باشد و سرانجام FN^6 نشان‌دهنده تعداد اشیایی است که دارای خوشه یکسان و کلاس‌های متفاوت می‌باشند.

۶-۲ الگوریتم‌های پیشنهادی

در این پژوهش دو الگوریتم بر پایه قطعه‌بندی با اندازه ثابت و خوشه‌بندی ترکیبی به منظور خوشه‌بندی سری‌های زمانی ارائه شده است.

۶-۲-۱ الگوریتم اول توسعه داده شده

این الگوریتم دارای ۵ گام به شرح زیر می‌باشد.

گام اول: یک مجموعه داده سری زمانی به قطعات با طول مساوی تقسیم می‌گردد. تعداد قطعات بین ۲ تا k قطعه می‌باشد.

گام دوم: برای هر یک از قطعات ایجاد شده و برای هر مقدار k با استفاده از الگوریتم خوشه‌بندی سلسله مراتبی، خوشه‌بندی انجام می‌گیرد. در این گام با توجه به اینکه مجموعه داده سری زمانی به ۲ تا

^۱ Dunns Indices

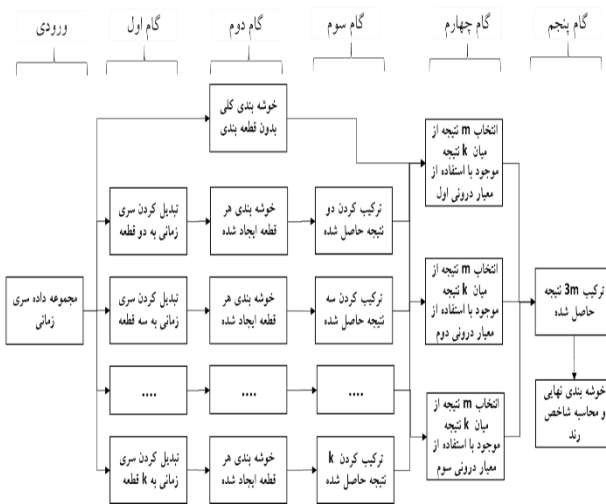
^۲ Calinski-Harabasz

^۳ True positive

^۴ True negative

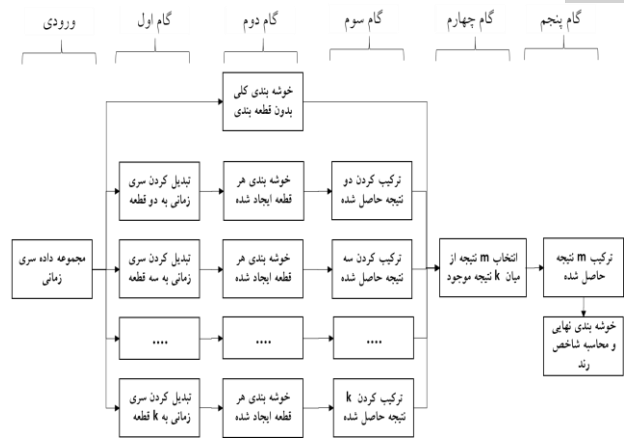
^۵ False positive

^۶ False negative



شکل ۶: الگوریتم دوم پیشنهادی

برای انتخاب m برچسب از k برچسب نهایی برای هر سری زمانی، از معیارهای درونی ارزیابی خوشه‌بندی می‌توان استفاده نمود. البته از آنجایی که این معیارها نیز متنوع می‌باشند، می‌توان از سه معیار واریانس بین‌گروهی، کالینسکی-هارباز و شاخص دان که در ادبیات موضوع بیشتر مورد استفاده قرار گرفته‌اند، استفاده نمود. با توجه به اینکه استفاده از واریانس بین گروهی در الگوریتم اول برای نتیجه‌گیری کفایت می‌کند از معیارهای دیگر استفاده نشده است. در الگوریتم اول صرفاً از یک معیار برای انتخاب m برچسب از k برچسب ایجاد شده، استفاده می‌گردد. الگوریتم اول پیشنهادی در شکل ۵ نمایش داده شده است.



شکل ۵: الگوریتم اول پیشنهادی

۳- اجرای مدل

به منظور بررسی کارایی الگوریتم‌های ارائه شده، این الگوریتم‌ها روی ۸۲ مجموعه داده از آرشیو UCR اجرا شده است [۴۵]. مجموعه داده مورد نظر زمینه‌های مختلفی مانند سلامت، مالی، رسانه و مهندسی را پوشش می‌دهد. با توجه به اینکه مجموعه داده‌های مورد نظر برای کلاس‌بندی سری‌های زمانی ارائه شده‌اند و دارای دو بخش آموزش و آزمون می‌باشند، برای خوشه‌بندی با یکدیگر جمع شده‌اند. پارامترهای ورودی برای هر الگوریتم شامل الگوریتم خوشه‌بندی استفاده شده، نوع پیوند، تعداد معیارهای درونی استفاده شده، معیار درونی استفاده شده، m و k می‌باشد. همچنین خروجی مدل نیز شامل سری‌های زمانی قرار گرفته در هر خوشه و همچنین شاخص رند نهایی می‌باشد.

شبه کد الگوریتم اول و دوم توسعه داده شده شامل ورودی‌ها، خروجی‌ها و روند اجرای هر الگوریتم توسط شکل ۷ نمایش داده شده است.

First Algorithm	Second Algorithm
Input: Time series dataset, m , k , internal measure	Input: Time series dataset, m , k , internal measure
Output: Best quality clustering, RI	Output: Best quality clustering, RI
Clustering dataset without segmentation	Clustering dataset without segmentation
for each k do	for each k do
Apply time series segmentation	Apply time series segmentation
for each segment do	for each segment do
Cluster the segment	Cluster the segment
end for	end for
Ensemble clustering all clustered segments	Ensemble clustering all clustered segments
end for	end for
for 3 internal measure do	for 3 internal measure do
Select m results from k results using an internal measure	Select m results from k results using an internal measure
Ensemble clustering m selected results	Ensemble clustering m selected results
Evaluate the goodness of the clustering	Evaluate the goodness of the clustering
Return Best quality clustering, RI	Return Best quality clustering, RI

شکل ۷: شبه کد الگوریتم اول و دوم

۲-۶ الگوریتم دوم توسعه داده شده

الگوریتم دوم نیز مانند الگوریتم اول دارای ۵ گام می‌باشد که سه گام نخست این الگوریتم نیز مانند الگوریتم اول می‌باشد. تفاوت این الگوریتم در گام سوم ایجاد می‌گردد که به نحوه انتخاب m نتیجه مرتبط می‌باشد. در گام چهارم این الگوریتم از سه معیار درونی معرفی شده در بخش ۲-۶-۱ به طور هم‌زمان به منظور انتخاب m نتیجه استفاده می‌گردد. در مرحله پنجم و نهایی نیز مطابق الگوریتم قبلی تعداد $3m$ برچسب نهایی ایجاد شده با استفاده از خوشه‌بندی ترکیبی با یکدیگر ترکیب می‌گردند. در مرحله نهایی در الگوریتم اول ما دارای m برچسب نهایی خواهیم بود اما در این الگوریتم تعداد برچسب‌ها $3m$ می‌باشد. باید توجه داشت که تعدادی از برچسب‌ها می‌تواند تکراری باشد. الگوریتم دوم در شکل ۶ نمایش داده شده است.

دوم نمایش می‌دهد. در تمام وضعیت‌ها به منظور خوشه‌بندی قطعات از الگوریتم سلسله مراتبی استفاده شده است.

برای الگوریتم اول سه وضعیت و برای الگوریتم دوم یک وضعیت با توجه به پارامترهای مختلف در نظر گرفته شده است. جدول ۲ پارامترهای استفاده شده برای سه وضعیت الگوریتم اول و الگوریتم

جدول ۲: پارامترهای الگوریتم‌های مختلف خوشه‌بندی

مقدار پارامتر m	معیار درونی	تعداد معیار درونی	نوع پیوند	الگوریتم خوشه‌بندی	نوع
۱	واریانس بین گروهی	۱	میانگین	سلسله‌مراتبی	وضعیت اول
۳	واریانس بین گروهی	۱	میانگین	سلسله‌مراتبی	وضعیت دوم
۳	واریانس بین گروهی	۱	حداکثر	سلسله‌مراتبی	وضعیت سوم
	واریانس بین گروهی، کالینسکی -				
۲	هارباسز، دان	۳	حداکثر	سلسله‌مراتبی	الگوریتم دوم

۱-۳ نتایج

رند و زمان اجرا برای میانگین ۱۰ تکرار در جدول ۳ گزارش شده است. مجموعه داده‌ها با سایز بزرگ، زمان‌های اجرای طولانی در مقایسه با سایر مجموعه داده‌ها دارند.

به منظور بررسی الگوریتم‌های ارائه شده نتایج حاصل روی ۸۲ مجموعه داده معرفی شده مورد مقایسه قرار گرفته است. الگوریتم‌ها برای هر یک از مجموعه داده‌ها ۱۰ بار اجرا شده است و نتایج شاخص

جدول ۳: نتایج شاخص رند و زمان اجرا برای دو الگوریتم ارائه شده

الگوریتم دوم	زمان اجرا (ثانیه)		شاخص رند					مجموعه داده
	الگوریتم اول وضعیت سوم	الگوریتم اول وضعیت دوم	الگوریتم اول وضعیت اول	الگوریتم دوم	الگوریتم اول وضعیت سوم	الگوریتم اول وضعیت دوم	الگوریتم اول وضعیت اول	
۷,۴	۸,۳۴	۴,۱۳	۳,۳۷	۵۰,۸۰٪	۵۰,۱۰٪	۵۱,۳۰٪	۵۰,۰۰٪	ITA
۵,۶۲	۳,۵۷	۲,۰۸	۱,۹۴	۸۲,۰۰٪	۸۱,۳۰٪	۷۹,۰۰٪	۸۰,۵۰٪	SYN
۹,۹	۱۰,۵۹	۶,۲۲	۴,۹۵	۶۳,۸۰٪	۶۷,۵۰٪	۵۴,۳۰٪	۵۴,۳۰٪	SO ₂
۴,۶۶	۳,۹۶	۲,۶۴	۲,۱۷	۵۰,۱۰٪	۵۲,۶۰٪	۵۰,۱۰٪	۴۹,۹۰٪	SO ₁
۳,۶۴	۲,۸۱	۱,۸۳	۱,۷۱	۷۲,۹۰٪	۷۳,۷۰٪	۷۱,۲۰٪	۷۱,۲۰٪	DPA
۷,۱۳	۷,۸۶	۴,۹۲	۴,۳۱	۵۳,۵۰٪	۵۲,۴۰٪	۵۳,۰۰٪	۵۳,۰۰٪	DPC
۳,۳۲	۳,۱۸	۱,۸۷	۱,۵۸	۸۷,۸۰٪	۸۶,۳۰٪	۸۱,۵۰٪	۸۴,۳۰٪	DPT
۴,۹۲	۳,۱	۲,۳۳	۱,۷۴	۷۳,۰۰٪	۷۲,۸۰٪	۷۰,۵۰٪	۷۰,۴۰٪	MPA
۸,۷۳	۷,۴۹	۵,۴۱	۴,۴۱	۵۰,۰۰٪	۵۰,۰۰٪	۵۳,۲۰٪	۵۳,۳۰٪	MPC
۴,۸۷	۲,۹۸	۱,۷۲	۱,۵	۸۳,۸۰٪	۸۳,۹۰٪	۸۱,۸۰٪	۸۲,۳۰٪	MPT
۷۷,۶۲	۸۴,۶۹	۴۹,۴۵	۴۵,۱۶	۵۴,۲۰٪	۵۲,۱۰٪	۵۴,۰۰٪	۵۴,۰۰٪	PHA
۵,۷۷	۳,۵۷	۲,۱۹	۲,۲۶	۷۸,۹۰٪	۷۹,۷۰٪	۷۷,۳۰٪	۷۷,۴۰٪	PPA
۸,۱۹	۷,۳۷	۴,۹۶	۴,۳۷	۵۲,۱۰٪	۵۷,۶۰٪	۵۳,۸۰٪	۵۳,۹۰٪	PPC
۵,۵۲	۳,۷۵	۲,۰۵	۱,۸۶	۷۸,۶۰٪	۷۹,۸۰٪	۸۳,۸۰٪	۸۳,۹۰٪	PPT
۱۴,۰۳	۱۴,۴۴	۸,۰۴	۷,۰۵	۵۰,۳۰٪	۵۰,۱۰٪	۵۰,۰۰٪	۵۰,۱۰٪	TWE
۱۷,۴۸	۲۰,۹۹	۱۰,۸۲	۹,۷۱	۶۱,۸۰٪	۶۳,۱۰٪	۵۰,۳۰٪	۵۰,۳۰٪	MOT
۱,۰۵	۱,۱۳	۰,۷۵	۰,۶۶	۶۲,۳۰٪	۶۱,۲۰٪	۵۴,۰۰٪	۶۲,۳۰٪	EC ₂
۲۱,۵۵	۱۵,۵	۱۰,۸۷	۱۰,۱۱	۶۳,۴۰٪	۶۴,۶۰٪	۶۱,۳۰٪	۶۰,۲۰٪	MED

10,94	10,55	6,9	6,25	66,00%	64,70%	63,10%	60,60%	CBF
23,42	18,67	13,6	12,06	87,20%	88,70%	37,90%	39,00%	SWE
361,46	383,85	258,36	240,73	62,80%	62,90%	62,10%	62,60%	TWP
63,1	62,1	75,5	39,63	83,30%	85,00%	81,60%	73,70%	FAA
10,55	10,77	7,95	6,53	50,00%	50,50%	50,30%	49,90%	ECF
328,33	463,83	278,71	250,06	74,10%	74,60%	90,60%	89,60%	EC0
2,99	1,87	1,05	1,05	94,00%	94,00%	92,50%	92,80%	PLA
1,6	1,27	1,98	0,66	53,70%	49,70%	50,90%	49,70%	GUN
1571,76	1395,94	618,98	580,69	53,40%	53,40%	53,40%	53,40%	WAF
281,72	290,34	208,31	210,29	53,10%	52,20%	51,50%	50,50%	CHL
11,56	11,82	8,25	7,57	92,20%	92,80%	48,70%	47,50%	ADI
99	98	55	50	66,87%	67,15%	62,52%	62,43%	كوتاه
5,37	3,82	2,52	0,93	49,90%	49,80%	49,90%	49,90%	WIN
24,01	22,91	15,13	13,68	50,40%	50,40%	51,50%	51,50%	STR
1,89	1,87	1,57	1,17	60,50%	66,80%	35,30%	35,80%	ARR
118,87	96,6	61,4	60,32	86,80%	86,00%	86,80%	86,50%	INS
32,85	22,63	9,42	10,23	94,50%	94,60%	94,40%	94,60%	W00
54,08	19,92	14,03	11,86	88,60%	88,70%	87,60%	83,90%	WOS
3,24	1,92	3,16	1,18	75,10%	75,00%	74,90%	75,00%	TRA
3,06	2,31	1,53	1,36	50,00%	49,90%	50,10%	49,90%	TO1
0,86	0,96	1,07	0,65	79,50%	70,10%	49,10%	49,80%	COF
14,79	13,29	9,4	9,48	83,40%	81,60%	77,90%	71,20%	CRX
17,26	13,15	9,92	9,32	84,00%	82,40%	77,40%	77,40%	CRY
15,7	13,48	9,35	8,48	84,30%	84,60%	77,10%	69,60%	CRZ
501,15	471,96	336,93	326,98	81,40%	81,90%	73,60%	71,50%	UWX
509,23	447,96	323,17	303,87	82,10%	82,20%	79,80%	80,50%	UWY
589,53	460,55	330,65	309,96	81,90%	81,90%	83,30%	79,90%	UWZ
2,78	2,14	1,35	1,15	75,60%	71,40%	75,80%	74,10%	LIV
2,02	1,59	1,05	0,92	51,10%	49,70%	49,90%	49,70%	TO2
4,15	3,19	2,57	2,39	30,60%	29,60%	30,60%	30,60%	DIA
1,6	1,44	1,5	1,02	69,50%	69,20%	73,80%	72,10%	FAF
40,45	27,47	18,86	19,63	89,20%	89,60%	89,90%	89,90%	SYM
481,22	420,37	295,01	272,77	50,00%	50,10%	50,00%	50,00%	YOG
12,04	6,87	5,35	4,89	73,20%	71,40%	62,70%	67,80%	OSU
7,13	2,49	2,05	1,68	51,40%	51,50%	50,20%	50,10%	HAM
2,62	1,66	1,23	1,12	72,00%	72,00%	77,00%	72,80%	MEA
6,44	5,48	4,29	3,78	71,40%	72,00%	71,60%	64,70%	FIS
1,52	3,59	1,15	1,04	63,80%	63,80%	65,50%	65,60%	BEE
1119,49	992,98	1048,86	796,02	50,00%	50,00%	50,20%	50,20%	FOA
958,58	807,43	694,94	659,74	50,00%	50,00%	50,00%	50,00%	FOB

۳,۱۴	۲,۵	۲,۰۷	۱,۶۷	۴۹,۸۰٪	۵۰,۰۰٪	۵۱,۴۰٪	۵۱,۴۰٪	SHS
۱۵۶	۱۳۴	۱۱۱	۹۸	۶۸,۲۸٪	۶۷,۸۰٪	۶۵,۴۲٪	۶۴,۳۴٪	متوسط
۱,۱۶	۱,۲۷	۰,۸۲	۰,۸۷	۵۰,۹۰٪	۵۱,۷۰٪	۵۶,۹۰٪	۵۵,۴۰٪	BFL
۱,۲	۱,۵۸	۰,۷۷	۰,۷۶	۵۳,۳۰٪	۴۸,۸۰٪	۴۹,۲۰٪	۴۸,۹۰٪	BIR
۱۰,۰۰۶	۹,۶	۷,۵۶	۶,۹۴	۵۲,۳۰٪	۵۲,۴۰٪	۵۶,۸۰٪	۵۶,۸۰٪	EAR
۱,۹۷	۱,۶۳	۱,۱۹	۱,۴۴	۵۰,۵۰٪	۵۰,۶۰٪	۵۱,۷۰٪	۵۱,۷۰٪	HER
۶۸,۸	۵۳,۳	۳۶,۹	۳۶,۲۱	۹۴,۸۰٪	۹۶,۱۰٪	۸۴,۴۰٪	۸۳,۵۰٪	SHA
۲,۳۶	۱,۶۴	۱,۰۶	۱,۱۴	۸۵,۸۰٪	۸۶,۰۰٪	۸۶,۱۰٪	۸۸,۷۰٪	OLI
۱,۸۵	۲,۳۳	۱,۴	۱,۲۶	۶۱,۳۰٪	۶۳,۳۰٪	۵۹,۵۰٪	۶۵,۳۰٪	CAR
۲,۱۸	۱,۷۷	۱,۳۷	۱,۲۳	۶۰,۰۰٪	۵۰,۴۰٪	۵۸,۵۰٪	۵۷,۹۰٪	LI۲
۱۴,۲	۱۲,۴۷	۱۰,۴۴	۱۰,۵۲	۵۰,۴۰٪	۵۰,۳۰٪	۴۹,۹۰٪	۴۹,۹۰٪	COM
۳۸,۶۹	۲۷,۳۷	۲۲,۰۱	۲۱,۱۹	۵۲,۲۰٪	۵۰,۹۰٪	۵۴,۲۰٪	۵۴,۱۰٪	LAR
۳۵,۷۸	۲۷,۱۸	۲۲,۸۴	۲۱,۶	۵۵,۲۰٪	۵۴,۹۰٪	۴۵,۵۰٪	۴۵,۸۰٪	REF
۳۷,۳۵	۲۹,۱۶	۲۲,۶۴	۲۲,۸	۵۵,۲۰٪	۵۵,۰۰٪	۴۰,۲۰٪	۴۰,۳۰٪	SCR
۳۸,۹	۳۱,۰۱	۲۳,۲۷	۲۲,۶۲	۴۳,۰۰٪	۴۶,۶۰٪	۳۴,۳۰٪	۳۳,۸۰٪	SMA
۹۹۲,۱۸	۶۳۰,۷۴	۵۵۹,۳۷	۵۳۲,۸۱	۹۴,۶۰٪	۹۵,۱۰٪	۸۵,۴۰٪	۸۶,۷۰٪	NO۱
۹۸۰,۵۲	۶۴۴,۱۹	۴۹۱,۲۶	۵۰۴,۸۸	۹۵,۶۰٪	۹۵,۹۰٪	۹۱,۱۰٪	۹۱,۳۰٪	NO۲
۱۲,۹	۶,۰۸	۵,۱۵	۵,۷۱	۶۳,۵۰٪	۶۳,۰۰٪	۵۹,۶۰٪	۶۲,۷۰٪	WOR
۸,۷	۵,۹۸	۴,۱۱	۴,۷۹	۴۹,۸۰٪	۴۹,۳۰٪	۴۹,۹۰٪	۴۹,۷۰٪	WOT
۱۷۵۳,۹۲	۱۱۸۰,۰۲	۹۰۷,۴۹	۸۶۵,۹۸	۸۵,۶۰٪	۸۶,۳۰٪	۸۵,۹۰٪	۸۷,۶۰٪	UWA
۵۰۹,۰۰۶	۳۳۶,۲۱	۲۵۳,۱۷	۲۴۰	۹۳,۹۰٪	۹۲,۲۰٪	۹۰,۲۰٪	۹۰,۳۰٪	MAL
۴۳۲,۷۵	۳۳۴,۱۲	۲۷۷,۲۸	۲۸۱,۷	۹۲,۶۰٪	۹۲,۰۰٪	۹۰,۱۰٪	۹۰,۱۰٪	PHO
۲۵,۶۴	۱۶,۵۱	۱۴,۸۷	۱۵,۷۹	۶۷,۱۰٪	۶۷,۱۰٪	۳۹,۴۰٪	۳۶,۶۰٪	HAP
۲۶۲,۵۶	۱۹۰,۵۷	۱۷۰,۳۹	۱۶۷,۹۳	۶۳,۸۰٪	۶۶,۱۰٪	۶۱,۸۰٪	۶۱,۴۰٪	CIN
۸۱,۸۸	۵۳,۸۵	۴۳,۱۸	۴۱,۹۷	۶۸,۰۰٪	۶۶,۳۰٪	۶۴,۶۰٪	۶۱,۷۰٪	INL
۵۵۲,۰۱	۴۳۹,۸۲	۳۶۳,۶۹	۳۴۴,۶	۵۶,۱۰٪	۵۰,۰۰٪	۵۴,۵۰٪	۵۷,۱۰٪	HAN
۲۴۴	۱۶۸	۱۳۵	۱۳۱	۶۶,۴۸٪	۶۵,۸۵٪	۶۲,۴۹٪	۶۲,۸۰٪	بلند
۱۶۱	۱۳۱,۱۳	۹۸,۲۱	۹۰,۹۳	۶۷,۲٪	۶۷,۰۰٪	۶۳,۵۴٪	۶۳,۲۲٪	متوسط

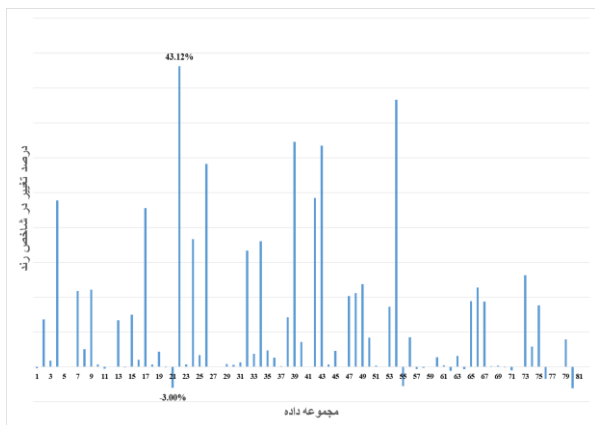
جدول ۴ به صورت مختصر و با توجه به معیارهای مختلف عملکرد دو الگوریتم ارائه شده را نمایش می‌دهد منظور از الگوریتم اولیه در این جدول، خوشه‌بندی بدون قطعه‌بندی با استفاده از الگوریتم سلسه‌مراتبی می‌باشد.

جدول ۴: نتایج الگوریتم‌های مختلف خوشه‌بندی

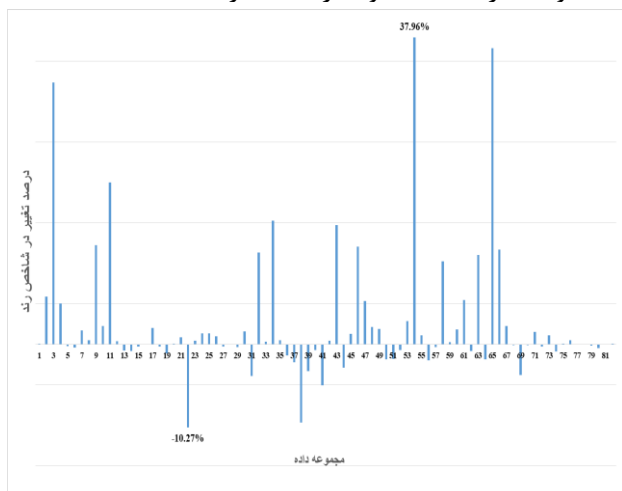
الگوریتم دوم	الگوریتم اول		معیار
	وضعیت سوم	وضعیت دوم	
۶۷,۲۵٪	۶۷,۰۰٪	٪۶۳,۵۴	میانگین شاخص رند
۰,۰۰۰۶۵	۰,۰۰۰۹۱	۰,۰۰۱۰۴	انحراف از معیار
۲,۹۲٪	۲,۶۹٪	٪۵,۷۴	میانگین بهبود شاخص رند نسبت به حالت اولیه
۳۷,۱۳٪	۳۷,۹۶٪	٪۴۳,۱۲	حداکثر مقدار بهبود در یک مجموعه داده
-۱۰,۷۸٪	-۱۰,۲۷٪	-۳,۰۰٪	حداکثر میزان تنزل شاخص رند در یک مجموعه داده
۵۹٪	۵۵٪	٪۷۳	نسبت مجموعه داده با بهبود شاخص رند
۳۰٪	۱۰٪	٪۱۰	نسبت مجموعه داده بدون تغییر شاخص رند

نسبت مجموعه داده با افت شاخص رند
نسبت ایجاد بهترین پاسخ در پژوهش های قبلی

۹٪
۱۸٪
۳۵٪
۲۰٪
۱۱٪
۳۳٪

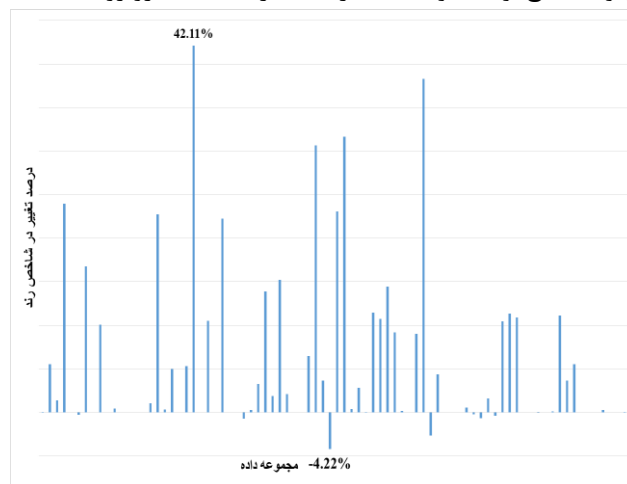


شکل ۹: اختلاف میانگین شاخص رند الگوریتم اول وضعیت دوم وضعیت سوم الگوریتم اول با پارامتر $m = 3$ و فاصله حداکثر برای الگوریتم خوشه بندی سلسله مراتبی اجرا شده است. با توجه به اطلاعات جدول ۴ می توان مشاهده نمود که میانگین شاخص رند در ۸۲ مجموعه داده به ۶۷ درصد رسیده است و نتایج اولیه خوشه بندی در این وضعیت ۲,۶۹ درصد به طور میانگین افزایش داشته است. همچنین در ۳۰ درصد حالات، وضعیت سوم الگوریتم اول توانسته است بهترین شاخص رند را در بین ادبیات موضوع داشته باشد. بر اساس اطلاعات شکل ۸، وضعیت سوم الگوریتم اول در بهترین و بدترین حالت به ترتیب توانسته است ۳۸ درصد نتیجه اولیه را بهبود بخشد و ۱۰ درصد نتیجه اولیه را ضعیف تر نماید.



شکل ۱۰: اختلاف میانگین شاخص رند الگوریتم اول وضعیت سوم در انتها الگوریتم دوم با پارامتر $m = 2$ و فاصله حداکثر برای الگوریتم خوشه بندی سلسله مراتبی اجرا شده است. در الگوریتم دوم همان گونه که قبلاً بیان شده است از سه معیار داخلی به صورت هم زمان استفاده شده است. با توجه به اطلاعات جدول ۴ الگوریتم

با توجه به اطلاعات جدول ۴ می توان مشاهده نمود که در الگوریتم اول وضعیت اول شاخص رند ۵,۴۴ درصد نسبت به الگوریتم اولیه بهبود پیدا نموده است. همین طور در ۳۰ درصد حالات، الگوریتم مورد اشاره توانسته است بهترین شاخص رند را در بین ادبیات موضوع داشته باشد. همچنین الگوریتم اول وضعیت اول توانسته است در ۵۹ درصد مجموعه داده ها دقت خوشه بندی را افزایش دهد و در ۳۲ درصد نیز دقت خوشه بندی تغییری نکرده است. با توجه به شکل ۸ در بهترین حالت، الگوریتم اول وضعیت اول می تواند تا ۴۲ درصد شاخص رند را برای یک مجموعه داده بهبود دهد و از سوی دیگر می تواند نتایج اولیه خوشه بندی را تا ۴ درصد با افت روبرو نماید.



شکل ۸: اختلاف میانگین شاخص رند الگوریتم اول وضعیت دوم وضعیت اول با مقدار پارامتر $m = 3$ اجرا شده است و مانند حالت قبل برای الگوریتم خوشه بندی سلسله مراتبی، فاصله برابر با میانگین در نظر گرفته شده است. با توجه به اطلاعات جدول ۴ می توان مشاهده نمود که وضعیت دوم الگوریتم اول میانگین شاخص رند برابر با ۶۳,۵۴ درصد می باشد و الگوریتم توانسته است دقت را در حدود ۵,۷۴ درصد بهبود دهد. مانند حالت قبل وضعیت دوم الگوریتم اول در ۲۰ درصد موارد بهترین جواب را داشته است و در ۷۳ درصد موارد نیز جواب خوشه بندی ابتدایی را بهبود داده است، همچنین در ۱۰ درصد مجموعه داده نیز این مقدار بدون تغییر باقی مانده است. بر اساس اطلاعات شکل ۷ در وضعیت دوم الگوریتم اول در بهترین حالت خود توانسته است تا ۴۳ درصد نتایج اولیه خوشه بندی را برای یک مجموعه داده خاص افزایش دهد و از طرفی در بدترین حالت نیز باعث کاهش ۳ درصدی شاخص رند برای یک مجموعه داده شده است.

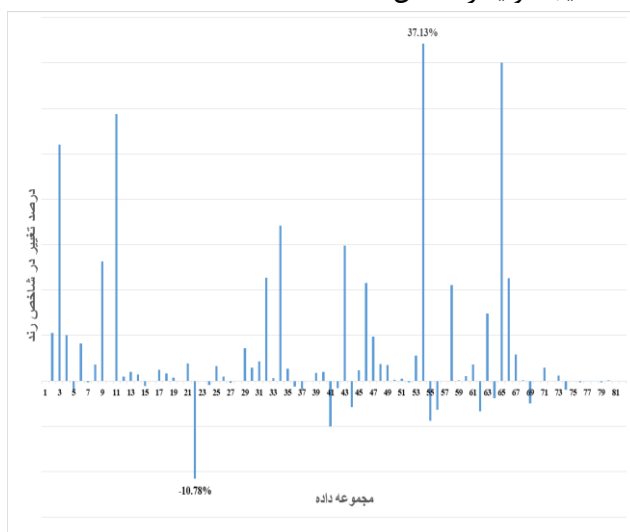
مقدار P-value در جدول ۵ می‌توان مشاهده نمود که شاخص رند برای الگوریتم اول وضعیت دوم از حالت اول با توجه به مقدار آلفای ۱۰ درصد بهتر می‌باشد. الگوریتم اول وضعیت سوم با توجه به مقدار آلفای ۵ درصد از هر دو الگوریتم اول وضعیت اول و دوم با توجه به آزمون بهتر می‌باشد و دارای شاخص رند بهتری می‌باشد. الگوریتم دوم با توجه به مقدار P-value به دست آمده و در سطح آلفای ۵ درصد از حالت اول و دوم الگوریتم اول بهتر می‌باشد اما در مورد مقایسه الگوریتم دوم و الگوریتم اول وضعیت سوم با توجه به مقدار p-value که برابر با ۰,۲۸۳ می‌باشد از لحاظ آماری مساوی بودن شاخص رند برای این دو وضعیت را نمی‌توان رد کرد. اما باید توجه داشت که میانگین شاخص رند برای الگوریتم دوم برابر با ۶۷,۲۵ و برای الگوریتم اول وضعیت سوم برابر با ۶۷ می‌باشد همچنین با توجه به دیگر معیارها در جدول ۴ می‌توان مشاهده نمود که الگوریتم دوم می‌تواند از الگوریتم اول وضعیت سوم عملکرد بهتری را داشته باشد.

جدول ۵: مقادیر P-value الگوریتم‌های ارائه شده

الگوریتم اول - وضعیت اول	الگوریتم اول - وضعیت دوم	الگوریتم اول - وضعیت سوم
۰,۰۸۸	-	-
۰,۰۰۱	۰,۰۰۲	-
۰,۰۰۰	۰,۰۰۰	۰,۲۸۳

به منظور بررسی دقت و سرعت الگوریتم‌های ارائه شده، بهترین نتایج که مربوط به الگوریتم دوم می‌باشد، با چهار الگوریتم در ادبیات موضوع بررسی شده است. چهار الگوریتم مذکور عبارت‌اند از DD_{DTW} که در آن از ترکیب دو معیار فاصله به همراه الگوریتم سلسله مراتبی استفاده شده است [۱۳]، الگوریتم KSC که در آن از گویهای زمانی مشخص به منظور خوشه‌بندی استفاده شده است [۴۷] و الگوریتم‌های TS^3C_{CH} و TS^3C_{MV} که بر اساس مشخصات قطعه‌بندی سری‌های زمانی توسعه داده شده‌اند [۲۲]. جدول ۶ شاخص رند و زمان اجرای مربوط به چهار الگوریتم منتخب و همچنین الگوریتم دوم را نمایش می‌دهد. با توجه به اطلاعات جدول ۶ و جدول ۷ می‌توان مشاهده نمود که با توجه به مقدار P-value و آلفای ۵ درصد، الگوریتم ارائه شده عملکرد بهتری نسبت به این الگوریتم‌ها داشته است. در مورد مقایسه هزینه زمانی الگوریتم ارائه شده نیز با توجه به اطلاعات جدول ۶ می‌توان مشاهده نمود که الگوریتم ارائه شده دارای زمان اجرای بسیار کمتری می‌باشد. نتایج جدول ۷ نیز این ادعا را تأیید می‌نماید. همچنین با توجه به اطلاعات جدول ۶ مشاهده می‌شود که الگوریتم ارائه شده در ۳۰ مورد از ۸۲ مجموعه داده‌ها توانسته است بهترین نتایج بین چهار الگوریتم را داشته باشد.

دوم توانسته است ۲,۹ درصد به صورت میانگین دقت خوشه‌بندی را افزایش دهد دقت شاخص رند را به ۶۷,۲۵ درصد برساند که بهترین نتیجه ممکن می‌باشد. در الگوریتم دوم نیز نتایج مربوط به ۶۰ درصد مجموعه داده‌ها به طور میانگین بهبود داشته است و در ۳۰ درصد موارد نیز تغییر نداشته است. همین‌طور در ۳۲ درصد حالات، الگوریتم دوم توانسته است بهترین شاخص رند را در بین ادبیات موضوع داشته باشد. بر اساس اطلاعات شکل ۹ الگوریتم دوم در بهترین حالت توانسته است ۳۷ درصد نتیجه شاخص رند را برای یک مجموعه داده بهبود ببخشد و در بدترین حالت نیز نزدیک به ۱۰ درصد نتیجه اولیه را کاهش دهد.



شکل ۱۱: اختلاف میانگین شاخص رند الگوریتم دوم

همچنین به منظور تحلیل دقیق‌تر نتایج به دست آمده مجموعه داده‌های مورد استفاده در سه کلاس با طول کوتاه (کمتر از ۲۰۰)، طول متوسط (بین ۲۰۰ تا ۵۰۰) و طول بلند (بیشتر از ۵۰۰) تقسیم شده‌اند. با توجه به نتایج جدول ۳ می‌توان مشاهده نمود که در حالت طول کوتاه الگوریتم اول - حالت سوم می‌تواند نتایج بهتری را داشته باشد با این وجود در دو کلاس متوسط و بلند الگوریتم دوم ارائه شده توانسته است نتایج بهتری را در برداشته باشد.

با توجه به اطلاعات جدول ۴ می‌توان مشاهده نمود که انحراف از معیار ۱۰ تکرار برای هر دو الگوریتم ارائه شده بسیار ناچیز می‌باشد که این نشان دهند پایداری خوب طرح ارائه شده می‌باشد.

۲-۳ مقایسه الگوریتم‌ها

در این قسمت آزمون آماری به منظور بررسی بهترین الگوریتم در دو مرحله انجام گرفته است در مرحله اول الگوریتم‌های ارائه شده با یکدیگر مقایسه شده‌اند و در گام بعدی بهترین الگوریتم ارائه شده با چهار الگوریتم ادبیات در این حوزه مقایسه آماری شده است. یکی از روش‌های پیشنهاد شده به منظور ارزیابی الگوریتم‌های خوشه‌بندی و طبقه‌بندی، استفاده از آزمون ویلکانسون می‌باشد [۴۶]. با توجه به

جدول ۶: مقایسه شاخص رند و زمان اجرای الگوریتم دوم با چهار الگوریتم منتخب

KSC	زمان اجرا				الگوریتم دوم	شاخص رند				مجموعه داده
	DD_{DTW}	TS^3C_{MV}	TS^3C_{CH}	الگوریتم دوم		KSC	DD_{DTW}	TS^3C_{MV}	TS^3C_{CH}	
۱۷۶۶۷	۱۰۴۲۸۵	۱۵۰۳	۱۴۷۸	۳۲,۸۵	۶۶٪	۹۲٪	۹۴٪	۹۴٪	۹۵٪	۵۰W
۵۳۷۴	۳۲۲۲۱	۹۰۴	۸۸۷	۱۱,۵۶	۹۵٪	۶۸٪	۹۲٪	۹۲٪	۹۲٪	ADI
۸۴	۴۶۰۳	۳۶۲	۳۶۱	۱,۸۹	۶۳٪	۳۵٪	۶۳٪	۶۲٪	۶۰٪	ARR
۶۳	۱۲۴۷	۱۸۸	۱۸۸	۱,۵۲	۷۱٪	۵۸٪	۶۸٪	۶۸٪	۶۴٪	BEE
۲۱	۶۴۲	۱۳۶	۱۳۶	۱,۱۶	۵۰٪	۵۹٪	۴۹٪	۴۹٪	۵۱٪	BFL
۱۹	۶۸۸	۱۱۷	۱۱۷	۱,۲۰	۵۴٪	۵۰٪	۴۹٪	۴۹٪	۵۳٪	BIR
۱۴۲	۸۳۷۱	۳۹۲	۳۹۲	۱,۸۵	۶۸٪	۵۰٪	۶۵٪	۶۵٪	۶۱٪	CAR
۵۸۴	۲۴۸۳۰	۹۷۸	۹۷۶	۱۰,۹۴	۵۶٪	۷۸٪	۶۷٪	۶۷٪	۶۶٪	CBF
۲۶۲۳	۷۶۳۵۸۷	۶۶۰۱	۶۵۵۶	۲۸۱,۷۲	۵۳٪	۴۰٪	۴۷٪	۴۹٪	۵۳٪	CHL
۲۷۷۷۲	۷۶۳۵۸۷	۱۱۶۵۳	۱۱۶۴۸	۲۶۲,۵۶	۶۹٪	۵۶٪	۶۴٪	۶۴٪	۶۴٪	CIN
۱۰	۳۸۵	۹۸	۹۸	۰,۸۶	۷۵٪	۴۹٪	۵۱٪	۵۱٪	۷۹٪	COF
۴۱۳	۲۱۰۱۱۳	۲۷۰۰	۲۶۹۹	۱۴,۲۰	۵۰٪	۵۰٪	۵۱٪	۵۰٪	۵۰٪	COM
۴۳۵۶	۸۲۹۷۰	۱۷۱۵	۱۷۰۸	۱۴,۷۹	۴۱٪	۷۸٪	۸۵٪	۸۵٪	۸۳٪	CRX
۴۴۸۸	۸۳۷۱۰	۱۷۵۹	۱۷۵۲	۱۷,۲۶	۵۳٪	۶۹٪	۸۴٪	۸۴٪	۸۴٪	CRY
۳۸۰۷	۷۸۵۶۳	۱۶۸۷	۱۶۸۰	۱۵,۷۰	۴۱٪	۷۱٪	۸۵٪	۸۴٪	۸۴٪	CRZ
۳۰۳	۲۰۴۳۳	۶۲۴	۶۲۳	۴,۱۵	۹۶٪	۳۰٪	۷۲٪	۷۲٪	۳۱٪	DIA
۸۳	۳۰۹۱	۳۱۴	۳۱۲	۳,۶۴	۷۲٪	۷۱٪	۶۰٪	۶۰٪	۷۳٪	DPA
۴۵	۸۵۳۶	۴۸۷	۴۸۴	۷,۱۳	۵۰٪	۵۳٪	۵۱٪	۵۱٪	۵۴٪	DPC
۱۹۴	۳۲۳۱	۳۰۶	۳۰۳	۳,۳۲	۶۶٪	۸۶٪	۶۶٪	۶۸٪	۸۸٪	DPT
۳۰۸	۹۳۲۶۲	۲۲۴۷	۲۲۴۶	۱۰,۰۶	۶۲٪	۵۴٪	۵۳٪	۵۳٪	۵۲٪	EAR
۱۸	۶۵۱	۱۲۲	۱۲۲	۱,۰۵	۶۱٪	۵۴٪	۵۰٪	۵۰٪	۶۲٪	EC۲
۱۶۲۱۴	۷۳۶۰۶۵	۵۰۸۸	۵۰۱۹	۳۲۸,۳۳	۵۹٪	۸۹٪	۶۰٪	۶۴٪	۷۴٪	EC۵
۱۳۶	۱۹۵۷۸	۹۴۵	۹۴۲	۱۰,۵۵	۸۱٪	۵۰٪	۵۰٪	۵۰٪	۵۰٪	EFC
۵۲۳۶	۷۶۳۵۸۷	۱۸۹۳	۱۸۷۶	۶۳,۱۰	۳۰٪	۶۰٪	۸۵٪	۸۵٪	۸۳٪	FAA
۹۰	۲۴۴۹	۲۷۰	۲۷۰	۱,۶۰	۳۸٪	۵۵٪	۵۷٪	۵۷٪	۶۹٪	FAF
۱۲۶۵	۴۲۷۰۷	۹۳۲	۹۳۰	۶,۴۴	۷۹٪	۱۸٪	۶۴٪	۷۳٪	۷۱٪	FIS
۳۶۸۲۳	۷۶۳۵۸۷	۱۶۳۱۶	۱۶۲۶۵	۱۱۱۹,۴۹	۵۰٪	۵۴٪	۵۱٪	۵۲٪	۵۰٪	FOA
۲۸۶۳۷	۷۶۳۵۸۷	۱۲۲۲۲	۱۲۱۸۷	۹۵۸,۵۸	۵۰٪	۵۰٪	۵۰٪	۵۰٪	۵۰٪	FOB
۱۰	۲۰۰۹	۱۳۴	۱۳۴	۱,۶۰	۵۱٪	۵۰٪	۵۴٪	۵۴٪	۵۴٪	GUN
۱۵۲	۱۴۰۳۹	۵۵۸	۵۵۷	۷,۱۳	۵۳٪	۵۰٪	۵۲٪	۵۲٪	۵۱٪	HAM
۳۴۰۵۲	۷۶۳۵۸۷	۲۰۴۰۸	۲۰۴۰۴	۵۵۲,۰۱	۶۹٪	۵۵٪	۵۵٪	۵۰٪	۵۶٪	HAN
۲۸۸۱	۷۶۳۵۸۷	۲۳۰۳	۲۳۰۲	۲۵,۶۴	۶۹٪	۳۹٪	۶۰٪	۶۰٪	۶۷٪	HAP
۳۸	۷۴۷۵	۲۷۳	۲۷۳	۱,۹۷	۵۰٪	۵۱٪	۵۰٪	۵۰٪	۵۰٪	HER
۱۱۲۳۴	۷۶۳۵۸۷	۵۶۰۰	۵۵۹۸	۸۱,۸۸	۷۴٪	۵۴٪	۷۱٪	۷۱٪	۶۸٪	INL
۱۰۴۷۳	۴۷۴۹۰۵	۳۶۸۵	۳۶۶۱	۱۱۸,۸۷	۶۹٪	۵۵٪	۸۱٪	۸۱٪	۸۷٪	INS
۲۰	۱۲۳۳	۲۴۶	۲۴۲	۷,۴۰	۶۴٪	۵۱٪	۵۰٪	۵۰٪	۵۱٪	ITA

711	369379	39.8	39.6	38,69	41%	34%	55%	55%	52%	LAR
128	9478	544	543	2,18	50%	50%	54%	50%	60%	LIY
180	3125	324	323	2,78	59%	60%	75%	75%	76%	LIY
18388	763587	14530	14510	509,06	92%	93%	80%	80%	94%	MAL
86	4718	279	279	2,62	76%	77%	40%	71%	72%	MEA
1062	19954	651	646	21,55	47%	64%	65%	65%	63%	MED
70	4086	251	250	4,92	73%	73%	56%	56%	73%	MPA
65	9306	396	395	8,73	50%	50%	51%	51%	50%	MPC
204	3319	289	286	4,87	81%	80%	82%	74%	84%	MPT
496	20456	721	717	17,48	58%	50%	50%	50%	62%	MOT
218650	763587	17226	17115	992,18	95%	70%	95%	94%	95%	NO1
208416	763587	13676	13594	980,52	97%	85%	95%	95%	96%	NO2
66	2444	139	139	2,36	85%	76%	77%	77%	86%	OLI
670	60765	1100	1097	12,04	29%	62%	73%	73%	73%	OSU
200	77558	1515	1499	77,62	51%	54%	51%	51%	54%	PHA
364703	763587	15525	15477	432,75	51%	45%	93%	93%	93%	PHO
53	2047	149	149	2,99	92%	100%	80%	83%	94%	PLA
82	4909	236	235	5,77	76%	78%	76%	76%	79%	PPA
32	8195	421	418	8,19	53%	54%	56%	56%	52%	PPC
221	3851	321	319	5,52	81%	88%	78%	78%	79%	PPT
900	422843	3972	3969	35,78	39%	35%	54%	56%	55%	REF
1358	355109	4024	4022	37,35	45%	35%	53%	53%	55%	SCR
111	13879	1006	1005	3,14	50%	50%	99%	99%	50%	SHS
51107	546585	3777	3738	68,80	63%	84%	97%	97%	95%	SHA
1487	379869	4050	4048	38,90	54%	34%	59%	59%	43%	SMA
74	2770	350	349	4,66	75%	50%	52%	51%	50%	SO1
149	9040	610	607	9,90	66%	53%	53%	60%	64%	SO2
366	9202	1334	1331	24,01	50%	50%	52%	50%	50%	STR
1620	36092	996	985	23,42	63%	35%	88%	88%	87%	SWE
2867	248454	2779	2774	40,45	60%	89%	81%	81%	89%	SYM
235	2112	365	361	5,62	38%	88%	78%	78%	82%	SYN
121	8315	528	527	3,06	53%	51%	51%	51%	50%	TO1
79	5258	339	338	2,02	53%	67%	50%	50%	51%	TO2
114	4986	325	325	3,24	72%	87%	84%	84%	75%	TRA
4602	581050	6532	6465	361,46	46%	85%	64%	64%	63%	TWP
240	11903	707	703	14,03	54%	50%	64%	64%	50%	TWE
47423	763587	7523	7470	501,15	51%	80%	75%	78%	81%	UWZ
48751	763587	9014	8951	509,23	54%	82%	76%	78%	82%	UWY
47554	763587	7315	7264	589,53	54%	74%	80%	80%	82%	UWZ
167079	763587	18799	18754	1753,92	45%	59%	76%	76%	86%	UWA

۱۶۸۲	۷۶۳۵۸۷	۴۱۵۸	۴۰۸۷	۱۵۷۱,۷۶	۵۹٪	۵۳٪	۶۶٪	۵۰٪	۵۳٪	WAF
۳۱	۱۲۷۳	۱۰۵	۱۰۵	۵,۳۷	۵۹٪	۵۰٪	۵۰٪	۵۷٪	۵۰٪	WIN
۸۲۱۲	۹۸۲۳۳	۱۴۷۶	۱۴۶۳	۵۴,۰۸	۵۰٪	۸۷٪	۸۷٪	۸۷٪	۸۹٪	WOS
۱۰۸۵	۸۲۵۲۷	۱۶۰۶	۱۶۰۵	۱۲,۹۰	۵۳٪	۶۲٪	۵۸٪	۶۰٪	۶۳٪	WOR
۶۷۶	۷۷۴۹۷	۱۶۲۱	۱۶۲۰	۸,۷۰	۵۰٪	۵۰٪	۵۱٪	۵۱٪	۵۰٪	WOT
۴۸۴۸	۷۶۳۵۸۷	۷۹۸۳	۷۹۵۹	۴۸۱,۲۲	۵۰٪	۵۰٪	۵۰٪	۵۱٪	۵۰٪	YOG
۱۷۳۶۵	۲۳۶۴۶۷	۳۲۸۴	۳۲۷۰	۱۶۲	۶۰,۲۶٪	۶۰,۵۵٪	۶۵,۶۸٪	۶۶,۱۰٪	۶۷,۲۵٪	متوسط

جدول ۷: نتایج آزمون آماری الگوریتم ارائه شده با چهار الگوریتم منتخب ادبیات در دو شاخص رند و زمان اجرا

$TS3C_{CH}$	$TS3C_{MV}$	KSC	DD_{DTW}	معیار مقایسه	الگوریتم
۰,۰۲۱	۰,۰۱۸	۰,۰۰۲	۰,۰۰۰	شاخص رند	p-value
۰,۰۰۰	۰,۰۰۰	۰,۰۰۰	۰,۰۰۰	زمان	p-value

۴- نتیجه گیری

است علاوه بر افزایش دقت خوشه‌بندی بتوان این کار را با حداقل هزینه زمانی انجام پذیرد. با توجه به نتایج ارائه شده در جدول ۶ به خوبی می‌توان مشاهده نمود که طرح ارائه شده علاوه بر افزایش دقت خوشه‌بندی می‌تواند این امر را در کمترین زمان ممکن انجام دهد. همچنین نشان داده شده است که طرح ارائه شده در ۳۲٪ از مجموعه داده‌های مورد استفاده توانسته است بهترین دقت را در بین الگوریتم‌های موجود داشته باشد. نتایج آزمون آماری به خوبی نشان داده است که الگوریتم دوم ارائه شده به لحاظ شاخص رند و همچنین زمان اجرا از همه الگوریتم‌های مقایسه شده کارایی بیشتری دارد. با توجه به دقت بالا و هزینه زمانی بسیار ناچیز طرح ارائه شده جهت خوشه‌بندی، الگوریتم مورد نظر می‌تواند به راحتی در مواردی مانند تشخیص الگوهای مصرف انواع حامل‌های انرژی، تشخیص الگوهای رفتاری بیماری‌های واگیردار و غیره مورد استفاده قرار گیرد. همچنین با توجه به اینکه یکی از کاربردهای اصلی خوشه‌بندی استفاده در پیش‌بینی سری‌های زمانی می‌باشد. با افزایش دقت خوشه‌بندی در زمان کم می‌توان عمل دقت پیش‌بینی این نوع داده‌ها را نیز افزایش داد.

در این مقاله یک رویکرد جدید به منظور خوشه‌بندی سری‌های زمانی بر مبنای قطعه‌بندی و خوشه‌بندی ترکیبی در ۵ گام ارائه شد. در این الگوریتم‌ها قطعاتی که می‌توانند نتایج بهتری را در خوشه‌بندی داشته باشند با استفاده از یک یا چند معیار داخلی انتخاب و در مرحله نهایی با یکدیگر ترکیب می‌گردند. الگوریتم اول بر اساس معیارهای داخلی و مقدار پارامتر m در سه وضعیت اجرا شده است. الگوریتم دوم ارائه شده تا گام سوم کاملاً شبیه الگوریتم اول می‌باشد، اما در گام سوم به صورت هم‌زمان از سه معیار داخلی به منظور انتخاب $3m$ برچسب استفاده می‌گردد و در گام پنجم این برچسب‌ها مجدداً به منظور نتیجه نهایی با یکدیگر ترکیب می‌گردند. شاخص رند در نتیجه ۱۰ تکرار برای سه وضعیت مختلف الگوریتم اول و یک وضعیت الگوریتم دوم به ترتیب برابر $۶۳,۵۲٪$ ، $۶۳,۵۲٪$ و $۶۷,۲۵٪$ برای ۸۲ مجموعه داده بوده است. همچنین آزمون آماری ویلکسون به همراه دیگر معیارها نشان داد که الگوریتم دوم می‌تواند عملکرد بهتری در خوشه‌بندی سری‌های زمانی داشته باشد. همچنین عملکرد الگوریتم دوم به لحاظ شاخص رند و متوسط زمان اجرای الگوریتم با ۴ الگوریتم DD_{DTW} ، KSC ، $TS3C_{MV}$ و $TS3C_{CH}$ مقایسه شده است.

۵- مطالعات آینده

طرح ارائه شده دارای پارامترها و مشخصات مختلفی می‌باشد که در این مطالعه صرفاً حالت‌ها اندکی از آن مورد بررسی و تحلیل قرار گرفته است. در مطالعات آینده می‌توان انواع دیگر الگوریتم‌های خوشه‌بندی، معیارهای درونی دیگر و همچنین پارامترهای ورودی دیگر را نیز مورد بررسی و تحلیل قرار داد. همچنین با توجه به کارایی بالای الگوریتم در زمان نسبتاً اندک از الگوریتم فوق می‌توان به منظور

در مطالعات اخیر سعی بر آن شده است با استفاده از تعریف معیارهای جدید فاصله و یا استفاده از الگوریتم‌های خاص کارایی خوشه‌بندی افزایش یابد. اما با این وجود نتایج مطالعات اخیر نشان داده است که استفاده از این روش‌ها علی‌رغم بهبود دقت خوشه‌بندی می‌تواند هزینه زمانی بسیار بالایی داشته باشد. به عنوان مثال در چهار الگوریتم موجود در ادبیات حتی می‌تواند به عددی نزدیک به ۲۳۶ هزار ثانیه برسد که عملاً استفاده از این الگوریتم‌ها در بخش صنعت و خدمات را ناممکن می‌نماید. در طرح ارائه شده سعی بر آن شده

- clustering: a comparative study," *Data Mining and Knowledge Discovery*, vol. ۳۶, no. ۱, pp. ۲۹-۸۱, ۲۰۲۲.
- [۹] S. Zolhavarieh, S. Aghabozorgi, and Y. W. Teh, "A review of subsequence time series clustering," *The Scientific World Journal*, vol. ۲۰۱۴, ۲۰۱۴.
- [۱۰] C. A. Ralanamahatana, J. Lin, D. Gunopulos, E. Keogh, M. Vlachos, and G. Das, "Mining time series data," in *Data mining and knowledge discovery handbook*: Springer, ۲۰۰۵, pp. ۱۰۶۹-۱۱۰۳.
- [۱۱] H. Kamalzadeh, A. Ahmadi, and S. Mansour, "Clustering time-series by a novel slope-based similarity measure considering particle swarm optimization," *Applied Soft Computing*, p. ۱۰۶۷۰۱, ۲۰۲۰.
- [۱۲] G. Soleimani and M. Abessi, "DLCSS: A new similarity measure for time series data mining," *Engineering Applications of Artificial Intelligence*, vol. ۹۲, p. ۱۰۳۶۶۴, ۲۰۲۰.
- [۱۳] M. Łuczak, "Hierarchical clustering of time series data with parametric derivative dynamic time warping," *Expert Systems with Applications*, vol. ۶۲, pp. ۱۱۶-۱۳۰, ۲۰۱۶.
- [۱۴] L. Rokach and O. Maimon, "Clustering methods," in *Data mining and knowledge discovery handbook*: Springer, ۲۰۰۵, pp. ۳۲۱-۳۵۲.
- [۱۵] M. A. Rahim Khan and M. Zakarya, "Longest common subsequence based algorithm for measuring similarity between time series: a new approach," *World Applied Sciences Journal*, vol. ۲۴, no. ۹, pp. ۱۱۹۲-۱۱۹۸, ۲۰۱۳.
- [۱۶] X. Wang, F. Yu, W. Pedrycz, and J. Wang, "Hierarchical clustering of unequal-length time series with area-based shape distance," *Soft Computing*, vol. ۲۳, no. ۱۵, pp. ۶۳۳۱-۶۳۴۳, ۲۰۱۹.
- [۱۷] K. K. W. Chu and M. H. Wong, "Fast time-series searching with scaling and shifting," in *Proceedings of the eighteenth ACM SIGMOD-SIGACT-SIGART*
- [۱] M. Maleki, H. Bidram, and D. Wraith, "Robust clustering of COVID-۱۹ cases across US counties using mixtures of asymmetric time series models with time varying and freely indexed covariates," *Journal of Applied Statistics*, pp. ۱-۱۵, ۲۰۲۲.
- [۲] E. Bair, "Semi-supervised clustering methods," *Wiley Interdisciplinary Reviews: Computational Statistics*, vol. ۵, no. ۵, pp. ۳۴۹-۳۶۱, ۲۰۱۳.
- [۳] H. A. Dau, N. Begum, and E. Keogh, "Semi-supervision dramatically improves time series clustering under dynamic time warping," in *Proceedings of the ۲۵th ACM International Conference on Information and Knowledge Management*, ۲۰۱۶, pp. ۹۹۹-۱۰۰۸.
- [۴] L. Alhusain and A. M. Hafez, "Cluster ensemble based on Random Forests for genetic data," *BioData Mining*, vol. ۱۰, no. ۱, p. ۳۷, ۲۰۱۷.
- [۵] R. Ma and R. Angryk, "Distance and density clustering for time series data," in *۲۰۱۷ IEEE International Conference on Data Mining Workshops (ICDMW)*, ۲۰۱۷: IEEE, pp. ۲۵-۳۲.
- [۶] S. Mehrmolaei and M. R. Keyvanpour, "A comparative study on weighting-based clustering techniques: Time series data," in *۲۰۱۸ ۸th Conference of AI & Robotics and ۱۰th RoboCup Iranopen International Symposium (IRANOPEN)*, ۲۰۱۸: IEEE, pp. ۶۵-۷۲.
- [۷] N. Tavakoli, S. Siami-Namini, M. Adl Khanghah, F. Mirza Soltani, and A. Siami Namin, "An autoencoder-based deep learning approach for clustering time series data," *SN Applied Sciences*, vol. ۲, pp. ۱-۲۵, ۲۰۲۰.
- [۸] B. Lafabregue, J. Weber, P. Gançarski, and G. Forestier, "End-to-end deep representation learning for time series

- Engineering (TCSET)*, 2020: IEEE, pp. 603-608.
- [26] R. J. Hyndman, E. Wang, and N. Laptev, "Large-scale unusual time series detection," in 2010 *IEEE international conference on data mining workshop (ICDMW)*, 2010: IEEE, pp. 1616-1619.
- [27] Y. Zou, R. V. Donner, N. Marwan, J. F. Donges, and J. Kurths, "Complex network approaches to nonlinear time series analysis," *Physics Reports*, vol. 487, pp. 1-97, 2019.
- [28] V. A. F. da Silva, "Time Series Analysis based on Complex Networks," 2018.
- [29] L. N. Ferreira and L. Zhao, "Time series clustering via community detection in networks," *Information Sciences*, vol. 326, pp. 227-242, 2016.
- [30] F. Bonacina, E. S. Miele, and A. Corsini, "Time Series Clustering: A Complex Network-Based Approach for Feature Selection in Multi-Sensor Data," *Modelling*, vol. 1, no. 1, pp. 1-21, 2020.
- [31] A. Koski, M. Juhola, and M. Meriste, "Syntactic recognition of ECG signals by attributed finite automata," *Pattern Recognition*, vol. 28, no. 12, pp. 1927-1940, 1995.
- [32] E. J. Keogh and M. J. Pazzani, "An Enhanced Representation of Time Series Which Allows Fast and Accurate Classification, Clustering and Relevance Feedback," in *Kdd*, 1998, vol. 98, pp. 239-243.
- [33] E. Keogh, S. Chu, D. Hart, and M. Pazzani, "Segmenting time series: A survey and novel approach," in *Data mining in time series databases*: World Scientific, 2004, pp. 1-21.
- [34] C. Faloutsos, M. Ranganathan, and Y. Manolopoulos, "Fast subsequence matching in time-series databases," *ACM Sigmod Record*, vol. 23, no. 2, pp. 419-429, 1994.
- [35] E. Keogh and C. A. Ratanamahatana, "Exact indexing of dynamic time warping," *Knowledge and information symposium on Principles of database systems*, 1999, pp. 237-248.
- [18] S. Gharghabi, S. Imani, A. Bagnall, A. Darvishzadeh, and E. Keogh, "Matrix profile xii: Mpdist: a novel time series distance measure to allow data mining in more challenging scenarios," in 2018 *IEEE International Conference on Data Mining (ICDM)*, 2018: IEEE, pp. 960-970.
- [19] G. Jiang, W. Wang, and W. Zhang, "A novel distance measure for time series: Maximum shifting correlation distance," *Pattern Recognition Letters*, vol. 117, pp. 58-65, 2019.
- [20] D. Hong, Q. Gu, and K. Whitehouse, "High-dimensional time series clustering via cross-predictability," in *Artificial Intelligence and Statistics*, 2017: PMLR, pp. 642-651.
- [21] T. Górecki, "Classification of time series using combination of DTW and LCSS dissimilarity measures," *Communications in Statistics-Simulation and Computation*, vol. 47, no. 1, pp. 263-276, 2018.
- [22] D. Guijo-Rubio, A. M. Durán-Rosal, P. A. Gutiérrez, A. Troncoso, and C. Hervás-Martínez, "Time-Series Clustering Based on the Characterization of Segment Typologies," *IEEE Transactions on Cybernetics*, 2020.
- [23] S. Aghabozorgi, T. Ying Wah, T. Herawan, H. A. Jalab, M. A. Shaygan, and A. Jalali, "A hybrid algorithm for clustering of time series data based on affinity search technique," *The Scientific World Journal*, vol. 2014, 2014.
- [24] X. Zhang, J. Liu, Y. Du, and T. Lv, "A novel clustering method on time series data," *Expert Systems with Applications*, vol. 38, no. 9, pp. 11891-11900, 2011.
- [25] N. Manakova and V. Tkachenko, "Two-stage time-series clustering approach under reducing time cost requirement," in 2020 *IEEE 10th International Conference on Advanced Trends in Radioelectronics, Telecommunications and Computer*

- [40] H. A. Dau *et al.*, "The UCR time series archive," *IEEE/CAA Journal of Automatica Sinica*, vol. 6, no. 6, pp. 1293-1300, 2019.
- [46] J. Demšar, "Statistical comparisons of classifiers over multiple data sets," *The Journal of Machine Learning Research*, vol. 9, pp. 1-30, 2006.
- [47] J. Yang and J. Leskovec, "Patterns of temporal variation in online media," in *Proceedings of the fourth ACM international conference on Web search and data mining*, 2011, pp. 177-186.
- [36] M. Djukanovic, G. R. Raidl, and C. Blum, "Finding Longest Common Subsequences: New anytime A* search results," *Applied Soft Computing*, vol. 90, p. 107499, 2020.
- [37] M. Paterson and V. Dančík, "Longest common subsequences," in *International symposium on mathematical foundations of computer science*, 1994: Springer, pp. 127-142.
- [38] R. Lin, A. King-Ip, and H. S. S. K. Shim, "Fast similarity search in the presence of noise, scaling, and translation in time-series databases," in *Proceeding of the 21th International Conference on Very Large Data Bases*, 1990: Citeseer, pp. 490-501.
- [39] M. Vlachos, G. Kollios, and D. Gunopulos, "Discovering similar multidimensional trajectories," in *Proceedings 14th international conference on data engineering*, 2002: IEEE, pp. 673-684.
- [40] T. Boongoen and N. Iam-On, "Cluster ensembles: A survey of approaches with recent extensions and applications," *Computer Science Review*, vol. 28, pp. 1-20, 2018.
- [41] D. Huang, C.-D. Wang, and J.-H. Lai, "Locally weighted ensemble clustering," *IEEE transactions on cybernetics*, vol. 48, no. 5, pp. 1470-1473, 2017.
- [42] Y. Liu, Z. Li, H. Xiong, X. Gao, and J. Wu, "Understanding of internal clustering validation measures," in 2010 *IEEE international conference on data mining*, 2010: IEEE, pp. 911-916.
- [43] J. C. Dunn, "Well-separated clusters and optimal fuzzy partitions," *Journal of cybernetics*, vol. 4, no. 1, pp. 90-104, 1974.
- [44] T. Caliński and J. Harabasz, "A dendrite method for cluster analysis," *Communications in Statistics-theory and Methods*, vol. 3, no. 1, pp. 1-27, 1974.