

# مقایسه‌ی عملکرد روش‌های اندازه‌گیری شباهت طولانی‌ترین زیردباله‌ی مشترک و چرخش زمانی پویا در داده‌کاوی سری‌های زمانی

غلام‌حسین سلیمانی (دانشجوی دکتری)

مسعود عابسی\* (استادیار)

گروه مهندسی صنایع، دانشکده فنی و مهندسی، دانشگاه بزد

تکنیک‌های داده‌کاوی به طور خاص برای داده‌های ثابت طراحی شده‌اند. لذا به کارگیری آنها برای داده‌های سری زمانی نیازمند اعمال تعییراتی (روش اندازه‌گیری شباهت) است. بر اساس تحقیقات اخیر، روش‌های طولانی‌ترین زیردباله‌ی مشترک و چرخش زمانی پویا، از پرکاربردترین و کارآترین این روش‌ها محسوب می‌شود. در این تحقیق، قصد داریم تا عملکرد این روش‌ها را در تکنیک‌های نزدیک‌ترین همسایگی و خوشه‌بندی کامدوید مورد ارزیابی و مقایسه قرار داده تا بتوان از آنها دقت بهتری در این تکنیک‌ها و در مسائلی نظری قسمت‌بندی مشترک‌یان، زمان‌بندی کارگاه و ... استفاده کرد. به همین منظور از ۶۳ مجموعه داده سری زمانی از بانک اطلاعاتی UCR، استفاده می‌شود. نتایج نشان می‌دهد که تأثیر آنها در دقت تشخیص درست دسته‌ی سری زمانی و دقت خوشه‌بندی، به طور معناداری تفاوت دارد، ولی تأثیر آنها در تعیین تعداد خوش و نخاینده خوش، تفاوت معناداری ندارد.

gholam\_soleimani@yahoo.com  
mabessi@gmail.com

واژگان کلیدی: داده‌کاوی سری‌های زمانی، خوشه‌بندی، نزدیک‌ترین همسایگی، طولانی‌ترین زیردباله‌ی مشترک، چرخش زمانی پویا.

## ۱. مقدمه

سری زمانی، منجر به ایجاد نتایج معتبه نخواهد شد. بنابراین به منظور قابل استفاده کردن این تکنیک‌ها برای داده‌های سری زمانی باستی تعییراتی در الگوریتم‌های داده‌کاوی اعمال کرد. یکی از این تعییرات، انتخاب و استفاده از یک روش اندازه‌گیری شباهت مناسب برای داده‌های سری زمانی است.<sup>[۱]</sup> از سوی دیگر، روش‌های بسیاری برای این منظور ارائه شده که هر کدام از یک سطح کیفی برخوردارند. از میان این روش‌ها، روش طولانی‌ترین زیردباله‌ی مشترک<sup>۱</sup> و چرخش زمانی پویا<sup>۲</sup> پرکاربردترین و کارآترین روش‌های اندازه‌گیری شباهت سری‌های زمانی هستند.<sup>[۳]</sup>

در این تحقیق قصد داریم به ارزیابی عملکرد این دو روش در مورد برخی از تکنیک‌های داده‌کاوی پردازیم تا مشخص شود که تأثیر آنها در کدام تکنیک بهتر است؟ به منظور این ارزیابی از تکنیک‌های نزدیک‌ترین همسایگی و خوشه‌بندی کامدوید استفاده می‌کنیم.

تکنیک نزدیک‌ترین همسایگی یکی از تکنیک‌های دسته‌بندی است که کاربرد گسترده‌بی دارد. این تکنیک در حوزه‌های مختلف نظری پژوهشی - در تعیین وضعیت بیماری مراجعین در مقایسه با وضعیت سلامتی بیماران قبلی - کاربرد دارد. در حوزه‌ی اقتصادی نیز این تکنیک برای تعیین وضعیت ارزش سهام یک شرکت در مقایسه با وضعیت ارزش سهام شرکت‌های موجود استفاده می‌شود. به طور ساده این تکنیک، عبارت است از جست‌وجو برای یافتن رده‌ی یک شیء در مقایسه با رده‌ی موجود از دیگر اشیاء. تکنیک نزدیک‌ترین همسایگی برای سری‌های زمانی نیز

داده‌های سری زمانی، دنباله‌بی از اعداد هستند که در فواصل زمانی غالباً یکسان اندازه‌گیری شده است.<sup>[۴]</sup> سری‌های زمانی تقریباً در همه حیطه‌ها وجود دارد (برای مثال: در حیطه‌ی پژوهشی به صورت نوار ضربان قلب، نوار تنفس و نوار مغز برای یک دوره‌ی زمانی؛ در حیطه‌ی آب و هوای نظیر دمای یک مکان یا میزان رطوبت یک مکان در فواصل زمانی مشخص؛ در حیطه‌ی فروش نظری ارزش فروش در فواصل زمانی مشخص و ...). سری‌های زمانی دارای دو ویژگی برجسته‌اند: ۱. طول سری زمانی؛ یک سری زمانی می‌تواند طولی برابر با صد هزارها باشد. این ویژگی در میزان اشتغال فضای حافظه و سخت شدن محاسبات تأثیر به سزاوی دارد. ۲. وابستگی بین اعداد؛ مقداری یک سری زمانی به هم وابسته‌اند و همین موضوع، ضرورت توجه به انتخاب ابزارهای مناسب برای انجام محاسبات و تحلیل آنها را مشکل‌تر می‌کند.<sup>[۵-۶]</sup>

برای کشف داشت از انبوی از داده، «داده‌کاوی» ابزاری مفید محسوب می‌شود. تکنیک‌های مختلفی برای داده‌کاوی وجود دارد، نظری دسته‌بندی، خوشه‌بندی، کشف الگوهای جذاب و تکراری و در نهایت کشف دورافتاده‌ها. جالب توجه است که همه این تکنیک‌ها برای داده‌های ثابت طراحی شده‌اند<sup>[۷]</sup> و به کارگیری آنها برای داده‌های

\* نویسنده مسئول  
تاریخ: دریافت ۱۳۹۹/۷/۹، اصلاحیه ۱۴۰۰/۲/۶، پذیرش ۱۴۰۰/۴/۲۳

کامل<sup>۳</sup> و انطباق زیربنایی.<sup>۴</sup> در رویکرد انطباق کامل، کل طول سری‌های زمانی در فرایند اندازه‌گیری شباهت استفاده می‌شود، به طوری که اگر طول هر دو سری زمانی برابر  $m$  باشد آنگاه همه‌ی  $m$  داده‌ی سری زمانی اول و همه‌ی  $m$  داده‌ی سری زمانی دوم در فرایند اندازه‌گیری شباهت این دو سری زمانی کاربرد دارد؛ نظیر روش اندازه‌گیری شباهت اقلیدسی، طولانی ترین زیربنایی مشترک و چرخش زمانی پویا.

در رویکرد انطباق زیربنایی، سری‌های زمانی می‌توانند طول‌های متفاوتی داشته باشند، به طوری که اگر طول یک سری زمانی برابر  $n$  و طول سری زمانی دیگر برابر  $m$  باشد و  $n > m$  باشد، آنگاه شباهت زیربنایی از سری زمانی بزرگ‌تر (با اعضای متواالی) با سری زمانی کوچک‌تر اندازه‌گیری می‌شود و شباهت دو سری زمانی برابر با حداکثر شباهت حاصل از فرایند فوق خواهد بود.<sup>[۲۰]</sup>

به طورکل معیارهای اندازه‌گیری مسافت (شباهت) را می‌توان به چهار گروه کلی تقسیم کرد:<sup>[۲۱]</sup>

۱. معیارهای مسافت مبتنی بر شکل<sup>۵</sup>;

۲. معیارهای مسافت مبتنی بر ویرایش<sup>۶</sup>؛

۳. معیارهای مسافت مبتنی بر ویژگی<sup>۷</sup>؛

۴. معیارهای مسافت مبتنی بر مدل.<sup>۸</sup>

در قسمت بعدی از میان تمامی روش‌های اندازه‌گیری شباهت موجود برای داده‌های سری زمانی، به دو روش مطرح DTW و LCSS اشاره می‌شود.

## ۲.۲. روش چرخش زمانی پویا (DTW)

در توضیح این روش باید عنوان کرد که گاهی اوقات برخی از سری‌های زمانی نسبتاً مشابه‌اند ولی این شباهت نه در زمان یکسان، بلکه با یک تأخیر زمانی اتفاق افتاده است. لذا برای شناسایی شباهت این چنین سری‌های زمانی باید روش متفاوتی ایجاد شود؛ این روش متفاوت همان روش چرخش زمانی پویا (DTW) است<sup>[۲۱]</sup> که می‌تواند شباهت را بین سری‌های زمانی با طول غیریکسان اندازه‌گیری کند. این روش با کشیدن یکی از سری‌های زمانی یا هر دو سری زمانی انجام می‌شود و سعی دارد با این کار، شکل دو سری زمانی را به یکدیگر مشابه کند. این روش که از مفهوم شباهت یک نقطه به چند نقطه و برعکس بهره می‌برد از جمله روش‌های شباهت مبتنی بر شکل محسوب می‌شود و در آن از انطباق کامل در تعیین شباهت دو سری زمانی استفاده می‌شود. در صورتی که  $TS_x = (x_1, x_2, \dots, x_n)$  (۱) است<sup>[۲۱]</sup>  $TS_y = (y_1, y_2, \dots, y_m)$  و بیان‌گر دو سری زمانی به ترتیب با طول‌های  $n$  و  $m$  باشند آن‌گاه مسافت بین دو سری زمانی طبق روش DTW به صورت  $TS_x, TS_y$  باشد  $DTW_{TS_x, TS_y} = M(n, m)$  به طوری که مقدار  $M(n, m)$  از رابطه‌ی برگشتی ۱ محاسبه می‌شود.

$$M(i, j) =$$

$$\left\{ \begin{array}{ll} (x_i - y_j)^r & ; i = 1, j = 1 \\ (x_i - y_j)^r + M(i, j - 1) & ; i = 1, j \geq 2 \\ (x_i - y_j)^r + M(i - 1, j) & ; i \geq 2, j = 1 \\ (x_i - y_j)^r + \min \left\{ \begin{array}{ll} M(i - 1, j) \\ M(i, j - 1) & ; i \geq 2, j \geq 2 \\ M(i - 1, j - 1) \end{array} \right. & \end{array} \right. \quad (1)$$

عبارت است از: «تعیین رده‌ی یک سری زمانی در مقایسه با رده‌ی سایر سری‌های زمانی موجود در یک مجموعه داده». عناصر موجود تکنیک خوش‌بندی یک روش یادگیری بدون نظرارت است و عبارت است از تقسیم‌بندی مجموعه‌ی از اعضا به گروه‌های موسوم به «خوش». عناصر موجود در هر خوش با یکدیگر بسیار شبیه‌اند ولی با عناصر سایر خوش‌ها شباهت کم‌تری دارند. طبق تعریف، خوش‌بندی سری زمانی به معنای تقسیم کردن یک مجموعه سری زمانی به چند خوش و با شرایط فوق است. باید توجه داشت که خوش‌بندی سری زمانی به چند گروه تقسیم می‌شوند: خوش‌بندی قسمت‌بندی، خوش‌بندی الگوریتم است. تکنیک‌های خوش‌بندی به چند گروه تقسیم کردند: خوش‌بندی مبتنی بر شبکه.<sup>[۱۷-۱۹]</sup> در اکثر این تکنیک‌ها، تعداد خوش توسط کار بر تعیین می‌شود و به عنوان ورودی الگوریتم است. در این تحقیق از الگوریتم خوش‌بندی کامدوید که زیرمجموعه‌ی تکنیک خوش‌بندی قسمت‌بندی است استفاده می‌شود. عمل انتخاب این الگوریتم نیز عدم نیاز به میانگین‌گیری در مراحل اجرای الگوریتم است.

در ادامه، ابتدا ضمن معرفی مفهوم شباهت در سری‌های زمانی، انواع رویکردهای اندازه‌گیری شباهت سری‌های زمانی معرفی می‌شود. سپس، روش‌های اندازه‌گیری شباهت DTW و LCSS تشریح خواهد شد. در قسمت سوم، تکنیک‌های نزدیک ترین همسایگی و خوش‌بندی کامدوید معرفی می‌شود. در قسمت چهارم، به تشریح ارزیابی به کار رفته در این تحقیق می‌پردازیم و در قسمت آخر، تابع اجرای تکنیک‌های ارزیابی به کار رفته در این تحقیق می‌پردازیم و در قسمت آخر، تابع اجرای تکنیک‌های فوق تحت روش‌های LCSS و DTW ارائه و مورد بحث قرار می‌گیرد.

## ۲. مرور ادبیات

### ۲.۱. تعریف شباهت و انواع معیارهای اندازه‌گیری شباهت

چنان که پیش‌تر نیز اشاره شد، یکی از مهم‌ترین موضوعات مطرح در داده‌کاوی سری زمانی، مسئله‌ی اندازه‌گیری شباهت سری‌های زمانی است. براساس تحقیقات مرتبط، شباهت در سری‌های زمانی به صورت شباهت نقطه به نقطه یا شباهت منعطف (یک نقطه به چند نقطه یا چند نقطه به یک نقطه) تعریف می‌شود. شباهت نقطه به نقطه به معنای آن است که هر داده از یک سری زمانی فقط می‌تواند با یک داده از سری زمانی دیگر شباهت داشته باشد و برعکس - نظری روش اندازه‌گیری شباهت اقلیدسی و طولانی ترین زیربنایی مشترک. شباهت یک نقطه به چند نقطه و برعکس نیز به معنای آن است که یک داده از یک سری زمانی می‌تواند با بیش از یک داده از سری زمانی دیگر شباهت داشته باشد و برعکس، نظری روش اندازه‌گیری شباهت چرخش زمانی پویا.<sup>[۱۹]</sup>

از منظر دیگر نیز می‌توان شباهت سری‌های زمانی را به صورت شباهت در زمان، شباهت در شکل و شباهت در مدل تعریف کرد. شباهت در زمان به معنای شباهت بین دو سری زمانی بر مبنای میزان شباهت داده‌ها در زمان‌های یکسان است. به عبارت دیگر در این تعریف، زمان و موقع اهمیت دارد. به عنوان مثال مقایسه‌ی سری زمانی فروش یک شرکت با سری زمانی فروش شرکت دیگر در همان دوره زمانی. شباهت در شکل به معنای شباهت بین دو سری زمانی بر مبنای شباهت بین زیربنایهای آنهاست. در این نوع شباهت، زمان اهمیت ندارد و فقط شکل سری زمانی اهمیت دارد. حتی ابعاد طولی و عرضی هم اهمیت ندارد. شباهت در مدل نیز به معنای میزان یکسان بودن پارامترها و مدل برآشش شده به سری‌های زمانی است.<sup>[۲۰]</sup> برای اندازه‌گیری شباهت دو سری زمانی، دو رویکرد کلی وجود دارد: انطباق:

از آن‌جا که نمی‌توان میزان تشابه دو دنباله‌ی  $S_x$  و  $S_y$  را با میزان تشابه دو دنباله‌ی دیگر نظری  $S_U$  و  $S_V$  (فرض کنید طول دنباله‌های  $S_U$  و  $S_V$  به ترتیب  $n_1$  و  $n_2$  باشد) فقط به استناد مقادیر  $LCSS(S_x, S_y)$  و  $LCSS(S_U, S_V)$  مشخص کرد (زیرا طول این دنباله‌ها می‌تواند با هم متفاوت باشد و هر مقدار که طول آنها بیشتر باشد  $n_1 + n_2 - LCSS(S_x, S_y)$  است). لذا از روابط  $sim(S_U, S_V) = \frac{LCSS(S_U, S_V)}{n_1 + n_2}$  و  $sim(S_x, S_y) = \frac{LCSS(S_x, S_y)}{n_1 + n_2}$  استفاده می‌شود. بدین معنا که مقدار  $sim(S_U, S_V)$  با مقدار  $sim(S_x, S_y)$  مقایسه می‌شود و هر کدام که بیشتر بود به معنای تشابه بیشتر آن دو دنباله نسبت به هم، در مقایسه با تشابه دو دنباله‌ی دیگر نسبت به هم است.

برای این که بتوان از این روش نیز برای داده‌های سری زمانی استفاده کرد، از حد مجاز شباخت ( $\infty$ ) برای تعیین شباخت دو داده از دو سری زمانی استفاده می‌شود.<sup>[۲۸]</sup> با این توضیحات، رابطه‌ی ۲ به صورت رابطه‌ی ۳ بازنویسی می‌شود.

$$M(i, j) =$$

$$\begin{cases} 0 & ; i = 0 \text{ or } j = 0 \\ 1 + M(i-1, j-1) ; |x_i - y_j| \leq \epsilon, i \geq 1 \text{ or } j \geq 1 \\ Max \begin{cases} M(i-1, j) \\ M(i, j-1) \end{cases} & ; < |x_i - y_j|, i \geq 1 \text{ or } j \geq 1 \end{cases} \quad (3)$$

چنان‌که مشاهده می‌شود، هرگاه قدر مطلق تفاصل دو داده از دو سری زمانی، کوچک‌تر یا مساوی با مقدار حد مجاز شباخت باشد آن‌گاه دو داده با یکدیگر مشابه‌اند. در این شرایط یک واحد طول به طول طولانی‌ترین زیردنباله‌ی مشترک آنها تا قبل از این دو داده اضافه خواهد شد. در غیر این صورت طول طولانی‌ترین زیردنباله‌ی مشترک آنها برایر با طول طولانی‌ترین زیردنباله‌ی مشترک آنها در شرایطی است که طول آن دنباله‌ها یکی کم‌تر باشد.

نتیجه‌ی حاصل از بهکارگیری این روش به شدت به مقدار حد مجاز شباخت وابسته است. به بیان واضح، اگر مقدار حد مجاز شباخت کوچک (سخت‌گیرانه) باشد، آنگاه طول طولانی‌ترین زیردنباله‌ی مشترک نیز کوتاه‌تر خواهد بود و اگر مقدار حد مجاز شباخت بزرگ (سهل‌گیرانه) باشد، آنگاه طول طولانی‌ترین زیردنباله‌ی مشترک نیز بلندتر خواهد بود. واقعیت آن است که به علت عدم آگاهی از ماهیت مجموعه داده، انتخاب بهترین مقدار حد مجاز شباخت به سیار دشوار است و این امر باعث می‌شود که در صورت انتخاب نامناسب مقدار حد مجاز شباخت، خطای نتیجه‌گیری نیز افزایش یابد.<sup>[۲۹]</sup>

به‌طورکلی، این روش نیز نسبت به روش‌های دیگر اندازه‌گیری شباخت سری‌های زمانی از عملکرد بهتری برخوردار است<sup>[۲۷]</sup> و در تحقیقات مختلف با اهداف متنوعی استفاده شده است. از این روش به منظور دسته‌بندی متن‌های داده‌یی نیز افزایش یابد.<sup>[۳۰]</sup> همچنین به منظور طراحی یک سیستم پیش‌بینی کارای داده‌های سری زمانی، از این روش با به کارگیری الگوریتم‌های انطباق‌گویی استفاده شده است.<sup>[۲۱]</sup> از این روش برای طراحی یک رویکرد ابتکاری برای بجینگ و زمان‌بندی ماشین‌های مغناطی به منظور کمینه‌سازی هزینه‌های راه‌اندازی<sup>[۲۲]</sup> و نیز برای طراحی یک الگوریتم ممتیک مبنی بر یادگیری برای مسئله‌ی برداشت و تحويل چندین وسایل باری بمنطبق با LIFO<sup>[۲۳]</sup> استفاده شده است. در برآورد یک الگوریتم توزیع برای کارگاه‌های فلوشاب به منظور کمینه‌سازی زمان جريان کل<sup>[۲۴]</sup> و نیز برای سرعت بخشیدن زنجیره‌های تأمین با بهینه سازی کلونی موره‌بی در طیف وسیعی از راه حل‌های ساخت‌افزاری،<sup>[۲۵]</sup> از این روش استفاده شده است.

این روش، دارای دو خروجی است. یکی از این خروجی‌ها، دنباله‌یی با عناصر زوجی و با طول  $r$  است به طوری که  $r$  بزرگ‌تر یا مساوی با  $\max(n, m)$  است. این دنباله‌ی زوجی را  $R$  نامگذاری می‌کنیم. خروجی دیگر مقدار مسافت (عدم شباخت) دو سری زمانی مورد نظر یا همان  $(TS_x, TS_y)$  است که بیان‌گر مجموع مربعات تفاصل عناصر زوجی دنباله‌ی  $R$  است.

از آن‌جا که نمی‌توان میزان عدم شباخت به سری زمانی  $TS_x$  و  $TS_y$  با میزان عدم  $DTW(TS_x, TS_y)$  را فقط به استناد مقادیر  $DTW(TS_x, TS_y)$  تعیین کرد (زیرا طول آین دنباله‌ی زمانی می‌تواند با هم متفاوت باشد و هر مقدار که طول آنها ممکن است مقدار  $DTW$  آنها نیز بزرگ‌تر باشد)، لذا از روابط  $dissim(TS_x, TS_y) = \sqrt{\frac{DTW(TS_x, TS_y)}{\|x\| \|y\|}}$  و  $dissim(TS_U, TS_V) = \sqrt{\frac{DTW(TS_U, TS_V)}{\|x\| \|y\|}}$  استفاده می‌شود، به طوری که  $r_1$  و  $r_2$  به ترتیب طول دنباله‌های ایجاد شده ناشی از اجرای روش برای  $DTW$  دو سری زمانی  $TS_x$  و  $TS_y$  و برای دو سری زمانی  $TS_U$  و  $TS_V$  مقدار  $dissim(TS_x, TS_y)$  با مقدار  $dissim(TS_U, TS_V)$  مقایسه می‌شود و هر کدام که کم‌تر بود به معنای تشابه بیشتر آن دو سری زمانی نسبت به یکدیگر در مقایسه با تشابه دو سری زمانی دیگر است.

از این روش اندازه‌گیری شباخت در تحقیق‌های بسیاری استفاده شده است، به عنوان مثال از روش  $DTW$  به عنوان روش اندازه‌گیری شباخت در بخش‌بندی بازار از طریق دادکاری،<sup>[۲۲]</sup> برای انتخاب ویژگی‌ها به منظور دسته‌بندی سرعی سری‌های زمانی از طریق الگوریتم ژنتیک کارا<sup>[۲۳]</sup> در تشخیص سیگنالهای زیست‌شناختی به عنوان روش اندازه‌گیری شباخت،<sup>[۲۴]</sup> برای خوشه‌بندی سری‌های زمانی چندمتغیر مبتنی بر امواج با بهکارگیری رویکرد SPCA<sup>[۲۵]</sup> برای بخش‌بندی صفحه<sup>[۲۶]</sup> استفاده شده است. به طور کلی این روش نسبت به روش‌های دیگر اندازه‌گیری شباخت سری‌های زمانی از عملکرد بهتری برخوردار است.<sup>[۲۷]</sup> ضمن این که به‌طور وسیع در تحقیقات مورد استفاده قرار می‌گیرد.

## ۳.۲. روش LCSS

این روش که مسئله‌یی کلاسیک در علم رایانه محسوب می‌شود، در ابتدا برای مقایسه دنباله‌یی از کاراکترها ایجاد شد. مهم‌ترین ویژگی این روش آن است که می‌تواند از مقادیر دورافتاده و نویزدار چشم‌پوشی کند. شباخت در این روش به صورت «عین هم بودن» دو کاراکتر از دو دنباله با طول‌های  $n$  و  $m$  باشند، آن‌گاه طولانی‌ترین زیردنباله‌ی مشترک آنها به صورت  $LCSS(S_x, S_y)$  نمایش داده می‌شود و برابر است با  $M(n, m)$  مقدار  $M(n, m)$  نیز از رابطه‌ی برگشتی<sup>[۲۷]</sup> محاسبه می‌شود و همواره عددی صحیح، بین صفر و  $\min(n, m)$  است.

$$M(i, j) = \begin{cases} 0 & ; i = 0 \text{ or } j = 0 \\ 1 + M(i-1, j-1) & ; x_i = y_j, i \geq 1 \text{ or } j \geq 1 \\ Max \begin{cases} M(i-1, j) \\ M(i, j-1) \end{cases} & ; x_i \neq y_j, i \geq 1 \text{ or } j \geq 1 \end{cases} \quad (2)$$

روش LCSS دارای دو خروجی است. خروجی اول، طولانی‌ترین زیردنباله‌ی مشترک دو دنباله است که با  $LCSS$  نشان داده می‌شود و خروجی دیگر، دنباله‌یی با طول  $LCSS$  که بیان‌گر عناصر مشترک بین دو دنباله مورد نظر است.

### ۳. تکنیک‌های داده‌کاوی

#### ۳.۱. تکنیک نزدیک‌ترین همسایگی

اگر  $TS$  بیان‌گر یک سری زمانی،  $DB$  بیان‌گر یک مجموعه داده سری زمانی و  $d$  بیان‌گر روش اندازه‌گیری شباهت باشد، تکنیک نزدیک‌ترین همسایگی، به دنبال یافتن یک سری زمانی از مجموعه  $DB$  است به‌گونه‌یی که دارای کمترین مقدار عدم شباهت (بیشترین شباهت) با سری زمانی  $TS$  باشد. البته این تکنیک را می‌توان به صورت «یافتن مجموعه‌یی از سری‌های زمانی از مجموعه  $DB$  که مقدار عدم شباهت آنها با سری زمانی  $TS$  از مقدار مشخصی کم‌تر باشد» نیز تعریف کرد. در تعریفی دیگر این تکنیک عبارت است از یافتن تعداد مشخصی از سری‌های زمانی از مجموعه  $DB$  که دارای کمترین مقدار عدم شباهت با سری زمانی  $TS$  باشند.

#### ۳.۲. تکنیک خوشه‌بندی با الگوریتم کامدوید

الگوریتم خوشه‌بندی کامدوید یکی از تکنیک‌های خوشه‌بندی قسمت‌بندی است. در این الگوریتم، مجموعه داده به چندین گروه تقسیم می‌شود، به‌طوری که عناصر هر گروه شباهت بالایی به هم دارند و با عناصر سایر گروه‌ها از شباهت کم‌تری برخوردارند. این الگوریتم در هر مرحله، از مدوید هر گروه به عنوان مرکز (نماینده) آن گروه استفاده می‌کند.

مشهورترین الگوریتم خوشه‌بندی کامدوید تحت عنوان الگوریتم قسمت‌بندی از مدوید<sup>۹</sup> یا PAM نام دارد. این الگوریتم از جستجوی حریصانه در فرایند قسمت‌بندی استفاده می‌کند و این امکان وجود دارد که به جواب بهینه منجر نشود بلکه به یک جواب بهینه محلی برسد، زیرا به جای جستجوی تمام فضای جواب، فقط بخشی از فضای جواب را جستجو می‌کند. ضمن این که همین امر موجب کاهش شدید زمان دست‌یابی به جواب هم خواهد شد. در این تحقیق از الگوریتم PAM استفاده شده است.

در صورتی که  $k$  بیان‌گر تعداد خوشه و  $n$  تعداد اعضای مجموعه داده باشد، مراحل اجرای این الگوریتم عبارت است از:

- گام ۱. به صورت تصادفی،  $k$  عضو از مجموعه داده انتخاب می‌شود. هر عضو تصادفی به عنوان یکی از مرکز خوشه خواهد بود؛
- گام ۲. سایر اعضای مجموعه داده به نزدیک‌ترین مرکز خوشه اختصاص داده می‌شود؛

گام ۳. یک عضو غیر مدوید به صورت تصادفی انتخاب می‌شود؛

گام ۴. هزینه‌ی کل ( $S$ ) که بیان‌گر هزینه‌ی ناشی از جایه‌جایی عضو غیر مدوید انتخاب شده با مدوید همان خوشه است، محاسبه می‌شود. یعنی عضو غیر مدوید به عنوان مرکز خوشه‌یی در نظر گرفته می‌شود که این عضو از آن خوشه انتخاب شده است. سپس سایر اعضاء به نزدیک‌ترین مرکز خوشه اختصاص می‌باشد. هزینه‌ی خوشه‌بندی وضعیت جدید محاسبه و از هزینه‌ی وضعیت قبلی خوشه‌بندی کسر و به عنوان هزینه‌ی کل در نظر گرفته می‌شود؛

گام ۵. اگر  $S$  منفی شده باشد، آن‌گاه جایه‌جایی عضو غیر مدوید به عنوان مدوید آن خوشه قطعی می‌شود، در غیر این صورت هیچ تغییری در ساختار خوشه‌ها اتفاق نمی‌افتد؛

گام ۶. اگر شرط توقف احراز شده باشد، عمل خوشه‌بندی خاتمه می‌پذیرد در غیر این صورت به گام ۲ مراجعه می‌شود.

شرط توقف می‌تواند یکی یا چند حالت از حالات زیر باشد:

۱. رسیدن به تعداد مشخصی جایه‌جایی؛
۲. رسیدن به تعداد مشخصی جایه‌جایی موفق؛
۳. رسیدن به درصد خاصی از کاهش در هزینه‌ی کل.

### ۴. رویکرد ارزیابی عملکرد

در این تحقیق، عملکرد روش‌های اندازه‌گیری شباهت LCSS و DTW در مورد تکنیک‌های نزدیک‌ترین همسایگی و خوشه‌بندی با الگوریتم کامدوید روی داده‌های سری زمانی محاسبه و مقایسه می‌شوند. در این تحقیق از ۶۳ مجموعه داده سری زمانی متعلق به بانک اطلاعاتی UCR (با آدرس زیر) که به‌طور تصادفی از میان آنها انتخاب شده‌اند، استفاده شده است.

([http://www.cs.ucr.edu/~eamonn/time\\_series\\_data\\_2018](http://www.cs.ucr.edu/~eamonn/time_series_data_2018))  
نام و مشخصات این مجموعه داده در جدول ۱ ارائه شده است. هر مجموعه داده دارای دو زیرمجموعه است، یکی تحت عنوان مجموعه داده آموزشی و دیگری تحت عنوان مجموعه داده آزمایشی. طول سری‌های زمانی، تعداد خوشه و کلاسه هر سری زمانی این مجموعه‌ها مشخص است. به عنوان نمونه، مجموعه داده «Statistical Control» دارای سری‌های زمانی با طول ۶۰، تعداد ۶ خوشه است و تعداد اعضای مجموعه داده آموزشی آزمایشی آن به ترتیب ۳۰۰ و ۳۰۵ سری زمانی است. از مجموعه داده آموزشی برای ارزیابی کیفیت الگوهای شناسایی شده استفاده می‌شود. مجموعه داده آزمایشی برای ارزیابی تأثیر بزرگ و تراکم - که به ترتیب بیان‌گر مراکز خوشه‌ها از بابت ویژگی‌های پراکنده و تراکم - به ترتیب بیان‌گر میزان دوری یا نزدیکی اعضای یک خوشه به یکدیگر و بیان‌گر میزان دوری یا نزدیکی مراکز خوشه‌ها با یکدیگر است - نیز از عوامل بسیار اثرگذار بر خروجی داده‌کاوی است. لذا در این تحقیق، ما نیز از گروه‌بندی این مجموعه‌های داده و در نتیجه، گروه‌بندی تابع حاصله از داده‌کاوی خودداری می‌کنیم و نتایج را به صورت کلی و از دید معیارهای ارزیابی داده‌کاوی بررسی می‌کنیم.

چنان‌که قبلاً نیز اشاره شد، در این تحقیق از تکنیک نزدیک‌ترین همسایگی به منظور شناسایی رده‌ی یک سری زمانی استفاده می‌شود. هدف از این کار ارزیابی تأثیر روش اندازه‌گیری شباهت، نتایج تکنیک چقدر و به چه سمتی تغییر می‌کنند. در این تکنیک فرض می‌شود رده‌ی سری‌های زمانی مجموعه آزمایشی مشخص نیستند و با مقایسه‌ی میزان شباهت آنها با سری‌های زمانی مجموعه آموزشی نسبت به تعیین رده‌ی آنها اقدام می‌شود. این فرایند در دو گام انجام می‌پذیرد:

- گام ۱. شباهت سری زمانی مجموعه آزمایشی ( $TS$ ) با تک تک سری‌های زمانی مجموعه آموزشی تحت روش اندازه‌گیری شباهت مورد نظر محاسبه می‌شود؛
- گام ۲. رده‌ی سری زمانی ( $TS$ ) تشخیص داده می‌شود که برابر است با رده‌ی شبیه‌ترین سری زمانی مجموعه آموزشی با آن تحت روش اندازه‌گیری شباهت مورد نظر.

جدول ۱. نام و مشخصات مجموعه داده‌های استفاده شده از بانک اطلاعاتی UCR.

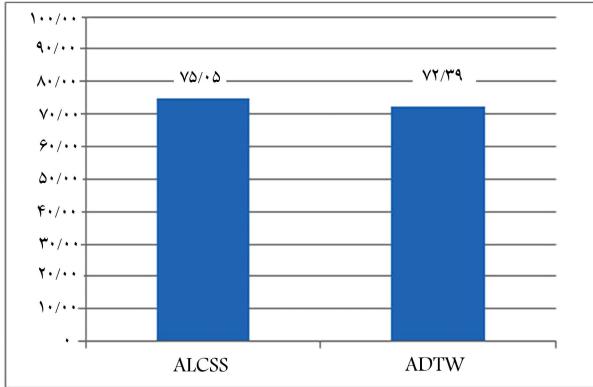
ردیف	نام مجموعه داده	تعداد خوشه	طول سری زمانی	تعداد اعضای آزمایشی	مجموعه داده آموزشی	ردیف	نام مجموعه داده	تعداد خوشه	طول سری زمانی	تعداد اعضای آزمایشی	مجموعه داده آموزشی
۱	Statistical Control	۶		۳۳	۳۰۰	۳۰۰	SonyII	۲	۶۵	۲۷	۹۵۳
۲	Gun-Point	۲		۳۴	۱۵۰	۵۰	sonySurface	۲	۷۰	۲۰	۶۰۱
۳	CBF	۳		۳۵	۹۰۰	۳۰	StarLightCurve	۳	۱۰۲۴	۱۰۰۰	۸۲۳۶
۴	ECG	۲		۳۶	۱۰۰	۱۰۰	Two Lead ECG	۲	۸۲	۲۲	۱۱۳۹
۵	Face4	۴		۳۷	۸۸	۲۴	Criket X	۱۲	۳۰۰	۳۹۰	۳۹۰
۶	Medical	۱۰		۳۸	۷۶۰	۲۸۱	Criket Y	۱۲	۳۰۰	۳۹۰	۳۹۰
۷	Sweedian	۱۵		۳۹	۶۲۵	۵۰۰	U Wave X	۸	۲۱۵	۸۹۶	۳۵۸۲
۸	OSU	۶		۴۰	۲۴۲	۲۰۰	U Wave Y	۸	۲۱۵	۸۹۶	۳۵۸۲
۹	Adiac	۳۷		۴۱	۳۹۱	۳۹۰	Insect Wing	۱۱	۲۵۶	۲۲۰	۱۹۸۰
۱۰	Beef	۵		۴۲	۳۰	۳۰	Arrow Head	۳	۲۵۱	۳۶	۱۷۵
۱۱	Lighting	۷		۴۳	۷۳	۷۰	Beetle Fly	۲	۵۱۲	۲۰	۲۰
۱۲	Fish	۷		۴۴	۱۷۵	۱۷۵	Bird Chicken	۲	۵۱۲	۲۰	۲۰
۱۳	50words	۵۰		۴۵	۴۵۵	۴۵۰	Ham	۲	۴۳۱	۱۰۹	۱۰۵
۱۴	Trace	۴		۴۶	۱۰۰	۱۰۰	Phalanges O-C	۲	۸۰	۱۸۰	۸۵۸
۱۵	Lighting7	۷		۴۷	۷۳	۷۰	Proximal POA	۳	۸۰	۴۰۰	۲۰۵
۱۶	Distal	۷		۴۸	۴۰۰	۱۳۹	Proximal POC	۲	۸۰	۶۰۰	۲۹۱
۱۷	Italy power demand	۲۴		۴۹	۱۰۲	۶۷	Proximal PT	۶	۸۰	۲۰۵	۴۰۰
۱۸	Middle-P-T	۷		۵۰	۲۹۹	۱۵۴	Toe Segmentation1	۲	۲۷۷	۴۰	۲۲۸
۱۹	Plane	۷		۵۱	۱۰۵	۱۰۵	Toe Segmentation2	۲	۲۴۳	۳۶	۱۸۰
۲۰	Car	۴		۵۲	۶۰	۶۰	Distal POA	۳	۸۰	۱۳۹	۴۰۰
۲۱	Olive Oil	۴		۵۳	۳۰	۳۰	Distal POC	۲	۸۰	۲۷۶	۶۰۰
۲۲	Diatom Size	۴		۵۴	۳۰۶	۱۶	Distal PT	۶	۸۰	۱۳۹	۴۰۰
۲۳	Symbol	۶		۵۵	۹۹۵	۲۵	Earth quakes	۲	۵۱۲	۱۳۹	۲۲۲
۲۴	Worms	۵		۵۶	۱۸۱	۷۷	Middle POA	۳	۸۰	۱۵۴	۴۰۰
۲۵	Two Pattern	۴		۵۷	۴۰۰	۱۰۰	Middle POC	۲	۸۰	۲۹۱	۶۰۰
۲۶	Wafer	۲		۵۸	۶۱۶	۱۰۰	Shapelet Sim	۲	۵۰۰	۲۰	۱۸۰
۲۷	Faceall	۱۴		۵۹	۱۶۹	۵۳۰	Wine	۲	۲۳۴	۵۷	۵۴
۲۸	Lighting2	۲		۶۰	۶۱	۶۰	Computers	۲	۷۲۰	۲۵۰	۲۵۰
۲۹	ECGFiveday	۲		۶۱	۸۶۱	۲۳	Meat	۳	۴۴۸	۶۰	۶۰
۳۰	Haptics	۵		۶۲	۳۰۸	۱۵۵	Refrigeration	۳	۷۲۰	۳۷۵	۳۷۵
۳۱	InLineSkate	۷		۶۳	۵۵۰	۱۰۰	Worm Two Class	۲	۹۰۰	۷۷	۱۸۱
۳۲	Motestrain	۷		۱۲۵	۲۰	۸۴					

پس از اجرای گام‌های فوق، شاخص «دقت» تکنیک محاسبه می‌شود. طبق تعریف، شاخص دقت عبارت است از نسبت تعداد سری‌های زمانی مجموعه آزمایشی که رده‌ی آنها به درستی تشخیص داده شده است به تعداد کل سری‌های زمانی مجموعه آزمایشی.

تکنیک دیگری که در این تحقیق از آن استفاده شده است، تکنیک خوشبندی است. هدف از این کار بررسی تأثیر روش‌های اندازه‌گیری شباهت در خوشبندی مناسب مجموعه داده است. بدین معناکه، تأثیر روش‌های اندازه‌گیری شباهت موردنظر در تعیین تعداد درست خوش، اعضا و نماینده‌ی آنها محاسبه و ارزیابی خواهد شد. از میان تکنیک‌های مختلف خوشبندی، از الگوریتم خوشبندی ندارد کامدوید استفاده می‌شود، زیرا این الگوریتم نیازی به انجام عمل میانگین‌گیری از اعضا خوشبندی در یک مرکز خوش در هر یک از مراحل خوشبندی ندارد. فرایند ارزیابی نتایج خوشبندی در دو گام انجام می‌شود.

گام ۱. خوشبندی مجموعه داده آموزشی: عمل خوشبندی با الگوریتم کامدوید

با مقادیر مختلف تعداد خوشه اجرا و برای هر مقدار از تعداد خوشه، ۵۰۰ عمل خوشبندی تکرار می‌شود و در هر بار تکرار نیز امکان جایه‌جایی نماینده خوشه به تعداد ۵۰۰ مرتبه وجود دارد. به عنوان مثال اگر تعداد خوشه برابر ۲ باشد، عمل خوشبندی به تعداد ۵۰۰ مرتبه اجرا می‌شود و در هر بار تکرار این امکان وجود دارد که برای بهتر شدن مقدار تابع هدف تا ۵۰۰ مرتبه نماینده‌های خوشه به هنگام شوند. بعد از اجرای خوشبندی، کمترین مقدار تابع هدف در هر حالت (مقدار تعداد خوشه) به عنوان بهترین وضعیت آن حالت در نظر گرفته می‌شود و مشخصات خوشبندی بهینه شامل نماینده خوشه‌ها و اعضا آنها ذخیره می‌شود. به عنوان مثال، مشخص می‌شود که بهترین مقدار تابع هدف خوشبندی در حالتی که تعداد خوشه برابر ۲ است در کدام وضعیت خوشبندی (نماینده و اعضای خوشه) ایجاد شده است. همین روال برای تعیین بهترین مقدار تابع هدف خوشبندی در حالتی که تعداد خوشه برابر ۲ است در کدام وضعیت خوشبندی کام می‌شود. سپس با بهکارگیری شاخص زاویه<sup>۱</sup> بهترین مقدار تعداد خوشه از میان تعداد خوشه‌های ممکن انتخاب خواهد شد. در انتهای نیز «دقت» خوشبندی که



شکل ۱. متوسط کل دقت تکنیک نزدیک‌ترین همسایگی تحت روش‌های LCSS و DTW

برای مقایسه دقت این تکنیک تحت روش‌های فوق الذکر از آزمون مقایسات جفتی  $t$  استفاده می‌کنیم تا از این طریق و دقت این تکنیک برای یک مجموعه داده سری زمانی تحت دو روش مورد نظر بتوانیم تشخیص دهیم که تأثیر کدامیک در تکنیک نزدیک‌ترین همسایگی بهتر است.

بنابراین اگر  $A_{LCSS}$  بیان‌گر دقت تکنیک نزدیک‌ترین همسایگی تحت روش DTW و  $A_{DTW}$  بیان‌گر دقت تکنیک نزدیک‌ترین همسایگی تحت روش LCSS باشد، آن‌گاه فرضیات این آزمون به شرح زیر خواهد بود.

$$A_{LCSS} = A_{DTW} : H_0$$

$$A_{LCSS} > A_{DTW} : H_1$$

نتایج این آزمون در جدول ۳ و برای سطح مختلطی از قابلیت اطمینان نمایش داده شده است. براساس این نتایج، می‌توان با  $92.5\%$  اطمینان ادعای کرد که دقت تکنیک نزدیک‌ترین همسایگی تحت روش LCSS از دقت این تکنیک تحت روش DTW بهتر است.

این مقایسه را می‌توان، از منظر آمار توصیفی نیز انجام داد که در جدول ۴ نمایش داده شده است. براساس اطلاعات متدرج در این جدول، دقت تکنیک نزدیک‌ترین همسایگی تحت روش LCSS نسبت به دقت این تکنیک تحت روش DTW در  $71.42\%$  حالات بهتر، در  $1.58\%$  از حالات برابر و فقط در  $26.98\%$  حالات بدتر است.

به طور خلاصه این بررسی‌ها نشان می‌دهد که روش LCSS نسبت به روش DTW تأثیر بهتری در دستیابی به دقت بالاتر دارد.

## ۲.۵. نتایج اجرای الگوریتم خوشبندی کامدوید

در این قسمت نتایج حاصل از اجرای الگوریتم خوشبندی کامدوید تحت روش‌های LCSS و DTW طبق فرایند ارزیابی توضیح داده شده در قسمت ۴ ارائه خواهد شد. این نتایج در دو قسمت، نتایج خوشبندی و نتایج گروه‌بندی ارائه می‌شود. هدف از این ارزیابی پاسخ به سؤالات زیر است:

۱. آیا دقت تکنیک خوشبندی تحت روش LCSS با دقت این تکنیک تحت روش DTW مقاوتمانند است؟

۲. آیا تأثیر روش LCSS با تأثیر روش DTW در تعیین تعداد خوشبندی حاصل از خوشبندی مقاوتمانند است؟

عبارت است از نسبت تعداد سری‌های زمانی که به درستی به خوشبندی درست خود اختصاص یافته‌اند به کل سری‌های زمانی مجموعه داده آموزشی، محاسبه می‌شود. خروجی این گام، تعیین بهترین مقدار تعداد خوشبندی برای مجموعه داده، تعیین دقت و مشخصات بهترین خوشبندی است:

گام ۲. گروه‌بندی مجموعه داده آزمایشی: در این گام از نتایج گام اول (تعداد خوشبندی و نماینده آنها) برای گروه‌بندی مجموعه داده آزمایشی استفاده می‌شود. بدین معنا که، شبیه‌ترین نماینده خوشبندی برای مجموعه داده آزمایشی تعیین می‌شود و سری زمانی مذکور به آن خوشبندی اختصاص داده می‌شود. در انتهای «دقت» گروه‌بندی که عبارت است از نسبت تعداد سری‌های زمانی بی که به درستی به خوشبندی درست خود اختصاص یافته‌اند به کل سری‌های زمانی مجموعه داده آزمایشی، محاسبه می‌شود. خروجی گام دوم، ارزیابی مناسب بودن کیفیت نماینده خوشبندی حاصل از گام اول خوشبندی است.

## ۵. نتایج ارزیابی و بحث

### ۵.۱. نتایج اجرای الگوریتم نزدیک‌ترین همسایگی

نتایج تکنیک نزدیک‌ترین همسایگی تحت مقادیر مختلف حد مجاز شباخت روش LCSS و DTW در جدول ۲ ارائه شده است. به عنوان مثال در مورد مجموعه داده‌ی «Statistical Control» وقتی حد مجاز شباخت ( $\epsilon$ ) برابر  $0.5$  است آن‌گاه رده‌ی  $69.67\%$  از سری‌های زمانی آزمایشی به درستی شناسایی شده است. دقت این تکنیک برای این مجموعه داده با تغییر مقدار حد مجاز شباخت تغییر می‌کند، به طوری که با افزایش مقدار آن از  $0.5$  تا  $0.35$ ، منجر به افزایش دقت تکنیک از  $69.67\%$  تا  $94.67\%$  می‌شود. متوسط و انحراف معیار مقادیر این دقت‌ها برای این مجموعه به ترتیب برابر  $95.95\%$  و  $8.63\%$  است. روند تأثیر تغییرات حد مجاز شباخت بر دقت این تکنیک برای تمامی مجموعه‌های داده به صورتی که اکنون توضیح داده شد نیست بلکه به عنوان مثال، این تأثیر در مورد مجموعه داده Gun با افزایش مقدار حد مجاز شباخت، در ابتدا افزایشی و سپس کاهشی خواهد بود، به طوری که بیشینه‌ی آن در مقدار حد مجاز شباخت  $0.15$  رخ می‌دهد و متوسط و انحراف معیار دقت آن نیز به ترتیب برابر  $93.33\%$  و  $5.05\%$  است. به طور کل، متوسط دقت و انحراف معیار این تکنیک برای همه مجموعه داده‌ها به ترتیب برابر  $93.41\%$  و  $4.80\%$  است، به طوری که این موضوع بیان‌گر وابستگی دقت این تکنیک به مقدار حد مجاز شباخت است. لذا انتخاب نادرست مقدار دقت مجاز می‌تواند منجر به نتایج نامناسبی شود. شایان ذکر است که بیشترین مقدار دقت متوسط کل در حد مجاز شباخت  $0.25$  رخ می‌دهد و برابر با  $75.05\%$  است.

در جدول ۲، دقت این تکنیک تحت روش DTW نیز ارائه شده است. به عنوان مثال در مورد مجموعه داده «Statistical Control»، با دقت  $0.33\%$ ، رده‌ی سری‌های زمانی آزمایشی به درستی شناسایی شده است. اما متأسفانه، این دقت برخی مجموعه‌های داده بسیار نامطلوب است، نظریه مجموعه داده‌ی OSU Middle-P-T به ترتیب با مقادیر دقت  $0.46\%$  و  $0.58\%$  به عبارت دیگر عملکرد روش DTW در برخی از مجموعه‌های داده بسیار خوب و در برخی دیگر کلیه مجموعه‌های داده تحت روش DTW برابر  $0.28\%$  است. نمودار ۱ بیان‌گر متوسط دقت تکنیک نزدیک‌ترین همسایگی تحت روش‌های LCSS و DTW است.

جدول ۲. دقت تکنیک نزدیکترین همسایگی تحت روش‌های LCSS و DTW (برحسب درصد).

ردیف	نام مجموعه داده	تحت روش LCSS با مقادیر مختلف حد مجاز شباهت									تحت روش DTW	
		حد مجاز شباهت ( $\epsilon$ )										
		۰,۳۵	۰,۳۰	۰,۲۵	۰,۲۰	۰,۱۵	۰,۱۰	۰,۰۵				
۱	statistical control	۹۷,۳۳	۹۳	۸,۶۳	۸۷,۹۶	۹۴,۶۷	۹۳	۹۳	۹۱,۶۷	۸۷	۸۶,۶۷	۶۹,۶۷
۲	Gun-Point	۹۰,۶۷	۹۲,۶۷	۵,۰,۰	۹۳,۳۳	۸۵,۳۳	۸۸,۶۷	۹۲,۶۷	۹۷,۳۳	۹۸,۶۷	۹۸	۹۲,۶۷
۳	CBF	۹۹	۹۹,۶۷	۱,۱۲	۹۹,۰,۰	۹۹,۸۹	۹۹,۷۸	۹۹,۶۷	۹۹,۴۴	۹۹,۰۶	۹۸,۱۱	۹۶,۸۹
۴	ECG	۷۹	۹۰	۴,۹۲	۸۶,۷۱	۹۱	۹۱	۹۰	۸۷	۸۶	۸۵	۷۷
۵	Face4	۸۱,۸۲	۹۴,۵	۴,۰,	۹۱,۱۰	۹۴,۳۲	۹۰,۰,۰	۹۴,۰	۹۳,۱۸	۹۲,۰,۰	۸۹,۷۷	۸۱,۸۲
۶	Medical	۶۵,۳۹	۶۲,۳۳	۲,۲۶	۶,۰,۱۲	۵۱,۷۱	۶۳,۱۶	۶۲,۳۳	۶۱,۱۲	۵۷,۷۶	۵۲,۹۲	
۷	Sweedian	۷۲,۱۶	۸۴,۶۷	۱۰,۱۱	۷۶,۶۶	۸۲,۷۲	۸۴,۹۶	۸۴,۶۲	۸۳,۰,۰	۸۲,۴	۷۵,۲	۴۳,۲
۸	OSU	۴۶,۲۸	۶۹,۸۳	۲,۴۸	۶۷,۷۱	۶۸,۱۸	۶۸,۱۸	۶۹,۸۳	۶۹,۴۲	۶۸,۶	۶۷,۳۶	۶۲,۴
۹	Adiac	۷۴,۰,۳	۴۲,۷	۱۸,۲۰	۰,۰,۱۷	۲۶,۲۳	۳۳,۳۸	۴۲,۷	۴۸,۸۶	۵۰,۱۹	۶۶,۸۸	۷۷,۹۲
۱۰	Beef	۶۳,۳۳	۷۳	۷,۸۹	۶۶,۶۲	۵۶,۶۷	۶۳,۳۳	۷۳	۷۶,۶۷	۷۰	۷۰	۵۶,۶۷
۱۱	Lighting	۶۸,۴۹	۷۵,۳۴	۷,۳۳	۶۸,۱۱	۶۸,۴۹	۷۵,۳۴	۷۵,۳۴	۷۱,۳۳	۶۹,۸۶	۵۶,۱۶	۶۰,۲۷
۱۲	Fish	۷۵,۴۳	۸۱,۶۹	۰,۷۷	۸۲,۵۳	۷۳,۱۴	۷۷,۱۴	۸۱,۶۹	۸۴,۵۷	۸۹,۱۴	۸۸	۸۴
۱۳	50words	۶۶,۳۹	۷۳,۹۵	۰,۶۱	۷۱,۶۹	۷۷,۰,۰	۷۵,۰,۰	۷۳,۹۵	۷۳,۵۱	۷۲,۱۱	۶۸,۳۵	۶۰,۰,۵
۱۴	Trace	۱۰۰	۹۸	۲,۸۲	۹۶,۵۷	۹۶,۵۷	۹۶,۵۷	۹۶,۵۷	۹۶,۵۷	۹۰	۹۷	۹۲
۱۵	Lighting7	۶۸,۴۹	۷۵,۳۴	۷,۳۳	۶۸,۱۱	۶۸,۴۹	۷۵,۳۴	۷۵,۳۴	۷۱,۲۳	۶۹,۸۶	۵۶,۱۶	۶۰,۲۷
۱۶	Distal	۷۰,۰,۰	۷۵,۰,۵	۲,۳۴	۷۴,۳۹	۷۶,۲۵	۷۶	۷۵,۰,۵	۷۵,۲۵	۷۴	۷۳,۰,۵	۶۹,۰,۵
۱۷	Italy power demand	۹۳,۹۷	۹۱,۷۴	۰,۷۳	۸۸,۰,۰	۹۲,۷۲	۹۲,۴۲	۹۱,۷۴	۹۰,۱۸	۸۷,۰,۰	۸۲	۷۸,۰,۰
۱۸	Middle-P-T	۵۸,۰,۰	۶۲,۴۱	۲,۱۱	۶۰,۰,۰	۶۱,۱۵	۶۱,۹	۶۲,۴۱	۵۹,۹	۵۶,۰,۰	۵۱,۰,۰	۵۱,۰,۰
۱۹	Plane	۱۰۰	۱۰۰	۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰	۱۰۰
۲۰	Car	۷۱,۶۷	۸۱,۶۷	۰,۰,۰	۸۱,۴۴	۷۲,۳۳	۷۲,۳۳	۸۱,۶۷	۸۳,۰,۰	۸۰	۸۰	۸۰,۰,۰
۲۱	Olive Oil	۸۳,۳۳	۴۵	۱۳,۰,۱	۵۰,۴۸	۴۰	۴۰	۴۵	۴۱,۰,۰	۵۰,۰,۰	۵۶,۰,۰	۵۶,۰,۰
۲۲	Diatom Size	۹۶,۴۱	۹۳,۴۵	۰,۱۹	۹۲,۶۷	۱۱,۳۷	۹۱,۸۳	۹۳,۴۵	۹۴,۴۴	۹۶,۴۱	۹۵,۷۵	۹۵,۰,۰
۲۳	Symbol	۹۵,۲۷	۹۵,۴۸	۱,۰,۰	۹۵,۱۰	۹۲,۴۶	۹۳,۷۷	۹۵,۴۸	۹۶,۸۸	۹۶,۲۷	۹۰,۰,۰	۹۰,۰,۰
۲۴	Worms	۴۵,۳۴	۵۱,۳۷	۳,۰,۰	۴۸,۶۲	۵۰,۰,۰	۵۱,۹۳	۵۱,۳۷	۴۸,۰,۰	۴۵,۲۹	۴۴,۱۸	
۲۵	Two Pattern	۱۰۰	۹۹,۲	۲,۰,۰	۹۸,۰,۰	۹۹,۰,۰	۹۹,۰,۰	۹۹,۰,۰	۹۹,۰,۰	۹۹,۰,۰	۹۳,۰,۰	۹۳,۰,۰
۲۶	Wafer	۹۸,۰,۰	۹۸,۷۲	۰,۰,۰	۹۸,۰,۰	۹۸,۰,۰	۹۸,۰,۰	۹۸,۰,۰	۹۸,۰,۰	۹۸,۰,۰	۹۸,۰,۰	۹۸,۰,۰
۲۷	Faceall	۸۴,۷۴	۸۸,۱۷	۱۲,۰,۰	۸۱,۰,۰	۹۳,۰,۰	۹۳,۰,۰	۸۸,۱۷	۸۴,۰,۰	۷۸,۰,۰	۷۳,۰,۰	۶۵,۰,۰
۲۸	Lighting2	۸۶,۰,۰	۷۸,۶۹	۰,۰,۰	۷۶,۰,۰	۷۷,۰,۰	۷۸,۶۹	۷۸,۶۹	۷۸,۶۹	۷۷,۰,۰	۷۷,۰,۰	۷۰,۰,۰
۲۹	ECGFiveday	۷۵,۰,۰	۸۱,۴۲	۰,۰,۰	۷۹,۰,۰	۸۶,۰,۰	۸۶,۰,۰	۸۱,۴۲	۸۰,۰,۰	۷۶,۰,۰	۷۶,۰,۰	۷۰,۰,۰
۳۰	Haptics	۴۱,۲۳	۴۲,۱۸	۲,۰,۰	۴۱,۷۹	۴۱,۲۳	۴۱,۲۳	۴۱,۱۸	۴۱,۱۸	۴۰,۰,۰	۳۹,۰,۰	۳۹,۰,۰
۳۱	InLineSkate	۳۸,۰,۰	۳۲,۲۷	۱,۰,۰	۳۷,۰,۰	۳۷,۰,۰	۳۷,۰,۰	۳۷,۰,۰	۳۷,۰,۰	۳۶,۰,۰	۳۵,۰,۰	۳۵,۰,۰
۳۲	Motestrain	۸۱,۰,۰	۹۲,۸۱	۰,۰,۰	۹۱,۷۸	۹۱,۰,۰	۹۱,۰,۰	۹۱,۰,۰	۹۱,۰,۰	۹۱,۰,۰	۹۱,۰,۰	۹۱,۰,۰
۳۳	SonyII	۸۴,۰,۰	۸۲,۴۸	۰,۰,۰	۸۴,۰,۰	۸۳,۰,۰	۸۳,۰,۰	۸۳,۰,۰	۸۳,۰,۰	۸۳,۰,۰	۸۳,۰,۰	۸۳,۰,۰
۳۴	sonySurface	۶۰,۰,۰	۶۶,۰,۰	۰,۰,۰	۶۷,۰,۰	۶۸,۰,۰	۶۸,۰,۰	۶۸,۰,۰	۶۸,۰,۰	۶۸,۰,۰	۶۸,۰,۰	۶۸,۰,۰
۳۵	StarLightCurve	۶۶,۰,۰	۷۸	۰,۰,۰	۷۷,۰,۰	۷۷	۷۷	۷۸	۷۸,۰,۰	۷۸,۰,۰	۷۸,۰,۰	۷۸,۰,۰
۳۶	Two Lead ECG	۸۹,۰,۰	۸۶,۸۳	۰,۰,۰	۸۵,۰,۰	۸۱,۰,۰	۸۱,۰,۰	۸۱,۰,۰	۸۱,۰,۰	۸۱,۰,۰	۸۱,۰,۰	۸۱,۰,۰
۳۷	Criket X	۶۹,۰,۰	۴۵,۰,۰	۰,۰,۰	۴۴,۰,۰	۴۵,۰,۰	۴۵,۰,۰	۴۵,۰,۰	۴۵,۰,۰	۴۵,۰,۰	۴۵,۰,۰	۴۵,۰,۰
۳۸	Criket Y	۵۶,۰,۰	۶۰,۰,۰	۰,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰
۳۹	U Wave X	۶۰,۰,۰	۵۹,۰,۰	۰,۰,۰	۵۸,۰,۰	۵۹,۰,۰	۵۹,۰,۰	۵۹,۰,۰	۵۹,۰,۰	۵۹,۰,۰	۵۹,۰,۰	۵۹,۰,۰
۴۰	U Wave Y	۴۰,۰,۰	۴۳,۱۹	۰,۰,۰	۴۴,۱۰	۴۲,۲۶	۴۰	۴۳,۱۹	۴۰,۰,۰	۴۰,۰,۰	۴۰,۰,۰	۴۰,۰,۰
۴۱	Insect Wing	۲۷,۰,۰	۵۰,۰,۰	۱,۰,۰	۴۷,۰,۰	۵۰,۰,۰	۴۷,۰,۰	۴۷,۰,۰	۴۷,۰,۰	۴۷,۰,۰	۴۷,۰,۰	۴۷,۰,۰
۴۲	Arrow Head	۷۰,۰,۰	۷۰,۰,۰	۰,۰,۰	۶۷,۰,۰	۶۷,۰,۰	۶۷,۰,۰	۶۷,۰,۰	۶۷,۰,۰	۶۷,۰,۰	۶۷,۰,۰	۶۷,۰,۰
۴۳	Beetle Fly	۷۰	۸۰	۰	۸۰	۸۰	۸۰	۸۰	۸۰	۸۰	۸۰	۸۰
۴۴	Bird Chicken	۷۵	۹۰	۰	۸۰	۸۰	۹۰	۸۰	۸۰	۸۰	۸۰	۸۰
۴۵	Ham	۴۲,۸۶	۷۴,۲۹	۰,۰,۰	۷۴,۳۳	۶۵,۷۱	۶۸,۰,۰	۷۴,۲۹	۷۱,۴۳	۶۵,۷۱	۶۵,۷۱	۶۰
۴۶	Phalanges O-C	۶۹,۳۳	۷۱,۰,۰	۱,۰,۰	۶۹,۲۳	۶۸,۱۸	۷۱,۰,۰	۶۹,۱۰	۶۹,۱۰	۶۸,۳۰	۶۶,۳۲	
۴۷	Proximal POA	۷۷,۶۴	۷۹,۰,۱	۰,۰,۰	۷۸,۰,۰	۷۷,۰,۰	۷۹,۰,۱	۷۹,۰,۱	۷۸,۰,۰	۷۸,۰,۰	۷۸,۰,۰	۷۸,۰,۰
۴۸	Proximal POC	۸۱,۶۳	۹۱,۰,۵	۰,۰,۰	۸۱,۷۲	۸۱,۰,۰	۸۱,۷۲	۸۱,۰,۰	۸۱,۰,۰	۸۱,۰,۰	۸۱,۰,۰	۸۱,۰,۰
۴۹	Proximal PT	۷۷,۰,۰	۷۸,۰,۰	۰,۰,۰	۷۸,۰,۰	۷۸,۰,۰	۷۸,۰,۰	۷۸,۰,۰	۷۸,۰,۰	۷۸,۰,۰	۷۸,۰,۰	۷۸,۰,۰
۵۰	Toe Segmentation1	۷۶,۰,۰	۸۷,۰,۰	۰,۰,۰	۸۰,۰,۰	۸۰,۰,۰	۸۰,۰,۰	۸۰,۰,۰	۸۰,۰,۰	۷۸,۰,۰	۷۸,۰,۰	۷۸,۰,۰
۵۱	Toe Segmentation2	۸۶,۰,۰	۹۶,۰,۰	۰,۰,۰	۹۵,۰,۰	۹۵,۰,۰	۹۵,۰,۰	۹۵,۰,۰	۹۵,۰,۰	۹۳,۰,۰	۹۲,۰,۰	۹۲,۰,۰
۵۲	Distal POA	۸۳,۰,۰	۷۷,۰,۰	۰,۰,۰	۷۷,۰,۰	۷۷,۰,۰	۷۷,۰,۰	۷۷,۰,۰	۷۷,۰,۰	۷۷,۰,۰	۷۷,۰,۰	۷۷,۰,۰
۵۳	Distal POC	۷۲,۰,۰	۷۱,۰,۰	۰,۰,۰	۷۰,۰,۰	۷۰,۰,۰	۷۰,۰,۰	۷۰,۰,۰	۷۰,۰,۰	۷۰,۰,۰	۷۰,۰,۰	۷۰,۰,۰
۵۴	Distal PT	۶۲,۰,۰	۶۱,۰,۰	۰,۰,۰	۶۰,۰,۰	۶۰,۰,۰	۶۰,۰,۰	۶۰,۰,۰	۶۰,۰,۰	۶۰,۰,۰	۶۰,۰,۰	۶۰,۰,۰
۵۵	Earth quakes	۵۷,۰,۰	۵۷	۰,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰
۵۶	Middle POA	۵۷,۰,۰	۵۷	۰,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰	۵۷,۰,۰
۵۷	Middle POC	۷۹	۸۰,۳۳	۰,۰,۰	۷۹,۰,۰	۸۲,۰,۰	۸۰,۰,۰	۸۰,۰,۰	۷۷,۱۷	۷۷,۱۷	۷۷,۱۷	۷۷,۱۷
۵۸	Shapelet Sim	۵۶,۱۱	۹۱,۰,۰	۱,۰,۰	۸۰,۰,۰	۹۲,۰,۰	۹۲,۰,۰	۹۱,۰,۰	۸۸,۰,۰	۸۶,۰,۰	۷۳,۰,۰	۷۳,۰,۰
۵۹	Wine	۵۷,۰,۰	۵۰	۰,۰,۰	۵۰	۰,۰,۰	۵۰	۴۸,۰,۰	۴۸,۰,۰	۴۸,۰,۰	۴۸,۰,۰	۴۸,۰,۰
۶۰	Computers	۶۲	۳,۰,۰	۰,۰,۰	۶۲	۶۲	۶۲	۶۲	۶۲	۶۲	۶۲	۶۲

جدول ۳. نتایج آزمون مقایسات جفتی  $t$  برای دقت تکنیک نزدیک‌ترین همسایگی تحت روش‌های LCSS و DTW.

Paired Differences		<i>T</i>	<i>Df</i>	sig. (2 - tailed)	Correlation
Mean	Std. Deviation				
$A_{LCSS} - A_{DTW}$	۲,۶۵۹۲۵	۱۲,۹۳۵۵۵	۱,۶۳۲	.۰۱۰	$A_{LCSS} & A_{DTW} = .۷۲۸$
Confidence Interval of the Difference					
		٪۹۰	٪۹۲,۵	٪۹۵	٪۹۷,۵
$A_{LCSS} - A_{DTW}$	Lower	۰,۵۴۸۱۸	۰,۲۸۳۷۸	-۰,۰۶۲۰۷	-۰,۰۵۹۸۵۲
					-۱,۲۳۲۵۵

بنابراین، اگر  $A_{LCSS}$  بیان‌گر دقت تکنیک خوشبندی تحت روش LCSS و  $A_{DTW}$  بیان‌گر دقت تکنیک خوشبندی تحت روش DTW باشد، آنگاه فرضیات این آزمون به شرح زیر خواهد بود.

$$A_{LCSS} = A_{DTW} : H_0$$

$$A_{LCSS} > A_{DTW} : H_1$$

نتایج این آزمون در جدول ۶ و برای سطوح مختلفی از قابلیت اطمینان نمایش داده شده است. بر اساس این نتایج، می‌توان حداقل با ٪۹۹ اطمینان ادعا کرد که دقت خوشبندی تحت روش LCSS از دقت این تکنیک تحت روش DTW بهتر است.

این مقایسه را می‌توان از منظر آمار توصیفی نیز انجام داد. نتایج این مقایسه در جدول ۷ نمایش داده شده است. بر اساس اطلاعات مندرج در این جدول، دقت خوشبندی تحت روش LCSS نسبت به دقت خوشبندی تحت روش DTW در ٪۵۳,۹۷ حالت بهتر، در ٪۱۲,۷۰ از حالات برابر و فقط در ٪۳۲,۳۳ حالت بدتر است. پهلو خلاصه این بررسی نشان می‌دهد که روش LCSS تأثیر بهتری در دست یابی به دقت بالاتر نسبت به روش DTW دارد.

در ادامه و در جدول ۸، تأثیر روش‌های LCSS و DTW در تعیین درست تعداد خوشبندی نشان داده شده است. از اطلاعات این جدول برای پاسخ به سؤال دوم استفاده می‌کنیم. بر اساس این اطلاعات، در خوشبندی تحت روش DTW و با بهکارگیری شاخص زاویه، می‌توان تعداد خوشبندی را برای ۲۸ مجموعه داده به درستی و برای ۱۷ مجموعه داده با اختلاف یک عدد تشخیص داد. این مقادیر به ترتیب معادل ٪۴۴,۴۴ و ٪۲۷ (جمعاً برابر ٪۷۱,۴۴) از کل مجموعه داده‌ها می‌شود. از سوی دیگر، در خوشبندی تحت روش LCSS و با بهکارگیری شاخص زاویه می‌توان تعداد خوشبندی را برای ۲۵ مجموعه داده به درستی و برای ۱۹ مجموعه داده با اختلاف یک عدد تشخیص داد. این مقادیر نیز به ترتیب معادل ٪۴۹,۶۸ و ٪۳۰,۱۵ (جمعاً برابر ٪۶۹,۸۳) از کل مجموعه داده‌ها هستند. این مطالب نشان از تأثیر نسبتاً برابر این روش‌ها در شناسایی تعداد درست خوشبندی دارد.

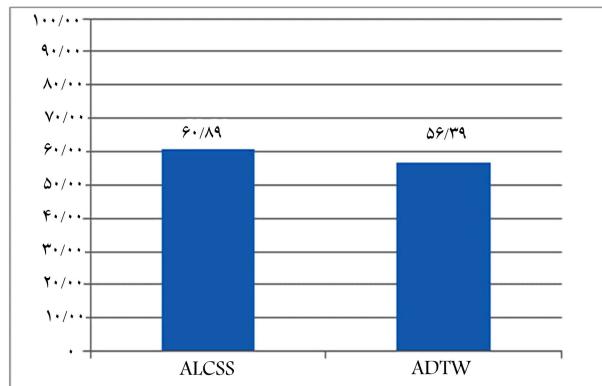
به طور خلاصه، روش LCSS تأثیر بهتری در تعیین اعضای هر خوشبندی به تأثیر روش DTW دارد، ولی از نظر تأثیر در تعیین درست تعداد خوشبندی، تقریباً از اثر یکسانی برخوردارند.

#### ۲.۵ نتایج گروه‌بندی مجموعه داده آزمایشی

چنان‌که قبلاً توضیح داده شد، پس از خوشبندی مجموعه داده آموزشی و تعیین بهترین تعداد خوشبندی و به تبع آن عناصر خوشبندی و نماینده هر خوشبندی، نسبت به گروه‌بندی مجموعه داده آزمایشی مبتنی بر این نتایج اقدام می‌شود. بدین معنا که اعضای

جدول ۴. مقایسه‌ی آماری دقت تکنیک نزدیک‌ترین همسایگی تحت روش‌های LCSS و DTW.

DTW vs LCSS			
ردیف	توضیحات	تعداد	درصد
۱	$A_{LCSS} > A_{DTW}$	۴۵	٪۷۱,۴۲
۲	$A_{LCSS} = A_{DTW}$	۱	٪۱,۵۸
۳	$A_{LCSS} < A_{DTW}$	۱۷	٪۲۶,۹۸



نمودار ۲. متوسط کل دقت خوشبندی مجموعه داده آموزشی تحت روش‌های LCSS و DTW

۳. آیا تأثیر روش LCSS با تأثیر روش DTW در تعیین نماینده خوشبندی حاصل از خوشبندی مقاومت دارد؟

#### ۲.۵ نتایج خوشبندی مجموعه آموزشی

جدول ۵ نتایج اجرای تکنیک خوشبندی کامدوید را تحت روش‌های DTW و LCSS نشان می‌دهد. بر اساس این نتایج و به عنوان مثال برای مجموعه داده‌ی «Statistical Control»، بهترین تعداد خوشبندی تحت روش DTW برابر ۶ خواهد بود که دارای دقت خوشبندی ٪۹۷,۶۷ است. بهترین تعداد خوشبندی تحت روش LCSS نیز برابر ۶ و با دقت خوشبندی ٪۸۷,۳۳ خواهد بود. یعنی این تکنیک تحت هر دو روش فوق توانسته است تعداد خوشبندی را به درستی تعیین کند.

با نگاه مجدد به نتایج مندرج در جدول ۵، مشاهده می‌شود که دقت خوشبندی تحت روش‌های فوق، متفاوت است که برای نتیجه‌گیری درست باید به طور آماری بررسی شود. از سوی دیگر، دقت تکنیک خوشبندی برای کلیه مجموعه‌های داده تحت روش DTW برابر ٪۵۶,۵۵ و تحت روش LCSS برابر ٪۵۶,۸۹ است. نمودار ۲ بیان‌گر متوسط دقت تکنیک نزدیک‌ترین همسایگی تحت روش‌های LCSS و DTW است. به منظور پاسخ به سؤال اول، از آزمون مقایسات زوجی  $t$  استفاده می‌شود.

جدول ۵. نتایج خوشبندی کامدوید (تعداد خوشبندی بهینه و دقت خوشبندی بر حسب درصد روی مجموعه آموزشی).

ردیف	نام مجموعه داده	تحت روش DTW				تحت روش LCSS				ردیف
		تعداد خوشبندی	دقت خوشبندی							
۱	statistical control	۹۷,۶۷	۶	۸۷,۲۳	۶	۹۷,۶۷	۶	۹۷,۶۰	۴	۷۸,۰۷
۲	Gun-Point	۵۶	۲	۵۶	۴	۵۶	۴	۶۰	۳	۶۰
۳	CBF	۹۶,۶۷	۳	۹۶,۶۷	۳	۹۶,۶۷	۳	۶۰	۴	۶۰
۴	ECG	۷۳	۲	۷۳	۲	۷۳	۲	۷۳	۴	۷۳
۵	Face4	۶۶,۶۷	۳	۸۳,۳۳	۵	۶۶,۶۷	۳	۶۰	۱۱	۶۰
۶	Medical	۳۲,۲۹	۶	۳۱,۷۶	۶	۳۲,۲۹	۶	۳۲,۲۹	۱۳	۳۲,۲۹
۷	Sweedian	۵۲,۴۰	۱۹	۶۶,۴۰	۱۹	۵۲,۴۰	۱۹	۵۲,۴۰	۷	۵۲,۴۰
۸	OSU	۴۷,۵	۵	۴۷,۵	۵	۴۷,۵	۵	۴۷,۵	۱۱	۴۷,۵
۹	Adiac	۴۲,۳۱	۳۰	۴۷,۱۷	۳۰	۴۲,۳۱	۳۰	۴۲,۳۱	۲	۴۲,۳۱
۱۰	Beef	۴۳,۳۳	۶	۴۶,۶۷	۶	۴۳,۳۳	۶	۴۳,۳۳	۲	۴۳,۳۳
۱۱	Lighting	۴۵,۷۱	۶	۴۵,۷۱	۶	۴۵,۷۱	۶	۴۵,۷۱	۱۱	۴۵,۷۱
۱۲	Fish	۴۹,۵۷	۹	۷۲,۰۷	۹	۴۹,۵۷	۹	۴۹,۵۷	۳	۴۹,۵۷
۱۳	50words	۴۸,۲۲	۷۱	۴۸,۲۲	۷۱	۴۸,۲۲	۷۱	۴۸,۲۲	۱۳	۴۸,۲۲
۱۴	Trace	۵۲	۳	۵۲	۲	۵۲	۲	۵۲	۴	۵۲
۱۵	Lighting7	۵۵,۷۱	۶	۵۵,۷۱	۶	۵۵,۷۱	۶	۵۵,۷۱	۱۵	۵۵,۷۱
۱۶	Distal	۶۷,۶۳	۳	۶۷,۶۳	۳	۶۷,۶۳	۳	۶۷,۶۳	۲	۶۷,۶۳
۱۷	Italy power demand	۶۷,۱۶	۳	۶۷,۱۶	۳	۶۷,۱۶	۳	۶۷,۱۶	۱۷	۶۷,۱۶
۱۸	Middle-P-T	۵۵,۱۵	۲	۵۵,۱۵	۲	۵۵,۱۵	۲	۵۵,۱۵	۱۸	۵۵,۱۵
۱۹	Plane	۱۰۰	۷	۱۰۰	۷	۱۰۰	۷	۱۰۰	۲	۱۰۰
۲۰	Car	۴۵	۴	۴۵	۴	۴۵	۴	۴۵	۲	۴۵
۲۱	Olive Oil	۴۳,۱۴	۴	۴۳,۱۴	۴	۴۳,۱۴	۴	۴۳,۱۴	۱۱	۴۳,۱۴
۲۲	Diatom Size	۶۲,۳۵	۲	۶۲,۳۵	۲	۶۲,۳۵	۲	۶۲,۳۵	۲	۶۲,۳۵
۲۳	Symbol	۱۰۰	۶	۱۰۰	۶	۱۰۰	۶	۱۰۰	۲	۱۰۰
۲۴	Worms	۲۲,۰۱	۲	۲۲,۰۱	۲	۲۲,۰۱	۲	۲۲,۰۱	۲	۲۲,۰۱
۲۵	Two Pattern	۶۰,۴۹	۴	۶۰,۴۹	۴	۶۰,۴۹	۴	۶۰,۴۹	۲	۶۰,۴۹
۲۶	Wafer	۶۰,۴۸	۲	۶۰,۴۸	۲	۶۰,۴۸	۲	۶۰,۴۸	۲	۶۰,۴۸
۲۷	Faceall	۷۸,۳۹	۱۵	۷۸,۳۹	۱۵	۷۸,۳۹	۱۵	۷۸,۳۹	۲	۷۸,۳۹
۲۸	Lighting2	۴۸,۳۳	۶	۴۸,۳۳	۶	۴۸,۳۳	۶	۴۸,۳۳	۲	۴۸,۳۳
۲۹	ECGFiveday	۶۰,۸۷	۲	۶۰,۸۷	۲	۶۰,۸۷	۲	۶۰,۸۷	۲	۶۰,۸۷
۳۰	Haptics	۴۲,۵۸	۲	۴۲,۵۸	۲	۴۲,۵۸	۲	۴۲,۵۸	۲	۴۲,۵۸
۳۱	InLineSkate	۴۱,۶۷	۳۱	۴۱,۶۷	۳۱	۴۱,۶۷	۳۱	۴۱,۶۷	۲	۴۱,۶۷
۳۲	Motestrain	۴۵	۳	۴۵	۳	۴۵	۳	۴۵	۲	۴۵

جدول ۶. نتایج آزمون مقایسات جفتی t برای دقت تکنیک خوشبندی تحت روش‌های LCSS و DTW.

	Paired Differences		T	Df	sig. (۲ - tailed)	Correlation
	Mean	Std.Deviation				
A <sub>LCSS</sub> – A <sub>DTW</sub>	۴,۳۴۶۳۵	۱۲,۴۳۴۴۰	۲,۷۷۴	۶۲	.۰۱۰	A <sub>LCSS</sub> &A <sub>DTW</sub> = .۰۷۸۹
A <sub>LCSS</sub> – A <sub>DTW</sub>	Lower	۲,۳۱۷۰۵	۲,۰۶۲۹۰	۱,۷۳۰۴۴	۱,۲۱۴۷۷	.۰۶۰۵۳۱
Confidence Interval of the Difference						
		.۹۰	.۹۲۵	.۹۵	.۹۷۵	.۹۹

جدول ۷. مقایسه‌ی آماری دقت خوشبندی تحت روش‌های LCSS و DTW.

ردیف	توضیحات	تعداد	درصد
۱	A <sub>LCSS</sub> > A <sub>DTW</sub>	۳۲	۵۲,۹۷
۲	A <sub>LCSS</sub> = A <sub>DTW</sub>	۸	۱۲,۷۰
۳	A <sub>LCSS</sub> < A <sub>DTW</sub>	۲۱	۳۲,۲۳

مجموعه آزمایشی به شبیه‌ترین نماینده خوشبندی اختصاص داده شده و دقت گروه‌بندی محسوبه می‌شود. نتایج این گروه‌بندی تحت روش‌های فوق در جدول ۹ ارائه شده است.

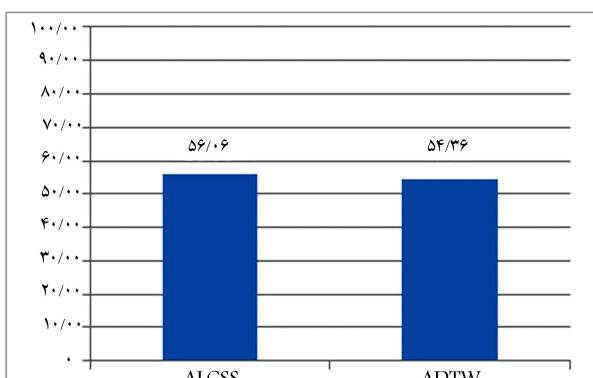
در این جدول و به عنوان مثال برای مجموعه داده‌ی «Statistical Control» در «Statistical Control» آگر تعداد و نماینده خوشبندی بهینه حاصل از خوشبندی مجموعه داده آموزشی این

جدول ۸. خلاصه تعداد حالات پیش‌بینی درست تعداد خوشه تحت روش‌های LCSS و DTW

ردیف	توضیحات	LCSS	DTW	ردیف
۱	تعداد حالاتی که تعداد خوشه به درستی پیش‌بینی شده است	۲۵	۲۸	
۲	تعداد حالاتی که تعداد خوشه، با اختلاف یک عدد پیش‌بینی شده است	۱۹	۱۷	

جدول ۹. دقت گروه‌بندی مجموعه داده‌های آزمایشی مبتنی بر نتایج خوشه‌بندی تحت روش‌های LCSS و DTW

ردیف	نام مجموعه داده	LCSS	DTW	ردیف	نام مجموعه داده	LCSS	DTW	ردیف	نام مجموعه داده	LCSS	DTW
۱	statistical control	۲۲	۸۵,۵۷	۹۶	Diatom Size	۴۳	۹۱,۸۲	۸۴,۶۴	Beetle Fly	۴۰	۴۵
۲	Gun-Point	۲۳	۴۶,۶۷	۴۸	Symbol	۴۴	۹۴,۱۷	۹۱,۷۶	Bird Chicken	۴۰	۴۵
۳	CBF	۲۴	۹۱	۹۳,۳۳	Worms	۴۵	۳۴,۲۵	۲۳,۷۶	Ham	۴۵,۷۲	۴۱,۹۱
۴	ECG	۲۵	۷۱	۵۴	Two Pattern	۴۶	۵۹,۳۲	۷۱,۱۷	Phalanges O-C	۴۴,۹۹	۴۱,۳۵
۵	Face4	۲۶	۸۶,۳۴	۴۴,۳۲	Wafer	۴۷	۶۸,۵۸	۵۹,۳۸	Proximal POA	۷۶,۱۰	۸۵,۳۷
۶	Medical	۲۷	۲۶,۴۵	۳۲,۷۶	Faceall	۴۸	۶۳,۶۷	۵۸,۸۱	Proximal POC	۳۶,۴۳	۳۶,۴۳
۷	Sweedian	۲۸	۶۳,۸۴	۵۵,۰۴	Lighting2	۴۹	۳۶,۰۷	۲۲,۷۹	Proximal PT	۵۲,۲۵	۷۲,۲۵
۸	OSU	۲۹	۳۹,۶۷	۳۵,۵۴	ECGFiveday	۵۰	۵۳,۶۶	۶۰,۸۶	Toe Segmentation1	۴۱,۶۷	۵۶,۱۴
۹	Adiac	۳۰	۳۵,۲۹	۳۷,۸۵	Haptics	۵۱	۲۹,۲۲	۱۹,۱۶	Toe Segmentation2	۸۲,۳۱	۷۴,۶۲
۱۰	Beef	۳۱	۵۰	۴۶,۶۷	InLineSkate	۵۲	۱۷,۲۶	۱۸	Distal POA	۶۵,۷۵	۶۹
۱۱	Lighting	۳۲	۵۰,۴۹	۵۱,۴۳	Motestrain	۵۳	۸۸,۵۸	۹۰,۰۷	Distal POC	۷۳,۱۷	۳۷,۱۷
۱۲	Fish	۳۳	۷۲,۵۷	۵۲	SonyII	۵۴	۷۴,۱۹	۶۶,۳۲	Distal PT	۷۳,۲۵	۷۷,۲۵
۱۳	50words	۳۴	۴۵,۴۹	۴۰,۶۶	sonySurface	۵۵	۴۳,۲۶	۴۴,۰۹	Earth quakes	۴۷,۰۲	۴۱,۳۱
۱۴	Trace	۳۵	۴۸	۷۲	StarLightCurve	۵۶	۷۵,۴۱	۶۷,۰۹	Middle POA	۶۸,۷۵	۲۵,۰
۱۵	Lighting7	۳۶	۴۹,۳۲	۵۲,۰۱	Two Lead ECG	۵۷	۶۶,۳۷	۶۰,۰۸	Middle POC	۴۴,۱۷	۴۸,۵
۱۶	Distal	۳۷	۷۸,۲۵	۷۲,۲۵	Criket X	۵۸	۱۹,۷۴	۳۸,۴۶	Shapelet Sim	۲۹,۴۴	۴۸,۸۹
۱۷	Italy power demand	۳۸	۶۴,۳۴	۷۳,۸۶	Criket Y	۵۹	۲۸,۲۰	۳۶,۶۷	Wine	۴۸,۱۵	۵۰
۱۸	Middle-P-T	۳۹	۶۱,۱۶	۶۱,۱۶	U Wave X	۶۰	۴۷,۶۳	۴۵,۱۰	Computers	۲۸,۴۰	۶۵,۲۰
۱۹	Plane	۴۰	۹۹,۰۵	۹۹,۰۵	U Wave Y	۶۱	۳۶,۶۸	۴۵,۰۹	Meat	۷۱,۶۷	۴۶,۶۷
۲۰	Car	۴۱	۵۱,۶۷	۴۳,۳۲	Insect Wing	۶۲	۴۲,۶۸	۱۸,۳۳	Refrigeration	۴۱,۶	۴۶,۱۳
۲۱	Olive Oil	۴۲	۵۶,۶۷	۸۶,۶۷	Arrow Head	۶۳	۵۲,۰۷	۴۱	Worm Two Class	۲۴,۲۵	۵۶,۹۱



نمودار ۳. متوسط کل دقت گروه‌بندی مجموعه داده آزمایشی تحت روش‌های DTW و LCSS

داده شده است. براساس این نتایج، می‌توان فرض صفر را پذیرفت و ادعا کرد که در سطوح مختلف اطمینان، دقت گروه‌بندی مجموعه داده آزمایشی تحت روش LCSS با دقت گروه‌بندی مجموعه داده آزمایشی تحت روش DTW تفاوتی ندارد. به عبارت دیگر، کیفیت نماینده خوشه‌های حاصله از خوشه‌بندی مجموعه داده‌های آموزشی تحت روش‌های یاد شده از کیفیت یکسانی برخوردار است و تفاوت معناداری با یکدیگر ندارند.

مجموعه داده تحت روش DTW را برای گروه‌بندی مجموعه داده‌ای آزمایشی خودش تحت روش DTW به کار ببریم، دقت گروه‌بندی برابر  $96\%$  خواهد بود. اگر با استفاده از تعداد و نماینده خوشه بهینه حاصل از خوشه‌بندی مجموعه داده آموزشی همین مجموعه داده تحت روش LCSS، عمل گروه‌بندی اعضايی مجموعه داده‌ای آزمایشی همین مجموعه داده را تحت روش LCSS انجام دهیم، آنگاه دقت گروه‌بندی برابر  $85,57\%$  خواهد بود. این عملیات برای سایر مجموعه داده‌ها نیز انجام شده است.

بر اساس این نتایج، متوسط دقت گروه‌بندی کل برای همه مجموعه‌های داده تحت روش DTW و LCSS به ترتیب برابر  $54,35\%$  و  $56,06\%$  است که در نمودار ۳ نمایش داده شده است.

با این توضیحات و به منظور پاسخ به سؤال سوم این تحقیق، مجدداً از آزمون مقایسات زوجی  $t$  استفاده می‌کنیم.

اگر  $ALCSS$  یا  $\bar{A}_{LCSS}$  دقت گروه‌بندی مجموعه داده آزمایشی تحت روش LCSS و  $ADTW$  یا  $\bar{A}_{DTW}$  دقت گروه‌بندی مجموعه داده آزمایشی تحت روش DTW باشد، آنگاه فرضیات این آزمون به شرح زیر خواهد بود.

$$ALCSS = ADTW : H,$$

$$ALCSS \neq ADTW : H,$$

نتایج این آزمون در جدول ۱۰ و برای سطوح مختلفی از قابلیت اطمینان نمایش

جدول ۱۰. نتایج آزمون مقایسات جفتی t برای دقت گروه‌بندی مجموعه داده آزمایشی تحت روش‌های LCSS و DTW

	Paired Differences		T	Df	sig. (۲ - tailed)	Correlation
	Mean	Std. Deviation				
A <sub>LCSS</sub> - A <sub>DTW</sub>	۱,۷۱۴۲۹	۱۴,۶۰۸۹۵	.۹۲۸	۶۲	.۳۵۷	$A_{LCSS} & A_{DTW} =$ .۷۳۲
Confidence Interval of the Difference						
	%۸۰	%۹۰	%۹۵	%۹۷,۵	%۹۹	
A <sub>LCSS</sub> - A <sub>DTW</sub>	Lower	-۰,۶۷۸۰۵	-۱,۳۶۹۶۰	-۱,۹۷۷۵۲	-۲,۵۲۸۲۰	-۳,۱۹۳۶۹
	Upper	۴,۱۰۶۶۲	۴,۷۹۸۱۷	۵,۴۰۶۰۹	۵,۹۵۶۷۷	۶,۶۲۲۲۶

## ۶. نتیجه‌گیری

در یک جمع‌بندی کلی، می‌توان ادعا کرد که دقت تکنیک نزدیک‌ترین همسایگی تحت روش LCSS با اطمینان ۹۲,۵٪ به نتایج بهتری نسبت حالتی خواهد داشت که این تکنیک تحت روش DTW اجرا شود. می‌توان ادعا کرد که با اطمینان ۹۹٪ دقت خوش‌بندی کامدوبید تحت روش LCSS نسبت به روش DTW به نتایج بهتری دست خواهیم یافت. تأثیر روش LCSS نسبت به تأثیر روش DTW در تعیین تعداد خوش‌های مجموعه داده به صورت آماری یکسان است. تأثیر روش LCSS با تأثیر روش DTW در تعیین نماینده خوش نیز تقریباً تفاوتی ندارد به طوری که گاه یکی از آنها برای برخی مجموعه داده‌های سری زمانی خوب و گاه دیگری برای سایر مجموعه داده‌های سری زمانی خوب است و این تأثیر حدوداً بینایین است. لازم به ذکر است که به منظور بررسی جامع‌تر تأثیر روش‌های یاد شده در داده کاری سری زمانی، می‌توان آنها را در تکنیک‌های دیگر داده‌کاری نیز مورد ارزیابی قرار داد یا از تکنیک‌های آماری دیگری برای ارزیابی و مقایسه‌ی آنها استفاده کرد.

جدول ۱۱. مقایسه‌ی آماری دقت خوش‌بندی تحت روش‌های LCSS و DTW

ردیف	توضیحات	درصد	تعداد
۱	$A_{LCSS} > A_{DTW}$	۵۵,۵۶	۳۵
۲	$A_{LCSS} = A_{DTW}$	۶,۳۵	۴
۳	$A_{LCSS} < A_{DTW}$	۳۸,۱۰	۲۴

این مقایسه را می‌توان از منظر آمار توصیفی نیز انجام داد. نتایج این مقایسه در جدول ۱۱ نمایش داده شده است. بر اساس اطلاعات مندرج در این جدول، دقت گروه‌بندی تحت روش LCSS نسبت به دقت این تکنیک تحت روش DTW، در ۵۵,۵۶٪ از حالات، بهتر، در ۶,۳۵٪ از حالات برابر و فقط در ۳۸,۱۰٪ از حالات، بدتر است. به طور خلاصه، این آمارها به همراه انحراف معیار اختلاف دقت گروه‌بندی تحت دو روش فوق نشان می‌دهد که تأثیر روش LCSS و تأثیر روش DTW در تعیین نماینده خوش تفاوتی ندارد.

## پانوشت‌ها

- longest common sub sequence (LCSS)
- dynamic time warping (DTW)
- whole matching
- subsequence matching
- shape based distance measure
- edit based distance measure
- feature based distance measure
- model based distance measure
- partitioning algorithm of medoid (PAM)
- elbow index

neering Applications of Artificial Intelligence, 24(1) pp. 164-181 (2011).

- Keogh, E. and Kasetty, S. "On the need for time series data mining benchmarks: a survey and empirical demonstration", *Data Mining and Knowledge Discovery*, 7(4), pp. 349-371 (2003).
- Sangeeta, R. and Geeta, S. "Recent techniques of clustering of time series data: a survey", *International Journal of Computer Applications*, 52(15), pp. 1-9 (2012).
- Lin, J., Vlachos, M., Keogh, E. and et al. "Iterative incremental clustering of time series", *International Conference on Extending Database Technology, Advances in Database Technology*, pp. 106-122 (2004).
- Gordon, A.D. "Clustering algorithms and cluster validity", *Computational Statistics*, PhysicaVerlag, pp. 497-512 (1994).
- Milligan, G.W. and Cooper, M.C. "An examination of procedures for determining the number of clusters in a data set", *Psychometrika*, 50, pp. 159-179 (1985).

## (References) مراجع

- Morris, B. and Trivedi, M. "Learning trajectory patterns by clustering: experimental studies and comparative evaluation", *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 312-319 (2009).
- Fu, T.C. "A review on time series data mining", *Engi-*

8. Aghabozorgi, S., Shirkhorshidi, A.S. and Wah, T.Y. "Time series clustering-A decade review", *Information Systems*, **53**, pp. 16-38 (2015).
9. Faloutsos, C. Ranganathan, M. and Manolopoulos, Y. "Fast subsequence matching in time-series databases", *In Proceedings of the ACM SIGMOD International Conference on Management of Data*, **23**, pp. 419-429 (1994).
10. Sakoe, H. and Chiba, S. "A dynamic programming approach to continuous speech recognition", *In Proceedings of the 7th International Congress on Acoustics*, **3**, pp. 65-69 (1971).
11. Sakoe, H. and Chiba, S. "Dynamic programming algorithm optimization for spoken word recognition", *IEEE Transactions on Acoustics, Speech, and Signal Processing*, **26**(1), pp. 43-49 (1987).
12. Chen, L. and Ng, R. "On the marriage of Lp-norms and edit distance", *In Proceedings of the 30th International Conference on Very Large Data Bases*, **30**, pp. 792-803 (2004).
13. Aßfalg, J., Kriegel, H.P., Kroger, P. and et al. "Similarity search on time series based on threshold queries", *In Advances in Database Technology*, pp. 276-294 (2006).
14. Vlachos, M., Kollios, G. and Gunopulos, D. "Discovering similar multidimensional trajectories", *In Proceedings of the 18th International Conference on Data Engineering*, pp. 673-684 (2002).
15. Banerjee, A. and Ghosh, J. "Clickstream clustering using weighted longest common subsequences", *In Proceedings of the Web Mining Workshop at the 1st SIAM Conference on Data Mining*, pp. 33-40 (2001).
16. Chen, L., Ozsu, M.T. and Oria, V. "Robust and fast similarity search for moving object trajectories", *In Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 491-502 (2005).
17. Keogh, E., Lonardi, S., Ratanamahatana, C.A. and et al. "Compression-based data mining of sequential data", *Data Mining and Knowledge Discovery*, **14**(1), pp. 99-129 (2007).
18. Han, J., Kamber, M. and Pei, J., *Data Mining: Concepts and Techniques*, 3rd ed (2011).
19. Rousseeuw, P.J. "Silhouettes: a graphical Aid to the interpretation and validation of cluster analysis", *Journal of Computational and Applied Mathematics*, **20**, pp. 53-65 (1987).
20. Hubert, L. and Arabie, P. "Comparing partitions", *Journal of Classification*, **2**, pp. 193-218 (1985).
21. Notredame, C., Higgins, D.G. and Heringa, J. "TCoffee: a novel method for fast and accurate multiple sequence alignment", *Journal of Molecular Biology*, **302**(1), pp. 205-217 (2000).
22. Murray, P.W., Agard, B. and Barajas, A. "Market segmentation through data mining: a method to extract behaviors from a noisy data set", *Computers&Industrial Engineering*, **109**, pp. 233-252 (2017).
23. Ahn, G. and Hur, S. "Efficient genetic algorithm for feature selection for early time series classification", *Computers & Industrial Engineering*, **142**, 106345 (2020).
24. Luzianin, I. and Krause, B. "Similarity measurement of biological signals using dynamic time warping algorithm", *Proc. of the 4th International Conference on Applied Innovations in IT*, (2016).
25. Barragan, J.F., Fontes, C.H. and Embicucu, M. "A wavelet-based clustering of multivariate time series using a multiscale SPCA approach", *Computers and Industrial Engineering* (2016). <https://doi.org/10.1016/j.cie.2016.03.003>.
26. Hong, Jae.Y., Park, Seung. H. and Baek, Jun-G. "SS-DTW:Shape segment dynamic time warping", *Expert Systems with Applications*, **150**, 113291 (2020).
27. Just, W. "Computational complexity of multiple sequence alignment with SP-score", *Journal of Computational Biology*, **8**(6), pp. 615-623 (2001).
28. Hubert, L. and Schultz, J. "Quadratic assignment as a general data-analysis strategy", *British Journal of Mathematical and Statistical Psychology*, **29**, pp. 190-241 (1987).
29. Chou, C.H., Su, C. and Lai, E. "A new cluster validity measure for clusters with different densities", *IASTED International Conference on Intelligence Systems and Control*, pp. 276-281 (2003).
30. Cho, Su.G. and Kim, Seoung. B. "A Data-driven document similarity measure based on classification algorithms", *International Journal of Industrial Engineering: Theory, Applications and Practice*, **24**(3), pp.328-339(2017).
31. Sridevi, S., Parthasarathy, S. and Rajaram, S. "An Effective prediction system for time series data using pattern matching algorithms", *International Journal of Industrial Engineering: Theory, Applications and Practice*, **25**(2), pp.123-136 (2018).
32. Agnetisa, A., Alfierib, A. and Nicosia, G. "An heuristics approach to batching and scheduling a singal machine to minimize setup costs", *Copmuters & Industrial Engineering*, **46**, pp. 793-802 (2004).
33. Penga, B., Zhangb, Y., Lub, Z. and et al. "A learning based memetic algorithm for the multiple vehicle pickup and delivery problem with LIFO loading", *Copmuters & Industrial Engineering*, **142**, 106241 (2020).
34. Zhang, Y. and Xiaoping, L. "Estimation of distribution algorithm for permutation flow shops with total flowtime minimization", *Copmuters & Industrial Engineering*, **60**, pp. 706-718 (2011).
35. Dzalbs, I. and Kalganova, T. "Accelerating supply chains with ant colony optimization across a range of hardware solutions", *Copmuters & Industrial Engineering*, **147**, 106610, (2020).