

# طراحی نمودار کنترل $T^2$ هتلینگ با استفاده از خوشه بندی

علیرضا پویا\* (استاد)

گروه مدیریت، دانشکده علوم اداری و اقتصادی، دانشگاه فردوسی مشهد

علی یگانه (دانشجوی دکتری)

گروه مهندسی صنایع، دانشکده مهندسی، دانشگاه فردوسی مشهد

سمیه فدائی (دانشجوی دکتری)

گروه مدیریت، دانشکده علوم اداری و اقتصادی، دانشگاه فردوسی مشهد

مهندسی صنایع و مدیریت شریف، تابستان ۱۴۰۰  
دوره ۱-۳۷، شماره ۱، ص. ۸۳-۷۸ (پژوهشی)

هدف اصلی تحقیقات در زمینه کنترل فرایند آماری چندمتغیره، در نظر گرفتن همبستگی بین چندین مشخصه کیفی برای یک مرحله از فرایند است. در فاز دوم رویه کنترل فرایند چندمتغیره با استفاده از حدود کنترلی به دست آمده از فاز اول و مشاهدات آتی، تحت کنترل بودن فرایند بررسی می شود. یافتن نقاط پرت فاز اول قبل از محاسبه حدود کنترلی برای حصول نتیجه مناسب اهمیت بسیار دارد. لذا در این تحقیق نمودار کنترل پیشنهادی با استفاده از روش خوشه بندی سلسله مراتبی، مشاهدات پرت را شناسایی می کند. اهمیت این روش در تعیین داده های پرت با استفاده از ضریب ناهمسانی و مجموعه بی از حدود کنترل متغیر است. در این روش ضریب ناهمسانی و مجموعه بی از حدود کنترل با استفاده از پارامترهای اندازه نمونه و تعداد متغیرها (شاخص های کیفی) تعیین می شود، در واقع فاصله بین مشاهدات به صورت خوشه بندی مدل می شود و داده های پرت با الگوریتم بازگشتی حذف می شوند و سپس میانگین و ماتریس واریانس و کوواریانس  $T^2$  بر اساس داده های باقیمانده تعیین می شود. در مرحله آخر با توجه به حد کنترل به دست آمده، آماری  $T^2$  تعیین می شود. برای ارزیابی عملکرد نمودار کنترل پیشنهادی و مقایسه آن با نمودار  $T^2$  هتلینگ معمولی در شناسایی نقاط پرت از شاخص عدم مرکزیت و روش الفارو و همکاران بر اساس تشخیص مجموعه داده های پرت استفاده شده است. همچنین برای بررسی و مقایسه بیشتر، دو نمودار از مجموعه داده های هاوکینز و فسفر نیز بررسی شده است. اهمیت و کارایی روش پیشنهادی در شرایط وجود تعداد داده های پرت زیاد در مجموعه داده ها، به طور محسوس مشاهده شد.

واژگان کلیدی: نمودار کنترل چندمتغیره، نمودار کنترل  $T^2$  هتلینگ، تحلیل خوشه بی سلسله مراتبی.

alirezapooya@um.ac.ir  
ali.yeganeh@mail.um.ac.ir  
somayehfadaei@mail.um.ac.ir

## ۱. مقدمه

امروزه کیفیت به عنوان راهبردی تجاری شناخته می شود و عنوان غول صنعتی نه به کشورهای با تولید انبوه، بلکه به کشورهایی با بالاترین درجه کیفیت در محصولات، صنایع و خدمات رسانی اطلاق می شود. حفظ دستاوردهای مرحله بهینه سازی قبل از ساخت در مرحله حین ساخت توسط فنون آماری به نام «نمودارهای کنترل» در مبحث کنترل آماری فرایند انجام می گیرد. کنترل فرایند آماری (SPC) به طور گسترده به عنوان یک روش قدرتمند برای اندازه گیری، طبقه بندی، تجزیه و تحلیل و تفسیر داده های فرایند برای بهبود کیفیت محصولات و خدمات با شناسایی

\* نویسنده مسئول

تاریخ: دریافت ۱۳۹۹/۳/۲۶، اصلاحیه ۱۴۰۰/۱/۱۷، پذیرش ۱۴۰۰/۲/۲۷

DOI:10.24200/J65.2021.55182.2089

از نمودارهای کنترل تک‌متغیره در فرایندهای چندمتغیره ممکن است منجر به نتایج نامطلوب شود. [۶۵]

هدف اصلی نمودار کنترل چندمتغیره، شناسایی وجود علل تغییر در ویژگی‌های کیفیت چندمتغیره است. به طور خاص، برای تجزیه و تحلیل فاز یک، از یک مجموعه داده‌ی تاریخی، دو هدف دنبال می‌شود: الف) برای شناسایی تغییرات در بردار میانگین که ممکن است برآورد میانگین بردار کنترل و ماتریس کوواریانس را تحریف کند؛ ب) شناسایی و از بین بردن نقاط پرت چندمتغیره. [۶]

در داده‌های یک‌متغیره، تشخیص داده‌ی پرت ساده است. در حالت دومتغیره، هنگام نمایش پراکندگی داده‌ها، یک داده ممکن است بدون وجود مشکل غیرضروری، پرت به نظر برسد. در داده‌های چندمتغیره، هنگامی که مجموعه داده بیش از دو بعد داشته باشد، شناسایی داده‌ی پرت توسط متدهای بصری امکان‌پذیر نیست. هتلینگ روش شناخته شده‌ی چندمتغیره‌ی است که به طور گسترده در نمودار کنترل  $T^2$  در مجموعه داده‌های چندبعدي به کار رفته است. [۸]

در نمودار کنترل  $T^2$  فرض بر این است که داده‌ها دارای توزیع نرمال چندمتغیره با میانگین  $\mu$  و واریانس  $\sigma$  است. اگر فرایند شامل  $p$  شاخص کیفی با تعداد  $m$  مشاهده‌ی  $n$  تایی باشد، آماره چنین تعریف می‌شود:

$$T^2 = n(x - \bar{x})' S^{-1} (x - \bar{x}) \quad (1)$$

که در آن،  $x = [x_1, x_2, \dots, x_p]^T$  و  $\bar{x}$  و  $S$  برآوردی از  $\mu$  و  $\sum$  است.  $S$  و  $\bar{x}$  مطابق رابطه‌ی ۲ و ۳ برآوردی از میانگین کل نمونه و ماتریس کوواریانس است.

$$\bar{X} = \frac{1}{n} \sum_{i=1}^m x_i \quad (2)$$

$$\sum = \frac{1}{n-1} \sum_{i=1}^m (X_i - \bar{X})(X_i - \bar{X})' \quad (3)$$

در حالتی که  $n = 1$  است، میسون و یانگ [۹] حدود کنترلی فاز ۱ را بر مبنای توزیع  $\beta$  به صورت رابطه‌ی ۴ برای نمودار کنترل  $T^2$  در نظر گرفته‌اند:

$$\begin{cases} UCL = \frac{p(m-1)^2}{m} \beta_{\alpha, p/2, m-p-1/2} \\ LCL = 0 \end{cases} \quad (4)$$

روش کلاسیک  $T^2$  برای تشخیص تغییر شیفیت میانگین و نقاط پرت که در رابطه‌ی ۱ توضیح داده شده است در صورتی که مجموعه داده فقط دارای یک داده پرت باشد، مناسب است. [۱۰] با این وجود، هنگامی که بیش از یک مشاهده پرت در مجموعه داده‌ها وجود داشته باشد، قدرت این تکنیک‌ها به طور قابل توجهی کاهش می‌یابد. [۱۱، ۱۲] به عبارتی هنگامی که بیش از یک داده‌ی پرت در فرایند وجود داشته باشد، به دلیل اثر درون‌آوری<sup>۱</sup> و برون‌بری<sup>۲</sup>، تشخیص نقاط پرت مشکل می‌شود. [۱۳] درون‌آوری هنگامی اتفاق می‌افتد که تشخیص مشاهدات پرت واقعی به‌عنوان پرت اعلام نمی‌شود. در حالی که برون‌بری هنگامی اتفاق می‌افتد که مشاهدات به‌صورت نادرست به‌عنوان پرت اعلام شود. [۱۴] داده‌های پرت، در صورت عدم کشف، قابلیت انحراف میانگین و ماتریس کوواریانس نمودارهای کنترل چندمتغیره‌ی  $T^2$  هتلینگ برای نظارت بر ویژگی‌های کیفیت فردی را دارند. تأثیر این تحریف این است که نمودار کنترل ساخته شده از آن غیرقابل اعتماد می‌شود، زیرا دارای اثر درون‌آوری و برون‌بری است. [۱۵] بنابراین، شناسایی نقاط پرت قبل از تجزیه و تحلیل داده‌ها بسیار مهم است به‌خصوص در نظارت بر کیفیت محصول فرایند تولید. [۸]

پدیده‌ی درون‌آوری و برون‌بری اتفاق می‌افتد چون میانگین و ماتریس واریانس و کوواریانس باثبات نیستند. [۱۶] یکی از راه‌های حل این مشکلات، استفاده از برآوردگرهای باثبات برای معادلات ۲ و ۳ به‌عنوان گزینه‌ی برای موقعیت و ماتریس کوواریانس است. این برآوردگرها شامل روش پیراسته<sup>۳</sup> با بهبود  $T^2$  هتلینگ، [۱۵، ۱۶] روش اختلاف پی‌درپی سالیوان و وودال [۱۷] است که  $SW_1$  نامیده می‌شوند و اصلاح طرح استلاکتیت<sup>۴</sup> اتکینسون و مولیرا [۱۸] به‌عنوان رویکرد دوم سالیوان و وودال، که آن را  $SW_2$  می‌گویند. روش‌های ذکر شده در واقع می‌توانند  $T^2$  هتلینگ را بهبود بخشند، اما هنوز در تشخیص تعداد زیادی از نقاط پرت موجود در داده‌ها مؤثر نیستند. [۸]

رسو<sup>۵</sup> [۱۹] دو روش حداقل حجم بیضی‌وار (MVE) و حداقل روش تعیین‌کننده‌ی کوواریانس (MCD) [۶] را برای مسئله‌ی تشخیص تعداد زیادی از نقاط پرت پیشنهاد کردند. با این حال، هر دو روش نیاز به محاسبات بسیار سنگین دارد. در دو دهه‌ی گذشته بحث‌های بسیاری در مورد تقریب MVE و MCD انجام شده است. اما در دهه‌ی اخیر، در برخی از مطالعات استفاده از روش‌های تحلیل خوشه‌ی برای بهبود نمودار کنترل چندمتغیره آغاز شده است. برای مثال اونگ و الوی [۱۰] از تحلیل خوشه‌ی رگرسیون، کانگ و بام کیم [۲۰] از تحلیل خوشه‌ی k-means، فان و همکارانش [۸] از تحلیل خوشه‌ی سلسله‌مراتبی برای نمودار کنترل چندمتغیره  $SW_2$  استفاده کرده‌اند. رویکردهای خوشه‌بندی معمولاً از لحاظ محاسباتی کارآمدتر از روش‌های نمودار کنترل است. [۸] لذا با توجه به این که تحلیل خوشه‌ی نسبت به داده‌های پرت حساس است، [۸] در این پژوهش از تکنیک خوشه‌بندی سلسله‌مراتبی در طراحی نمودار کنترل در فاز اول، برای به دست آوردن مجموعه‌ی باثبات از مشاهدات استفاده می‌شود و تمرکز بر استفاده از رویکردهای خوشه‌ی به‌عنوان فیلتر اولیه‌ی داده‌های پرت است.

یکی از عوامل رایج خارج از کنترل شدن نمودار، رخ دادن چندین داده‌ی پرت است. برآورد پارامترهای  $\mu$  و  $\sum$  مخصوصاً در نمودار کنترل  $T^2$  هتلینگ، نسبت به یک مشاهده‌ی پرت حساس‌اند اما آماره‌های به دست آمده از نمودار کنترل مذکور، بر اساس این برآوردگرها وقتی که چندین مشاهده پرت وجود داشته باشد، بسیار ضعیف عمل می‌کند. [۲۱] همانطور که گفته شد در پژوهش‌های زیادی برای بهبود نمودار کنترل  $T^2$  هتلینگ، اقدام شده است، اما هنوز در تشخیص تعداد زیادی از نقاط پرت موجود در داده‌ها مؤثر نیستند. [۸] لذا دلیل استفاده از تحلیل خوشه‌ی سلسله‌مراتبی در تحقیق حاضر این است که خوشه‌بندی سلسله‌مراتبی امکان استفاده از فاصله‌های مختلف در بین داده‌ها را دارد. همچنین چون تمامی داده‌ها با یکدیگر مقایسه می‌شوند، امکان شناسایی نقاط پرت با الگوهای مختلف فراهم می‌شود. [۱۰]

برای ارزیابی عملکرد نمودار کنترل پیشنهادی ابتدا بر اساس تشخیص یک داده‌ی پرت و با استفاده از شاخص عدم مرکزیت وارگاس [۲۲] استفاده می‌شود. شاخص عدم مرکزیت (ncp) نشان‌دهنده‌ی شدت تغییر بردار میانگین خارج از کنترل ( $\mu_1$ ) از بردار میانگین تحت کنترل میانگین ( $\mu_0$ ) مطابق رابطه‌ی ۱ است. وارگاس معیار دقت روش را تشخیص دست‌کم یک سیگنال در داده‌ها دانسته است.

$$ncp = (\mu_1 - \mu_0)' \sum^{-1} (\mu_1 - \mu_0) \quad (5)$$

در این تحقیق علاوه بر روش ارزیابی وارگاس، از ارزیابی عملکرد بر اساس تشخیص مجموعه‌ی داده‌های پرت الفارو و ارتگا [۱۱] نیز استفاده می‌شود. در این

وارگاس<sup>[۲۲]</sup> در یک مطالعه‌ی مقایسه‌ی  $T^2$  با استفاده از تکنیک‌های شبیه‌سازی عملکرد نمودار کنترلی را بر اساس تشخیص یک داده‌پرت و با استفاده از شاخص عدم مرکزیت (ncp) بررسی کرد. نتایج حاصله و ارزیابی عملکرد نشان داد که نمودار کنترل  $T^2$  با استفاده از برآوردگرهای حداقل حجم بیضی (MVE) در شناسایی تعداد نقاط پرت مؤثر است و روش SW۱ در شناسایی تعداد زیادی از نقاط انحرافی مؤثر نیست. روش SW۲ نیز در مقابل داده‌هایی با نقاط دورافتاده زیاد، آسیب‌پذیر است و همچنین به نمونه‌ی تصادفی اولیه وابسته است.

الفارو و ارتگا<sup>[۱۶]</sup> یک برآوردگر باثبات به جای  $T^2$  با استفاده از روش پیراسته پیشنهاد کردند. این نمودار کنترل جدید از ساختار مشابه  $T^2$  برخوردار است، اما بردار موقعیت و برآورد ماتریس کوواریانس با استفاده از پیرایش به دست می‌آید. آن‌ها رفتار و عملکرد این دو نمودار را با استفاده از یک مثال و یک مطالعه شبیه‌سازی بر اساس تشخیص مجموعه‌ی داده‌های پرت مقایسه کردند. یک سال بعد این دو پژوهش‌گر<sup>[۱۸]</sup> یک مطالعه‌ی شبیه‌سازی برای تحلیل عملکرد برآوردگرهای MVE، MCD، MCD، مازموزون و برآوردگر پیراسته را در موقعیت‌های مختلف انجام دادند. احسان و همکاران<sup>[۱۴]</sup> نمودار  $T^2$  مبتنی بر ترکیب PCA با حدود کنترل چگالی کرنل برای داده‌های پیوسته و دسته‌ی مختلط ارائه کردند.

پیشینه‌ی مطالعات نشان می‌دهد که در دهه‌ی اخیر استفاده از روش‌های تحلیل خوشه‌ی برای بهبود نمودار کنترل چندمتغیره برای شناسایی نقاط پرت و غیرمعمول چندمتغیره، افزایش یافته است. انکارا و ویل<sup>[۲۸]</sup> نمودار کنترل  $R$  اصلاح شده را با استفاده از روش خوشه‌ی تک پیوند ارائه دادند. جوب و پاکجوبی<sup>[۲۹]</sup> یک روش مبتنی بر خوشه‌ی چندمرحله‌ی رایانه‌ی را پیشنهاد کردند. آنان چند سال بعد نیز برای غلبه بر این اثر درون‌آوری، یک رویکرد مبتنی بر خوشه رایانه‌ی پیشنهاد کردند<sup>[۳۰]</sup> که آماره‌ی  $mcd$  مازموزون را با یک الگوریتم مبتنی بر خوشه‌ی چندمرحله‌ی ترکیب می‌کند. فان و همکارانش<sup>[۸]</sup> نمودار کنترل چندمتغیره SW۲ برای تشخیص پرت‌ها با استفاده از درخت خوشه‌ی سلسله‌مراتبی ارائه دادند که به طور مؤثر بتواند داده‌های پرت شناسایی کند. کانگ و بام کیم<sup>[۳۰]</sup> نمودار کنترل مبتنی بر الگوریتم خوشه‌ی  $k$ -means برای فرایندهای TFT-LCD با توزیع غیرهمگن ارائه کردند. اونگ و الی<sup>[۱۰]</sup> نمودار کنترل را بر اساس تنظیم خوشه‌ی رگرسیون برای نظارت بر خصوصیات کیفیت فردی در یک محیط چندمتغیره ارائه کرده‌اند. هوانگ و همکاران<sup>[۳۱]</sup> یک الگوریتم تشخیص خوشه‌ی پرت به نام ROCF را بر اساس مفهوم گراف همسایه متقابل و با این عقیده که معمولاً اندازه‌ی خوشه‌های دورافتاده بسیار کوچک‌تر از خوشه‌های معمولی است، پیشنهاد کردند. بررسی پیشینه نشان می‌دهد که تحقیقات متعددی با روش‌های متفاوت اقدام به بهبود و توسعه‌ی نمودار کنترل  $T^2$  هتلینگ پرداخته‌اند و در دهه‌ی اخیر ترکیب روش‌های تحلیل خوشه‌ی با نمودارهای چندمتغیره برای بهبود نمودار کنترل در شناسایی نقاط پرت، افزایش یافته است. برای مثال فان و همکارانش<sup>[۸]</sup> از خوشه‌بندی برای بهبود برآوردگر سالیوان ۲ (SW۲) استفاده کرده‌اند. جوب و پاکجوبی<sup>[۳۰]</sup> از الگوریتم خوشه‌بندی چندمرحله‌ی برای بهبود MCD رسو استفاده کردند. همچنین اونگ و الی<sup>[۱۰]</sup> از رگرسیون خوشه‌ی برای بهبود نمودار کنترل  $T^2$  استفاده کرده‌اند. مشاهده می‌شود در این پژوهش‌ها از خوشه‌بندی برای بهبود نوع متفاوتی از نمودار کنترل  $T^2$  اقدام شده است؛ در حالی که در پژوهش حاضر از تحلیل خوشه‌ی سلسله‌مراتبی برای بهبود برآوردگر  $T^2$  استفاده شده است. در این روش به‌صورت بازگشتی نقاط شناسایی شده پرت توسط الگوریتم ترکیبی پیشنهاد شده، شناسایی

روش اگر تمامی داده‌های پرت شناسایی شوند، دقت روش مورد تأیید است. سپس با استفاده از دو مجموعه داده‌ی واقعی فسفر و هاوکینز<sup>[۸]</sup> نمودار کنترل  $T^2$  معمولی با نمودار  $T^2$  خوشه‌بندی مقایسه می‌شود.

در ادامه‌ی این نوشتار، در بخش دوم پیشینه‌ی مطالعاتی نمودار کنترل  $T^2$  و تلاش‌هایی که برای بهبود این نمودار انجام شده بررسی می‌شود. در بخش سوم روش پیشنهادی تحقیق حاضر معرفی خواهد شد. در بخش چهارم ارزیابی عملکرد نمودار کنترل پیشنهادی ( $T^2$  خوشه‌بندی) با استفاده از روش وارگاس و روش الفارو و ارتگا بررسی و مقایسه می‌شود؛ سپس با استفاده از مجموعه داده‌های واقعی فسفر و هاوکینز، نسبت به ارزیابی نمودار کنترل  $T^2$  خوشه‌بندی و مقایسه‌ی آن با نمودار  $T^2$  هتلینگ معمول اقدام می‌شود. در بخش پایانی به بحث و نتیجه‌گیری مطالب به دست آمده خواهیم پرداخت.

## ۲. پیشینه‌ی تحقیق

تاکنون روش‌های برآورد باثبات متفاوتی برای شناسایی نقاط پرت و غیرمعمول چندمتغیره در فاز اول نمودارهای کنترلی مطرح شده است که مهم‌ترین و کاراترین آنها در ادامه معرفی می‌شود. رسو<sup>[۱۹]</sup> دو روش حداقل حجم بیضی‌وار (MVE) و حداقل دترمینان کوواریانس (MCD) را برای شناسایی نقاط پرت در داده‌های چندمتغیره پیشنهاد کرد. برآوردگر حداقل حجم بیضی‌وار (MVE) اولین برآوردگر باثبات بر اساس موقعیت و پراکندگی چندمتغیره است. در حالی که ایده‌ی MVE بسیار شهودی است، یافتن برآوردگر MVE در عمل می‌تواند بسیار دشوار باشد. با افزایش اندازه‌ی نمونه و ابعاد داده‌ها محاسبات به طور نمایی افزایش می‌یابد؛ به همین دلیل رسو و لئوری<sup>[۲۳]</sup> یک روش تقریبی برای یافتن برآوردگرهای MVE با استفاده از الگوریتم زیرنمونه‌گیری<sup>۷</sup> را پیشنهاد دادند. دیگر برآوردگر بردار میانگین و ماتریس کوواریانس که می‌تواند تا ۵۰٪ نقاط پرت را در مجموعه داده اصلاح کند کم‌ترین برآوردگر تعیین کوواریانس (MCD) است. روسو و ون‌درایسین<sup>[۲۴]</sup> الگوریتمی سریع برای محاسبه‌ی برآوردگر تعیین کوواریانس پیشنهاد دادند و آن را FastMCD نامیدند؛ این الگوریتم قابلیت کنترل مجموعه داده‌های بزرگ در یک زمان معقول را دارد. جنسن و همکاران<sup>[۲۵]</sup> دو روش MVE و MCD را در داده‌های چندمتغیره‌ی زیادی مقایسه کردند و نتیجه گرفتند هر دو روش در تشخیص نقاط پرت مؤثرند، اما عملکرد این دو رویکرد به اندازه‌ی نمونه و نسبت نقاط پرت موجود بستگی دارد. ویلمز و همکارانش<sup>[۲۶]</sup> الگوریتم ترکیبی حداقل دترمینان کوواریانس مازموزون (RMCD)<sup>۸</sup> را پیشنهاد کردند که از جایگزینی برآوردگرهای موقعیت و پراکندگی حداقل دترمینان کوواریانس مازموزون با میانگین و واریانس کلاسیک آماره‌ی  $T^2$  هتلینگ به دست می‌آید. به طور مشابه وریات و وتاتور<sup>[۲۷]</sup> برآوردگرهای باثبات بر اساس حداقل دترمینان کوواریانس مازموزون (RMCD) و حداقل حجم بیضی‌وار مازموزون (RMVE) برای نظارت بر مشاهدات چندمتغیره در داده‌های فاز I پیشنهاد کردند.

سالیوان و وودال<sup>[۱۷]</sup> دو الگوریتم با دو دیدگاه متفاوت ارائه کردند. آن‌ها نشان دادند که نمودار  $T^2$  هتلینگ با استفاده از ماتریس کوواریانس نمونه، در شناسایی تغییرات بردار میانگین مؤثر نیست، از این رو آن‌ها بر اساس نمودار استلاکنیت اتکینسون و مولیرا، پیشنهاد استفاده از اختلاف برداری بین مشاهدات پی‌درپی برای برآورد ماتریس کوواریانس در کنترل روند را دادند.<sup>[۱۸]</sup> با این وجود

می‌شوند. لذا استفاده از این ساختار ترکیبی در شناسایی داده‌های پرت و استفاده از حدود کنترل متغیر متناسب با تعداد داده‌های پرت انتخاب شده، نوآوری تحقیق حاضر است.

### ۳. معرفی روش پیشنهادی

#### ۱.۳. مبانی روش پیشنهادی

اساس روش پیشنهادی در این پژوهش، خوشه‌بندی سلسله‌مراتبی مشابه روش فن و همکاران<sup>[۸]</sup> است. در خوشه‌بندی سلسله‌مراتبی داده‌ها بر اساس فاصله تفکیک می‌شوند و در نهایت خروجی خوشه‌بندی سلسله‌مراتبی، نمودار دندوگرام داده‌ها خواهد بود. بر اساس دندوگرام داده‌ها می‌توان تعداد خوشه‌ها را مشخص کرد. در خوشه‌بندی سلسله‌مراتبی نوع فاصله و نحوه‌ی محاسبه‌ی فاصله‌ی خوشه‌ها باید تعیین شود. در این مدل از فاصله‌ی اقلیدسی و فاصله‌ی بین خوشه‌ی مینیمم (single) استفاده شده است.

برای تعیین وضعیت خوشه‌ها شاخص‌های متفاوتی ارائه شده که یکی از پرکاربردترین شاخص‌ها، ضریب ناهمسانی خوشه‌هاست. این ضریب به‌وسیله‌ی ماتریس ناهمسانی که دارای  $m - 1$  سطر ( $m$  تعداد داده‌هاست) و ۴ ستون است بیان می‌شود.<sup>[۳۲]</sup> ستون اول این ماتریس میانگین فاصله‌ی داده‌های بین دو خوشه در بالاترین سطح است. ستون دوم انحراف معیار فاصله‌ها و ستون سوم تعداد فاصله‌ی موجود بین دو خوشه و ستون چهارم ضریب ناهمسانی است. هر چه این ضریب مقدار کم‌تری باشد یعنی شباهت دو خوشه کم‌تر است و خوشه‌بندی مناسبی صورت پذیرفته است. برای توضیحات بیشتر می‌توان به مرجع مربوطه مراجعه کرد.<sup>[۳۲]</sup>

در این پژوهش پس از انجام خوشه‌بندی سلسله‌مراتبی داده‌ها به دو خوشه تقسیم می‌شوند. سپس ضریب ناهمسانی برای دو خوشه‌ی نهایی (سطر آخر ماتریس ناهمسانی) محاسبه و با حد مجاز آن که بر اساس سیگنال خطا تعیین می‌شود مقایسه می‌شود. اگر از حد مجاز کم‌تر باشد خوشه‌بندی تأیید می‌شود و در غیر این صورت داده‌ی پرت وجود ندارد. در صورت تأیید خوشه‌بندی، خوشه با تعداد داده‌ی کم‌تر، معادل داده‌های پرت در نظر گرفته شده و خوشه‌ی داده‌های اولیه‌ی پرت نامیده می‌شود. در گام دوم برای تمامی داده‌ها آماره‌ی  $T^2$  محاسبه می‌شود و داده‌ی دارای بیشینه مقدار  $T^2$  به‌عنوان داده‌ی انتخاب شده‌ی پرت در نظر گرفته می‌شود. اگر این داده در خوشه‌ی داده‌های اولیه‌ی پرت وجود داشت داده‌ی پرت ثانویه در نظر گرفته می‌شود و از مجموعه‌ی داده‌ها کنار گذاشته می‌شود. الگوریتم فوق دوباره تکرار می‌شود تا داده‌ی پرت ثانویه جدید ایجاد نشود. در انتهای مرحله‌ی فوق، داده‌ها به دو دسته‌ی داده‌های پرت ثانویه و داده‌های اصلی محاسباتی تقسیم می‌شوند. داده‌های پرت ثانویه از مجموعه‌ی داده‌ها کنار گذاشته می‌شود و با توجه به داده‌های اصلی محاسباتی میانگین و ماتریس واریانس کوواریانس مطابق رابطه‌ی ۲ و ۳ ایجاد می‌شود.

در مرحله‌ی بعدی با توجه به میانگین و ماتریس کوواریانس ایجاد شده بر اساس داده‌های اصلی محاسباتی، آماره‌ی  $T^2$  برای تمامی داده‌ها محاسبه و با حد مجاز آن مقایسه می‌شود. اگر از حد مجاز آن بیش‌تر شود، آن داده پرت قطعی شناخته می‌شود. در قسمت طراحی پارامترهای مدل نحوه‌ی تعیین حدود در روش تشریح شده است.

در نهایت ملاک نهایی تصمیم‌گیری برای داده‌ی پرت در این روش، داده‌ی پرت قطعی خواهد بود. بر این اساس اگر داده‌ی در رده پرت ثانویه تعیین شد ولی پس از

محاسبه‌ی آماره‌ی مربوط به آن، پرت قطعی نبود (از حد ایجاد شده‌ی نمودار کم‌تر بود) آن داده پرت محسوب نمی‌شود.

#### ۲.۳. طراحی پارامترهای مدل

دو پارامتر ضریب ناهمسانی و حدود کنترل باید در این روش طراحی شود. ابتدا باید پارامتر ضریب ناهمسانی مجاز را برای تأیید خوشه‌بندی تعیین کرد. برای تعیین این پارامتر برای هر حالت  $m$  و  $p$  مشخص، داده‌های نرمال با میانگین صفر و ماتریس واریانس - کوواریانس همانی شبیه‌سازی می‌شود. میانگین ضریب ناهمسانی برای خوشه‌های نهایی (سطر آخر ماتریس ناهمسانی) به‌عنوان این حد تعیین می‌شود.

در مرحله‌ی دوم باید حدود کنترل تنظیم شود. در این روش از حدود مجاز متغیر برای داده‌ها استفاده شده است. این حدود بر اساس تعداد داده‌ی موجود در خوشه‌ی داده‌های پرت ثانویه تغییر می‌کند. به‌طور کلی می‌توان به‌صورت نظری بر مبنای توزیع احتمال یا با شبیه‌سازی، حدود کنترل (LCL و UCL) را تعیین کرد.<sup>[۲۷]</sup> حد پایین در روش  $T^2$  صفر در نظر گرفته می‌شود، اما برای حد بالا در این روش به‌ازای هر حالت تعداد داده‌ی موجود در خوشه‌ی داده‌های اولیه یک مقدار در نظر گرفته می‌شود به‌طوری که نسبت تعداد سیگنال در هر حالت به تعداد هر حالت از داده‌ی موجود در خوشه‌ی اولیه مساوی شود. به‌عنوان مثال ورگاس برای آماره‌ی  $T^2$  با پارامترهای  $m = 30$  و  $p = 2$  حد بالای  $10/5$  را با شبیه‌سازی برای نسبت سیگنال خطای  $0/5$  محاسبه کرده است.<sup>[۲۷]</sup> این بدان معناست که در صورت شبیه‌سازی در حالت میانگین  $0$  و ماتریس واریانس همانی در  $5$  درصد مواقع نقاط بالای حد کنترل قرار می‌گیرد و سیگنال خطای اشتباه صادر می‌شود.

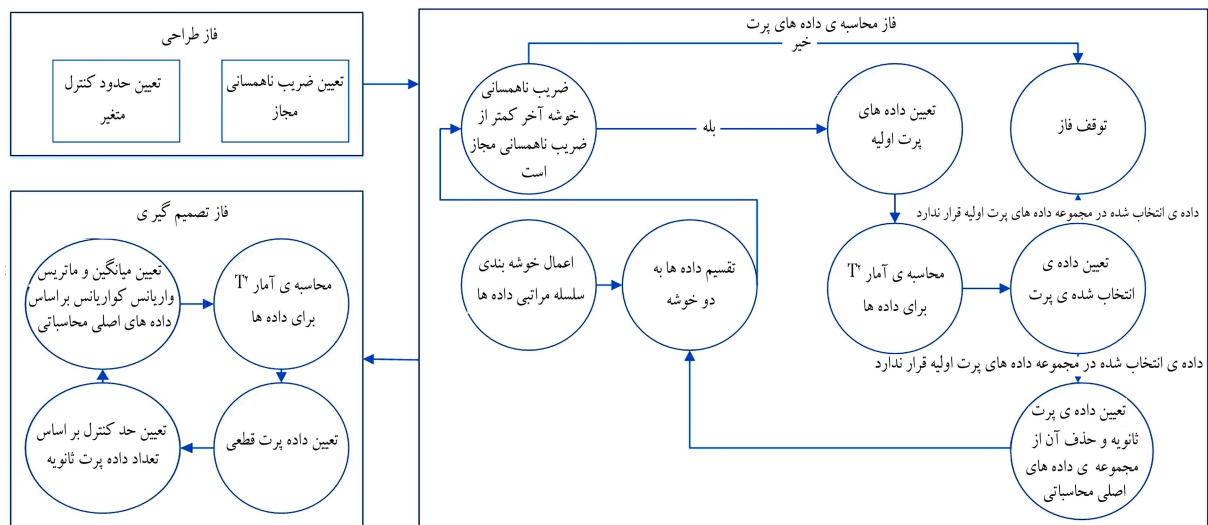
در حقیقت در این روش بر اساس تعداد داده‌ی ثانویه‌ی پرت شناسایی شده (به‌وسیله‌ی خوشه‌بندی سلسله‌مراتبی) حد کنترل بالای نمودار  $T^2$  تغییر می‌کند. برای حالت  $m = 30$  و  $p = 2$  در جدول ۱ محاسبات تعیین حدود کنترل متغیر بر اساس ضریب ناهمسانی  $1/2$  بر اساس  $5000$  شبیه‌سازی ارائه شده است. در سطر اول نشان داده شده که در هر مرحله از شبیه‌سازی چند داده‌ی پرت ثانویه شناسایی شده و بر اساس آن آماره‌ی  $T^2$  ساخته شده است. در سطر دوم در صورتی که بخواهیم  $0/5$  خطا داشته باشیم، تعداد سیگنال مجاز بیان شده است (از ضرب سطر اول در  $0/5$  سطر دوم حاصل می‌شود). در حقیقت این نسبت بیان می‌کند که از هر حالت شبیه‌سازی چندبار باید سیگنال اشتباه صادر شود.

به‌عنوان مثال در  $5000$  شبیه‌سازی، در  $1768$  حالت روش هیچ داده‌ی پرتی را در حالت اولیه شناسایی نمی‌کند و برای این حالت حد کنترل  $8/87$ ، در  $80$  حالت نمودار سیگنال خارج کنترل برای داده‌های با  $m = 30$  و  $p = 2$  صادر می‌کند. در حالت ایده‌آل باید این اعداد بر روی  $884$  تنظیم شود که امکان‌پذیر نیست. دقت شود که سایر حالت‌های بزرگ‌تر از  $8$  داده‌ی پرت اولیه به‌دلیل کم بودن احتمال رخ دادن در نظر گرفته نمی‌شود و حد کنترل آن‌ها نیز معادل حد کنترل  $28/5$  در نظر گرفته می‌شود. با این رویه به‌ازای حالت‌های مختلف از داده اولیه پرت، حدود کنترل متغیر تعریف می‌شود. توجه شود که تعداد حالات شبیه‌سازی در سطر اول باید دقیقاً برابر  $5000$  شود که سایر حالات که بیش از  $8$  داده پرت اولیه در خوشه‌بندی شناخته شده است. برای اختصار جدول ارائه نشده است رویه‌ی فوق برای هر یک از حالات  $m$  و  $p$  باید انجام گیرد و حدکنترل مورد نظر تعیین شود.

این ساختار ترکیبی و استفاده از حدود کنترل متغیر متناسب با تعداد داده‌های پرت انتخاب شده نوآوری و مزیت اصلی این پژوهش است. بر این اساس، در این روش هدف اصلی این است که با توجه به خوشه‌بندی صورت گرفته، به تعدادی

جدول ۱. مراحل طراحی حد کنترل به ازای نسبت سیگنال اشتباهی ۰/۰۵ برای  $m$  و  $p$  برابر ۳۰ و ۲ و ضریب ناهمسانی ۱/۲.

جمع کل	تعداد داده‌های پرت ثانویه									
	۸	۷	۶	۵	۴	۳	۲	۱	۰	
۴۹۶۱	۲۵	۴۷	۷۸	۱۲۸	۲۶۳	۴۸۹	۸۳۸	۱۳۲۵	۱۷۶۸	تعداد حالت ایجاد شده در شبیه‌سازی
۲۴۸,۰۵	۱,۲۵	۲,۳۵	۳,۹	۶,۴	۱۳,۱۵	۲۴,۴۵	۴۱,۹	۶۶,۲۵	۸۸,۴	تعداد سیگنال خطای مجاز
۲۵۴	۷	۱۶	۸	۱۹	۱۰	۲۵	۲۸	۶۱	۸۰	تعداد سیگنال خطا بر اساس حدود تعیین شده
	۲۸,۵	۲۸,۵	۳۳,۵	۲۸,۵	۳۰,۵	۲۷	۲۴,۳۹	۱۸,۲	۸,۸۷	حد بالای کنترل



شکل ۱. فلوجارت روش پیشنهادی برای تعیین داده پرت.

کم‌تری است به‌عنوان داده‌های پرت اولیه در نظر گرفته می‌شود. سپس داده‌های پرت اولیه بر این اساس مشخص می‌شود و آماره‌ی  $T^2$  برای آن‌ها محاسبه می‌شود. داده‌ی پرت انتخاب شده با توجه به این آماره و حدکنترل متغیر به دست می‌آید. در صورتی که این داده در مجموعه‌ی پرت اولیه باشد، به‌عنوان داده‌ی پرت قطعی شناخته می‌شود. پس از مشخص شدن داده‌های پرت از مجموعه داده‌ها حذف می‌شوند.

در فاز تصمیم‌گیری، میانگین و ماتریس واریانس کواریانس  $T^2$  بر اساس داده‌های باقیمانده در فاز دوم محاسبه خواهد شد، سپس مقادیر آماره‌ی  $T^2$  با توجه به حدکنترل به دست آمده در فاز اول برای تمامی داده‌ها محاسبه می‌شود.

### ۳.۳. مثال مفهومی

برای تشریح بهتر الگوریتم فوق از یک مثال عددی استفاده می‌شود. داده‌های جدول ۳ شامل ۲۵ داده ابتدایی ایجاد شده با توزیع نرمال دومتغیره با میانگین صفر و واریانس همانی ۱ است و ۵ داده‌ی آخر بر اساس  $nCP$  برابر ۲۰ (با میانگین ۳/۱۶) ایجاد شده است (در حقیقت این ۵ داده پرت محسوب می‌شوند).

داده‌ی پرت ابتدایی با انتخاب شده دست یافت و سپس براساس آن‌ها، آماره‌ی  $T^2$  محاسبه می‌شود که این آماره با حدکنترل مخصوص به خود که با توجه به تعداد داده‌ی پرت انتخاب شده مشخص شده، مقایسه می‌شود. روش پیشنهادی این پژوهش را می‌توان در ۳ فاز مطابق فلوجارت شکل ۱ طبقه‌بندی کرد:

۱. طراحی؛

۲. محاسبه‌ی داده‌ی پرت؛

۳. تصمیم‌گیری.

در فاز اول با استفاده از پارامترهای اندازه نمونه ( $m$ ) و تعداد متغیرها (تعداد شاخص‌های کیفی) ( $p$ ) ضریب ناهمسانی تعیین می‌شود. سپس با استفاده از پارامتر تعیین شده (ضریب ناهمسانی مجاز) و پارامتر نسبت سیگنال خطا، حدود کنترل متغیر مطابق جدول ۱ تعیین می‌شود.

در فاز دوم محاسبه‌ی داده پرت با توجه به ضریب ناهمسانی تعیین شده در فاز اول به دست می‌آید. محاسبه‌ی داده پرت در هر مرحله از تکرار الگوریتم خوشه‌بندی سلسله‌مراتبی صورت می‌گیرد و داده‌ها به دو خوشه تحت کنترل و خارج از کنترل تقسیم می‌شود (با توجه به ضریب ناهمسانی). خوشه‌ی که دارای تعداد اعضای

جدول ۲. مراحل خوشه‌بندی و آماره‌های محاسبه شده بر روی داده‌ها.

مرحله	داده‌های پرت اولیه	مقدار ضریب ناهمسازی	داده با بیشترین مقدار آماری $T^2$	مقدار آماری $T^2$	داده پرت انتخاب شده	داده پرت قطعی
۱	۳۰ - ۲۹ - ۲۸ - ۲۷ - ۲۶	۱/۱	۲۶	۹,۶۴۶۳۴۹۰۱۵	۲۶	۲۶
۲	۲۹ - ۲۸ - ۲۷ - ۲۶	۱/۱	۲۹	۹,۰۱۲۶۹۶۵۱۹	۲۹	۲۹(۳۰)
۳	۲۸ - ۲۷ - ۲۶	۱/۱	۲۸	۸,۹۷۲۸۱۱۷۶۵	۲۸	۲۸(۲۹)
۴	۲۷ - ۲۶	۱/۱۲	۲۷	۹,۰۱۰۰۲۰۹۱۷	۲۷	۲۷(۲۸)
۵	۲۶	۰,۷۰۷	۲۶	۱۳,۸۳۱۳۲۴۰۲	۲۶	۲۶(۲۷)
۶	۱۵	۰,۷۰۷	۱۶	۵,۴۵۰۰۳۴۷۳۱	-	-
بردار میانگین داده‌ها پس از حذف داده‌های پرت		ماتریس واریانس کوواریانس داده‌ها پس از حذف داده‌های پرت				
۰,۴۱۵	۰,۰۰۷۵	۰,۱۸۷	-	۰,۱۸۷	۰,۶۴	-
۰,۴۱۵	۰,۰۰۷۵	۰,۵۸۷۶	۰,۱۸۷	-	۰,۱۸۷	-

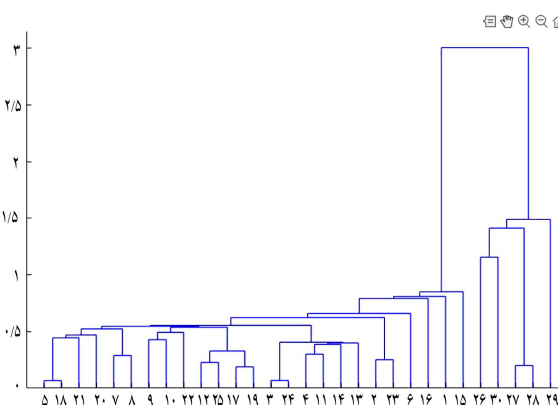
این مرحله آماری  $T^2$  برای ۲۹ داده در ستون چهارم جدول ۳ محاسبه شده است و داده‌ی ۲۹ فعلی و ۳۰ اصلی دارای بیشترین مقدار آماره است که در خوشه‌ی داده‌های پرت نیز قرار دارد و از داده‌ها حذف می‌شود (داده‌ی پرت ثانویه). مرحله‌ی سوم و چهارم و پنجم نیز به همین شکل انجام می‌گیرد و داده‌های ۲۸ فعلی (۲۹ اصلی) و ۲۷ فعلی (۲۸ اصلی) و ۲۶ فعلی (۲۷ اصلی) از مجموعه‌ی داده‌ها حذف می‌شوند. در مرحله‌ی ششم دندوگرام داده‌ی ۱۵ فعلی را داده‌ی پرت اولیه تشخیص می‌دهد ولی آماری  $T^2$  برای ۲۵ داده‌ی باقیمانده در داده‌ی ۱۶ دارای بیشترین مقدار است (ستون هشتم جدول ۳). در این قسمت فاز محاسبه‌ی داده‌ی پرت به پایان می‌رسد.

در فاز تصمیم‌گیری باید ابتدا بر مبنای ۲۵ داده‌ی باقیمانده میانگین و ماتریس واریانس - کوواریانس را محاسبه کرد. این مقادیر در جدول ۲ گزارش شده است (سطر انتهایی). در مرحله‌ی نهایی بر مبنای این مقادیر میانگین و ماتریس واریانس - کوواریانس آماری  $T^2$  برای کل داده‌ها محاسبه می‌شود که در ستون نهم جدول ۳ گزارش شده است. به دلیل این که در این الگوریتم ۵ داده‌ی پرت اولیه شناسایی شد، در این حالت حد کنترل باید (۲۸/۵) حد کنترل به‌ازای ۵ داده‌ی ثانویه پرت در جدول ۱) در نظر گرفته شود که با توجه به این حد کنترل هر ۵ داده‌ی انتهایی پرت محسوب و با رنگ قرمز مشخص شده است. این ۵ داده به‌عنوان داده‌ی پرت قطعی در نظر گرفته می‌شود که به‌طور کاملاً صحیح توسط روش شناسایی شده‌اند.

#### ۴. ارزیابی عملکرد

##### ۱.۴. عملکرد روش بر اساس تشخیص یک داده‌ی پرت

ورگاس برای تشخیص داده‌ی پرت از شاخص ncp بر اساس رابطه‌ی ۵ برای ایجاد  $k$  داده‌ی پرت استفاده کرد.<sup>[۲۴]</sup> او بر اساس شبیه‌سازی، سیگنال خطا را در حد ۰,۵ تنظیم کرد و ملاک دقت روش را تشخیص حداقل یک سیگنال در داده‌ها قرار داد. در نمودار کنترل پیشنهادی این پژوهش، به‌ازای هر حالت  $m$  و  $p$  حدود کنترل (حد کنترل بالا) تعیین می‌شود و سپس در ۵۰۰ شبیه‌سازی حالت‌هایی که سیگنال داده شود



شکل ۲. دندوگرام داده‌ها در اولین مرحله.

با توجه به مقادیر آماری  $T^2$  در سومین ستون جدول ۳، روش  $T^2$  معمولی با توجه به حدود کنترل شبیه‌سازی شده ۱۰/۵ برای  $m = 3$  و  $p = 2$  هیچ داده‌ی پرتی را در این مجموعه نشان نمی‌دهد.

در نمودار کنترل پیشنهادی، اولین مرحله با خوشه‌بندی سلسله‌مراتبی با نرم‌افزار متلب دندوگرام شکل ۲ به دست می‌آید. در تقسیم داده‌ها به دو خوشه، مطابق شکل ۲ داده‌های ۲۶، ۲۷، ۲۸، ۲۹ و ۳۰ پرت اولیه هستند. به‌علت این که ضریب ناهمسازی خوشه‌ی آخر بیش‌تر از ۱/۲ است این خوشه‌بندی تأیید می‌شود. سپس آماری  $T^2$  معمولی برای کل داده‌ها محاسبه شده است (نتایج مربوط به این مرحله، در ستون سوم جدول ۳ گزارش شده است). مشاهده می‌شود، داده‌ی ۲۶ دارای بیشترین آماری  $T^2$  است و در مجموعه‌ی داده‌های پرت اولیه قرار دارد و به‌عنوان داده‌ی پرت ثانویه از داده‌ها حذف می‌شود. در جدول ۳ این داده با رنگ زرد مشخص شده است (مقدار آماره ۹,۶۴۶ است).

در مرحله‌ی دوم برای ۲۹ داده‌ی بعدی دندوگرام رسم می‌شود (به‌منظور اختصار سایر دندوگرام‌ها نشان داده نشده است). در این مرحله داده‌های ۲۶، ۲۷، ۲۸ و ۲۹ در خوشه‌ی داده‌های پرت اولیه قرار دارند (دقت شود که در داده‌های اصلی این داده‌ها ۲۷، ۲۸، ۲۹ و ۳۰ هستند و در جدول ۲ در پارتیز نشان داده شده‌اند). در

جدول ۳. داده‌های مثال تشریحی به همراه مقادیر آماره‌های محاسبه شده برای هر داده در هر مرحله.

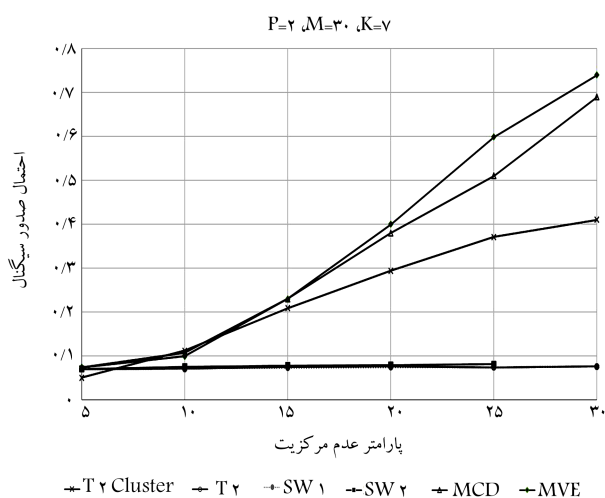
آماره	آماره $T^2$ در مرحله ۵	آماره $T^2$ در مرحله ۴	آماره $T^2$ در مرحله ۳	آماره $T^2$ در مرحله ۲	آماره $T^2$ در مرحله ۱	آماره معمولی $T^2$	$X_2$	$X_1$	
۵.۱۳۰	۵.۱۳۰	۲.۴۵۷	۱.۸۵۷	۱.۱۱۰	۱.۱۶۸	۱.۲۲۸	۰.۱۵۷	۱.۸۱۲	۱
۲.۲۵۶	۲.۲۵۶	۱.۵۴۴	۱.۳۷۴	۱.۴۶۱	۱.۰۷۸	۰.۷۰۸	۱.۱۶۱	-۰.۲۵۴	۲
۱.۴۸۶	۱.۴۸۶	-۰.۹۲۱	-۰.۸۰۸	-۰.۷۷۹	-۰.۷۶۶	-۰.۸۰۰	-۰.۴۲۲	-۰.۸۲۵	۳
۱.۲۲۷	۱.۲۲۷	۱.۰۱۹	۰.۹۹۸	۱.۰۷۲	۰.۹۴۰	۰.۸۳۱	-۰.۷۴۳	-۰.۳۰۵	۴
۱.۶۲۳	۱.۶۲۳	۱.۳۵۳	۱.۳۶۷	۱.۱۳۳	۱.۱۴۳	۱.۰۶۳	-۰.۴۶۴	-۰.۸۹۶	۵
۴.۳۰۹	۴.۳۰۹	۴.۲۳۷	۴.۳۹۶	۳.۹۲۹	۳.۸۱۶	۳.۴۲۹	۰.۹۸۱	-۱.۳۶۳	۶
۲.۷۷۸	۲.۷۷۸	۲.۵۹۵	۲.۶۰۷	۲.۶۹۳	۲.۲۸۴	۱.۸۴۲	-۱.۲۳۳	-۰.۰۴۴	۷
۲.۳۹۷	۲.۳۹۷	۲.۴۴۶	۲.۵۳۵	۲.۵۵۰	۲.۲۲۸	۱.۸۲۷	-۱.۱۳۲	-۰.۲۲۷	۸
۰.۳۶۵	۰.۳۶۵	۰.۱۰۵	۰.۰۵۷	۰.۰۸۲	۰.۰۶۱	۰.۰۵۶	۰.۴۰۵	۰.۲۹۵	۹
۰.۴۰۳	۰.۴۰۳	۰.۱۵۰	۰.۱۱۷	۰.۰۷۴	۰.۱۱۰	۰.۱۴۵	۰.۰۳۹	-۰.۵۱۵	۱۰
۰.۵۶۵	۰.۵۶۵	۰.۴۷۸	۰.۴۸۶	۰.۵۳۹	۰.۵۱۷	۰.۵۱۲	-۰.۴۴۴	-۰.۳۰۱	۱۱
۰.۷۵۷	۰.۷۵۷	۰.۷۲۰	۰.۷۴۷	۰.۷۸۵	۰.۶۸۸	۰.۵۶۱	۰.۶۹۳	-۰.۱۶۰	۱۲
۰.۵۳۴	۰.۵۳۴	۰.۴۱۳	۰.۴۱۸	۰.۴۲۲	۰.۴۶۳	۰.۵۱۳	-۰.۱۲۸	-۰.۵۴۵	۱۳
۰.۰۶۱	۰.۰۶۱	۰.۰۹۹	۰.۱۳۹	۰.۱۸۰	۰.۲۱۶	۰.۲۵۴	-۰.۱۰۶	-۰.۱۱۱	۱۴
۴.۳۸۹	۴.۳۸۹	۲.۶۸۷	۲.۳۳۶	۱.۸۸۲	۱.۹۶۰	۲.۰۵۹	-۰.۱۹۲	-۱.۶۴۴	۱۵
۵.۴۵۰	۵.۴۵۰	۲.۷۷۴	۲.۰۴۴	۲.۰۹۶	۱.۳۲۷	۰.۷۳۲	۱.۵۵۷	-۰.۹۵۳	۱۶
۱.۶۶۱	۱.۶۶۱	۱.۷۴۴	۱.۸۵۰	۱.۷۳۶	۱.۶۶۹	۱.۴۸۸	۰.۷۳۰	-۰.۷۵۴	۱۷
۱.۴۴۵	۱.۴۴۵	۱.۲۳۲	۱.۲۵۶	۱.۰۵۶	۱.۰۶۲	۰.۹۸۶	-۰.۴۶۲	-۰.۸۲۹	۱۸
۱.۱۵۵	۱.۱۵۵	۱.۲۰۸	۱.۲۸۹	۱.۲۱۳	۱.۱۹۸	۱.۱۰۶	۰.۵۵۷	-۰.۶۸۰	۱۹
۲.۲۵۰	۲.۲۵۰	۲.۳۳۶	۲.۴۶۳	۲.۲۹۴	۲.۱۳۱	۱.۸۲۹	-۰.۹۱۵	-۰.۷۰۳	۲۰
۰.۸۲۸	۰.۸۲۸	-۰.۸۸۲	-۰.۹۵۳	-۰.۹۱۵	-۰.۸۷۶	-۰.۷۸۷	-۰.۵۶۰	-۰.۳۹۵	۲۱
۱.۵۶۳	۱.۵۶۳	۰.۶۹۵	۰.۵۱۷	۰.۲۹۰	۰.۳۳۱	۰.۳۷۰	۰.۰۷۱	۱.۰۰۷	۲۲
۳.۲۱۲	۳.۲۱۲	۲.۴۲۷	۲.۲۴۷	۲.۳۵۶	۱.۷۹۴	۱.۲۲۲	۱.۴۰۰	-۰.۱۷۶	۲۳
۱.۴۹۸	۱.۴۹۸	۰.۹۳۱	۰.۸۱۶	۰.۸۰۴	۰.۷۷۲	۰.۷۹۲	-۰.۴۸۰	-۰.۷۹۰	۲۴
۰.۶۵۷	۰.۶۵۷	۰.۷۱۳	۰.۷۷۳	۰.۷۸۱	۰.۷۴۰	۰.۶۵۷	۰.۵۷۰	-۰.۳۵۱	۲۵
۵۳.۰۴۷	-	-	-	-	-	۹.۶۴۶	۴.۷۰۵	۳.۰۶۳	۲۶
۳۳.۸۲۲	-	۱۳.۸۳۱	۸.۵۴۰	۶.۱۰۴	۴.۲۸۳	۳.۸۲۰	۲.۲۲۴	۳.۹۹۶	۲۷
۳۵.۵۴۸	-	-	۹.۰۱۰	۶.۶۹۳	۴.۶۴۷	۳.۸۷۱	۲.۴۲۴	۴.۰۰۳	۲۸
۴۳.۷۰۸	-	-	-	۸.۹۷۳	۸.۶۱۹	۸.۹۴۱	۱.۲۳۵	۵.۱۱۱	۲۹
۴۸.۴۱۸	-	-	-	-	۹.۰۱۳	۵.۹۲۶	۳.۸۲۴	۳.۸۱۳	۳۰

به عنوان حالت صحیح شناخته می‌شود و نسبت تعداد حالت‌هایی که روش داده پرت تشخیص می‌دهد به تعداد کل شبیه‌سازی ثبت می‌شود. در شکل ۳ تا ۵ مقایسه‌ی ارزیابی عملکرد نمودار  $T^2$  خوشه‌بندی در شناسایی نقاط پرت با MCD، MVE،  $SW_1$ ،  $SW_2$  و  $T^2$  معمولی بر اساس  $m = 30$  و  $p = 2$  ارائه شده است. در شکل‌های ۳ تا ۵ ارزیابی عملکرد به ترتیب با ۳، ۵ و ۷ داده‌ی پرت آورده شده است. مقایسه‌ی عملکرد نمودار  $T^2$  خوشه‌بندی با نمودار  $T^2$  معمولی در جداول ۴ تا ۷ آمده است. در جدول ۴ مقایسه‌ی عملکرد نمودار  $T^2$  با نمودار کنترل  $T^2$  خوشه‌بندی از نظر درصد تشخیص داده‌ی پرت به ازای  $P = 3$  آورده شده است. هر دو نمودار در شاخص عدم مرکزیت برابر با ۵ ( $nep = 5$ ) در تعداد نقاط پرت

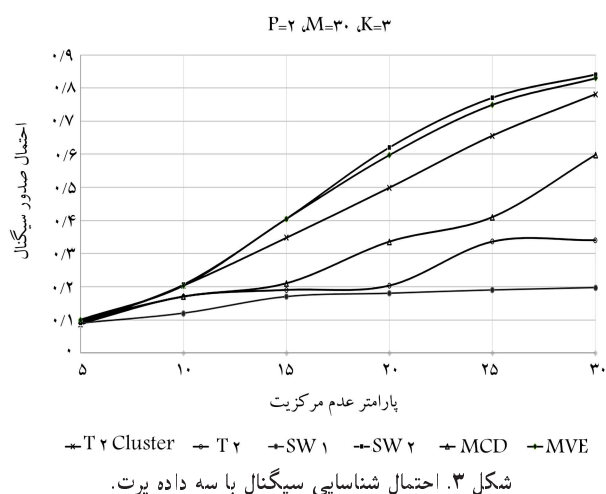
مختلف و با تعداد داده‌ها متفاوت، عملکرد تقریباً مشابه دارند اما با افزایش شاخص عدم مرکزیت عملکرد نمودار  $T^2$  خوشه‌بندی بهتر می‌شود. در جدول ۵ نیز درصد تشخیص دو نمودار با ۵ متغیر آورده شده است. در این جدول نیز با افزایش مقدار nep عملکرد نمودار خوشه‌بندی به مراتب بهتر بوده است. در جدول ۶ درصد تشخیص داده پرت، به ازای  $p = 10$  نشان داده شده است. عملکرد هر دو نمودار مانند جداول ۴ و ۵ در ( $nep = 5$ ) تقریباً مشابه است و با افزایش مقدار شاخص عدم مرکزیت عملکرد نمودار خوشه‌بندی افزایش یافته‌است. نکته قابل مشاهده دیگر اینست که به طور کلی با افزایش تعداد متغیرها عملکرد هر دو نمودار کاهش می‌یابد.

جدول ۴. درصد تشخیص داده پرت به ازای  $p = 3$  در نمودار  $T^2$  و نمودار  $T^2$  خوشه بندی به ازای مقادیر مختلف شاخص عدم مرکزیت.

ncp	$p = 3$								Methods	
	$m = 30$			$m = 50$			$m = 100$			
	$K = 2$	4	6	2	5	10	5	10		20
5	0.076	0.076	0.05	0.093	0.082	0.059	0.121	0.084	0.045	$T^2$
	0.089	0.076	0.059	0.092	0.09	0.063	0.133	0.107	0.064	$T^2$ Cluster
15	0.263	0.077	0.044	0.447	0.137	0.057	0.441	0.137	0.039	$T^2$
	0.42	0.257	0.118	0.406	0.356	0.104	0.586	0.284	0.0588	$T^2$ Cluster
25	0.449	0.08	0.051	0.761	0.182	0.038	0.688	0.171	0.045	$T^2$
	0.741	0.497	0.246	0.772	0.701	0.212	0.896	0.566	0.098	$T^2$ Cluster



شکل ۵. احتمال شناسایی سیگنال با هفت داده پرت.



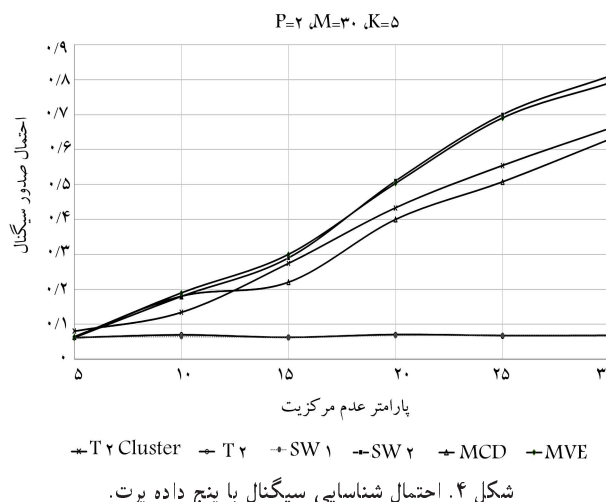
شکل ۳. احتمال شناسایی سیگنال با سه داده پرت.

اصلی و ۹/۰ بر روی سایر قسمت‌ها ایجاد می‌شوند. داده‌های پرت با میانگین ۵ ایجاد می‌شوند. حدود سیگنال خطا در دو روش بر روی حد مشخصی با شبیه‌سازی تنظیم می‌شوند. در این روش اگر تمامی داده‌های پرت شناسایی شوند، دقت روش مورد تأیید است. در جدول ۷ دو عدد برای ستون‌های  $T^2$  و  $T^2$  خوشه‌بندی ذکر شده است. عدد اول نسبت تعداد دفعاتی است که روش توانسته است کلیه داده‌های پرت را شناسایی نماید و عدد دوم نرخ سیگنال خطا در آن حالت است. سیگنال خطا هنگامی حساب می‌شود که روش در داده‌های غیرپرت (تفاوتی ندارد که در یک داده یا چند داده از یک مجموعه‌ی داده سیگنال صادر شود) سیگنال دهد.

### ۳.۴. عملکرد روش بر اساس مجموعه داده فسفر

در این مجموعه داده  $m = 18$  و  $p = 2$  است. [۳۴،۳۳] مطابق جدول ۱ ابتدا باید حدود کنترل برای سیگنال ۰/۰۵ خطا تعیین شود.

در این مجموعه داده حدود کنترل در سطح خطای ۰/۰۵ به صورت [۲۰/۵، ۳۳/۵، ۵۳/۵، ۴۴/۵، ۴۹، ۵۷/۳۹، ۴۷/۲، ۱۰/۰۷] است. در اجرای الگوریتم، مجموعه داده‌های پرت ثانویه به ترتیب ۱۷، ۶، ۱، ۸، ۴ است که باید حد



شکل ۴. احتمال شناسایی سیگنال با پنج داده پرت.

### ۲.۴. عملکرد روش بر اساس تشخیص مجموعه‌ی داده‌های پرت

در این قسمت بر اساس روش آلفارو و همکاران عمل می‌شود. [۲] در این بخش داده‌های تحت کنترل با میانگین ۰ و ماتریس واریانس کوواریانس ۱ بر روی قطر



جدول ۵. درصد تشخیص داده پرت به ازای  $p = 5$  در نمودار  $T^2$  و نمودار  $T^2$  خوشه بندی به ازای مقادیر مختلف شاخص عدم مرکزیت.

پارامتر عدم مرکزیت	$p = 5$									روش
	$m = 30$			$m = 50$			$m = 100$			
	۲	۴	۶	۲	۵	۱۰	۵	۱۰	۲۰	
۵	۰٫۰۶۸	۰٫۰۷۵	۰٫۰۶۲	۰٫۰۷۶	۰٫۰۷۱	۰٫۰۶۲	۰٫۰۹۴	۰٫۰۷۴	۰٫۰۴۵	$T^2$
	۰٫۰۸۴۴	۰٫۰۶۳	۰٫۰۶۹	۰٫۰۶۳۶	۰٫۰۶۰۸	۰٫۰۴۲	۰٫۰۹	۰٫۰۷۹۲	۰٫۰۶۵	$T^2$ Cluster
۱۵	۰٫۱۶۹	۰٫۰۷۶	۰٫۰۴۹	۰٫۲۷	۰٫۱۱	۰٫۰۵	۰٫۲۶۷	۰٫۰۹۳	۰٫۰۵	$T^2$
	۰٫۲۳۳	۰٫۱۴۹	۰٫۰۸۷	۰٫۲۲۶	۰٫۱۵۸	۰٫۰۵۳	۰٫۳۵۵	۰٫۱۴۸	۰٫۰۵۹	$T^2$ Cluster
۲۵	۰٫۲۸	۰٫۰۸۴	۰٫۰۵۳	۰٫۵۳	۰٫۱۱۸	۰٫۰۴۶	۰٫۴۵۶	۰٫۱۵۵	۰٫۰۶	$T^2$
	۰٫۴۹۹	۰٫۳۰۹	۰٫۱۲۴	۰٫۵۰۹	۰٫۳۷۴	۰٫۰۷۶	۰٫۷۲۷	۰٫۲۷۲	۰٫۰۶۷	$T^2$ Cluster

جدول ۶. درصد تشخیص داده پرت به ازای  $p = 10$  در نمودار  $T^2$  و نمودار  $T^2$  خوشه بندی به ازای مقادیر مختلف شاخص عدم مرکزیت.

پارامتر عدم مرکزیت	$p = 10$									روش
	$m = 30$			$m = 50$			$m = 100$			
	۲	۴	۶	۲	۵	۱۰	۵	۱۰	۲۰	
۵	۰٫۰۶۱	۰٫۰۴۹	۰٫۰۴۸	۰٫۰۵۶	۰٫۰۵۸	۰٫۰۶۳	۰٫۰۶۸	۰٫۰۷۶	۰٫۰۵۸	$T^2$
	۰٫۰۵۹	۰٫۰۵۹	۰٫۰۵۲	۰٫۰۶۶	۰٫۰۵۸	۰٫۰۵۹	۰٫۰۶۵	۰٫۰۶۱	۰٫۰۵۴	$T^2$ Cluster
۱۵	۰٫۰۶۲	۰٫۰۵۶	۰٫۰۵۵	۰٫۱۳۳	۰٫۰۷۱	۰٫۰۵۳	۰٫۱۳۸	۰٫۰۷	۰٫۰۵۲	$T^2$
	۰٫۱۲۹	۰٫۰۶۹	۰٫۰۵۴	۰٫۱۷۳	۰٫۰۹۸	۰٫۰۵۹	۰٫۱۵۹	۰٫۰۹۵	۰٫۰۶	$T^2$ Cluster
۲۵	۰٫۱	۰٫۰۵۷	۰٫۰۵۷	۰٫۲۷۵	۰٫۰۸۲	۰٫۰۵۹	۰٫۲۲۹	۰٫۰۷۲	۰٫۰۴۸	$T^2$
	۰٫۲۶	۰٫۰۸۹	۰٫۰۶۴	۰٫۴۱	۰٫۱۴۶	۰٫۰۶۱	۰٫۳۵۷	۰٫۱۱۷	۰٫۰۵۹	$T^2$ Cluster

جدول ۷. درصد تشخیص صحیح تمامی داده های پرت و سیگنال خطا در  $m = 25$  برای دو نمودار  $T^2$  و  $T^2$  خوشه بندی.

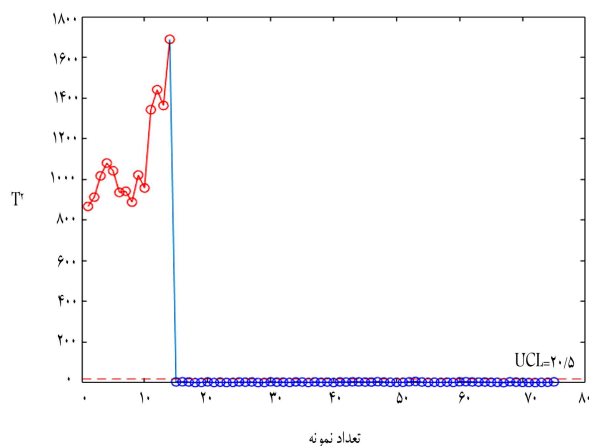
$p$	$k$	$T^2$		$T^2$ Cluster	
		درصد تشخیص	سیگنال خطا	درصد تشخیص	سیگنال خطا
۲	۰	-	۱٫۱	-	۱٫۵۶
۲	۱	۶۳٫۱۲	۰٫۳۴	۵۱٫۹۴	۱٫۷۶
۲	۲	۰٫۱۸	۰٫۲۴	۳۱٫۶۶	۱٫۵۶
۲	۳	۰	۰٫۰۲	۲۳٫۱۶	۱٫۶۴
۳	۰	-	۰٫۲۲	-	۰٫۳
۳	۱	۳۱٫۸	۰٫۱	۲۰٫۸۲	۰٫۳۲
۳	۲	۰٫۰۴	۰٫۰۸	۷	۰٫۳۸
۳	۳	۰	۰٫۲	۳٫۹۸	۰٫۶

#### ۴.۴. عملکرد روش بر اساس مجموعه داده ها و کینز

این مجموعه داده شامل ۷۵ مشاهده با چهار متغیر است، [۲۵] در این مجموعه داده ۱۴ داده ای ابتدایی پرت هستند. در شکل ۸ مشاهده می شود نمودار  $T^2$  معمولی از ۱۴ داده پرت داده ای ۱۲ و ۱۴ را شناسایی کرده است در صورتی که نمودار  $T^2$  خوشه بندی در شکل ۹ داده های ۱، ۲ تا ۱۴ را پرت ثانویه تشخیص می دهد و

کنترل آن را ۵۳٫۵ در نظر گرفت. نمودار آماره های  $T^2$  معمولی و خوشه بندی مطابق شکل ۶ و ۷ است.

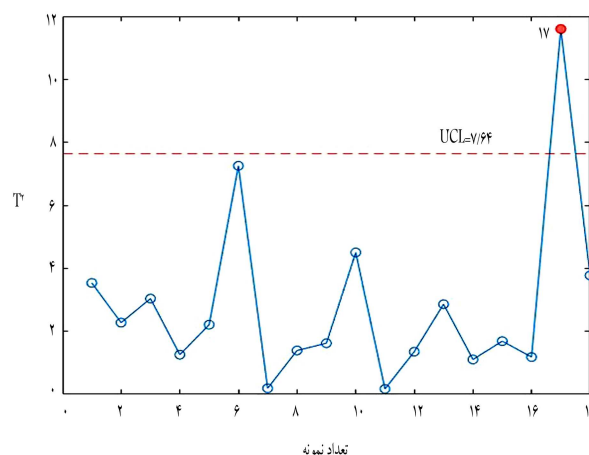
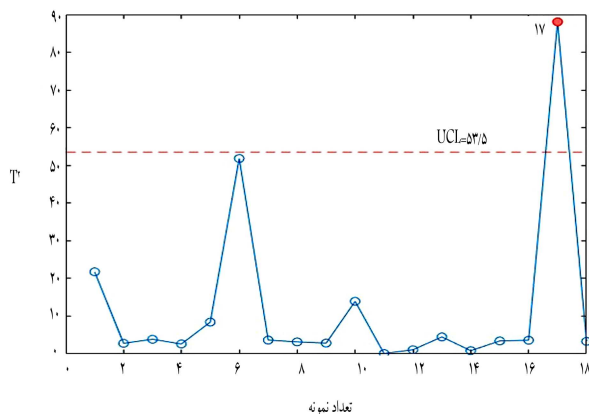
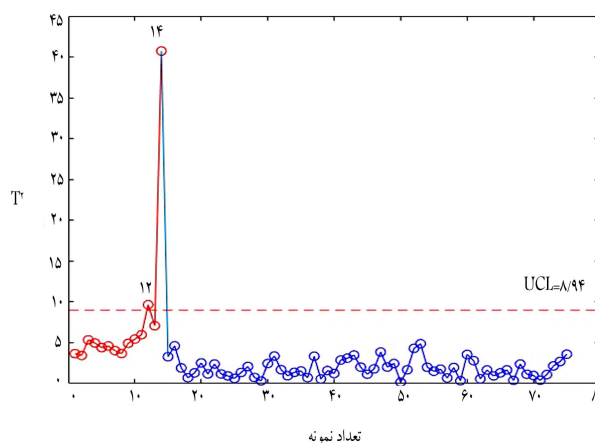
در این مجموعه داده یک داده ای پرت وجود دارد. نمودار  $T^2$  معمولی در صورتی که مجموعه داده فقط دارای یک داده پرت باشد، عالی عمل می کند. [۱۰] با توجه به شکل ۶ و ۷ نیز مشاهده می شود عملکرد هر نمودار مشابه هم است و هر دو نمودار داده ای پرت ۱۷ را شناسایی کرده اند.

شکل ۹. نمودار  $T^2$  خوشه‌بندی در مجموعه داده‌ی هاوکینز.

انحرافی در نمونه یا داده‌های نامتعارف قرار نمی‌گیرد و نقاط پرت موجود در فاز اول نمودارهای کنترلی چندمتغیره را شناسایی و حذف می‌کند. روش پیشنهادی برای بهبود عملکرد نمودار کنترل  $T^2$  در مثال مفهومی تشریح شد. برای تأیید عملکرد روش پیشنهادی نسبت به سایر روش‌های موجود در ادبیات که برای بهبود نمودار کنترل  $T^2$  معرفی شده است، شبیه‌سازی و مقایسه انجام شد. بنابراین ابتدا عملکرد نمودار کنترل پیشنهادی در شناسایی نقاط پرت با روش حجم بیضی‌وار (MVE) و حداقل دترمینان کوواریانس (MCD) رسن<sup>[۱۸]</sup>، برآوردگرهای SW۱، SW۲ سالیوان و وودال<sup>[۱۷]</sup> و نمودار  $T^2$  معمولی بر اساس  $m = 3$  و  $p = 2$  مقایسه شده است. در حالتی که سه نقطه‌ی پرت وجود داشته باشد، مطابق شکل ۳ مشاهده می‌شود، نمودار کنترل پیشنهادی ( $T^2$  خوشه‌بندی) از نمودار  $T^2$  معمولی، SW۱ سالیوان و وودال<sup>[۱۷]</sup> و MCD رسن<sup>[۱۸]</sup> به مراتب عملکرد بهتری دارد و عملکرد تقریباً مشابهی با SW۲ سالیوان و وودال<sup>[۱۷]</sup> و MVE رسن<sup>[۱۸]</sup> دارد. با این حال این روش‌ها در مقایسه با روش خوشه‌بندی، به لحاظ عملی بسیار وقت‌گیر است.

چنان‌که گفته شد هدف پژوهش حاضر استفاده از خوشه‌بندی سلسه‌مراتبی برای بهبود نمودار کنترل  $T^2$  معمولی است، لذا نمودار کنترل پیشنهادی ( $T^2$  خوشه‌بندی) با نمودار  $T^2$  معمولی در حالت‌های مختلف شبیه‌سازی و مقایسه شد. در جدول ۴ نشان می‌دهد به ازای  $n_{cp} = 5$  با دو نقطه‌ی پرت، درصد تشخیص داده‌ی پرت در نمودار کنترل  $T^2$  خوشه‌بندی نزدیک به نمودار  $T^2$  معمولی است که با افزایش اندازه‌ی ( $m$ ) این اختلاف کم‌تر نیز شده است، اما با افزایش نقاط پرت درصد تشخیص داده‌ی پرت در نمودار  $T^2$  معمولی به شدت کاهش داشته است. همچنین مشاهده می‌شود با افزایش شاخص عدم مرکزیت ( $n_{cp}$ ) قدرت نمودار کنترل  $T^2$  خوشه‌بندی نسبت به نمودار  $T^2$  معمولی به مراتب بهتر است. در جدول ۵ تعداد متغیرها از ۳ به ۵ افزایش داده شد؛ مشاهده می‌شود با افزایش تعداد متغیرها عملکرد هر دو نمودار افت دارد اما در نمودار کنترل  $T^2$  خوشه‌بندی کاهش شدیدی حاصل نشده است.

برای تحلیل بیشتر، عملکرد نمودار کنترل  $T^2$  خوشه‌بندی با عملکرد نمودار  $T^2$  معمولی بر اساس تشخیص مجموعه داده‌های پرت، شبیه‌سازی و مقایسه شد. همان‌طور که در جدول ۷ مشاهده می‌شود با افزایش تعداد نقاط پرت عملکرد نمودار  $T^2$  خوشه‌بندی در شناسایی صحیح تمامی نقاط پرت به مراتب خیلی بهتر از نمودار  $T^2$  معمولی است. اما نرخ سیگنال خطا در نمودار  $T^2$  خوشه‌بندی بالاتر بوده که با

شکل ۶. نمودار  $T^2$  معمولی در مجموعه داده فسفر.شکل ۷. نمودار  $T^2$  خوشه‌بندی در مجموعه داده فسفر.شکل ۸. نمودار  $T^2$  معمولی در مجموعه داده‌ی هاوکینز.

به‌ازای ۱۴ مقدار داده در مجموعه داده‌های پرت ثانویه این روش مقدار UCL را ۲۰/۵ گزارش می‌دهد.

## ۵. بحث

با توجه به حساسیت پایین نمودار کنترل  $T^2$ ، برای بهبود آن، نمودار کنترل پیشنهادی بر مبنای روش خوشه‌بندی سلسه‌مراتبی ارائه شده است که تحت تأثیر داده‌های

نیازمند تلاش‌های محاسباتی فوق‌العاده زیادی هستند. به همین علت، مطلوب است که نمودارهای کنترل قوی برای موقعیت‌های چندمتغیره ایجاد شود که همچنین تعادل خوبی بین قابلیت تشخیص و تلاش محاسبات را حفظ کند. لذا در این تحقیق نیز برای این منظور از تکنیک خوشه‌بندی سلسه‌مراتبی در طراحی نمودار کنترل  $T^2$  هتلینگ در فاز اول استفاده شد. به عبارتی در این تحقیق با استفاده از رویکردهای خوشه‌بندی توان نمودار را در تشخیص داده‌های پرت افزایش داده شد. عملکرد روش پیشنهادی به دو صورت مورد ارزیابی قرار گرفت. با استفاده از شاخص عدم مرکزیت نمودار کنترل  $T^2$  معمولی و  $T^2$  خوشه‌بندی در ابعاد ۳، ۵ و ۱۰ و با تعداد داده‌های ۳۰، ۵۰ و ۱۰۰ مقایسه شد. نتایج به طور کلی نشان داد که با افزایش شاخص عدم مرکزیت عملکرد نمودار کنترل  $T^2$  خوشه‌بندی افزایش یافته است. همچنین نتایج ارزیابی نیز با استفاده از روش الفارو همکاران نشان داد که نمودار کنترل  $T^2$  خوشه‌بندی در شناسایی تغییرات بسیار کارا تر بوده اما نرخ خطا بالاست که نرخ خطا هم با افزایش تعداد متغیرها کاهش یافته است. مطابق جدول ۷ نمودار کنترل  $T^2$  خوشه‌بندی در ۲۵ داده‌ی ۳ متغیره با ۳ داده پرت، نسبت تعداد دفعاتی است که روش توانسته است کلیه‌ی داده‌های پرت با نرخ خطای ۰/۶ را شناسایی کند معادل ۳/۹۸ بوده است در حالی که نمودار  $T^2$  معمولی قادر به شناسایی نقاط پرت نبوده است. در مرحله‌ی دوم با استفاده از یک مجموعه داده‌ی فسفر و هاوکینز نمودار کنترل  $T^2$  خوشه‌بندی طراحی و با نتایج به دست آمده از نمودار کنترل  $T^2$  معمولی مقایسه شد. نتایج اجرای نمودارهای کنترل نشان می‌دهد که نمودار کنترل  $T^2$  خوشه‌بندی در شناسایی نقاط پرت کارا تر است.

## ۷. ملاحظات

این مقاله مستخرج از طرح شماره ۴۹۴۰۹ مورخه ۹۷/۱۲/۲۷ با حمایت معاونت پژوهشی دانشگاه فردوسی مشهد است.

## پانویس‌ها

1. masking
2. swamping
3. trimming
4. stalactite
5. minimum volume ellipsoid (MVE)
6. minimum covariance determinant
7. Subsampling
8. reweighted minimum covariance determinant

## منابع (References)

1. Fazel Zarandi, M.H., Alaeddini, A. and Turksen, I.B. "A hybrid fuzzy adaptive sampling-run rules for Shewhart control charts", *Information Sciences*, **178**, pp. 1152-1170 (2008).

افزایش تعدادمتغیرها ( $P$ ) سیگنال خطا کاهش یافته است.

برای تحلیل عملکرد نمودار  $T^2$  خوشه‌بندی، مقایسه‌ی عملکرد آن با نمودار  $T^2$  معمولی روی داده‌های واقعی فسفر و هاوکینز نیز انجام شد. در مجموعه داده‌ی فسفر یک داده پرت وجود دارد و عملکرد نمودار  $T^2$  معمولی مشابه نمودار کنترل  $T^2$  خوشه‌بندی است، در حالی که در مجموعه داده‌ی هاوکینز که ۱۴ داده‌های پرت وجود دارد نمودار  $T^2$  معمولی فقط سه داده‌ی پرت را شناسایی کرده است اما نمودار کنترل  $T^2$  خوشه‌بندی هر ۱۴ داده پرت شناسایی کرده است، بنابراین با توجه به مقایسات با نمودارهای دیگر در بخش اول ارزیابی و مقایسات مختلف نمودار با نمودار  $T^2$  معمولی می‌توان نتیجه گرفت که عملکرد نمودار کنترل  $T^2$  خوشه‌بندی کارا تر از نمودار  $T^2$  معمولی است.

## ۶. نتیجه‌گیری

در بسیاری از شرایط عملی، اغلب لازم است به طور همزمان ویژگی‌های چند کیفیت کنترل شود. شناسایی پرت‌های بالقوه یا مشاهدات مؤثر در مجموعه داده‌های چندگانه، قبل از شروع به تحلیل‌های پیشرفته‌تر آماری، یک وظیفه‌ی بسیار مهم است. نمودار کنترل چندمتغیره رایج‌ترین ابزار آماری مورد استفاده در چنین شرایط است، اما نمودار کنترل چندمتغیره‌ی کلاسیک نمی‌تواند به طور مؤثری داده‌های پرت را تشخیص دهد. تاکنون روش‌های متعددی برای بهبود عملکرد نمودارهای کنترل چندمتغیره ارائه شده است. از جمله می‌توان به برآوردگر باثبات حداقل حجم بیضی‌وار، حداقل دترمینان واریانس، برآوردگر سالیوان - وودال و حداقل دترمینان واریانس بازموزون اشاره کرد. همه‌ی این روش‌ها به دنبال به دست آوردن برآوردی از میانگین و ماتریس کوواریانس مجموعه داده‌های چندمتغیره‌اند که نسبت به نقاط غیرمعمول و دورافتاده، مقاوم بوده و تحت تأثیر این نقاط قرار نگیرند. در حالی که این روش‌ها یعنی MVE و MCD طبق اظهارنظر فان و همکارانش (۲۰۱۲)

2. De Vries, A. and Conlin, B.J. "A comparison of the performance of statistical quality control charts in a dairy production system through stochastic simulation", *Agricultural Systems*, **85**, pp. 317-341 (2005).
3. Shewhart, W.A., *Statistical Methods From the Viewpoint of Quality Control*, Republished in 1986 by Dover Publications, New York, NY (1939).
4. Runger, G.C. "Multivariate statistical process control for autocorrelated processes", *International Journal of Production Research*, **34**, pp. 1715-1724 (1996).
5. Mason, R.L., Champ, C.W., Tracy, N.D. and et al. "Assessment of multivariate process control techniques", *Journal of Quality Technology*, **29**(2), pp. 140-143 (1997a).
6. Mason, R.L., Tracy, N.D. and Young, J.C. "A practical approach for interpreting multivariate T2 control chart signals", *Journal of Quality Technology*, **29**(4), pp. 396-406 (1997b).

7. Williams, J.D., Woodall, W.H., Birch, J.B. and et al. "Distribution of hotelling's T2 statistic based on the successive differences estimator", *Journal of Quality Technology*, **38**, pp. 217-229 (2006).
8. Fan, S.k., Huang, H.-K. and Chang, Y.-J. "Robust multivariate control chart for outlier detection using hierarchical cluster tree in SW2", *Quality and Reliability Engineering* 29(7), (2012). DOI: 10.1002/qre.1448.
9. Mason, R.L. and Young, J.C., *Multivariate Statistical Process Control With Industrial Applications American Statistical Association and the Society for Industrial and Applied Mathematics*, Philadelphia, PA (2002).
10. Ong, H.C. and Alih, E. "A control chart based on cluster-regression adjustment for retrospective monitoring of individual characteristics", *PLoS ONE*, **10**(4), e0125835 (2015). doi:10.1371/journal.pone.0125835.
11. Alfaro, J.L. and Ortega, J.F. "A comparison of robust alternatives to hotelling's T 2 control chart", *Journal of Applied Statistics*, **36**(12), pp. 1385-1396 (2009).
12. Chenouri, S.E., Steiner, S.H. and Variyath, A.M. "A multivariate robust control chart for individual observations", *Journal of Quality Technology*, **41**(3), pp. 259-271 (2009).
13. Rousseeuw, P.J. and Van Zomeren, B.C. "Unmasking multivariate outliers and leverage points", *Journal of the American Statistical Association*, **85**(411), pp. 633-639 (1990).
14. Ahsan, M., Mashuri, M., Kuswanto, H. and et al. "Outlier detection using PCA mix based T 2 control chart for continuous and categorical data", *Communications in Statistics-Simulation and Computation*, pp. 1-28 (2019).
15. Rocke, D.M. "And RQ charts: robust control charts", *Journal of the Royal Statistical Society: Series D (The Statistician)*, **41**(1), pp. 97-104 (1992).
16. Alfaro, J.L. and Ortega, J.F. "A robust alternative to Hotelling's T2 control chart using trimmed estimators", *Quality and Reliability Engineering International*, **24**(5), pp. 601-611 (2008).
17. Sullivan, J.H. and Woodall, W.H. "A Comparison of multivariate control charts for individual observations", *Journal of Quality Technology*, **28**(4), pp. 398-408 (1996).
18. Atkinson, A.C. and Mulira, H.M. "The stalactite plot for the detection of multivariate outliers", *Statistics and Computing*, **3**(1), pp. 27-35 (1993).
19. Rousseeuw, P.J. "Least median of squares regression", *Journal of the American Statistical Association*, **79**(388), pp. 871-880 (1984).
20. Kang, JI.H. and Bum Kim, S. "A clustering algorithm-based control chart for inhomogeneously distributed TFT-LCD processes", *International Journal of Production Research*, **51**(18), pp. 5645-5657 (2013).
21. Hajabbasi. "Multivariate quality control charts", Unpublished master's Thesis. Shahid Bahonar University of Kerman. Kerman (2013).
22. Vargas, J.A.N. "Robust estimation in multivariate control charts for individual observations", *Journal of Quality Technology*, **35**(4), pp. 367-376 (2003).
23. Rousseeuw, P.J. and Leroy, A.M., *Robust Regression and Outlier Detection*, **589**, John wiley & sons (2005).
24. Rousseeuw, P.J. and Driessen, K.V. "A fast algorithm for the minimum covariance determinant estimator", *Technometrics*, **41**(3), pp. 212-223 (1999).
25. Jensen, W.A., Birch, J.B. and Woodall, W.H. "High breakdown estimation methods for phase I multivariate control charts", *Quality and Reliability Engineering International*, **23**(5), pp. 615-629 (2007).
26. Willems, G., Pison, G., Rousseeuw, P.J. and et al. "A robust Hotelling test", *Metrika*, **55**(1-2), pp. 125-138 (2002).
27. Variyath, A.M. and Vattathoor, J. "Robust control charts for monitoring process mean of phase-I multivariate individual observations", *Journal of Quality and Reliability Engineering Volume 2013, Article ID 542305*, 14 pages (2013). doi.org/10.1155/2013/542305.
28. Ankara, H. and Yerel, S. "Determination of sampling errors in natural stone plates through single linkage cluster method", *Journal of Materials Processing Technology*, pp. 2483-2487 (2008). doi:10.1016/j.jmatprotec.2008.05.048.
29. Jobe, J.M. and Pokojovy, M. "A multistep, cluster-based multivariate chart for retrospective monitoring of individuals", *Journal of Quality Technology*, **41**(4), pp. 323-339 (2009).
30. Jobe, J.M. and Pokojovy, M. "A cluster-based outlier detection scheme for multivariate data", *Journal of the American Statistical Association*, **110**(512), pp. 1543-1551 (2015).
31. Huang, J., Zhu, Q., Yang, L. and et al. "A novel outlier cluster detection algorithm without top-n parameter", *Knowledge-Based Systems*, **121**, pp. 32-40 (2017).
32. <https://www.mathworks.com/help/stats/hierarchical-clustering.html>.
33. <https://rdrr.io/cran/robustbase/man/phosphor.html>.
34. Snedecor, G.W. and Cochran, W.G., *Statistical Methods*, 6th Edition. The Iowa State University Press: Ames, Iowa (1967).
35. Hawkins, D.M., Bradu, D. and Kass, G.V. "Location of several outliers in multiple-regression data using elemental sets", *Technometrics*, **26**, pp. 197-208 (1984).