

بهبود کیفیت مدل سازی و پیش بینی سری های زمانی با استفاده از تحلیل مجموعه مقادیر تکین چندمتغیره استوار

طاهره امینی

دانشجوی دکتری آمار دانشگاه پیام نور، تهران، ایران. tahereamini90@yahoo.com

مسعود یارمحمدی*

(نویسنده ی مسؤل)، دانشیار گروه آمار دانشگاه پیام نور، تهران، ایران. masyar@pnu.ac.ir

علی شادروخ

دانشیار گروه آمار دانشگاه پیام نور، تهران، ایران. a.shadrokh@pnu.ac.ir

مهدی کلانتری

استادیار گروه آمار دانشگاه پیام نور، تهران، ایران. kalantarimahdi@pnu.ac.ir

چکیده: در تحلیل سری های زمانی نادیده گرفتن نقاط دورافتاده منجر به شناسایی نادرست مدل، برآورد اریب پارامترها و در نتیجه پیش بینی های ضعیف و به عبارتی کاهش کیفیت و دقت مدل سازی می شود. یکی از روش های ناپارامتری معتبر در پیش بینی و بهبود کیفیت مدل سازی سری های زمانی چند متغیره، روش مجموعه مقادیر تکین چند متغیره^۱ (MSSA) است که نیازمند هیچ گونه فرض اولیه ای نیست. از آنجایی که وجود نقاط دورافتاده کارایی روش MSSA را کاهش داده و نرم ماتریسی فروبنیوس به کار رفته در آن را متأثر ساخته و به عبارتی غیر استوار می سازد، در این تحقیق، نسخه ی جدیدی از روش MSSA بر اساس نرم L_1 پیشنهاد می شود. در ادامه با استفاده از مطالعات شبیه سازی و نیز استفاده از داده های واقعی، عملکرد روش تحلیل مجموعه مقادیر تکین چندمتغیره بر اساس هر دو نرم مورد مقایسه قرار می گیرد. معیارهای مورد استفاده شامل ریشه ی میانگین توان دوم خطاها و میانگین قدرمطلق خطاها، برتری روش تحلیل مجموعه مقادیر تکین چندمتغیره بر اساس نرم L_1 را در بازسازی و پیش بینی سری زمانی نشان می دهند. **واژه های کلیدی:** تحلیل مجموعه مقادیر تکین چند متغیره، استوارسازی، نقاط دورافتاده، نرم L_1 ، نسبت میانگین توان دوم خطا، نسبت میانگین قدرمطلق خطا.

۱- مقدمه

و توانایی پیش بینی صحیح جهت تغییر^۲ (DC)، بستگی دارد [۱]. از این رو روش های متفاوتی برای پیش بینی سری های زمانی وجود دارد که با توجه به نوع و ماهیت داده ها می توان از آنها استفاده کرد. البته اغلب به منظور پیش بینی داده ها با استفاده از روش های کلاسیکی مانند ARIMA، فرض هایی خطی بودن مدل و نرمال بودن باقیمانده ها محدودیت هایی را برای رسیدن به نتایج مطلوب فراهم می کنند. یکی از روش های ناپارامتری مورد استفاده توسط محققین برای تحلیل و پیش بینی سری های زمانی به خصوص در زمینه های مهندسی و بهینه سازی مدل های خطی، روش تحلیل مجموعه مقادیر تکین^۳ (SSA) بوده که به دلیل عدم نیاز به برقراری فرض های محدودکننده و اولیه ای مانند مانایی داده ها، نرمال بودن باقیمانده ها و خطی بودن مدل، مورد توجه بسیاری از محققین در این حوزه قرار گرفته است. روش تحلیل مجموعه

به دلیل عدم اطمینان روزافزونی که با آن مواجه هستیم، تحلیل و پیش بینی سرهای زمانی در تمامی حوزه های صنعتی، کشاورزی، اقتصادی و... از اهمیت بسیار بالایی برخوردار است. علت این امر آن است که مدل های پیش بینی دقیق، توانایی ما را در تصمیم گیری، برنامه ریزی و مدیریت ریسک افزایش می دهند. رسیدن به چنین منظوری تا حد زیادی به ماهیت روش پیش بینی (از دیدگاه پارامتری یا ناپارامتری)، دقت پیش بینی ها

تاریخ دریافت: ۱۴۰۲/۰۵/۲۶ تاریخ پذیرش: ۱۴۰۲/۰۹/۰۷

دوره ۱۲ / شماره ۴

صفحات ۴۹۵-۵۲۰

*Corresponding author: masyar@pnu.ac.ir

¹ Multivariate Singular Spectrum Analysis (MSSA)

² Direction of Change (DC)

³ Singular Spectrum Analysis (SSA)

صنعت [۱۲]، [۸]، [۱۳] و [۱۴] و پیش‌بینی مرگ و میر [۱۵]، می‌توان نام برد. از طرف دیگر روش MSSA کارایی زیادی در بازسازی و بهینه‌سازی مدل‌های آماری و پیش‌بینی داده‌های مالی و اقتصادی، صنعتی و تجاری دارد و برتری این روش در برابر روش SSA و حتی روش‌های کلاسیک پارامتری و اقتصادسنجی در مطالعاتی نشان داده شده است [۱۳]، [۱۵]، [۱۶]، [۱۷] و [۱۸]. از این‌رو در این مقاله به موضوع تحلیل مجموعه مقادیر تکین چندمتغیره پرداخته خواهد شد.

تحلیل سری‌های زمانی همراه با نقاط دورافتاده یک بحث اساسی و مهم در علم آمار است. این داده‌ها اغلب به دلیل عواملی خارجی، مانند جنگ‌ها و بحران‌های سیاسی و اقتصادی، خرابی دستگاه‌ها، وقوع جریان‌های سیاسی، اعتصاب‌ها و... به وجود می‌آیند و باعث گسست‌های ساختاری در سری‌های زمانی نیز می‌شوند. وجود داده‌های دورافتاده بر تحلیل سری‌های زمانی و نتایج پیش‌بینی تأثیرات منفی بسیار زیادی گذاشته، می‌تواند باعث انحراف نتایج شده و در نتیجه باعث مشکلات عدیده‌ای شوند. اگر چه بر اساس مطالعات صورت‌گرفته، حساسیت مدل‌های ناپارامتری به داده‌های دورافتاده کمتر از مدل‌های پارامتری است اما به دلیل کارایی بیشتر مدل‌های ناپارامتری در سری‌های زمانی، بررسی و مطالعه در زمینه‌ی کاهش اثر منفی داده‌های دورافتاده بر این مدل‌ها می‌تواند قابل تأمل باشد [۱۹]. با توجه به این‌که ساختار روش مجموعه مقادیر تکین چندمتغیره بر اساس نرم فروبنیوس است، از این‌رو نسبت به داده‌های دورافتاده حساس بوده و کارایی مناسبی ندارد. به‌منظور کاهش اثرات نقاط دورافتاده بر روی این مدل، در این تحقیق نرم ماتریسی L_1 به جای نرم فروبنیوس جایگزین و نسخه‌ی جدیدی از MSSA معرفی شده که در ادامه‌ی متن با نماد $L_1 - MSSA$ نمایش داده می‌شود. در این راستا ابتدا روشی برای محاسبه‌ی ماتریس سیگنال بلوکی بر اساس نرم L_1 معرفی شده و سپس روش هنکل‌سازی با استفاده از این نرم برای حالت عمودی و افقی ماتریس مسیر به‌صورت یک قضیه بیان خواهد شد. البته لازم است تا ماتریس مسیر هر دو حالت افقی و عمودی به دقت مورد واکاوی قرار گرفته و تغییرات لازم در شیوه‌ی محاسبه‌ی ماتریس سیگنال برای ساخت ماتریس مسیر بلوک‌بندی هنکل‌سازی شده، ارائه شود. در ادامه با استفاده

مقادیر تکین برای اولین‌بار به وسیله برومهد و کینگ^۴ معرفی شد. بعد از آن گیل و همکاران^۵ و گولیاندینا و استیانف^۶ ویرایش‌های دیگری از روش‌های تحلیل مجموعه مقادیر تکین یک و چند متغیره را برای پیش‌بینی مورد بحث و بررسی قرار دادند [۲]، [۳] و [۴].

روش تحلیل مجموعه مقادیر تکین چندمتغیره (MSSA) و یک متغیره (SSA) دارای الگوریتمی دو مرحله‌ای و هر مرحله شامل دو گام است. در مرحله‌ی تجزیه، سری زمانی به‌صورت یک ماتریس خاص که به آن ماتریس مسیر^۷ گفته می‌شود در می‌آید. سپس در مرحله‌ی بازسازی مشاهدات سری زمانی در ساده‌ترین حالت تجزیه خود به دو مؤلفه‌ی نوفه و سیگنال تقسیم شده و متناسب با سری بازسازی‌شده از مؤلفه‌ی سیگنال، مقادیر پیش‌بینی با استفاده از یک رابطه‌ی خطی بازگشتی^۸ (LRR) محاسبه می‌شود.

در هر بار استفاده از روش‌های تحلیل مجموعه مقادیر تکین یک و چندمتغیره و قبل از تعیین ماتریس مسیر، کاربر ملزم به تعیین «پارامتر»، طول پنجره^۹ [۵] و پس از آن در گام دوم، تعیین پارامتر بازسازی یا تعداد سه‌تایی‌های ویژه^{۱۰} مورد استفاده برای بازسازی است [۶]. تعیین پارامتر بازسازی، متناسب با ساختار ماتریس مسیر سری زمانی بوده که در حالت چند متغیره بنا به تشخیص کاربر، ماتریس مسیر به دو حالت افقی و عمودی بلوک‌بندی می‌شود. برای توضیحات کامل‌تر به [۵] مراجعه شود. اگر چه هر دو روش SSA و MSSA رویکردی بهینه برای مدل‌سازی و پیش‌بینی ساختارهای متفاوت سری زمانی اعم از خطی یا غیر خطی، مانا یا نامانا و همچنین سری‌های زمانی با ساختارهایی ناشناخته و پنهان دارند [۷]، اما بر خلاف روش SSA، درباره‌ی جنبه‌های نظری MSSA تحقیقات کمتری انجام شده است، MSSA بدون داشتن پیچیدگی‌های خاصی، همانند SSA کاربردهای زیادی مانند نوفه‌زدایی، استخراج روند، علیت و پیش‌بینی دارد [۸] و [۹]. از طرفی در سال‌های اخیر به دلیل پیشرفت فناوری و قابلیت‌های ذخیره‌سازی ماشین‌های روزمره، توجه بیشتری به سمت سری‌های زمانی چندمتغیره جلب شده است. همچنین موارد بسیاری از کاربردهای MSSA را در زمینه‌هایی مانند تجارت [۱۰]، ترافیک خطوط هوایی [۱۱]،

⁸ Linear Recurrence Relations (LRR)

⁹ Window length

¹⁰ eigentriples

⁴ Broomhead, D. S. and King, G. P. (1986 b).

⁵ Ghil, M., R. M. Allen, M. D. Dettinger, K. Ide, D.

Kondrashov, et al. (2002)

⁶ Golyandina, N., and Stepanov, D. (2005)

⁷ Trajectory matrix

که وجود نقاط دور افتاده در یک سری زمانی می‌تواند مرتبه ماتریس مسیر، بزرگی مقادیر ویژه، بردارهای ویژه و در نتیجه تجزیه، بازسازی و پیش‌بینی مقادیر آینده سری زمانی را تحت تأثیر قرار دهند [۱۹]. لذا بر خلاف تمام مزیت‌هایی که برای روش MSSA بیان شد، این روش نسبت به نقاط دور افتاده استوار نبوده و یک نقطه ضعف برای آن محسوب می‌شود. حال با توجه به این که وجود نقاط دورافتاده در سری‌های زمانی بر تحلیل‌ها از جمله پیش‌بینی مقادیر آینده سری زمانی و برآورد پارامترهای مدل تأثیر منفی دارد، لزوم شناسایی و کاهش اثرات آنها از اهمیت بسیاری برخوردار است. از طرفی دیگر به دلیل کارا بودن روش MSSA نسبت به روش SSA و به منظور استوارسازی MSSA کلاسیک با استفاده از نرم L_1 به جای L_2 در روش MSSA کلاسیک، نسخه جدیدی از این روش ارائه می‌شود که استفاده از این نرم می‌تواند یکی از راه‌های استوارسازی و بهبود کارایی در این روش باشد.

همانند روش کلاسیک MSSA، روش $L_1 - MSSA$ نیز دارای دو مرحله‌ی مکمل بوده و هر مرحله هم شامل دو قدم اصلی است، گام‌های نشانیدن و گروه‌بندی در روش $L_1 - MSSA$ ، همانند روش کلاسیک MSSA است اما گام‌های تجزیه‌ی مقادیر تکین و هنکل‌سازی در روش $L_1 - MSSA$ تفاوتی اساسی با این گام‌ها در روش کلاسیک MSSA دارند که در ادامه شرح داده خواهند شد.

۲-۱- تجزیه

این مرحله شامل دو گام است: نشانیدن و تجزیه‌ی مقادیر

تکین

الف) نشانیدن

در روش SSA تعبیه کردن به‌عنوان یک نگاشت برای تبدیل یک سری زمانی مانند $Y_N = (y_1, y_2, \dots, y_N)$ به طول N به سری‌های زمانی X_1, X_2, \dots, X_k به صورت بردارهای $X_i = (y_1, y_2, \dots, y_{i+L-1})^T \in R^L$ در نظر گرفته می‌شود. در واقع در این گام سری زمانی Y_N به k زیرسری تبدیل می‌شود، وقتی که عدد صحیح L به‌عنوان طول پنجره تعریف شود ($2 \leq L \leq N$) و $K = N - L + 1$. دقت کنید بردار X_i یک زیرسری متشکل از یک بردار ستونی با L مؤلفه است و گاهی

از مطالعات شبیه‌سازی عملکرد روش تحلیل مجموعه مقادیر تکین چندمتغیره بر اساس هر دو نرم مورد مقایسه قرار می‌گیرد و با به کارگیری این روش برای سری‌های زمانی واقعی، معیارهای مورد استفاده شامل ریشه‌ی میانگین توان دوم خطاها و میانگین قدرمطلق خطاها، برتری روش تحلیل مجموعه مقادیر تکین چندمتغیره بر اساس نرم L_1 را نسبت به نسخه‌ی قبلی، در بازسازی و پیش‌بینی سری زمانی نشان می‌دهد. در ادامه در بخش دوم علاوه بر معرفی اجمالی MSSA، راهکارهای محاسباتی روش تحلیل مجموعه مقادیر تکین چند متغیره بر اساس نرم L_1 ارائه می‌شود پس از آن با استفاده از نتایج شبیه‌سازی و به کارگیری داده‌های واقعی، کارایی روش جدید مورد ارزیابی قرار داده می‌شود.

۲- روش‌شناسی

تحلیل مجموعه مقادیر تکین چند متغیره (MSSA) یکی از روش‌های ناپارامتری برای تحلیل سری‌های زمانی چندمتغیره بوده و مانند حالت یک متغیره (SSA) شامل دو مرحله تجزیه^{۱۱} و بازسازی^{۱۲} است، که هر مرحله خود به دو گام نشانیدن^{۱۳} و تجزیه‌ی مقادیر تکین (SVD)^{۱۴} در مرحله‌ی اول و گروه‌بندی^{۱۵} و میانگین‌گیری قطری^{۱۶} در مرحله‌ی دوم دسته‌بندی می‌شوند. علاقه‌مندان برای مطالعه‌ی جزئیات مدل MSSA می‌توانند به [۹] مراجعه کنند. در روش MSSA، تعریف ماتریس مسیر به دو حالت مختلف افقی و عمودی است. در مرحله‌ی تجزیه هر متغیر که یک سری زمانی می‌باشد، به طور جداگانه به چند مؤلفه‌ی قابل تفسیر همچون روند، مؤلفه‌ی نوسانی و نوفه تجزیه می‌شود و بعد از گروه‌بندی مناسب، مؤلفه‌های بازسازی‌شده مربوط به هر سری در یک ماتریس معین، درست به همان ترتیب و چیدمان قبل از تجزیه در کنار هم قرار می‌گیرند. به این ترتیب مؤلفه‌های بازسازی شده برای پیش‌بینی مشاهدات جدید مربوط به هر سری مورد استفاده قرار می‌گیرند. روش MSSA همانند SSA به انتخاب دو پارامتر اساسی طول پنجره و پارامتر ساختاری مرتبط با نحوه گروه‌بندی بستگی دارد و ساختار این روش بر اساس نرم L_2 بوده که نسبت به نقاط دورافتاده استوار نیست. حسنی و همکاران (۲۰۱۴)، نشان دادند

¹⁴ Singular Value Decomposition (SVD)

¹⁵ Grouping

¹⁶ Diagonal Averaging

¹¹ Decomposition

¹² Reconstruction

¹³ Embedding

نمایش داده می شود در رابطه ی (۴) نحوه ی شکل گیری ماتریس مسیر عمودی نمایش داده شده است:

$$X_V = \begin{bmatrix} X^{(1)} \\ \vdots \\ X^{(p)} \end{bmatrix} = \begin{bmatrix} y_1^1 & y_2^1 & y_3^1 & \dots & y_K^1 \\ y_2^1 & y_3^1 & y_4^1 & \dots & y_{K+1}^1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ y_L^1 & y_{L+1}^1 & y_{L+2}^1 & \dots & y_N^1 \\ y_1^2 & y_2^2 & y_3^2 & \dots & y_K^2 \\ y_2^2 & y_3^2 & y_4^2 & \dots & y_{K+1}^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ y_L^2 & y_{L+1}^2 & y_{L+2}^2 & \dots & y_N^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ y_1^p & y_2^p & y_3^p & \dots & y_K^p \\ y_2^p & y_3^p & y_4^p & \dots & y_{K+1}^p \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ y_L^p & y_{L+1}^p & y_{L+2}^p & \dots & y_N^p \end{bmatrix} \quad (4)$$

دقت به این نکته حائز اهمیت است که عناصر روی قطرهای فرعی هر ماتریس $X^{(i)}$ برای هر بلوک با هم برابرند. چنین ماتریسی را ماتریس مسیر هنکل می نامند. از طرفی واضح است که با در اختیار داشتن ماتریس $X^{(i)}$ می توان سری زمانی $Y_N^{(i)}$ را به دست آورد. در واقع با کنار هم قرار دادن ستون اول و سطر آخر (یا سطر اول و ستون آخر) در هر بلوک، ماتریس $X^{(i)}$ ، سری زمانی $Y_N^{(i)}$ حاصل می شود. همچنین یک تناظر یک به یک بین ماتریس مسیر (که هنکل نیز می باشد) و سری زمانی وجود دارد. در جدول ۱ تفاوت بین تعیین پارامترهای بازسازی و طول پنجره متناسب با تعداد داده های سری زمانی و عمودی یا افقی^{۱۹} بودن ماتریس مسیر بیان می شود.

جدول ۱- نحوه محاسبات پارامترها برای دو حالت افقی و عمودی مجموعه مقادیر تکین چند متغیره

تعداد مقادیر ویژه	K	L	طول سری ها	روش
$L_{sum} = \sum_{i=1}^p L_i$	یکسان	مختلف	مختلف	$L_1 - VMSSA$
L	مختلف	یکسان	مختلف	$L_1 - HMSSA$

ب) تجزیه ی مقادیر تکین (SVD) ماتریس مسیر در حالت افقی

دومین گام از مرحله ی اول ساختن مدل $L_1 - MSSA$ تجزیه ی ماتریس مسیر بر اساس تجزیه ی مقادیر تکین است. با استفاده از روش SVD، ماتریس مسیر به صورت مجموع ماتریس های

^{۱۹} برای نمایش حالت افقی یا عمودی روش تحلیل مجموعه مقادیر تکین چندمتغیره به ترتیب از نمادهای VMSSA و HMSSA استفاده می شود.

اوقات بردار L -تاخیر^{۱۷} نامیده می شود. نتیجه ی مرحله ی نشاندن، ماتریس $X = [X_1, X_2, \dots, X_K]$ است که تشکیل یک ماتریس $L \times K$ را می دهد و به صورت رابطه ی (۱) نمایش داده می شود.

$$X = (x_{ij})_{i,j=1}^{L,K} = \begin{pmatrix} y_1 & y_2 & y_3 & \dots & y_K \\ y_2 & y_3 & y_4 & \dots & y_{K+1} \\ y_3 & y_4 & y_5 & \dots & y_{K+2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ y_L & y_{L+1} & y_{L+2} & \dots & y_{L+K-1} \end{pmatrix} \quad (1)$$

همچنین: $X_{ij} = y_{i+j-1}$. ماتریس مسیر ساخته شده (X) در این مرحله را ماتریس هنکل^{۱۸} نیز می نامند، زیرا عناصر روی قطرهای فرعی این ماتریس با هم برابرند [۲۰]. یکی از تفاوت های مراحل ساخت SSA با MSSA و $L_1 - MSSA$ در ساختن ماتریس مسیر است. روش های ساخت ماتریس مسیر در مدل $L_1 - MSSA$ به دو حالت افقی و عمودی است که در ادامه به صورت زیر بیان خواهد شد.

فرض کنید $y_j = \{y_j^{(1)}, y_j^{(2)}, \dots, y_j^{(p)}\}$ مشاهده ی j ام یک نمونه ی به طول N از یک سری زمانی p -متغیره بوده $X^{(i)}$ ($i = 1, \dots, p$)، ماتریس مسیر متناظر با سری زمانی i ام مسیر چندمتغیره را می توان با کنار هم قراردادن ماتریس های مسیر هر سری زمانی (به صورت افقی) یا زیر هم قراردادن آنها (به صورت عمودی) به دست آورد. در مرحله ی تعبیه کردن ماتریس ها به ترتیب در حالت افقی و عمودی (از چپ به راست) به صورت رابطه (۲) در خواهند آمد:

$$X_V = \begin{bmatrix} X^{(1)} \\ \vdots \\ X^{(p)} \end{bmatrix} \text{ یا } X_H = [X^{(1)} : X^{(2)} : \dots : X^{(p)}] \quad (2)$$

چنانچه همه ی p - تا ماتریس های مسیر به صورت افقی کنار هم قرار داده شوند ماتریس مسیر به حالت افقی در می آید که با X_H نامیده می شود. بنا بر این در رابطه (۳) نحوه ی شکل گیری ماتریس مسیر افقی نمایش داده شده است.

$$X_H = [X^{(1)} : X^{(2)} : \dots : X^{(p)}] = \begin{bmatrix} y_1^1 & y_2^1 & y_3^1 & \dots & y_K^1 & y_1^2 & y_2^2 & y_3^2 & \dots & y_K^2 & \dots & y_1^p & y_2^p & y_3^p & \dots & y_K^p \\ y_2^1 & y_3^1 & y_4^1 & \dots & y_{K+1}^1 & y_2^2 & y_3^2 & y_4^2 & \dots & y_{K+1}^2 & \dots & y_2^p & y_3^p & y_4^p & \dots & y_{K+1}^p \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \dots & \vdots & \vdots & \vdots & \ddots & \vdots \\ y_L^1 & y_{L+1}^1 & y_{L+2}^1 & \dots & y_N^1 & y_L^2 & y_{L+1}^2 & y_{L+2}^2 & \dots & y_N^2 & \dots & y_L^p & y_{L+1}^p & y_{L+2}^p & \dots & y_N^p \end{bmatrix} \quad (3)$$

و اگر همه ی p - تا ماتریس های مسیر به صورت عمودی زیر هم قرار داده شوند ماتریس مسیر حالت عمودی حاصل شده و با X_V

^{۱۷} L-Lagged
^{۱۸} Hankel Matrix

$$Q = X_V X_V^T$$

$$Q = \begin{bmatrix} X_V^{(1)} X_V^{(1)T} & X_V^{(1)} X_V^{(2)T} & \dots & X_V^{(1)} X_V^{(p)T} \\ X_V^{(2)} X_V^{(1)T} & X_V^{(2)} X_V^{(2)T} & \dots & X_V^{(2)} X_V^{(p)T} \\ \vdots & \vdots & \ddots & \vdots \\ X_V^{(p)} X_V^{(1)T} & X_V^{(p)} X_V^{(2)T} & \dots & X_V^{(p)} X_V^{(p)T} \end{bmatrix}$$

لازم به توضیح است که SVD ماتریس مسیبر X_V را می‌توان به صورت زیر نوشت.

$$X_V = U_V \Sigma_V V_V^T$$

$$U_V = [U_{V_1}, U_{V_2}, \dots, U_{V_{L_{sum}}}] \in R^{L_{sum} \times L_{sum}}$$

$$V_V = [V_{V_1}, V_{V_2}, \dots, V_{V_k}] \in R^{K \times K}$$

$$\Sigma = \text{diag}(\sqrt{\lambda_{V_1}} \geq \sqrt{\lambda_{V_2}} \geq \dots \geq \sqrt{\lambda_{V_{L_{sum}}}})$$

$L_{sum} = \sum_{i=1}^p L_i$ مقادیر ویژه $\lambda_{V_1}, \lambda_{V_2}, \dots, \lambda_{V_{L_{sum}}}$ از ماتریس Q به دست می‌آید که $\lambda_{V_1} \geq \lambda_{V_2} \geq \dots \geq \lambda_{V_{L_{sum}}}$ و همچنین $U_{V_1}, U_{V_2}, \dots, U_{V_{L_{sum}}}$ بردارهای ویژه‌ی مرتبط با این مقادیر ویژه هستند.

در روش $L_1 - \text{MSSA}$ تجزیه مقادیر تکین هر کدام از ماتریس مسیبرها به ترتیب در حالت عمودی و افقی به صورت زیر به دست می‌آید.

$$X_V = X_{V_1} + X_{V_2} + \dots + X_{V_{L_{sum}}} = \sum_{i=1}^{d_V} \sqrt{\lambda_{V_i}} U_{V_i} W_i V_{V_i}^T$$

$$X_H = X_{H_1} + X_{H_2} + \dots + X_{H_d} = \sum_{i=1}^{d_H} \sqrt{\lambda_{H_i}} U_{H_i} W_i V_{H_i}^T \quad (V)$$

وقتی که $X_{H_i} = \sqrt{\lambda_{H_i}} U_{H_i} W_i V_{H_i}^T$ و $X_{V_i} = \sqrt{\lambda_{V_i}} U_{V_i} W_i V_{V_i}^T$ هستند. همچنین d_H و d_V به ترتیب برابر با تعداد مقادیر ویژه‌ی مثبت یا رتبه‌ی ماتریس مسیبر در حالت عمودی و افقی است.

بنا بر خاصیت ساختاری روش مجموعه مقادیر تکین چندمتغیره، در حالت افقی یا حالت عمودی ماتریس سیگنال $Z = U \Sigma_r V^T$ در بین تمام ماتریس‌های سیگنال با رتبه‌ی $r < d$ کمترین فاصله را از ماتریس مسیبر دارد یعنی $\|X - U \Sigma_r V^T\|_F$ حداقل مقدار را دارد. لذا از این ویژگی مهم می‌توان در محاسبه‌ی ماتریس سیگنال بر پایه‌ی نرم L_1 استفاده کرد. برای جلوگیری از تکرار مطالب و نیز به دلیل تغییرات اساسی ابعاد و ساختار ماتریس مسیبر در حالت عمودی، در این قسمت تنها به بررسی یک حالت از ماتریس سیگنال بر اساس نرم L_1 پرداخته شده و ویژگی و ساختار ماتریس وزن در ساخت نزدیکترین ماتریس سیگنال به ماتریس مسیبر حالت عمودی تشریح خواهد شد. حالت افقی شبیه

مقدماتی متعامد (هر کدام با رتبه‌ی ۱) تجزیه می‌شوند. در نسخه‌ی $L_1 - \text{MSSA}$ همانند روش کلاسیک MSSA با توجه ساختار متفاوت ماتریس مسیبر برای هر دو حالت افقی یا عمودی، ابعاد ماتریس XX^T به دلیل نوع قرار گرفتن ماتریس مسیبر هر سری زمانی در کنار هم تحت تأثیر قرار می‌گیرد و لذا تفاوتی در تعداد مقادیر ویژه و بردارهای ویژه ظاهر می‌شود. اما شیوه‌ی محاسبه ماتریس‌های مقدماتی متعامد با حالت کلاسیک MSSA متفاوت است. بدین صورت که برای استوارسازی، با معرفی ماتریس قطری W به عنوان ماتریس وزن و ضریب قرار دادن این ماتریس در محاسبه‌ی ماتریس‌های مقدماتی متعامد، تجزیه‌ی مقادیر تکین به صورتی متفاوت خواهد بود. لذا در ابتدا علاوه بر توضیح کامل روش تجزیه‌ی مقادیر تکین برای روش تحلیل مجموعه مقادیر تکین کلاسیک در هر دو حالت افقی و عمودی، ماتریس قطری وزن W معرفی شده و پس از آن این مرحله برای روش $L_1 - \text{MSSA}$ بیان خواهد شد.

برای حالت HMSSA و به منظور اجرای روش SVD بر روی ماتریس مسیبر X_H ، سه تایی‌های ویژه $(\lambda_{H_i}, U_{H_i}, V_{H_i})$ از ماتریس $X_H X_H^T$ به دست می‌آیند. فرض کنید $\lambda_{H_1}, \dots, \lambda_{H_L}$ مقادیر ویژه‌ی ماتریس $X_H X_H^T$ بوده که به صورت نزولی مرتب شده‌اند $(\lambda_{H_1} \geq \dots \geq \lambda_{H_L} \geq 0)$ و نیز بردارهای ویژه یک‌متعامد مقادیر ویژه‌ی $\lambda_{H_1}, \dots, \lambda_{H_L}$ باشند.

$$X_H X_H^T = [X^{(1)} : X^{(2)} : \dots : X^{(p)}] \begin{bmatrix} X^{(1)} \\ X^{(2)} \\ \dots \\ X^{(p)} \end{bmatrix} \quad (5)$$

$$X_H X_H^T = X^{(1)} X^{(1)T} + X^{(2)} X^{(2)T} + \dots + X^{(p)} X^{(p)T}$$

در این صورت SVD ماتریس X_H به صورت زیر تعریف می‌شود:

$$X_H = X_{H_1} + X_{H_2} + \dots + X_{H_L} \quad (6)$$

وقتی که $X_{H_i} = \sqrt{\lambda_{H_i}} U_{H_i} V_{H_i}^T$ به عنوان ماتریس‌های مقدماتی متعامد $V_{H_i} = X_H^T U_{H_i} / \sqrt{\lambda_{H_i}}$ به ازای $i = 1, \dots, L$ ، سه تایی $(\sqrt{\lambda_{H_i}}, U_{H_i}, V_{H_i})$ را سه تایی ویژه و $\sqrt{\lambda_{H_i}}$ را مقادیر تکین می‌نامند.

برای حالت VMSSA و به منظور اجرای SVD بر روی ماتریس مسیبر X_V ، سه تایی ویژه $(\lambda_{V_i}, U_{V_i}, V_{V_i})$ را از ماتریس $X_V X_V^T$ به دست می‌آید.

$$X_V X_V^T = \begin{bmatrix} X^{(1)} \\ X^{(2)} \\ \dots \\ X^{(p)} \end{bmatrix} [X^{(1)} : X^{(2)} : \dots : X^{(p)}]$$

به حالت تک‌متغیره می‌باشد که علاقمندان می‌توانند برای اطلاعات بیشتر به [۲۱]، مراجعه نمایند.

فرض کنید ماتریس مسییر حالت عمودی $X_V = X_{V_1} + X_{V_2} + \dots + X_{V_{L_{sum}}}$ دارای رتبه‌ی قطری $d_V = r = rank(X) < L < K$ باشد. در این قسمت ماتریس W به‌عنوان ماتریس وزن معرفی می‌شود.

$$\min_{A_r} \|X_V - A_r B_1\|_{L_1} \quad (۱۱)$$

برای حل (۱۱) فرض کنید $S = A_r B_1$ باشد آن‌گاه $S^T = B_1^T A_r^T$

$$S_j^T = B_1^T A_j^T \quad (۱۲)$$

که در آن A_j^T و S_j^T به ترتیب زامین ستون ماتریس‌های A_r^T و S^T هستند. رابطه‌ی (۱۳)، با در نظر گرفتن خاصیت تعریف L_1 برای یک

$$\begin{aligned} \|X_V - S\|_{L_1} &= \|(X_V - S)^T\|_{L_1} \\ &= \|X_V^T - S^T\|_{L_1} \\ &= \sum_{j=1}^{L_{sum}} \|X_{V_j}^T - S_j^T\|_{L_1} \\ &= \sum_{j=1}^{L_{sum}} \|X_{V_j}^T - B_1^T A_j^T\|_{L_1} \quad (۱۳) \end{aligned}$$

که $X_{V_j}^T$ به‌عنوان زامین ستون ماتریس X_V^T در نظر گرفته شده است. رابطه (۱۳) نشان می‌دهد برای کمینه کردن عبارت $\|X_V - S\|_{L_1}$ کافی است عبارت $\|X_{V_j}^T - B_1^T A_j^T\|_{L_1}$ نسبت به A_j^T کمینه شود.

با توجه به این‌که مسئله‌ی فوق یک مسئله‌ی کمینه‌سازی شناخته‌شده در رگرسیون کمترین قدرمطلق انحرافات^{۲۰} (LAD) بوده و حل آن به‌صورت عددی وجود دارد؛ لذا با فرض این‌که $X_{V_j}^T$ ، B_1^T و A_j^T را به ترتیب به‌عنوان بردار متغیر وابسته، ماتریس مشاهدات متغیرهای مستقل و بردار ضریب مجهول در نظر گرفته شوند، بردار A_j^T را می‌توان به کمک روش‌های عددی تکراری، مشابه آن‌چه که در [۲۲] و [۲۳] آمده است محاسبه کرد. بنا بر این با داشتن بردارهای A_j^T ، $j = 1, 2, \dots, L_{sum}$ ، ماتریس A_r به‌دست می‌آید. برای انجام محاسبات مربوط به رگرسیون LAD از تابع `rq` در بسته‌ی `quantreg` نرم‌افزار R با آرگومان $\tau = 0.5$ استفاده شده است.

۲-۲- بازسازی

این مرحله شامل دو گام است: گروه‌بندی و میانگین‌گیری قطری.

الف) گروه‌بندی

گروه‌بندی یکی از مراحل مهم در تجزیه و بازسازی سری‌های زمانی به روش‌های SSA و MSSA است. در حال حاضر

فرض کنید ماتریس مسییر حالت عمودی $X_V = X_{V_1} + X_{V_2} + \dots + X_{V_{L_{sum}}}$ دارای رتبه‌ی قطری $d_V = r = rank(X) < L < K$ باشد. در این قسمت ماتریس W به‌عنوان ماتریس وزن معرفی می‌شود.

$$W = \text{diag}(\overbrace{w_1, w_2, \dots, w_r}^r, \overbrace{0, 0, \dots, 0}^{L_{sum}-r}) \in R^{L_{sum} \times L_{sum}} \quad (۸)$$

توجه داشته باشید که

$$w = \begin{pmatrix} w_1 & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & w_2 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & w_r & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \end{pmatrix} = \begin{pmatrix} W_r & 0 \\ 0 & 0 \end{pmatrix},$$

که در آن $W_r = \text{diag}(w_1, w_2, \dots, w_r) \in R^{r \times r}$

برای محاسبه‌ی ماتریس سیگنال بر پایه نرْم L_1 ، لازم است ماتریس وزن قطری $W_{L_{sum} \times L_{sum}}$ را به‌گونه‌ای یافت که $\|X_V - U_V W \Sigma_d V_V^T\|_{L_1}$ کمترین مقدار را داشته باشد. به طور دقیق‌تر، مسأله کمینه‌سازی زیر باید حل شود:

$$\min_w \|X_V - U_V W \Sigma_d V_V^T\|_{L_1}$$

توجه داشته باشید که

$$\begin{aligned} U_V W \Sigma_d V_V^T &= [U_V, U_{V_1}] \begin{bmatrix} W_r & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_{V_1}^T \\ V_V^T \end{bmatrix} \\ &= [U_V, W_r, 0] \begin{bmatrix} \Sigma_1 V_{V_1}^T \\ 0 \end{bmatrix} = U_V W_r \Sigma_1 V_{V_1}^T \quad (۹) \end{aligned}$$

که در آن $\Sigma_1 \in R^{r \times r}$ و $V_{V_1} \in R^{k \times r}$ ، $U_V \in R^{L_{sum} \times r}$ آسانی می‌توان نشان داد که

$$U_{V_1} W_r \Sigma_1 V_{V_1}^T = \sum_{i=1}^r \sqrt{\lambda_{V_i}} w_i U_{V_i} V_{V_i}^T \quad (۱۰)$$

با استفاده از رابطه (۸) و (۹) مسئله کمینه‌سازی زیر حل می‌شود:

$$\min_w \|X_V - U_{V_1} W_r \Sigma_1 V_{V_1}^T\|_{L_1}$$

به همین منظور، فرض کنید $A_r = U_{V_1} W_r$ و $B_1 = \Sigma_1 V_{V_1}^T$ که در آن $B_1 \in R^{k \times L_{sum}}$ بر اساس مشاهدات ماتریس‌های B_1 و U_{V_1} معلوم هستند. هدف یافتن ماتریس A_r است به‌طوری‌که

²⁰ Least Absolute Deviation (LAD)

حالت افقی است با این تفاوت که ماتریس مسیر ترکیبی از سری‌های زمانی است که به صورت عمودی کنار هم قرار گرفته‌اند.

ب) میانگین‌گیری قطری

در روش MSSA، ماتریس مسیر ترکیبی از ماتریس‌هایی است که به صورت بلوک کنار هم قرار گرفته‌اند. از این‌رو در این مرحله میانگین‌گیری قطری برای هر بلوک به صورت جداگانه محاسبه می‌شود. در بخش میانگین‌گیری قطری ماتریس‌های بازسازی‌شده در مرحله‌ی گروه‌بندی خاصیت هنکلی ندارند و چنانچه ماتریس مسیر هنکل نباشد نمی‌توان به سری متناظر با آن ماتریس دست یافت. از این‌رو هدف اصلی در این مرحله تبدیل ماتریس بازسازی‌شده به یک ماتریس مسیر هنکل است. در واقع برای رسیدن به این هدف به جای عناصر روی قطر فرعی میانگین آنها قرار داده می‌شود. در آخر به ازای تعداد سری‌های زمانی‌ای (تعداد بلوک‌ها) که در مسئله وجود دارد، به همان تعداد هم ماتریس هنکل‌سازی شده در ماتریس مسیر وجود خواهد داشت. از این‌رو در حالت چندمتغیره ماتریس سیگنال، انباشتی از ماتریس‌های هنکل‌سازی شده برای هر سری زمانی (هر بلوک) به تنهایی است. با این فرض که p تعداد بلوک‌ها یا تعداد متغیرهای سری زمانی باشد، در حالت کلی می‌توان ماتریس مسیر افقی (X_H) و ماتریس مسیر عمودی (X_V) به ترتیب با نمادهای $J_{MSSA}(X_H)$ و $J_{MSSA}(X_V)$ نمایش داده و آن‌ها را به صورت رابطه‌ی (۱۴) نوشت.

$$J_{MSSA}(X_H) = X_H = [X^{(1)} : \dots : X^{(p)}]$$

$$J_{MSSA}(X_V) = X_V = \begin{bmatrix} X^{(1)} \\ X^{(2)} \\ \vdots \\ X^{(p)} \end{bmatrix} \quad (14)$$

بنا بر این می‌توان ماتریس مسیر هنکل‌سازی شده را به صورت زیر در نظر گرفت، به طوری که برای حالت افقی

$$\begin{bmatrix} X^{(1)} \\ X^{(2)} \\ \vdots \\ X^{(p)} \end{bmatrix} \quad \text{و} \quad [X^{(1)} : \dots : X^{(p)}]$$

به ترتیب بلوک‌های تشکیل‌دهنده ماتریس مسیر بوده و

$$J_{MSSA}^{-1}(X_H) = [J_{SSA}^{-1}(X^{(1)}) : \dots : J_{SSA}^{-1}(X^{(p)})]$$

$$J_{MSSA}^{-1}(X_V) = \begin{bmatrix} J_{SSA}^{-1}(X^{(1)}) \\ J_{SSA}^{-1}(X^{(2)}) \\ \vdots \\ J_{SSA}^{-1}(X^{(p)}) \end{bmatrix} \quad (15)$$

متداول‌ترین روش‌ها برای گروه‌بندی رسم نمودار مقادیر تکین ماتریس مسیر و استفاده از مفهوم تفکیک‌پذیری و محاسبه‌ی همبستگی بین زیرسری‌هاست. در خصوص نمودار مقادیر تکین، چنانچه λ_i ها مقادیر نزدیک به هم داشته باشند در یک گروه قرار می‌گیرند.

گروه‌بندی برای ماتریس مسیر افقی

مرحله‌ی گروه‌بندی مربوط به افزایش ماتریس‌های $X_{H_1}, X_{H_2}, \dots, X_{H_L}$ به چندین گروه جدا از هم و جمع‌بندی ماتریس‌های درون هر گروه است. در مرحله‌ی گروه‌بندی با توجه به ماتریس مسیر حاصل‌شده، مجموعه‌ی $\{1, \dots, L\}$ به m زیرمجموعه‌ی $\{I_1, I_2, \dots, I_m\}$ جدا از هم افزایش می‌شوند در این صورت رابطه‌ی $X_H = X_{H_1} + X_{H_2} + \dots + X_{H_L}$ را می‌توان در گام گروه‌بندی به صورت $X_H = X_{I_1} + X_{I_2} + \dots + X_{I_m}$ نوشت. شیوه‌ی انتخاب مجموعه‌های I_1, I_2, \dots, I_m گروه‌بندی نامیده می‌شود. در واقع در این گام هدف یافتن سه‌تایی‌های ویژه مربوط به هر یک از اجزای سری زمانی نظیر روند، مؤلفه‌های فصلی، نوفه و غیره است. ساده‌ترین حالت گروه‌بندی برای زمانی است که تنها دو گروه وجود داشته باشند که این دو گروه می‌تواند شامل سیگنال و باقیمانده‌ها باشد $I_1 = \{1, 2, \dots, r\}$ در واقع r تا از بزرگ‌ترین مقادیر تکین و بردارهای ویژه متناظر برای تقریب سری اصلی انتخاب شده و $I_2 = \{r+1, r+2, \dots, L\}$ بقیه مقادیر تکین به‌عنوان مؤلفه‌ی نوفه در نظر گرفته می‌شود. در متون علمی موجود از r به‌عنوان پارامتر بازسازی نیز اسم برده شده است. بنا بر این چنانچه در متن مقاله از پارامتر بازسازی نام برده شود منظور همان r (بزرگ‌ترین مقدارهای تکین و بردارهای ویژه متناظر برای تقریب سری اصلی انتخاب شده) است.

گروه‌بندی برای ماتریس مسیر عمودی

در مرحله‌ی گروه‌بندی با توجه به ماتریس مسیر حاصل‌شده، مجموعه‌ی $\{1, 2, \dots, L_{sum}\}$ به m زیرمجموعه‌ی $\{I_1, I_2, \dots, I_m\}$ جدا از هم افزایش می‌شود. در این صورت رابطه‌ی $X_V = X_{V_1} + X_{V_2} + \dots + X_{V_{L_{sum}}}$ را می‌توان در گام گروه‌بندی به صورت $X_V = X_{I_1} + X_{I_2} + \dots + X_{I_m}$ در نظر گرفت. مانند حالت افقی شیوه‌ی انتخاب مجموعه‌های I_1, \dots, I_m گروه‌بندی نامیده می‌شود. سایر توضیحات این قسمت همانند

نشریه مهندسی و مدیریت کیفیت

مسیر بازسازی‌شده هنکل‌سازی شوند. همچنین چون اختلاف با هر بلوک با ماتریس مسیر هنکل‌سازی شده کمینه مقدار است بنا بر این ماتریس مسیر تولید شده در مرحله آخر نزدیک‌ترین داده‌ها را به سری زمانی تجزیه و بازسازی شده تولید خواهد کرد و داده‌های هر سری زمانی قابل دستیابی است. قضیه‌ی فوق بیان می‌کند که برای هنکل‌سازی ماتریس A بر اساس نرم L_1 کافی است به جای درایه‌های قطر فرعی در هر بلوک مقدار میانه‌ی درایه‌های آن قطر را جایگزین درایه‌های موجود در آن قطر کرد.

۲-۳- پیش‌بینی بازگشتی برای سری‌های زمانی به روش $L_1 - MSA$

برای پیش‌بینی سری‌های زمانی به روش تحلیل مجموعه مقادیر تکین چندمتغیره دو شیوه‌ی پیش‌بینی به روش‌های بازگشتی و برداری معرفی شده است که می‌توان برای مطالعه‌ی بیشتر به [۲۰] و [۲۴] مراجعه نمود. در این مقاله و در بخش‌های شبیه‌سازی برای $HMSSA$ و $VMSSA$ پیش‌بینی به روش بازگشتی خطی^{۲۱} (LRF) به کار برده خواهد شد و ضرایب رابطه‌ی خطی بازگشتی از بردارهای ویژه‌ی ماتریس مسیره‌ها به دست می‌آیند.

سطر آخر r ستون اول، از ماتریس بردارهای ویژه که در مرحله SVD روی ماتریس $X_H X_H^T$ به دست آمده بودند را بر $L - 1$ سطر اول رگرسیون کرده و بردار ضرایب را به شکل $\mathcal{R} = [\widehat{a}_1, \widehat{a}_2, \dots, \widehat{a}_{L-1}]$ تعریف می‌شوند. با فرض آن که $U_{H_j}^V$ نشان‌دهنده‌ی $L - 1$ مؤلفه‌ی اول بردار U_{H_j} و نماد Π_{H_j} آخرین مؤلفه بردار ویژه‌ی U_{H_j} ($j = 1, 2, \dots, r$) باشد، می‌توان بردار ضرایب خطی \mathcal{R} را به صورت رابطه‌ی (۱۶) تعریف کرد که در آن

$$\mathcal{R} = [\widehat{a}_1, \widehat{a}_2, \dots, \widehat{a}_{L-1}] = \frac{1}{1 - \nu^2} \sum_{j=1}^r \Pi_{H_j} U_{H_j}^V \quad (16)$$

اگر $\nu^2 < 1$ ، آن‌گاه h گام جلوتر (پیش‌بینی h گام جلوتر) از سری زمانی چندمتغیره به روش $HMSSA$ با استفاده از رابطه (۱۷) به دست می‌آید.

$$[\widehat{y}_j^{(1)}, \widehat{y}_j^{(2)}, \dots, \widehat{y}_j^{(p)}]^T = \begin{cases} [\widehat{y}_j^{(1)}, \widehat{y}_j^{(2)}, \dots, \widehat{y}_j^{(p)}] & j = 1, 2, \dots, N_i \\ \mathcal{R}^T \mathbf{Z}_h & j = N_i + 1, \dots, N_i + h \end{cases} \quad (17)$$

در حالی که؛

نماد \mathcal{T}^{-1} هنکل شدن یک ماتریس مسیر است.

در ابتدای این بحث، با توضیح مرحله‌ی میانگین‌گیری قطری که در روش $MSSA$ کلاسیک به آن اشاره شد به جای عناصر روی قطر فرعی، میانگین آنها قرار داده می‌شود. از طرفی دیگر، یکی از روش‌های هنکل‌سازی استفاده از ملاک کمترین فاصله بر اساس نرم فروبنیوس است. بر اساس این ملاک، ماتریس هنکل‌سازی شده‌ی A که با نماد $\mathcal{H}A$ نمایش داده می‌شود طوری تعیین می‌شود که $\|A - \mathcal{H}A\|_F$ کمترین مقدار را داشته باشد، اما در روش $MSSA - L_1$ هنکل‌سازی بر اساس نرم L_1 ارائه می‌شود. همانگونه که در حالت یک متغیره ثابت شده است [۲۱]، در حالت چندمتغیره نیز با استفاده از قضیه زیر می‌توان ثابت کرد که برای هنکل‌سازی ماتریس A بر اساس نرم L_1 کافی است به جای درایه‌های قطر فرعی در هر بلوک، مقدار میانه‌ی درایه‌های آن قطر را جایگزین درایه‌های موجود در آن قطر کرد.

قضیه:

فرض کنید A_I یک ماتریس $L_I \times K$ (بلوک l ام از ماتریس X_V) بوده $S_I = (i_1, \dots, p)$ ، $2 \leq S_I \leq L_I \times K$ و $S_I = i_l + j_l$ اگر درایه‌های $\tilde{a}_{i_l j_l}$ ماتریس $\mathcal{H}A_I$ به صورت عبارت زیر باشد،

$$\tilde{a}_{i_l j_l} = \underset{(L, k) \in A_S}{\text{median}} a_{l k}$$

آن‌گاه وقتی که

$$(A_S)_I = \{(L_i, K): L_i \times K = S_I, 1 \leq l_i \leq L_i, 1 \leq k_i \leq K_i\}$$

$\|A_I - \mathcal{H}A_I\|_{L_1}$ کمترین مقدار را خواهد داشت.

طبق تعریف درایه‌های $(b_{ij})_I$ ماتریس هنکل $B = \mathcal{H}A_I$ باید به ازای $S_I = i_l + j_l$ و برخی مقادیر $(g_S)_I$ در شرط $(b_{ij})_I = (g_S)_I$ صدق کند، از طرفی با در نظر گرفتن تعریف نرم L_1 یک ماتریس داریم $(I = 1, 2, \dots, p)$ ؛

$$\|A_I - \mathcal{H}A_I\|_{L_1} = \|A_I - B\|_{L_1}$$

$$\sum_{i_l=1}^{L_i} \sum_{j_l=1}^{K_i} |(a_{ij})_I - (b_{ij})_I| = \sum_{S_I=1}^{L_i+K_i} \sum_{S_I=i_l+j_l} |(a_{ij})_I - (g_S)_I|$$

بنابر نامساوی میانه، مقادیر $(g_S)_I$ که سمت راست رابطه‌ی فوق را کمینه می‌سازند برابرند با میانه‌ی اعضای مجموعه‌ی $(A_S)_I$. علاوه بر این باید دقت کرد که لندیس I بیانگر تعداد متغیرهای (سری‌های زمانی) به کار برده شده در مدل است و باید هنکل‌سازی بر اساس نرم L_1 برای همه‌ی p متغیر صورت پذیرد تا به این صورت تمام بلوک‌ها با کمترین اختلاف با ماتریس

²¹ Linear Recurrent Formula

۳- کاربرد

$$\mathbf{z}_h = [\mathbf{z}_h^{(1)}, \mathbf{z}_h^{(2)}, \dots, \mathbf{z}_h^{(p)}]^T$$

$$\mathbf{z}_h^{(i)} = [\hat{y}_{N_i-L_i+h+1}^{(i)}, \dots, \hat{y}_{N_i+h-1}^{(i)}] \quad i = 1, 2, \dots, p$$

در این بخش کارایی روش‌های کلاسیک MSSA و VMSSA با استفاده از سری‌های زمانی شبیه‌سازی شده و داده‌های واقعی مورد مقایسه قرار می‌گیرند. به همین منظور مجموعه‌ی $\{1, 2, \dots, L_{sum}\}$ به دو زیرمجموعه داده تقسیم می‌شود. برای مقایسه‌ی عملکرد هر دو روش، در مرحله‌ی بازسازی و پیش‌بینی از دو معیار ریشه‌ی میانگین توان دوم خطا^{۲۲} (RMSE) و میانگین قدرمطلق خطا^{۲۳} (MAE) و همچنین از نسبت معیارهای نام‌برده شده که به صورت رابطه (۱۹) تعریف می‌شوند استفاده شد.

$$RRMSE = \frac{L_1 - MSSA \text{ با استفاده از } RMSE}{MSSA \text{ با استفاده از روش کلاسیک } RMSE}$$

$$RMAE = \frac{L_1 - MSSA \text{ با استفاده از } MAE}{MSSA \text{ با استفاده از روش کلاسیک } MAE} \quad (19)$$

اگر $(RMAE < 1)$ یا $(RRMSE < 1)$ ، آن‌گاه می‌توان گفت روش $MSSA - L_1$ نسبت به روش کلاسیک $MSSA$ از کارایی بیشتری برخوردار است.

۳-۱- مطالعات شبیه‌سازی

در مطالعات شبیه‌سازی، ۲۰۰ داده با استفاده از مدل‌های مختلف ساختاری تولید شده است. از ۱۹۰ مشاهده‌ی اول برای بازسازی و ۱۰ مشاهده‌ی بعدی برای پیش‌بینی استفاده شده است. در تمام روش‌های مطرح شده نکته‌ی قابل توجه تعیین پارامترهای طول پنجره و بازسازی است که مقدار آنها در گروه‌بندی و در نتیجه بازسازی سری‌های زمانی اهمیت زیادی داشته و انتخاب آنها به ساختار داده‌ها و نیز هدف تجزیه و تحلیل بستگی دارد. شبیه‌سازی برای هر یک از مدل‌ها در ۱۰۰۰ تکرار انجام شده و دو معیار RRMSE و RMAE محاسبه می‌شوند. نتایج شبیه‌سازی نشان می‌دهد که دقت روش MSSA به مقدار زیادی به این پارامترها بستگی دارد. انتخاب نامناسب پارامترهای طول پنجره و بازسازی، باعث تغییر ساختار ماتریس مسیر و تغییر در نحوه‌ی گروه‌بندی‌ها شده و نتایج غیر قابل اتکایی پدید خواهد آمد [۲۵].

برای حالت عمودی و پیش‌بینی h گام جلوتر از سری زمانی چندمتغیره به روش VMSSA نیز می‌توان با رگرسیون کردن سطر آخر r ستون اول، از هر کدام از ماتریس بردارهای ویژه بر $L - 1$ سطرها قبل از سطر آخر، که در واقع به صورت بلوک برای هر سری زمانی در مرحله SVD روی ماتریس $X_V X_V^T$ حاصل می‌شود، به دست آورد. در این حالت، $U_j^{(i)\nabla}$ به عنوان $L - 1$ مولفه اول بردار $U_j^{(i)}$ و نماد $\Pi_j^{(i)}$ به عنوان آخرین مؤلفه بردار

$$U_j^{(i)} \text{ در نظر گرفته شده و زمانی که } U_j^{(i)\nabla} = \begin{bmatrix} U_j^{(1)\nabla} \\ U_j^{(2)\nabla} \\ \vdots \\ U_j^{(p)\nabla} \end{bmatrix} \text{ باشد،}$$

ماتریس $U^{p\nabla}$ به صورت $U^{p\nabla} = (U_1^{p\nabla}, U_2^{p\nabla}, \dots, U_r^{p\nabla})$ تعریف می‌شود.

$$\Psi = \begin{bmatrix} \pi_1^{(1)} & \pi_1^{(2)} & \dots & \pi_1^{(r)} \\ \pi_1^{(1)} & \pi_1^{(2)} & \dots & \pi_1^{(r)} \\ \vdots & \vdots & \ddots & \vdots \\ \pi_1^{(p)} & \pi_1^{(p)} & \dots & \pi_1^{(p)} \end{bmatrix} \text{ همچنین ماتریس } \Psi \text{ به صورت}$$

تعریف می‌شود.

اگر ماتریس $(I_{p \times p} - \Psi \Psi^T)^{-1}$ وجود داشته باشد و $r \leq L_{sum} - p$ آن‌گاه h مرحله جلوتر VMSSA که همان پیش‌بینی h گام جلوتر است بر اساس رابطه (۱۸) به دست می‌آید.

$$\begin{bmatrix} \hat{y}_{j_1}^{(1)}, \hat{y}_{j_1}^{(2)}, \dots, \hat{y}_{j_p}^{(p)} \end{bmatrix}^T = \begin{cases} [\hat{y}_{j_1}^{(1)}, \hat{y}_{j_1}^{(2)}, \dots, \hat{y}_{j_p}^{(p)}] & j = 1, 2, \dots, N_i \\ (I_{p \times p} - \Psi \Psi^T)^{-1} \Psi U^{(p)\nabla T} \mathbf{z}_h & j = N_i + 1, \dots, N_i + h \end{cases} \quad (18)$$

در حالی که

$$\mathbf{z}_h = [\mathbf{z}_h^{(1)}, \mathbf{z}_h^{(2)}, \dots, \mathbf{z}_h^{(p)}]^T$$

$$\mathbf{z}_h^{(i)} = [\hat{y}_{N_i-L_i+h+1}^{(i)}, \dots, \hat{y}_{N_i+h-1}^{(i)}] \quad i = 1, 2, \dots, p$$

²² Root Mean Square Error (RMSE)

²³ Mean Absolute Error (MAE)

برای حجم نمونه‌های جایگزین $n = 5$ یا $n = 10$ رسم شده است. ضمن این‌که به هر دو سری زمانی خطا (ε) با توزیع نرمال استاندارد به‌عنوان عامل نوفه اضافه شده است. نمودار هر خط (روند) نشان‌دهنده‌ی نسبت ریشه‌ی میانگین توان دوم خطاها برای هر کدام از روش $MSSA - L_1$ و روش کلاسیک $MSSA$ است که برای طول پنجره‌های متفاوت با پارامترهای بازسازی اختصاص داده شده به آنها با یک حجم خاص و مقدار δ تعیین شده است.

مثال اول: سری‌های زیر به‌عنوان یک سری زمانی دو متغیره برای شبیه‌سازی در نظر گرفته شده‌اند.

$$y_1 = 3 \sin\left(\frac{2t\pi}{13}\right) + \varepsilon$$

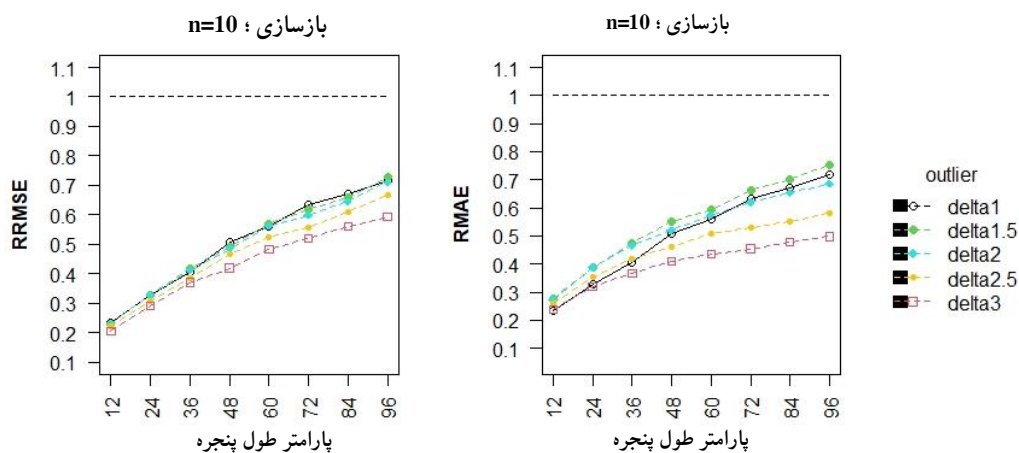
$$y_2 = 2 \sin\left(\frac{2t\pi}{13} + \frac{\pi}{4}\right) + \varepsilon$$

با توجه به این‌که مرتبه‌ی این سری زمانی برای حالت افقی ۲ در نظر گرفته شده [۴]، بنا بر این از ۲ سه‌تایی ویژه اول برای بازسازی و پیش‌بینی این سری زمانی دومتغیره در حالت افقی استفاده شده است. همچنین با توجه به راهنمایی‌های کلی در مورد نحوه‌ی انتخاب پارامتر طول پنجره که در قسمت قبل ارائه شده و نیز با توجه به این‌که پارامتر طول پنجره ضرایبی از دوره‌ی تناوب سری زمانی است، طول پنجره به‌صورت ضرایبی از دوره‌ی تناوب این سری زمانی که ۱۲ است در نظر گرفته شده است.

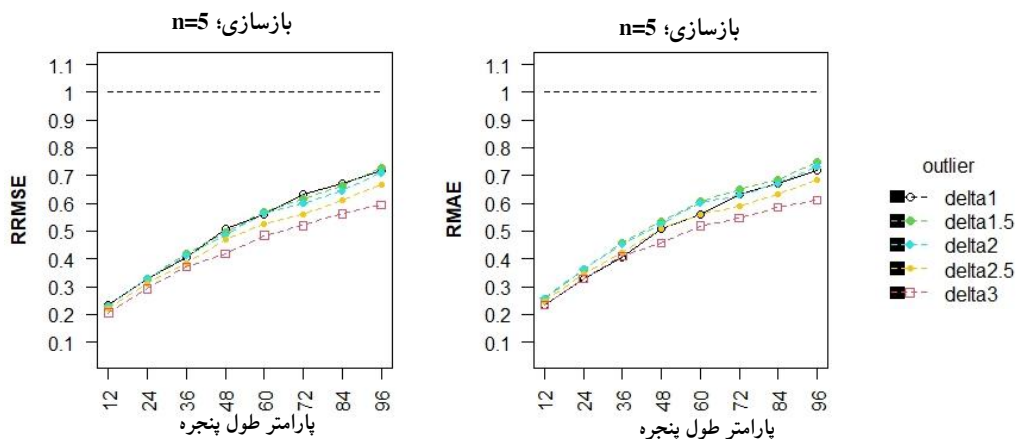
هدف از محاسبه‌ی پارامتر بازسازی در گروه‌بندی جدا کردن نوفه و سیگنال است [۸]، که در واقع همان $-r$ تا r یا $-r$ تا r بزرگ‌ترین مقادیرهای تکین و بردارهای ویژه متناظر برای تقریب سری اصلی انتخاب شده هستند. برای تعیین پارامتر طول پنجره مقادیرهای نزدیک به نصف طول سری زمانی می‌تواند انتخاب مناسبی برای پارامتر طول پنجره باشد. چنان‌چه سری زمانی دارای دوره‌ی تناوب یا مؤلفه‌ی متناوب به‌صورت یک عدد صحیح باشد طول پنجره متناسب با این دوره انتخاب می‌شود [۲۱]. همچنین در روش $MSSA$ برای حالت $HMSSA$ و $VMSSA$ به‌منظور دستیابی به بهترین مقدار برای پارامتر طول پنجره به ترتیب از روابط $L = \left\lfloor \frac{p(N+1)}{p+1} \right\rfloor$ و $L = \left\lfloor \frac{N+1}{p+1} \right\rfloor$ استفاده می‌شود [۱۳] و [۲۶]. در این روابط N طول سری زمانی و p تعداد متغیرهای موجود در سری زمانی است.

در مطالعات شبیه‌سازی شده در این مقاله طول پنجره متناسب با دوره‌ی تناوب سری زمانی و به‌صورت مضربی از آن در نظر گرفته شده است، بنا بر این برای هر ۳ مثال بیان شده مقدار طول پنجره ضرایبی از ۱۲ در نظر گرفته شده‌اند $L = 12, 24, 36, 48, 60, 72, 84, 96$.

از طرفی با توجه به هدف مورد بررسی در مقاله یعنی مقایسه‌ی روش $MSSA - L_1$ با روش کلاسیک $MSSA$ در برابر داده‌های دورافتاده از حضور داده‌های دورافتاده در سری‌های زمانی استفاده شده است. برای تولید داده‌های دورافتاده، ابتدا به تعداد n مشاهده ($n=5, 10$) از سری زمانی شبیه‌سازی شده به‌طور تصادفی انتخاب کرده و سپس به جای آنها، δ برابر آنها را جایگزین شد. در تمامی مثال‌هایی که ارائه خواهند شد مقدار ضرایب δ برابر با $\{1, 1.5, 2, 2.5, 3\}$ است. لازم به توضیح است در تمامی شکل‌ها هر نمودار به‌طور مشخص



شکل ۱- نمودارهای RMAE و RRMSE برای بازسازی با اندازه نمونه ۱۰



شکل ۲- نمودارهای RRMSE و RMAE برای بازسازی با اندازه نمونه ۵

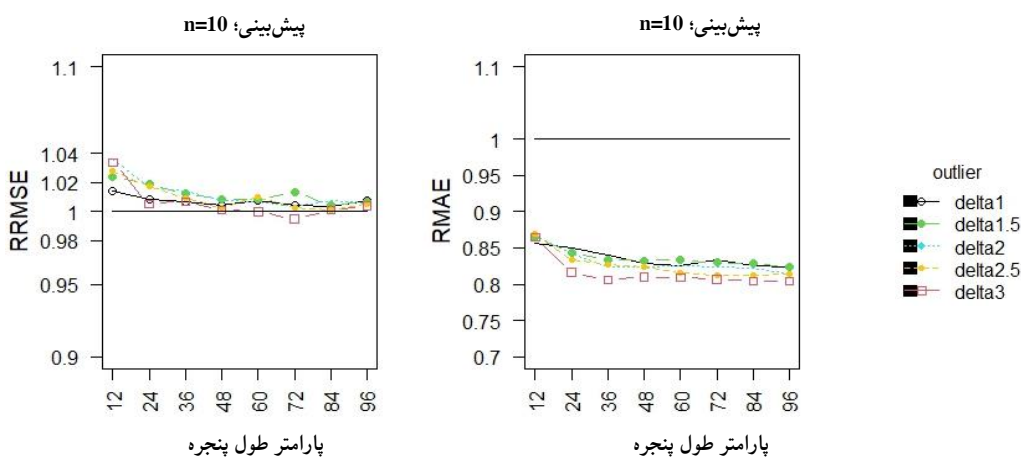
زمانی داده‌های دورافتاده با ضرایب $\delta = \{1, 1/5, 2, 2/5, 3\}$ به صورت تصادفی به ترتیب با حجم نمونه ۵ و ۱۰- تایی اضافه شده است.

در هر دو حالتی که حجم نمونه‌ی انتخابی ۱۰ یا ۵ بوده، با افزایش ضریب داده‌های دور افتاده، $RRMSE < 1$ شده و عملکرد روش $L_1 - MSSA$ بهتر از روش کلاسیک مجموعه مقادیر تکین چندمتغیره است.

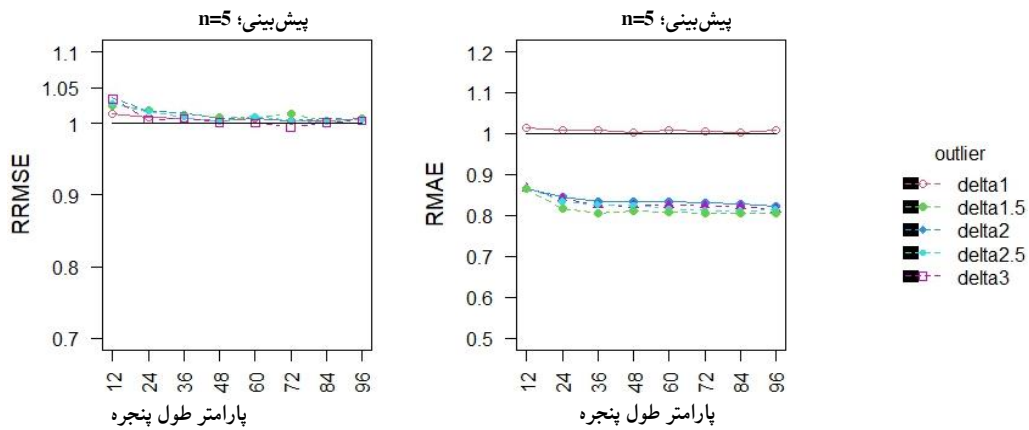
اما برای نمودار RMAE با انتخاب نمونه ۱۰- تایی در تمام ضرایب داده‌های دورافتاده، عملکرد روش $L_1 - MSSA$ بهتر از روش کلاسیک مجموعه مقادیر تکین چندمتغیره است. بر اساس همین معیار در پیش‌بینی، زمانی که حجم نمونه ۵- تایی باشد، در حالی که مشاهده‌ی داده‌ی دورافتاده‌ای وجود نداشته باشد عملکرد هر دو روش یکسان است.

شکل‌های ۱ و ۲ نمودارهایی هستند که برای مقایسه بازسازی سری‌های زمانی به دو روش $L_1 - MSSA$ و روش کلاسیک $MSSA$ رسم شده‌اند. در شکل ۱ به ازای $\delta = \{1, 1/5, 2, 2/5, 3\}$ انتخاب یک نمونه ۱۰- تایی برای ساختن داده‌های دورافتاده و در شکل ۲ به ازای همین ضرایب یک نمونه ۵- تایی برای ساختن داده‌های دورافتاده در نظر گرفته شده‌اند. بنا بر این همان‌طور که در این شکل‌ها مشاهده می‌شود برای هر دو نمونه‌ی انتخابی نسبت خطاها یعنی RRMSE و RMAE کمتر از یک شده و این بدان معناست که با حضور داده‌های دور افتاده در سری زمانی عملکرد روش $L_1 - MSSA$ در بازسازی سری‌های زمانی بهتر از روش کلاسیک است.

شکل‌های ۳ و ۴ نمودارهایی هستند که برای مقایسه پیش‌بینی سری‌های زمانی به دو روش کلاسیک $MSSA$ و $L_1 - MSSA$ ارائه شده‌اند. همانند قسمت بازسازی در این شکل‌ها به سری‌های



شکل ۳- نمودارهای RRMSE و RMAE برای پیش‌بینی با اندازه نمونه ۱۰



شکل ۴- نمودارهای RRMSE و RMAE برای پیش‌بینی با اندازه نمونه ۵

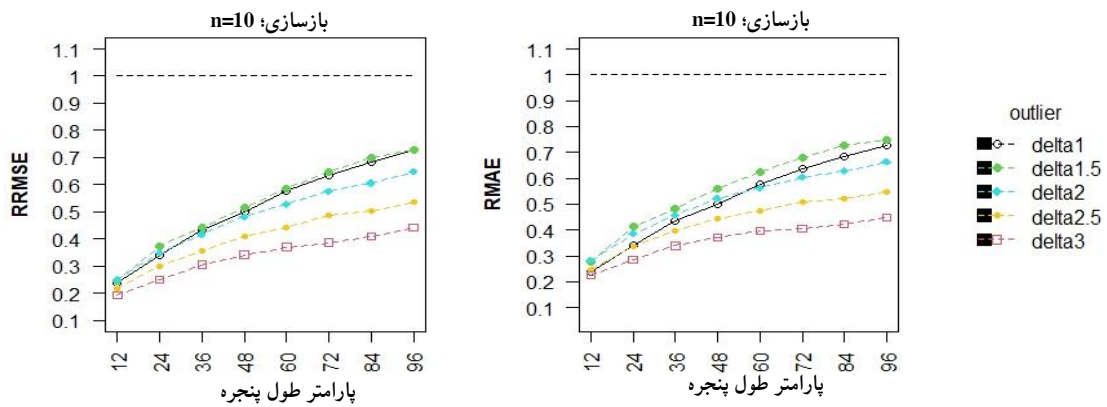
بازسازی و پیش‌بینی این سری زمانی دومتغیره در حالت افقی استفاده شده است. همچنین طول پنجره را به‌عنوان ضرایبی از دوره‌ی تناوب سری زمانی که ۱۲ است در نظر گرفته شده است.

مثال دوم: سری‌های زیر به‌عنوان یک سری زمانی دو متغیره برای شبیه‌سازی در نظر گرفته شده‌اند:

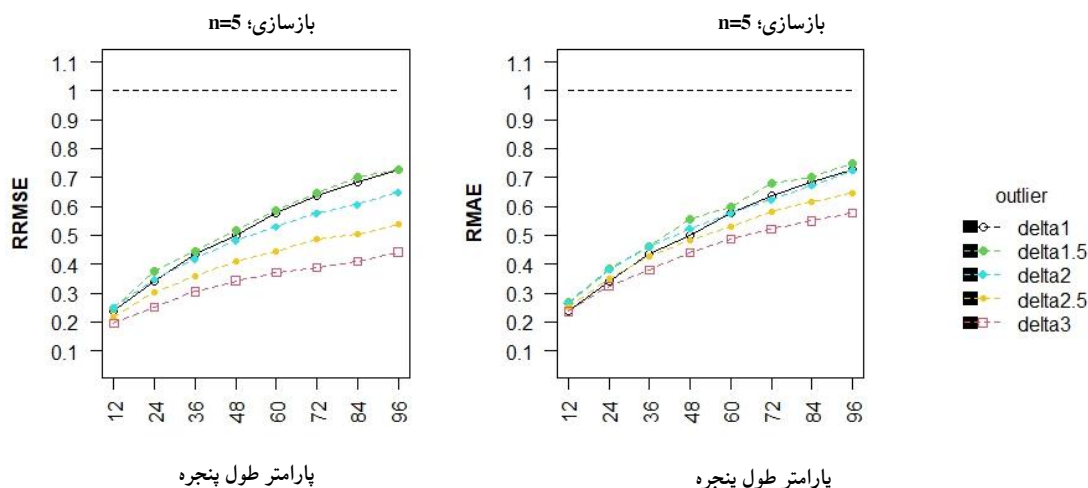
$$y_1 = 3 \cos\left(\frac{2t\pi}{12}\right) + \varepsilon$$

$$y_2 = 3 \cos\left(\frac{2t\pi}{12} + \frac{\pi}{2}\right) + \varepsilon$$

رتبه‌ی ماتریس مسیر این سری‌های زمانی برای حالت افقی ۲ در نظر گرفته شده است [۴]. بنا بر این از ۲، سه‌تایی ویژه اول برای



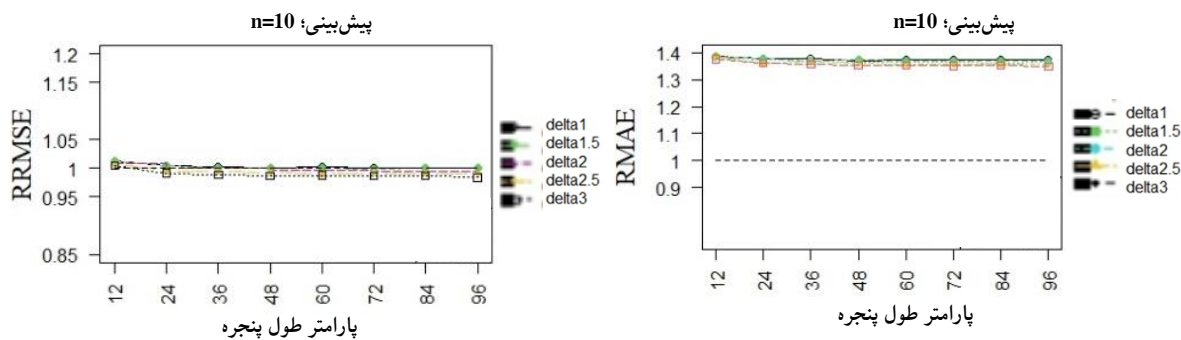
شکل ۵- نمودارهای RRMSE و RMAE برای بازسازی با اندازه نمونه ۱۰



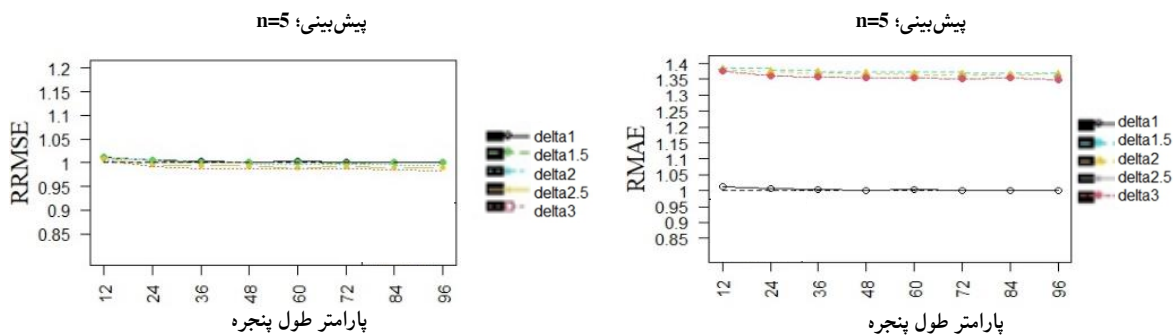
شکل ۶- نمودارهای RRMSE و RMAE برای بازسازی با اندازه نمونه ۵

نمودارهای RRMSE و RMAE در شکل ۷ برای مقایسه پیش‌بینی سری‌های زمانی به دو روش کلاسیک $MSSA$ و $L_1 - MSSA$ مورد استفاده قرار گرفته‌اند. در نمودار $RRMSE$ نسبت خطای بین دو روش کمتر از یک شده و در نتیجه عملکرد روش $L_1 - MSSA$ بهتر از عملکرد روش کلاسیک مجموعه مقادیر تکین چندمتغیره است. اما برای $RMAE$ با انتخاب نمونه ۱۰- تایی در تمام ضرایب داده‌های دورافتاده، عملکرد روش کلاسیک $MSSA$ بهتر از روش $L_1 - MSSA$ است.

شکل‌های ۵ و ۶ نمودارهایی هستند که برای مقایسه بازسازی سری‌های زمانی به دو روش $L_1 - MSSA$ و کلاسیک $MSSA$ رسم شده‌اند. برای هر دو حالت حجم نمونه ۵ و ۱۰- تایی و تمام ضرایب تولید داده‌های دورافتاده، نسبت خطاها یعنی $RRMSE$ و $RMAE$ کمتر از یک شده و این بدان معناست که با حضور داده‌های دور افتاده در سری زمانی عملکرد روش $L_1 - MSSA$ در بازسازی سری‌های زمانی بهتر از حالت روش کلاسیک $MSSA$ است.



شکل ۷- نمودارهای RRMSE و RMAE برای پیش‌بینی با اندازه نمونه ۱۰



شکل ۸- نمودارهای RMAE و RRMSE برای پیش‌بینی با اندازه نمونه ۵

همچنین طول پنجره به‌عنوان ضرایبی از دوره‌ی تناوب سری زمانی که ۱۲ است در نظر گرفته شده است.

در شکل ۹ نمودارهای RRMSE و RMAE در حالت بازسازی برای مثال سوم نتایج متفاوت‌تر با خروجی دو مثال دیگر دارند. در این مثال نمودار RRMSE برای مقایسه‌ی عملکرد $L_1 - MSSA$ و روش کلاسیک $MSSA$ در حالت بازسازی با حضور ۱۰ داده‌ی دورافتاده با ضرایب $\{1.5, 2, 2.5, 3\}$ را نشان می‌دهد. با توجه به این نمودار متوجه می‌شویم چنان‌چه ضرایب تولید داده‌های دورافتاده بزرگتر شوند عملکرد روش $L_1 - MSSA$ بهتر از عملکرد روش کلاسیک $MSSA$ خواهد بود. بنا بر این تنها در دو حالت که ضرایب تولید داده دورافتاده ۲ و ۳ باشد، $RRMSE < 1$ خواهد بود. نمودار RMAE تفسیری شبیه نمودار RRMSE دارد و عملکرد $L_1 - MSSA$ زمانی که داده‌های دورافتاده بزرگتری در مشاهدات باشد در مقایسه با روش کلاسیک $MSSA$ بهتر خواهد بود و $RMAE < 1$ است.

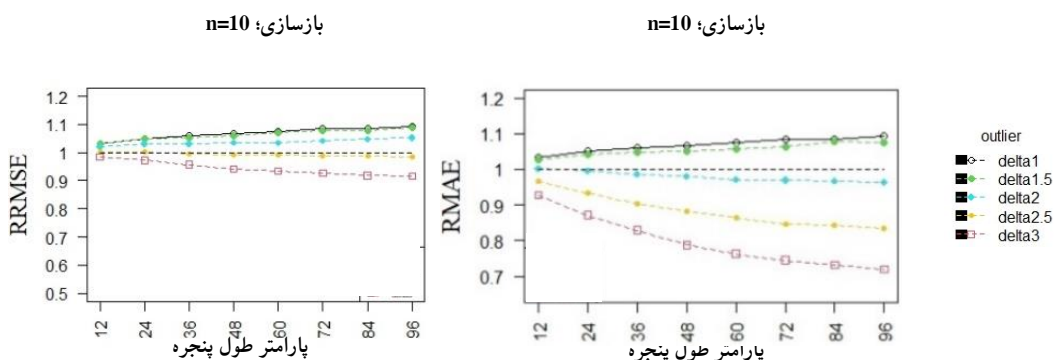
شکل ۸ که شامل نمودارهای RRMSE و RMAE برای پیش‌بینی سری‌های زمانی دو متغیره‌ی مثال دوم است، نشان می‌دهد چنان‌چه اندازه نمونه‌ی انتخابی برای اضافه کردن داده‌های دورافتاده ۵ باشد با افزایش ضریب داده‌های دورافتاده، با توجه به معیار RRMSE عملکرد روش $L_1 - MSSA$ بهتر از عملکرد روش کلاسیک $MSSA$ است.

مثال سوم: سری‌های زیر به‌عنوان یک سری زمانی دو متغیره برای شبیه‌سازی در نظر گرفته شده‌اند.

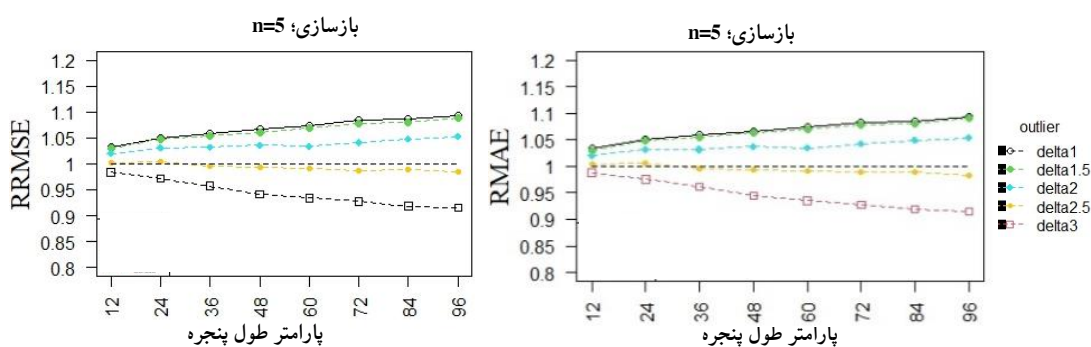
$$y_1 = 3 \cos\left(\frac{y_1 t \pi}{12}\right) + \varepsilon$$

$$y_2 = 2 \cos\left(\frac{y_2 t \pi}{12} + \frac{\pi}{4}\right) + \varepsilon$$

رتبه‌ی ماتریس مسیر این سری زمانی برای حالت افقی ۲ در نظر گرفته شده است [۴]. بنا بر این از ۲، سه‌تایی ویژه اول برای بازسازی و پیش‌بینی این سری زمانی دومتغیره در حالت افقی استفاده شده است.



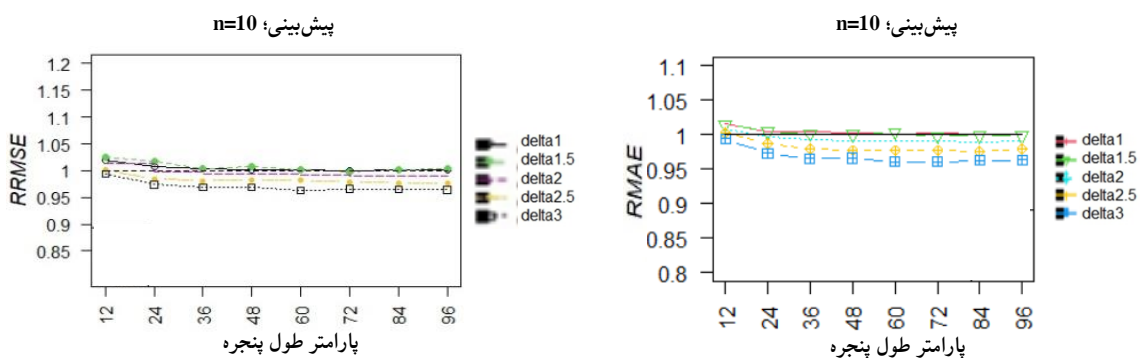
شکل ۹- نمودارهای RMAE و RRMSE برای بازسازی با اندازه نمونه ۱۰



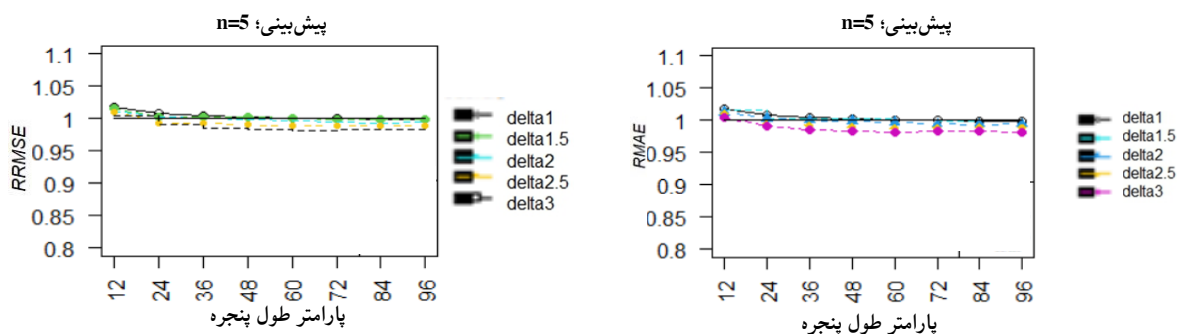
شکل ۱۰- نمودارهای RRMSE و RMAE برای بازسازی با اندازه نمونه ۵

شکل ۱۰ نمودارهای RRMSE و RMAE عملکرد دو روش $L_1 - MSSA$ و کلاسیک $MSSA$ را در مرحله‌ی بازسازی نشان می‌دهد. با توجه به این شکل زمانی که ضریب داده‌های دورافتاده افزایش می‌یابد نسبت خطاها کاهش یافته تا جایی که برای $\delta = \{2/5, 3\}$ در هر دو مورد مقادیرهای RRMSE و RMAE کمتر از یک شده و در نتیجه عملکرد $L_1 - MSSA$ در مقایسه با روش کلاسیک $MSSA$ در مرحله‌ی بازسازی بهتر است.

شکل ۱۰ نمودارهای RRMSE و RMAE عملکرد دو روش $L_1 - MSSA$ و کلاسیک $MSSA$ را در مرحله‌ی بازسازی نشان می‌دهد. با توجه به این شکل زمانی که ضریب داده‌های دورافتاده افزایش می‌یابد نسبت خطاها کاهش یافته تا جایی که برای $\delta = \{2/5, 3\}$ در هر دو مورد مقادیرهای RRMSE و RMAE کمتر از یک شده و در نتیجه عملکرد $L_1 - MSSA$ در مقایسه با روش کلاسیک $MSSA$ در مرحله‌ی بازسازی بهتر است.



شکل ۱۱- نمودارهای RRMSE و RMAE برای پیش‌بینی با اندازه نمونه ۱۰



شکل ۱۲- نمودارهای RRMSE و RMAE برای پیش‌بینی با اندازه نمونه ۵

شکل‌های ۱۱ و ۱۲ به ترتیب نمودارهای RRMSE و RMAE برای حضور ۱۰ و ۵ داده‌ی دورافتاده در مشاهدات سری‌های زمانی چندمتغیره، این نمودارها بهتر بودن عملکرد $L_1 - MSSA$ در مرحله‌ی پیش‌بینی را نسبت به عملکرد روش کلاسیک $MSSA$ نشان می‌دهند. همچنین از این نمودارها می‌توان به این نتیجه دست یافت که در مرحله‌ی پیش‌بینی

شکل‌های ۱۱ و ۱۲ به ترتیب نمودارهای RRMSE و RMAE برای حضور ۱۰ و ۵ داده‌ی دورافتاده در مشاهدات سری‌های زمانی چندمتغیره، این نمودارها بهتر بودن عملکرد $L_1 - MSSA$ در مرحله‌ی پیش‌بینی را نسبت به عملکرد روش کلاسیک $MSSA$ نشان می‌دهند. همچنین از این نمودارها می‌توان به این نتیجه دست یافت که در مرحله‌ی پیش‌بینی

روش ۶ مقدار اول مقادیر ویژه برای بازسازی و پیش‌بینی در نظر گرفته می‌شود.

با بررسی قدم به قدم نمودارهای «بردارهای ویژه جفتی» و «W-همبستگی» و وجود دوره‌ی تناوب ۳ و ۴ برای داده‌ها، پارامتر طول پنجره را مقداری بین ۳۶ (ضریبی از دوره تناوب داده‌ها) به‌عنوان کمترین و ۴۶۸ به‌عنوان بیشترین طول پنجره‌ای که می‌توان انتخاب کرد، در نظر گرفته شده است. از ۳۶ تا ۴۶۸ در هر مرحله ضریبی از ۳۶ را به مقدار طول پنجره اضافه کرده و در محاسبات، به دست آوردن پارامتر بازسازی در نظر گرفته شد. با توجه به این مطلب که برای تعیین پارامتر بازسازی یا نحوه گروه‌بندی داده‌های واقعی روش ریاضی و فرموله‌شده خاصی وجود ندارد بنا بر این به ازای هر طول پنجره، نمودارهای «بردارهای ویژه»، «بردارهای ویژه جفتی» و «W-همبستگی» مجدداً رسم شده و بر اساس اطلاعات حاصل شده از نمودارها و تجربه محقق، پارامتر بازسازی متناسب با هر طول پنجره تعیین شد. با تعیین پارامتر طول پنجره و پارامتر بازسازی، دو گام گروه‌بندی و میانگین‌گیری از مرحله‌ی بازسازی داده‌ها انجام می‌شود. به همین خاطر متناسب با طول‌های پنجره متفاوت، مقادیر پارامتر بازسازی نیز تعیین شده و خروجی‌های برنامه متناسب با آنها به دست آمد، سپس بر اساس معیارهای RMAE و RRMSE، عملکرد دو روش مورد مقایسه قرار گرفت. نتیجه این مقایسه‌ها نیز به‌صورت نمودارهای ۱۵ برای بازسازی و ۱۶ برای پیش‌بینی رسم شده است.

چنانچه تعداد داده‌های دورافتاده بیشتر و بزرگتر باشند عملکرد $L_1 - MSSA$ بهتر از عملکرد روش کلاسیک $MSSA$ است. با مشاهده‌ی همه‌ی نمودارهای $RRMSE$ و $RMAE$ در هر دو حالت نمونه‌های ۵ و ۱۰ تایی، این‌گونه به نظر می‌رسد که انتخاب مناسب طول پنجره در رسیدن به نتیجه مطلوب تأثیرگذار خواهد بود.

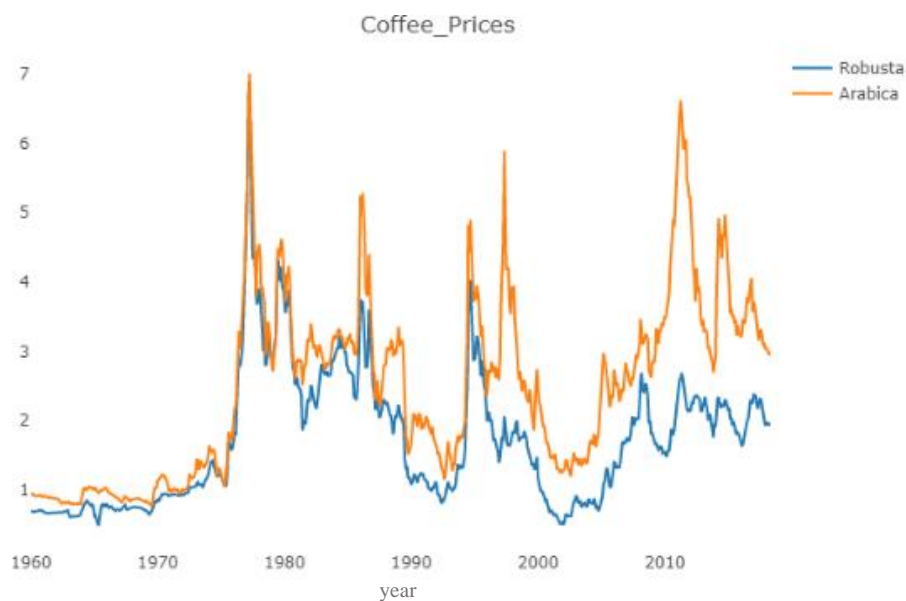
۳-۲- داده‌های واقعی

در این بخش برای بررسی عملکرد کارایی روش تحلیل مجموعه مقادیر تکین چندمتغیره‌ی استوار در حالت افقی با روش کلاسیک $MSSA$ در حالت افقی از داده‌های قیمت قهوه که یک سری زمانی دو متغیره با ۷۰۱ مشاهده به‌صورت ماهانه از January سال ۱۹۶۰ تا May سال ۲۰۱۸ می‌باشد، استفاده شده است. این داده‌ها در نرم‌افزار R و در بسته‌ی TSstudio با نام Coffee_Prices قابل دسترس هستند. دلیل استفاده از این داده‌ها اولاً داشتن نقاط دورافتاده در نمودار سری زمانی و دوماً دومتغیره بودن آنها می‌باشد (شکل ۱۳). از ۶۶۶ مشاهده‌ی اول این سری زمانی برای بازسازی و ۳۵ مشاهده‌ی بعدی برای پیش‌بینی استفاده شده است.

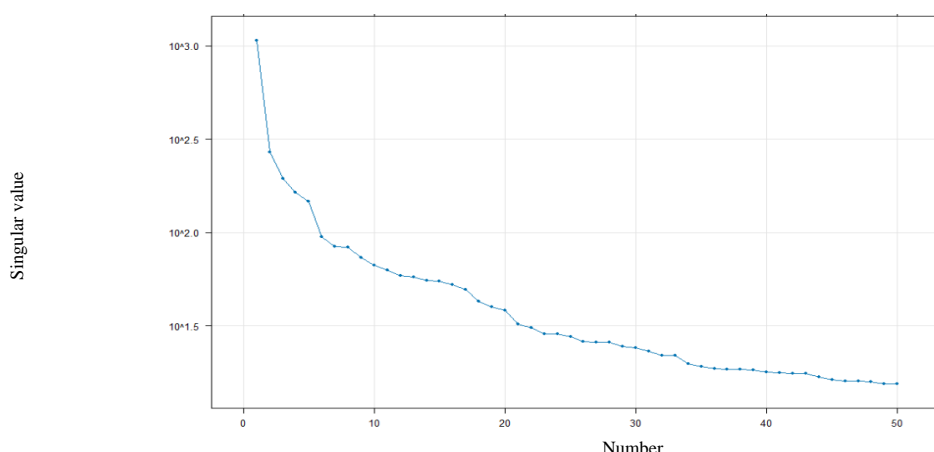
با توجه به این‌که برای انجام محاسبات لازم، ماتریس مسیر را به‌صورت افقی در نظر گرفته‌ایم لذا با استفاده از رفتار مقادیر ویژه ماتریس مسیر برای طول پنجره ۴۶۸ که در شکل ۱۴ نمایش داده شده و با توجه به مطالب و دلایل ارائه شده [۹] در هر دو

$$L = \frac{p}{p+1} (N + 1) \quad (۲۴) \text{ (بر اساس فرمول)}$$

نشریه مهندسی و مدیریت کیفیت



شکل ۱۳- سری‌های زمانی فروش قهوه از ماه اول سال ۱۹۶۰ تا ماه پنجم سال ۲۰۱۸

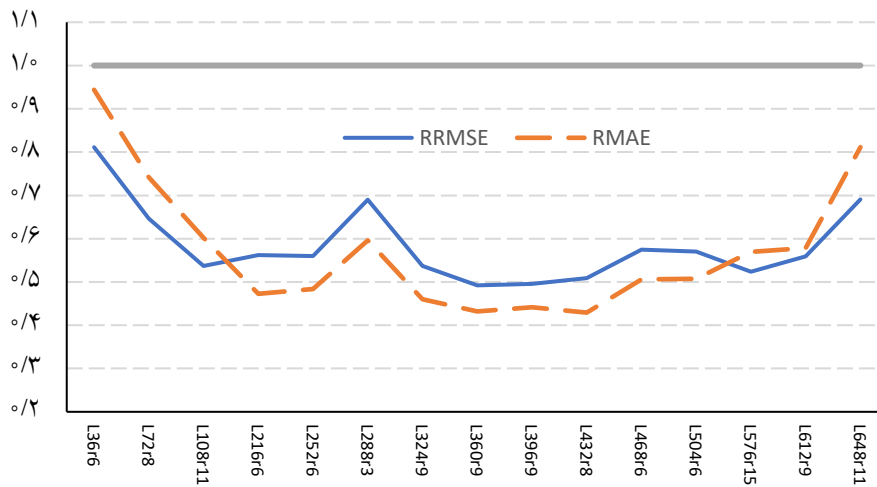


شکل ۱۴- نمودار رفتار مقادیر ویژه ماتریس مسیر برای طول پنجره ۴۶۸

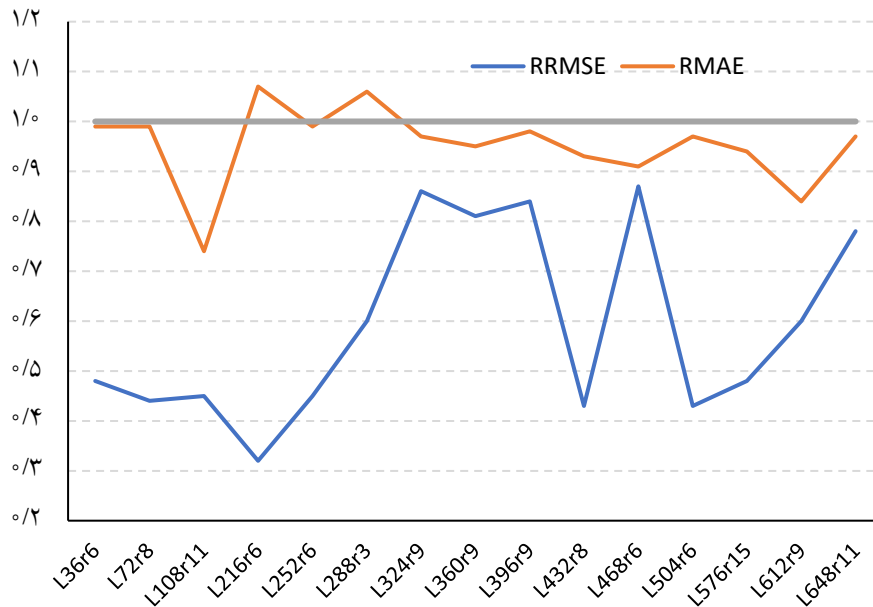
واقعیست می‌باشد که تقریباً با افزایش طول پنجره و پارامتر بازسازی، کارایی روش $L_1 - MSSA$ بهتر از روش کلاسیک $MSSA$ است و مقدار $RRMSE < 1$ و $RMAE < 1$ است.

شکل ۱۵ نشان می‌دهد که تقریباً با افزایش طول پنجره و پارامتر بازسازی، کارایی روش $L_1 - MSSA$ بهتر از روش کلاسیک $MSSA$ است و مقدار $RRMSE < 1$ و $RMAE < 1$ است.

شکل ۱۶ که مقایسه عملکرد دو روش $L_1 - MSSA$ و روش کلاسیک $MSSA$ را در مرحله‌ی پیش‌بینی نشان می‌دهد گویای این



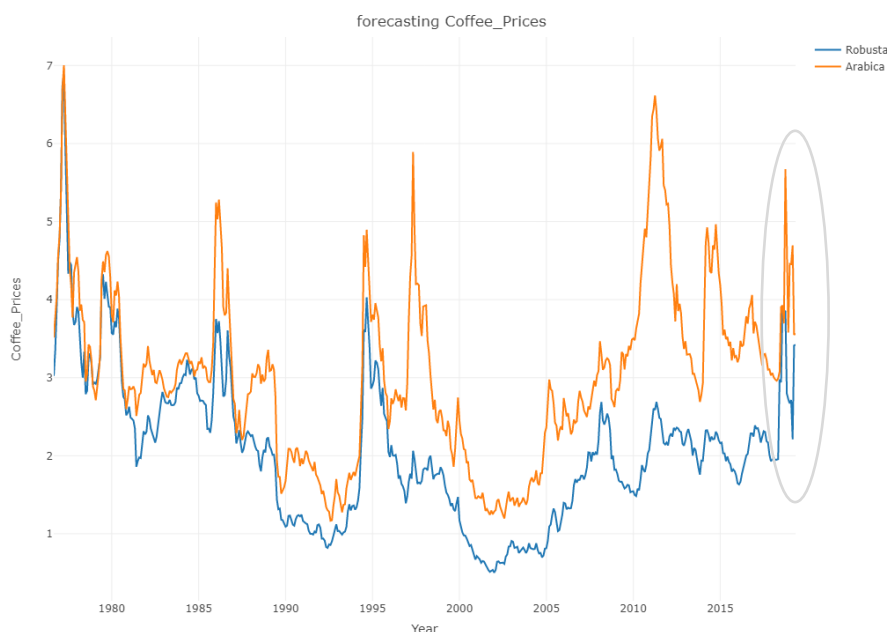
شکل ۱۵- نمودارهای RMAE و RRMSE برای مقایسه عملکرد دو روش $L_1 - MSSA$ و کلاسیک $MSSA$ در بازسازی داده‌های واقعی



شکل ۱۶- نمودارهای RMAE و RRMSE برای مقایسه عملکرد روش $L_1 - MSSA$ و روش کلاسیک $MSSA$ در پیش‌بینی داده‌های واقعی

در این نمودار مقادیر پیش‌بینی با کادری که به دور آنها کشیده شده مشخص شده است.

در شکل ۱۷، پیش‌بینی یک‌ساله برای سری‌های زمانی فروش قهوه با استفاده از روش بازگشتی خطی (LRF) و در نظر گرفتن پارامتر طول پنجره‌ی ۴۶۸ و پارامتر بازسازی ۶ انجام شده است.



شکل ۱۷- نمودار پیش‌بینی سری‌های زمانی فروش قهوه از ماه ششم سال ۲۰۱۸ تا ماه پنجم سال ۲۰۱۹

برای ارزیابی عملکرد نسخه‌ی جدید پیشنهادی ($L_1 - MSSA$) در روش کلاسیک $MSSA$ از معیارهای $RRMSE$ و $RMAE$ در مطالعات شبیه‌سازی و داده‌های واقعی استفاده شد. در این ارزیابی با اضافه کردن داده‌های دورافتاده با حجم و ضرایب متفاوت در مرحله‌ی شبیه‌سازی عملکرد دو روش بیان شده در مرحله‌ی بازسازی و پیش‌بینی مورد مقایسه گرفتند. براساس مشاهدات حاصل از نتایج تحلیل شبیه‌سازی و داده‌های واقعی نشان داده شد که اگر داده‌های دورافتاده در سری زمانی وجود داشته باشند، تقریباً در تمام موارد کارایی روش ناپارامتری مجموعه مقادیر تکین استوار شده چند متغیره $L_1 - MSSA$ بهتر از حالت کلاسیک این روش است.

۵- منابع

- [1] Hassani, H., Kalantari, M., and Yarmohammadi, M. (2017). An improved SSA forecasting result based on a filtered recurrent forecasting algorithm. *Comptes Rendus. Mathématique*, 355(9), 1026-1036.
- [2] Broomhead, D. S., and King, G. P. (1986 b). On the Qualitative Analysis of Experimental Dynamical Systems. *Nonlinear Phenomena and Chaos*, 113, 114.

۴- نتیجه‌گیری

در این مقاله به منظور بهینه‌سازی کیفیت مدل‌سازی و پیش‌بینی به دست آمده از روش کلاسیک $MSSA$ با حضور داده‌های دورافتاده، روش استوارسازی بر اساس تغییر نُرم به کار رفته در ماتریس مسیر پیشنهاد شد. این روش که یک نسخه‌ی جدید از روش کلاسیک $MSSA$ بر اساس نُرم L_1 ، و تغییر ساختار ماتریس مسیر در مرحله‌ی بازسازی است به عنوان شیوه‌ای نوین در استوارسازی سری‌های زمانی برای بهبود مدل‌سازی شناخته می‌شود. لذا این روش علاوه بر بهینه‌سازی کیفیت مدل پیشنهادی، باعث افزایش کارایی تحلیل و پیش‌بینی سری‌های زمانی نیز می‌شود.

برخی محققان برای افزایش قابلیت‌های بازسازی و پیش‌بینی در روش SSA ، شناسایی و حذف نقاط دورافتاده را پیشنهاد داده‌اند زیرا SSA مبتنی بر نُرم L_2 نسبت به نقاط دورافتاده حساس‌تر است [۱۹]. در حالی که نتایج مقاله حاضر نشان می‌دهد که در صورت مواجهه با سری‌های زمانی دارای نقاط دورافتاده، می‌توان با تغییر نُرم ماتریسی در $MSSA$ از L_2 به L_1 به نتایج بهتری دست پیدا کرد زیرا نیازی به تبدیل داده وجود نداشته و تضمین می‌کند که هیچ اطلاعاتی از دست نرود.

- and Reliability Engineering International, 32, 2127–2137.
- [15] Mahmoudvand, R.; Alehosseini, F.; and Rodrigues, P. C. (2015). Mortality Forecasting with Singular Spectrum Analysis. *RevStat - Statistical Journal*, 13, 193–206.
- [16] Mahmoudvand, R., Konstantinides, D. and Rodrigues, P.C. (2017). Forecasting Mortality Rate by Multivariate Singular Spectrum Analysis. *Applied Stochastic Models in Business and Industry*. DOI: 10.1002/asmb.2274.
- [17] Hassani, H., Soofi, A., and Avazalipour, M. S. (2011). Forecasting GDP with aggregated and sectoral data. *Fluctuation and Noise Letters*, 10 (3), 249-265.
- [18] Sirimal Silva, E., Hassani, H. and Heravi, S. (2018). Modeling European industrial production with multivariate singular spectrum analysis: A cross-industry analysis, *Journal of Forecasting*, <https://doi.org/10.1002/for.2508>
- [19] Mahmoudvand, R. (2012). Some Theoretical Development of the Singular Spectrum Analysis. Ph.D. Thesis (In Persian), Shahid Beheshti University, Tehran, Iran.
- [20] Kalantari, M., Yarmohammadi, M. and Hassani, H. (2016). Singular Spectrum Analysis Based on L_1 -norm. *Fluctuation and Noise Letters*, 15(01), 1650009.
- [21] Bloomfield, P. and W. Steiger (1983). *Least Absolute Deviations: Theory, Applications and Algorithms*. Boston: Birkhäuser.
- [22] Birkes, D., & Dodge, Y. (1993). *Alternative Methods of Regression*. John Wiley&Sons. Inc., New York.
- [23] Golyandina, N. and Zhigljavsky, A. (2018). *Singular Spectrum Analysis for Time Series* Second Edition., Springer Briefs in Statistics.
- [24] Golyandina, N., Nekrutkin, V. and Zhigljavsky, A. (2001). *Analysis of Time Series Structure: SSA and Related Techniques*. London: Chapman & Hall/CRC.
- [25] Rodrigues, P. C., & Mahmoudvand, R. (2018). The benefits of multivariate singular spectrum analysis over the univariate version. *Journal of the Franklin Institute*, 355(1), 544-564.
- [3] Ghil, M., Allen, M. R., Dettinger, M. D., Ide, K., Kondrashov, D., Mann, M. E., ... and Yiou, P. (2002). Advanced spectral methods for climatic time series. *Reviews of geophysics*, 40(1), 3-1.
- [4] Golyandina, N. E. and Stepanov, D. (2005). SSA-based approaches to analysis and forecast of multidimensional time series. *Proceedings of the 5th St.Petersburg Workshop on Simulation*, 293–298.
- [5] Wang, R., Ma, H.-G., Liu, G.-Q. and Zuo, D.-G. (2015). Selection of window length for singular spectrum analysis, *Journal of the Franklin Institute*, 352, 1541–1560.
- [6] Alharbia, N. and Hassani, H. (2016). A new approach for selecting the number of the eigenvalues in singular spectrum analysis, *Journal of the Franklin Institute*, 353, 1-16.
- [7] Hassani, H. and Rua, A. and Silva, E.S. and Thomakos, D. (2019). Monthly forecasting of GDP with mixed-frequency multivariate singular spectrum analysis. *International Journal of Forecasting*, 35 (4). pp. 1263-1272.
- [8] Patterson, K., Hassani, H., Heravi, S., and Zhigljavsky, A. (2011). Forecasting the final vintage of the industrial production series, *Journal of Applied Statistics*, 38, 2183–2211.
- [9] Hassani, H. and Mahmoudvand, R. (2013). Multivariate singular spectrum analysis: A general view and new vector forecasting approach, *Int. J. Energy Stat.* 1(1) (2013) 55–83.
- [10] de Carvalho, M., Rodrigues, P.C. and Rua, A. (2012). Tracking the US business cycle with a singular spectrum analysis, *Economics Letters*, 114, 32–35.
- [11] Rodrigues, P.C. and de Carvalho, M. (2013). Spectral modeling of time series with missing data, *Applied Mathematical Modelling*, 37, 4676–4684.
- [12] Hasani, H., Heravi, S., and Zhigljavsky, A. (2009). Forecasting European industrial production with singular spectrum analysis. *International journal of forecasting*, 25(1), 103-118.
- [13] Mahmoudvand, R. and Rodrigues, P.C. (2016). Missing value imputation in time series using Singular Spectrum Analysis. *International Journal of Energy and Statistics*, 4(01), 1650005.
- [14] Rodrigues, P.C., and Mahmoudvand, R. (2016). Correlation analysis in contaminated data by singular spectrum analysis, *Quality*

Improving The Quality of Time Series Modeling and Forecasting Using Robust Multivariate Singular Spectrum Analysis

Tahere Amini

Ph.D. Student in Statistics, Payame Noor University, Tehran, Iran.
tahreamini90@yahoo.com

Masoud Yarmohammadi¹

(Corresponding Author): Associate Professor of Statistics, Payame Noor University, Tehran, Iran. masyar@pnu.ac.ir

Ali Shadrokh

Associate Professor of Statistics, Payame Noor University, Tehran, Iran.
a.shadrokh@pnu.ac.ir

Mahdi Kalantari

Assistant Professor of Statistics, Payame Noor University, Tehran, Iran.
kalantarimahdi@pnu.ac.ir

Abstract

In time series analysis, ignoring outliers leads to misidentification of the model, biased estimation of parameters, and poor predictions. One of the reliable non-parametric methods in predicting and improving the quality of multivariate time series modeling is the multivariate Singular Spectrum Analysis (MSSA) technique, which does not require any initial assumptions. The presence of outliers affects the Frobenius norm of matrix and reduces the efficiency of the MSSA method. In this research, a new version of MSSA based on the L_1 -norm is proposed. Then the performance of this method is compared with basic MSSA using simulation studies and real data.

Key Words: Multivariate Singular Spectrum Analysis (MSSA), Robustification, Outliers, L_1 -norm, Ratio of Root Mean Square Errors (RRMSE), Ratio of Mean Absolute Errors (RMAE).

1- Introduction

Multivariate Singular Spectrum Analysis (MSSA) is a nonparametric method to analyze signals without any assumptions of the underlying system. It can decompose the original time series into a set of components, which are recognized as either a trend, periodic or quasi-periodic signal or noise. The matrix norm used in MSSA is the Frobenius norm which is not robust to the outliers. In order to robustify the method, the Frobenius norm is replaced by the L_1 -norm and a new version of MSSA (L_1 -MSSA) is proposed. In this approach; a method for computing the signal matrix and the hankelization method using the L_1 -norm are proposed. Using real-world data and the simulation studies the

¹ Corresponding Author: masyar@pnu.ac.ir

performance of L_1 -MSSA and basic MSSA are compared. Applying the Root Mean Square Errors (RMSE) and Mean Absolute Error (MAE) criteria, show that the L_1 -MSSA outperforms basic MSSA in reconstruction and forecasting.

2- Methodology

The Multivariate Singular Spectrum Analysis (MSSA) technique consists of two complementary stages: Decomposition and Reconstruction, and both of them include two separate steps. In the first stage, we decompose the series to enable signal extraction and noise reduction. In the second stage, we reconstruct a less noisy series and use the reconstructed series for forecasting new data points [1]. In the MSSA method there are two forms for the trajectory matrix: vertical and horizontal. The vertical and the horizontal forms are called VMSSA and HMSSA, respectively. The matrix norm used in basic MSSA is the Frobenius norm which is not robust in the presence of outliers. In this paper, in order to robustify the method, the Frobenius norm is replaced by the L_1 norm and a new version of MSSA (L_1 -MSSA) is proposed. In L_1 -MSSA, the negative effects of outliers is reduced by utilizing the median instead of averaging for trajectory matrix hankelization in the last step of time series reconstruction (averaging) and the performance of L_1 -MSSA and basic MSSA in reconstruction and forecasting are evaluated using the real-world and simulated time series that are contaminated with outliers. To measure the accuracy of forecasting results, we use the commonly adopted forecasting performance evaluation measures of Root Mean Squared Errors (RMSE) and Mean Absolute Errors (MAE). Moreover, the effects of outliers on forecasting and reconstruction for two procedures are evaluated. The following ratios are used for comparing the established and newly proposed MSSA methods:

$$RRMSE = \frac{RMSE \text{ based on } L_1 - MSSA}{RMSE \text{ based on basic MSSA}}$$

$$RMAE = \frac{MAE \text{ based on } L_1 - MSSA}{MAE \text{ based on basic MSSA}}$$

If $RRMSE < 1$ and $RMAE < 1$, then the L_1 -MSSA procedure outperforms basic MSSA. Alternatively, when $RRMSE > 1$ and $RMAE > 1$, it would indicate that the performance of the $L_1 - MSSA$ procedure is worse than basic MSSA.

The purpose of determining the reconstruction parameter in the grouping step of MSSA is to separate noise and signal [6]. This parameter is equal to the first r singular values and the corresponding eigenvectors are chosen to approximate the original series. Selection of the proper window length (denoted by L) depends on the structure of time series and the preliminary information of the data. Theoretical results indicate that L should be large enough but not greater than $\left\lfloor \frac{N}{2} \right\rfloor$. Furthermore, if we know that the time series may have a periodic component with an integer period then to get better separability of this periodic component it is advisable to take the window length proportional to that period [3]. Regarding the window length, Hassani and Mahmoudvand (2013) showed that a value close to the $L = \left\lfloor \frac{p(N+1)}{p+1} \right\rfloor$ and $L = \left\lfloor \frac{(N+1)}{p+1} \right\rfloor$

is optimal for HMSSA and VMSSA, respectively. In these formulas, N is the length of the time series and p is the number of variables in the time series [4, 5].

3- Results

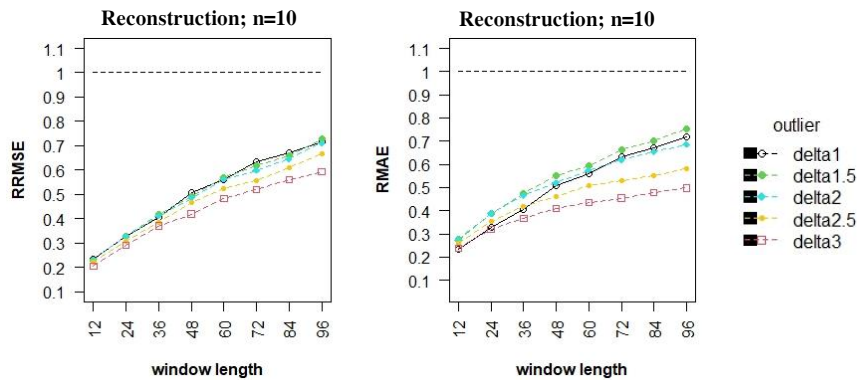
In this section, the results of the paper are presented. To generate outliers, first, n observations ($n=5,10$) are randomly selected from the simulated time series and then a multiplication (δ) of them is substituted for each selected observation. In all examples, the values of δ are equal to 1, 1.5, 2, 2.5, and 3. In this section, we present an example to declare the proposed method. As an example, we consider the following Sine series:

$$y_1 = 3 \sin\left(\frac{2t\pi}{12}\right) + \varepsilon$$

$$y_2 = 2 \sin\left(\frac{2t\pi}{12} + \frac{\pi}{4}\right) + \varepsilon$$

$t=1, 2, \dots, 200,$

where $\varepsilon \sim N(0, 1)$ is the noise component. The first 190 observations were considered as in-sample (reconstruction) and the rest as out-of-sample series (forecasting). Since the rank of the trajectory matrix for this model (HMSSA) is equal to 2 we choose $r=2$. Using the recommendations given in Sec.2, we take $L=12, 24, 36, 48, 60, 72, 84,$ and 96 . In figures 1 and 2, the plots of RRMSE and RMAE for reconstruction are depicted for $n=5$ and 10 . The plots of RRMSE and RMAE show that L_1 -MSSA is better than basic MSSA for all values of L and δ ($RRMSE < 1$ and $RMAE < 1$).



Figs. 1- The plots RRMSE and RMAE for reconstruction and sample size 10

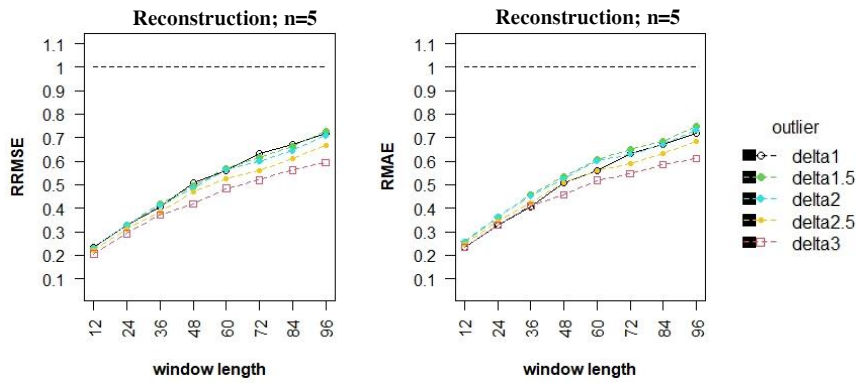


Fig. 2- The plots RRMSE and RMAE for reconstruction and sample size 5

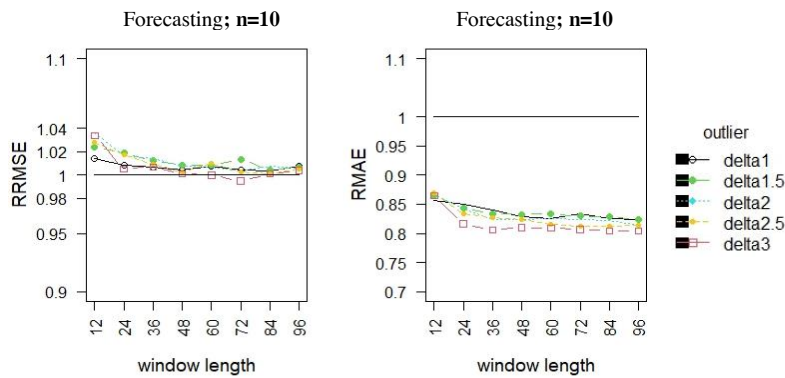


Fig. 3- The plots RRMSE and RMAE for forecasting and sample size 10

In figures 3 and 4, the plots of RRMSE and RMAE for forecasting out of samples are depicted for $n = 5, 10$. The plot of RRMSE shows that L_1 -MSSA is better than basic MSSA for $L = 72$ and $\delta = 3$.

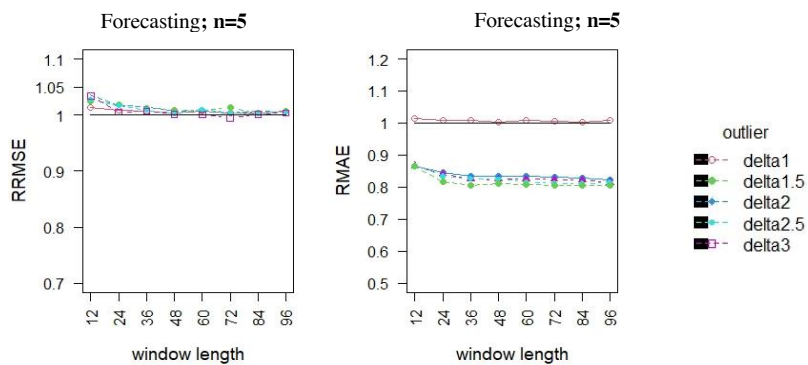


Fig.4- The plots RRMSE and RMAE for forecasting and sample size 5

For $n = 5$ and 10 , according to RRMSE, by increasing δ , the efficiency of $L_1 - MSSA$ tends to increase. But based on RMAE, the performance of the L_1 -MSSA method is better than the basic MSSA for all values of n , L , and δ . When $\delta = 1$ (i.e. no outliers), the performance of both methods (L_1 -MSSA and basic MSSA) are the same in forecasting.

4- Discussion and Conclusion

In this paper, to optimize the quality of modeling and forecasting obtained from the basic MSSA with outliers, a robust method was proposed by replacing Frobenius norm with the $L_1 - norm$. This method is a new way of using robust time series to improve modeling. In addition to optimizing the quality of the proposed model, this technique increases the efficiency of time series analysis and forecasting. In order to increase the reconstruction and forecasting performance of the SSA method, identifying and removing outliers have been suggested by some researchers, because SSA based on the L_2 norm is more sensitive to outliers [2]. The results of this paper show that in the case of encountering time series with outliers, better results can be obtained by changing the matrix norm in MSSA from L_2 to L_1 because there is no need to transform the data and it ensures that no information is lost. This study considered both RMSE and MAE criteria. In brief, the comparison of forecasting results showed that L_1 -MSSA is more accurate than the basic MSSA version, further confirming the results obtained via theoretical results.

5- Reference

- [1] Hassani, H., Webster, A. E., Silva, S. and Heravi, S. (2015). Forecasting U.S. tourist arrivals using optimal singular spectrum analysis, *Tour. Manage.* 46, 322–335.
- [2] Hassani, H. and Mahmoudvand. R. (2013). Multivariate singular spectrum analysis: A general view and new vector forecasting approach, *Int. J. Energy Stat.* 1(1), 55–83.
- [3] Hassani, H. (2007). Singular spectrum analysis: Methodology and comparison, *J. Data Sci.* 5(2), 239–257.
- [4] Mahmoudvand, R. and Rodrigues, P.C. (2016). Missing value imputation in time series using Singular Spectrum Analysis. *International Journal of Energy and Statistics* 4. DOI: 10.1142/S2335680416500058
- [5] Mahmoudvand, R. (2012). Some Theoretical Development of the Singular Spectrum Analysis. Ph.D. Thesis (In Persian), Shahid Beheshti University, Tehran, Iran.
- [6] Patterson, K., Hassani, H., Heravi, S., and Zhigljavsky, A. (2011). Forecasting the final vintage of the industrial production series, *Journal of Applied Statistics*, 38, 2183–2211.

