

شناسایی نقطه‌ی تغییر در سری‌های زمانی با استفاده از روش تحلیل مجموعه مقادیر تکین

مسعود یار محمدی^۱، مهدی کلانتری^۲

^۱گروه آمار، دانشگاه پیام نور، صندوق پستی ۴۶۹۷-۱۹۳۹۵، تهران، ایران

چکیده در تحلیل سری‌های زمانی اگر مشاهدات جمع آوری شده شامل بخش‌هایی با توزیع‌های متفاوت باشد، یعنی نقاطی وجود داشته باشد که وضعیت ظاهری سری به دلیل مختلف از جمله تغییر در توزیع داده‌ها دستخوش تغییر شود، به این نقاط از دیدگاه آماری نقاط تغییر گویند. بنابراین در نقاط تغییر، داده‌ها به بخش‌های همگن و متجانس مجزا تقسیم می‌گردند. در عمل تعداد نقاط تغییر و محل آنها نامعلوم بوده و توانایی کشف و شناسایی آنها از حیثگاه کاربردی به ویژه مدل بندی و پیش بینی سری زمانی از اهمیت به سزایی برخوردار می‌باشد. در این مقاله نهمی شناسایی نقاط تغییر در رگرسیون‌های سری زمانی یا استفاده از یک روش ناپارامتری به نام تحلیل مجموعه‌ی مقادیر تکین مورد بررسی قرار می‌گیرد. کارایی این روش با توجه به عدم نیاز به هیچگونه فرض اولیه در مورد سری زمانی، در مقایسه با سایر روش‌های متداول لایل توجه می‌باشد.

واژه‌های کلیدی: سری زمانی، نقطه‌ی تغییر، تحلیل مجموعه مقادیر تکین.

کد موضوعی: ریاضی (۲۰۱۰): 91B84, 62M10, 37M10

۱ مقدمه

فرآیندها غالباً طی یک مدت زمان نسبتاً طولانی در حالت پایدار و بدون تغییر خود به سر می‌برند. بدیهی است هیچ فرآیندی به طور دائمی پایدار نیست و تدریجاً تغییرات با دلیل (معمولاً به صورت اتفاقی) ظاهر شده و باعث می‌شوند تا فرآیند با خواستهای مورد نظر انطباق نداشته باشد. این نقاط را که باعث ایجاد تغییر در ساختار توزیع داده‌ها در یک فرآیند تصادفی یا سری زمانی می‌شود را نقطه‌ی تغییر می‌نامند. نقاط تغییر به دلیل مختلفی از جمله تغییر در مکانیسم تولید داده‌ها (توزیع احتمال داده‌ها) رخ می‌دهند. روش‌های شناسایی نقطه‌ی تغییر یکی از مباحث مهم در تحلیل سری‌های زمانی است و در سال‌های اخیر مورد توجه بسیاری از پژوهشگران قرار گرفته است. تحلیل نقطه‌ی تغییر در زمینه‌های مختلفی نظیر کنترل فرآیند، بیچ (۱۹۵۲)، بررسی‌های هواشناسی، واتارد و جیل (۱۹۸۹)، اقتصاد،

^۱مسعود یار محمدی، masyar@pnu.ac.ir

چن و گویتا (۱۹۹۲) و تحلیل رگردهای نوار مغزی، پرودسکی و همکاران (۱۹۹۹) ، استفاده شده است.

چرنوف و زاکس (۱۹۶۳) یک آزمون بیزی برای تشخیص تغییرات میانگین در مشاهدات نرمال و شناسایی قطعی تغییر به کار بردند. روش‌های بیزی متعددی توسط چتون و کیم (۲۰۱۰) ، کیم و چتون (۲۰۱۰) برای شناسایی نقاط تغییر در مدل‌های مختلف از جمله نرمال، دوجمله‌ای، نمایی و پواسن به کار گرفته شده است. اولین روش ناپارامتری برای تشخیص قطعی تغییر توسط پاتاچاریا و فریرسون (۱۹۸۱) ارائه شده است. روش آنها مشابه روش چرنوف و زاکس (۱۹۶۳) می‌باشد، با این تفاوت که از رتبه‌ی داده‌ها به جای داده‌های اصلی استفاده کرده‌اند. در رابطه با شناسایی قطعی تغییر می‌توان گفت روش‌های ناپارامتری نسبت به روش‌های پارامتری به دلیل آنکه نیازمند هیچ گونه فرض اولیه در مورد داده‌ها نیستند، از کارایی خاصی برخوردارند. مهمترین روش در این زمینه بر اساس روش تحلیل مجموعه مقادیر تکین^۱ (SSA) است.

SSA بر خلاف مدل‌های کلاسیک باگس- چتکیت روشی ناپارامتری در حوزه تحلیل سری‌های زمانی است و نیازی به برقراری شرط مانایی سری زمانی و نرمال بودن خطاها ندارد. از طرف دیگر، کم بودن تعداد مشاهدات محدودیت جدی برای آن ایجاد نمی‌کند. به دلیل دارا بودن چنین مزیت‌هایی، این روش کاربردهای وسیعی در بسیاری از شاخه‌های علوم یافته است. SSA سعی در تجزیه سری زمانی به اجزایی نظیر روند، مولفه‌های فصلی (با دوره تناوب‌های مختلف) و خطای تصادفی (نویز) دارد. در بخش دوم، این روش به اختصار معرفی خواهد شد. برای کسب اطلاعات بیشتر در مورد SSA به گولیانینا و همکاران (۲۰۰۱) ، گولیانینا و ژینگنجاوسکی (۲۰۱۳) و حسنی (۲۰۰۷) مراجعه کنید. در بخش سوم، الگوریتمی برای شناسایی نقاط تغییر در سری زمانی ارائه خواهد شد که مبتنی بر روش SSA است و نیازمند هیچ گونه فرض اولیه در مورد مشاهدات سری زمانی نیست. در بخش چهارم نیز نحوه‌ی به کارگیری الگوریتم ارائه شده را به کمک یک سری زمانی شبیه سازی شده نشان خواهیم داد.

۲ تحلیل مجموعه مقادیر تکین

روش SSA از دو مرحله تشکیل شده است: تجزیه و بازسازی. هر کدام از این مراحل نیز شامل دو گام هستند. فرض کنید x_1, x_2, \dots, x_N یک سری زمانی مشاهده شده به طول N بوده و M نیز عددی صحیح، که طول پنجره نامیده می‌شود، باشد به طوری که $1 < M < N$. مراحل SSA عبارتند از:

۱.۲ تجزیه

این مرحله شامل دو گام است: نشانیدن^۲ و تجزیه مقدار تکین^۳ (SVD).

۱.۱.۲ نشانیدن

در این گام، ابتدا سری زمانی را به K زیر سری تبدیل می‌کنیم، وقتی که $K = N - M + 1$. زیر سری k ام به ازای $k = 1, \dots, K$ به صورت $X_k = (x_k, \dots, x_{k+M-1})^T$ تعریف می‌شود. حقت کنید که زیر سری k ام یک بردار ستونی با M مولفه است و گاهی اوقات بردار M -تأخیری نامیده می‌شود. سپس این بردارهای M -تأخیری را به صورت ستونی در کنار هم می‌نشانیم تا تشکیل یک

^۱ Singular Spectrum Analysis (SSA)

^۲ Embedding

^۳ Singular Value Decomposition (SVD)

ماتریس $M \times K$ دهنده، این ماتریس که با X نشان داده می‌شود، ماتریس مسیر^۱ می‌نامیم. داریم:

$$X = [X_1 : \dots : X_K] = (x_{ij})_{i,j=1}^{M,K} = \begin{pmatrix} x_{11} & x_{12} & x_{13} & \dots & x_{1K} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2K} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{M1} & x_{M2} & x_{M3} & \dots & x_{MK} \end{pmatrix}$$

دقت کنید که عناصر روی قطرهای قرصی ماتریس X با هم برابرند. چنین ماتریسی را ماتریس هنگل^۲ می‌نامند. از طرفی واضح است که با در اختیار داشتن ماتریس X ، می‌توان سری زمانی را به دست آورد. در واقع یا کنار هم قرار دادن ستون اول و سطر آخر (یا سطر اول و ستون آخر) ماتریس X ، سری زمانی حاصل می‌شود. بنابراین یک تناظر یک به یک بین ماتریس مسیر (که هنگل نیز می‌باشد) و سری زمانی وجود دارد.

۲.۱.۲ تجزیه مقدار تکین

در این گام، تجزیه مقدار تکین ماتریس مسیر را به دست می‌آوریم. فرض کنید $\lambda_1, \dots, \lambda_m$ مقادیر ویژه ماتریس XX^T بوده که به صورت نزولی (از بزرگ به کوچک) مرتب شده‌اند ($\lambda_1 \geq \dots \geq \lambda_m \geq 0$) و U_1, \dots, U_m نیز بردارهای یکا متعامد ویژه متناظر با مقادیر ویژه $\lambda_1, \dots, \lambda_m$ باشند. تعریف می‌کنیم $d = \max\{i, \lambda_i > 0\}$ (تعداد مقادیر ویژه مثبت) و $V_i = X^T U_i / \sqrt{\lambda_i}$ در این صورت SVD ماتریس X به صورت $X = X_1 + \dots + X_d$ نوشته می‌شود وقتی که $X_i = \sqrt{\lambda_i} U_i V_i^T$ به ازای $i = 1, \dots, d$ سه تایی $(\sqrt{\lambda_i}, U_i, V_i)$ را سه تایی ویژه^۳ (ET) می‌نامند.

۲.۲ بازسازی

این مرحله شامل دو گام است: گروه بندی و میانگین‌گیری قطری.

۱.۲.۲ گروه بندی

هدف از اجرای این مرحله یافتن سه تایی‌های ویژه مربوط به هر یک از اجزای سری زمانی نظیر روند، مولفه‌های فصلی، توفه و غیره است. در این گام، پس از اینکه SVD ماتریس مسیر به دست آمد، مجموعه‌ی انجیس‌های $\{1, \dots, d\}$ را به m زیر مجموعه‌ی I_1, \dots, I_m افزایش می‌دهیم. در واقع این امر را کردن معادل با این است که اجزای سری زمانی را به کمک سه تایی‌های ویژه تشخیص داده و از هم تفکیک کنیم. فرض کنید $I = \{1, \dots, d\}$. در این صورت ماتریس X_I متناظر با گروه I به صورت $X_I = X_{I_1} + \dots + X_{I_m}$ تعریف می‌شود. مثلاً اگر $I = \{1, 2, 7\}$ آنگاه $X_I = X_1 + X_2 + X_7$. بدین ترتیب می‌توان X_I را به ازای $I = I_1, \dots, I_m$ به دست آورد. در این صورت از SVD ماتریس X نتیجه می‌گیریم که $X = X_{I_1} + \dots + X_{I_m}$. شیوه‌ی انتخاب مجموعه‌های I_1, \dots, I_m ، گروه بندی سه تایی ویژه نامیده می‌شود.

^۱ Trajectory Matrix

^۲ Hankel Matrix

^۳ Eigen Triple

۲.۲.۲ میانگین‌گیری قطری

هدف اصلی در این گام، تبدیل هر ماتریس X_{T_i} ، $i = 1, \dots, m$ ، از گام گروه بندی به یک سری زمانی به طول N است. همان طور که در گام نشانیدن اشاره شد، با در اختیار داشتن یک ماتریس متکل می‌توان سری زمانی متناظر با آن را به دست آورد؛ ولی ماتریس‌های X_{T_i} که در مرحله‌ی گروه بندی به دست می‌آیند دارای خاصیت متکلی نیستند. متکل سازی ماتریس X_{T_i} به وسیله‌ی میانگین‌گیری روی عناصر قطرهای فرعی انجام می‌شود بدین معنی که همه‌ی عناصر روی یک قطر فرعی را با میانگین عناصر همین قطر فرعی جایگزین می‌کنیم. بدین ترتیب ماتریس X_{T_i} به یک ماتریس متکل تبدیل می‌شود. پس از این تبدیل، زیر سری بازسازی شده‌ی $\{\tilde{x}_1^{(i)}, \dots, \tilde{x}_N^{(i)}\}$ به طول N به دست می‌آید. بنابراین با توجه به SVD ماتریس X ، سری زمانی اصلی به صورت مجموع m زیر سری بازسازی شده به شکل زیر بازسازی می‌شود:

$$x_t = \sum_{i=1}^m \tilde{x}_t^{(i)}, \quad t = 1, 2, \dots, N.$$

۳ تشخیص نقطه تغییر بر اساس SSA

موسکونا و زیگلنجاوسکی (۲۰۰۲) الگوریتمی را برای تشخیص نقطه تغییر ارائه دادند که مبتنی بر روش SSA است. این الگوریتم شامل سه گام است. فرض کنید x_1, x_2, \dots یک سری زمانی بوده و N, M, l, p و q اعداد صحیحی باشند به طوری که $0 \leq p < q$ و $l < M \leq N/2$ به ازای هر $n = 0, 1, \dots$ گام‌های زیر را اجرا می‌کنیم:

گام ۱ در این گام سه مرحله‌ی اول روش SSA، یعنی نشانیدن، SVD و گروه بندی را روی سری زمانی در بازه‌ی $[n+1, n+N]$ اجرا می‌کنیم. به عبارت دیگر مراحل زیر را به ترتیب طی می‌کنیم:

۱- یک ماتریس مسیر به نام ماتریس پایه به صورت زیر می‌سازیم:

$$X_B^{(n)} = \begin{pmatrix} x_{n+1} & x_{n+2} & x_{n+3} & \dots & x_{n+K} \\ x_{n+2} & x_{n+3} & x_{n+4} & \dots & x_{n+K+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n+M} & x_{n+M+1} & x_{n+M+2} & \dots & x_{n+N} \end{pmatrix}$$

وقتی که $K = N - M + 1$ ، توجه کنید که ستون‌های ماتریس پایه^۱ را بردارهای $x_j^{(n)} = (x_{n+j}, \dots, x_{n+j+M-1})^T$ و $j = 1, \dots, K$ تشکیل می‌دهند.

۲- SVD ماتریس $R_n = X_B^{(n)}(X_B^{(n)})^T$ را به دست می‌آوریم. در این مرحله M بردار ویژه حاصل می‌شود.

۳- اولین l بردار ویژه از M بردار ویژه‌ی مرحله قبل را انتخاب می‌کنیم.

^۱Basic Matrix

گام ۲ در این گام یک ماتریس به نام ماتریس آزمون^۸ با بعد $M \times Q$ می‌سازیم. داریم:

$$X_j^{(n)} = \begin{pmatrix} x_{n+p+1} & x_{n+p+2} & x_{n+p+3} & \dots & x_{n+q} \\ x_{n+p+2} & x_{n+p+3} & x_{n+p+4} & \dots & x_{n+q+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n+p+M} & x_{n+p+M+1} & x_{n+p+M+2} & \dots & x_{n+q+M-1} \end{pmatrix}$$

وقتی که $q = p + Q$. توجه کنید که ستون‌های این ماتریس را بردارهای Q و $j = p + 1, \dots, p + Q$ تشکیل می‌دهند.

گام ۳ در این گام آمارهای مربوط به تشخیص قطعی تغییر محاسبه می‌شوند. این آمارها عبارتند از:

$$D_{n,l,p,q} = \sum_{j=p+1}^l ((X_j^{(n)})^T X_j^{(n)} - (X_j^{(n)})^T U U^T X_j^{(n)}) \quad \bullet$$

ستون‌های آن را بردارهای ویژه‌ای که در مرحله سوم از گام اول انتخاب کردیم، تشکیل می‌دهند.

$$S_n = \bar{D}_{n,l,p,q} / \mu_{n,l} \quad \bullet$$

وقتی که $S_n = \bar{D}_{n,l,p,q} / \mu_{n,l}$ برآوردگر $\bar{D}_{n,l,p,q}$ در بازه زمانی $[n + m]$

جمع انباشته (CUSUM)

$$W_1 = S_1, \quad W_{n+1} = \max(0, W_n + S_{n+1} - S_n - k/\sqrt{MQ}), \quad n \geq 1$$

k یک مقدار ثابت نامنفی کوچکی است و موسکونتا (۲۰۰۱) مقدار آن را برابر $1/3\sqrt{MQ}$ پیشنهاد کرده است.

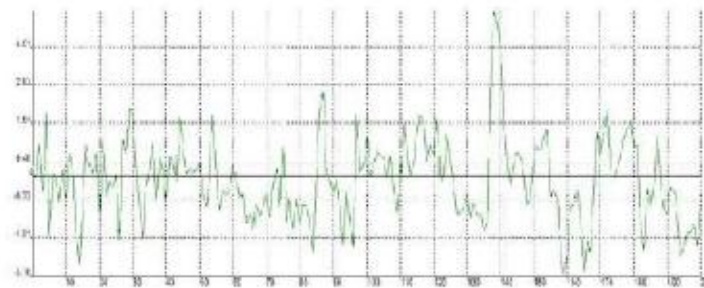
مقادیر بزرگ $D_{n,l,p,q}$ ، S_n و W_n نشان دهنده‌ی یک تغییر در ساختار سری زمانی است. این الگوریتم یک تغییر در زمان n را نشان می‌دهد هر گاه $W_n > h$ وقتی که $h = \frac{z_{\alpha}}{\sqrt{MQ}} \sqrt{MQ(3MQ - Q^2 + 1)}$ و z_{α} چننگ $(1 - \alpha)$ توزیع نرمال استاندارد است. برای کسب اطلاعات بیشتر در مورد نحوه‌ی تعیین پارامترهای M ، l ، p و q به موسکونتا و زیگلجاوسکی (۲۰۰۳) مراجعه کنید.

۴ شبیه سازی

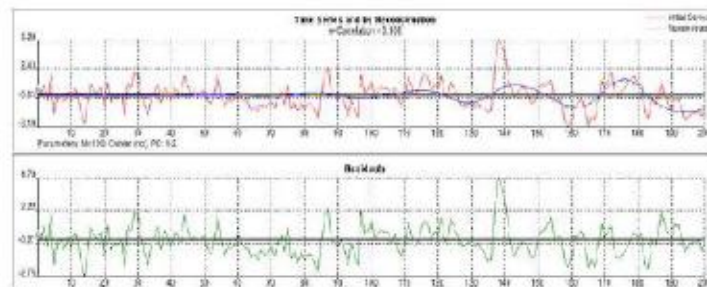
در این بخش الگوریتم ارائه شده را روی یک مثال شبیه سازی شده پیاده خواهیم کرد. ۲۰۰ مشاهده از فرآیند $AR(1)$ تولید شده است ($N = 200$). ۱۰۰ مشاهده‌ی اول با پارامتر $\phi = 0.7$ و ۱۰۰ مشاهده‌ی بعدی با پارامتر $\phi = 0.7$. توزیع توفه‌ها هم نرمال استاندارد در نظر گرفته شده است. بدین ترتیب انتظار داریم که در زمان ۱۰۱ یک تغییر رخ دهد. شکل ۱ نمودار سری زمانی داده‌های شبیه سازی شده را نشان می‌دهد.

پارامتر طول پنجره برابر $M = \frac{N}{4} = 100$ در نظر گرفته شده است. پس از اجرای مرحله‌ی SVD و به دست آوردن مقادیر ویژه و بردارهای ویژه، دو بردار ویژه اول برای بازسازی سری زمانی در نظر گرفته شد. شکل ۲ نمودار سری زمانی بازسازی شده به وسیله‌ی دو بردار ویژه اول به همراه مانده‌ها را نشان می‌دهد. با انتخاب مقادیر $l = 2$ ، $p = 27$ و $q = 75$ مقدار h برابر با 0.7334 به دست می‌آید. در شکل ۳ نمودار آماری تغییر نشان داده شده است. این شکل به وضوح وجود یک تغییر در زمان ۱۰۱ را نشان می‌دهد.

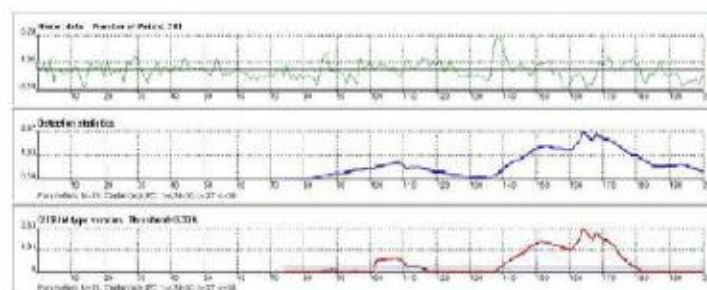
^۸Test Matrix



شکل ۱: نمودار سری زمانی داده‌های شبیه‌سازی شده



شکل ۲: نمودار سری زمانی بازسازی شده



شکل ۳: نمودار آماری تغییر

بحث و نتیجه‌گیری

شناسایی نقاط تغییر در سری‌های زمانی از دیدگاه مدل بندی و پیش بینی از اهمیت ویژه‌ای برخوردار است. در این مقاله کاربرد روش تحلیل مجموعه مقادیر تکین جهت شناسایی نقاط تغییر بیان گردید. در ادامه کارایی این روش برای شناسایی یک نقطه تغییر مشخص، در یک سری زمانی شبیه‌سازی شده مورد بحث و بررسی قرار گرفت. عدم نیاز این روش به فرض‌های حائلی سری زمانی، تعداد مشاهدات و نرمال بودن خطاها را می‌توان به عنوان مزیت‌های این روش محسوب نمود.

- Bhattacharya P. K. and Frierson J. R. (1981), A Nonparametric Control Chart for Detecting Small Disorders, *The Annals of Statistics*, 9, 3, 544-554.
- Brodsky B.E., Durkbovsky, B.S., Kaplan A.Ya., Shishkin S.L. (1999), A nonparametric method for the segmentation of the EEG, *Computer Methods and Programs in Biomedicine*, 60, 93-106.
- Chen J. and Gupta A.K. (1997), Testing and locating variance change points with application to stock prices, *Journal of the American Statistical Association*, 92, 739-747.
- Cheon S. and Kim J. (2010), Multiple change-point detection of multivariate mean vectors with the bayesian approach. *Computational Statistics and Data Analysis*, 54, 406-415.
- Chernoff H. and Zacks S. (1964), Estimating the current mean of a normal distribution which is subject to changes in time. *Annals of Mathematical Statistics*, 35, 999-1018.
- Golyandina N., Nekrutkin V., Zhigljavsky A., (2001), *Analysis of time series structure: SSA and related techniques*, Boca Raton, Chapman and Hall/CRC.
- Golyandina N. and Zhigljavsky A. (2013), *Singular Spectrum Analysis for Time Series*, London, Springer.
- Hassani H. (2007), Singular Spectrum Analysis: Methodology and Comparison, *Journal of Data Science*, 5(2), 239-257.
- Kim J. and Cheon S. (2010), Bayesian multiple change-point estimation with annealing stochastic approximation monte carlo, *Computational Statistics* 25, 215-239.
- Moakvina V. and Zhigljavsky A. (2003), An algorithm based on singular spectrum analysis for change-point detection, *Communication in Statistics: Simulation and Computation*, 32(2), 319-352.
- Moakvina V. (2001), Application of the Singular-Spectrum Analysis for Change-Point Detection in Time Series. Ph.D. thesis, School of Mathematics, Cardiff: Cardiff University.
- Page E.S. (1954), Continuous inspection scheme, *Biometrika*, 1, 100-115.
- Vautard R. and Ghil M. (1989), Singular spectrum analysis in nonlinear dynamics, with applications to paleoclimatic time series, *Physica D*, 35, 395-424.