

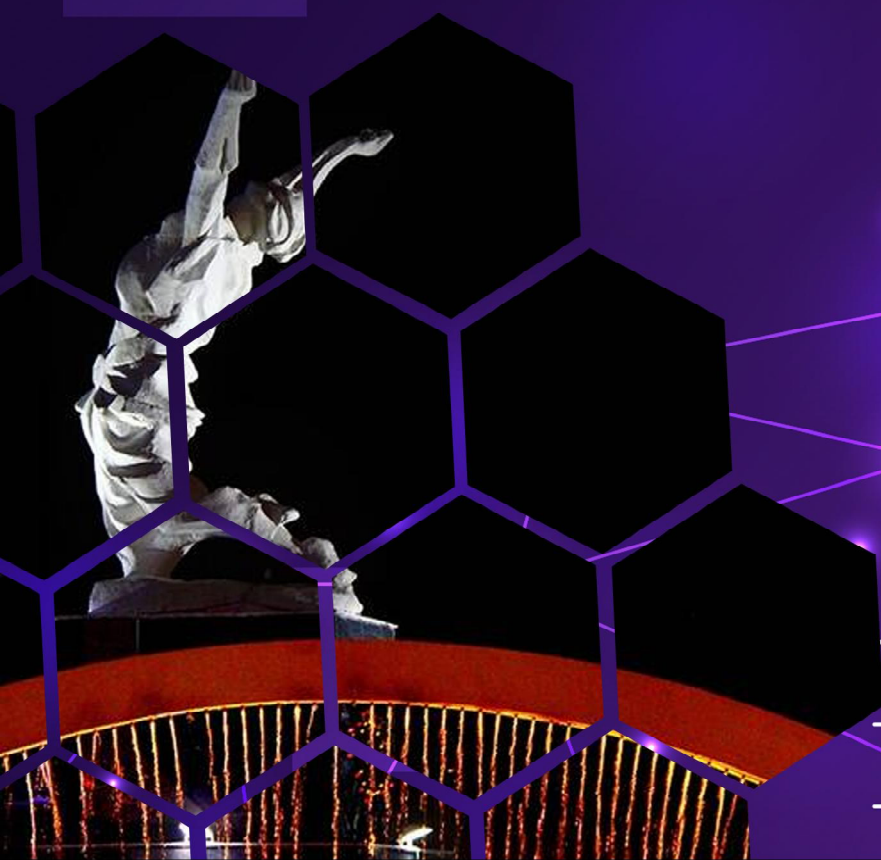
بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

مجموعه مقالات پانزدهمین سمینار

احتمال و فرآیندهای تصادفی

● ۸ و ۹ شهریور ۱۴۰۴

● دانشگاه کردستان



مجموعه مقالات پانزدهمین سمینار احتمال و فرایندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴ دانشگاه کردستان

سنندج

شناسنامه کتاب

عنوان: مجموعه مقالات پانزدهمین سمینار احتمال و فرایندهای تصادفی

موضوع: آمار، احتمال، فرایندهای تصادفی

گردآوری: دکتر سیروس فتحی منش، دکتر هادی امامی، دکتر هادی احمدی

تاریخ برگزاری: هشتم و نهم شهریورماه ۱۴۰۴

محل برگزاری: دانشگاه کردستان، سنندج

به نام خداوند جان و خرد

پیشگفتار

خداوند متعال را سپاسگزاریم که توفیق برگزاری پانزدهمین سمینار احتمال و فرایندهای تصادفی را در دانشگاه کردستان و با همکاری انجمن آمار ایران، به منظور خدمت به جامعه علمی کشور، نصیب ما نمود. این سمینار با هدف گردهمایی صاحب نظران و پژوهشگران حوزه احتمال و فرایندهای تصادفی برگزار شد و فرصت مغتنمی برای تبادل نظر و ارائه تحقیقات و دستاوردهای نوین این عرصه فراهم آورد.

مجموعه حاضر، شامل مقالات فارسی پذیرفته شده در این سمینار است که امیدواریم گشاینده افق‌های تازه‌ای پیش روی پژوهشگران، به ویژه دانشجویان تحصیلات تکمیلی، باشد. در این دوره از سمینار، در مجموع ۹۸ مقاله به زبان‌های فارسی و انگلیسی به دبیرخانه ارسال گردید که پس از طی فرایند داوری، ۷۵ مقاله برای ارائه و چاپ به صورت کامل یا خلاصه پذیرفته شد. در فرایند ارزیابی مقالات، همکاران ارجمندی از دانشگاه‌های سراسر کشور ما را یاری نمودند. بیشترین تعداد مقالات پذیرفته شده با ۱۵ مقاله، متعلق به محور سری‌های زمانی و کاربردها بود.

علاوه بر ارائه مقالات، هفت سخنرانی کلیدی توسط اساتید مدعو، شش نشست تخصصی در زمینه‌های احتمال، فرایندهای تصادفی، سری‌های زمانی، نظریه توزیع‌ها، آنالیز تصادفی، ریاضیات مالی، استنباط آماری برای فرایندهای تصادفی و همچنین کاربرد و نقش فرایندهای تصادفی در هوش مصنوعی و علم داده، و یک کارگاه آموزشی در حوزه ریاضیات مالی برگزار گردید. مقالات مندرج در این مجموعه، بر اساس محورهای اصلی سمینار تدوین شده‌اند. این مقالات در پایگاه استنادی علوم جهان اسلام (ISC) و پایگاه سیویلیکا ثبت شده‌اند. امیدواریم این اثر، گامی هرچند کوچک در جهت پیشبرد مرزهای دانش در حوزه تخصصی خود باشد و مورد توجه و استفاده جامعه علمی کشور قرار گیرد. از خوانندگان گرامی دعوت می‌شود تا با مطالعه این مقالات، از یافته‌های نوین پژوهشگران این عرصه آگاه شوند.

برگزاری این رویداد علمی بدون مساعدت و راهنمایی‌های ارزشمند همکاران محترم در دانشگاه‌های مختلف کشور میسر نبود. همکاری این عزیزان در ارسال و داوری مقالات، ارائه راهکارهای سودمند و مشارکت فعال، سهمی بسزا در موفقیت سمینار داشت. همچنین از پشتیبانی و زحمات بی‌دریغ ریاست محترم دانشگاه و رئیس محترم دانشکده علوم پایه دانشگاه کردستان، ریاست محترم انجمن آمار ایران، واحد فناوری اطلاعات دانشگاه کردستان و کلیه

کارکنان و دانشجویانی که در برگزاری هرچه بهتر این سمینار یاری گر ما بودند، صمیمانه سپاسگزاریم. در پایان، حمایت‌های مالی سازمان برنامه و بودجه استان کردستان که مشوقی ارزشمند در برگزاری این سمینار بود، را ارج می‌نهم.

کمیته‌های علمی و اجرایی پانزدهمین سمینار احتمال و فرایندهای تصادفی

فهرست

۱ مجریان و حامیان
۲ اعضای کمیته اجرایی
۳ اعضای کمیته علمی
۴ هیأت داوران
۵ مقالات ارائه شده

مجریان و حامیان :



دانشگاه کردستان



انجمن آمار ایران



سازمان مدیریت و برنامه ریزی
استان کردستان

اعضای کمیته اجرایی سمینار:

- دکتر هادی احمدی
- دکتر هادی امامی (دبیر اجرایی سمینار)
- دکتر حسین بیورانی
- خانم سمیه جعفر رمشتی
- دکتر کوروش دادخواه
- دکتر شاهو زارعی
- خانم سمیه زند
- دکتر سیروس فتحی منش
- دکتر اقبال قادری

اعضای کمیته علمی سمینار :

- دکتر محمد امینی (دانشگاه فردوسی مشهد)
- دکتر حسین بیورانی (دانشگاه کردستان)
- دکتر امید خوارزمی (دانشگاه ولی عصر رفسنجان)
- دکتر زهرا رضایی قهرودی (دانشگاه تهران)
- دکتر شاهوزارعی (دانشگاه کردستان)
- دکتر احمد رضا سلطانی (دانشگاه شیراز)
- دکتر سیروس فتحی منش (دانشگاه کردستان - دبیر علمی سمینار)
- دکتر علیرضا نعمت الهی (دانشگاه شیراز)

اسامی داوران :

دانشگاه اصفهان	دکتر بهاره اختری
دانشگاه کردستان	دکتر هادی امامی
دانشگاه فردوسی	دکتر محمد امینی
دانشگاه بو علی	دکتر ابراهیم امینی سرشت
دانشگاه کردستان	دکتر سلمان ایزدخواه
دانشگاه رازی	دکتر محی الدین ایزدی
دانشگاه کردستان	دکتر حسین بیورانی
دانشگاه اصفهان	دکتر افشین پرورده
دانشگاه رازی	دکتر عبدالله جلیلیان
دانشگاه زنجان	دکتر امید خادم نوع
دانشگاه تربیت مدرس	دکتر مجید خالدي
دانشگاه ولی عصر (عج) رفسنجان	دکتر امید خوارزمی
دانشگاه شیراز	دکتر کاووس خورشیدیان
دانشگاه تهران	دکتر زهرا رضایی قهرودی
دانشگاه کردستان	دکتر شاهره زارعی
دانشگاه صنعتی شریف	دکتر شیوا زمانی
دانشگاه امیرکبیر	دکتر عرفان صلواتی
دانشگاه کردستان	دکتر محمد فتحی
دانشگاه کردستان	دکتر سیروس فتحی منش
دانشگاه پیام نور زنجان	دکتر مهدی کلانتری
دانشگاه مازندران	دکتر مهرناز محمدپور
دانشگاه زنجان	دکتر علی محمدیان مصمم
دانشگاه بوعلی	دکتر رحیم محمودوند

مقالات ارائه شده

خوشه‌بندی دوطرفه فازی با استفاده از تابع فاصله اطلاعاتی،

احمدی، ز. و زارعی، ر. ۹

یک روش تنک‌سازی فرایندهای نقطه‌ای مبتنی بر مجموعه‌های تصادفی بولی،

اسدی، ر. و خزایی، م. و گنجعلی، م. ۱۸

آزمون نیکویی برازش تقارن توزیع‌های پیوسته تحت داده‌های از چپ بریده شده و از راست سانسور شده،

اکبری، م. ۳۱

یک کلاس از مدل‌بندی سری‌زمانی دوخطی گسسته‌مقدار با کاربردهای آن،

بامدادی، ر. و محمودپور، م. ۳۹

مدل بندی قابلیت اطمینان سیستم های تنش-مقاومت بر اساس رویکرد تابع مولد عام،

تصدیقی، ب. و زارعی، ر. و فتحی، ب. ۴۸

اخلاق در آمار و احتمال،

دولتی، ع. ۵۷

نقش فرآیندهای تصادفی در هوش مصنوعی و علم داده: یک مطالعه موردی،

راستین، ا. ۷۰

فرایند INAR(۱)-PJ: یک جایگزین جدید فرایند INAR (۱) پواسونی،

۷۸ پاسخ، م

بررسی انشعاب‌های تصادفی مدل رشد لجستیک جمعیت،

۸۸ مهموئی، س. و ربیعی مطلق، ا. و محمدی نژاد، ح. م.

تشخیص اعداد دست‌نویس با رهیافت یادگیری ماشین،

۹۶ رضایی طالقانی، ن. و حاج‌رجبی، آ.

مقایسه دو رویکرد نیم‌رخ و نامقید در برازش مدل رگرسیون بر اساس نمونه داده‌های ترکیبی با تواتر متفاوت،

۱۰۴ رنجبر، ف. و آقامحمدی، ع. و ادیب، م.

برخی نتایج مجانبی در مورد اختلاف میانگین دو جامعه در داده‌های تابعی جزئی مشاهده‌شده،

۱۱۲ سلیمی فر، پ. و حسینی نسب، س. م. و خادم نوع، ا.

مدل خودبازگشتی صحیح مقدار مرتبه اول فصلی با توزیع حاشیه‌ای دلاپورت،

۱۱۹ شالباف، م. و پرهام، غ.

رویکرد فازی خوشه‌بندی داده‌ها بر پایه الگوریتم،

۱۲۸ میاهی، م. و طارمی، ب. و میاهی، م.

بررسی رفتار مجانبی رهیافت PTE در حضور هم‌خطی چندگانه،

۱۳۴ طباطبایی شیرازی، س. ا. ح. و عمادی، م. و آرشی، م. و سیف‌اللهی، س.

چگونه نظریه فرآیندهای تصادفی به مدل گسترش یافته در سنتز مدرن تکامل کمک کرد؟

عریضی، ح. ۱۴۳

بررسی معیار Covar بر اساس رویکرد COUPLE-ARCH و توزیع GED،

علیزاده، ف. و امینی، م. و محتشمی برزادران، غ. ۱۵۳

مروری بر استفاده از خانواده توزیع‌های فاز-نوع در مدل شوک‌های تعمیم یافته،

فتحی‌منش، س. و ایزدی، م. ۱۶۲

بررسی عملکرد مجانبی برآوردگر ترکیبی در مدل‌های خطی با خطا در اندازه‌گیری،

قیانی، ف. ۱۷۴

مقایسه روش‌های کلاسیک و بیزی در برآورد پارامتر قابلیت اعتماد تنش-مقاومت برای توزیع گمپتر تعمیم یافته یکه،

کاراندیش مروتی، ا. و ارمز، ا. و بصیرت، م. ۱۸۲

کاربرد روش استوار مونت کارلوی بهینه در استنباط بیزی بدون درست‌نمایی برای داده‌های وابسته،

کریمی، ا. و حسینی، ف. ۱۹۲

مفصل‌های دایره‌ای-خطی: ویژگی‌ها و کاربردها،

کریمی، ن. و دست‌برآورده، ع. ۲۰۰

مدل‌سازی فرایندهای سری زمانی شبکه عصبی بیزی برای شناسایی پیش‌زلزله با ناهنجاری‌های یونوسفری زلزله خاش،

کیکاووسی، م. و کریمی، ف. و حسینی، ا. ۲۱۳

پیش‌بینی زمان حوادث ترافیکی با استفاده از شبکه عصبی بقا،

محمدی، ف. و آرشی، م. و حبیبی راد، آ. و محمدزاده، ا. ۲۲۱

مدل‌بندی تصادفات جاده‌ای استان خراسان جنوبی با استفاده از مدل‌های اتورگرسیو صحیح مقدار،

نخعی‌زاده، ز. و جمهوری، س. ۲۳۶

فرایند چوب‌شکنی رگرسیون-بتا با کوواریانس ناماننا وابسته به متغیر کمکی،

یارعلی، ا. و ریواز، ف. و جعفری خالیدی، م. ۲۴۵



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30-31, 2025
University of Kurdistan



خوشه‌بندی دوطرفه فازی با استفاده از تابع فاصله اطلاعاتی

زهرا احمدی^۱، دکتر رضا زارعی^۲

^۱استادیار گروه آمار، دانشکده علوم ریاضی، دانشگاه گیلان

^۲فارغ التحصیل کارشناسی ارشد آمار ریاضی، گروه آمار، دانشکده علوم ریاضی، دانشگاه گیلان

چکیده: با گسترش روزافزون داده‌ها در حوزه‌های مختلف علمی، نیاز به روش‌های مؤثر برای تحلیل و استخراج الگوهای پنهان در داده‌ها بیش از پیش احساس می‌شود. خوشه‌بندی یکی از روش‌های مهم در داده‌کاوی غیرنظارتی است که با هدف شناسایی ساختارهای پنهان و تقسیم داده‌ها به گروه‌های مشابه به‌کار می‌رود. در این میان، خوشه‌بندی فازی به دلیل انعطاف‌پذیری بالا در برخورد با داده‌های نادقیق، جایگاه ویژه‌ای یافته است. در این مقاله ضمن معرفی رویکرد خوشه‌بندی فازی دوطرفه، یک الگوریتم جدید در این حوزه مورد بررسی قرار می‌گیرد. جزئیات محاسباتی این الگوریتم با استفاده از مجموعه داده‌های واقعی در حوزه زیست‌پزشکی تشریح شده و عملکرد آن با چند الگوریتم مشابه مقایسه شده است. نتایج نشان‌دهنده برتری نسبی الگوریتم مذکور از نظر دقت، تفسیرپذیری و معیارهای ارزیابی خوشه‌بندی است.

واژه‌های کلیدی: خوشه‌بندی فازی، خوشه‌بندی دوطرفه، تابع هدف، توابع عضویت

کد موضوع‌بندی ریاضی (۲۰۲۰): 68T09, 03B52, 62H30

۱ مقدمه

تحلیل و دسته‌بندی داده‌های حجیم و پیچیده یکی از چالش‌های اساسی در علوم داده و آمار است. خوشه‌بندی^۱ به عنوان یکی از روش‌های پایه در داده‌کاوی، نقش مؤثری در کشف ساختار درون داده‌ها دارد. هدف اصلی خوشه‌بندی، تقسیم داده‌ها به گروه‌هایی است که درون هر گروه، داده‌ها بیشترین شباهت و بین گروه‌ها بیشترین تفاوت را داشته باشند. در روش‌های سنتی خوشه‌بندی یا خوشه‌بندی سخت^۲، هر داده تنها به یک خوشه اختصاص می‌یابد. با این حال، بسیاری از داده‌های دنیای واقعی دارای مرزهای نامشخص یا هم‌پوشانی هستند.

^۱ سخنران، ahmadizahra19.77@gmail.com

^۱Clustering

^۲Hard Clustering

در چنین شرایطی، خوشه‌بندی فازی^۳ گزینه‌ای مناسب‌تر است زیرا به هر داده اجازه می‌دهد تا به چند خوشه با درجات عضویت متفاوت تعلق گیرد.

در بسیاری از کاربردها، مانند تحلیل بیان ژن، خوشه‌بندی همزمان نمونه‌ها (ردیف‌ها) و ویژگی‌ها (ستون‌ها) اهمیت فزاینده‌ای پیدا کرده است. این رویکرد که به خوشه‌بندی دوطرفه موسوم است، قادر است زیرگروه‌هایی از نمونه‌ها را همراه با زیرگروه‌هایی از ویژگی‌ها که رفتاری مشابه از خود نشان می‌دهند، شناسایی کند. علی‌رغم پیشرفت‌های قابل توجه در خوشه‌بندی فازی دوطرفه، همچنان چالش‌هایی اساسی در مواجهه با داده‌های بسیار پیچیده و نویزی و به ویژه در استخراج خوشه‌های با مرزهای هم‌پوشان و نتایج با قابلیت تفسیرپذیری بالا وجود دارد. روش‌های موجود اغلب در مدل‌سازی دقیق وابستگی‌های پیچیده میان نمونه‌ها و ویژگی‌ها و همچنین مدیریت بهینه ابهام و عدم قطعیت در داده‌ها با محدودیت‌هایی مواجه‌اند. در این راستا، این مقاله به معرفی و بررسی یک الگوریتم جدید مبتنی بر اطلاعات می‌پردازد. نوآوری اصلی این الگوریتم در استفاده از یک تابع فاصله اطلاعاتی و ترکیب دو عبارت آنتروپی مجزا برای کنترل پراکندگی عضویت‌ها در سطرها و ستون‌ها است که همگرایی پایدارتر، مدیریت بهینه عدم قطعیت و در نتیجه نتایج دقیق‌تر و قابل تفسیرتری ارائه می‌دهد.

باقیمانده این مقاله به شرح زیر سازماندهی شده است: در بخش دوم، مبانی نظری خوشه‌بندی دوطرفه فازی به کوتاهی مرور می‌شود. بخش سوم به معرفی جامع و جزئیات محاسباتی الگوریتم پیشنهادی اختصاص دارد. در بخش چهارم، نتایج آزمایش‌های تجربی روی مجموعه داده‌های واقعی و مقایسه عملکرد الگوریتم تحت بررسی با الگوریتم‌های موجود در این زمینه ارائه می‌شود. در نهایت، بخش پنجم شامل بحث و نتیجه‌گیری کلی از یافته‌های پژوهش است.

۲ خوشه‌بندی دوطرفه

خوشه‌بندی کلاسیک معمولاً تنها بر سطرها (نمونه‌ها) یا ستون‌ها (ویژگی‌ها) انجام می‌شود. اما داده‌های پیچیده‌تر، مانند ماتریس‌های با ابعاد بزرگ و دارای ساختارهای پنهان، نیازمند روشی هستند که بتواند همزمان خوشه‌بندی روی نمونه‌ها و ویژگی‌ها انجام دهد. این روش با عنوان خوشه‌بندی دوطرفه^۴ شناخته می‌شود که اولین بار توسط بوریس میرکین مطرح شد. در این روش، هدف شناسایی زیرماتریس‌هایی از داده‌هاست که در آن‌ها گروهی از ویژگی‌ها با گروهی از نمونه‌ها به صورت همزمان الگوهای مشترکی دارند. خوشه‌بندی دوطرفه کاربردهای گسترده‌ای در متن‌کاوی، زیست‌فناوری، داده‌های پزشکی و سیستم‌های توصیه‌گر دارد و نسبت به روش‌های سنتی، نتایج دقیق‌تر و تفسیرپذیرتری ارائه می‌دهد. تفاوت این روش با خوشه‌بندی کلاسیک در این است که خوشه‌بندی دوطرفه همزمان روی ردیف‌ها و ستون‌ها کار می‌کند و خروجی آن زیرماتریس‌هایی با الگوهای مشترک است و انعطاف‌پذیری زیادی دارد یعنی هر آیت می‌تواند در چند خوشه باشد این در حالی است که در خوشه‌بندی کلاسیک خوشه‌بندی فقط روی ردیف یا ستون انجام می‌شود و خروجی آن گروهی از ردیف‌ها یا ستون‌ها است و انعطاف‌پذیری کمی دارد. فرارو، ام.بی. (۲۰۲۴)

فرض کنید X یک ماتریس $n \times p$ دادگان باشد. هدف الگوریتم خوشه‌بندی دوطرفه تقسیم‌بندی همزمان n سطر به k خوشه و p ستون به h خوشه است. k -میانگین دوطرفه توسط ویچی و همکارانش (۲۰۰۱) مطرح شد. نسخه فازی آن توسط فرارو و همکارانش (۲۰۱۵) معرفی شد و سپس در سال ۲۰۲۱ عمیقاً مورد بررسی قرار گرفت.

^۳Fuzzy Clustering

^۴Co-Clustering

تابع هدف الگوریتم k -میانگین دوطرفه فازی به صورت زیر مدل‌بندی می‌شود

$$J(\mathbf{X}, U, C) = \sum_{j=1}^k \sum_{i=1}^n \sum_{f=1}^h \sum_{g=1}^p (x_i - c_j)^2 (u_{ij})^m (v_{gf})^l$$

که در آن

$$\sum_{j=1}^k u_{ij} = 1 \quad \sum_{f=1}^h v_{gf} = 1, \quad u_{ij}, v_{gf} \in [0, 1].$$

همچنین U ماتریسی $n \times k$ حاوی درجات عضویت سطرها و V ماتریسی $p \times h$ حاوی درجات عضویت ستون‌ها می‌باشند. l و m مشخص‌کننده درجات فازی و اعداد حقیقی و بزرگتر از یک هستند.

۳ الگوریتم خوشه‌بندی دوطرفه فازی مبتنی بر اطلاعات

الگوریتم خوشه‌بندی دوطرفه فازی مبتنی بر اطلاعات^۵ (به کوتاهی $ibFCC$) یکی از روش‌های نوین در زمینه خوشه‌بندی فازی دوطرفه است که با هدف بهبود دقت خوشه‌بندی و تفسیر بهتر خوشه‌ها طراحی شده است. این الگوریتم از تابع هدفی بهره می‌برد که شامل سه جزء اصلی است: (۱) تابع فاصله اطلاعاتی بین داده‌ها و مراکز خوشه‌ها، (۲) آنتروپی توزیع عضویت نمونه‌ها در خوشه‌ها، و (۳) آنتروپی توزیع عضویت ویژگی‌ها در خوشه‌ها.

در این مدل، دو ماتریس عضویت به نام‌های U و V تعریف می‌شوند که به ترتیب درجات عضویت نمونه‌ها و ویژگی‌ها را در خوشه‌های مربوطه تعیین می‌کنند. پارامترهای T_u و T_v برای تنظیم میزان پراکندگی عضویت‌ها استفاده می‌شوند. مزیت اصلی $ibFCC$ نسبت به روش‌های مشابه در استفاده از فاصله اطلاعاتی به جای فاصله اقلیدسی است که باعث می‌شود ارتباط میان ویژگی‌ها و نمونه‌ها بهتر مدل‌سازی شود.

هدف $ibFCC$ به حداقل رساندن تابع هدف معادله (۱.۳) با توجه به محدودیت‌های معادله زیر است

$$J_{ibFCC} = \sum_{c=1}^C \sum_{i=1}^n \sum_{j=1}^k u_{ci} v_{cj} d_{cij} + T_u \sum_{c=1}^C \sum_{i=1}^n u_{ci} \ln u_{ci} + T_v \sum_{c=1}^C \sum_{j=1}^k v_{cj} \ln v_{cj}. \quad (1.3)$$

$$\sum_{c=1}^C u_{ci} = 1 \quad \text{و} \quad \sum_{j=1}^k v_{cj} = 1 \quad i = 1, \dots, n \quad \text{و} \quad c = 1, \dots, C$$

عبارت اول معادله (۱.۳) درجه تطابق است که باید طی فرایند خوشه‌بندی دوطرفه به حداقل برسد تا افراد و ویژگی‌های بسیار مرتبط را باهم خوشه‌بندی دوطرفه کند. u_{ci} و v_{cj} دو تابع عضویت هستند که به ترتیب عضویت اسناد و ویژگی‌ها را نشان می‌دهند. عبارت دوم و سوم عوامل تنظیم آنتروپی هستند که u_{ci} و v_{cj} را جداگانه ترکیب می‌کنند. T_u و T_v پارامترهای وزنی هستند که درجه فازی را در خوشه‌های نهایی کنترل می‌کنند. معادلات به‌روزرسانی برای عضویت سند و ویژگی که d_{cij} فاصله بین نقطه داده ویژگی و مرکز خوشه ویژگی است، به صورت زیر بدست می‌آید

$$u_{ci} = \frac{\exp\left(-\frac{\sum_{j=1}^k v_{cj} d_{cij}}{T_u}\right)}{\sum_{c=1}^C \exp\left(-\frac{\sum_{j=1}^k v_{cj} d_{cij}}{T_u}\right)} \quad \text{و} \quad v_{cj} = \frac{\exp\left(-\frac{\sum_{i=1}^n u_{ci} d_{cij}}{T_v}\right)}{\sum_{j=1}^k \exp\left(-\frac{\sum_{i=1}^n u_{ci} d_{cij}}{T_v}\right)}. \quad (2.3)$$

⁵Information-based Fuzzy Co-Clustering

برای محاسبه اطلاعات فاصله مجموعه دادگان از رابطه زیر استفاده می‌کنیم

$$d_{cij} = \frac{1}{n} x_{ij} \log\left(\frac{x_{ij}}{t_{cij}}\right) + \frac{|c|}{n} p_{cj} \log\left(\frac{p_{cj}}{t_{cij}}\right), \quad (3.3)$$

که $t_{cij} = \frac{(x_{ij} + |c| * p_{cj})}{(1 + |c|)}$ ، تعداد اسناد در c امین خوشه است. تعریف مقدار $|c|$ در خوشه‌بندی فازی نسبت به خوشه‌بندی سخت کمی پیچیده‌تر است زیرا باید عملیات فازی‌سازی روی ماتریس عضویت انجام شود. پس از فازی‌سازی به راحتی می‌توان مقدار $|c|$ را در خوشه‌بندی سخت بدست آورد.

در الگوریتم *ibFCC* بدست آوردن p_{cj} به طور صریح دشوار است حتی اگر مقدار p_{cj} از نظر تئوری به صورت u_{ci} و v_{cj} محاسبه شود ممکن است فرایند محاسباتی بالایی داشته باشد. بنابراین یک رویکرد جایگزین انتخاب می‌شود که از روش میانگین وزنی استفاده می‌کند. در خوشه‌بندی فازی مرکز یک خوشه میانگین تمام نقاط است که براساس درجه تعلق آن‌ها به خوشه وزن می‌شود. در نتیجه معادله به‌روزرسانی نرمال شده p_{cj} به صورت زیر بدست می‌آید

$$p_{cj} = \frac{\sum_{i=1}^n u_{ci} x_{ij}}{\sum_{i=1}^{Nn} u_{ci}}. \quad (4.3)$$

۱.۳ مراحل اجرای الگوریتم *ibFCC*

مرحله	شرح عملیات
۱	مقداردهی اولیه به پارامترها: تعداد خوشه‌ها C ، بیشینه تکرارها τ_{\max} ، آستانه همگرایی ε ، ضرایب تنظیم آنتروپی T_u و T_v .
۲	تنظیم $\tau = 1$.
۳	مقداردهی اولیه به درجات عضویت u_{ci} و v_{cj} به صورت تصادفی.
۴	شروع تکرارها
۵	محاسبه مراکز خوشه ویژگی‌ها با استفاده از رابطه (۴.۳).
۶	محاسبه فاصله اطلاعاتی بین اسناد و ویژگی‌ها با استفاده از رابطه (۳.۳).
۷	به‌روزرسانی درجات عضویت با استفاده از رابطه (۲.۳).
۸	افزایش شمارنده تکرار: $\tau \leftarrow \tau + 1$.
۹	بررسی شرط توقف: $\tau = \tau_{\max}$ یا $\max u_{ci}^{(\tau)} - u_{ci}^{(\tau-1)} \leq \varepsilon$.

یکی از تفاوت‌های کلیدی بین الگوریتم $ibFCC$ و الگوریتم‌های فازی دوطرفه کلاسیک، ساختار تابع هدف آن‌هاست. در روش‌های سنتی مانند فازی k -میانگین دوطرفه، تابع هدف بر اساس فاصله اقلیدسی و ترکیب درجه‌های عضویت سطرها و ستون‌ها تعریف می‌شود و تمرکز آن صرفاً بر کمینه‌سازی فاصله بین داده‌ها و مراکز خوشه است. اما این رویکرد در مواجهه با داده‌های پیچیده، نویزی یا دارای هم‌پوشانی بالا، کارایی محدودی دارد. در مقابل، تابع هدف الگوریتم $ibFCC$ با بهره‌گیری از یک تابع فاصله اطلاعاتی (اطلاعات متقابل بین ویژگی و سند) و نیز دو عبارت آنتروپی مجزا برای سطرها و ستون‌ها، رویکردی غنی‌تر و منعطف‌تر ارائه می‌دهد. پارامترهای کنترلی T_u و T_v امکان تنظیم میزان پراکندگی عضویت‌ها را فراهم کرده و از همگرایی سریع و نامناسب جلوگیری می‌کنند. این طراحی باعث می‌شود الگوریتم $ibFCC$ بتواند ساختارهای پنهان و الگوهای معنادارتر را در داده‌ها شناسایی کرده و خوشه‌بندی دقیق‌تری ارائه دهد. به‌طور خلاصه، تفاوت‌های موجود در تابع هدف موجب افزایش تطبیق‌پذیری، دقت، و تفسیرپذیری خروجی الگوریتم $ibFCC$ نسبت به روش‌های سنتی شده‌اند، که در بخش بعدی از طریق ارزیابی عددی در داده‌های واقعی به آن پرداخته می‌شود.

۴ نتایج عددی

برای ارزیابی عملکرد الگوریتم $ibFCC$ ، آزمایش‌هایی روی پنج مجموعه داده واقعی در حوزه زیست‌پزشکی انجام شد. ویژگی‌های این مجموعه داده‌ها در جدول زیر آمده است

جدول ۱: جزئیات مجموعه داده‌های مورد بررسی

نام مجموعه داده	تعداد نمونه	تعداد ویژگی	تعداد دسته‌بندی‌ها	توضیحات
Ohsumed	۲۰۰۰۰ (سند از ۵۰۲۱۶)	۵۰۰	۱۵	۲۳ دسته بیماری‌های موضوعی پزشکی (MeSH).
سرطان ریه (LC)	۲۷	۵۶	۳	برای شناخت انواع پاتولوژیک سرطان ریه.
بافت سینه (BT)	۱۰۶	۹	۶	یش‌بینی طبقه‌بندی با ادغام کلاس‌های فیبروآدنوم، ماستوپاتی و غده‌ای.
کاردیوتوکوگرافی (CTG)	۲۱۲۶	۲۱	۱۰	۲۱۲۶ کاردیوتوکوگرام جنین، ویژگی‌های تشخیصی اندازه‌گیری شده.
بیان پروتئین موش (MPE)	۱۰۷۶	۶۸	۸	۱۰۷۶ اندازه‌گیری در هر پروتئین، با ۸ دسته از موش‌ها.

هر ۵ الگوریتم به صورت تصادفی مقداره‌ی و ده بار اجرا شدند تا تاثیر بهینه‌سازی‌های موضعی کاهش یابد. از بین پنج الگوریتم، خوشه‌بندی c -میانگین (به کوتاه‌ی FCM) یک الگوریتم خوشه‌بندی فازی استاندارد است و بقیه، الگوریتم‌های خوشه‌بندی دوطرفه فازی هستند. در ادامه، الگوریتم $ibFCC$ با چهار الگوریتم دیگر شامل

FCM ، خوشه‌بندی فازی برای داده‌های چندمتغیره رسته‌ای (به کوتاه‌ی $FCCM$)، خوشه‌بندی دوطرفه استوار فازی (به کوتاه‌ی $RFCC$) و خوشه‌بندی دوطرفه فازی برای تصاویر (به کوتاه‌ی $FCCI$) مقایسه شد. معیارهای ارزیابی شامل آنتروپی، شاخص F و $Log(CS)$ بودند و نتایج نشان داد که الگوریتم $ibFCC$ در اکثر موارد عملکرد بهتری داشته و به‌ویژه در داده‌های BT و MPE

جدول ۲: مقایسه الگوریتم‌ها بر اساس معیارهای F ، آنتروپی و $\log(CS)$

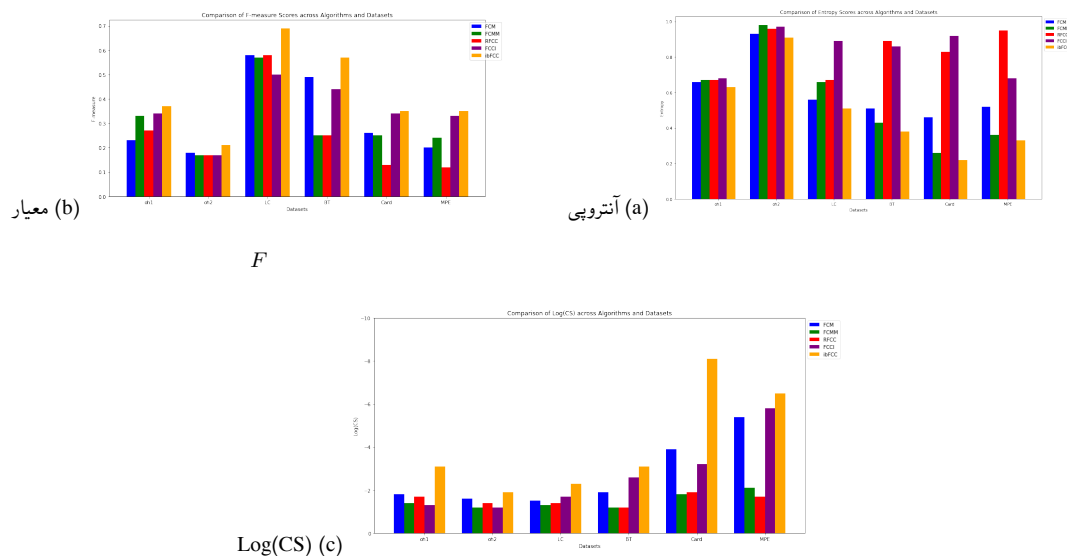
الگوریتم	نویسندگان	داده‌ها	F	آنتروپی	$\log(CS)$
FCM	یزدک (۱۹۷۳)	Oh1	۰/۲۳	۰/۶۶	-۱/۸
		Oh2	۰/۱۸	۰/۹۳	-۱/۶
		LC	۰/۵۸	۰/۵۶	-۱/۵
		BT	۰/۴۹	۰/۵۱	-۱/۹
		Card	۰/۲۶	۰/۴۲	-۳/۹
		MPE	۰/۳۵	۰/۲۲	-۵/۱
FCCM	اوسا، ج و همکاران (۲۰۰۱)	Oh1	۰/۳۳	۰/۶۷	-۱/۴
		Oh2	۰/۲۷	۰/۹۶	-۱/۳
		LC	۰/۵۵	۰/۶۶	-۱/۵
		BT	۰/۴۶	۰/۶۰	-۱/۸
		Card	۰/۲۵	۰/۳۶	-۲/۱
		MPE	۰/۳۴	۰/۲۶	-۳/۹
RFCC	تجی، دلیو، سی. و همکاران (۲۰۰۷)	Oh1	۰/۲۷	۰/۶۷	-۱/۷
		Oh2	۰/۲۰	۰/۹۲	-۱/۴
		LC	۰/۵۸	۰/۹۶	-۱/۳
		BT	۰/۵۲	۰/۸۹	-۱/۲
		Card	۰/۲۵	۰/۸۵	-۱/۹
		MPE	۰/۳۴	۰/۹۵	-۳/۷
FCCI	هانماندلوام، و همکاران (۲۰۱۳)	Oh1	۰/۳۴	۰/۶۸	-۱/۳
		Oh2	۰/۲۸	۰/۹۰	-۱/۲
		LC	۰/۵۷	۰/۸۵	-۱/۲
		BT	۰/۴۹	۰/۸۰	-۱/۴
		Card	۰/۲۷	۰/۷۸	-۳/۵
		MPE	۰/۳۶	۰/۶۵	-۵/۸
ibFCC	لیو و همکاران (۲۰۱۷)	Oh1	۰/۳۵	۰/۷۰	-۱/۳
		Oh2	۰/۲۹	۰/۹۱	-۱/۲
		LC	۰/۵۷	۰/۸۵	-۱/۲
		BT	۰/۴۹	۰/۸۰	-۱/۴
		Card	۰/۳۵	۰/۲۲	-۶/۵
		MPE	۰/۳۵	۰/۳۳	-۵/۳

توانسته دقت خوشه‌بندی را به‌طور چشم‌گیری افزایش دهد. بعد از $ibFCC$ الگوریتم $FCCI$ عملکرد بهتری از FCM ، $FCCM$ و $RFCC$ دارد.

بحث و نتیجه‌گیری

در این مطالعه الگوریتم $ibFCC$ به عنوان یک روش جدید برای خوشه‌بندی فازی دوطرفه مورد مطالعه و بررسی قرار گرفت. این الگوریتم با در نظر گرفتن اطلاعات متقابل بین ویژگی‌ها و نمونه‌ها، نسبت به روش‌های موجود که عمدتاً بر فاصله اقلیدسی تمرکز دارند، برتری محسوسی دارد. علاوه بر این، $ibFCC$ با استفاده از دو عبارت آنتروپی مجزا برای سطرها (نمونه‌ها) و ستون‌ها (ویژگی‌ها)، انعطاف‌پذیری بیشتری در مدل‌سازی ساختارهای پنهان داده‌ها ارائه می‌دهد. وجود پارامترهای کنترلی در این الگوریتم، امکان تنظیم میزان پراکندگی عضویت‌ها را فراهم کرده و به همگرایی سریع و جلوگیری از همگرایی نامناسب کمک کرده است.

نتایج ارزیابی بر روی پنج مجموعه داده واقعی نشان داد که الگوریتم $ibFCC$ در بیشتر موارد عملکرد بهتری داشته و دقت خوشه‌بندی را به‌طور چشم‌گیری افزایش داده است. این برتری عملکردی، $ibFCC$ به دلیل توانایی آن در مدل‌سازی بهتر ارتباطات پیچیده میان



شکل ۱: مقایسه سه معیار F ، آنتروپی و $\text{Log}(CS)$ برای الگوریتم‌های جدول نتایج روی داده‌ها

ویژگی‌ها و نمونه‌ها، و همچنین انعطاف‌پذیری بالای آن در مواجهه با داده‌های نادقیق و هم‌پوشان است. در نهایت، با توجه به نتایج به‌دست‌آمده از ارزیابی‌های تجربی، می‌توان نتیجه گرفت که الگوریتم ibFCC به‌عنوان یک ابزار مؤثر و کارآمد برای تحلیل داده‌های پیچیده و دارای ساختار پنهان محسوب می‌شود. قابلیت‌های پیشرفته این الگوریتم در کشف ساختارهای نهفته و ارائه نتایج دقیق‌تر و تفسیرپذیرتر، آن را به گزینه‌ای مناسب برای کاربردهای مختلف در علوم داده، به‌ویژه در حوزه‌هایی مانند زیست‌پزشکی که با داده‌های حجیم و پیچیده سروکار دارند، تبدیل می‌کند.

مراجع

طیبی، م. و زارعی، ر. (۱۴۰۱)، رویکرد بیزی در خوشه‌بندی تحت شرایط نادقیق، پایان نامه کارشناسی ارشد آمار ریاضی، دانشگاه گیلان.

Bezdek J. C. (1973). *Fuzzy Mathematics in Pattern Classification*, John Wiley Sons, Hoboken, New Jersey.

Ferraro M.B., Vichi M., Grzegorzewski P., Gagolewski M., Hryniewicz O. and Gil M. (2015), Fuzzy double clustering: a robust proposal. , *Advances in Intelligent systems and computing* 315. Springer, Cham. 225-232.

Ferraro M.B. (2024), Fuzzy k-means: history and applications, *Elsevier*; 2024: 110-123.

Hanmandlu M., Verma O.P., Susan S. and Madasu V.K. (2013), Color segmentation by fuzzy co-clustering of chrominance color features, *Neurocomputing* 120: 235–249.

Liu Y., Wu. S., Liu Z. and Chao H. (2017), A fuzzy co-clustering algorithm for biomedical data., *PLoS ONE* 12(4): e0176536

Oh CH., Honda K., and Ichihashi H. (2001), Fuzzy clustering for categorical multivariate data, 2154–2159 vol. 2154.

Tjhi W.C. and Chen L., (2007), Robust fuzzy Co-clustering algorithm, 1-5.

Vichi M., Borra S., Rocci R., Vichi M. and Schader M., (2001), Double kmeans clustering for simultaneous classification of objects and variables, *Advances in classification and data analysis. studies in classification, data analysis, and knowledge organization*. Springer, Berlin, Heidelberg

Fuzzy Co-Clustering using Information Distance

Zahra Ahmadi¹, Reza Zarei²

¹Department of Statistics, University of Guilan

²Department of Statistics, University of Guilan

Abstract: With the rapid growth of data in various scientific domains, the need for effective methods to analyze and extract hidden patterns has become increasingly important. Among clustering techniques, fuzzy clustering has gained significant attention due to its high flexibility in handling ambiguous and overlapping data. This paper introduces and investigates an advanced algorithm called interval-based bidirectional fuzzy co-clustering (ibFCC) designed for imprecise environments. Using real-world biomedical datasets, the performance of this algorithm is compared with several related methods. The results indicate its relative superiority in terms of accuracy, interpretability, and clustering evaluation metrics.

Keywords: Fuzzy clustering, Co-clustering, Membership functions, Objective function

Mathematics Subject Classification (2020): 62H30, 62B52, 68T09.



پانزدهمین سمینار احتمال
و فرایندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30-31, 2025
University of Kurdistan



یک روش جدید تنک‌سازی فرایندهای نقطه‌ای مبتنی بر مجموعه‌های تصادفی بولی

رضا اسدی^۱، مجتبی خزائی^۱، مجتبی گنجعلی^۱

^۱گروه آمار، دانشگاه شهید بهشتی

چکیده: فرایندهای نقطه‌ای فضایی از مدل‌های آماری برای تحلیل الگوهای نقطه‌ای اند. یکی از روش‌های ساخت فرایندهای نقطه‌ای جدید تنک‌سازی فرایندهای نقطه‌ای موجود است. در این مقاله یک الگوریتم تنک‌سازی جدید مبتنی بر یک مجموعه تصادفی بولی ارائه می‌شود. این الگوریتم را می‌توان تعمیمی از روش تنک‌سازی ارائه شده توسط اشتویان در معرفی فرایندهای نقطه‌ای وابسته دانست. بعضی ویژگی‌های فرایند نقطه‌ای تنک‌شده و فرایند نقاط حذف شده حاصل از اعمال این الگوریتم تنک‌سازی بر یک فرایند نقطه‌ای پواسون همگن بررسی شده و بعضی از آماره‌های خلاصه آن‌ها محاسبه شده است. از روش درستنمایی پالم برای برآورد پارامترها استفاده شده است. در نهایت، مدل‌های فوق برای مدل‌بندی دو مجموعه داده واقعی استفاده می‌شود.

واژه‌های کلیدی: فرایند نقطه‌ای پواسون؛ فرایندهای کاکس؛ مجموعه تصادفی بولی؛ تنک‌سازی وابسته؛ درستنمایی پالم.

کد موضوع‌بندی ریاضی (۲۰۲۰): xxCxx, xxBxx, xxAxx.

۱ مقدمه

فرایندهای نقطه‌ای فضایی از مدل‌های آماری برای مدل‌بندی الگوهای نقطه‌ای اند. به عنوان مثال، از فرایندهای نقطه‌ای فضایی برای مدل‌بندی محل انواع درختان جنگلی، محل وقوع زمین لرزه‌ها، لانه پرندگان و ... می‌توان استفاده کرد. رفتار الگوهای نقطه‌ای به سه دسته رفتار خوشه‌ای، رفتار تصادفی محض و رفتار منظم یا دفعی تقسیم‌بندی می‌شود. متناسب با رفتار هر الگوی نقطه‌ای دسته‌ای از مدل‌ها برای مدل‌بندی الگوهای نقطه‌ای ارائه شده است. روش‌های زیادی از جمله اعمال تبدیلات و یا اصلاح فرایندهای نقطه‌ای برای ساخت فرایند نقطه‌ای جدید وجود دارند. چهار رویکرد کلی برای این کار عبارتند از: نگاشت، برهم نهی، خوشه‌بندی و تنک‌سازی.

در رویکرد تنک‌سازی نقاط فرایند اولیه براساس یک الگوریتم حذف می‌شوند. در صورتی که احتمال حذف نقاط فرایند اولیه از هم مستقل باشند، تنک‌سازی مستقل و در صورتی که احتمال حذف نقاط فرایند اولیه به هم وابسته باشند، تنک‌سازی وابسته نامده می‌شوند.

^۱ سخنران، rezaasadi220@gmail.com

در این مقاله ابتدا یک الگوریتم تنک‌سازی معرفی می‌شود که در آن احتمال‌های حذف براساس تحقیقی از یک مدل بولی مشخص می‌شوند. سپس فرایندهای تنک‌شده و نقاط حذف شده حاصل از اعمال این الگوریتم روی یک فرایند اولیه پواسون همگن بررسی می‌شود. الگوریتم تنک‌سازی ارائه شده در این مقاله را می‌توان تعمیمی از الگوریتم تنک‌سازی ارائه شده توسط **اشتویان (۱۹۷۹)** و **مولر و لوانسیر (۲۰۱۶)** دانست با این تفاوت که در اینجا احتمال حذف درون تحقق مجموعه تصادفی بولی می‌تواند هر مقداری بین صفر و یک باشد در حالی که **اشتویان (۱۹۷۹)** و **مولر و لوانسیر (۲۰۱۶)** تنها احتمال حذف نقاط خارج از مجموعه تصادفی بولی را برابر یک در نظر گرفته‌اند. یک مجموعه تصادفی بولی به صورت زیر تعریف می‌شود.

تعریف ۱.۱. (**اشتویان ، ۱۹۷۹**) یک مجموعه تصادفی بولی از قرار دادن تحقق‌های یک مجموعه تصادفی بسته، معروف به دانه، روی نقاط یک فرایند نقطه‌ای پواسون، معروف به ریشه، حاصل می‌شود. در صورتی که Z فرایند پواسون مانا روی \mathbb{R}^2 با شدت θ و مستقل از M_s ها؛ یعنی مجموعه‌های تصادفی بسته مستقل و هم توزیع با M_o به‌طوری‌که $E\{\sup\|x\|, x \in M_o\} < \infty$ باشد، آنگاه مجموعه تصادفی بولی حاصل از دانه بسته تصادفی M_o و فرایند نقطه‌ای Z برابر است به $\Xi = \bigcup_{s \in Z} (s + M_s)$

الگوریتم‌های تنک‌سازی متنوعی که فرایند تنک‌شده حاصل از اعمال آن‌ها بر فرایند نقطه‌ای فضایی پواسون مورد بررسی قرار گرفته است معرفی شده است. به عنوان مثال فرایندهای هسته سخت مترن، در فرایندهای هسته سخت نقاط فرایند اولیه براساس تابعی از فاصله نقاط از یکدیگر حذف می‌شوند لکن در الگوریتم تنک‌سازی معرفی شده در این مقاله نقاط فرایند اولیه براساس تابعی از فرایند Z که مستقل از فرایند اولیه است حذف می‌شوند.

۲ معرفی الگوریتم تنک‌سازی پیشنهادی و فرایندهای مربوطه

فرایند نقطه‌ای اولیه Y و فرایند نقطه‌ای Z به عنوان فرایندی که بر اساس آن Y تنک می‌شود را در نظر بگیرید. در صورتی که فرایندهای Y و Z بر \mathbb{R}^2 تعریف شده و از هم مستقل باشند. حول هر نقطه از فرایند نقطه‌ای فضایی Z دایره‌ای با شعاع r را در نظر گرفته و احتمال حذف نقطه‌ای از فرایند اولیه درون هر دایره را مستقل از سایر دایره برابر با q در نظر می‌گیریم. به این ترتیب یک الگوریتم تنک‌سازی حاصل می‌شود که در آن:

۱- در ناحیه خارج از اجتماع دایره احتمال حذف صفر است.

۲- در نواحی درون یک دایره احتمال حذف $q = 1 - p$ است.

۳- در نواحی مشترک بین k دایره احتمال حذف $1 - p^k$ است.

در صورتی که Z فرایند پواسون همگن و همسانگرد با تابع شدت θ باشد و $\Xi = \bigcup_{s \in Z} b(s, r)$ تعریف شود، یک مجموعه تصادفی بولی با دانه‌های دایره‌ای شکل به شعاع r است. الگوریتم تنک‌سازی فوق را می‌توان به صورت یک الگوریتم تنک‌سازی روی تحقق‌های این مجموعه بولی به صورت زیر تعریف کرد:

۱- در نواحی که خارج از مجموعه تصادفی بولی Ξ است نقاط Y حذف نمی‌شوند.

۲- در نواحی که محل تحقق فقط یک دانه از مجموعه تصادفی بولی Ξ است نقاط Y با احتمال $q = 1 - p$ حذف می‌شوند.

۳- در نواحی که محل اشتراک k دانه از مجموعه تصادفی بولی Ξ است نقاط فرایند اولیه Y با احتمال $q = 1 - p^k$ حذف می‌شوند.

در قضیه زیر توزیع حاشیه‌ای و توام احتمال بقای نقاط محاسبه می‌شود. اثبات قضیه‌های ارائه شده در این مقاله در ضمایم ارائه شده است.

قضیه ۱.۲. در الگوریتم تنک‌سازی پیشنهادی اگر $\Pi(\zeta)$ احتمال بقا در \mathbb{R}^2 باشد آن گاه:

(i) اگر $A_\zeta = b(\zeta, r)$ باشد برای هر $\zeta \in \mathbb{R}^2$:

$$P[\Pi(\zeta) = p^k] = P[N(A_\zeta) = k] = \frac{e^{-\theta\pi r^2} (\theta\pi r^2)^k}{k!}, \quad k = 0, 1, 2, \dots, \quad (1.2)$$

(ii) توزیع توام احتمال بقا برای هر $\zeta_1, \zeta_2 \in \mathbb{R}^2$ برابر است به:

$$P[\Pi(\zeta_1) = p^n, \Pi(\zeta_2) = p^m] = \begin{cases} P[\Pi(\zeta_1) = p^n]P[\Pi(\zeta_2) = p^m], & \|\zeta_1 - \zeta_2\| > 2r, \\ C[\Pi(\zeta_1, \zeta_2)], & \|\zeta_1 - \zeta_2\| \leq 2r, \end{cases} \quad (2.2)$$

که در آن

$$C[\Pi(\zeta_1, \zeta_2)] = e^{-\theta(\pi r^2 - A(t, r))} [\theta(\pi r^2 - A(t, r))]^{m+n} \times \sum_{k=0}^{\min(m, n)} \frac{1}{k!(n-k)!(m-k)!} \left[\frac{A(t, r)}{\theta(\pi r^2 - A(t, r))} \right]^k, \quad (3.2)$$

و $m, n \in \mathbb{N} \cup \{0\}, r > 0$

$$A(t, r) = 2r^2 \text{Arc cos} \left(\frac{t}{2r} \right) - \frac{t}{2} [4r^2 - t^2]^{\frac{1}{2}}, \quad (4.2)$$

مساحت ناحیه مشترک بین دو دایره به شعاع r با فاصله t ، $(t < 2r)$ است.

در صورتی که فرایند اولیه، Y ، یک فرایند نقطه‌ای پواسون با شدت δ باشد از اعمال الگوریتم فوق بر Y دو فرایند تنک‌شده‌ی X, Y ، و نقاط حذف شده از $Y - X$ ، فرایندهای نقطه‌ای تعریف می‌کنند که به پارامترهای (δ, θ, r, p) بستگی دارد. در ادامه ویژگی‌ها و رفتار این فرایندها بررسی می‌شود.

۳ آماره‌های خلاصه

آماره‌های خلاصه به عنوان ابزاری برای تحلیل اکتشافی داده‌ها، آزمون و بررسی اعتبار مدل توسط آماردان‌ها مورد استفاده قرار می‌گیرند. در تحلیل فرایندهای نقطه‌ای فضایی ویژگی‌های مرتبه اول و دوم متغیر شمارشی $N(B); B \subseteq \mathbb{R}^2$ اطلاعات سودمندی در مورد رفتار فرایند نقطه‌ای به تحلیل‌گر می‌دهند. تابع شدت مرتبه اول و دوم این ویژگی‌ها را تبیین می‌کنند. اگر $\mu(B) = E[N(B)]; B \subseteq \mathbb{R}^2$ ، $\mu(\cdot)$ اندازه شدت نامیده می‌شود. در صورتی که بتوان نوشت $\mu(B) = \int_B \lambda(\zeta) d\zeta; B \subseteq \mathbb{R}^2$ نوشت که در آن یک تابع نامنفی است، $\lambda(\cdot)$ تابع شدت نامیده می‌شود. اگر $\lambda^{(2)}(C) = E[\sum_{\zeta, \eta \in C} I[(\zeta, \eta) \in C]]$ ، $\lambda^{(2)}(\cdot)$ ، $\alpha^{(2)}(\cdot)$ اندازه شدت فاکتوری مرتبه دوم نامیده می‌شود. بعلاوه اگر $\alpha^{(2)}(C) = \int \int I[(\zeta, \eta) \in C] \lambda^{(2)}(\zeta, \eta) d\zeta d\eta$ که در آن $\lambda^{(2)}(\cdot, \cdot)$ یک تابع

نامنفی است به $(.,.)^{(2)} \lambda$ تابع شدت مرتبه دوم گفته می‌شود. در این‌جا با توجه به این‌که فرایند اولیه‌ی Y توسط یک میدان تصادفی تنک می‌شود لذا هر دو فرایند X و $Y - X$ فرایندهای کاکس هستند (مولر و واگه‌پیتسن، ۲۰۰۳).

قضیه ۱.۳. برای فرایند تنک‌شده‌ی X ، تابع شدت $\lambda(.)$ و تابع شدت مرتبه دوم $\lambda^{(2)}(.,.)$ برابر است به:

$$(i) \lambda(u) = \delta e^{-\theta(1-p)\pi r^2}.$$

$$(ii) \lambda^{(2)}(u, v) = \begin{cases} \delta^2 e^{-2\pi r^2 \theta(1-p) + \theta A(t, r)(1-p)^2}, & \|t\| \leq 2r, \\ \delta^2 e^{-2(1-p)\theta \pi r^2}, & \|t\| > 2r, \end{cases}$$

که در آن $t = \|u - v\|$ است.

تابع همبستگی زوجی یکی دیگر از آماره‌هایی است که ویژگی‌های مرتبه دوم فرایندهای نقطه‌ای را بررسی می‌کند. تابع همبستگی زوجی به صورت $g(\zeta, \eta) = \frac{\lambda^{(2)}(\zeta, \eta)}{\lambda(\zeta)\lambda(\eta)}$ تعریف می‌شود که در آن برای $a \geq 0$ قرار می‌دهیم $\frac{a}{\pi} = 0$. برای فرایندهای مانا می‌توان نوشت $g(\zeta, \eta) = g(\zeta - \eta)$ و برای فرایندهای مانا و همسانگرد $g(\zeta, \eta) = g(t)$ که در آن $t = \|\zeta - \eta\|$. در حالت کلی اگر $t \rightarrow \infty$ آن‌گاه $g(t) \rightarrow 1$. ولی برای t های کوچک، $g(t) < 1$ نشان‌دهنده فرایند با رفتار دفعی، $g(t) = 1$ نشان‌دهنده فرایند با رفتار تصادفی محض و $g(t) > 1$ نشان‌دهنده فرایند با رفتار خوشه‌ای است.

قضیه ۲.۳. تابع همبستگی زوجی برای فرایندهای X و $Y - X$ برابر است به:

$$(i) g_X(t) = \begin{cases} e^{\theta A(t, r)(1-p)^2}, & t \leq 2r, \\ 1, & t > 2r, \end{cases}$$

$$(ii) g_{Y-X}(t) = 1 + \frac{a^2}{(1-a)^2} (g_X(t) - 1),$$

که در آن $t = \|u - v\|$, $a = E[\Pi(o)]$.

(i) نشان می‌دهد که $g_X(t) \geq 1$ و از (ii) نتیجه می‌شود که $g_{Y-X}(t) \geq 1$. نتایج فوق بیانگر این است که هر دو فرایند X و $Y - X$ دارای رفتار خوشه‌ای هستند. تابع شدت پالم یکی دیگر از آماره‌هایی است که در استنباط آماری فرایندهای نقطه‌ای کاربرد زیادی دارد. در صورتی که x نقطه‌ای دلخواه در \mathbb{R}^2 با فاصله t از مبدأ، o ، باشد. آن‌گاه نرخ تحقق در x به شرط آن‌که o یک نقطه از فرایند باشد برای مجموعه لبگ با اندازه بسیار کوچک dx برابر است با: $1|N\{o\} = 1|$ به $\lambda_o(x)dx = P[N(dx) = 1|N\{o\} = 1]$ به $\lambda_o(x)$ شدت پالم در نقطه x گفته می‌شود. در صورتی‌که فرایند مانا و همسانگرد باشد تابع شدت پالم صرفاً به t وابسته است و می‌توان آن را به صورت $\lambda_o(t)$ نوشت. همچنین برای فرایندهای مانا و همسانگرد $\lambda_o(t) = \lambda g(t)$ (تاناکا و همکاران، ۲۰۰۸).

قضیه ۳.۳. تابع شدت پالم برای فرایندهای X و $Y - X$ برابر است به:

$$(i) \lambda_{o, X}(t) = \begin{cases} \delta e^{-\theta(1-p)[\pi r^2 - A(t, r)(1-p)]}, & t \leq 2r, \\ \delta e^{-\theta(1-p)\pi r^2}, & t > 2r, \end{cases}$$

$$(ii) \lambda_{o, Y-X}(t) = \delta[1 - e^{-\theta q \pi r^2}]\{1 + \frac{a^2}{(1-a)^2} (g_X(t) - 1)\}.$$

۴ برآورد پارامترها و ارزیابی مدل

روش‌های زیادی برای برآورد پارامترهای فرایندهای نقطه‌ای وجود دارد، از جمله حداقل تضادها، ماکسیم درستنمایی، شبه درستنمایی و درستنمایی پالم. در این مقاله برای برآورد پارامترهای فرایندهای X و $Y - X$ از روش درستنمایی پالم استفاده می‌شود. از روش‌های متداول بررسی نیکویی برازش در فرایندهای نقطه‌ای بررسی پوشش‌ها است. این روش مبتنی بر محاسبه و رسم نمودار چندک‌های آماره خلاصه به ازای t های متفاوت است. از روش بوت استرپ برای محاسبه چندک‌ها می‌توان استفاده کرد. این روش مبتنی بر این است که اگر مدل برازش داده شده صحیح باشد بروردگر ناپارامتریک آماره‌های خلاصه الگوی مشاهده شده به ازای t های متفاوت در اغلب اوقات رفتاری مشابه مدل مورد نظر (با پارامترهای برآورد شده) دارد.

۱.۴ درستنمایی پالم

درستنمایی پالم یک روش تقریبی برآورد پارامترهای فرایندهای نقطه‌ای است. این روش مبتنی بر فرایند تفاضل است. برای فرایند نقطه‌ای P که بر پنجره محدب W مشاهده شده است، فرایند تفاضل به صورت $D_{P,W} = \{x - y; x, y \in P_W, \text{ s.t. } x \neq y\}$ تعریف می‌شود. برای فرایند مانا و همسانگرد P ، فرایند تفاضل یک فرایند همسانگرد با تابع شدت $N(W)\lambda_o(t)$ است که در آن $\lambda_o(t)$ تابع شدت پالم فرایند P است. همچنین در صورتی که فرض شود فرایند P به خوبی توسط یک فرایند پواسون تقریب زده می‌شود آن‌گاه تابع درستنمایی فرایند برابر است با:

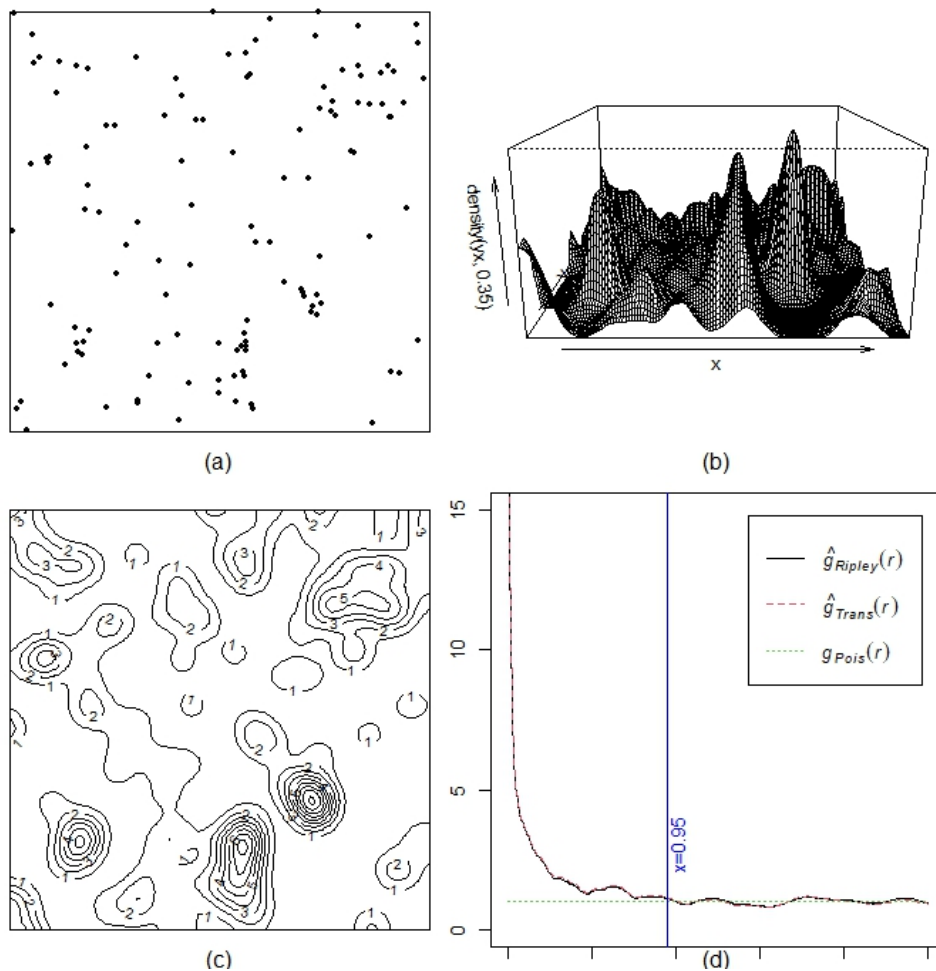
$$\log L(\delta, \theta, p, r, t) = \sum_{\{i,j; i \neq j; t_{ij} < T\}} \log\{N(W)\lambda_o(t_{ij})\} - N(W) \int_0^T \lambda_o(t) 2\pi t dt. \quad (1.4)$$

که در آن t_{ij} فاصله بین نقاط p_i و p_j است و T آستانه‌ای است که فرض می‌شود نقاط با فاصله بیشتر از آن از هم مستقل هستند. به $(\hat{\delta}, \hat{\theta}, \hat{r}, \hat{p})$ که عبارت فوق را ماکسیم کند برآوردگر ماکسیم درستنمایی پالم گفته می‌شود. برای آشنایی بیشتر به **تاناکا و همکاران (۲۰۰۸)** به مراجعه شود.

۵ کاربرد

در این بخش دو مجموعه داده واقعی بررسی می‌شود و فرایندهای X و $Y - X$ به آن‌ها برازش داده می‌شود.

مثال ۱.۵. الگوی نقطه‌ای نهال‌های کاج که در بسته *spatstat* نرم افزار \mathbb{R} تحت عنوان *finpines* موجود است، محل قرار گرفتن ۱۲۶ نهال کاج به همراه قطر و ارتفاع آن‌ها در جنگلی در فنلاند بر پنجره $[-8, 2] \times [-5, 5]$ را ثبت کرده است. در این‌جا صرفاً محل قرار گرفتن نهال‌ها بررسی و مدل‌بندی می‌شود. در شکل ۱ نمودار (a) الگوی نقطه‌ای، (b) نمودار ۳ بعدی تابع شدت، (c) منحنی‌های تراز تابع شدت، (d) نمودار تابع همبستگی زوجی، نهال‌های کاج را نشان می‌دهند. نمودار تابع همبستگی زوجی نشان‌دهنده رفتار خوشه‌ای الگوی نقطه‌ای است همچنین منحنی‌های تراز و نمودار ۳ بعدی تابع شدت نشان‌دهنده این است که الگوی نقطه‌ای تمایل به تجمع در نواحی دایره‌ای شکل دارند. لذا برازش فرایند $Y - X$ به الگوی نقطه‌ای فوق پیشنهاد می‌شود. برای برازش فرایند $Y - X$ به الگوی نقطه‌ای فوق براساس روش درستنمایی پالم پارامترهای مدل فوق برآورد می‌گردد. با توجه به این‌که برآورد ناپارامتریک تابع همبستگی زوجی برای $t \geq 0.95$ کوچکتر از ۱.۱ است (شکل ۱ مشاهده شود) لذا T برابر با ۰.۹۵ در نظر گرفته می‌شود. همچنین



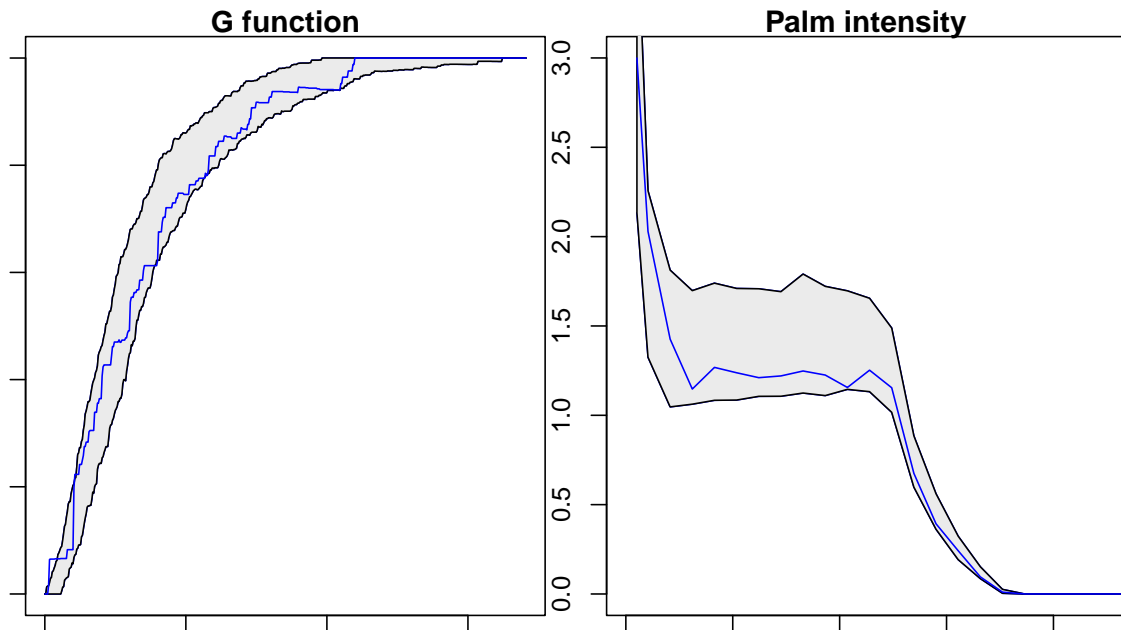
شکل ۱: (a) الگوی نقطه‌ای، (b) نمودار ۳ بعدی تابع شدت، (c) منحنی‌های تراز تابع شدت، (d) تابع همبستگی زوجی، نهال‌های کاج را نشان می‌دهد.

با جایگذاری $\lambda_{o,Y-X}(t)$ در (۱.۴) تابع درستنمایی پالم حاصل می‌شود. برآوردگر ماکسیمم درستنمایی پالم فرایند فوق در جدول ۱ ارائه شده است.

پارامتر	δ	θ	p	r
برآورد	۶۸۳/۵	۷۱۳/۱	۰۴۶/۰	۲۳۸۸/۰

جدول ۱: برآوردهای پارامترهای حاصل از برازش فرایند $Y - X$ به الگوی نهال‌های کاج

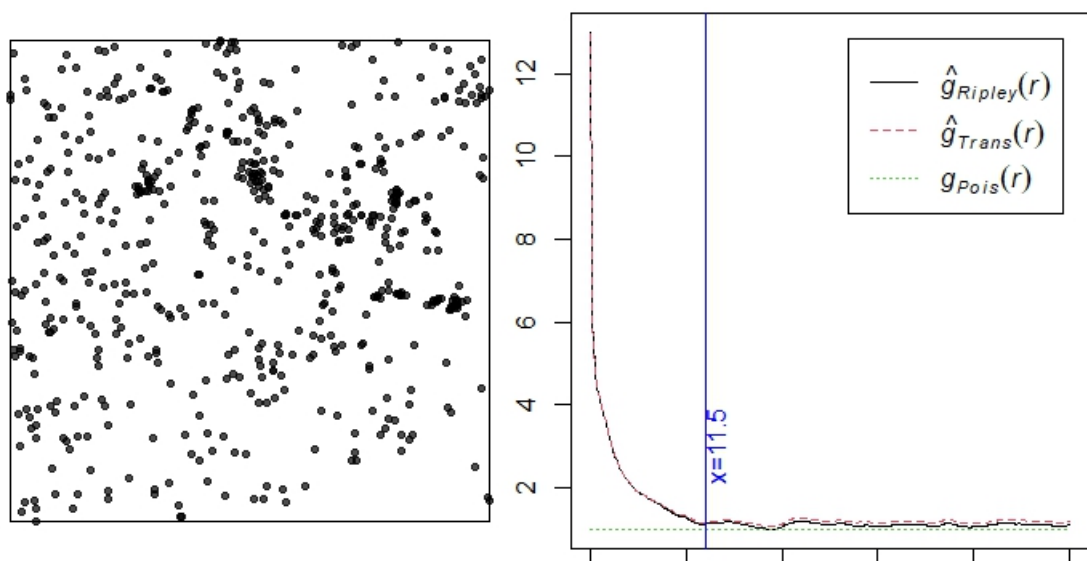
برای بررسی نیکویی برازش مدل برازش داده شده پوشش آماره‌های خلاصه بررسی می‌شود. از آن‌جایی که همه پارامترهای مدل در ساختار تابع شدت پالم وجود دارد پوشش تابع شدت پالم بررسی می‌شود همچنین پوشش تابع G که یکی از آماره‌های خلاصه مبتنی بر فاصله بین نقاط است نیز بررسی می‌شود. در شکل ۲ نمودار پوشش تابع شدت پالم رسم شده است که موید این است که فرایند برازش داده شده به خوبی الگوی نقطه‌ای فوق را مدل‌بندی می‌نماید. همچنین نمودار پوشش تابع G نیز رسم شده است. از آن‌جایی که روش



شکل ۲: نمودار پوشش تابع شدت پالم و تابع G برای مدل برازش داده شده

درست‌نمایی پالم یک روش تقریبی برآورد پارامترها است و این‌که، در اغلب نقاط تابع G مشاهده شده در فاصله اطمینان 0.95 قرار گرفته است لذا می‌توان نتیجه گرفت که مدل برازش داده شده با پارامترهای فوق برازش مناسبی به الگوی فوق دارد.

مثال ۲.۵. الگوی نقطه‌ای درختان کاج برگ بلند در ناحیه‌ای در جنوب جورجیا (ایالات متحده آمریکا) بر پنجره $[0, 200] \times [0, 200]$ محل قرار گرفتن و قطر 584 درخت کاج برگ بلند را نشان می‌دهد. در این‌جا صرفاً محل قرار گرفتن درخت‌های فوق بررسی می‌شود. این داده‌ها در پکیج *spatstat* نرم افزار \mathbb{R} با عنوان *longleaf* موجود است.



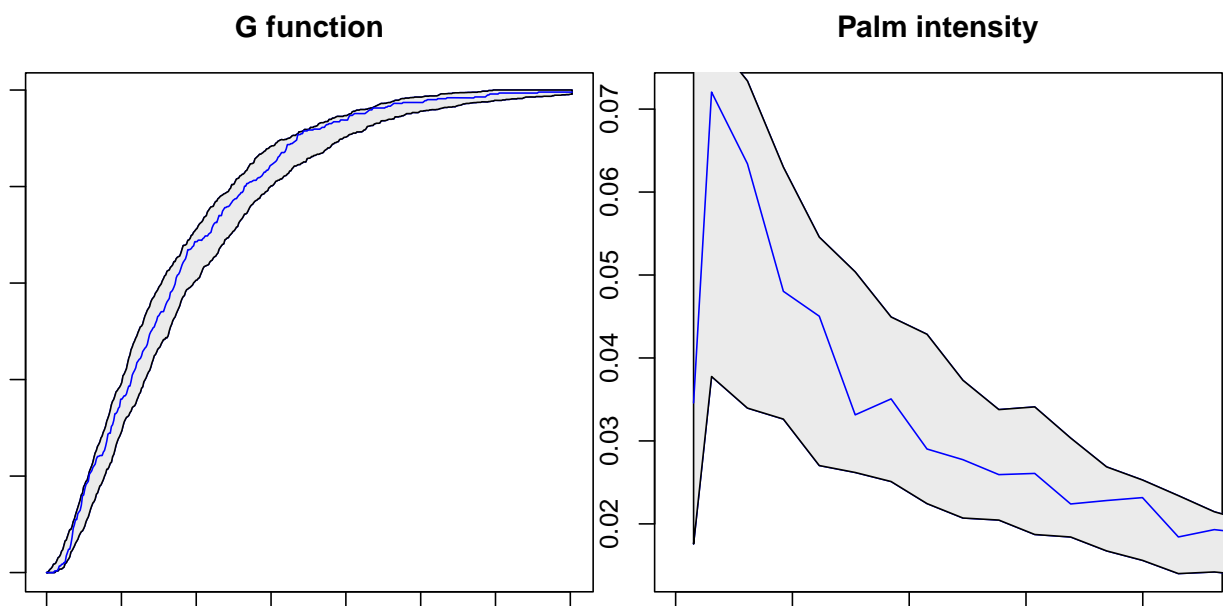
شکل ۳: نمودار الگوی نقطه‌ای محل قرار گرفتن درختان کاج برگ بلند و تابع همبستگی زوجی آن

در شکل ۳ نمودار محل قرار گرفتن درختان کاج و تابع همبستگی زوجی آن‌ها نشان داده شده. در شکل فوق به نظر می‌رسد بعضی نواحی دایره‌ای شکل تنک‌تر هستند. لذا برازش فرایند X به الگوی فوق پیشنهاد می‌شود. برای برازش برآورد پارامترهای مدل توسط روش درستنمایی پالم با توجه به رفتار تابع همبستگی زوجی در شکل ۳ و از آنجایی که برای $t \geq 11/5$ تابع همبستگی زوجی در اغلب نقاط کمتر از $1/1$ است لذا T در رابطه (۱.۴) برابر با $11/5$ در نظر می‌گیریم. با جایگذاری تابع شدت پالم فرایند X ، $\lambda_{o,X}(t)$ ، از ۳.۳ در تابع درستنمایی پالم (۱.۴) برآورد درستنمایی پارامترهای فرایند که در جدول ۲ ارائه شده اند بدست می‌آید.

پارامتر	δ	θ	p	r
برآورد	۵۴۴/۰	۱۰۹/۰	۳۴۹/۰	۵۱۱/۵

جدول ۲: برآوردگر درستنمایی پالم حاصل از برازش فرایند X به الگوی درختان کاج برگ بلند

از طرفی در فرایند X با توجه به رفتار تابع همبستگی زوجی براساس رابطه $\{g(2t) > 0 \text{ \& } g(2t) < 1 + \varepsilon\}$ $\hat{r}_T = \min_t$ پارامتر r را نیز می‌توان برآورد کرد. برای بهبود و تعدیل برآوردگر پارامتر r ، آن را می‌توان به کمک رابطه $\hat{r}_{imp} = a\hat{r}_{Palm} + (1-a)\hat{r}_T$ برآورد نمود. با در نظر گرفتن $\varepsilon = 1/1$ و $a = 4/5$ آن‌گاه $\hat{r}_{imp} = 5/654$



شکل ۴: نمودار پوشش برای تابع شدت پالم و تابع G

برای بررسی نیکویی برازش مدل برازش داده شده، پوشش آماره‌های خلاصه بررسی می‌شود. برای این منظور پوشش تابع شدت پالم و تابع G بررسی می‌شود. در شکل ۴ نمودار پوشش برای تابع شدت پالم و تابع G رسم شده است. پوشش آماره‌های خلاصه نشان‌دهنده این است که دلیلی برای عدم برازش مناسب مدل برازش داده شده وجود ندارد و مدل برازش داده شده عملکرد مناسبی دارد.

۶ بحث و نتیجه‌گیری

در این مقاله یک الگوریتم جدید تنک‌سازی ارائه شد. همچنین فرایند تنک‌شده و نقاط تنک شده حاصل از اعمال الگوریتم فوق بر فرایند پواسون همگن بررسی شد و روشی برای برازش مدل‌های فوق ارائه شد. برای مطالعات آتی می‌توان به مطالعه فرایند تنک‌شده حاصل از اعمال الگوریتم فوق بر سایر فرایندهای اولیه اقدام نمود. همچنین پیشنهاد می‌شود سایر روش‌های برآورد پارامترهای مدل بررسی و با روش درست‌نمایی پالم مقایسه گردد.

مراجع

- Stoyan, D. (1979), Interrupted point processes, *Biometrical Journal*, **21(7)**, 607-610.
- Lavancier, F. and Møller, J. (2016), Modelling aggregation on the large scale and regularity on the small scale in spatial point pattern datasets, *Scandinavian Journal of Statistics*, **43(2)**, 587-609.
- Møller, J., *Statistical inference and simulation for spatial point processes*, CRC Press.
- Tanaka, U., Y. Ogata, and D. Stoyan (2008), Parameter estimation and model selection for neyman-scott point processes, *Biometrical Journal: Journal of Mathematical Methods in Bio- sciences* , **50(1)**, 43-57.
- Neyman, J. and E. L. Scott (1958), Statistical approach to problems of cosmology, of the Royal Statistical Society: Series B (Methodological), **20(1)**, 1-29.

۷ ضمایم

قضیه ۱.۲ اثبات. (i) با استفاده از اینکه Z فرایند نقطه‌ای پواسون است به دست می‌آید. (ii) تساوی به ازای $\| \zeta_1 - \zeta_2 \| > 2r$ با استفاده از مستقل بودن تعداد نقاط در ناحیه‌های مجزا در فرایند نقطه‌ای پواسون بدست می‌آید. در حالتی که $\| \zeta_1 - \zeta_2 \| \leq 2r$ داریم:

$$\begin{aligned}
 P[\Pi(\zeta_1) = p^n, \Pi(\zeta_2) = p^m] &= P[N(A_{\zeta_1}) = n, N(A_{\zeta_2}) = m], \\
 &= \sum_{k=0}^{\min(m,n)} P[N(A_{\zeta_1} - A_{\zeta_2}) = n - k, N(A_{\zeta_1} \cap A_{\zeta_2}) = k, N(A_{\zeta_2} - A_{\zeta_1}) = m - k], \\
 &= \sum_{k=0}^{\min(m,n)} P[N(A_{\zeta_1} - A_{\zeta_2}) = n - k] P[N(A_{\zeta_1} \cap A_{\zeta_2}) = k] P[N(A_{\zeta_2} - A_{\zeta_1}) = m - k], \\
 &= \sum_{k=0}^{\min(m,n)} \frac{e^{-\theta \|A_{\zeta_1} - A_{\zeta_2}\|} [\theta \|A_{\zeta_1} - A_{\zeta_2}\|]^{n-k}}{(n-k)!} \\
 &\times \frac{e^{-\theta \|A_{\zeta_1} \cap A_{\zeta_2}\|} [\theta \|A_{\zeta_1} \cap A_{\zeta_2}\|]^k}{k!} \times \frac{e^{-\theta \|A_{\zeta_2} - A_{\zeta_1}\|} [\theta \|A_{\zeta_2} - A_{\zeta_1}\|]^{m-k}}{(m-k)!}, \\
 &= e^{-\theta [\pi r^2 - A(t,r)]} [\theta (\pi r^2 - A(t,r))]^{m+n} \sum_{k=0}^{\min(m,n)} \frac{1}{k!(n-k)!(m-k)!} \left[\frac{A(t,r)}{\theta (\pi r^2 - A(t,r))^2} \right]^k.
 \end{aligned} \tag{۱.۷}$$

□

تساوی آخر با کمی محاسبات جبری بدست می‌آید.

قضیه ۱.۳ اثبات. برای (i) با توجه به کاکس بودن فرایند X داریم:

$$\begin{aligned}
 \lambda(u) &= E[\delta \Pi(u)] = \delta E[\Pi(u)] = \delta \sum_{k=0}^{\infty} p^k P(\Pi(u) = p^k) = \delta \sum_{k=0}^{\infty} p^k \frac{e^{-\mu(A_u)} \mu(A_u)^k}{k!} \\
 &= \delta e^{-\mu(A_u)} \sum_{k=0}^{\infty} \frac{[p\mu(A_u)]^k}{k!} = \delta e^{-\mu(A_u)(1-p)} = \delta e^{-(1-p)\pi\theta r^2}
 \end{aligned} \tag{۲.۷}$$

در تساوی چهارم از رابطه (۱.۲) استفاده شده و $A_u = B(u, r)$. همچنین با توجه به اینکه $Z \sim \text{Poisson}(\theta)$ لذا $\mu(A_u) = \theta \pi r^2$ است.

برای (ii) با توجه به کاکس بودن فرایند X داریم:

$$\lambda^{(2)}(u, v) = E[\delta \Pi(u) \delta \Pi(v)] = \delta^2 E[\Pi(u) \Pi(v)]$$

در صورتی که $\| \zeta_1 - \zeta_2 \| \geq 2r$ از استقلال ناشی از الگوریتم تنک‌سازی داریم:

$$E[\Pi(\zeta_1) \Pi(\zeta_2)] = E[\Pi(\zeta_1)] E[\Pi(\zeta_2)].$$

و با توجه به قسمت (i) نتیجه حاصل می‌شود.

برای $\| \zeta_1 - \zeta_2 \| < 2r$ با استفاده از رابطه (۲.۲) داریم:

$$\begin{aligned}
 E[\Pi(\zeta_{\backslash})\Pi(\zeta_{\Upsilon})] &= \sum_{m,n=\circ}^{\infty} p^{n+m} P[\Pi(\zeta_{\backslash}) = p^n, \Pi(\zeta_{\Upsilon}) = p^m] \\
 &= \sum_{m,n=\circ}^{\infty} p^{n+m} e^{-\theta[\Upsilon\pi r^{\Upsilon}-A(t,r)]} [\theta(\pi r^{\Upsilon} - A(t,r))]^{m+n} \sum_{k=\circ}^{\min(m,n)} \frac{1}{k!(n-k)!(m-k)!} \left[\frac{A(t,r)}{\theta(\pi r^{\Upsilon} - A(t,r))^{\Upsilon}} \right]^k \\
 &= e^{-\theta[\Upsilon\pi r^{\Upsilon}-A(t,r)]} \sum_{m,n=\circ}^{\infty} p^{n+m} [\theta(\pi r^{\Upsilon} - A(t,r))]^{m+n} \sum_{k=\circ}^{\min(m,n)} \frac{1}{k!(n-k)!(m-k)!} \left[\frac{A(t,r)}{\theta(\pi r^{\Upsilon} - A(t,r))^{\Upsilon}} \right]^k \\
 &= e^{-\theta[\Upsilon\pi r^{\Upsilon}-A(t,r)]} \sum_{k=\circ}^{\infty} \sum_{n=k}^{\infty} \sum_{m=n}^{\infty} \frac{[p\theta(\pi r^{\Upsilon} - A(t,r))]^{m+n}}{k!(n-k)!(m-k)!} \left[\frac{A(t,r)}{\theta(\pi r^{\Upsilon} - A(t,r))^{\Upsilon}} \right]^k \\
 &\quad + e^{-\theta[\Upsilon\pi r^{\Upsilon}-A(t,r)]} \sum_{k=\circ}^{\infty} \sum_{m=k}^{\infty} \sum_{n=m+1}^{\infty} \frac{[p\theta(\pi r^{\Upsilon} - A(t,r))]^{m+n}}{k!(n-k)!(m-k)!} \left[\frac{A(t,r)}{\theta(\pi r^{\Upsilon} - A(t,r))^{\Upsilon}} \right]^k \\
 &= e^{-\theta[\Upsilon\pi r^{\Upsilon}-A(t,r)]} \sum_{k=\circ}^{\infty} \frac{1}{k!} [\theta p^{\Upsilon} A(t,r)]^k \sum_{n=\circ}^{\infty} \frac{[p\theta(\pi r^{\Upsilon} - A(t,r))]^n}{n!} \sum_{m=n}^{\infty} \frac{[p\theta(\pi r^{\Upsilon} - A(t,r))]^m}{m!} \\
 &\quad + e^{-\theta[\Upsilon\pi r^{\Upsilon}-A(t,r)]} \sum_{k=\circ}^{\infty} \frac{1}{k!} [\theta p^{\Upsilon} A(t,r)]^k \sum_{m=\circ}^{\infty} \frac{[p\theta(\pi r^{\Upsilon} - A(t,r))]^m}{m!} \sum_{n=m+1}^{\infty} \frac{[p\theta(\pi r^{\Upsilon} - A(t,r))]^n}{n!} \\
 &= e^{-\Upsilon\pi r^{\Upsilon}\theta(\backslash-p)+\theta A(t,r)(\backslash-p)^{\Upsilon}} \sum_{n=\circ}^{\infty} \frac{e^{-p\theta[\pi r^{\Upsilon}-A(t,r)]} [p\theta(\pi r^{\Upsilon} - A(t,r))]^n}{n!} \\
 &\quad \times \sum_{m=n}^{\infty} \frac{e^{-p\theta[\pi r^{\Upsilon}-A(t,r)]} [p\theta(\pi r^{\Upsilon} - A(t,r))]^m}{m!} \\
 &\quad + e^{-\Upsilon\pi r^{\Upsilon}\theta(\backslash-p)+\theta A(t,r)(\backslash-p)^{\Upsilon}} \sum_{m=\circ}^{\infty} \frac{e^{-p\theta[\pi r^{\Upsilon}-A(t,r)]} [p\theta(\pi r^{\Upsilon} - A(t,r))]^m}{m!} \\
 &\quad \times \sum_{n=m+1}^{\infty} \frac{e^{-p\theta[\pi r^{\Upsilon}-A(t,r)]} [p\theta(\pi r^{\Upsilon} - A(t,r))]^n}{n!},
 \end{aligned}$$

اگر $T_{\backslash}, T_{\Upsilon}$ متغیرهای تصادفی مستقل و هم توزیع با $pos(p\theta(\pi r^{\Upsilon} - A(t,r)))$ باشند آنگاه :

$$\begin{aligned}
 E[\Pi(\zeta_{\backslash})\Pi(\zeta_{\Upsilon})] &= e^{-\Upsilon\pi r^{\Upsilon}\theta(\backslash-p)+\theta A(t,r)(\backslash-p)^{\Upsilon}} P(T_{\backslash} \leq T_{\Upsilon}) + e^{-\Upsilon\pi r^{\Upsilon}\theta(\backslash-p)+\theta A(t,r)(\backslash-p)^{\Upsilon}} P(T_{\backslash} > T_{\Upsilon}), \\
 &= e^{-\Upsilon\pi r^{\Upsilon}\theta(\backslash-p)+\theta A(t,r)(\backslash-p)^{\Upsilon}} [P(T_{\backslash} \leq T_{\Upsilon}) + P(T_{\backslash} > T_{\Upsilon})] = e^{-\Upsilon\pi r^{\Upsilon}\theta(\backslash-p)+\theta A(t,r)(\backslash-p)^{\Upsilon}},
 \end{aligned}$$

و در نتیجه:

$$\lambda^{(\Upsilon)}(u, v) = \begin{cases} \delta^{\Upsilon} e^{-\Upsilon\pi r^{\Upsilon}\theta(\backslash-p)+\theta A(t,r)(\backslash-p)^{\Upsilon}}, & ||u - v|| \leq \Upsilon r, \\ \delta^{\Upsilon} e^{-\Upsilon(\backslash-p)\theta\pi r^{\Upsilon}}, & ||u - v|| > \Upsilon r. \end{cases}$$

□

قضیه ۲.۳ اثبات. قسمت (i) از قضیه ۱.۳ به دست می‌آید. برای (ii) داریم:

$$\begin{aligned} g_{Y-X}(u, v) &= \frac{E\{\delta[\mathbf{1} - \Pi(u)]\delta[\mathbf{1} - \Pi(v)]\}}{E\{\delta[\mathbf{1} - \Pi(u)]\}E\{\delta[\mathbf{1} - \Pi(v)]\}} = \frac{\mathbf{1} - \mathfrak{z}E[\Pi(u)] + E[\Pi(u)\Pi(v)]}{\{\mathbf{1} - E[\Pi(u)]\}^{\mathfrak{z}}} \\ &= \frac{\mathbf{1} - \mathfrak{z}E[\Pi(u)] + E[\Pi(u)]^{\mathfrak{z}} + E[\Pi(u)\Pi(v)] - E[\Pi(u)]^{\mathfrak{z}}}{\{\mathbf{1} - E[\Pi(u)]\}^{\mathfrak{z}}} = \mathbf{1} + \frac{E[\Pi(u)\Pi(v)] - E[\Pi(u)]E[\Pi(v)]}{\{\mathbf{1} - E[\Pi(u)]\}^{\mathfrak{z}}} \\ &= \mathbf{1} + \frac{E[\Pi(u)]^{\mathfrak{z}}}{\{\mathbf{1} - E[\Pi(u)]\}^{\mathfrak{z}}}[g_X(u, v) - \mathbf{1}] = \mathbf{1} + \frac{a^{\mathfrak{z}}}{(\mathbf{1} - a)^{\mathfrak{z}}}[g_X(t) - \mathbf{1}] \end{aligned}$$

□

A New Thinning Method of Spatial Point Processes Based on using Boolean Random Set Models

Reza Asadi¹, Mojtaba Khazaei¹, Mojtaba Ganjali¹

¹Department of Statistics, Shahid Beheshti University, Tehran, Iran

Abstract: Spatial point processes are statistical models for the analysis of point patterns. One of the methods of making new point processes is thinning the pre-existing point processes. In this paper, a new dependent thinning algorithm based on a Boolean random set is proposed. This algorithm can be considered as a generalization of the Stoyan thinning algorithm in the interrupted point process. Some properties of the thinned processes and eliminated point processes resulted from the application of the proposed thinning algorithm on the homogeneous Poisson point process, have been studied. Some summary statistics of the both thinned and eliminated point processes are calculated. Palm likelihood method is used to estimate parameters of these models. Finally, the proposed models are used to model two real datasets.

Keywords: Homogeneous Poisson point process; Cox process; Boolean random set; Dependent thinning; Cox process; Palm likelihood.

Mathematics Subject Classification (2020): xxAxx, xxBxx, xxCxx.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30-31, 2025
University of Kurdistan



آزمون نیکویی برازش تقارن توزیع‌های پیوسته تحت داده‌های از چپ بریده شده و از راست سانسور شده

معصومه اکبری^۱

گروه آمار، دانشگاه مازندران، بابلسر، ایران.

چکیده: در این مقاله، ابتدا یک آماره آزمون برای متقارن بودن توزیع‌های پیوسته تحت داده‌های از چپ بریده شده و از راست سانسور شده معرفی می‌شود. نشان داده می‌شود که این آماره آزمون مجاناً نرمال است. در یک مطالعه شبیه سازی گسترده، توان آزمون پیشنهادی مورد بررسی قرار گرفته است. نتایج گویای رضایت بخش بودن عملکرد آزمون است. در نهایت کاربردی بودن آماره آزمون در یک مثال واقعی نشان داده می‌شود.

واژه‌های کلیدی: آزمون نیکویی برازش، توزیع متقارن، داده‌های سانسور شده از راست و بریده شده از چپ.

کد موضوع‌بندی ریاضی (۲۰۲۰): 62G10، 62N01.

۱ مقدمه

فرض کنید X_1, X_2, \dots, X_n یک نمونه تصادفی به حجم n ، از یک تابع چگالی ناشناخته f با میانگین متناهی μ باشد. گفته می‌شود چگالی f نسبت به μ متقارن است اگر و فقط اگر

$$H_0: f(\mu - x) = f(\mu + x), \quad \forall x \in \mathbb{R}. \quad (1.1)$$

بررسی فرض تقارن یکی از موضوعات اساسی در آمار ناپارامتری است و کاربردهای گسترده‌ای در شاخه‌های مختلف علمی دارد. توجه به تقارن توزیع داده‌ها از دیرباز مورد توجه پژوهشگران بوده، چراکه بسیاری از آزمون‌ها و روش‌های آماری کلاسیک بر مبنای این فرض توسعه یافته‌اند. اولین مطالعات در زمینه آزمون تقارن به دهه‌های میانی قرن بیستم بازمی‌گردد، زمانی که فرضیه تقارن به عنوان معیاری برای اعتبارسنجی مدل‌های نرمال در داده‌های تجربی مورد استفاده قرار گرفت. در آن دوران، آزمون‌هایی مبتنی بر مرتب‌سازی داده‌ها، توزیع رتبه‌ها، و مقایسه میانه و میانگین برای آزمون تقارن توسعه یافتند.

^۱ سخنران، m.akbari@umz.ac.ir

در سال‌های بعد، با پیچیده‌تر شدن ساختار داده‌ها، روش‌های آماری نیز تکامل یافتند و آزمون‌های متنوع‌تری برای بررسی تقارن توسعه داده شد. به عنوان نمونه، آزمون‌های مبتنی بر آمار ناپارامتری، آزمون‌های استوار، و روش‌هایی مبتنی بر فاصله و تبدیل داده‌ها پیشنهاد شدند. این آزمون‌ها نه تنها در تحلیل داده‌های کلاسیک، بلکه در زمینه‌هایی مانند اقتصاد، ژنتیک، پزشکی و مهندسی نیز به کار گرفته شده‌اند.

در مطالعات مالی، فرض تقارن نقش کلیدی در مدل‌سازی بازده دارایی‌ها و قیمت‌گذاری ابزارهای مالی ایفا می‌کند؛ به عنوان نمونه می‌توان به مدل قیمت‌گذاری دارایی‌های سرمایه‌ای شارپ-لینتر و مدل بلک-شولز برای قیمت‌گذاری اختیار معامله اشاره کرد که در آن‌ها تقارن توزیع بازده دارایی‌ها یا خطاها فرض می‌شود. از سوی دیگر، در روش‌های پارامتری و استوار آماری، تقارن اغلب به عنوان فرض پایه‌ای برای خطاها یا متغیرهای پاسخ در مدل‌هایی چون رگرسیون خطی و تحلیل واریانس تلقی می‌شود.

با توجه به اهمیت این موضوع، پژوهش‌های متعددی به آزمون فرضیه H_0 در معادله (۱.۱) پرداخته‌اند. از جمله می‌توان به مطالعات رندلز و همکاران (۱۹۸۰)، مک ویلیامز (۱۹۹۰)، مدرس و گستیرت (۱۹۹۰)، گینز و چاکرابتی (۱۹۹۲)، تاجدین (۱۹۹۴)، باکلیزی (۲۰۰۳، ۲۰۰۷ و ۲۰۰۸)، چنگ و بالاکریشن (۲۰۰۴)، کورزو و باباتیوا (۲۰۱۳) و شیانگ و همکاران (۲۰۲۱) اشاره کرد که همگی آزمون‌هایی برای داده‌های کامل پیشنهاد داده‌اند.

افزون بر مطالعات کلاسیک، در سال‌های اخیر تلاش‌های متعددی برای توسعه آزمون‌های تقارن با ویژگی‌های آماری بهتر یا سازگار با شرایط خاص داده‌ها صورت گرفته است. به عنوان نمونه، ژو و ونگ (۲۰۱۵) آزمونی مبتنی بر نسبت درست‌نمایی تجربی ارائه دادند که ویژگی‌های توان مناسبی دارد. جوریکوا و کالینا (۲۰۱۶) آزمون‌هایی برای تحلیل تقارن در داده‌های اقتصادی ارائه کردند. همچنین گوش و سن (۲۰۱۷) نیز طبقه‌ای از آزمون‌های ناپارامتری مبتنی بر U -آماره‌ها برای بررسی تقارن توسعه دادند. در شرایطی که داده‌ها ناقص باشند، مطالعات محدودی در این حوزه انجام شده‌است. به عنوان مثال، شیانگ و همکاران (۲۰۱۹) آزمونی برای تقارن در حضور داده‌های گمشده پیشنهاد کردند و اخیراً کائو و ژو (۲۰۲۱) آزمون‌های ناپارامتری جدیدی برای داده‌های سانسور شده تصادفی توسعه داده‌اند. اکبری و زمانی (۱۴۰۳) آزمونی برای سنجش تقارن در توزیع‌های پیوسته در حضور داده‌های سانسور شده از راست معرفی کرده‌اند.

در این مقاله، به گسترش این موضوع پرداخته و یک آزمون جدید برای بررسی تقارن توزیع پیوسته در شرایط وجود داده‌های بریده شده از چپ و سانسور شده از راست پیشنهاد می‌دهیم. ساختار مقاله به صورت زیر است: در بخش ۲، آماره آزمون معرفی می‌شود. از آنجا که توزیع دقیق این آماره به راحتی در دسترس نیست، در بخش ۳ با استفاده از روش مونت کارلو مقادیر بحرانی و توان آزمون به صورت تقریبی محاسبه می‌شود. نهایتاً، در بخش ۴ کاربرد عملی آزمون با استفاده از یک مجموعه داده واقعی توضیح داده می‌شود.

۲ آماره آزمون

فرض کنید X_1, X_2, \dots, X_n یک نمونه تصادفی از توزیع پیوسته F با میانه نامعلوم θ باشد. اگر i, j و k سه مقدار مجزا متعلق به مجموعه $\{1, 2, \dots, n\}$ باشد، در این صورت مشاهدات سه تایی (X_i, X_j, X_k) تشکیل یک سه تایی راست‌گرا^۱ را می‌دهند، بدین

^۱right triple

معنا که توزیع بنظر می‌رسد که به راست چوله است اگر مشاهده میانی به کوچکترین مشاهده نزدیکتر باشد. بطور مشخص تر فرض کنید

$$f^*(X_i, X_j, X_k) = \frac{1}{3} \{ \text{sign}(X_i + X_j - 2X_k) + \text{sign}(X_i + X_k - 2X_j) + \text{sign}(X_j + X_k - 2X_i) \}, \quad (1.2)$$

که در آن $\text{sign}(u) = -1, 0, 1$ به ترتیب برای مقادیر $u < 0, u = 0, u > 0$. با توجه به معادله (۱.۲)، (X_i, X_j, X_k) تشکیل یک سه تایی راست گرا را می‌دهد اگر $f^*(X_i, X_j, X_k) = \frac{1}{3}$ و بطور مشابه تشکیل یک سه تایی چپ گرا را می‌دهد اگر $f^*(X_i, X_j, X_k) = -\frac{1}{3}$ و زمانی که $f^*(X_i, X_j, X_k) = 0$ باشد بیانگر آن است که سه تایی نه راست گرا و نه چپ گراست. قابل توجه است که $f^*(X_i, X_j, X_k)$ فقط می‌تواند مقادیر $\frac{1}{3}, 0, -\frac{1}{3}$ را اختیار کند. بر همین اساس رندلز و همکاران (۱۹۸۰) یک آزمون بر اساس U -آماره

$$\hat{\eta} = \frac{1}{\binom{n}{3}} \sum_{i < j < k} f^*(X_i, X_j, X_k), \quad (2.2)$$

پیشنهاد کردند که مقادیر بزرگ $|\hat{\eta}|$ منجر به رد فرضیه تقارن می‌شود. در ادامه نحوه لحاظ کردن مشاهدات از چپ بریده شده و از راست سانسور شده را در آماره آزمون (۲.۲) مورد بررسی قرار می‌دهیم.

فرض کنید که X, C, L به ترتیب با تابع توزیع‌های F, G, H بیانگر متغیرهای تصادفی طول عمر، زمان سانسور و زمان برش از چپ باشند و $\bar{F}(x) = 1 - F(x), \bar{G}(x) = 1 - G(x)$ و $\bar{H}(x) = 1 - H(x)$. تحت سانسور از راست و برش از چپ، n نمونه مستقل و هم توزیع $(T_i \epsilon_i, \delta_i)$ از $(T \epsilon, \delta)$ در نظر گرفته می‌شود که $T = \min(X, C)$ و $\delta = I(X < C)$ و $\epsilon = I(T > L)$. به عبارت دیگر، δ نشانه سانسور از راست و ϵ برای مشخص کردن برش استفاده شده است. واضح است که T_i یک مشاهده برای i مین نمونه است اگر $T_i > L_i$ باشد. بر اساس داتا (۲۰۱۰) می‌توان براساس داده‌های از راست سانسور شده و از چپ بریده شده، یک نسخه از آماره (۲.۲) بصورت زیر در نظر گرفت.

$$\hat{\eta}_{LTRC} = \frac{1}{\binom{n}{3}} \sum_{1 \leq i < j < k \leq n} \frac{f^*(T_i, T_j, T_k) \delta_i \delta_j \delta_k}{\hat{K}_C(T_i) \hat{K}_C(T_j) \hat{K}_C(T_k)},$$

که در آن $\hat{K}_C(\tau)$ براوردگری از تابع بقا زمان سانسور C تحت برش چپ است و بصورت زیر می‌باشد.

$$\hat{K}_C(\tau) = \prod_{t \leq \tau} \left(1 - \frac{dN^C(t)}{Y(t)} \right),$$

که در آن

$$dN_i^C(t) = I(T_i = t, \delta_i = 0),$$

$$Y(t) = \sum_{i=1}^n Y_i(t), \quad \text{and} \quad Y_i(t) = I(T_i \geq t \geq L_i). \quad (3.2)$$

فرض H_0 یعنی متقارن بودن توزیع پیوسته، به ازای مقادیر بزرگ $|\hat{\eta}_{LTRC}|$ رد می‌شود. تحت شرایطی می‌توان نشان داد که آماره آزمون $\hat{\eta}_{LTRC}$ مجاناً نرمال است.

۳ مقادیر بحرانی و توان آزمون

در این بخش، بر اساس شبیه‌سازی مونت کارلو، عملکرد آزمون پیشنهادی مورد بررسی قرار می‌گیرد. همه محاسبات بر اساس نرم افزار R انجام شده است. بدین منظور در ابتدا، مقادیر بحرانی به ازای مقادیر $n = ۲۰, ۵۰, ۷۰$ تقریب زده شده است در حالی که، دو مقدار $p = ۰/۱, ۰/۴$ برای سانسور از راست یعنی $P(X > C) = p$ و $C \sim Exp(\lambda)$ و مستقل از X ، در نظر گرفته شده است. علاوه بر این توزیع متغیر تصادفی برش از چپ بصورت، $L \sim Exp(\mu)$ فرض شده به قسمی که $P(L > X) = ۰/۲$ گردد. نتایج در سطح خطای نوع اول، $\alpha = ۰/۰۵$ ، بصورت زیر در جدول ۱ گزارش شده است. برای بررسی توان آزمون پیشنهاد شده، توزیع‌های

جدول ۱: مقادیر بحرانی آزمون $\hat{\eta}_{LTRC}$

n			
۷۰	۵۰	۲۰	p
$۰/۰۳۲$	$۰/۰۳۹$	$۰/۰۶۸$	$۰/۱$
$۰/۰۱۸$	$۰/۰۲۳$	$۰/۰۳۵$	$۰/۴$

مقابل از خانواده توزیع لامبدای تعمیم یافته (GLD^2) با تابع چنک زیر در نظر گرفته شده است.

$$F^{-1}(u) = \lambda_1 + \frac{u^{\lambda_2} - (1-u)^{\lambda_2}}{\lambda_2}, \quad 0 < u < 1.$$

در جدول ۲، ۹ حالت از GLD به ازای پارامترهای مختلف آن بعنوان ۹ توزیع مقابل نشان داده شده است. جدول ۳ توان آزمون

جدول ۲: ۹ حالت از GLD

λ_1	λ_2	λ_3	λ_4	Case
$۰/۰۰۰۰$	$۰/۱۹۷۵$	$۰/۱۳۴۹$	$۰/۱۳۴۹$	۱
$-۰/۱۱۶۷$	$-۰/۳۵۱۷$	$-۰/۱۳۰۰$	$-۰/۱۶۰۰$	۲
$۰/۰۰۰۰$	$-۱/۰۰۰۰$	$-۰/۱۰۰۰$	$-۰/۱۸۰۰$	۳
$۳/۵۸۶۵$	$۰/۰۴۳۱$	$۰/۰۲۵۲$	$۰/۰۹۴۰$	۴
$۰/۰۰۰۰$	$-۱/۰۰۰۰$	$-۰/۰۰۷۵$	$۰/۳۰۰۰$	۵
$۰/۰۰۰۰$	$۱/۰۰۰۰$	$۱/۴۰۰۰$	$۰/۰۲۵۰۰$	۶
$۰/۰۰۰۰$	$۱/۰۰۰۰$	$۰/۰۰۰۱$	$۰/۱۰۰۰$	۷
$۰/۰۰۰۰$	$-۱/۰۰۰۰$	$-۰/۰۰۱۰$	$-۰/۱۳۰۰$	۸
$۰/۰۰۰۰$	$-۱/۰۰۰۰$	$-۰/۰۰۰۱$	$-۰/۱۷۰۰$	۹

محاسبه شده $\hat{\eta}_{LTRC}$ را گزارش می‌دهد. همان‌طور که مشاهده می‌شود به ازای افزایش حجم نمونه، توان آزمون‌ها افزایش یافته و توان‌ها زمانی که درجه سانسور کوچک است ($P = ۰/۱$) به نسبت بزرگتر می‌باشد. در مجموع توان آزمون قابل قبول بوده و در اکثر موارد جز در توزیع مقابل حالت دوم، توان خوبی بدست آمده است.

از طرفی حالت اول از توزیع مقابل، همان توزیع نرمال استاندارد است لذا توان آزمون در این مورد، همان خطای نوع اول را نشان می‌دهد که مقادیر گزارش شده حول مقدار $۰/۰۵$ همان خطای نوع اول اسمی در نظر گرفته شده است. بنابراین می‌توان گفت برآوردهای خطای نوع اول نیز تحت کنترل و قابل قبول است.

²generalized lambda distribution

جدول ۳: توان آزمون									
Case ۹	Case ۸	Case ۷	Case ۶	Case ۵	Case ۴	Case ۳	Case ۲	Case ۱	n
$p = ۰/۱$									
۰/۶۳۴۸	۰/۹۰۰۸	۰/۵۸۳۴	۰/۲۱۶۴	۰/۲۵۶۲	۰/۱۷۰۴	۰/۱۰۶۶	۰/۰۷۸	۰/۰۵۰۴	۲۰
۰/۹۸۶	۰/۹۹۹۶	۰/۹۶۹۷	۰/۵۱۶	۰/۶۲۶۸	-۰/۴۳۹۸	-۰/۲۰۰۴	-۰/۰۶۹	-۰/۰۴۴۸	۵۰
۰/۹۹۵	۰/۹۹۹۸	۰/۹۸۶۸	۰/۶۵۵	۰/۷۵۸۸	۰/۵۹۵۸	۰/۳۱۱	۰/۰۸۱	۰/۰۶۳۶	۷۰
$p = ۰/۴$									
۰/۲۴۷۶	۰/۲۴۲۶	۰/۲۲۰۴	۰/۳۳۰۶	۰/۰۵	۰/۰۵۷۶	۰/۰۴۴۸	۰/۱۲۷۲	۰/۰۴۹۴	۲۰
۰/۵۲۹۲	۰/۵۰۳۲	۰/۴۵۶۲	۰/۴۶۲۸	۰/۰۳۶۸	۰/۰۵۷۶	۰/۰۲۱۶	۰/۲۰۸۸	۰/۰۵۴	۵۰
۰/۷۶۸۴	۰/۷۳۳۶	۰/۶۷۵۴	۰/۵۹۱۸	۰/۰۶۲	۰/۱۰۸۲	۰/۰۳۸۶	۰/۳۶۹۴	۰/۰۹۷۶	۷۰

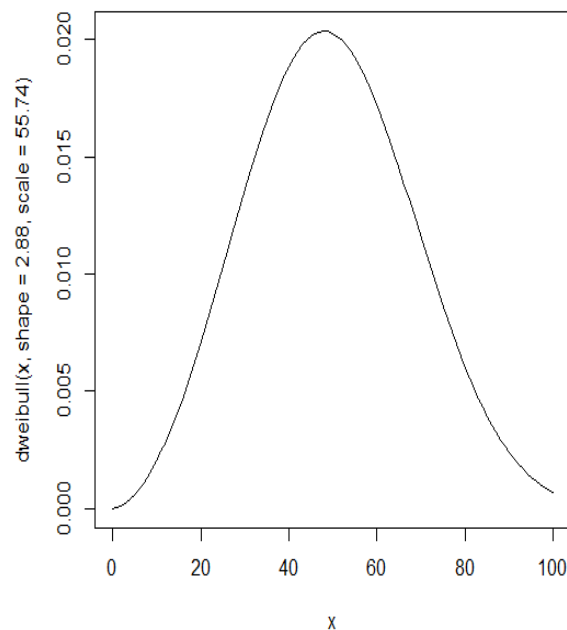
۴ مثال واقعی

برای تشریح نحوه بکار بردن آزمون پیشنهادی، داده‌های از چپ بریده شده و از راست سانسور شده **هانگ و همکاران (۲۰۰۹)** را بکار می‌بریم که این مجموعه داده شامل طول عمر ۷۱۰ ترانسفورماتور یک شرکت انرژی با ۶۲ شکست است. اگرچه داده‌های اصلی موجود نیست اما یک زیر مجموعه ۲۸۶ تایی با ۳۹ شکست در **امورا و شیو (۲۰۱۶)** آمده است. داده‌ها بین سال‌های ۱۹۸۰ الی ۲۰۰۸ جمع‌آوری شده است. از این رو ترانسفورماتورهایی که قبل از ۱۹۸۰ نصب شده‌اند بعنوان داده‌های از چپ بریده شده در نظر گرفته می‌شوند. بنابراین داده‌ها از چپ بریده شده و از راست سانسور شده است. نرخ سانسور راست و برش از چپ به ترتیب ۸۶/۴ و ۵۸/۳۹ درصد هستند. **هانگ و همکاران (۲۰۰۹)** یک توزیع وایبل با پارامتر شکل ۲/۸۸ و پارامتر مقیاس ۵۵/۷۴ به این مجموعه داده‌ها برازش داده‌اند. نمودار چگالی این توزیع در شکل زیر آمده است و به نوعی نشان می‌دهد که توزیع مجموعه داده موردنظر، متقارن است. آماره آزمون پیشنهادی مشاهده شده بصورت $\hat{\eta}_{LTRC} = ۰/۰۶۲۵$ بدست آمده است و برای بدست آوردن p -مقدار یک نمونه به حجم $n = ۵۰$ از توزیع وایبل با پارامترهای ذکر شده تولید کرده و $L \sim Exp(\mu)$ و $C \sim Exp(\lambda)$ در نظر می‌گیریم. مقادیر μ و λ را طوری تعیین می‌کنیم که به نرخ‌های به ترتیب ۵۸/۲۹ و ۸۶/۴ منتج شود. بر اساس هر تکرار یک مقدار از آماره $\hat{\eta}_{LTRC}$ بدست خواهد آمد و براساس ۵۰۰۰ تکرار، میانگین تعداد دفعاتی که $\hat{\eta}_{LTRC}$ بیشتر از آماره آزمون مشاهده شده $۰/۰۶۲۵$ باشد را به عنوان p -مقدار در نظر خواهیم گرفت. با توجه به مراحل مذکور، p -مقدار بدست آمده برابر با ۰/۷ بدست آمد که حاکی از آن است که توزیع مجموعه داده موردنظر متقارن است.

مراجع

اکبری، م. و زمانی، ز. (۱۴۰۳)، معرفی آزمون جدید برای متقارن بودن یک توزیع در حضور داده‌های سانسور شده، هشتمین همایش ریاضیات و علوم انسانی ریاضیات مالی، ۲۳ و ۲۴ آبان ۱۴۰۳، دانشگاه علامه طباطبائی، تهران.

Baklizi A. (2003), A conditional distribution free runs test for symmetry, *Journal of Nonparametric Statistics*, **15**, 713–718.



شکل ۱: چگالی توزیع وایبل

Baklizi A. (2007), Testing symmetry using a trimmed longest run statistics, *Australian and New Zealand Journal of Statistics*, **49**, 339–347.

Baklizi A. (2008), Improving the power of the Hybrid test, *International Journal of Contemporary Mathematical Sciences*, **3**, 497–499.

Cao Y. and Zhu W. (2021). Nonparametric tests for symmetry under random censoring. *Statistics & Probability Letters*, **169**, 108–967.

Cheng W. and Balakrishnan N. (2004), A modified sign test for symmetry, *Communications in Statistics-Simulation and Computation*, **33**, 703–709.

Corzo J. and Babativa G. (2013), A modified runs test for symmetry, *Journal of Statistical Computation and Simulation*, **83**, 984–991.

Datta S., Bandyopadhyay D., and Satten, G. A. (2010), Inverse probability of censoring weighted U-statistics for right-censored data with an application to testing testing hypotheses, *Scandinavian Journal of Statistics*, **37**, 680–700.

Emura T., and Shiu v. (2016), Estimation and model selection for left-truncated and right-censored lifetime

- data with application to electric power transformers analysis, *Communication in Statistics: Simulation and Computation*, **45**, 171–318.
- Gibbons C., and Chakraborti S. (1992), *Nonparametric Statistical Inference*, New York: Marcel Dekker.
- Ghosh A. K., and Sen A. (2017). A general class of tests for symmetry based on U-statistics. *Journal of Nonparametric Statistics*, **29**, 1–18.
- Hong Y., Meeker W.Q. and McCalley J.D. (2009), Prediction of remaining life of power transformers based on left truncated and right censored lifetime data, *Annals of Applied Statistics*, **3** 3, 857–879.
- Jurečková J., and Kalina J. (2016). Tests for symmetry of distributions with application in economics. *Statistical Papers*, **57**, 197–214.
- McWilliams P. (1990), A distribution-free test for symmetry based on a runs statistic, *Journal of the American Statistical Association*, **85**, 1130–1133.
- Modarres R. and Gastwirth J. (1998), Hybrid test for the hypothesis of symmetry, *Journal of Applied Statistics*, **25**, 777-783.
- Randles R. H., Fligner M. A., Policello G. E., and Wolfe D. A. (1980), An asymptotically distribution-free test for symmetry versus asymmetry, *Journal of the American Statistical Association*, **75**, 168-172.
- Tajuddin I. (1994), Distribution-free test for symmetry based on Wilcoxon two-sample test, *Journal of Applied Statistics*, **21**, 409–415.
- Xiong J., Zhou X., and Xu J. (2019). A test for symmetry of a distribution with missing data. *Communications in Statistics - Theory and Methods*, **48**, 1682–1695.
- Xiong P., Zhuang W., and Qiu G. (2021), Testing symmetry based on the entropy of record values, *Journal of Nonparametric Statistics*, **33**, 134–155.
- Zou C., and Wang Z. (2015). A new test for symmetry based on the empirical likelihood ratio. *Journal of Statistical Computation and Simulation*, **85**, 348–358.

A Symmetry Goodness-of-Fit Test for Continuous Distributions under Left Truncation and Right Censoring

Masoumeh Akbari

Department of Statistics, University of Mazandaran, Babolsar, Iran.

Abstract: In this paper, a test statistic is first introduced for assessing the symmetry of continuous distributions under right-censored and left-truncated data. It is shown that the proposed test statistic is asymptotically normal. In an extensive simulation study, the power of the proposed test is examined. The results indicate that the test performs satisfactorily. Finally, the applicability of the test statistic is demonstrated through a real data example.

Keywords: Goodness-of-fit test, Symmetric distribution, Left-truncated and right-censored data.

Mathematics Subject Classification (2020): 62N01, 62G10.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



یک کلاس از مدل‌بندی سری زمانی دوخطی گسسته مقدار با کاربردهای آن

رعنا بامدادی^۱ و مهرناز محمدپور

گروه آمار، دانشگاه مازندران

چکیده: در این مطالعه، یک مدل تعدیل یافته دوخطی گسسته مقدار مرتبه دو برای تحلیل سری‌های زمانی غیرخطی معرفی می‌شود. این مدل از انعطاف‌پذیری بالایی برخوردار است به گونه‌ای که ورودی‌های جدید می‌توانند به طور مستقل رفتار کنند. مدل‌بندی ارائه شده برای فرایندهای بیش‌پراکنده به کار می‌رود که بر یک چارچوب غیرخطی تکیه دارند. ویژگی‌های اصلی مدل پیشنهادی به طور نظری بررسی شده است. همچنین، به منظور استفاده عملی از مدل، چند روش برای برآورد پارامترها ارائه و به منظور ارزیابی عملکرد این روش‌ها شبیه‌سازی انجام شده است. در پایان، برای نشان دادن کاربرد مدل، به تحلیل داده‌ها پرداخته شده است. واژه‌های کلیدی: پیش‌بینی، مدل تعدیل یافته دوخطی گسسته مقدار، تقریب نقطه زینی، روش ویتل. کد موضوع بندی ریاضی (۲۰۲۰): 60G25, 62M10, 62P10.

۱ مقدمه

در سال‌های اخیر، مدل‌بندی سری‌های زمانی شمارشی به دلیل کاربردهای وسیع آن مورد توجه قرار گرفته و توسعه‌ی چشمگیری یافته است. این نوع مدل‌بندی در بسیاری از حوزه‌های علمی، از جمله بیمه، اقتصاد، داروسازی، اپیدمیولوژی، سیستم‌های صف، ارتباطات، هواشناسی، مخابرات و غیره کاربردهای گسترده‌ای دارند و ضرورت دارد که برازش مناسب بر روی این داده‌ها صورت گیرد. برای مدل‌سازی سری‌های زمانی شمارشی، به کارگیری عملگرهای رقیق‌ساز^۱ به عنوان روشی مؤثر شناخته شده‌اند. عملگرهای رقیق‌ساز مختلفی توسط ویب (۲۰۰۸) و اسکات و همکاران (۲۰۱۵) معرفی شده‌اند. جهت مدل‌بندی داده‌هایی با این خصوصیت ژانگ و همکاران (۲۰۱۰) عملگر رقیق‌ساز سری توانی تعمیم یافته علامت^۲ به صورت زیر

$$\alpha \circledast Y = \text{sign}(\alpha) \text{sign}(Y) \sum_{i=1}^{|Y|} W_i,$$

^۱ سخنران، bamdadirana@gmail.com

^۱ Thinning operator

^۲ Signed Generalized Power Series Thinning Operator

معرفی کرده‌اند، که در آن، دنباله‌ی شمارشی W_i از متغیرهای تصادفی مستقل و هم‌توزیع با توزیع توان تعمیم‌یافته^۳ در نظر گرفته می‌شود. با اینکه در حوزه مدل‌سازی سری‌های زمانی شمارشی پیشرفت‌های زیادی حاصل شده، ولی هنوز مدل‌های غیرخطی در این زمینه کمتر مورد توجه و بررسی قرار گرفته‌اند. در یکی از تلاش‌های برجسته در این حوزه، **دوقان و همکاران (۲۰۰۶)** مدل‌های دوخطی را با استفاده از عملگر رقیق‌ساز گسترش دادند. یک رویکرد جایگزین برای این مدل‌ها توسط **بامدادی و همکاران (۲۰۲۴)** ارائه شده است. آن‌ها مدل **دوقان و همکاران (۲۰۰۶)** را اصلاح کرده و مدل دوخطی را معرفی کردند که انعطاف‌پذیری بیشتری برای ورودی‌های جدید ارائه می‌دهد. این مدل از این مزیت برخوردار است که افراد جدید می‌توانند به‌طور مستقل رفتار کنند. در این مقاله، مدل فوق با مرتبه ۲ برای فرایندهایی که تابع خودهمبستگی جزئی آن در تاخیر ۲ به‌طور معنی‌داری از صفر اختلاف دارد بررسی می‌شود.

۲ معرفی مدل و برخی ویژگی‌های آن

فرآیند $\{X_t\}_{t \in \mathbb{Z}}$ ، را به‌صورت زیر در نظر بگیرید

$$X_t = a \circledast X_{t-2} + (a \circledast X_{t-2})(b \circledast \varepsilon_{t-1}) + \varepsilon_t, \quad (1.2)$$

که در آن \circledast عملگر رقیق‌ساز سری توانی تعمیم‌یافته علامت می‌باشد، $a, b \in [-1, 1]$ و $\{\varepsilon_t\}$ دنباله‌ای از متغیرهای تصادفی گسسته i.i.d. با میانگین μ_ε و واریانس σ_ε^2 هستند که از $\{X_s\}_{s < t}$ مستقل می‌باشند. همچنین، $\{W_i\}$ و $\{\tilde{W}_i\}$ که در عملگرهای $a \circledast$ و $b \circledast$ به‌کار رفته‌اند، دارای توزیع سری توانی با میانگین‌های $|a|$ و $|b|$ و واریانس‌های α و β هستند. این مدل به عنوان مدل تعدیل‌یافته دوخطی گسسته مقدار مثبتی بر عملگر رقیق‌ساز از مرتبه دو یا به اختصار $\text{MINBL}(2, 0, 2, 1)$ شناخته می‌شود.

گزاره ۱.۲. اگر $|a| + |ab|\mu_{|\varepsilon|} + |ab|^2\sigma_\varepsilon^2 + |a|^2\beta\mu_{|\varepsilon|} < 1$ ، آنگاه فرآیند $\{X_t\}_{t \in \mathbb{Z}}$ در **(۱.۲)** ایستا می‌باشد.

گزاره ۲.۲. فرض کنید $\{X_t\}_{t \in \mathbb{Z}}$ فرآیند $\text{MINBL}(2, 0, 2, 1)$ معرفی شده در **(۱.۲)** باشد که شرایط گزاره **۱.۲** برقرار باشد. در این صورت، امیدریاضی و تابع خودکواریانس فرآیند به‌صورت زیر می‌باشد

$$E(X_t) = \frac{\mu_\varepsilon}{1 - (a + ab\mu_\varepsilon)}, \quad (2.2)$$

و

$$\gamma_X(1) = \frac{ab\sigma_\varepsilon^2 E(X_t)}{1 - (a + ab\mu_\varepsilon)}, \quad \gamma_X(2k) = (a + ab\mu_\varepsilon)^k \gamma_X(0), \quad \gamma_X(2k+1) = (a + ab\mu_\varepsilon)^k \gamma_X(1), \quad k \geq 1.$$

گزاره ۳.۲. میانگین شرطی فرآیند $\{X_t\}$ برای k -گام جلوتر به‌صورت زیر است

$$\begin{aligned} E(X_{t+1}|t) &= (a + ab\varepsilon_t)X_{t-1} + \mu_\varepsilon, \\ E(X_{t+2}|t) &= (a + ab\mu_\varepsilon)X_t + \mu_\varepsilon, \\ E(X_{t+k}|t) &= (a + ab\mu_\varepsilon)^{\frac{k}{2}}X_t + \mu_\varepsilon \sum_{i=1}^{\frac{k}{2}} (a + ab\mu_\varepsilon)^{i-1}, \quad k = 2h, \end{aligned} \quad (3.2)$$

$$E(X_{t+k}|t) = (a + ab\mu_\varepsilon)^{\lceil \frac{k}{2} \rceil} E(X_{t+1}|t) + \mu_\varepsilon \sum_{i=1}^{\lceil \frac{k}{2} \rceil} (a + ab\mu_\varepsilon)^{i-1}, \quad k = 2h + 1. \quad (4.2)$$

³ Generalized Power Series

گزاره ۴.۲. واریانس شرطی یک گام به جلو فرآیند $\{X_t\}$ به صورت

$$Var(X_{t+1}|t) = a^2\beta|\varepsilon_t|X_{t-1}^2 + (\alpha + 2\alpha b\varepsilon_t + \alpha b^2\varepsilon_t^2 + \alpha\beta|\varepsilon_t|)|X_{t-1}| + \sigma_\varepsilon^2$$

است.

۳ برآوردیابی و شبیه‌سازی

در این بخش، به برآورد پارامترهای مدل $MINBL(2, 0, 2, 1)$ با استفاده از چهار روش مختلف می‌پردازیم. برای ارزیابی کارایی این روش‌ها، از داده‌های شبیه‌سازی شده استفاده شده و نتایج حاصل از هر روش با یکدیگر مقایسه می‌شوند. در روند برآورد، توزیع $\{\varepsilon_t\}$ پواسن با پارامتر μ_ε در نظر گرفته شده است.

۱.۳ روش یول‌واکر

با به‌کارگیری گزاره ۲.۲، امید ریاضی و تابع خودکواریانس مدل تحت فرض توزیع پواسن خطاها به صورت زیر بیان می‌شوند

$$E(X_t) = \frac{\mu_\varepsilon}{1 - (a + ab\mu_\varepsilon)}, \quad \gamma_X(1) = \frac{ab\mu_\varepsilon E(X_t)}{1 - (a + ab\mu_\varepsilon)},$$

$$\gamma_X(2k) = (a + ab\mu_\varepsilon)^k \gamma_X(0), \quad \gamma_X(2k+1) = (a + ab\mu_\varepsilon)^k \gamma_X(1).$$

فرض کنید $A = a + ab\mu_\varepsilon$ و $B = ab\mu_\varepsilon$. با برابر قرار دادن گشتاورهای مرتبه اول و توابع خودکواریانس فرآیند با میانگین نمونه‌ای و تابع خودکواریانس نمونه‌ای برآورد پارامترها به صورت

$$\hat{a}_{YW} = \hat{A} - \hat{B}, \quad \hat{b}_{YW} = \frac{\hat{B}}{\hat{a}_{YW}\hat{\mu}_{YW}}, \quad \hat{\mu}_{YW} = \bar{X}(1 - \hat{A}),$$

می‌باشد، که در آن $\hat{B} = \frac{\gamma_X(1)(1-\hat{A})}{\bar{X}}$ و $\hat{A} = \frac{\gamma_X(2)}{\gamma_X(0)}$.

۲.۳ روش کمترین مربعات شرطی

برآوردگرهای کمترین مربعات شرطی پارامترهای $\Theta = (a, b, \mu_\varepsilon)$ ، از طریق مینیمم‌سازی تابع هدف

$$Q(\Theta) = \sum_{t=3}^n (X_t - E(X_t|t-1))^2 = \sum_{t=3}^n (X_t - (a + ab\varepsilon_{t-1})X_{t-2} - \mu_\varepsilon)^2,$$

حاصل می‌شود، که در آن $\varepsilon_t = X_t - aX_{t-2} - abX_{t-2}\varepsilon_{t-1}$. مینیمم کردن تابع $Q(\Theta)$ با استفاده از الگوریتم نیوتن-رافسون انجام می‌شود.

۳.۳ روش ویتل

در این بخش، روشی برای برآورد پارامترها در حوزه فرکانس با استفاده از معیار ویتل^۴ معرفی می‌شود. در بسیاری از موارد، محاسبه تابع چگالی طیفی، به‌ویژه برای فرآیندهای پیچیده‌تر، ساده‌تر و سریع‌تر از محاسبه تابع احتمال دقیق است. برآوردگرهای ویتل برای پارامترهای

⁴Whittle

$\Theta = (a, b, \mu_\varepsilon)$ ، با مینیمم سازی تابع زیر به دست می آیند

$$\hat{l}_n(\Theta) = \frac{1}{n} \sum_{k=1}^{\lfloor \frac{n}{\gamma} \rfloor} (\log f_X(\omega_k) + \frac{I_n(\omega_k)}{f_X(\omega_k)}),$$

که در آن

$$f_X(\omega) = \frac{1}{\sqrt{\pi}} \frac{(1 - (a + ab\mu_\varepsilon))(\gamma_X(\circ))(1 + (a + ab\mu_\varepsilon)) + \sqrt{\gamma_X(1)} \cos \omega}{1 + (a + ab\mu_\varepsilon)^2 - \sqrt{\gamma_X(1)} \cos \omega}$$

تابع چگالی طیفی فرآیند $\text{MINBL}(\sqrt{2}, \circ, \sqrt{2}, 1)$ در فرکانس فوریه $\omega_k = \frac{\sqrt{2}\pi k}{n}$ و $I_n(\cdot)$ دوره نگار فرآیند می باشند.

۴.۳ روش درستنمایی ماکزیمم نقطه زینی

روش درستنمایی ماکزیمم به دلیل پیچیدگی های مرتبط با عملگر رقیق ساز در محاسبه تابع درستنمایی

$$L_n(\Theta) = \sum_{t=1}^n \log f_{X_t|t-1}(x_t), \quad (1.3)$$

چالش برانگیز است. در ادامه با استفاده از تقریب نقطه زینی به تقریب تابع درستنمایی می پردازیم. لگاریتم تابع مولد تجمعی شرطی X_t را به صورت $K_t(u) = \log E(e^{uX_t} | t-1)$ در نظر بگیرید. تقریب نقطه زینی برای تابع چگالی شرطی (۱.۳) به صورت زیر می باشد

$$\tilde{f}_{X_t|t}(x_t) = (\sqrt{2}\pi K_t''(\tilde{u}_t))^{-\frac{1}{\sqrt{2}}} \exp \{K_t(\tilde{u}_t) - \tilde{u}_t x_t\}, \quad (2.3)$$

به طوری که مقدار u که در معادله نقطه زینی $K_t'(u) = x_t$ صدق می کند، با \tilde{u}_t نشان داده می شود و K_t' و K_t'' به ترتیب مشتقات اول و دوم K_t نسبت به u هستند. محاسبه K_t مدل (۱.۲) در (۲.۳) به راحتی امکان پذیر نیست و با بکارگیری بسط تیلور مرتبه اول و دوم $K_t(u)$ در $u = \circ$ و اندکی محاسبات می توان نشان داد تقریب نقطه زینی تابع چگالی شرطی به صورت زیر می باشد

$$\tilde{f}_{X_t|t-1}(X_t) = (\sqrt{2}\pi\sigma_t^{\sqrt{2}}(\Theta))^{-\frac{1}{\sqrt{2}}} \exp \left\{ \frac{-(x_t - \mu_t(\Theta))^2}{\sqrt{2}\sigma_t^{\sqrt{2}}(\Theta)} \right\}.$$

و لذا برآورد درستنمایی ماکزیمم نقطه زینی پارامترها از طریق ماکزیمم کردن تابع زیر بدست می آید

$$\tilde{L}_n(\Theta) := \sum_{t=1}^n \log \tilde{f}_{X_t|t-1}(x_t) = - \sum_{t=1}^n \frac{1}{\sqrt{2}} \log(\sqrt{2}\pi\sigma_t^{\sqrt{2}}(\Theta)) - \sum_{t=1}^n \frac{(x_t - \mu_t(\Theta))^2}{\sqrt{2}\sigma_t^{\sqrt{2}}(\Theta)}.$$

که در آن μ_t و $\sigma_t^{\sqrt{2}}(\Theta)$ به ترتیب میانگین شرطی و واریانس شرطی X_t هستند.

۵.۳ شبیه سازی

در این بخش، با استفاده از شبیه سازی به بررسی عملکرد چهار روش برآوردیابی می پردازیم. هدف این مطالعه، ارزیابی عملکرد هر یک از این روش ها در برآورد پارامترهای مدل $\text{MINBL}(\sqrt{2}, \circ, \sqrt{2}, 1)$ است. مطالعات شبیه سازی برای پارامترهای زیر ارائه شده اند:

$$A1: (a, b, \mu) = (-\circ/1, \circ/9, \sqrt{2}),$$

$$A2: (a, b, \mu) = (\circ/2, -\circ/5, 1),$$

$$A3: (a, b, \mu) = (-\circ/2, -\circ/4, \sqrt{2}),$$

$$A4: (a, b, \mu) = (\circ/3, \circ/4, \sqrt{2}).$$

جدول ۱: اریبی و میانگین مربعات خطا برآوردگرهای W, YW, CLS و SPML.

n	YW			W			CLS			SPML		
	$\hat{\mu}$	$\hat{\delta}$	\hat{a}	$\hat{\mu}$	$\hat{\delta}$	\hat{a}	$\hat{\mu}$	$\hat{\delta}$	\hat{a}	$\hat{\mu}$	$\hat{\delta}$	\hat{a}
$(a, b, \mu) = (-0.1, 0.8, 7)$												
۱۰۰	۰/۲۹۱۹	۰/۷۶۰۷	۰/۳۱۱۴	۰/۳۰۱۰۰	۰/۰۰۰۷	۰/۰۰۰۵	۰/۰۰۰۶۲	۰/۰۰۰۲۴	۰/۰۰۰۳۱	۰/۰۰۰۰۱	۰/۰۰۰۱۳	۰/۰۰۰۰۵
	MSE	۰/۵۹۷۸	۰/۱۲۳۹	۰/۰۰۱۰	۰/۰۰۰۴	۰/۰۰۰۸	۰/۰۰۰۰۴	۰/۰۰۰۰۱	۰/۰۰۰۰۵	۰/۰۰۰۰۱	۰/۰۰۰۰۶	۰/۰۰۰۰۳
۳۰۰	۰/۳۰۰۹	۰/۷۲۹۷	۰/۲۶۵۱	۰/۰۰۱۰۰	۰/۰۰۰۷	۰/۰۰۰۵	۰/۰۰۰۰۷	۰/۰۰۰۰۷	۰/۰۰۰۰۸	۰/۰۰۰۰۵	۰/۰۰۰۰۹	۰/۰۰۰۰۵
	MSE	۰/۵۶۱۰	۰/۰۹۱۴	۰/۰۰۱۰	۰/۰۰۰۷	۰/۰۰۰۸	۰/۰۰۰۰۸	۰/۰۰۰۰۳	۰/۰۰۰۰۹	۰/۰۰۰۰۵	۰/۰۰۰۰۶	۰/۰۰۰۰۱
۵۰۰	۰/۳۰۳۳	۰/۶۲۴۰	۰/۲۱۲۶	۰/۰۰۰۹۹	۰/۰۰۰۷	۰/۰۰۰۵	۰/۰۰۰۰۵	۰/۰۰۰۰۳	۰/۰۰۰۰۹	۰/۰۰۰۰۵	۰/۰۰۰۰۹	۰/۰۰۰۰۱
	MSE	۰/۴۲۸۳	۰/۰۷۸۱	۰/۰۰۱۰	۰/۰۰۰۶	۰/۰۰۰۸	۰/۰۰۰۰۱	۰/۰۰۰۰۷	۰/۰۰۰۰۵	۰/۰۰۰۰۹	۰/۰۰۰۰۸	۰/۰۰۰۰۱
$(a, b, \mu) = (0.2, -0.5, 8)$												
۱۰۰	۰/۹۴۷۵	۰/۴۲۲۶	۰/۲۰۹۸	۰/۰۰۱۹۶	۰/۰۰۰۷۷	۰/۰۰۰۵	۰/۰۰۰۲۷	۰/۰۰۰۶۸	۰/۰۰۰۵۰	۰/۰۰۰۱۷۱	۰/۰۰۰۲۵۹	۰/۰۰۰۰۲
	MSE	۰/۰۰۲۱	۰/۰۰۵۷۵	۰/۰۰۰۳۹	۰/۰۰۰۰۶	۰/۰۰۰۰۶	۰/۰۰۰۰۵	۰/۰۰۰۰۷	۰/۰۰۰۰۸	۰/۰۰۰۰۲	۰/۰۰۰۰۵	۰/۰۰۰۰۷
۳۰۰	۰/۰۰۵۷۹	۰/۰۰۱۵۲۹	۰/۰۰۰۹۶۶	۰/۰۰۱۹۵	۰/۰۰۰۷۴	۰/۰۰۰۴	۰/۰۰۰۱۰	۰/۰۰۰۰۲	۰/۰۰۰۰۷	۰/۰۰۰۰۹	۰/۰۰۰۰۵	۰/۰۰۰۰۷
	MSE	۰/۰۰۲۸	۰/۰۰۳۲۴	۰/۰۰۰۳۸	۰/۰۰۰۰۶	۰/۰۰۰۰۶	۰/۰۰۰۰۵	۰/۰۰۰۰۵	۰/۰۰۰۰۹	۰/۰۰۰۰۶	۰/۰۰۰۰۵	۰/۰۰۰۰۱
۵۰۰	۰/۰۰۱۹۳	۰/۰۰۰۹۳۴	۰/۰۰۰۴۸۸	۰/۰۰۱۹۳	۰/۰۰۰۷۴	۰/۰۰۰۴	۰/۰۰۰۰۵	۰/۰۰۰۰۶	۰/۰۰۰۰۱	۰/۰۰۰۰۵	۰/۰۰۰۰۹	۰/۰۰۰۰۱
	MSE	۰/۰۰۰۷	۰/۰۰۰۳۳۶	۰/۰۰۰۳۷	۰/۰۰۰۰۵	۰/۰۰۰۰۶	۰/۰۰۰۰۸	۰/۰۰۰۰۵	۰/۰۰۰۰۷	۰/۰۰۰۰۵	۰/۰۰۰۰۷	۰/۰۰۰۰۱
$(a, b, \mu) = (0.2, -0.7, 4)$												
۱۰۰	۰/۸۱۹۰	۰/۰۰۰۲۲	۰/۰۰۰۷۳	۰/۰۰۰۸۱	۰/۰۰۰۴۰	۰/۰۰۰۷	۰/۰۰۰۱۱	۰/۰۰۰۰۶	۰/۰۰۰۰۳	۰/۰۰۰۰۳	۰/۰۰۰۰۳	۰/۰۰۰۰۳
	MSE	۰/۰۰۰۸۲۱	۰/۰۰۰۳۲۶	۰/۰۰۰۱۲	۰/۰۰۰۰۳	۰/۰۰۰۰۶	۰/۰۰۰۰۵	۰/۰۰۰۰۱	۰/۰۰۰۰۳	۰/۰۰۰۰۱	۰/۰۰۰۰۱	۰/۰۰۰۰۱
۳۰۰	۰/۰۰۰۵۹	۰/۰۰۰۲۰۰	۰/۰۰۰۵۱۹	۰/۰۰۰۰۲	۰/۰۰۰۰۲	۰/۰۰۰۰۲	۰/۰۰۰۰۲	۰/۰۰۰۰۷	۰/۰۰۰۰۵	۰/۰۰۰۰۵	۰/۰۰۰۰۲	۰/۰۰۰۰۸
	MSE	۰/۰۰۰۲۸۳	۰/۰۰۰۲۸۶	۰/۰۰۰۰۲	۰/۰۰۰۰۳	۰/۰۰۰۰۶	۰/۰۰۰۰۵	۰/۰۰۰۰۷	۰/۰۰۰۰۵	۰/۰۰۰۰۵	۰/۰۰۰۰۶	۰/۰۰۰۰۱
۵۰۰	۰/۰۰۰۸۸۲	۰/۰۰۰۳۳۹	۰/۰۰۰۶۲۰	۰/۰۰۰۰۶	۰/۰۰۰۰۳	۰/۰۰۰۰۳	۰/۰۰۰۰۲	۰/۰۰۰۰۷	۰/۰۰۰۰۵	۰/۰۰۰۰۷	۰/۰۰۰۰۵	۰/۰۰۰۰۷
	MSE	۰/۰۰۰۲۸۷	۰/۰۰۰۲۸۹	۰/۰۰۰۰۵	۰/۰۰۰۰۶	۰/۰۰۰۰۷	۰/۰۰۰۰۶	۰/۰۰۰۰۵	۰/۰۰۰۰۵	۰/۰۰۰۰۷	۰/۰۰۰۰۷	۰/۰۰۰۰۱
$(a, b, \mu) = (0.2, 0.4, 7)$												
۱۰۰	۰/۰۰۸۲۸	۰/۰۰۰۸۲۶	۰/۰۰۰۱۲۷۵	۰/۰۰۰۰۷	۰/۰۰۰۱۲	۰/۰۰۰۰۶	۰/۰۰۰۰۶	۰/۰۰۰۰۶	۰/۰۰۰۰۸	۰/۰۰۰۰۳	۰/۰۰۰۰۳	۰/۰۰۰۰۳
	MSE	۰/۰۰۰۱۲۸	۰/۰۰۰۱۸۳	۰/۰۰۰۰۲۵	۰/۰۰۰۰۵	۰/۰۰۰۰۵	۰/۰۰۰۰۶	۰/۰۰۰۰۸	۰/۰۰۰۰۵	۰/۰۰۰۰۳	۰/۰۰۰۰۳	۰/۰۰۰۰۳
۳۰۰	۰/۰۰۰۱۸۷	۰/۰۰۰۰۸۵	۰/۰۰۰۰۴۷	۰/۰۰۰۰۵	۰/۰۰۰۰۸	۰/۰۰۰۰۶	۰/۰۰۰۰۶	۰/۰۰۰۰۶	۰/۰۰۰۰۵	۰/۰۰۰۰۵	۰/۰۰۰۰۳	۰/۰۰۰۰۳
	MSE	۰/۰۰۰۰۱۷	۰/۰۰۰۰۰۸	۰/۰۰۰۰۴۷	۰/۰۰۰۰۶	۰/۰۰۰۰۵	۰/۰۰۰۰۸	۰/۰۰۰۰۴	۰/۰۰۰۰۵	۰/۰۰۰۰۵	۰/۰۰۰۰۳	۰/۰۰۰۰۳
۵۰۰	۰/۰۰۰۱۸۴	۰/۰۰۰۱۲۹	۰/۰۰۰۰۸۶۵	۰/۰۰۰۰۳	۰/۰۰۰۰۵	۰/۰۰۰۰۶	۰/۰۰۰۰۵	۰/۰۰۰۰۵	۰/۰۰۰۰۵	۰/۰۰۰۰۵	۰/۰۰۰۰۳	۰/۰۰۰۰۳
	MSE	۰/۰۰۰۰۱۴	۰/۰۰۰۰۰۷	۰/۰۰۰۰۱۳	۰/۰۰۰۰۷	۰/۰۰۰۰۵	۰/۰۰۰۰۶	۰/۰۰۰۰۴	۰/۰۰۰۰۵	۰/۰۰۰۰۵	۰/۰۰۰۰۳	۰/۰۰۰۰۳

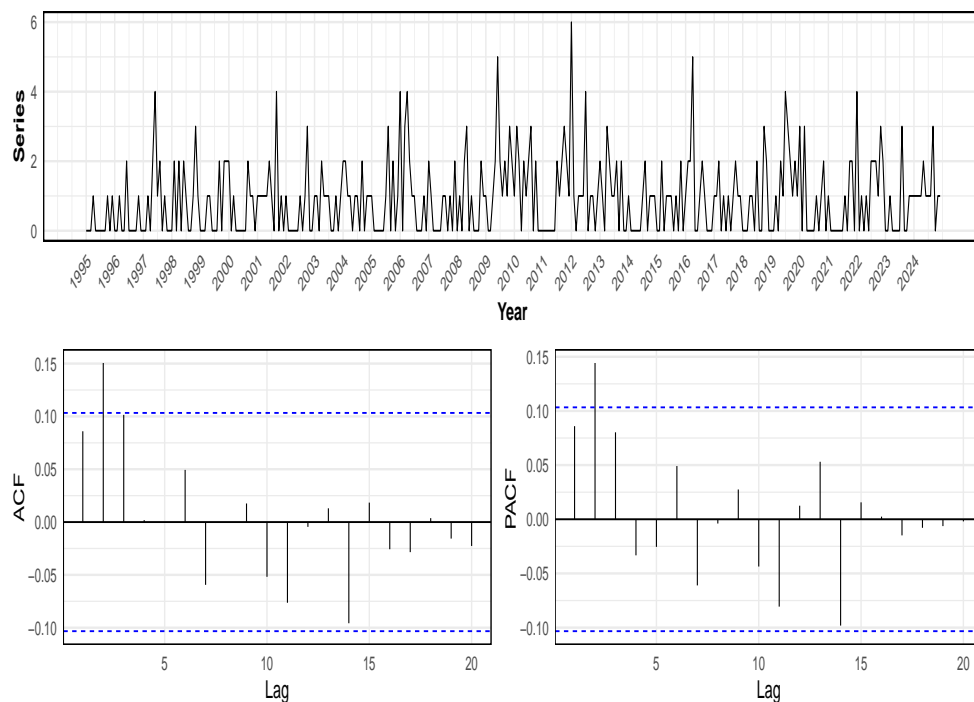
عملگرهای $a \otimes$ و $b \otimes$ عملگر رقیق‌ساز دو جمله‌ای تعمیم‌یافته علامت و سری‌های شمارشی در $a \otimes$ و $b \otimes$ دارای توزیع برنولی به ترتیب با پارامترهای $|a|$ و $|b|$ و همچنین ε_t دارای توزیع پواسن با پارامتر μ می‌باشد. جهت بررسی کارایی چهار روش برآوردیابی، شبیه‌سازی‌هایی با اندازه‌های نمونه‌ای ۱۰۰، ۳۰۰، ۵۰۰ از مدل با پارامترهای A۱ - A۴ انجام و عملکرد برآوردها با $s = 100$ تکرار بررسی شده است. نتایج شبیه‌سازی بر اساس معیارهای Bias و MSE در جدول ۱ ارائه شده است. همان‌طور که در جدول ۱ مشاهده می‌شود، هر چهار روش برآورد از نظر اریبی و میانگین مربعات خطا عملکرد رضایت‌بخشی دارند و برآوردهای ویتل بهترین عملکرد را دارند.

۴ تحلیل داده‌های واقعی

در این بخش، به کاربرد مدل $\text{MINBL}(2, 0, 2, 1)$ پرداخته می‌شود. داده‌ها مربوط به تعداد گزارش‌های ماهانه پلیس در مورد خرید و فروش اموال سرقت‌شده در گاندا^۵، استرالیا است. این اطلاعات به‌صورت ماهانه از ژانویه ۱۹۹۵ تا دسامبر ۲۰۲۴ جمع‌آوری شده است و داده‌ها از سایت bocsar^۶، گرفته شده است. در نمودار ۱ مسیر نمونه‌ای، ACF و PACF برای داده‌ها نمایش داده شده است. در ادامه،

⁵Gunnedah

⁶ <http://www.bocsar.nsw.gov.au>



شکل ۱: مسیر نمونه‌ای، تابع خودهمبستگی و تابع خودهمبستگی جزئی داده‌های مربوط به خرید و فروش اموال سرقت شده.

ایستایی داده‌ها با استفاده از آزمون دیکی-فولر افزایشی تأیید شده است که p -مقادیرهای آن برابر با ۰/۰۱ بوده است. این نتیجه بیان می‌کند که داده‌ها ایستا هستند. همچنین، آزمون غیرخطی کینان که p -مقادیرهای ۰/۰۰۰ را نشان می‌دهد، غیرخطی بودن داده‌ها را تأیید می‌کند. از سوی دیگر، داده فوق با شاخص‌های پراکندگی $\hat{I}_X = ۱/۶۵$ بیش‌پراکنده می‌باشند. مدل $\text{MINBL}(2,0,2,1)$ به‌عنوان مدلی مناسب برای برازش این مجموعه داده‌ها معرفی می‌شود، همان‌طور که در نمودار ۱ مشاهده می‌شود. در ادامه، پارامترهای مدل با در نظر گرفتن توزیع‌های مختلف از جمله پواسن، پواسن-لیندلی، هندسی و هندسی اصلاح‌شده برآورد شده‌اند. جدول ۲ شامل برآوردهای ویتل به همراه مقادیر ریشه میانگین مربعات خطا (RMSE)، میانگین قدرمطلق خطا (MAE)، میانگین مربعات خطای پیش‌بینی (FMSE) و میانگین قدرمطلق خطای پیش‌بینی (FMAE) هستند. برای محاسبه FMSE و FMAE، از پیش‌بینی درون‌نمونه‌ای ۱۰-گام به جلو استفاده شده است. نتایج جدول ۲ نشان می‌دهند که مدل PL-MINBL مقادیر کمتری از RMSE، MAE، FMSE و FMAE را به‌دست می‌آورد، که حاکی از عملکرد قوی این مدل برای داده است.

۱.۴ پیش‌بینی

در این بخش به پیش‌بینی مقادیر آینده می‌پردازیم. همان‌طور که قبلاً اشاره شد، پیش‌بینی کلاسیک k -گام به جلو از طریق امیدریاضی شرطی (۳.۲) و (۴.۲) به‌صورت زیر محاسبه می‌شود

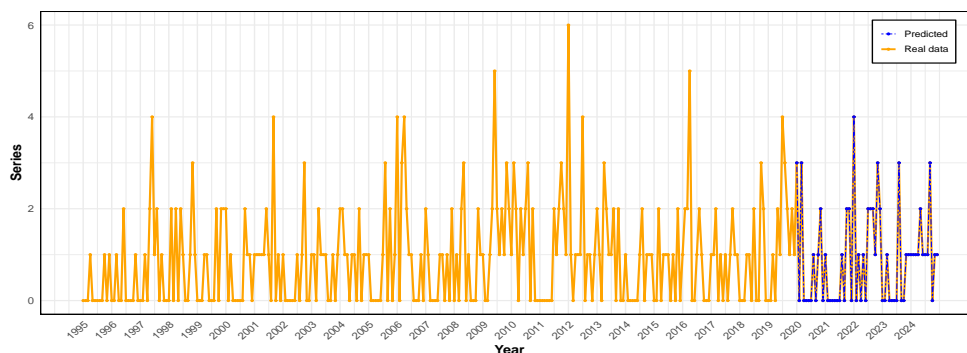
$$\hat{X}_{t+1} = (a + ab\varepsilon_t)X_{t-1} + \mu_\varepsilon, \quad \hat{X}_{t+k} = (a + ab\mu_\varepsilon)X_{t+k-2} + \mu_\varepsilon.$$

که در عمل پارامترهای a ، b و μ_ε با برآوردهای مربوطه جایگزین می‌گردد.

پیش‌بینی بوت‌استرپ

جدول ۲: برآورد پارامترها، RMSE، MAE، FMSE، و FMAE داده‌های مربوط به اموال سرقت شده.

FMAE	FMSE	MAE	RMSE	برآورد ویتل	مدل
۰.۵۲۰۲۹۰۷	۰.۷۰۰۲۷۴۷	۰.۸۰۳۹۰۱۸	۱.۰۵۰۷۷۵	$\hat{a} = ۰.۰۳۱$ $\hat{b} = ۱.۸۷۲$ $\hat{\mu} = ۰.۷۶۵$	P-MINBL(۲, ۰, ۲, ۱)
۰.۵۱۳۱۴۱۴	۰.۶۸۹۷۵۱	۰.۷۹۹۵۴۶۶	۱.۰۴۹۱۸۳	$\hat{a} = ۰.۰۶۰$ $\hat{b} = ۰.۳۷۶$ $\hat{\theta} = ۱.۷۷۶$	PL-MINBL(۲, ۰, ۲, ۱)
۰.۵۱۶۱۹۴۴	۰.۶۹۳۳۱۸۲	۰.۸۰۰۵۲۱	۱.۰۴۹۳۵۱	$\hat{a} = ۰.۰۵۱$ $\hat{b} = ۰.۶۷۲$ $\hat{\mu} = ۰.۷۶۵$	NG-MINBL(۲, ۰, ۲, ۱)
۰.۵۳۰۲۱۶۶	۰.۷۰۵۴۰۵۶	۰.۸۰۱۶۲۲۹	۱.۰۵۰۵۶۹	$\hat{a} = ۰.۰۵۲$ $\hat{b} = ۰.۶۷۲$ $\hat{p} = ۰.۵۷۳$	G-MINBL(۲, ۰, ۲, ۱)



شکل ۲: مقادیر واقعی در مقابل مقادیر پیش‌بینی‌شده به روش بوت‌استرپ مربوط به داده‌های اموال سرقت شده در پنج سال گذشته.

در این قسمت از تکنیک بوت‌استرپ اصلاح‌شده برای پیش‌بینی استفاده می‌کنیم. این روش به‌طور خاص برای مدل $\text{MINBL}(۲, ۰, ۲, ۱)$ با تنظیمات مشخص اعمال شده و در الگوریتم ۱ به تفصیل شرح داده شده است. برای ارزیابی قابلیت پیش‌بینی این روش در طول زمان، از رویکرد بوت‌استرپ برای پیش‌بینی پنج سال پایانی از داده‌های موجود استفاده شد. نتایج نشان می‌دهند که مقادیر پیش‌بینی‌شده شباهت زیادی به سری‌های اصلی دارند و این موضوع تأیید می‌کند که مدل $\text{MINBL}(۲, ۰, ۲, ۱)$ برای پیش‌بینی‌های دقیق در مجموعه داده مؤثر است.

الگوریتم ۱

گام اول: پارامترهای a, b و μ_ε را با استفاده از روش ویتل برآورد کنید.

گام دوم: باقیمانده‌های فرآیند $\hat{\varepsilon}_t = X_t - \hat{a}X_{t-1} - \hat{b}X_{t-2}\hat{\varepsilon}_{t-1}$ را برای $t = 3, \dots, n$ تعیین کنید.

گام سوم: توزیع تجربی از باقیمانده‌های تعدیل شده $\tilde{\varepsilon}_t = [\hat{\varepsilon}_t]$ زمانی که $\hat{\varepsilon}_t > 0$ ، در غیر این صورت $\tilde{\varepsilon}_t = 0$.

گام چهارم: داده‌های بوت‌استرپ X_t^b برای $b = 1, \dots, B$ به صورت زیر شبیه‌سازی شده است.

$$X_t^b = \hat{a} \oplus X_{t-1}^b + (\hat{a} \oplus X_{t-2}^b)(\hat{b} \oplus \tilde{\varepsilon}_{t-1}^b) + \tilde{\varepsilon}_t^b,$$

که در آن $\tilde{\varepsilon}_t^b$ برای $t = 1, 2, \dots, n$ در گام سوم شبیه‌سازی شده است.

گام پنجم: بر اساس داده‌های شبیه‌سازی شده X_t^b پارامترهای \hat{a}^b, \hat{b}^b و $\hat{\mu}_\varepsilon^b$ مشابه گام اول محاسبه می‌شود.

گام ششم: پارامترهای a, b و μ_ε برای B تکرار بر اساس میانگین‌های نمونه‌ای به صورت زیر برآورد شود:

$$\hat{a}^* = \frac{\sum_{b=1}^B \hat{a}^b}{B}, \quad \hat{b}^* = \frac{\sum_{b=1}^B \hat{b}^b}{B}, \quad \hat{\mu}_\varepsilon^* = \frac{\sum_{b=1}^B \hat{\mu}_\varepsilon^b}{B}.$$

گام هفتم: پیش‌بینی k -گام به جلو بر اساس معادلات بازگشتی زیر محاسبه می‌گردد

$$X_{t+k}^b = \hat{a}^* \oplus X_{t+k-1}^b + (\hat{a}^* \oplus X_{t+k-2}^b)(\hat{b}^* \oplus \tilde{\varepsilon}_{t+k-1}^b) + \tilde{\varepsilon}_{t+k}^b, \quad k \geq 1.$$

این طرح با استفاده از اعتبار ویژه پژوهشی دانشگاه مازندران اجرا گردید.

مراجع

- Bamdadi, R., Mohammadpour, M., and Ramezani, S. (2024). A bilinear modeling in counts time series with applications, *Communications in Nonlinear Science and Numerical Simulation*, **139**, 108282.
- Doukhan, P., Latour, A., and Oraichi, D. (2006). A Simple Integer-Valued Bilinear Time Series Model, *Advances in Applied Probability*, **38**, 559–578.
- Scotto, M. G., Weiß, C. H., and Gouveia, S. (2015). Thinning-based models in the analysis of integer-valued time series: a review, *Statistical Modelling*, **15**, 590–618.
- Weiß, C. H. (2008). Thinning operations for modeling time series of counts—a survey, *Advances in Statistical Analysis*, **92**, 319–341.
- Zhang, H., Wang, D., and Zhu, D. (2010). Inference for INAR(p) processes with signed generalized power series thinning operator, *Journal of Statistical Planning and Inference*, **140**, 667–683.

A class of integer-valued bilinear time series model with applications

Rana Bamdadi and Mehrnaz Mohammadpour

Department of Statistics, University of Mazandaran

Abstract: This study presents a specific class of integer-valued bilinear model of order 2 for nonlinear time series. These models possess the flexibility to incorporate the new individuals into the model in a way that allows them to behave independently. The modeling methodology applies to processes characterized by overdispersion, relying on a non-linear framework. The main characteristics of the proposed model are established. For computational purposes, some approaches for parameter estimation are carried out. To evaluate the performance of them, simulation experiments are conducted. To illustrate the modeling strategy, a social science cases is analyzed.

Keywords: Forecasting, Integer-valued bilinear process, Saddlepoint approximation, Whittle method.

Mathematics Subject Classification (2020): 62P10· 62M10· 60G25.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



مدل بندى قابليت اطمینان سیستم های تنش - مقاومت بر اساس رویکرد تابع مولد عام

بهناز تصدیقی، دکتر رضا زارعی، دکتر بهروز فتحی^۱

گروه آمار، دانشکده علوم ریاضی، دانشگاه گیلان، رشت، ایران

چکیده: مدل تداخل تنش-مقاومت پیوسته (SSI) تنش و مقاومت را به عنوان متغیرهای تصادفی پیوسته با تابع چگالی احتمال شناخته شده در نظر می‌گیرد. این امر تا حدی منجر به محدودیت در کاربرد آن می‌شود. در این مقاله، تنش و مقاومت به عنوان متغیرهای تصادفی گسسته در نظر گرفته می‌شوند و یک مدل SSI گسسته با استفاده از روش تابع مولد عام ارائه می‌شود. در نهایت، مطالعات موردی اعتبار مدل گسسته را در شرایط مختلف نشان می‌دهد، که در آن تنش و مقاومت را می‌توان با متغیرهای تصادفی پیوسته، متغیرهای تصادفی گسسته یا دو گروه از داده‌های تجربی نشان داد.

واژه‌های کلیدی: قابليت اطمینان؛ تابع مولد عام؛ مدل تداخل تنش-مقاومت؛ مدل گسسته.

کد موضوع بندی ریاضی (۲۰۲۰): 62N05, 60E10, 90B25.

۱ مقدمه

مدل تداخل تنش-مقاومت (SSI) به‌طور گسترده‌ای برای طراحی قابليت اطمینان اجزای مکانیکی استفاده شده‌است. در این مدل، اگر تنش روی یک جزء و مقاومت یک جزء به ترتیب با S_1 و S_2 نشان داده شود، قابليت اطمینان جزء که با R نشان داده شده‌است، به صورت زیر تعریف می‌شود.

$$R = \Pr(S_2 > S_1) \quad (1.1)$$

معادله (۱.۱) اساسی ترین بیان مدل تداخل تنش-مقاومت است، به این معنی که قابليت اطمینان جزء به عنوان احتمال بزرگتر بودن مقاومت از تنش در نظر گرفته می‌شود. علاوه بر این، اگر هر دو تنش-مقاومت به عنوان متغیرهای تصادفی پیوسته و توابع چگالی

^۱ سخنران، behnaztasdighi77@email.com

احتمال آنها که به ترتیب با $f_1(S_1)$ و $f_2(S_2)$ نشان داده می‌شوند، در نظر گرفته شوند، معادله (۱۰۱) را می‌توان به صورت فرمول‌های زیر بازنویسی کرد

$$R = \int_{-\infty}^{\infty} f_1(S_1) \left[\int_{S_1}^{\infty} f_2(S_2) dS_2 \right] dS_1 \quad (12.1)$$

$$R = \int_{-\infty}^{\infty} f_2(S_2) \left[\int_{-\infty}^{S_2} f_1(S_1) dS_1 \right] dS_2 \quad (20.1)$$

برای وضوح، معادله (۲۰۱) را می‌توان مدل تداخل تنش-مقاومت پیوسته نامید.

از نظر تئوری، اگر تابع چگالی احتمال تنش و مقاومت موجود باشد، قابلیت اطمینان مؤلفه را می‌توان به صورت تحلیلی یا عددی محاسبه کرد. اما در عمل، اغلب توزیع دقیق تنش و مقاومت معلوم نیست و تنها داده‌های آزمایشی محدودی درباره آن‌ها در دسترس است. بنابراین، بررسی روش‌های تقریبی برای محاسبه قابلیت اطمینان مؤلفه‌ها ضروری بوده و پژوهش‌های زیادی در این زمینه انجام شده است.

کاپور (۱) رویکردی برای تعیین مرزهای غیرقابل اطمینان دقیق ابداع کرد و این رویکرد فقط به اطلاعاتی در مورد احتمالات زیربازهای در یک منطقه تداخلی نیاز داشت. پارک و کلارک (۲) برای بهبود دقت محاسبات، فرمول کاپور را در مورد مسئله برنامه‌ریزی درجه دوم اصلاح کردند و راه حلی برای این مشکل ارائه کردند. وانگ و لیو (۳) رویکردی را برای محاسبه عدم اطمینان فازی یک جزء/سیستم ارائه کردند. علاوه بر این، کوتز و همکاران (۴) با خلاصه کردن نتایج تحقیق از رشته‌های مختلف، مدل تنش-مقاومت را تعمیم دادند و روش‌های محاسباتی را بر اساس تخمین حداکثر احتمال ارائه کردند.

در این مقاله، برخلاف مدل‌های پیوسته، تنش و مقاومت به صورت متغیر تصادفی گسسته در نظر گرفته شده و توابع جرم احتمال آن‌ها با استفاده از توابع مولد عام (UGF) نمایش داده می‌شوند. با توجه به تعریف قابلیت اطمینان مؤلفه، یک مدل تداخل تنش-مقاومت گسسته ارائه شده که برای محاسبه قابلیت اطمینان در سه حالت زیر کاربرد دارد:

۱. تنش و مقاومت متغیر تصادفی گسسته باشند،

۲. تنش و مقاومت متغیر تصادفی پیوسته باشند،

۳. توزیع تنش و مقاومت در دسترس نباشد اما توزیع فراوانی بر اساس داده‌های آزمایش موجود باشد.

در ادامه در بخش دوم روش UGF مورد بررسی قرار می‌گیرد، در بخش سوم مدل گسسته SSI معرفی می‌شود و در بخش چهارم یک مثال کاربردی ارائه می‌شود.

۲ رویکرد تابع مولد عام در متغیرهای تصادفی گسسته

تعریف ۱.۲. فرض کنید که یک متغیر تصادفی گسسته X دارای یک تابع جرم احتمال باشد که با بردار x متشکل از مقادیر ممکن X و بردار p متشکل از احتمالات مربوطه مشخص می‌شود که می‌توان با عبارات زیر فرمول‌بندی کرد:

$$x = (x_1, x_2, \dots, x_k), \quad p = (p_1, p_2, \dots, p_k), \quad p_i = \Pr(X = x_i), \quad i = 1, 2, \dots, k$$

بر اساس اصل اساسی روش UGF، تابع جرم احتمال متغیر تصادفی گسسته X را می‌توان با یک تابع چند جمله‌ای از متغیر z ، $u_X(z)$ نشان داد که مقادیر ممکن X را به احتمالات مربوطه مرتبط می‌کند.

$$u_X(z) = p_1 z^{x_1} + p_2 z^{x_2} + \dots + p_k z^{x_k} = \sum_{i=1}^k p_i z^{x_i} \quad (1.2)$$

لازم به ذکر است که برای یک متغیر تصادفی گسسته دلخواه، UGF آن منحصرأ توسط تابع جرم احتمال آن تعیین می‌شود. این بدان معنی است که یک مکاتبه یک به یک بین تابع جرم احتمال و UGF یک متغیر تصادفی گسسته وجود دارد.

بدون از دست دادن کلیت، n متغیر تصادفی گسسته مستقل X_1, X_2, \dots, X_n و یک تابع دلخواه $f(X_1, X_2, \dots, X_n)$ را می‌توان در نظر گرفت.

تعریف ۲.۲. فرض کنید تعداد مقادیر ممکن هر متغیر تصادفی به ترتیب k_1, k_2, \dots, k_n باشد. با توجه به معادله (۱.۲) UGF متغیر تصادفی تکی را می‌توان به صورت زیر به دست آورد:

$$\begin{aligned} u_{X_1}(z) &= \sum_{j_1=1}^{k_1} p_{1j_1} z^{x_{1j_1}}, \\ u_{X_2}(z) &= \sum_{j_2=1}^{k_2} p_{2j_2} z^{x_{2j_2}}, \\ &\vdots \\ u_{X_n}(z) &= \sum_{j_n=1}^{k_n} p_{nj_n} z^{x_{nj_n}}. \end{aligned}$$

تعریف ۳.۲. برای به دست آوردن UGF تابع $f(X_1, X_2, \dots, X_n)$ ، عملگر ترکیب \otimes به صورت زیر تعریف می‌شود:

$$\begin{aligned} \otimes \left(\sum_{j_1=1}^{k_1} p_{1j_1} z^{x_{1j_1}}, \sum_{j_2=1}^{k_2} p_{2j_2} z^{x_{2j_2}}, \dots, \sum_{j_n=1}^{k_n} p_{nj_n} z^{x_{nj_n}} \right) \\ = \sum_{j_1=1}^{k_1} \sum_{j_2=1}^{k_2} \dots \sum_{j_n=1}^{k_n} \left(\prod_{i=1}^n p_{ij_i} z^{f(x_{1j_1}, x_{2j_2}, \dots, x_{nj_n})} \right). \end{aligned}$$

در واقع، عملگر ترکیب \otimes یک قانون عملیات را نشان می‌دهد که به شدت به بیان تابع $f(X_1, X_2, \dots, X_n)$ بستگی دارد. بر اساس این قاعده می‌توان UGF تابع تصادفی $f(X_1, X_2, \dots, X_n)$ را به دست آورد که شکلی از تابع چند جمله‌ای نیز دارد و با تابع جرم احتمال تابع $f(X_1, X_2, \dots, X_n)$ مطابقت دارد. لازم به ذکر است که UGF ها علی‌رغم شباهت به چند جمله‌ای‌ها، چند جمله‌ای‌های منظم نیستند. با این حال، UGF ها خواص اساسی چند جمله‌ای‌های منظم را به ارث می‌برند. به عنوان مثال، در عملیات UGF ، اصطلاحات مشابهی را می‌توان جمع آوری کرد، و قانون جایگزین و قانون انجمنی قابل اجرا هستند:

$$\begin{aligned}
u_f(z) &= \bigotimes (u_{X_1}(z), \dots, u_{X_i}(z), u_{X_{i+1}}(z), \dots, u_{X_n}(z)) \\
&= \bigotimes (u_{X_1}(z), \dots, u_{X_{i+1}}(z), u_{X_i}(z), \dots, u_{X_n}(z)),
\end{aligned}$$

$$\begin{aligned}
u_f(z) &= \bigotimes (u_{X_1}(z), \dots, u_{X_i}(z), u_{X_{i+1}}(z), \dots, u_{X_n}(z)) \\
&= \bigotimes (u_{X_1}(z), \dots, u_{X_i}(z)), \bigotimes (u_{X_{i+1}}(z), \dots, u_{X_n}(z)).
\end{aligned}$$

هنگامی که UGF تابع تصادفی $f(X_1, X_2, \dots, X_n)$ به دست آمد، می‌توان آن را به عنوان یک متغیر تصادفی جدید در نظر گرفت و ویژگی‌های آماری آن را تحلیل کرد.

۳ مدل گسسته SSI

تعریف ۱.۳. فرض کنید که تنش‌های وارد بر یک جزء و مقاومت یک جزء دو متغیر تصادفی گسسته مستقل هستند که به ترتیب با S_1 و S_2 نشان داده می‌شوند. اگر تابع جرم احتمال تنش-مقاومت به صورت زیر شناخته شوند:

$$S_1 = (S_{11}, S_{12}, \dots, S_{1k_1}), P_1 = (P_{11}, P_{12}, \dots, P_{1k_1})$$

$$S_2 = (S_{21}, S_{22}, \dots, S_{2k_2}), P_2 = (P_{21}, P_{22}, \dots, P_{2k_2})$$

که در آن k_1 و k_2 تعداد مقادیر ممکن هستند که به ترتیب S_1 و S_2 می‌توانند دریافت کنند، سپس، با توجه به معادله (۱.۲) UGF تنش-مقاومت را می‌توان به صورت زیر به دست آورد:

$$u_{S_1}(z) = \sum_{j_1=1}^{k_1} p_{1j_1} z^{S_{1j_1}}, \quad u_{S_2}(z) = \sum_{j_2=1}^{k_2} p_{2j_2} z^{S_{2j_2}}$$

یک تابع $f(S_1, S_2)$ از متغیر تصادفی تنش-مقاومت به صورت زیر ساخته می‌شود

$$f(S_1, S_2) = S_2 - S_1 \quad (۱.۳)$$

بر اساس روش UGF معرفی شده در بخش ۲، UGF تابع گسسته $f(S_1, S_2)$ را می‌توان به صورت زیر به دست آورد

$$\begin{aligned}
u_f(z) &= \bigotimes (u_{S_1}(z), u_{S_2}(z)) = \bigotimes \left(\sum_{j_1=1}^{k_1} p_{1j_1} z^{S_{1j_1}}, \sum_{j_2=1}^{k_2} p_{2j_2} z^{S_{2j_2}} \right) \\
&= \sum_{j_1=1}^{k_1} \sum_{j_2=1}^{k_2} \left(\prod_{i=1}^2 p_{ij_i} z^{f(S_{1j_1}, S_{2j_2})} \right) = \sum_{j_1=1}^{k_1} \sum_{j_2=1}^{k_2} \left(\prod_{i=1}^2 p_{ij_i} z^{(S_{2j_2} - S_{1j_1})} \right) \quad (۲.۳)
\end{aligned}$$

به عنوان یک حالت قطعی، عملگر \otimes در معادله (۲.۳) به عنوان تفریق توان‌های متناظر در چندجمله‌ای‌های ضرب تعریف می‌شود. درک اینکه شکل نهایی UGF تابع $f(S_1, S_2)$ یک تابع چند جمله‌ای است که شامل $K \leq k_1 \times k_2$ عبارت است (تعداد کل عبارت‌ها می‌تواند کمتر از $k_1 \times k_2$ باشد پس از جمع‌آوری عبارت‌های مشابه) دشوار نیست. بنابراین، معادله (۲.۳) را می‌توان به صورت زیر بازنویسی کرد

$$u_f(z) = \sum_{j=1}^k P_j z^{f_j} \quad (3.3)$$

که در آن f_j و P_j ($j = 1, 2, \dots, k$) به ترتیب مقادیر ممکن تابع $f(S_1, S_2)$ و احتمالات مربوطه هستند. همانطور که با معادله (۱.۱) بیان می‌شود، قابلیت اطمینان جزء به عنوان احتمال بزرگتر بودن مقاومت از تنش تعریف می‌شود. با تبدیل معادله (۱.۱)، می‌توان به دست آورد

$$R = \Pr(S_2 - S_1 > 0) \quad (4.3)$$

جایگزینی معادله (۱.۳) به معادله (۴.۳) نتیجه می‌دهد

$$R = \Pr(f(S_1, S_2) > 0) \quad (5.3)$$

که در آن $f(S_1, S_2)$ اساساً، یک متغیر تصادفی گسسته جدید است و خواص توزیع آن را می‌توان با UGF آن نشان داد. برای محاسبه احتمال بیان شده توسط معادله (۵.۳) بر اساس UGF تابع $f(S_1, S_2)$ می‌توان یک تابع با ارزش دودویی با دامنه روی مجموعه مقادیر ممکن تابع $f(S_1, S_2)$ تعریف کرد.

$$\alpha(f_i) = \begin{cases} 1, & f_i > 0, \\ 0, & f_i \leq 0. \end{cases}$$

سپس، بر اساس معادله (۵.۳)، قابلیت اطمینان جزء را می‌توان به صورت زیر محاسبه کرد

$$R = \Pr(f(S_1, S_2) > 0) = \sum_{j=1}^K P_j \alpha(f_j) \quad (6.3)$$

برای مقایسه، معادله (۶.۳) را می‌توان مدل تداخل تنش-مقاومت گسسته نامید.

۴ مثال کاربردی

فرض کنید مدل تداخل تنش-مقاومت گسسته ارائه شده در بخش ۳ می‌تواند برای ارزیابی قابلیت اطمینان یک جزء در چندین مورد استفاده شود. بدیهی است که وقتی تنش-مقاومت متغیر تصادفی گسسته هستند، قابلیت اطمینان جزء را می‌توان مستقیماً مطابق با معادلات (۲.۳)، (۳.۳) و (۶.۳) به دست آورد. در این بخش دو مورد در نظر گرفته شده است. مورد ۱ نشان می‌دهد که تنش-مقاومت متغیر تصادفی پیوسته با تابع چگالی احتمال شناخته شده هستند و مورد ۲ نشان می‌دهد که تنها توزیع فراوانی تنش-مقاومت بر اساس تجزیه و تحلیل آماری داده‌ها در دسترس است.

مورد ۱. مثال (۳) را در نظر بگیرید: تنش بر روی یک جزء، S_1 ، به صورت نمایی با میانگین $\mu_1 = 50 MPa$ توزیع می‌شود، و مقاومت مولفه، S_2 ، به طور s -نرمال با میانگین $\mu_2 = 100 MPa$ و انحراف استاندارد $\sigma_2 = 10 MPa$ توزیع می‌شود. مقدار دقیق قابلیت اطمینان این جزء برابر با ۸۶۱۹۴٪ است.

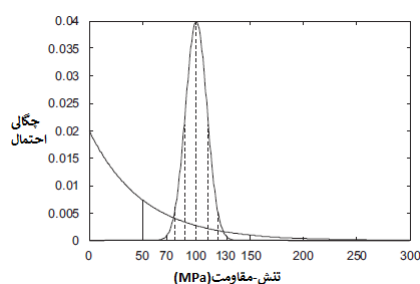
در این روش، برای محاسبه قابلیت اطمینان مؤلفه، متغیر تصادفی پیوسته با تابع چگالی احتمال شناخته‌شده به یک متغیر تصادفی گسسته با تابع جرم احتمال تبدیل می‌شود. در مرحله اول، با توجه به شرایط عملیاتی قطعه، محدوده تقریبی مقادیر ممکن برای تنش و مقاومت تعیین می‌شود که به ترتیب با فواصل $\langle S_{1min}, S_{1max} \rangle$ و $\langle S_{2min}, S_{2max} \rangle$ مشخص می‌شوند. سپس فواصل $\langle S_{1min}, S_{1max} \rangle$ و $\langle S_{2min}, S_{2max} \rangle$ به ترتیب به زیر بازه‌های m و n تقسیم می‌شوند ($m = n$ مجاز است). مقادیر نقطه میانی هر زیر بازه به عنوان مقادیر ممکن متغیر تصادفی و مساحت هر زیر بازه به عنوان احتمال آن در نظر گرفته می‌شوند. با این روش، دو متغیر تصادفی مجزا برای تنش و مقاومت با تابع جرم احتمال مشخص به دست می‌آید. سپس متغیر تصادفی گسسته تنش-مقاومت با استفاده از روش UGF نمایش داده شده و قابلیت اطمینان جزء با مدل تداخل تنش-مقاومت گسسته ارائه شده در بخش ۳ محاسبه می‌شود.

بگذارید دامنه تنش-مقاومت به ترتیب $\langle 0, 300 \rangle MPa$ و $\langle 0, 6 \rangle MPa$ باشد و هر دو بازه را به شش زیر بازه تقسیم کنید. نتیجه تقسیم بندی بازه‌ای در شکل (۱) نشان داده شده است. در این مورد، مقادیر مساحت زیربازه‌ها با استفاده از $MATLAB$ ۱۰.۷ محاسبه می‌شود. بر اساس مقادیر نقطه میانی و مقادیر مساحت همه زیربازه‌ها، می‌توان متغیر تصادفی گسسته تنش-مقاومت را به ترتیب به صورت زیر به دست آورد:

$$S_1 = (25, 75, 125, 175, 225, 275), S_2 = (75, 85, 95, 105, 115, 125)$$

$$p_1 = (0.6321, 0.2325, 0.0855, 0.0315, 0.0116, 0.0043)$$

$$p_2 = (0.0214, 0.1359, 0.3413, 0.1359, 0.0214)$$



شکل ۱: تابع چگالی احتمال منحنی‌ها و تقسیم بندی بازه‌ای

با توجه به روش ارائه شده در بخش ۳، قابلیت اطمینان این مولفه به صورت $R = 8625\%$ به دست می‌آید. خطای نسبی در مقایسه با مقدار دقیق قابلیت اطمینان ۶۵٪ است. به منظور نشان دادن تأثیر تقسیم بندی بازه‌ای بر دقت محاسبه، می‌توان بازه تنش را به ۱۲ زیربازه تقسیم کرد در حالی که تعداد زیربازه‌های مقاومت را بدون تغییر نگه داشت. بنابراین، می‌توان توصیف دیگری از تنش به دست آورد:

$$S_1 = (12,537,5, 62,5, 87,5, 112,5, 137,5, 162,5, 187,5, 212,5, 237,5, 262,5, 287,5)$$

$$P_1 = (0,3953, 0,2387, 0,1447, 0,0878, 0,0533, 0,0323, 0,0196, 0,0119, 0,0072,$$

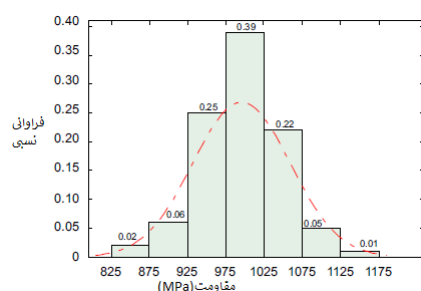
$$0,0044, 0,0027, 0,0016)$$

به طور مشابه، پایایی این مؤلفه را می توان به صورت $R = 0,86163$ محاسبه کرد و خطای نسبی برابر با $0,36\%$ است. این نتیجه نشان می دهد که کاهش طول زیربازه می تواند دقت محاسباتی را زمانی که دامنه تنش-مقاومت ثابت است، بهبود بخشد.

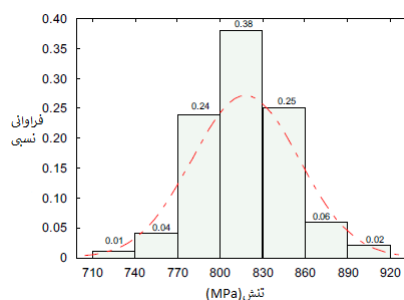
مورد ۲. برای محاسبه قابلیت اطمینان یک جزء، اگر مدل تداخل تنش-مقاومت پیوسته استفاده شود، تنش و مقاومت به صورت متغیر تصادفی پیوسته در نظر گرفته شده و توزیع آن ها با روش های برازش توزیع و تخمین پارامتر به دست می آید. اما در مدل گسسته، نیازی به توجه به توزیع واقعی تنش-مقاومت نیست و می توان با استفاده از دو متغیر تصادفی جداگانه اطلاعات آماری را مستقیماً برای محاسبه قابلیت اطمینان به کار برد. فرض کنید تعداد داده های موجود در هر گروه برابر با 100 است و پس از یک پردازش ساده می توان داده ها را با هیستوگرام توصیف کرد. همانطور که در شکل های ۲(آ) و ۲(ب) نشان داده شده است، فواصل کلاسی داده ها و فرکانس های نسبی متناظر آن ها به دست آمده است. به طور مشابه، مقادیر نقطه میانی هر بازه کلاس به عنوان مقادیر ممکن متغیر تصادفی تنش-مقاومت در نظر گرفته می شود و فرکانس های نسبی هر بازه کلاس به عنوان احتمالات مربوطه در نظر گرفته می شوند. بنابراین، دو متغیر تصادفی جدید تنش-مقاومت با تابع جرم احتمال شناخته شده به شرح زیر به دست می آید:

$$S_1 = (725, 755, 785, 815, 845, 875, 905), S_2 = (850, 900, 950, 1000, 1050, 1100, 1150)$$

$$p_1 = (0,01, 0,04, 0,24, 0,38, 0,25, 0,06, 0,02), p_2 = (0,02, 0,06, 0,25, 0,39, 0,22, 0,05, 0,01)$$



(ب) هیستوگرام داده های مقاومت



(آ) هیستوگرام داده های تنش

با استفاده از مدل تداخل تنش-مقاومت گسسته ارائه شده در بخش ۳، قابلیت اطمینان این جزء به صورت $R = 0,9972$ به دست می آید.

همانطور که در بخش ۱ ذکر شد، روش ارائه شده توسط (۴) می تواند برای حل این مشکل استفاده شود که نوع توزیع احتمال تنش-مقاومت مشخص باشد. فرض کنید هر دو تنش-مقاومت در این حالت متغیر تصادفی نرمال مستقل با پارامترهای ناشناخته هستند. دو گروه از داده ها را می توان به عنوان مشاهدات تنش-مقاومت در نظر گرفت. با استفاده از روش برآورد حداکثر درست نمایی،

پایایی این مؤلفه را می‌توان $R = 0.9962$ محاسبه کرد. اگر این نتیجه به عنوان مقدار دقیق قابلیت اطمینان در نظر گرفته شود، خطای نسبی حاصل از مدل تداخل تنش-مقاومت گسسته برابر با ۰/۱ است.

بحث و نتیجه‌گیری

در این مقاله، یک رویکرد جدید بر اساس روش تابع مولد عام برای مدل‌سازی قابلیت اطمینان سیستم‌های تنش-مقاومت بر اساس متغیرهای تصادفی گسسته مورد بررسی و مطالعه قرار گرفت که در آن یک مدل تداخل تنش-مقاومت گسسته برای محاسبه قابلیت اطمینان مؤلفه را تحت شرایط مختلف، از جمله زمانی که داده‌ها به صورت پیوسته، گسسته، یا تنها در قالب دو گروه از داده‌های تجربی موجود باشند، داراست. یکی از مهمترین مزایای مدل پیشنهادی، عدم نیاز به اطلاعات قبلی در مورد توزیع دقیق تنش و مقاومت است که به طور قابل توجهی محدودیت‌های موجود در مدل‌های پیوسته SSI را برطرف می‌کند. در بسیاری از کاربردهای عملی، دسترسی به توزیع دقیق داده‌ها دشوار است و تنها داده‌های آزمایشی محدودی در دسترس قرار می‌گیرد. نتایج عددی نشان دهنده برتری روش مورد مطالعه است. برای داده‌های پیوسته که به متغیرهای گسسته تبدیل شدند، مدل دقت بالایی بوده و نتایج نشان داد که کاهش طول زیربازه در فرآیند گسسته‌سازی، دقت محاسبات را به طور محسوسی افزایش می‌دهد. علاوه بر این، مدل SSI گسسته کارایی مناسبی را در شرایطی که تنها دو گروه از داده‌های تنش و مقاومت موجود بود، ارائه نمود.

مراجع

- [1] Kapur, K.C., 1975, *Reliability bounds in probabilistic design*, IEEE Transactions on Reliability, 24(3), pp.193-195.
- [2] Park, J.W. and Clark, G.M., 1986, *A computational algorithm for reliability bounds in probabilistic design*, IEEE transactions on reliability, 35(1), pp.30-31.
- [3] Wang, J.D. and Liu, T.S., 2002, *Fuzzy reliability using a discrete stress-strength interference model*, IEEE transactions on reliability, 45(1), pp.145-149.
- [4] Kotz, S., Lumelskii, Y. and Pensky, M., 2003, *The stress-strength model and its generalizations: theory and applications*.

A discrete stress–strength interference model based on universal generating function

Behnaz Tsadighi, Dr.Reza Zarei, Dr.Behroz Fathi ¹

¹Department of Statistics, Faculty of Mathematical Sciences, University of Guilan, Rasht, Iran

Abstract: Continuous stress–strength interference (SSI) model regards stress and strength as continuous random variables with known probability density function. This, to some extent, results in a limitation of its application. In this paper, stress and strength are treated as discrete random variables, and a discrete SSI model is presented by using the universal generating function (UGF) method. Finally, case studies demonstrate the validity of the discrete model in a variety of circumstances, in which stress and strength can be represented by continuous random variables, discrete random variables, or two groups of experimental data.

Keywords: Reliability; Universal generating function; Stress–strength interference ;Discrete model.

Mathematics Subject Classification (2020): 90B25, 60E10, 62N05.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



اخلاق در آمار و احتمال

علی دولتی^۱

گروه آمار دانشگاه یزد

چکیده: امروزه با گسترش استفاده از داده‌ها و روش‌های آماری در تصمیم‌سازی‌ها، اقدامهای غیر اخلاقی مانند سوءاستفاده از داده‌ها، ارائه تفسیرهای مغرضانه و نادرست از نتیجه تحلیل‌ها و عدم شفافیت در ارائه گزارش‌ها، به صورت آگاهانه یا ناآگاهانه در انجام کارهای آماری وجود دارد. از این رو، برای تقویت اعتماد عمومی به آمار و حفظ سلامت و اعتبار این حوزه، آموزش و ترویج اصول اخلاقی، امری حیاتی است. هدف این مقاله ارائه الگویی برای گنجاندن اصول اخلاقی در سرفصل‌های درسی آمار است تا دانشجویان علاوه بر یادگیری مفاهیم نظری و کاربردی آمار، به مسئولیت‌پذیری اجتماعی و حرفه‌ای خود نیز آگاهی پیدا کنند.

واژه‌های کلیدی: آموزش آمار، اخلاق در آمار، برنامه‌درسی، مسئولیت‌پذیری حرفه‌ای، حریم خصوصی داده‌ها.

کد موضوع‌بندی ریاضی (۲۰۲۰): ۹۷ک۸۰، ۹۷ک۹۹.

۱ مقدمه

اهمیت ملاحظات اخلاقی در یک رشته علمی یا حوزه تخصصی، زمانی برجسته می‌شود که نتایج آن تأثیرات گسترده و قابل توجهی بر جهان پیرامون خود داشته باشد. در رشته‌هایی نظیر پزشکی، مهندسی، روانشناسی، علوم انسانی و اجتماعی، عموماً این باور وجود دارد که کار و تخصص آنها می‌تواند پیامدهای اخلاقی گسترده و فراگیری برای جامعه داشته باشد و این پیامدها بایستی مورد توجه و ارزیابی قرار گیرند. در این حوزه‌ها سازمان‌ها و نهادهای تخصصی گوناگونی مانند نظام مهندسی، نظام پزشکی و انجمن‌های مختلفی وجود دارند که با تدوین شیوه‌نامه‌ها و دستورالعمل‌هایی به‌طور جدی بر اهمیت مسئولیت‌های اجتماعی و ملاحظات اخلاقی متخصصان در زمینه‌های کاری خود تأکید می‌ورزند و آن را جزو لاینفک استانداردهای حرفه‌ای خود می‌دانند. این موضوع نشان دهنده این است که در این رشته‌ها، توجه به پیامدهای اخلاقی کارها، بخش مهمی از هویت حرفه‌ای محسوب می‌شود. در ایران نیز مانند سایر کشورهای دنیا در

^۱ سخنران، adolati@yazd.ac.ir

زمینه اخلاق علمی و حرفه‌ای، قوانین، دستورالعمل‌ها، شیوه‌نامه‌ها، نظام‌نامه‌ها و راهنماهای مختلفی در حوزه‌هایی مانند وکالت، مهندسی، روان‌شناسی و مشاوره، حسابداری و پزشکی، وجود دارد (۷؛ ۸؛ ۹؛ ۱۰؛ ۱۱؛ ۱۲) در مورد ضرورت آموزش اخلاق در رشته‌های دانشگاهی مقالات زیادی نوشته شده است. به‌عنوان نمونه، مراجع (۱؛ ۲؛ ۳؛ ۴؛ ۵؛ ۶) را ملاحظه کنید. در بسیاری از کشورها، وکلا، پزشکان، زیست‌شناسان، مهندسان، فیزیک‌دانان، متخصصان علوم کامپیوتر، آماردانان و حسابداران خبره اخلاق تخصصی مرتبط با رشته خود را فرامی‌گیرند، زیرا در مقام متخصص، با تصمیم‌گیری‌های اخلاقی در حرفه خود روبرو خواهند شد. فیسler و همکاران (۹) در سال ۲۰۲۰ به تحلیل محتوای بیش از صد عنوان درسی با محتوای آموزش اخلاق در رشته‌های مختلف پرداخته‌اند که در برخی از دانشگاه‌های دنیا تدریس می‌شود. این بررسی، طیف وسیعی از رویکردهای موجود در زمینه آموزش اخلاق در رشته‌های مختلف را در طول دهه‌ها نشان می‌دهد. در ایران نیز در برنامه‌های درسی برخی از رشته‌های مهندسی و علوم پزشکی، دروسی با محتوای اخلاق علمی و حرفه‌ای وجود دارد. رجبعلی پور و همکاران (۵) به تحلیل محتوای دروس ۲۰۵ رشته در حوزه پزشکی و پیراپزشکی پرداخته و دروسی را که محتوای اخلاق حرفه‌ای داشته‌اند، شناسایی کرده‌اند. در حوزه مهندسی نیز دروسی با محتوای اخلاق مهندسی در برنامه درسی برخی از رشته‌های دانشگاهی کشور وجود دارد (۴).

در عصر انفجار داده‌ها و تحولات سریع فناوری در حوزه تحلیل داده‌ها، آمار به عنوان ابزاری حیاتی برای تصمیم‌گیری‌های حساس در حوزه‌های مختلف علمی، پزشکی، اقتصادی و اجتماعی مطرح است. با این حال، قدرت فزاینده روش‌های آماری همراه با مسئولیت اخلاقی روزافزونی است که متخصصان این حوزه باید به آن توجه کنند. اخلاقیات آماری به اصول و دستورالعمل‌هایی اشاره دارد که مسئولیت‌پذیری و رفتار اخلاقی در فعالیت‌های آماری را تنظیم می‌کنند. این اصول تمام مراحل کار آماری، از جمع‌آوری و تحلیل داده‌ها تا تفسیر و گزارش‌دهی را در بر می‌گیرد. رعایت اخلاق در آمار برای حفظ صداقت پژوهش، تقویت اعتماد عمومی و اطمینان از استفاده مسئولانه از داده‌ها در تصمیم‌گیری‌های حیاتی است. موارد متعددی مانند سوءاستفاده از تکنیک‌های آماری، نداشتن تخصص و صلاحیت انجام کار آماری، تحریف عمدی نتایج، نقض حریم خصوصی داده‌ها و تفسیرهای گمراه‌کننده از نتایج، توسط انجام دهندگان کارهای آماری، لزوم توجه جدی‌تر به آموزش اصول اخلاقی در کنار مباحث فنی آمار را آشکار ساخته است. آموزش رسمی دانشجویان رشته آمار عمدتاً بر جنبه‌های محاسباتی و تکنیکی متمرکز بوده و کمتر به ابعاد اخلاقی این حرفه پرداخته است. این در حالی است که یک تصمیم آماری نادرست یا غیراخلاقی می‌تواند پیامدهای گسترده‌ای در سطح جامعه داشته باشد. از خطاهای پزشکی ناشی از تحلیل‌های نادرست گرفته تا سیاست‌گذاری‌های اقتصادی اشتباه بر اساس داده‌های تحریف شده، همگی بر اهمیت رعایت اصول اخلاقی در این حوزه تأکید دارند. نبود آموزش‌های مرتبط با این حوزه، به‌ویژه زمانی که با سطح آموزش‌های اخلاقی در سایر رشته‌ها مانند پزشکی و مهندسی مقایسه می‌شود، نمود بیشتری پیدا می‌کند. این مقاله با توجه به این ضرورت، به بررسی راهکارهایی برای گنجاندن آموزش اخلاق در برنامه‌های درسی آمار می‌پردازد.

۲ تاریخچه شیوه‌نامه‌های اخلاقی آمار

توجه به ملاحظات اخلاقی در تحقیقات سابقه‌ای طولانی دارد. پس از جنگ جهانی دوم و فجایعی مانند آزمایش‌های نازی‌ها بر روی انسان‌ها و برگزاری دادگاه‌های نورنبرگ مجموعه‌ای از اصول اخلاقی که به کد نورنبرگ (۱۷) معروف است تدوین شد و به‌عنوان یکی از پایه‌های مهم اخلاق پزشکی و تحقیقات انسانی شناخته می‌شود. مهم‌ترین اصول کد نورنبرگ عبارت‌اند از: رضایت داوطلبانه و مشارکت فرد با آگاهی کامل از خطرات و اهداف تحقیق و ممنوعیت اجبار، فریب یا فشار غیراخلاقی، طراحی آزمایش‌ها بر پایه دانش علمی و

با اهداف مفید برای جامعه، نرساندن رنج یا آسیب جسمی و روانی به افراد، حق انصراف شرکت‌کننده در هر مرحله از تحقیق و نهایتاً برتری قابل‌توجه منافع تحقیق بر خطرات آن. این اصول بعدها در بیانیه هلسینکی (۲۱) و سایر قوانین اخلاق پزشکی جهانی گنجانده شد. آمار به‌عنوان یکی از ابزارهای علمی برای فهم واقعیت‌های اجتماعی، اقتصادی و پزشکی، همواره با داده‌هایی سروکار دارد که گاه به طور مستقیم با زندگی انسان‌ها گره‌خورده‌اند. به همین دلیل، توجه به جنبه‌های اخلاقی در انجام مطالعات آماری، امری حیاتی و اجتناب‌ناپذیر است. سوءاستفاده‌های آماری در پژوهش‌ها (مانند دست‌کاری داده‌ها یا گزارش انتخابی نتایج) منجر به توجه بیشتر به اخلاق در انتشار یافته‌های علمی بر پایه روشهای آمار شد. اخلاقیات آماری، به عنوان مجموعه اصول راهنمای انجام مسئولانه کارهای آماری، در طول زمان تکامل یافته است. از دهه ۸۰ میلادی، انجمن‌های آماری بین‌المللی دستورالعمل‌هایی برای رفتار حرفه‌ای تدوین کردند. انجمن آمار آمریکا^۱ نقش کلیدی در این تحول ایفا کرده و طی چند دهه، «راهنمای اخلاقی برای انجام کارهای آماری» (۱۴) خود را تدوین و اصلاح کرده است. تاریخچه بیش از سه دهه فعالیت انجمن آمار آمریکا برای تدوین این شیوه‌نامه را در مراجع (۱۹؛ ۲۲) ببینید. با گسترش فناوریهای نوظهور و جمع‌آوری داده‌های بزرگ و افزایش چالش‌ها و نگرانی‌ها درباره حریم خصوصی و سوءاستفاده از داده‌ها، استانداردهای اخلاقی جامع‌تری در انجام کارهای آماری لازم گردید و سازمانهای دیگری مانند موسسه بین‌المللی آمار^۲ فدراسیون انجمنهای ملی آمار اروپا^۳، انجمن آکچواری (بیمسنجی) آمریکا^۴، انجمن آمار سلطنتی^۵، اداره آمار بریتانیا^۶، برای ترویج اخلاق در انجام کارهای آماری دستورالعمل‌های خود را تدوین و منتشر نمودند. کمیسیون آمار سازمان ملل متحد در سال ۱۹۹۴ اصول اساسی آمار رسمی را تصویب کرد. نسخه تجدیدنظر شده این اصول در سال ۲۰۱۳ منتشر شد^۷. اخیراً انجمن آمار ایران نیز با تشکیل «کمیته اخلاق آماری» و تصویب «شیوه‌نامه اخلاقی برای کارهای آماری» (۱۲) در این زمینه فعالیت‌های خود را آغاز نموده است.

۳ اصول اخلاق آماری

در این بخش به اصول اخلاقی انجام کارهای آماری پرداخته می‌شود. نخست برخی تعریف‌های لازم یادآوری می‌شود. در شیوه‌نامه اخلاقی برای کارهای آماری که توسط انجمن آمار ایران (۱۲) و با بهره‌گیری از شیوه‌نامه‌های معتبر منتشر شده در این زمینه توسط انجمن آمار آمریکا و موسسه بین‌المللی آمار تدوین شده است، «کار آماری»، شامل فعالیت‌هایی مانند نمونه‌گیری، طراحی آزمایش‌ها، خلاصه‌سازی، پردازش، تصویری‌سازی، تحلیل و تفسیر داده‌ها، ایجاد یا بکارگیری مدل یا توسعه و تعمیم الگوریتم‌ها است. «انجام دهندگان کارهای آماری» همه کسانی هستند که بدون توجه به عنوان شغل، حرفه، میزان اطلاعات یا رشته تحصیلی، به این فعالیت‌ها می‌پردازند. کسانی که در یک کار آماری سرمایه‌گذاری می‌کنند، به آن کمک می‌کنند، از آن استفاده می‌کنند، یا به گونه‌ای تحت تأثیر آن قرار می‌گیرند، «ذینفعان» محسوب می‌شوند. سازمان‌ها، مؤسسات، شرکت‌ها، دانشگاه‌ها، نهادهای دولتی یا غیردولتی که به نحوی در انجام، حمایت، تأمین مالی، اجرای پروژه‌های آماری یا بهره‌برداری از نتایج آن نقش دارند، به عنوان «کارفرما، ناظر، حامی مالی، یا استفاده‌کننده از نتایج آماری» شناخته می‌شوند. به فردی که وظیفه نظارت، هدایت و تضمین کیفیت در اجرای کار آماری را بر عهده دارد و ممکن است به عنوان استاد راهنما، مدیر پروژه، سرپرست تحلیل آماری، پژوهشگر ارشد، مربی یا هر عنوان دیگری فعالیت کند و مسئولیت اجرای صحیح، شفافیت

¹American Statistical Association (ASA)

²International Statistical Institute (ISI)

³The Federation of European Statistical National Societies (FENStatS)

⁴American Academy of Actuaries

⁵Royal Statistical Society (RSS)

⁶UK Statistics Authority

⁷United Nations Statistical Commission

و رعایت اصول اخلاقی در انجام کار آماری را بر عهده دارد یا در آن شریک است، «راهنمای علمی و فنی در کار آماری» گفته می‌شود. اصول اخلاق در انجام کارهای آماری شامل موارد زیر است (۱۲) :

(اصل الف) صداقت حرفه‌ای و مسئولیت‌پذیری : انجام‌دهنده کارهای آماری اخلاق‌مدار، برای ایجاد اعتماد بین خود و ذینفعان، توانمندی‌ها و فعالیت‌های خود را به‌طور صادقانه عرضه می‌کند، با دیگران با احترام برخورد می‌نماید و مسئولیت کارهای خود را می‌پذیرد.

(اصل ب) درستی و یکپارچگی داده‌ها و روش‌ها:

انجام دهنده کارهای آماری اخلاق‌مدار در حالی که در صدد کشف و رفع محدودیت‌ها، نقص‌ها یا سوگیری‌های شناخته شده یا مشکوک در داده‌ها یا روش‌ها است، پیامدهای بالقوه آنها را بر تفسیر، نتیجه‌گیری، توصیه‌ها، تصمیم‌گیری‌ها یا سایر نتایج کارهای آماری اطلاع‌رسانی می‌کند.

(اصل پ) مسئولیت‌های حرفه‌ای: انجام‌دهندگان کارهای آماری اخلاق‌مدار به اصول علمی و فنی تکیه می‌کنند و در قبال پیامدهای اجتماعی، اخلاقی و انسانی تصمیمات آماری خود نیز مسئولیت‌پذیر هستند. این مسئولیت‌ها به‌طور مستقیم به نحوه تعامل با ذینفعان، داده‌ها، اعضای تیم‌های چندرشته‌ای و همکاران مربوط می‌شود. با رعایت دقیق و شفاف این مسئولیت‌ها، می‌توان اطمینان حاصل کرد که فرایندهای آماری در راستای منافع عمومی، احترام به حقوق افراد و ارتقای کیفیت تصمیم‌سازیهای مبتنی بر داده‌ها پیش می‌روند.

(اصل پ ۱) مسئولیت در قبال ذینفعان: انجام‌دهنده کارهای آماری اخلاق‌مدار، انتظارات ذینفعان از هر پروژه خاص را شناسایی و تبیین می‌کند و به منافع و حقوق آنها احترام می‌گذارد.

(اصل پ ۲) مسئولیت در قبال داده‌ها و آزمودنیهای تحت تأثیر کارهای آماری: انجام دهنده کارهای آماری اخلاق‌مدار برای تولید و استفاده از داده‌ها، به گونه‌ای عمل می‌کند که موجب آسیب به افراد یا گروه‌ها نشود و همواره به حقوق انسان‌ها، حیوانات و محیط زیست احترام می‌گذارد. در عین حال، امکان سوءاستفاده‌های دیگران از داده‌ها را نیز نادیده نمی‌گیرد. مسئولیت‌پذیر است و به عواقب کارهای آماری خود، به ویژه کارهایی که ممکن است تأثیر مستقیم بر زندگی دیگران و محیط زیست داشته باشد، توجه دارد.

(اصل پ ۳) مسئولیت در قبال همکاران و اعضای تیم‌های چندرشته‌ای: کارهای آماری معمولاً در حوزه‌های مختلف و تیم‌های چندرشته‌ای انجام می‌شود که افراد با تخصص‌ها و دیدگاه‌های متفاوت در آن مشارکت دارند. این تنوع، فارغ از عنوان شغلی یا میزان تحصیلات، می‌تواند منجر به هم‌افزایی و بهبود کیفیت کارها شود، اما در عین حال نیازمند رعایت اصول اخلاقی در تعاملات و همکاری‌ها است. انجام‌دهندگان کارهای آماری باید در محیط‌های گروهی و تیمی به‌طور اخلاقی و حرفه‌ای عمل کنند و همواره اصول احترام متقابل و شفافیت را رعایت نمایند.

(اصل ت) الزامات و تعهدات اشخاص حقوقی، استادان راهنما، راهبران، سرپرستان و مربیان: در فرآیند انجام کارهای آماری، علاوه بر فرد انجام‌دهنده، نقش استادان راهنما، راهبران، سرپرستان، مشاوران و سازمانها نیز از اهمیت ویژه‌ای برخوردار است. این افراد و سازمان‌ها باید با رعایت اصول اخلاقی و حرفه‌ای، ترویج‌دهندگان فرهنگ آماری اخلاق‌مدار باشند و از انجام کارهای آماری با کیفیت و متکی به شیوه‌نامه‌های اخلاقی حمایت کنند. در این راستا، اقدامات مشخصی لازم است که هم‌راستایی سازمان‌ها و نهادها با شیوه‌نامه‌های اخلاقی و حرفه‌ای ایجاد شود، برای این که محیطی امن، سازنده و احترام‌آمیز برای تمامی اعضای تیم‌های آماری فراهم شود.

(اصل ت ۱) **استادان راهنما، راهبران، سرپرستان و مربیان:** استادان راهنما، رهبران، سرپرستان و مربیان، انجام‌دهندگان اصلی، ناظران و مشاوران کارهای آماری را به پیروی، ترویج و حمایت از اصول اخلاقی متعهد می‌کنند و بر انجام کارهای آماری اخلاق‌مدار توسط کارگروه مجری تاکید دارند.

(اصل ت ۲) **سازمان‌ها، موسسات و نهادها:** سازمان‌ها، مؤسسات و نهادهای درگیر در فعالیتهای آماری (مانند جمع‌آوری، خلاصه‌سازی، پردازش، تحلیل، تفسیر، یا ارائه نتایج آماری با استفاده از انواع داده‌ها)، مسئولیت دارند که کارهای آماری اخلاق‌مدار را ترویج نمایند.

(اصل ت ۳) **تعهدات اخلاقی در مواجهه با تخلفات آماری و حرفه‌ای:** همواره این امکان وجود دارد که انجام‌دهنده کارهای آماری اخلاق‌مدار با شبهاتی در ارتباط با تخلفات احتمالی در زمینه کارهای آماری، شیوه‌های علمی یا عملکرد حرفه‌ای مواجه شود. گاهی اوقات، یک انجام‌دهنده ممکن است فردی را به تخلف متهم کند یا توسط دیگران به تخلفی متهم شود. در برخی مواقع نیز، یک انجام‌دهنده ممکن است در رسیدگی به تخلفات دیگران مشارکت داشته باشد. انجام دهنده کارهای آماری اخلاق‌مدار، تعاریف و شیوه‌های مربوط به اتهام‌های تخلف در حوزه‌های آماری در محیط سازمانی خود را می‌داند، تفاوت رفتار غیراخلاقی با اختلاف نظر و اشتباه صادقانه را تشخیص می‌دهد و قبل از اظهارنظر درباره تخلف دیگران به دنبال کشف حقایق و انگیزه می‌گردد.

هر کدام از اصول بالا، دارای چندین بند است که توضیحات آنها در شیوه‌نامه اخلاقی کارهای آماری انجمن آمار (۱۲) آمده است. هدف اصلی این شیوه‌نامه و شیوه‌نامه‌های مشابه، تبیین استانداردهای اخلاقی و ایجاد فرهنگ پاسخگویی و مسئولیت‌پذیری در میان افرادی است که به نحوی با داده‌ها، روش‌ها و تحلیل‌های آماری و تفسیر نتایج سر و کار دارند. این شیوه‌نامه‌ها به دانشجویان، پژوهشگران و تمامی افرادی که در حوزه‌های مختلف از روش‌های آماری استفاده می‌کنند، کمک می‌کند تا مسئولانه و با رعایت اصول اخلاقی به انجام کارهای آماری بپردازند. پایبندی به اصول اخلاقی، برای تضمین کیفیت، اعتبار علمی و اعتماد عمومی به نتایج کارهای آماری ضروری است. موضوع اخلاق در آمار، جنبه‌های متفاوتی قابل بحث و بررسی است. برای درک بیشتر اهمیت موضوع، می‌توان به کتاب‌هایی که در این زمینه نوشته شده است مراجعه کرد (۱۸؛ ۲۴؛ ۲۶).

۴ اخلاق و آموزش آمار

همان گونه که در دانشگاه طیف وسیعی از مباحث آمار به دانشجویان آموزش داده می‌شود تا برای فرصت‌های تحصیلی و حرفه‌ای مختلف آماده شوند، منطقاً باید طیف وسیعی از موقعیتهای اخلاقی مرتبط با انجام کارهای آماری در زمینه‌های مختلف نیز آموزش داده شود. ترویج اخلاق در انجام فعالیت‌های مرتبط با آمار، بایستی از طریق گنجاندن آموزش اخلاق در برنامه‌های درسی انجام شود تا بتواند در شکل‌دهی به رفتار حرفه‌ای دانش‌آموختگان موثر باشد. گنجاندن اصول اخلاق در آموزش آمار را می‌توان با استفاده از رویکردهای متنوعی مانند: گنجاندن فعالیت‌ها یا تکلیف‌های خاصی در خلال تدریس برخی از دروس، توزیع مفاهیم اخلاقی در محتوای برخی از دروس، برگزاری کارگاه‌های آموزشی خارج از برنامه درسی رسمی، واگذاری مسئولیت آموزش اخلاق، به دوره‌های ارائه شده توسط متخصصین رشته‌های مرتبط و نهایتاً گنجاندن درسی مستقل در برنامه درسی انجام داد. در بسیاری از دانشگاه‌ها و مؤسسات آموزشی دنیا دروس یا دوره‌هایی با عنوان «اخلاق علمی» تدریس می‌شود که به مباحث عمومی مرتبط با ملاحظات اخلاقی پژوهش و انتشار نتایج آن اختصاص دارد. از جمله:

- آگاهی در مورد تخلفات علمی مانند سرقت علمی و ادبی و تشریح پیامدهای آن برای اعتبار فرد و جامعه علمی
 - آگاهی از تفاوت تخلف علمی و اشتباه صادقانه و نحوه برخورد با تخلفات علمی
 - نحوه عملکرد در اصلاح اشتباهات اساسی، پس از انتشار یک پژوهش
 - آموزش نحوه ارجاع به مقالات و منابع علمی در قالب استاندارد
 - اهمیت صداقت، شفافیت در گزارش و انتشار نتایج پژوهش‌ها و مسئولیت‌پذیری
 - اصول اخلاقی در همکاری‌های علمی گروهی
 - برخورد صحیح با همکاران در صورت وجود اختلاف نظر علمی
 - تقسیم عادلانه کار و رعایت انصاف در نوشتن مقالات مشترک و سهم اعضای گروه
 - اصول اخلاقی و قوانین و مقررات مربوط به نوشتن گزارش‌های علمی
 - آگاهی از فرآیند بررسی نتایج پژوهش توسط هم‌تایان و نحوه تعامل با داوران.
- در برخی از رشته‌ها در کنار آموزش ملاحظات عمومی اخلاق علمی در انتشار نتایج پژوهش‌ها، دوره‌ها یا دروسی با عنوان «اخلاق حرفه‌ای» نیز ارائه می‌شود. در مرجع (۲۰) محتوای ۱۱۵ عنوان درس در زمینه آموزش اخلاق در رشته‌های مختلف، در برخی از دانشگاه‌های دنیا (امریکا، کانادا، اروپا و آسیا) مورد بررسی و تحلیل و دسته‌بندی موضوعی قرار گرفته است. این دسته‌بندی طیف وسیعی از ملاحظات اخلاقی را شامل می‌شود. برخی از این موارد مختص رشته‌های مربوط و برخی نیز عمومی هستند. مسائل اخلاقی مرتبط با هوش مصنوعی و فناوری‌های نوین در اغلب دروس ظاهر می‌شود. از جمله:
- حفظ حریم خصوصی افراد و قوانین مربوط به محرمانگی و امنیت داده‌ها
 - ناشناس‌سازی اطلاعات و ردپای افراد در کلان داده‌ها (داده‌های بزرگ) و داده‌های دیجیتال
 - توجه به محدودیت‌های مربوط به تلقی داده‌های دیجیتال به عنوان دارایی شرکت‌ها
 - رعایت مالکیت معنوی و قوانین کپی رایت و نسخه‌برداری
 - دوری از انواع تبعیض و نابرابری اجتماعی و آزار و اذیت افراد
 - حمایت از آزادی بیان، حقوق بشر و حفظ ارزش‌های دموکراسی
 - مسئولیت در قبال محیط زیست و حیوانات
 - مسئولیت در قبال اخبار جعلی، مطالب نادرست و شبه علم.
- در زمینه ضرورت و چگونگی گنجانیدن آموزش ملاحظات اخلاقی در برنامه درسی رشته آمار و همچنین علوم ریاضی، مقالات زیادی نوشته شده است و نیز در برخی از دانشگاه‌های دنیا، دروسی نیز تعریف شده است. به عنوان نمونه مراجع (۱۳؛ ۱۶؛ ۲۳؛ ۲۷؛ ۲۸؛ ۲۹)

را ملاحظه کنید. در اینجا به محتوای برخی از این دروس، به عنوان نمونه اشاره می‌کنیم. چیدو و مولیر (۱۵) در گزارشی که اخیراً منتشر نموده‌اند، رویکردی برای گنجانیدن آموزش اخلاق در علوم ریاضی در برنامه دوره سطح کارشناسی ارائه می‌دهند. این درس برای دانشجویان رشته‌های مهندسی، علوم کامپیوتر، آمار، فیزیک و اقتصاد نیز قابل تدریس است. مواد آموزشی درس به سه نوع تقسیم شده است: تمرین‌ها، پروژه‌ها و مطالب درسی. تمرین‌های انتخابی که برای آموزش ملاحظات اخلاقی، استفاده شده است، در حد آمار مقدماتی و ریاضیات عمومی سال اول و دوم مقطع کارشناسی است. برخی از تمرین‌های اساسی‌تر، به عنوان پروژه‌های مستقل به دانشجویان داده می‌شود تا ملاحظات اخلاقی آنها را به صورت گروهی، به طور عمیق‌تر بررسی کنند. این گونه مسائل، معمولاً به تحقیق و مطالعه بیشتری نیاز دارند و از آنها به عنوان بخشی از ارزیابی دوره استفاده می‌شود. فهرستی از مطالب مربوط به اخلاق در آمار، ریاضی و سایر علوم در این گزارش ارائه شده است که مدرس، با توجه به مخاطبان خود، می‌تواند از آنها استفاده کند. تمرین‌ها به طور موازی هم برای آموزش مطالب نظری و هم برای آگاهی جنبه اخلاقی آنها طراحی شده‌اند و حوزه‌های مختلفی را پوشش می‌دهند. هر تمرین دارای برجستگی است که موضوع رشته مرتبط و یک رکن اخلاقی را توضیح می‌دهد. برجستگی‌هایی که برای رشته‌های مرتبط استفاده شده است شامل آمار و احتمال، علوم کامپیوتر، اقتصاد، مهندسی و علوم طبیعی است. برجستگی‌های مربوط به حوزه‌های اخلاقی استفاده شده در تمرین‌ها، شامل ۱۰ مورد به شرح زیر است (۱۵):

- ۱ - توجیه آغاز یک کار: توجیه شما برای انجام یک کار آماری چیست؟ آیا اساساً انجام آن موجه است؟
- ۲ - بهره‌گیری از تنوع دیدگاه و اجتناب از سوگیری‌ها: آیا از تنوع دیدگاه کافی در بین همکاران و مدیران برخوردارید؟ آیا محدودیت‌ها و سوگیری‌های شناختی خود را درک می‌کنید؟
- ۳ - مدیریت مسئولانه و اثربخش داده‌ها: آیا منحصراً از مجموعه داده‌های مجاز و به‌دست‌آمده به شیوه‌ای اخلاقی استفاده می‌کنید؟
- ۴ - پردازش و استنتاج داده‌ها با دقت و صحت: آیا از تخصص کافی برای تحلیل صحیح و اخلاقی داده‌ها و تضمین کیفیت نتایج برخوردارید؟
- ۵ - مدل‌سازی ریاضی مسئله و پیامدهای واقعی آن: اهداف و محدودیت‌های مدل انتخابی شما چیست؟ پیامدهای واقعی آن برای ذینفعان مختلف کدامند؟
- ۶ - برقراری ارتباط شفاف و مستندسازی دقیق فرایند کار: چگونه کار خود را به طور شفاف شرح داده و مستند می‌کنید؟ چگونه نتایج را به ذینفعان مربوطه منتقل می‌کنید؟
- ۷ - قابلیت ابطال و ایجاد حلقه بازخورد مستمر: آیا کار شما قابل ابطال است؟ چگونه تأثیرات گسترده و سازوکارهای بازخورد ایجاد شده را مدیریت می‌کنید؟
- ۸ - قابلیت تفسیر و ایمنی مدل: آیا خروجی مدل شما قابل تفسیر است؟ آیا سازوکارهای نظارت و نگهداری مناسب و توسعه آن پیش‌بینی شده است؟
- ۹ - آگاهی از جنبه‌های سیاسی و اجتماعی مرتبط با کار: آیا از جنبه‌های غیر فنی و ماهیت سیاسی و اجتماعی کار خود آگاه هستید؟ برای جلب اعتماد ذینفعان به کار و محصول خود چه تدابیری اندیشیده‌اید؟

۱۰ - تدوین استراتژی‌های واکنش اضطراری مناسب: آیا استراتژی دادن پاسخ غیرفنی مناسب، برای مواجهه با شرایط پیش‌بینی نشده را دارید؟ آیا یک شبکه پشتیبانی متشکل از همکاران آگاه و حامی در اختیار دارید؟

نمونه‌هایی از تمرین‌ها و پروژه‌های این درس به همراه راه حل‌های پیشنهادی، به طور مفصل در گزارش (۱۵) آمده است. در وبگاه انجمن علم داده (۳۰)^۸، وبگاه مرکز اخلاق مهندسی و علوم (۳۱)^۹، وبگاه نشریه انجمن اخلاق در برنامه درسی (۳۲)^{۱۰} و همچنین وبگاه مرکز ملی آموزش آمار (۳۳)^{۱۱} می‌توان مطالب مفیدی در زمینه آموزش اخلاق در آمار و ریاضی جستجو نمود.

۱.۴ مثالی از آموزش ملاحظات اخلاقی در احتمال

شما با دوستان عباس در تعطیلات به کوه نوردی رفته‌اید. متأسفانه، عباس روی یک سنگ می‌افتد و پای راستش می‌شکند و به یک بیمارستان کوچک شهر نزدیک منتقل می‌شود و مجبور است یک شب دیگر در بیمارستان بماند. برق بیمارستان هر ساعت به صورت مستقل با احتمال p قطع می‌شود.

(الف) فرمولی برای تعداد ساعات مورد انتظار تا قطعی برق بیمارستان استخراج کنید.

(ب) مقدار p باید چقدر کوچک باشد تا تعداد مورد انتظار قطعی برق در یک سال کمتر از ۱ باشد؟

(پ) عباس باید x ساعت دیگر در بیمارستان بماند. اگر برق قطع نشود، او به موقع مرخص می‌شود و شما دو نفر هنوز می‌توانید به

پرواز برگشتتان برسید. احتمال اینکه برق در $x - ۱$ ساعت اول قطع شود چقدر است؟

(ت) عباس از شما می‌پرسد که چه مقدار از p شما را آسوده‌خاطر می‌کند. چگونه تصمیم می‌گیرید؟

پاسخ مسئله

(الف) این مسئله از توزیع هندسی پیروی می‌کند. امیدریاضی متغیر تصادفی X ، تعداد ساعات تا اولین قطعی برق برابر است با

$$E(X) = \frac{1}{p}$$

(ب) می‌خواهیم مقدار p را طوری تعیین کنیم که تعداد قطعی‌های مورد انتظار در یک سال کمتر از ۱ باشد. تعداد ساعات در یک سال

$۸۷۶۰ \times ۲۴ = ۳۶۵$ ساعت و تعداد مورد انتظار قطعی‌ها در سال $۸۷۶۰p$ است. شرط مسئله:

$$۸۷۶۰p < ۱ \implies p < \frac{1}{۸۷۶۰} \approx ۰/۰۰۰۱۱۴$$

(پ) احتمال قطع نشدن برق در یک ساعت $1 - p$ ، احتمال قطع نشدن برق در $x - ۱$ ساعت متوالی $(1 - p)^{x-1}$ است. بنابراین

احتمال حداقل یک قطعی

$$P = 1 - (1 - p)^{x-1}$$

⁸Data Science Association

⁹Ethics Center for Engineering and Science website

¹⁰Society of Ethics Across the Curriculum website

¹¹National Center for Education Statistics (NCES)

(ت) برای انتخاب مقدار مناسب p آیا دانشجوی فقط به دوست خود (عباس) فکر می‌کند یا بیماران دیگر را نیز در نظر می‌گیرد

- دیدگاه فردی: اگر فقط به پرواز بازگشت فکر کنیم: احتمال قطع برق در x ساعت

$$1 - (1 - p)^x \leq 0.01$$

برای $x = 24$ ساعت، باید $p \approx 0.0004$ باشد.

- دیدگاه اجتماعی و ملاحظه اخلاقی: با در نظر گرفتن ایمنی همه بیماران مقدار p باید بسیار کوچک باشد مثلاً:

$$p < 0.0001.$$

- نتیجه‌گیری: مقدار p باید آنقدر کوچک باشد که برای عباس و سایر بیماران ایمنی قابل قبولی فراهم کند. پیشنهاد می‌شود:

$$p \leq 0.0001$$

این مثال، به موارد ۲ و ۵ و ۶ از ملاحظات اخلاقی ده گانه بالا اشاره دارد.

آیا آماردانان باید اخلاق را تدریس کنند؟

استادان و اعضای هیئت علمی خود باید به عنوان الگوهای اخلاقی برای دانشجویان خود عمل کنند. آن‌ها باید در انجام تحقیقات علمی و نشر مقالات به اصول اخلاقی پایبند باشند تا دانشجویان از آن‌ها الگو بگیرند. در نهایت، آموزش اخلاق در آمار می‌تواند دانشجویان را قادر سازد تا در مسیر حرفه‌ای خود با مسئولیت‌پذیری، دقت و صداقت عمل کنند و به تقویت اعتبار و سلامت علمی جامعه آماری کمک نمایند. این سؤال پیش می‌آید که آیا آماردانان باید به طور رسمی محتوای درس اخلاق در آمار را آموزش دهند؟ از آنجا که پیامدهای اخلاقی کارهای آماری را تنها متخصصین آموزش دیده در آمار می‌توانند درک کنند، خود آماردانان علاقه‌مند بهترین گزینه برای آموزش ملاحظات اخلاقی در آمار هستند. این مسائل را نمی‌توان صرفاً به متخصصان اخلاق و فلسفه واگذار نمود. برون‌سپاری آن به عنوان یک درس سرویسی یا تدریس مشترک با متخصصین اخلاق، راهکار دیگری است. البته دانشجویان آمار می‌توانند از طریق شرکت در دوره‌های ارائه شده توسط سایر گروه‌های آموزشی مرتبط، در مورد اخلاق حرفه‌ای آموزش ببینند. شرکت در چنین دوره‌هایی، اگرچه تمرکز اصلی بر کاربردهای آماری نیست، اما می‌تواند دانشجویان آمار را با اصول اخلاقی ارزشمندی آشنا کند.

بحث و نتیجه‌گیری

همان گونه که در این مقاله اشاره شد، در ایران در زمینه اخلاق علمی و حرفه‌ای، قوانین، دستورالعمل‌ها، شیوه‌نامه‌ها، نظام‌نامه‌ها و راهنماهای مختلفی در حوزه‌هایی مانند وکالت، مهندسی، روان‌شناسی و مشاوره، حسابداری و پزشکی، وجود دارد. در مورد ضرورت آموزش اخلاق در رشته‌های دانشگاهی نیز مقالاتی نوشته شده است. در برخی از رشته‌های مهندسی و علوم پزشکی دروسی با محتوای اخلاق علمی و حرفه‌ای وجود دارد. در زمینه اخلاق در آمار و آموزش آن، در برنامه‌های درسی دانشگاه‌های

کشور تاکنون کاری انجام نشده است. اخیراً انجمن آمار ایران با تشکیل «کارگروه اخلاق آماری»، «شیوه‌نامه اخلاقی برای کارهای آماری» خود را تدوین و منتشر نموده است. آموزش ملاحظات اخلاقی در برنامه درسی آمار، امری ضروری است که نباید در برابر پیچیدگی‌های فنی و نظری این رشته نادیده گرفته شود. به‌ویژه با توجه به کاربردهای گسترده آمار در حوزه‌های مختلف علمی، مهندسی، پزشکی، اقتصاد و تصمیم‌گیری‌های اجتماعی، توجه به اصول اخلاقی می‌تواند مانع از سوءاستفاده‌ها و خطاهای فاحش شود. جا دارد انجمن آمار ایران، با قرار دادن عنوان «اخلاق در آمار» در محورهای همایش‌های علمی انجمن در جهت ترویج اخلاق در آمار و جلب توجه آماردانان علاقه‌مند کشور به پژوهش در این زمینه، گام بردارد. گام‌های بعدی، می‌تواند تهیه محتوای آموزشی برای افزودن به برنامه درسی رشته آمار باشد. ترویج اصول اخلاقی در جامعه آمار و آموزش آن، به تقویت مسئولیت‌پذیری حرفه‌ای دانش‌آموختگان در انجام کارهای آماری کمک خواهد که به نوبه خود باعث ارتقای کیفیت و اعتبار نتایج حاصل از کارهای آماری آماری در جامعه و صنعت خواهد شد.

قدردانی و تشکر

این مقاله به پیشنهاد کارگروه اخلاق انجمن آمار ایران (متشکل از آقایان دکتر محمد آرشی، دکتر علی دولتی، دکتر علی رجالی، دکتر محمدرضا فقیهی، دکتر رحیم محمودوند و دکتر حمید نیلی ثانی) توسط نویسنده تهیه شده است. از نظرات ارزشمند این عزیزان در تدوین مقاله قدردانی می‌کنم.

مراجع

- [۱] امین خندقی، م. و پاک مهر، ح. (۱۳۹۱). آموزش معیارهای اخلاق پژوهش: ضرورتی انکارناپذیر در برنامه‌های درسی آموزش عالی، فصلنامه اخلاق در علوم و فناوری، شماره ۴، صفحات ۱۳-۱۰.
- [۲] بکروی، م. و همکاران (۱۳۹۹). اخلاق حرفه‌ای در علوم کامپیوتر، انتشارات دانشگاه آزاد.
- [۳] ثقه‌الاسلامی، ع. (۱۴۰۰). برنامه‌ریزی درسی و تدوین سرفصل‌های آموزشی اخلاق پژوهش، فصلنامه اخلاق در علوم و فناوری، شماره ۳، صفحات ۱۰۱-۹۲.
- [۴] خوشدست، ح. و سام، ع. (۱۳۸۸). ارائه الگویی برای آموزش مؤثر اخلاق مهندسی در دوره کارشناسی، فصلنامه آموزش مهندسی ایران، سال ۱۱، شماره ۴۳، صفحه ۱۰۸-۹۹.
- [۵] رجبعلی‌پور، م.، رستگاری، ف. و قنبری‌زاده، ف. (۱۴۰۱). بررسی دروس مرتبط با اخلاق حرفه‌ای در برنامه درسی رشته‌های مختلف علوم پزشکی در ایران، دو فصلنامه آموزش و اخلاق در پرستاری، دوره ۱۱، شماره ۳ و ۴، صفحات ۴۱-۳۲.
- [۶] فراستخواه، م. (۱۳۸۵). اخلاق علمی رمز ارتقای آموزش عالی، فصلنامه اخلاق در علوم و فناوری، شماره ۱، صفحات ۲۷-۱۳.

- [۷] منشور اخلاق حرفه‌ای وکالت (۱۳۹۹). هیات مدیره کانون وکلای دادگستری.
- [۸] نظامنامه رفتار حرفه‌ای اخلاقی در مهندسی ساختمان (۱۳۹۵). مصوب وزارت راه و شهرسازی.
- [۹] نظامنامه اخلاق حرفه‌ای سازمان نظام روان‌شناسی (۱۳۸۶). سازمان نظام روان‌شناسی و مشاوره.
- [۱۰] آئین رفتار حرفه‌ای برای حسابداران (۱۳۹۶). جامعه حسابداران رسمی ایران.
- [۱۱] قوانین، دستورالعمل‌ها و راهنماهای اخلاق در پژوهش‌های زیست پزشکی (۱۴۰۱). وزارت بهداشت و درمان و آموزش پزشکی.
- [۱۲] شیوه‌نامه اخلاقی برای کارهای آماری (۱۴۰۴)، کارگروه اخلاق انجمن آمار ایران.
- [13] Alayont, F. (2022). A case for ethics in the mathematics major curriculum. *Journal of Humanistic Mathematics*, 12(2), 160-177.
- [14] American Statistical Association (ASA). (2018). *Ethical Guidelines for Statistical Practice*.
- [15] Chiodo, M. and Müller, D. (2025). Teaching resources for embedding ethics in mathematics: Exercises, projects, and handouts. arXiv preprint arXiv:2310.08467.
- [16] Ernest, P. (2024). *Ethics and Mathematics Education*. Cham, Switzerland: Springer.
- [17] Code, N. (1949). The Nuremberg code. *Trials of war criminals before the Nuremberg military tribunals under control council law*, 10(1949), 181-2.
- [18] Doosti, H. (Ed). (2024). *Ethics in Statistics : Opportunities and Challenges* , Ethics International Press.
- [19] Ellenberg, J. H. (1983). Ethical guidelines for statistical practice: A historical perspective. *The American Statistician*, 37(1), 1-4.
- [20] Fiesler, C., Garrett, N. and Beard, N. (2020). What do we teach when we teach tech ethics? A syllabi analysis, pages 289–295 in *Proceedings of the 51st ACM Technical Symposium on Computer Science Education*.
- [21] Declaration of Helsinki (1964). <https://www.britannica.com/topic/Declaration-of-Helsinki>.
- [22] Hurwitz, S. and Gardenier, J. S. (2012). Ethical guidelines for statistical practice: The first 60 years and beyond. *The American Statistician*, 66(2), 99-103.
- [23] Lesser, L. M. and Nordenhaug, E. (2004). Ethical statistics and statistical ethics: Making an interdisciplinary module. *Journal of Statistics Education*, 12(3).

- [24] Panter, A. T. (Ed.). (2011). Handbook of ethics in quantitative methodology. Taylor Francis.
- [25] Royal Statistical Society. 2014. Code of Conduct, RSS, Accessed 04.07.2024. <https://rss.org.uk/about/policy-and-guidelines/code-of-conduct/>.
- [26] Smeyers, P. (2010). Educational research-the ethics and aesthetics of statistics. Springer Science, Business Media.
- [27] Tractenberg, R. E. (2016). Why and How the ASA Ethical Guidelines should be integrated into every quantitative course. In Proceedings of the 2016 Joint Statistical Meetings, Chicago, IL, USA (pp. 517-535).
- [28] Tractenberg, R. (2023). Ethical Practice of Statistics and Data Science: Ethical Practice of Statistics and Data Science. Ethics International Press.
- [29] Vardeman, S. B. and Morris, M. D. (2003). Statistics and ethics: some advice for young statisticians. The American Statistician, 57(1), 21-26.
- [30] Data Science Association <https://www.datascienceassn.org>
- [31] Ethics Center for Engineering and Science website, <https://onlineethics.org>.
- [32] Society of Ethics Across the Curriculum website, <https://www.seac-online.org/resources/>
- [33] National Center for Education Statistics (NCES) <https://nces.ed.gov/forum/dataethics/>

Ethics in Statistics and Probability

Ali Dolati¹

¹Yazd University, Department of Statistics

Abstract: Today, with the expansion of the use of statistical data and methods in waste, there are unethical practices such as the misuse of data, the presentation of biased and inaccurate content of the results of analyses, and the lack of transparency in the presentation of reports, in an informed or biased manner. Therefore, it is vital to strengthen public confidence in statistics and to maintain the health and credibility of the field, to educate and promote ethics. The purpose of this paper is to provide a model for the inclusion of ethics in statistics courses, so that students can learn about social and professional development in addition to learning the theoretical and applied concepts of statistics.

Keywords: Statistics education, ethics in statistics, curriculum, professional training, data privacy.

Mathematics Subject Classification (2020): 97k80, 97k99.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30-31, 2025
University of Kurdistan



نقش فرآیندهای تصادفی در هوش مصنوعی و علم داده: یک مطالعه موردی

اعظم راستین^۱

دکترای آمار، دانشگاه شهید چمران اهواز

چکیده:

چکیده

در این مقاله، نقش تاثیرگذار فرآیندهای تصادفی در هوش مصنوعی و علم داده با تمرکز ویژه بر یادگیری تقویتی در محیط *Gridworld* به همراه یک مطالعه موردی مورد بررسی قرار می‌گیرد. فرآیندهای تصادفی، به‌ویژه مدل‌سازی‌های مارکوفی، چارچوبی ریاضی و مفهومی برای تحلیل و پیش‌بینی رفتار عامل‌های هوشمند در محیط‌های نامعین فراهم می‌آورند. در این راستا، الگوریتم $Q - Learning$ در سه سناریوی مختلف (محیط قطعی، محیط با نویز تصادفی، و محیط با پاداش‌های تصادفی) پیاده‌سازی و تحلیل شده است. نتایج عددی، نمودارهای همگرایی، و مسیرهای یادگیری عامل‌ها نشان می‌دهد که چگونه فرآیندهای تصادفی درک دقیق‌تری از تعامل عامل با محیط فراهم کرده و بر کیفیت یادگیری تأثیر می‌گذارند. این مقاله نشان می‌دهد که ادغام تئوری احتمال و فرآیندهای تصادفی در ساختارهای هوش مصنوعی، مسیر توسعه الگوریتم‌های پایدارتر و تطبیق‌پذیرتر را هموار می‌سازد. واژه‌های کلیدی: فرآیند تصادفی، یادگیری تقویتی، $Gridworld$ ، MDP ، یادگیری سیاست، ارزش‌گذاری، مدل‌سازی محیط کد موضوع‌بندی ریاضی (۲۰۲۰): 68T20، 68T05، 60J20.

۱ مقدمه

فرآیندهای تصادفی یکی از ابزارهای بنیادی در مدل‌سازی پدیده‌هایی هستند که در آن‌ها عدم قطعیت و تصادف نقش دارند. در دهه‌های اخیر، اهمیت این فرآیندها در حوزه‌هایی مانند هوش مصنوعی و علم داده نیز به طور چشمگیری افزایش یافته است. الگوریتم‌های یادگیری ماشین، تحلیل سری‌های زمانی، یادگیری تقویتی و بسیاری دیگر از شاخه‌های هوش مصنوعی برای مدل‌سازی رفتارهای پیچیده و پیش‌بینی‌های دقیق، به نظریه‌های تصادفی متکی هستند. فرآیندهای تصادفی امکان مدل‌سازی

^۱ سخنران، rastinstat@gmail.com

پویایی‌هایی را فراهم می‌کنند که در آن‌ها آینده سیستم به صورت احتمالی و بر اساس گذشته یا وضعیت فعلی تعیین می‌شود. به ویژه، زنجیره‌های مارکوف، فرایندهای تصمیم‌گیری مارکوف، و فرایندهای گاوسی، ابزارهایی کلیدی در توسعه سیستم‌های هوشمند و تحلیل داده‌های پیچیده هستند. در سال‌های اخیر، هوش مصنوعی و علم داده به شدت به مفاهیم احتمالاتی و فرایندهای تصادفی متکی شده‌اند تا بتوانند رفتار پیچیده و غیرقطعی سیستم‌های واقعی را مدل‌سازی کنند. یکی از حوزه‌هایی که به طور مستقیم از این ابزارها بهره‌مند است، یادگیری تقویتی است. یادگیری تقویتی یکی از شاخه‌های مهم یادگیری ماشینی است که به طور خاص برای تصمیم‌گیری در محیط‌های پویا و تعاملی طراحی شده است. در این روش، عامل از طریق تعامل با محیط و بر اساس بازخورد پاداش، سیاستی بهینه برای بیشینه‌سازی پاداش بلندمدت فرا می‌گیرد (سوتون و باروت، ۲۰۱۸).

در سال‌های اخیر، یادگیری تقویتی به عنوان یکی از زیرشاخه‌های مهم هوش مصنوعی، نقش چشم‌گیری در حل مسائل تصمیم‌گیری در محیط‌های نامطمئن و تصادفی ایفا کرده است. مطالعات متعددی بر اهمیت پایداری الگوریتم‌های یادگیری تقویتی در مواجهه با تغییرات دینامیکی محیط تأکید داشته‌اند (لادوز و همکاران، ۲۰۲۲). با گسترش کاربردهای یادگیری تقویتی در محیط‌های بزرگ و با فضای کنش گسترده، روش‌های جدیدی برای مقابله با پیچیدگی‌های محاسباتی ارائه شده‌اند. به عنوان نمونه، فورانی و همکاران (۲۰۲۴) یک نسخه تصادفی از الگوریتم $Q-learning$ را معرفی کردند که امکان یادگیری در محیط‌هایی با فضای کنش بسیار بزرگ را فراهم می‌کند.

فرایندهای تصادفی نقشی کلیدی در مدل‌سازی عدم قطعیت در یادگیری تقویتی دارند. مطالعات پیشین نشان داده‌اند که مدل‌سازی دقیق نویزها و تصادفی بودن پاداش یا انتقال حالت می‌تواند به طراحی الگوریتم‌های پایداری و واقع‌گرایانه‌تر منجر شود (مراجعه شود به ترون و همکاران (۱۹۹۵) و کرنز و سینگز (۲۰۰۲)).

در این مقاله، ما محیط $Gridworld$ را به عنوان بستری ساده ولی گویا برای بررسی تأثیر فرایندهای تصادفی در مدل‌سازی عدم قطعیت انتخاب کرده‌ایم. با استفاده از الگوریتم $Q-learning$ و در سه سناریوی مجزا-محیط قطعی، محیط با نویز در انتقال، و محیط با پاداش تصادفی-روند یادگیری و عملکرد عامل مورد تحلیل قرار گرفته‌اند.

۲ مفاهیم پایه

فرایند تصادفی به توالی از متغیرهای تصادفی، $\{X_t\}_{t \in T}$ ، گفته می‌شود که به صورت تابعی از پارامتر t تعریف شده‌اند که در آن T یک مجموعه اندیس گذار بوده و می‌تواند گسسته (مثلاً زمان‌های طبیعی) یا پیوسته (مثلاً بازه‌ای از اعداد حقیقی) باشد. یک مورد خاص از فرایندهای تصادفی، زنجیره مارکوف است که دارای ویژگی زیر می‌باشد:

$$P(X_{t+1} = x | X_t = x_t, X_{t-1} = x_{t-1}, \dots) = P(X_{t+1} = x | X_t = x_t), t \in T$$

به این خاصیت، ویژگی مارکوفی یا حافظه کوتاه گفته می‌شود. زنجیره‌های مارکوف به طور گسترده‌ای در الگوریتم‌های یادگیری، مدل‌سازی تصمیم‌گیری، و تحلیل سری‌های زمانی به کار می‌روند. فرایندهای تصادفی مبنای ریاضی مدل‌سازی رفتار سیستم‌هایی هستند که با گذر زمان و تحت تأثیر عدم قطعیت تغییر می‌کنند. در یادگیری تقویتی، به ویژه از فرایندهای تصمیم‌گیری مارکوفی (MDP) استفاده می‌شود (پورتن، ۲۰۰۵).

تعریف ۱.۲. فرایند تصمیم‌گیری مارکوف (MDP): یک MDP که به اختصار آن را با $M = (S, A, P, R, \gamma)$ نشان می‌دهیم،

شامل مجموعه‌ای از حالات ممکن S ، اعمال ممکن A ، تابع احتمال انتقال $P(s' | s, a)$ (احتمال انتقال از حالت s به حالت s' با انتخاب عمل a) و تابع پاداش $R(s, a)$ (پاداش دریافتی برای اجرای عمل a در حالت s) و ضریب تنزیل آینده $\gamma \in [0, 1]$ می‌باشد.

توجه کنید که طبق خاصیت مارکوفی، $P(s_{t+1} | s_t, a_t, s_{t-1}, \dots, s_0) = P(s_{t+1} | s_t, a_t)$. در واقع، این خاصیت تضمین می‌کند که وضعیت آینده فقط به وضعیت و عمل فعلی بستگی دارد، نه به مسیر طی شده قبلی.

تعریف ۲.۲. سیاست ($Policy$): تابعی که مشخص می‌کند عامل در هر حالت چه عملی انجام دهد. سیاست بهینه π^* باعث بیشینه شدن امید ریاضی پاداش تجمعی می‌شود.

تعریف ۳.۲. تابع Q : بیانگر ارزش حالت-عمل تحت یک سیاست خاص است:

$$Q(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

هدف اصلی در MDP ، یافتن سیاست بهینه π^* است به طوری که تابع ارزش Q بیشینه شود. یادگیری تقویتی نوعی یادگیری ماشینی است که در آن عامل با تعامل با محیط، سیاست بهینه را از طریق آزمون و خطا می‌آموزد. در این تحقیق از الگوریتم $Q-Learning$ برای این منظور استفاده شده است. $Q-Learning$ یکی از الگوریتم‌های پایه در یادگیری تقویتی است که بدون مدل محیط و به صورت زیر کار می‌کند:

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

که در آن α, r, s' و γ به ترتیب بیانگر نرخ یادگیری، پاداش دریافتی در انتقال، حالت بعدی پس از انجام a در s و ضریب تنزیل است. با گذر زمان، مقدارهای Q به سمت مقدارهای بهینه همگرا می‌شوند (واتکینز و دایان، ۱۹۹۲).

۳ مطالعه موردی: یادگیری تقویتی در محیط $Gridworld$

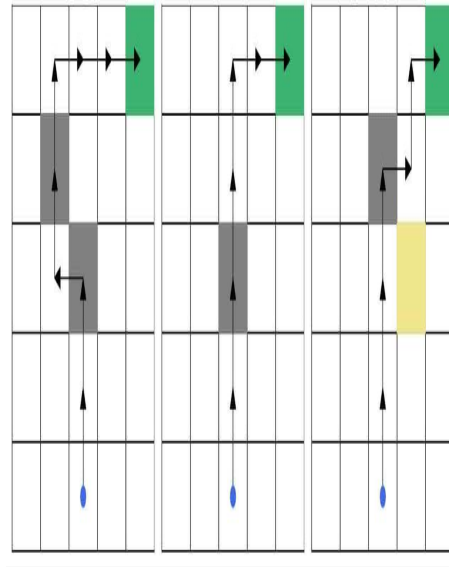
محیط $Gridworld$ یک شبکه دوبعدی است که هر خانه از آن یک حالت را نشان می‌دهد. در این شبکه یک عامل ($Agent$) می‌تواند به چهار جهت اصلی (بالا، پایین، چپ، راست) حرکت کند. بعضی خانه‌ها می‌توانند موانع یا چاله‌هایی باشند که مانع حرکت عامل می‌شوند. هدف عامل، در حالی که از موانع اجتناب و پاداش حداکثری دریافت کند، رسیدن به خانه هدف می‌باشد. در این مطالعه سه نوع محیط تعریف شده است:

۱. محیط قطعی: انتقال دقیق و بدون نویز به جهت دلخواه

۲. محیط نویزدار: احتمال ۸۰٪ برای حرکت مورد نظر و ۲۰٪ برای انحراف به جهت دیگر

۳. محیط با پاداش‌های تصادفی: مقدار پاداش در برخی خانه‌ها به صورت تصادفی تغییر می‌کند

در شکل ۱ این سه محیط متفاوت نمایش داده شده است. این تصویر سه محیط مختلف را در قالب شبکه 5×5 نشان می‌دهد که در هر کدام شرایط متفاوتی حاکم است: در محیط قطعی، حرکت عامل با دقت کامل در جهت اعمال شده صورت می‌گیرد و پاداش‌ها ثابت هستند. در محیط نویزی، عامل با احتمال ۰.۸ در جهت اعمال شده حرکت می‌کند و با احتمال ۰.۲



شکل ۱: نمایش سه محیط Gridworld با ویژگی‌های مختلف: محیط قطعی، محیط با نویز، محیط با پاداش تصادفی.

به جهت‌های دیگر می‌رود. و در محیط با پاداش تصادفی، مقصد حرکت مشخص است، اما مقدار پاداش از یک توزیع تصادفی (مثلاً یکنواخت یا نرمال) گرفته می‌شود.

در این مطالعه، محیط Gridworld شامل یک شبکه 5×5 است که عامل باید از خانه شروع به خانه هدف برسد. عامل در هر مرحله می‌تواند یکی از چهار جهت را انتخاب کند. هدف یافتن سیاست بهینه برای بیشینه کردن پاداش است. پس از اجرای الگوریتم $Q - Learning$ ، سیاست بهینه استخراج شده نشان‌دهنده بهترین جهت حرکت از هر خانه به سمت هدف است. به عنوان مثال، اگر در خانه‌ای فلش به سمت راست باشد (\rightarrow)، به معنی آن است که حرکت به سمت راست بیشترین ارزش را دارد. این سیاست به عامل کمک می‌کند تا با پیمودن مسیر بهینه، در کوتاه‌ترین زمان و بیشترین پاداش به هدف برسد. پس از آموزش، سیاست بهینه به صورت جهت حرکت از هر خانه به سوی هدف نمایش داده می‌شود. این نتایج نشان می‌دهند که فرآیندهای تصادفی چگونه می‌توانند مسیرهای بهینه را در محیط‌های نامعین بیاموزند.

برای هر یک از سه سناریو، الگوریتم $Q - learning$ به فرم زیر

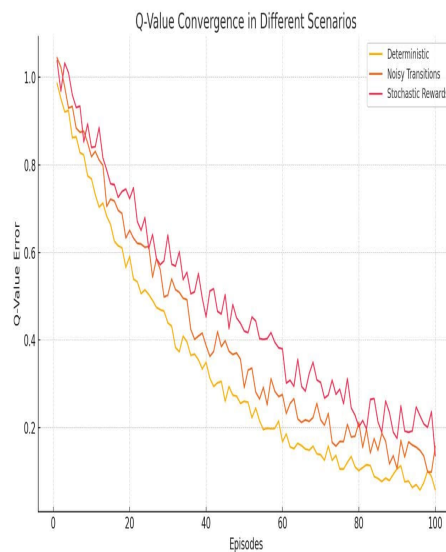
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

در صد اپیزود و با پارامترهای $\alpha = 0.1$ و $\gamma = 0.95$ اجرا شده و نمودار همگرایی و مسیر یادگیری استخراج شده است. در شکل ۲ همگرایی مقدار Q در سه سناریو نمایش داده شده است. همان‌طور که مشاهده می‌شود، مقدار خطای Q در محیط قطعی سریع‌تر کاهش می‌یابد. در محیط‌های نویزی و دارای پاداش تصادفی، نرخ همگرایی کندتر و همراه با نوسان است.

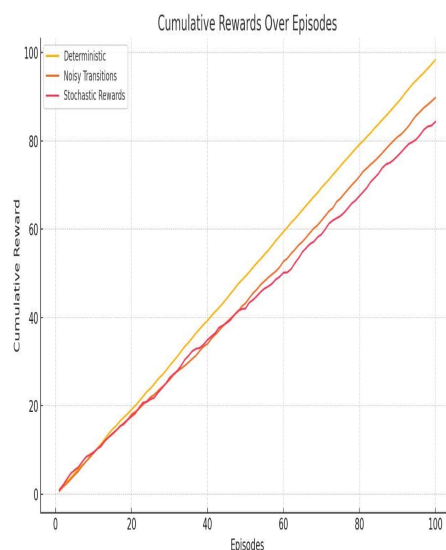
در شکل ۳ پاداش تجمعی در طول اپیزودها برای سه محیط نمایش داده شده است. در محیط قطعی، پاداش تجمعی با شیب مثبت مشخص رشد می‌کند، ولی در محیط‌های تصادفی، پاداش‌ها بی‌ثبات‌ترند.

در شکل ۴ مسیر عامل در شبکه نهایی نمایش داده شده است. عامل در طی فرآیند یادگیری، مسیری بهینه برای رسیدن به هدف پیدا کرده است.

تجزیه و تحلیل این نمودارها نشان می‌دهد که تصادفی بودن در پاداش یا انتقال، موجب کاهش سرعت یادگیری و همچنین



شکل ۲: نمودار همگرایی مقدار Q برای سه حالت

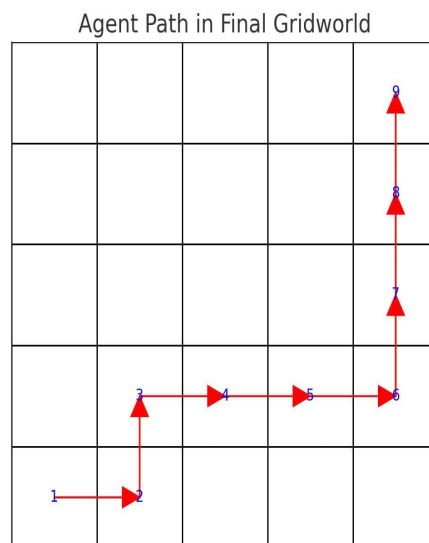


شکل ۳: نمودار پاداش تجمعی در طول اپیزودها برای سه حالت

افزایش نوسان در عملکرد عامل می‌شود. بنابراین، استفاده از تکنیک‌های مقاوم در برابر نویز و سیاست‌های اکتشافی مناسب، اهمیت زیادی دارند.

تحلیل نتایج بخش اصلی این مقاله است که نقش فرآیندهای تصادفی را آشکار می‌کند.

- محیط قطعی: در این حالت، عامل به سرعت سیاست بهینه را می‌یابد. همگرایی الگوریتم بسیار سریع است. مسیر عامل کاملاً مستقیم و بهینه است. پاداش تجمعی در اپیزودهای پایانی به حداکثر مقدار ممکن می‌رسد.
- محیط با نویز حرکتی: عامل مجبور است سیاست محافظه‌کارانه‌تری اتخاذ کند، زیرا هر حرکت ممکن است منجر به حالت ناخواسته شود. این موضوع باعث نوسانات بیشتر در مقادیر Q و تأخیر در همگرایی می‌شود. با این حال،



شکل ۴: مسیر طی شده عامل در شبکه Gridworld نهایی

الگوریتم همچنان قادر به یادگیری سیاست نسبتاً خوب است.

- محیط با پاداش‌های تصادفی: در این حالت عامل باید با تغییرات تصادفی پاداش‌ها سازگار شود. مقدار Q دچار نوسانات دائمی می‌شود. مسیر عامل نیز گاه به گاه تغییر می‌کند، بسته به اینکه کدام مسیر اخیراً پاداش بیشتری داشته است. همگرایی کندتر و نوسانی‌تر از سایر سناریوها است.

بحث و نتیجه‌گیری

مطالعه حاضر نشان داد که فرآیندهای تصادفی نقش کلیدی در تصمیم‌گیری هوشمند دارند. حتی در محیط‌های ساده‌ای مانند *Gridworld*، وجود نویز و تصادفی بودن پاداش می‌تواند یادگیری عامل را با چالش روبرو کند. Q -learning با وجود این چالش‌ها، قادر است در بلندمدت سیاست‌های کارآمدی یاد بگیرد. مطالعه موردی ارائه شده به خوبی نشان می‌دهد که چگونه مفاهیم ریاضی می‌توانند به راه‌حل‌های عملی در مسائل هوش مصنوعی منجر شوند. برای کارهای آینده، استفاده از الگوریتم‌های پیشرفته‌تر مانند *DQN* و بررسی اثر حافظه و شبکه‌های عصبی پیشنهاد می‌شود.

مراجع

Bertsekas, D. P. (2017). *Dynamic Programming and Optimal Control*. Athena Scientific.

Durrett, R. (2019). *Probability: Theory and Examples*. Cambridge University Press.

- Kearns, M., and Singh, S. (2002). Near-optimal reinforcement learning in polynomial time. *Machine Learning*, 49(2), 209-232.
- Ladosz, P. , Slowik, M., and Dębowski, M. (2022). Exploration in deep reinforcement learning: A survey, *Information Fusion*, vol. 89, pp. 1–23, 2022.
- Fourati, F., Aggarwal, V., and Alouini, M. S. (2024). Stochastic Q-learning for Large Discrete Action Spaces. In *Proceedings of the 41st ICML*, PMLR 235:13734–13759.
- Mnih, V., Kavukcuoglu, K. and Silver, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529–533.
- Puterman, M. L. (2005). *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley.
- Sutton, R. S., Barto, A. G. (2018). *Reinforcement Learning: An Introduction* (2nd ed.). MIT Press.
- Thrun, S, Burgard, W., and Fox, D. (1995). *Probabilistic robotics*. MIT press.
- Watkins, C. J. C. H., and Dayan, P. (1992). "Q-learning". *Machine Learning*, 8(3-4), 279–292.

The Role of Stochastic Processes in Artificial Intelligence and Data Science: A Case Study

Azam Rastin¹

¹Shahid Chamran university of Ahvaz

Abstract:

Stochastic processes are a foundational component in modeling uncertainty and dynamics within Artificial Intelligence (AI) and Data Science. This paper presents a theoretical and empirical investigation into the role of stochastic processes in Reinforcement Learning (RL), focusing on a case study using the classical Gridworld environment. We formalize the decision-making process as a Markov Decision Process (MDP), and examine how stochastic transition dynamics influence learning policies and value estimation. Experimental results demonstrate how varying the stochasticity of the environment impacts policy convergence, optimality, and stability. These findings provide critical insights into the importance of accurate modeling of environment dynamics for robust AI systems.

Keywords: Most important keywords of article (3-5 items) should appear here.

Mathematics Subject Classification (2020): 60J20, 68T05, 68T20.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30-31, 2025
University of Kurdistan



فرایند INAR(1)-PJ : یک جایگزین جدید فرایند INAR(1) پواسونی

مهدی راسخی^۱

گروه آمار، دانشکده علوم ریاضی و آمار، دانشگاه ملایر، ملایر، ایران

چکیده: در این مقاله، یک مدل خودبازگشتی مرتبه اول با مقادیر صحیح نامنفی معرفی می‌شود که نوفه فرایند از یک توزیع یک پارامتری جدید پیروی می‌کند. برای برآورد پارامترهای این فرایند، دو روش مختلف برآوردیابی در نظر گرفته شده است: حداکثر درستنمایی شرطی و روش یول-والکر. همچنین در یک مطالعه شبیه‌سازی عملکرد و کارایی این روش‌ها ارزیابی شده است. علاوه بر این، کاربرد فرایند پیشنهادی از طریق تحلیل یک مجموعه داده واقعی ارائه شده و در این قسمت برتری آن نسبت به مدل INAR(1) پواسونی نشان داده شده است.

واژه‌های کلیدی: توزیع احتمال جوجز، داده‌های شمارشی، بیش‌پراکنش، فرایند خودبازگشتی عدد صحیح مرتبه اول.

کد موضوع بندی ریاضی (۲۰۲۰): 62E15, 62F10, 62M10.

۱ مقدمه

داده‌های سری‌های زمانی با مقادیر اعداد صحیح نامنفی در حوزه‌های گوناگونی همچون علوم اجتماعی، اپیدمیولوژی، بیمه، مهندسی صنایع و بسیاری دیگر از زمینه‌ها به طور گسترده‌ای دیده می‌شوند. برای مدل‌سازی بهینه این نوع داده‌ها که دارای وابستگی زمانی نیز هستند، مدل خودبازگشتی عدد صحیح مرتبه اول (INAR(1)) پیشنهاد شده است که مکنزی (۱۹۸۶) و الاش و الزید (۱۹۸۷) طور مستقل به معرفی آن پرداختند. ویژگی مشترکی که اغلب در این نوع سری‌های زمانی دیده می‌شود، بیش‌پراکنش است؛ یعنی واریانس داده‌ها بیشتر از میانگین آنهاست. با توجه به اینکه توزیع پواسون اغلب برای مدل‌سازی داده‌های شمارشی با بیش‌پراکنش مناسب نیست، برخی از پژوهشگران مدل‌های INAR(1) مبتنی بر توزیع‌های نوفه انعطاف‌پذیرتری را پیشنهاد کرده‌اند. از جمله این مدل‌ها می‌توان به موارد زیر اشاره کرد: مدل INAR(1) گسسته هندسی توسط همکاران و جازی (۲۰۱۲)، فرایند INAR(1) مبتنی بر خانواده توزیع‌های کاتس که توسط کیم و لی (۲۰۱۷) معرفی شد. در راستای پیشرفت‌های اخیر در مدل‌سازی داده‌های شمارشی، در این مقاله توزیع

^۱ سخنران، m.rasekhi@malayeru.ac.ir

پواسون-جوچز (PJ) معرفی شده است: یک توزیع گسسته جدید با یک پارامتر که از ترکیب توزیع‌های پواسون و جوچز (اکسپری و امبگو ۲۰۲۲) حاصل شده است. توزیع استفاده شده در ترکیب فوق (توزیع جوچز) در مقایسه با توزیع گاما دارای تعداد پارامتر کمتری

است اما از لحاظ انعطاف پذیری در شکل تابع چگالی و میزان شاخص‌های چولگی و کشیدگی همانند توزیع گاما عمل می‌کند. توزیع PJ به ویژه در مدل‌بندی داده‌ها با ویژگی بیش‌پراکنش، کارا است و دارای فرم‌های بسته و ساده‌ای برای تابع جرم احتمال و گشتاورهای آن می‌باشد. با استفاده از این توزیع به عنوان مکانیزم نوفه، یک فرآیند خودبازگشتی عدد صحیح مرتبه اول را توسعه داده و بصورت INAR(1)-PJ نامگذاری شده است. فرآیند پیشنهادی روی یک مجموعه داده واقعی اعمال شده است تا کاربرپذیری و انعطاف‌پذیری آن نشان داده شود. همچنین، مطالعات شبیه‌سازی برای ارزیابی دقت و کارایی روش‌های برآوردیابی پارامترها در چارچوب پیشنهادی انجام شده است. مقاله به شرح زیر سازمان‌دهی گردیده است: در بخش ۲، توزیع PJ معرفی و ویژگی‌های آماری آن به طور کامل بررسی شده است. در بخش ۳، یک فرایند خودبازگشتی عدد صحیح مبتنی بر نوفه‌های PJ پیشنهاد شده و خواص ساختاری مهم آن ارائه شده است. در بخش ۴، دو روش برآوردیابی برای پارامترهای مدل INAR(1)-PJ مورد بحث قرار می‌گیرد: روش حداکثر درستنمایی شرطی و روش یول-والکر. در بخش ۵، نتایج یک مطالعه شبیه‌سازی به منظور ارزیابی دقت و کارایی این روش‌های برآوردیابی گزارش شده‌اند. در بخش ۶، فرآیند پیشنهادی روی یک مجموعه داده واقعی اعمال شده است تا کاربرد عملی آنها نشان داده شود.

۲ توزیع پواسن-جوچز

تابع جرم احتمال (PMF) توزیع جوچز به صورت زیر تعریف می‌شود:

$$f(x; \theta) = \frac{\theta^x}{\theta^3 + \theta^2 + 6} (1 + x + x^2) e^{-\theta x} \quad x > 0, \quad (1.2)$$

که $\theta > 0$. تابع توزیع تجمعی (CDF) متناظر با توزیع جوچز به صورت زیر بیان می‌شود:

$$F(x; \theta) = 1 - \left(1 + \frac{\theta x(\theta^2 + \theta^2 x^2 + 3\theta x + 6)}{\theta^3 + \theta^2 + 6} \right) e^{-\theta x} \quad x \geq 0.$$

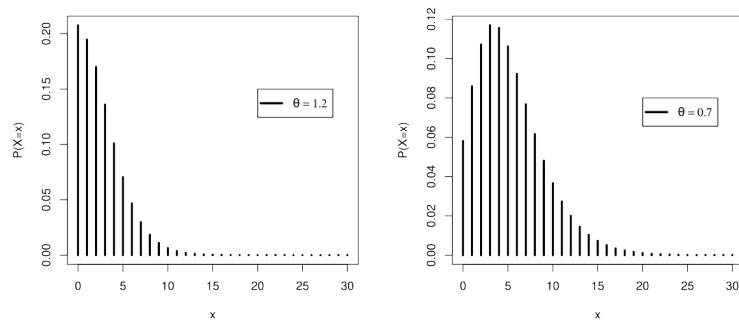
در ادامه، یک توزیع پواسون ترکیبی جدید یک پارامتری را تعریف می‌کنیم که از ترکیب توزیع پواسون با توزیع جوچز حاصل می‌شود. فرض کنید متغیر تصادفی λ دارای توزیع جوچز با پارامتر θ باشد که در برابری (۱.۲) تعریف شده است. همچنین فرض کنید که توزیع شرطی X به شرط λ ، یک توزیع پواسون با میانگین λ باشد. در این صورت، توزیع حاشیه‌ای X که از حذف λ به دست می‌آید، مربوط به توزیع پواسون-جوچز (PJ) است و دارای تابع جرم احتمال (PMF) به صورت زیر می‌باشد:

(۲.۲)

$$p(x; \theta) = \left(\frac{\theta^x}{\theta^3 + \theta^2 + 6} \right) \left(\frac{1}{(\theta + 1)^{x+1}} + \frac{(x+1)}{(\theta + 1)^{x+2}} + \frac{(x+3)(x+2)(x+1)}{(\theta + 1)^{x+4}} \right), \quad x = 0, 1, 2, \dots$$

فرض کنید X نشان‌دهنده متغیر تصادفی با تابع چگالی (۲.۲) باشد؛ در این صورت می‌نویسیم $X \sim PJ(\theta)$. شکل ۱ انواع مختلفی از شکل‌های تابع جرم احتمال (PMF) توزیع PJ را تحت مقادیر مختلف پارامتر نشان می‌دهد. تابع جرم احتمال این توزیع به ازای مقادیر متفاوت پارامتر، دارای ساختار نزولی و چوله به راست است. دو گشتاور اول توزیع PJ به ترتیب به فرم زیر هستند:

$$\begin{aligned} E(X) &= \left(\frac{\theta^x}{\theta^3 + \theta^2 + 6} \right) \left\{ \sum_{x=0}^{\infty} \frac{x}{(\theta+1)^{x+1}} + \sum_{x=0}^{\infty} \frac{x(x+1)}{(\theta+1)^{x+2}} + \sum_{x=0}^{\infty} \frac{x(x+3)(x+2)(x+1)}{(\theta+1)^{x+4}} \right\} \\ &= \left(\frac{\theta^x}{\theta^3 + \theta^2 + 6} \right) \left\{ \frac{1}{\theta^2} + \frac{2}{\theta^3} + \frac{24}{\theta^5} \right\} = \frac{\theta^2 + 2\theta^3 + 24}{\theta(\theta^3 + \theta^2 + 6)}, \end{aligned} \quad (3.2)$$



شکل ۱: نمودار تابع جرم احتمال توزیع PJ تحت مقادیر مختلف پارامتر

$$E(X^2) = \left(\frac{\theta^4}{\theta^3 + \theta^2 + 6} \right) \left\{ \sum_{x=0}^{\infty} \frac{x^2}{(\theta+1)^{x+1}} + \sum_{x=0}^{\infty} \frac{x^2(x+1)}{(\theta+1)^{x+2}} + \sum_{x=0}^{\infty} \frac{x^2(x+2)(x+1)}{(\theta+1)^{x+3}} \right\} \quad (4.2)$$

$$= \left(\frac{\theta^4}{\theta^3 + \theta^2 + 6} \right) \left\{ \frac{\theta+2}{\theta^3} + \frac{2(\theta+3)}{\theta^4} + \frac{24(\theta+5)}{\theta^5} \right\} = \frac{\theta^5 + 4\theta^6 + 6\theta^5 + 32\theta^4 + 120\theta^3 + 60\theta^2 + 144\theta + 144}{\theta^2(\theta^3 + \theta^2 + 6)^2}.$$

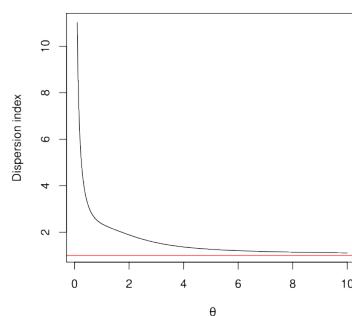
بر اساس روابط (۳.۲) و (۴.۲) واریانس توزیع PJ عبارتست از

$$V(X) = \frac{\theta^5 + 4\theta^6 + 6\theta^5 + 32\theta^4 + 120\theta^3 + 60\theta^2 + 144\theta + 144}{\theta^2(\theta^3 + \theta^2 + 6)^2}.$$

همچنین شاخص پراکنش به فرم زیر حاصل می شود.

$$DI(X) = \frac{\theta^5 + 4\theta^6 + 6\theta^5 + 32\theta^4 + 120\theta^3 + 60\theta^2 + 144\theta + 144}{\theta(\theta^3 + \theta^2 + 6)(\theta^3 + 2\theta^2 + 24)}.$$

شکل ۲ نشان می دهد که برای تمامی مقادیر پارامتر θ ، شاخص پراکنشی بزرگتر از ۱ است. مقدار این شاخص به ازای مقادیر بزرگ پارامتر θ به عدد ۱ نزدیک می شود. این رفتار، کارایی این توزیع در مدل سازی داده های شمارشی با بیش پراکنش را تأیید می کند.



شکل ۲: نمودار شاخص پراکنش توزیع PJ تحت مقادیر مختلف پارامتر

الگوریتم زیر می تواند برای تولید نمونه تصادفی از توزیع PJ با بهره بردن از روش معکوس تابع توزیع و بر اساس توزیع جوجز استفاده

شود: (۱) تولید U_i ، $i = 1, 2, \dots, n$ ، جاییکه $U_i \sim U(0, 1)$.

(۲) یافتن ریشه برابری زیر بر اساس λ_i به کمک دستور uniroot در نرم افزار R.

$$\ln \left(1 + \frac{\theta \lambda_i (\theta^2 + \theta^2 \lambda_i^2 + 3\theta \lambda_i + 6)}{\theta^3 + \theta^2 + 6} \right) - \theta \lambda_i - \ln(1 - U_i) = 0$$

(۳) تولید X_i از توزیع $Poisson(\lambda_i)$.

۳ فرایند INAR(1) مبتنی بر نوفه‌های PJ

فرض کنید $\{X_t; t \in \mathbb{Z}\}$ نشان دهنده فرایند INAR(1) بر اساس تعریف زیر باشد.

$$X_t = p \circ X_{t-1} + \varepsilon_t, \quad t \in \mathbb{Z}, \quad (1.3)$$

که در آن $0 \leq \alpha < 1$ و $\{\varepsilon_t; t \in \mathbb{Z}\}$ دنباله ای از متغیرهای تصادفی هم توزیع و مستقل با مقادیر صحیح نامنفی و میانگین $E(\varepsilon_t) = \mu_\varepsilon$ و واریانس $V(\varepsilon_t) = \sigma_\varepsilon^2$ است. فرض کنید نوفه $\{\varepsilon_t; t \in \mathbb{Z}\}$ مستقل از مشاهدات گذشته $k \geq 1$ است. X_{t-k} عملگر \circ نشان دهنده عملگر رقیق‌ساز دوجمله ای است که به فرم زیر تعریف می‌شود:

$$p \circ X_{t-1} := \sum_{j=1}^{X_{t-1}} W_j$$

که در آن $\{W_j; j \geq 1\}$ نشان دهنده دنباله ای از متغیرهای تصادفی هم توزیع و مستقل با توزیع برنولی با احتمال موفقیت p است. اگر $0 \leq p < 1$ این فرایند ایستا و در غیر این صورت فرایند نایستا است. برای تشخیص ساختار حرکت در فرایند، احتمال انتقال یک مرحله ای به شرح زیر است:

$$P(X_t = k | X_{t-1} = l) = \sum_{i=1}^{\min(k,l)} \binom{l}{i} p^i (1-p)^{l-i} P(\varepsilon_t = k-i), \quad k, l \geq 0$$

که در آن $0 < p < 1$. برای ایجاد یک فرایند تصادفی INAR(1) جدید، دنباله متغیرهای تصادفی $\{\varepsilon_t; t \in \mathbb{Z}\}$ که از توزیع PJ پیروی می‌کند را به عنوان متغیرهای نوفه فرایند در برابری (۱.۳) لحاظ می‌شود. در نتیجه احتمال انتقال یک مرحله فرایند تصادفی جدید INAR(1)-PJ به فرم زیر حاصل می‌شود.

$$P(X_t = k | X_{t-1} = l) = \left(\frac{\theta^*}{\theta^3 + \theta^2 + 6} \right) \sum_{i=1}^{\min(k,l)} \binom{l}{i} p^i (1-p)^{l-i} \left(\frac{1}{(\theta+1)^{k-i+1}} + \frac{(k-i+1)}{(\theta+1)^{k-i+2}} + \frac{(k-i+3)(k-i+2)(k-i+1)}{(\theta+1)^{k-i+3}} \right). \quad (2.3)$$

از این پس، فرایند توصیف شده در معادله (۲.۳) فرایند مربوطه را با INAR(1)-PJ نامگذاری می‌کنیم. بر اساس روش ویب (۲۰۱۸) میانگین، واریانس و شاخص پراکنش فرایند INAR(1)-PJ به ترتیب به صورت زیر بیان می‌شوند:

$$\mu_X = \frac{\theta^3 + 2\theta^2 + 24}{\theta(\theta^3 + \theta^2 + 6)(1-p)},$$

$$\sigma_X^2 = \frac{\left((p+1)\theta^7 + (3p+4)\theta^6 + (2p+6)\theta^5 + (30p+32)\theta^4 + (36p+120)\theta^3 + 60\theta^2 + (144p+144)\theta + 144 \right)}{(1-p)^2 \theta^2 (\theta^3 + \theta^2 + 6)^2},$$

$$DI_X = \frac{\frac{\theta^7 + 4\theta^6 + 6\theta^5 + 32\theta^4 + 120\theta^3 + 60\theta^2 + 144\theta + 144}{\theta(\theta^3 + \theta^2 + 6)} + p}{1 + p}.$$

تحت فرضیات مدل، بر اساس روش ارائه شده در [الاش و الزید \(۱۹۸۸\)](#)، میانگین و واریانس شرطی فرآیند به ترتیب به شرح زیر حاصل می‌شود:

$$E(X_t | X_{t-1}) = pX_{t-1} + \frac{\theta^3 + 2\theta^2 + 24}{\theta(\theta^3 + \theta^2 + 6)},$$

و

$$V(X_t | X_{t-1}) = p(1-p)X_{t-1} + \left(\frac{\left(\begin{array}{c} \theta^7 + 4\theta^6 + 6\theta^5 + 32\theta^4 \\ + 120\theta^3 + 60\theta^2 + 144\theta + 144 \end{array} \right)}{\theta^2(\theta^3 + \theta^2 + 6)^2} \right).$$

۴ برآوردیابی پارامترهای فرآیند

۱.۴ روش ماکسیمم درستنمایی شرطی

برای انجام روش ماکسیمم درستنمایی شرطی تحت فرآیند INAR(1)-PJ، تابع درستنمایی شرطی زیر تعریف شده است.

$$\begin{aligned} \ell(p, \theta) &= \sum_{t=2}^T \ln(P(X_t = k | X_{t-1} = l)) \\ &= \sum_{t=2}^T \ln \left(\left(\frac{\theta^k}{\theta^k + \theta^l + 6} \right)^{\min(X_t, X_{t-1})} \sum_{i=1}^{X_{t-1}} \binom{X_{t-1}}{i} p^i (1-p)^{X_{t-1}-i} \left(\frac{1}{(\theta+1)^{X_t-i+1}} \right. \right. \\ &\quad \left. \left. + \frac{(X_t-i+1)}{(\theta+1)^{X_t-i+2}} + \frac{(X_t-i+2)(X_t-i+1)}{(\theta+1)^{X_t-i+3}} \right) \right). \end{aligned} \quad (1.4)$$

برآوردگر ماکسیمم درستنمایی شرطی (CML) فرآیند INAR(1)-PJ، که با این $\hat{p}_{CLS}, \hat{\theta}_{CLS}$ نشان داده می‌شود، با ماکسیمم کردن تابع درستنمایی شرطی برابری (۱.۴) حاصل می‌شود. برآوردیابی پارامترها را به کمک تابع optim در نرم‌افزار R انجام می‌دهیم. برای مدیریت دامنه پارامترهای p و θ از الگوریتم L-BFGS-B استفاده می‌شود. ماتریس اطلاع مشاهده شده در پارامترهای برآورد شده را با استفاده از خروجی ماتریس هسیان حاصل از بهینه‌سازی محاسبه می‌شود. خروجی حاصل امکان محاسبه خطاهای استاندارد تقریبی برای برآوردگرها را فراهم می‌نماید. تحت شرایط نظم، برآوردگرهای CML دارای سازگاری، توزیع نرمال مجانبی و کارایی هستند. [جو \(۱۹۹۷\)](#)

۲.۴ روش یول-والکر

روش برآوردیابی یول-والکر براساس برابری توابع خودکواریانس و میانگین فرآیند با حالت نمونه‌ای آنها و برآورد پارامترها با حل دستگاه معادلات ایجاد شده حاصل می‌شود. برای فرآیند INAR(1) تابع خود همبستگی نمونه ای (ACF) عبارتست از $\rho_X(h) = p^h$ که

$h \geq 1$ و $p \in (0, 1)$ عملگر نازل فرآیند است. بنابراین برآورد یول-والکر p عبارتست از

$$\hat{p}_{YW} = \frac{\sum_{t=2}^T (X_t - \bar{X})(X_{t-1} - \bar{X})}{\sum_{t=2}^T (X_t - \bar{X})^2}.$$

برآوردگر یول-والکر $\hat{\theta}_{YW}$ از برابر قرار دادن میانگین نظری فرایند PJ-INAR(1) با مقدار نمونه‌ای آن به دست می‌آید. معادله برآوردیابی با قرار دادن تفاضل این دو کمیت برابر با صفر تشکیل شده و سپس $\hat{\theta}_{YW}$ از حل این معادله نسبت به θ حاصل می‌شود.

$$((1 - \hat{p}_{YW})\bar{X})\theta^4 + ((1 - \hat{p}_{YW})\bar{X} - 1)\theta^3 - 2\theta^2 + (6(1 - \hat{p}_{YW})\bar{X})\theta - 24 = 0 \quad (2.4)$$

برابری (2.4) به صورت عددی به کمک تابع uniroot در نرم افزار R بدست می‌آید.

۵ مطالعات شبیه سازی

یک مطالعه شبیه سازی به منظور ارزیابی رفتار برآوردگرهای CML و YW بر مبنای حجم متفاوت نمونه ای برای پارامترهای مدل INAR(1)-PJ انجام شده و این شبیه سازی بر اساس ۱۰۰۰ تکرار برای حجم های نمونه ۱۰۰، ۳۰۰، ۵۰۰ اجرا شده است. نمونه های تصادفی با توزیع PJ با استفاده از الگوریتم معرفی شده در بخش ۱ تولید شده اند. یک مجموعه انتخابی از مقادیر پارامتر برای انجام شبیه سازی عبارت است از $\theta = 1$ ، $p = 0.5$. عملکرد برآوردگرهای CML و YW با استفاده از سه معیار انحراف (Bias)، میانگین مربعات خطا (MSE) و میانگین خطای نسبی (MRE) ارزیابی می شود. این معیارها به صورت زیر تعریف می شوند:

$$Bias = \sum_{j=1}^N \frac{\hat{\gamma}_{i,j} - \gamma_i}{N}, \quad MSE = \sum_{j=1}^N \frac{(\hat{\gamma}_{i,j} - \gamma_i)^2}{N}, \quad MRE = \sum_{j=1}^N \frac{\hat{\gamma}_{i,j}/\gamma_i}{N},$$

برای $i = 1, 2$ و $\gamma = (p, \theta)$.

مقادیر محاسبه شده برای Bias و MSE که در جدول ۱ گزارش شده اند، به طور کلی نزدیک به صفر هستند و این موضوع میزان اطمینان هر دو روش برآوردیابی را در حجم های محدود تأیید می کند. با افزایش حجم نمونه، میانگین خطاهای نسبی (MREs) به سمت یک همگرا شده، که نشان از بهبود دقت برآوردیابی با افزایش حجم داده ها دارد. این نتایج نشان می دهند که روش CML از لحاظ هر دو معیار Bias و MSE در تمامی حجم های نمونه ای، عملکرد بهتری نسبت به روش YW دارد. علاوه بر این، مقادیر MRE محاسبه شده نیز برتری روش CML را تأیید می کنند، زیرا این مقادیر به طور مداوم به مقدار ایده آل یک نزدیک تر هستند.

۶ کاربرد فرآیند پیشنهادی

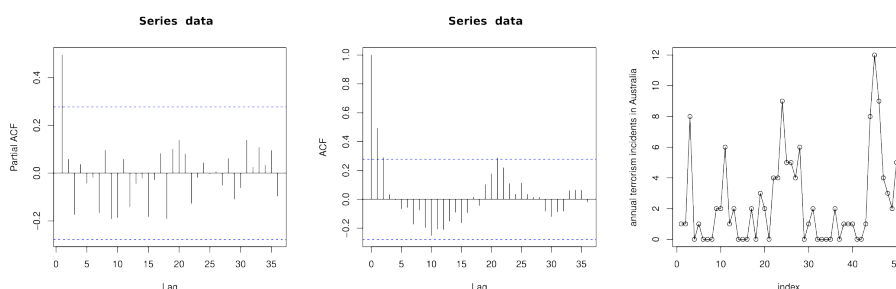
در این بخش، یک مجموعه داده واقعی برای بررسی کاربرد عملی فرایند پیشنهادی INAR(1)-PJ در مقایسه با مدل استاندارد INAR(1) پواسونی (INAR(1)-P) مورد تحلیل قرار می گیرد. برای مقایسه مدل های رقیب، از معیار AIC و معیار BIC استفاده می کنیم، به طوری که مقادیر کمتر نشان دهنده برازش بهتر مدل به داده ها است. ما یک سری زمانی سالانه از وقایع تروریستی در استرالیا که از سال ۱۹۷۰ تا ۲۰۲۰ ادامه داشته است، مورد تحلیل قرار می دهیم. داده ها از بسته Ecdat در نرم افزار R به دست آمده اند. آمار توصیفی در این مجموعه داده عبارتست از میانگین ۴/۲، واریانس ۴۸/۸ و شاخص پراکنش ۵۳/۳، که به خوبی رفتار بیش پراکنش را در داده های شمارشی نشان

جدول ۱: نتایج شبیه سازی برای فرآیند INAR(1)-PJ: اولین مجموعه انتخابی پارامترها

$\theta = 1$		$p = 0.5$			n
YW	CML	YW	CML		
-0.246/0	0.058/0	-0.316/0	-0.032/0	Bias	100
0.188/0	0.093/0	0.085/0	0.023/0	MSE	
9491/0	9790/0	9112/0	9662/0	MRE	
-0.069/0	0.024/0	-0.103/0	-0.015/0	Bias	300
0.063/0	0.035/0	0.026/0	0.009/0	MSE	
9827/0	9910/0	9687/0	9852/0	MRE	
-0.027/0	0.019/0	-0.049/0	-0.002/0	Bias	500
0.040/0	0.019/0	0.015/0	0.005/0	MSE	
9862/0	9909/0	9790/0	9884/0	MRE	

می‌دهند و لزوم استفاده از مدل‌هایی قوی تر از مدل استاندارد INAR(1) پواسونی را مطرح می‌کند. روش ارائه شده در **شویر و ویب (۱۹۸۶)** یک آزمون برای وجود بیش‌پراکنش در فرایندهای خودبازگشتی با مقادیر عدد صحیح نامنفی ارائه دادند. با اعمال این روش روی مجموعه داده حاضر، آماره آزمون برابر با $37/11$ و مقدار p آزمون کمتر از 0.01 به دست می‌آید که منجر به رد فرضیه صفر هم‌پراکنش می‌شود. این نتیجه برجسته می‌کند که توزیع نوفه مدل باید قادر به مدل‌سازی تغییرات بیشتر از میانگین باشد؛ بنابراین، ممکن است مدل مبتنی بر پواسون برای این داده‌ها مناسب نباشند.

شکل ۳ نمودار سری زمانی تعداد سالانه حملات تروریستی در استرالیا را به همراه توابع خودهمبستگی (ACF) و خودهمبستگی جزئی (PACF) برآورد شده را نشان می‌دهد. تابع PACF پس از گام اول قطع می‌شود، که نشان می‌دهد ساختار خودبازگشتی مرتبه اول به خوبی وابستگی دنباله‌ای داده‌ها را مدل‌سازی می‌کند. بنابراین، مدلی از نوع $AR(1)$ برای مدل‌سازی این مجموعه داده مناسب است.



شکل ۳: نمودار سری زمانی، تابع خودهمبستگی (ACF) و تابع خودهمبستگی جزئی (PACF) برای تعداد سالانه حملات تروریستی در استرالیا

در جدول ۲، پارامترهای برآورد شده مدل‌ها، خطاهای استاندارد مرتبط به آنها و معیارهای انتخاب مدل ارائه شده است. مدل INAR(1)-PJ از نظر هر دو معیار AIC و BIC مقادیر پایین‌تری را نسبت به مدل INAR(1)-P ارائه می‌دهد. این نتایج حاکی از عملکرد بهتر مدل پیشنهادی در توصیف داده‌های مربوط به حملات تروریستی سالانه در استرالیا دارد.

جدول ۲: برآورد ماکسیمم درستنمایی شرطی پارامترهای مدل، خطای استاندارد، AIC، BIC و آماره های مرتبط با تعداد سالانه حملات

تروریستی در استرالیا

مدل	پارامتر	برآورد	خطای استاندارد	AIC	BIC	μ_X	σ_X^2	DI_X
(INAR-PJ\1)	p	۲۳۷/۰	۰۱۱۰/۰	۸۴۹/۲۰۷	۶۳۳/۲۱۱	۳۶/۲	۴۵/۴	۸۸/۱
	θ	۵۱۴/۱	۰۳۱/۰					
(INAR-P\1)	p	۳۳۱/۰	۰۰۵/۰	۳۳۹/۲۳۹	۱۲۳/۲۴۳	۴۶/۲	۴۶/۲	۱
	λ	۶۵۰/۱	۰۵۱/۰					
تجربی								
						۴/۲	۴۸/۸	۵۳/۳

۷ نتیجه گیری

در این مقاله، یک مدل جدید خودبازگشتی عدد صحیح نامنفی مرتبه اول با نوفه های PJ با نام (INAR(1)-PJ)، برای مدل سازی سری های زمانی شمارشی با بیش پراکنش پیشنهاد شده است. به منظور برآورد پارامترهای مدل، دو روش مختلف مورد بررسی قرار گرفته و عملکرد آنها از طریق شبیه سازی مقایسه شده است. مدل پیشنهادی روی یک مجموعه داده واقعی اعمال شده و نتایج آن با مدل استاندارد INAR(1) بواسونی مقایسه گردیده است. یافته های تجربی نشان می دهد که مدل INAR(1)-PJ در مواجهه با بیش پراکنش عملکرد بهتری دارد. بنابراین مدل پیشنهادی برای مدلبندی داده هایی با چولگی مثبت و شاخص پراکنش بزرگتر از ۱ گزینه مناسبتری محسوب می شود.

مراجع

- Al-Osh, M and Alzaid, A (1987). First-order integer-valued autoregressive (INAR (1)) process. *J Time Ser Anal*, **8(3)**: 261-275.
- Alzaid, A, Al-Osh, M (1988). First-order integer-valued autoregressive (INAR (1)) process: distributional and regression properties. *Stat Neerland*, **42(1)**: 53-61.
- Echebiri, U.V., Mbegbu, J.I. (2022). Juchez Probability Distribution: Properties and Applications. *Asian Journal of Probability and Statistics*, **20(2)**: 56-71.
- Jazi, M., Jones, G. and Lai, C.D. (2012). Integer valued AR (1) with geometric innovations. *J. Iran. Stat. Soc. (JIRSS)*, **11**: 173-190.
- Joe, H (1997). Multivariate models and multivariate dependence concepts. *Chapman and Hall, London*.
- Kim, H. and Lee, S. (2017). On first-order integer-valued autoregressive process with Katz family innovations. *J Stat Comput Simul*, **87**: 546-562.

- McKenzie, E (1986). Autoregressive moving-average processes with negative binomial and geometric marginal distributions. *Adv Appl Probabx*, **18**: 679-705.
- Schweer, S, WeiB, CH (2014). Compound Poisson INAR (1) processes: stochastic properties and testing for overdispersion. *Comput Stat Data Anal*, **77**: 267-284.
- WeiB, C.H.(2018). An introduction to discrete-valued time series. *Wiley*.

The INAR(1)-PJ Process: A New Alternative to the Poisson INAR(1) Process

Mahdi Rasekhi¹

¹Department of Statistics, Faculty of Mathematical Sciences and Statistics,
Malayer University, Malayer, Iran

Abstract: In this paper, we introduce a first-order non-negative integer-valued autoregressive model, in which the innovation process follows a newly proposed one-parameter distribution. To estimate the unknown parameters of the model, two estimation approaches are considered: conditional maximum likelihood and Yule-Walker methods. A simulation study is conducted to evaluate the performance and efficiency of these estimation techniques. Furthermore, the practical applicability of the proposed model is demonstrated through the analysis of a real-world data sets, showing that it outperforms the traditional Poisson INAR(1) model.

Keywords: Juchez probability distribution, count data, over-dispersion, INAR(1) process.

Mathematics Subject Classification (2020): 62M10, 62F10, 62E15.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



بررسی انشعاب‌های تصادفی مدل رشد لجستیک جمعیت

سمیرا مهموئی^۱، امید ربیعی مطلق^۲، حاجی محمد محمدی نژاد^۳

^{۱،۲،۳} دانشکده علوم ریاضی و آمار دانشگاه بیرجند

چکیده: تحلیل رفتار جمعیت و فرآیندهای زیستی از مسائل کلیدی در علوم ریاضی و زیست‌شناسی به‌شمار می‌آید. در این زمینه، مدل رشد لجستیک به‌عنوان ابزاری مؤثر برای بررسی رفتار جمعیت‌ها و فرآیندهای بیولوژیکی مورد استفاده قرار می‌گیرد. با توجه به افزایش جمعیت و چالش‌های محیطی، مدیریت پایدار منابع امری ضروری است. این مقاله به بررسی انشعاب‌های معادله رشد لجستیک تحت تأثیر حرکت وینر و عوامل تصادفی می‌پردازد. نتایج به‌دست آمده نشان‌دهنده تغییرات ناگهانی در امید ریاضی و وقوع P -انشعاب در رفتار سیستم است. هدف ارائه چارچوبی جامع برای درک بهتر دینامیک‌های سیستم‌های جمعیتی در شرایط ناپایدار است.

واژه‌های کلیدی: مدل رشد لجستیک جمعیت، انشعاب تصادفی، معادلات تصادفی، تابع چگالی احتمال.

کد موضوع بندی ریاضی (۲۰۲۰): 37N25, 34F05, 37G10.

۱ مقدمه

مدل رشد لجستیک^۱ در ساختارهای ریاضی متنوعی مانند مدل‌سازی رشد جمعیت باکتری‌ها، پیاده‌سازی شبکه‌های عصبی مصنوعی، تخمین پویایی ۱۹ - COVID برای پیشگیری از شیوع، تحلیل مدل‌های رشد جمعیت تومورها به منظور درک رفتار فرآیند آن‌ها و ... کاربرد دارد (Tsagkis, 2023; Postnikov, 2020; Mansour, 2022; Koyama, 2021). با توجه به روند افزایشی جمعیت جهان، نیاز برای مدیریت رشد پایدار و اطمینان از اینکه رشد جمعیت در محدوده منابع موجود باقی می‌ماند، ضروری است. همانطور که می‌دانیم، چالش‌هایی مانند تخریب محیط زیست با افزایش جمعیت تشدید می‌شود. اندازه جمعیت به‌طور قابل توجهی بر توسعه اجتماعی، اقتصادی، سیاسی و محیط زیستی تأثیر می‌گذارد.

^۱ سخنران، s.mahmoiy@birjand.ac.ir

^۱ Logistic Growth Model

با توجه به تأثیر عوامل متعددی بر نرخ رشد جمعیت، ورهولست^۲ این عوامل را در قالب مدل لجستیک فرمول‌بندی کرده است که نشان‌دهنده رشد محدود شونده‌ی یک جمعیت زیستی است. این معادله غالباً به عنوان نمونه استاندارد در نمایش رشد جمعیت مورد استفاده قرار می‌گیرد، که در آن نرخ تولد مستقیماً با اندازه‌ی کنونی جمعیت و منابع موجود مرتبط است. در ادامه، از توسیع مدل ورهولست به شکل زیر استفاده می‌شود:

$$\begin{aligned} dx &= \left[x \left(1 - \frac{x}{y} \right) - ax^2 \right] dt, \\ dy &= \left[-h \frac{x(y-x)}{y(c+x)^2} - bxy + \alpha y \right] dt, \end{aligned} \quad (1.1)$$

$$x(\circ) = x_0, \quad y(\circ) = y_0,$$

در این مدل، $x(t)$ نشان‌دهنده‌ی جمعیتی است در لحظه t که محیط زیست خودش را تخریب می‌کند و در عین حال رقابت درون‌گروهی هم دارد، $y(t)$ نقش ظرفیت برد سیستم را ایفا می‌کند که نشان‌دهنده‌ی حداکثر جمعیتی است که محیط قادر است در زمان t پشتیبانی کند. برخلاف مدل کلاسیک ورهولست که ظرفیت برد مقداری ثابت است، در اینجا $y(t)$ به صورت تابعی پویا مدل شده که تحت تأثیر چندین عامل زیست‌محیطی تغییر می‌کند. عبارت $-bxy$ در معادله‌ی dy نشان می‌دهد که افزایش جمعیت x منجر به تخریب سریع‌تر محیط از طریق استفاده بی‌رویه از منابع یا آلودگی آن می‌شود؛ بنابراین، افزایش پارامتر b باعث کاهش سریع‌تر $y(t)$ می‌شود. در مقابل، عبارت $+\alpha y$ بیانگر بازسازی طبیعی ظرفیت محیط است که به صورت تدریجی موجب بهبود شرایط زیستی می‌شود. افزایش α رشد سریع‌تر $y(t)$ را در پی دارد و بنابراین به پایداری بیشتر ظرفیت برد منجر می‌شود. افزون‌بر این، عبارت غیرخطی $-\frac{hx(y-x)}{y(c+x)^2}$ نوعی اصطکاک پویا میان جمعیت و ظرفیت محیط را مدل‌سازی می‌کند. پارامترهای a, b, c, h, α همگی مثبت بوده و دینامیک کلی سیستم زیستی را تعیین می‌کنند.

در سیستم‌های دینامیکی غیرتصادفی، انشعاب زمانی رخ می‌دهد که تغییر در پارامترهای سیستم باعث تغییر ناگهانی در رفتار کیفی جواب‌های آن شود. برای مثال، از دست رفتن پایداری یک نقطه تعادل و یا ظهور چرخه‌های حدی یا تغییر در تعداد نقاط تعادل نمونه‌های از وقوع یک انشعاب در دستگاه‌های معادلات دیفرانسیل غیرتصادفی هستند. در محیط‌های تصادفی، این انشعابات تحت تأثیر نویز قرار می‌گیرند، که می‌تواند رفتار انشعاب را تغییر دهد یا پدیده‌های جدیدی ایجاد کند که در سیستم‌های قطعی مشاهده نمی‌شوند. نظریه انشعاب تصادفی^۳ تغییرات کیفی را در رفتار مجانبی سیستم‌های دینامیکی تصادفی هنگامی که پارامترهای سیستم تغییر می‌کنند، بررسی می‌کند. انشعاب‌های تصادفی در دو قالب کلی دسته‌بندی می‌شوند:

۱- (D -انشعاب^۴ یا انشعاب دینامیکی^۵): این نوع انشعاب وابسته به تغییر رفتار کیفی جواب‌ها است. در صورت نبودن فاکتور تصادف در معادله، این انشعاب همان انشعاب معمولی در معادلات دیفرانسیل است.

۲- (P -انشعاب^۶ یا انشعاب پدیده‌ای^۷): این انشعاب وابسته به تغییر ماهیت تصادفی جواب‌ها است (به عنوان مثال وقوع یک انشعاب یا تغییر کیفی در تابع چگالی احتمال جواب‌ها، یک P -انشعاب محسوب می‌شود).

^۲ مدل رشد لجستیک اولین بار توسط پیر فرانسوا ورهولست *Pierre – Franois Verhulst* در سال ۱۸۳۸ معرفی شد.

^۳ stochastic bifurcation

^۴ D-bifurcation

^۵ dynamical bifurcation

^۶ P-bifurcation

^۷ phenomenological bifurcation

تحلیل P - انشعاب‌ها معمولاً با مطالعه‌ی معادله معروف فوکر-پلانک^۸ انجام می‌گیرد که توزیع احتمال مسیرهای سیستم را توصیف می‌کند؛ در این‌جا ابتدا به معرفی این معادله می‌پردازیم تا بستر لازم برای مطالعه‌ی پدیده‌ی انشعاب در سیستم‌های تصادفی فراهم شود. فرض کنید $f: \mathbb{R}^n \times [t_0, T] \rightarrow \mathbb{R}^{n \times d}$ و $g: \mathbb{R}^n \times [t_0, T] \rightarrow \mathbb{R}^{n \times d}$ توابعی هموار و $\{W(t)\}_{t \geq 0}$ یک حرکت براونی d -بعدی باشد. معادله دیفرانسیل تصادفی زیر را در نظر بگیرید:

$$dx(t) = f(x(t), t) dt + g(x(t), t) dW(t),$$

فرض کنیم $p \equiv p(x, t)$ تابع چگالی احتمال جواب $x(t)$ باشد، در این صورت می‌توان نشان داد که p در معادله فوکر-پلانک متناظر زیر صدق می‌کند (برای ملاحظه جزئیات، مرجع (Arnold, 1974, p. 44) را ملاحظه کنید):

$$\partial_t p = - \sum_i \partial_i [f_i(x(t), t)p] + \frac{1}{2} \sum_{i,j} \partial_i \partial_j \left\{ [g(x(t), t)g^{tr}(x(t), t)]_{ij} p \right\},$$

این معادله تکامل زمانی توزیع احتمال جواب‌های تصادفی سیستم را توصیف می‌کند. مطالعه‌ی دقیق این معادله نه تنها امکان پیش‌بینی رفتار سیستم را در آینده فراهم می‌کند، بلکه به ما اجازه می‌دهد تا رفتار سیستم در حضور نویز را تحلیل کنیم.

۲ بررسی انشعاب‌های تصادفی مدل رشد لجستیک جمعیت

با توجه به تأثیر عوامل تصادفی و غیرقابل پیش‌بینی بر رفتار جمعیت‌ها، بررسی این پدیده در مدل‌های ریاضی اهمیت زیادی دارد. در این بخش، انشعاب‌های ماهیتی معادله رشد لجستیک جمعیت در حضور حرکت وینر مورد مطالعه قرار می‌گیرد تا اثرات نویز تصادفی بر پایداری و رفتار بلندمدت سیستم تحلیل و مدل‌سازی شود. هدف اصلی این بخش، ارائه چارچوبی جامع برای درک عمیق‌تر رفتار سیستم‌های پویا در محیط‌های تصادفی و تحلیل تأثیر تغییرات پارامتری بر ویژگی‌های دینامیکی آن‌ها است. مدل رشد لجستیک جمعیت ۱.۱ را در نظر بگیرید. پس از افزودن فاکتور تصادف به سیستم داریم:

$$\begin{aligned} dx &= \overbrace{\left[x \left(1 - \frac{x}{y} \right) - ax^2 \right]}^{=L_1(x,y)} dt + \sigma x dw_1(t) \\ dy &= \underbrace{\left[-h \frac{x(y-x)}{y(c+x)^2} - bxy + \alpha y \right]}_{=L_2(x,y)} dt + \sigma y dw_2(t), \end{aligned} \quad (1.2)$$

$$x(0) = x_0, \quad y(0) = y_0,$$

که در آن $dw_1(t)$ تغییرات تصادفی جمعیت و $dw_2(t)$ تغییرات تصادفی محیطی را مدل‌سازی می‌کنند. نقطه $E = (x^*, y^*)$ نقطه ثابت دستگاه غیرتصادفی است که در آن $y^* = x^* / (1 - ax^*) > 0$ و $x^* > 0$ جواب معادله زیر است:

$$U(x) = x^2 (a^2(-h) - \alpha + 2bc) + x (ah + bc^2 - 2\alpha c) + bx^3 - \alpha c^2.$$

با استفاده از مقادیر پارامترهای جدول ۱، نقطه تعادل به صورت تقریبی برابر با $E = (x^*, y^*) = (0.0825, 0.0899)$ می‌باشد.

⁸Fokker-Planck equation

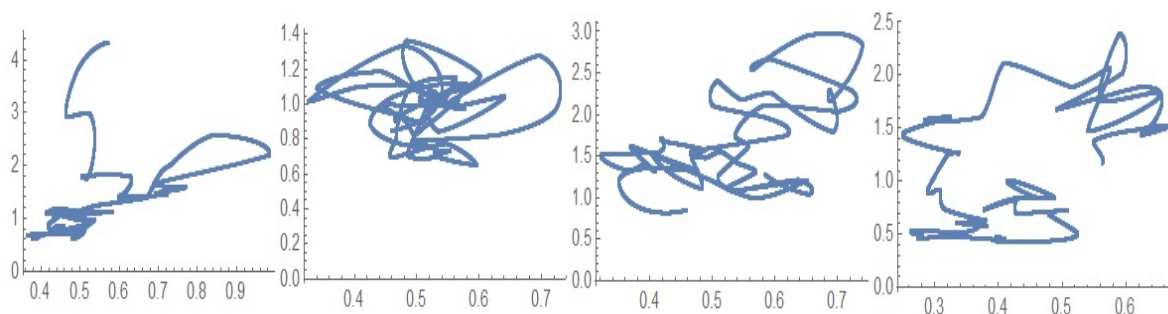
ماتریس قسمت خطی سیستم به صورت زیر است:

$$\begin{bmatrix} 1 - \frac{2x(1+y)}{y} & \frac{x^2}{y^2} \\ -\frac{2}{(2+x)^3} - y + \frac{x(4+y)}{(2+x)^3 y} & 0.1 + x \left(-1 - \frac{x}{(2+x)^2 y^2} \right) \end{bmatrix}$$

محاسبه‌ی مقادیر ویژه این ماتریس در نقطه‌ی E منجر به دو مقدار $\lambda_1 = -1/0.95$, $\lambda_2 = -0.812$ می‌شود. از آن‌جا که هر دو مقدار ویژه دارای قسمت حقیقی منفی و غیرصفر هستند، نتیجه می‌گیریم که نقطه‌ی E یک نقطه‌ی تعادل هذلولوی است. این موضوع، نشانگر ثبات و پایداری موضعی نقطه‌ی E است. بنابراین، در همسایگی این نقطه هیچ‌گونه D -انشعاب (که معمولاً در نقاط ناپایدار رخ می‌دهد) روی نمی‌دهد. شکل ۱ جواب‌های سیستم (۱.۲) را در اطراف E به ازای مقادیر پارامتر جدول ۱ نشان می‌دهد.

σ	T	$y(^{\circ})$	$x(^{\circ})$	b	c	h	α	a
۰.۳	۱۰	۰.۰۸۹۹	۰.۰۸۲۵	۱	۲	۱	۰.۵	۱

جدول ۱: مقادیر پارامتر سیستم (۱.۲)



شکل ۱: جواب‌های تصادفی معادله (۱.۲)

به منظور بررسی P -انشعاب دستگاه ۱.۲، فرض کنید $P \equiv P(t, x, y)$ تابع چکالی احتمال جواب این سیستم باشد، معادله فوکر-پلانک وابسته به این معادله را به شکل زیر در نظر بگیرید:

$$\begin{aligned} \frac{\partial P}{\partial t} &= -\frac{\partial}{\partial x} [L_1(x, y)P] - \frac{\partial}{\partial y} [L_2(x, y)P] + \frac{\sigma^2}{2} \frac{\partial^2}{\partial x^2} [x^2 P] + \frac{\sigma^2}{2} \frac{\partial^2}{\partial y^2} [y^2 P] \\ &= -P \frac{\partial L_1(x, y)}{\partial x} - L_1(x, y) \frac{\partial P}{\partial x} - P \frac{\partial L_2(x, y)}{\partial y} - L_2(x, y) \frac{\partial P}{\partial y} \\ &\quad + \frac{\sigma^2}{2} \left(2P + 2x \frac{\partial P}{\partial x} + x^2 \frac{\partial^2 P}{\partial x^2} + 2P + 2y \frac{\partial P}{\partial y} + y^2 \frac{\partial^2 P}{\partial y^2} \right) \end{aligned}$$

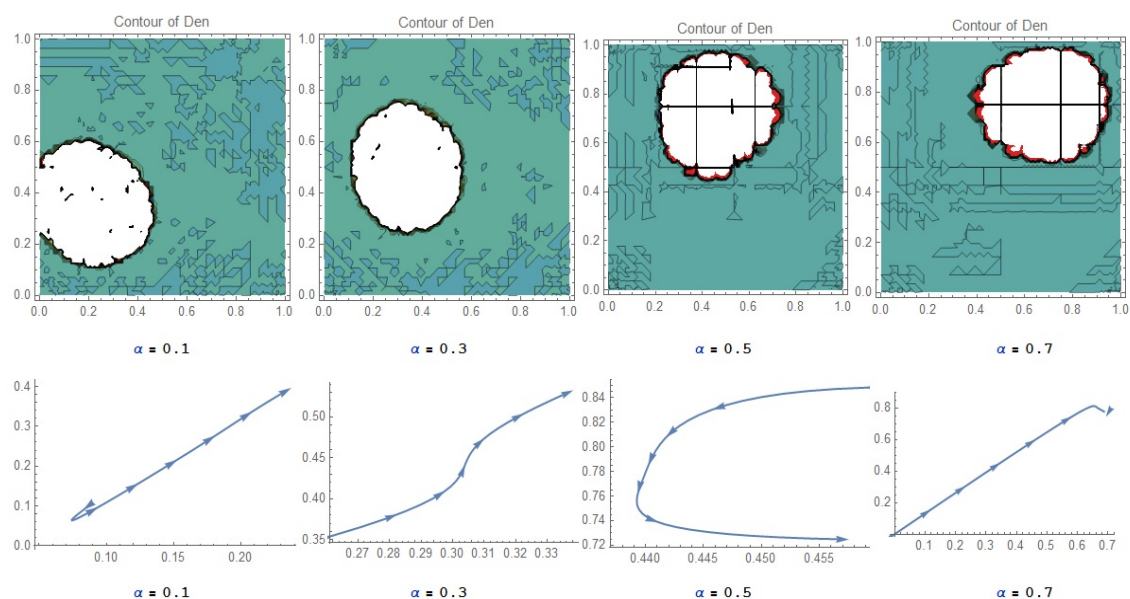
$$\begin{aligned} \frac{\partial P}{\partial t} &= \frac{\partial P}{\partial x} \left[2\sigma^2 x - L_1(x, y) \right] + \frac{\partial P}{\partial y} \left[2\sigma^2 y - L_2(x, y) \right] \\ &\quad + P \left[2\sigma^2 - \frac{\partial L_1(x, y)}{\partial x} - \frac{\partial L_2(x, y)}{\partial y} \right] + \frac{\partial^2 P}{\partial x^2} \left[\frac{\sigma^2 x^2}{2} \right] + \frac{\partial^2 P}{\partial y^2} \left[\frac{\sigma^2 y^2}{2} \right] \end{aligned}$$

با توجه به شرط اولیه زیر،

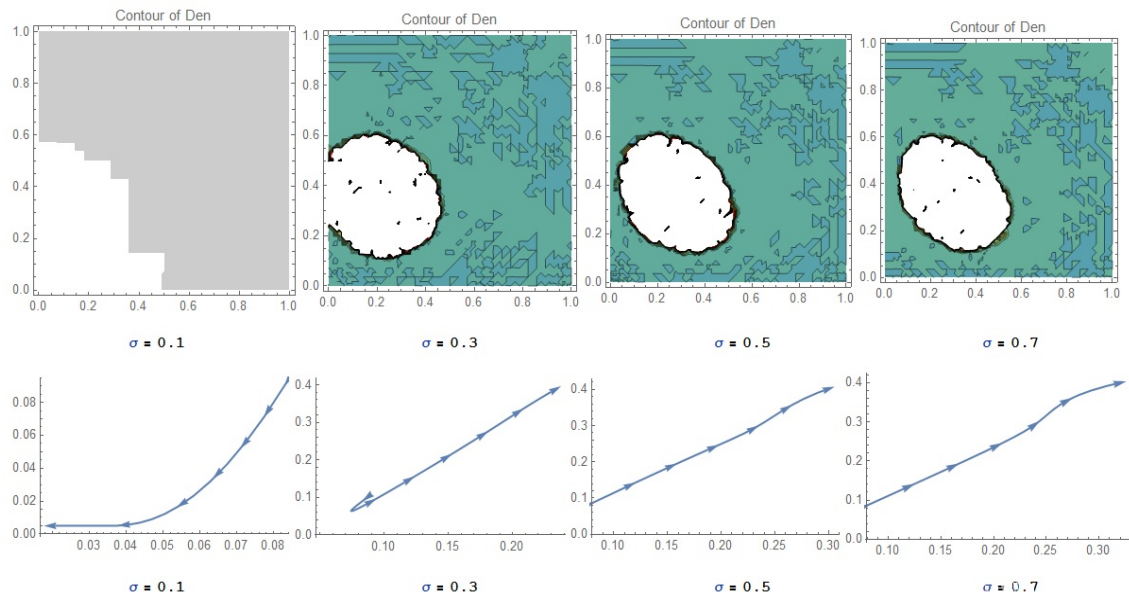
$$P(\circ, x, y) = \Delta(x, y, \epsilon) = \frac{1}{\pi\epsilon} \exp\left(-\frac{(x-x_\circ)^2 + (y-y_\circ)^2}{\epsilon}\right), \quad \circ < \epsilon \ll 1$$

$$\lim_{\epsilon \rightarrow \circ} \Delta(x, y, \epsilon) = \delta(x, y)$$

که در آن $\delta(x, y)$ تابع دلتای دیراک است. نمودارهای ۲ و ۳ تغییرات جواب معادله فوکر-پلانک در صفحه را برحسب تغییرات به ترتیب σ و α نشان می‌دهند و نمودارهای قسمت پایین هر کدام از شکل‌ها امید ریاضی جواب دستگاه (۱.۲) را نشان می‌دهند. همان‌طور که ملاحظه می‌شود نمودار امید ریاضی دچار تغییرات ناگهانی شده که تغییر رفتار سیستم (۱.۲) را نشان می‌دهد و بنابراین این سیستم دچار P -انشعاب شده است.



شکل ۲: برای $\sigma = 0.3; \epsilon = 0.002; a = 1; b = 1; c = 2; h = 1$



شکل ۳: برای $\alpha = 0.1$; $h = 1$; $c = 2$; $b = 1$; $a = 1$; $\epsilon = 0.02$

بحث و نتیجه‌گیری

در این مطالعه، با تکیه بر مدل رشد لجستیک و توسعه آن در بستر سیستم‌های تصادفی، به بررسی پایداری دینامیکی جمعیت در حضور نویزهای تصادفی پرداخته شد. تحلیل معادله رشد لجستیک با در نظر گرفتن حرکت وینر نشان داد که اعمال اختلالات تصادفی می‌تواند منجر به بروز رفتارهای پیچیده و انشعابات تصادفی در سیستم شود. نتایج حاصل از محاسبه مقادیر ویژه در نقطه تعادل E نشان داد که این نقطه هذلولوی بوده و در آن D -انشعاب رخ نمی‌دهد، در حالی که بررسی P -انشعاب در چارچوب معادله فوکر-پلانک، تصویر روشن‌تری از ناپایداری‌های ممکن در سیستم ارائه می‌دهد و نشان می‌دهد که حتی تغییرات کوچک در پارامترهای محیطی می‌توانند منجر به تحولات بنیادی در پاسخ سیستم شوند.

این نتایج اهمیت در نظر گرفتن نویزهای تصادفی در مدل‌سازی پدیده‌های واقعی، به‌ویژه در حوزه‌هایی مانند رشد جمعیت، انتشار بیماری‌ها و گسترش تومورها را برجسته می‌سازد. چارچوب ارائه‌شده در این پژوهش می‌تواند به‌عنوان مبنایی برای تحلیل پایداری و پاسخ سیستم‌های زیستی و اجتماعی در برابر ناپایداری‌های محیطی مورد استفاده قرار گیرد. در مجموع، بررسی مدل رشد لجستیک در بستر تصادفی نه تنها درک عمیق‌تری از رفتار بلندمدت سیستم‌های پویا فراهم می‌سازد، بلکه راهگشای توسعه سیاست‌ها و راهکارهای مبتنی بر داده برای مدیریت پایداری در مواجهه با عدم قطعیت‌های دنیای واقعی خواهد بود.

مراجع

Arnold L. (1974), *Stochastic Differential Equations: Theory and Applications*, Wiley Interscience.

Tsagkis P., Bakogiannis E. and Nikitas A. (2023), *Analysing urban growth using machine learning and*

open data: An artificial neural network modelled case study of five Greek cities, *Sustainable Cities and Society*, **89**, 104337.

Postnikov E. B. (2020), Estimation of COVID-19 dynamics “on a back-of-envelope”: Does the simplest SIR model provide quantitative parameters and predictions?, *Chaos, Solitons and Fractals*, **135**, 109841.

Mansour M. B. A. and Abobakr A. H. (2022), Stochastic differential equation models for tumor population growth, *Chaos, Solitons and Fractals*, **164**, 112738.

Koyama K., Hiura S., Abe H. and Koseki S. (2021), Application of growth rate from kinetic model to calculate stochastic growth of a bacteria population at low contamination level, *Journal of Theoretical Biology*, **525**, 110758.

Investigation of Stochastic Bifurcations in the Logistic Growth Model of Populations

Samira Mahmooee¹, Omid RabieiMotlagh², Hajimohammad Mohammadinejad³

^{1,2,3} Dept. of Applied Math., University of Birjand, Birjand, Iran

Abstract: The logistic growth model is employed as a key tool in analyzing the behavior of populations and biological processes. Given the increasing population and environmental challenges, sustainable resource management is essential. This paper investigates the bifurcations of the logistic growth equation under the influence of Wiener motion and stochastic factors. The results obtained reveal sudden changes in the expected value and the occurrence of P -bifurcation in the system's behavior. The objective is to provide a comprehensive framework for a better understanding of the dynamics of population systems under unstable conditions.

Keywords: Logistic growth model, Stochastic bifurcation, Stochastic equations, Probability density function.

Mathematics Subject Classification (2020): 37G10, 34F05, 37N25.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



تشخیص اعداد دست‌نویس با رهیافت یادگیری ماشین

نوشین رضائی طالقانی^۱، آرزو حاج رجبی
گروه آمار، دانشگاه بین‌المللی امام خمینی (ره)

چکیده:

در سال‌های اخیر با گسترش داده‌های تصویری و افزایش نیاز به پردازش و تحلیل خودکار آن‌ها، مسئله‌ی تشخیص اعداد دست‌نویس به عنوان یکی از مسائل مهم و کاربردی در حوزه یادگیری ماشین و بینایی ماشین مطرح شده است. در این پژوهش مسئله‌ی تشخیص اعداد دست‌نویس با استفاده از الگوریتم‌های یادگیری ماشین مانند جنگل تصادفی و درخت تصمیم بررسی و پیاده‌سازی گردیده است و برای این منظور از پایگاه‌های داده فارسی Hoda بهره گرفته شده است. عملکرد این مدل‌ها با معیارهای آماری مناسب ارزیابی شده و نتایج، حاکی از کارایی بالای جنگل تصادفی در مقایسه با درخت تصمیم می‌باشد.

واژه‌های کلیدی: یادگیری ماشین، تشخیص اعداد دست‌نویس، جنگل تصادفی، درخت تصمیم.

۱ مقدمه

در سال‌های اخیر با گسترش فناوری‌های مبتنی بر هوش مصنوعی، مسائلی نظیر شناسایی و دسته‌بندی الگوهای نوشتاری در تصاویر به‌ویژه تشخیص اعداد دست‌نویس، توجه بسیاری از پژوهشگران را به خود جلب کرده است. تشخیص اعداد دست‌نویس یکی از مسائل مهم در حوزه شناسایی الگوها و پردازش تصویر است که در سامانه‌های بانکی، مراکز پستی، آرشیو اسناد و سایر نهادهای اداری کاربرد فراوان دارد. به دلیل تنوع در سبک نگارش افراد و ناهم‌هنگی‌های مختلف در تصاویر، این مسئله همواره به‌عنوان یکی از چالش‌های اصلی در این حوزه مطرح بوده است. در زمینه بازشناسی اعداد و حروف دست‌نویس پژوهش‌های بسیاری انجام شده است، به عنوان مثال با به‌کارگیری شبکه عصبی و سه مرحله پیش‌پردازش، استخراج ویژگی و دسته‌بندی، بر روی پایگاه Hoda با دقت ۹۵/۳۷٪ بازشناسی شده است (اسدی و دیگران، ۱۳۹۵). **زربافی و همکاران (۱۴۰۲)** با استفاده از شبکه‌های عصبی کانولوشنی و مدل ترکیبی شامل حافظه کوتاه مدت و بلند مدت به اعتبار سنجی ۹۸/۰۳ روی پایگاه داده‌ی اعداد فارسی دست یافتند.

همچنین در پژوهشی، عملکرد سه الگوریتم SVM، درخت تصمیم و شبکه عصبی در تشخیص اعداد دست‌نویس فارسی بر روی پایگاه داده Hoda مقایسه شده است. **(آقای و اخوت، ۱۴۰۳)**

^۱ سخنران، rezaeinrt@gmail.com

در این مقاله برای بازشناسی اعداد دست‌نویس از پایگاه داده فارسی Hoda و از الگوریتم درخت تصمیم و جنگل تصادفی استفاده شده است و با استفاده از معیارهای ارزیابی نشان داده شده که الگوریتم جنگل تصادفی دارای دقت بیشتری نسبت به درخت تصمیم است. در بخش دوم به یادگیری ماشین و مفاهیم مقدماتی پرداخته شده است. الگوریتم جنگل تصادفی و درخت تصمیم در بخش سوم بررسی شده است. در بخش چهارم به مطالعه‌ی روش‌های پیشنهادی بر روی پایگاه داده Hoda پرداخته شده است.

۲ یادگیری ماشین

یادگیری ماشین شاخه‌ای از هوش مصنوعی است که به ماشین‌ها امکان می‌دهد از داده‌ها، الگو بیاموزند و پیش‌بینی کنند. این حوزه در مسائل گوناگون از جمله طبقه‌بندی تصاویر و پردازش متن، کاربرد گسترده‌ای دارد. یکی از مسائل مهم، شناسایی و طبقه‌بندی اعداد دست‌نویس است که به دلیل تنوع در سبک نگارش، چالشی پیچیده به شمار می‌رود. روش‌های یادگیری ماشین به چهار نوع اصلی نظارت‌شده،^۱ بدون نظارت،^۲ نیمه‌نظارت‌شده^۳ و تقویتی^۴ تقسیم می‌شوند. در این پژوهش، تمرکز بر روش‌های نظارت‌شده است که این نوع یادگیری یکی از شاخه‌های اصلی یادگیری ماشین است که در آن مدل با استفاده از مجموعه داده‌های برچسب‌دار آموزش می‌بیند تا بتواند پیش‌بینی‌های دقیقی برای داده‌های جدید ارائه دهد. در این روش، داده‌ها شامل جفت‌های ورودی-خروجی هستند که ورودی‌ها (ویژگی‌ها) و خروجی‌های مورد انتظار (برچسب‌ها) را مشخص می‌کنند. هدف، یادگیری یک تابع نگاشت است که بتواند ورودی‌ها را به برچسب‌های صحیح مرتبط سازد. این رویکرد به‌ویژه در مسائل طبقه‌بندی، مانند تشخیص اعداد دست‌نویس، و رگرسیون کاربرد گسترده‌ای دارد. یکی از نکات کلیدی در یادگیری ماشین، انتخاب ویژگی‌های مناسب از داده‌هاست. ویژگی‌هایی مانند شدت پیکسل‌ها در تصاویر یا الگوهای هندسی می‌توانند به شناسایی بهتر اعداد کمک کنند. پیش‌پردازش داده‌ها مانند نرمال‌سازی یا کاهش نویز دقت مدل‌ها را بهبود می‌بخشد. برای مثال در مسائل طبقه‌بندی تصاویر، تبدیل داده‌ها به مقیاس‌های استاندارد یا استفاده از روش‌های کاهش

جدول ۱: دسته‌بندی الگوریتم‌های یادگیری ماشین

ردیف	نوع یادگیری	الگوریتم‌های پرکاربرد
۱	یادگیری نظارت‌شده	رگرسیون خطی، رگرسیون لجستیک، درخت تصمیم، جنگل تصادفی، نزدیک ترین همسایه
۲	یادگیری بدون نظارت	خوشه‌بندی، تحلیل مؤلفه‌های اصلی
۳	یادگیری نیمه‌نظارت‌شده	خودآموز، ماشین بردار پشتیبان نیمه‌نظارتی
۴	یادگیری تقویتی	روش ترکیبی شبکه عمیق و ماشین بردار پشتیبان

ابعاد مانند تحلیل مؤلفه‌های اصلی، سرعت و کارایی مدل‌ها را افزایش می‌دهد. در جدول ۱ به انواع یادگیری ماشین اشاره شده است که نتایج نشان می‌دهد که انتخاب درست الگوریتم نقش مهمی را در برازش مدل مناسب به داده‌ها ایفا می‌کند.

¹Supervised

²Unsupervised

³Semi-Supervised

⁴reinforcement

۳ جنگل تصادفی و درخت تصمیم

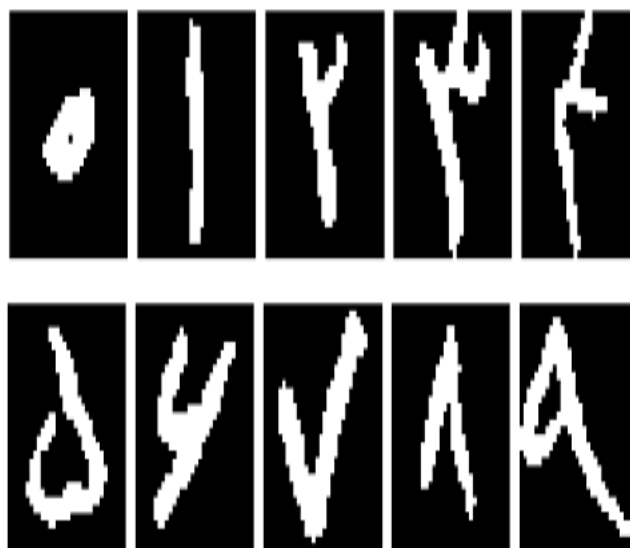
درخت تصمیم یک الگوریتم یادگیری ماشین است که برای شناسایی دست‌نوشته‌های اعداد به کار می‌رود. این روش داده‌های ورودی، مانند ویژگی‌های پیکسلی یا الگوهای هندسی تصاویر دست‌نوشته را در یک ساختار درختی سازمان‌دهی می‌کند. هر گره در درخت یک ویژگی خاص (مثلاً شدت رنگ پیکسل در یک ناحیه) را بررسی می‌کند. و شاخه‌ها، داده‌ها را بر اساس معیارهای تصمیم به زیرمجموعه‌های کوچک‌تر تقسیم می‌کنند تا در نهایت به یک برگ برسند که نشان‌دهنده یک عدد خاص (مثلاً ۰ تا ۹) است. این روش به دلیل سادگی و قابلیت تفسیر آسان، برای مسائل شناسایی اعداد مناسب است. با این حال، اگر داده‌ها دارای خطا باشند یا تنوع زیادی در سبک‌های دست‌خط وجود داشته باشد، درخت تصمیم ممکن است بیش از حد به داده‌های آموزشی وابسته شود و بیش‌برازش بوجود آید و عملکردش روی داده‌های جدید کاهش یابد. جنگل تصادفی یک الگوریتم پیشرفته‌تر است که با ترکیب چندین درخت تصمیم، دقت شناسایی دست‌نوشته‌های اعداد را افزایش می‌دهد. این روش با انتخاب تصادفی زیرمجموعه‌هایی از داده‌ها (مانند تصاویر دست‌نوشته) و ویژگی‌ها (مثلاً بخشی از پیکسل‌ها)، تعداد زیادی درخت تصمیم مستقل می‌سازد و سپس، برای پیش‌بینی نهایی نتایج همه درخت‌ها را ترکیب می‌کند. این رویکرد باعث می‌شود جنگل تصادفی در برابر خطا و تغییرات در دست‌خط‌ها مقاوم‌تر باشد و معمولاً دقت بالاتری نسبت به یک درخت تصمیم تنها ارائه دهد. با این حال، به دلیل استفاده از چندین درخت، محاسبات آن سنگین‌تر است و تفسیر مدل به اندازه یک درخت تصمیم ساده نیست. برای جزییات بیشتر به (Geron, ۲۰۱۹) مراجعه شود.

۴ تحلیل داده‌ی واقعی

در این پژوهش از مجموعه داده Hoda استفاده شده است که اولین مجموعه بزرگ اعداد دست‌نویس فارسی است. این مجموعه شامل ۱۰۲۳۵۳ تصاویر خاکستری از اعداد ۰ تا ۹ است که در شکل ۱ قابل مشاهده است. این مجموعه داده‌ها با توجه به حجم کم، ساختار ساده و برچسب‌گذاری دقیق، گزینه‌ای ایده‌آل برای آزمایش سریع و مقایسه عملکرد الگوریتم‌ها محسوب می‌شوند. جدول ۲، تعداد نمونه‌ها در هر کلاس مشخص شده است.

جدول ۲: تعداد نمونه‌ها در هر کلاس

کلاس	۰	۱	۲	۳	۴	۵	۶	۷	۸	۹
تعداد نمونه	۱۰۰۷۰	۱۰۳۳۰	۹۹۲۳	۱۰۳۳۴	۱۰۳۳۳	۱۰۱۱۰	۱۰۲۵۴	۱۰۳۶۳	۱۰۲۶۴	۱۰۳۷۱



شکل ۱: تصویر اعداد ۰ تا ۹

در گام اول تصاویر 32×32 پیکسل از مجموعه داده Hoda به آرایه‌های یک‌بعدی 10×24 تایی تبدیل شده است، که هر عنصر نشان‌دهنده شدت رنگ خاکستری یک پیکسل (بین ۰ تا ۲۵۵) می‌باشد، این مقادیر بر ۲۵۵ تقسیم شده تا در بازه $[0, 1]$ نرمال‌سازی شوند و برای مدل‌های یادگیری ماشین مناسب‌تر گردند. سپس ۷۰ درصد داده‌ها را به مجموعه‌های آموزشی و ۳۰ درصد را به مجموعه آزمایشی تقسیم کرده و پس از آن مجموعه آموزشی را به مدل‌های یادگیری نظارت‌شده مانند درخت تصمیم و یا جنگل تصادفی می‌دهیم. مدل با داده‌های آموزشی و برچسب‌های مربوطه (اعداد ۰ تا ۹) آموزش می‌بیند، در مدل درخت تصمیم با استفاده از شاخص جینی به‌عنوان معیار پیش‌فرض تقسیم‌بندی، گره‌ها و شاخه‌ها تشکیل شده است. شاخص جینی میزان ناخالصی گره‌ها را با فرمول (۱.۴)

$$Gini = 1 - \sum_{i=1}^K p_i^2 \quad (1.4)$$

محاسبه می‌کند، که در آن p_i احتمال حضور در کلاس i است. مدل جنگل تصادفی از مجموعه‌ای از درخت‌های تصمیم مستقل تشکیل شده که هر یک با شاخص جینی آموزش دیده و پیش‌بینی‌های آن‌ها از طریق رأی‌گیری اکثریت ترکیب می‌شوند. برای ارزیابی عملکرد مدل‌ها از معیارهای زیر استفاده شده است:

۱. دقت^۷: این معیار نسبت نمونه‌های درست پیش‌بینی‌شده به کل نمونه‌ها را نشان می‌دهد.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

• TP^8 : نمونه‌های مثبت که به‌درستی مثبت پیش‌بینی شده‌اند (مثبت درست).

• TN^9 : نمونه‌های منفی که به‌درستی منفی پیش‌بینی شده‌اند (منفی درست).

⁵Normalization

⁶Gini Impurity

⁷Accuracy

⁸True Positive

⁹True Negative

• FP^0 : نمونه‌های منفی که به اشتباه مثبت پیش‌بینی شده‌اند (مثبت نادرست).

• FN^1 : نمونه‌های مثبت که به اشتباه منفی پیش‌بینی شده‌اند (منفی نادرست).

۲. **صحت^{۱۲}** : این معیار نسبت پیش‌بینی‌های درست مثبت به کل پیش‌بینی‌های مثبت را نشان می‌دهد و برای کاهش خطاهای مثبت نادرست (FP) اهمیت دارد.

$$Precision = \frac{TP}{TP + FP}$$

۳. **بازخوانی یا حساسیت^{۱۳}** : این معیار نسبت نمونه‌های مثبت درست شناسایی شده به کل نمونه‌های مثبت واقعی را نشان می‌دهد و برای کاهش خطاهای منفی نادرست (FN) حیاتی است.

$$Recall = \frac{TP}{TP + FN}$$

۴. **امتیاز $F1$ ^{۱۴}** : میانگین هارمونیک صحت و حساسیت است که برای داده‌های نامتوازن مناسب بوده و تعادل بین این دو معیار را برقرار می‌کند.

$$F1 \text{ Score} = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

معیارها به صورت درصدی در بازه $[0\%, 100\%]$ گزارش شده‌اند، که 100% بهترین و 0% بدترین عملکرد را نشان می‌دهد. نتایج ارزیابی مدل‌های درخت تصمیم در جدول ۳ و جنگل تصادفی در جدول ۴ ارائه شده است.

ماتریس‌های درهم‌ریختگی ارائه شده، عملکرد دو مدل یادگیری ماشین، یعنی جنگل تصادفی و درخت تصمیم، را بر روی داده‌های آزمایشی برای تشخیص اعداد دست‌نویس (۰ تا ۹) نشان می‌دهند. ماتریس درهم‌ریختگی جنگل تصادفی نشان‌دهنده دقت بسیار بالایی این مدل است. مقادیر قطری (پیش‌بینی‌های درست) برای هر کلاس (اعداد ۰ تا ۹) نشان می‌دهند که اکثر نمونه‌ها به درستی طبقه‌بندی شده‌اند. به عنوان مثال، برای عدد ۰، مدل در ۲۹۶۴ مورد از مجموع نمونه‌ها پیش‌بینی درست داشته و تنها در ۴ مورد خطا کرده است (۱) و ماتریس درهم‌ریختگی درخت تصمیم عملکرد ضعیف‌تری را نسبت به جنگل تصادفی نشان می‌دهد. به طور کلی، جنگل تصادفی به دلیل استفاده از چندین درخت تصمیم و کاهش واریانس، عملکرد بهتری نسبت به درخت تصمیم در تشخیص اعداد دست‌نویس دارد. تعداد بالای پیش‌بینی‌های درست و خطاهای کمتر در ماتریس درهم‌ریختگی جنگل تصادفی، برتری این مدل را تأیید می‌کند.

ماتریس درهم‌ریختگی درخت تصمیم (دو کلاس اول)

$$\begin{pmatrix} 2908 & 3 & 1 & 5 & 15 & 19 & 4 & 8 & 5 & 0 \\ 1 & 3005 & 10 & 1 & 7 & 3 & 12 & 2 & 3 & 22 \end{pmatrix}$$

ماتریس درهم‌ریختگی جنگل تصادفی (دو کلاس اول)

$$\begin{pmatrix} 2964 & 0 & 0 & 0 & 1 & 3 & 0 & 0 & 0 & 0 \\ 2 & 3054 & 1 & 0 & 3 & 0 & 3 & 0 & 0 & 3 \end{pmatrix}$$

¹⁰False Positive

¹¹False Negative

¹²Precision

¹³Recall

¹⁴Score- $F1$

جدول ۴: معیارهای ارزیابی جنگل تصادفی

امتیاز F۱	بازخوانی یا حساسیت	صحت	کلاس
%۱۰۰	%۱۰۰	%۹۹	۰
%۹۹	%۱۰۰	%۹۹	۱
%۹۸	%۹۸	%۹۷	۲
%۹۷	%۹۷	%۹۸	۳
%۹۷	%۹۷	%۹۷	۴
%۹۹	%۹۹	%۹۹	۵
%۹۸	%۹۷	%۹۹	۶
%۹۹	%۹۹	%۹۹	۷
%۱۰۰	%۱۰۰	%۱۰۰	۸
%۹۸	%۹۹	%۹۸	۹
دقت : %۹۸			

جدول ۳: معیارهای ارزیابی درخت تصمیم

امتیاز F۱	بازخوانی یا حساسیت	صحت	کلاس
%۹۸	%۹۸	%۹۷	۰
%۹۷	%۹۸	%۹۶	۱
%۹۱	%۹۲	%۸۹	۲
%۹۱	%۹۱	%۹۲	۳
%۹۰	%۸۹	%۹۱	۴
%۹۴	%۹۲	%۹۵	۵
%۹۱	%۹۰	%۹۱	۶
%۹۶	%۹۶	%۹۵	۷
%۹۷	%۹۷	%۹۷	۸
%۹۴	%۹۳	%۹۴	۹
دقت : %۹۴			

بحث و نتیجه‌گیری

این پژوهش به بررسی تشخیص اعداد دست‌نویس با استفاده از روش‌های یادگیری ماشین پرداخت. با بهره‌گیری از مجموعه داده Hoda و اعمال پیش‌پردازش‌های لازم، الگوریتم‌های جنگل تصادفی و درخت تصمیم آزمایش شدند. بوسیله ماتریس‌های درهم‌ریختگی و معیارهای ارزیابی نشان داده شد که جنگل تصادفی خطاهای کمتر و دقت پیش‌بینی بهتری نسبت به درخت تصمیم دارد. این برتری به دلیل ترکیب چندین درخت تصمیم و کاهش واریانس در جنگل تصادفی است. در نهایت، این پژوهش نشان داد که انتخاب الگوریتم مناسب نقش کلیدی در بهبود دقت مدل‌های یادگیری ماشین دارد.

مراجع

اسدی، فرشید و صیدی پیری، رسول و نوری، زینب و لطفی پور، امین (۱۳۹۵)، تشخیص حروف دست‌نویس فارسی با استفاده از حسگر شتاب‌سنج و الگوریتم‌های یادگیری ماشین، چهارمین کنفرانس بین‌المللی مهندسی برق و کامپیوتر، تهران، ایران.

زربافی، م. و همکاران (۱۴۰۲)، تشخیص اعداد گفتاری با استفاده از شبکه‌های یادگیری عمیق، مجله علمی-پژوهشی فناوری اطلاعات و ارتباطات، ۱۰، ۴۵-۵۶.

آقای، امیرعلی و اخوت، راضیه سادات (۱۴۰۳)، مقایسه عملکرد ماشین بردار پشتیبان، درخت تصمیم و شبکه عصبی در تشخیص اعداد دست‌نویس فارسی، نوزدهمین کنفرانس سیستم‌های هوشمند ایران، سیرجان، ایران.

Geron, A. (2019), *Hands on Machine Learning with Scikit Learn, Keras and TensorFlow*, Second Edition, O'Reilly Media, Inc.

Handwritten Digit Recognition Using a Machine Learning Approach

Arezo Hajrajabi¹, Nooshin Rezaei Taleghani²

¹Department of Statistics, Imam Khomeini International University, Qazvin

²Department of Statistics, Imam Khomeini International University, Qazvin

Abstract:

In recent years, with the expansion of visual data and the growing need for its automatic processing and analysis, the issue of handwritten digit recognition has emerged as a significant and practical problem in the fields of machine learning and computer vision. In this research, the problem of handwritten digit recognition has been investigated and implemented using machine learning algorithms such as Random Forest and Decision Tree. For this purpose, the Persian Hoda dataset has been utilized. The performance of these models has been evaluated using appropriate statistical metrics, and the results indicate the superior efficiency of Random Forest compared to Decision Tree.

Keywords: Machine Learning, Handwritten Digit Recognition, Random Forest, Decision Tree.

Mathematics Subject Classification (2020): 68T07, 68T10, 62H30.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30-31, 2025
University of Kurdistan



مقایسه دو رویکرد نیم‌رخی و نامقید در برازش مدل رگرسیون بر اساس نمونه داده‌های ترکیبی با تواتر متفاوت

فرزین رنجبر^۱، علی آقامحمدی^۲، مجید ادیب^۳

^۱گروه آمار، دانشگاه زنجان

^۳گروه ریاضی، دانشگاه زنجان

چکیده: این مقاله به مقایسه روش نیم‌رخی و نامقید در برآورد پارامترهای مدل رگرسیون بر اساس نمونه داده‌های ترکیبی با تواتر متفاوت (MIDAS) می‌پردازد. روش نیم‌رخی با ترکیب حداقل مربعات معمولی (OLS) برای برآورد ضرایب رگرسیونی و بهینه‌سازی پارامترهای وزن‌دهی چندجمله‌ای، کارایی محاسباتی بهتر و خطای پیش‌بینی کمتری نسبت به روش برآورد پارامتر نامقید نشان می‌دهد.

واژه‌های کلیدی: مدل رگرسیونی نمونه داده‌های ترکیبی با تواتر متفاوت، روش برآورد پارامتر نامقید، روش برآورد پارامتر نیم‌رخی.
کد موضوع‌بندی ریاضی (۲۰۲۰): 62M10، 91B84.

۱ مقدمه

در اقتصادسنجی مدرن، تحلیل داده‌های سری زمانی با تواتر متفاوت به یکی از چالش‌های اساسی محققان تبدیل شده است. مدل‌های رگرسیونی بر اساس نمونه داده‌های ترکیبی با تواتر متفاوت (MIDAS)^۲ که توسط گیلز و همکاران (۲۰۰۴) و گیلز و همکاران (۲۰۰۶) معرفی شدند، ابزاری قدرتمند برای مواجهه با این چالش محسوب می‌شوند. این مدل‌ها قابلیت ترکیب اطلاعات متغیرهای اقتصادی با تواترهای متفاوت نظیر داده‌های روزانه، ماهانه و فصلی را فراهم می‌آورند. اهمیت این تحقیق در چندین جنبه نمایان می‌شود. نخست، مدل‌های MIDAS سنتی نیازمند استفاده از روش‌های حداقل مربعات غیرخطی هستند که از لحاظ محاسباتی پیچیده و گاهی دارای مشکلات همگرایی می‌باشند آندرو و همکاران (۲۰۱۰).

در ادامه برای کاهش پیچیدگی محاسباتی و کاهش تعداد پارامترها، روش استنباطی نیم‌رخی ارائه شده است که ابتدا به مفهوم آن می‌پردازیم.

^۱ Farzin.r@znu.ac.ir
^۲ Mixed Data Sampling

مفهوم تابع درست‌نمایی نیم‌رخشی^۱ ریشه‌ای عمیق در ادبیات آماری دارد و توسط محققانی چون پاتینفلد (۱۹۷۷) توسعه یافته است. فرض کنید متغیرهای تصادفی y_t مستقل و هم‌توزیع با تابع چگالی احتمال $f(y; \delta, \theta)$ باشند و هدف برآورد پارامترهای δ, θ باشد. با یک نمونه به اندازه T لگاریتم تابع درست‌نمایی به صورت زیر تعریف می‌شود.

$$\mathcal{L}(\delta, \theta) = \sum_{t=1}^T \log f(y_t; \delta, \theta)$$

در بسیاری مواقع، بیشینه‌سازی این تابع نسبت به تمامی پارامترها دشوار است. اما اگر θ را ثابت در نظر گرفته شود، یعنی $\theta = \bar{\theta}$ ، بیشینه‌سازی نسبت به δ برای θ ثابت ساده‌تر است. در چنین مواردی می‌توان لگاریتم تابع درست‌نمایی به شکل $L_{T,\delta}(\theta)$ بازنویسی کرد که به آن لگاریتم تابع درست‌نمایی نیم‌رخشی گفته می‌شود و نسبت به δ برای یک θ معین بهینه می‌شود. بنابراین مسئله بهینه‌سازی به $(\hat{\delta}, \hat{\theta}) = \arg \max_{\theta} L_{T,\delta}(\theta)$ تبدیل می‌شود. این توابع عمدتاً زمانی مورد استفاده قرار می‌گیرند که δ پارامتری با بعد کم و θ پارامتری با بعد بالا یا مزاحم محسوب شود. مورفی (۲۰۰۰) توجیه کلی برای استفاده از تابع درست‌نمایی نیم‌رخشی را به عنوان یک رهیافت استنباطی ارائه دادند.

در این پژوهش، رویکرد تابع درست‌نمایی نیم‌رخشی برای مدل‌های رگرسیون MIDAS ارائه می‌شود که در آن θ نشان‌دهنده پارامترهای وزن چندجمله‌ای MIDAS و δ شامل پارامترهای عرض از مبدأ و شیب مدل رگرسیون است. برای مقدار ثابت $\theta = \bar{\theta}$ ، رویکرد نیم‌رخشی مسئله برآورد را به مجموعه‌ای از برآوردهای حداقل مربعات معمولی (OLS) مدل‌های رگرسیون خطی نسبتاً ساده $\hat{\delta}(\bar{\theta})$ تقلیل می‌دهد. از آنجا که پارامترهای چندجمله‌ای MIDAS یعنی θ دارای ابعاد پایین هستند - گاهی اوقات حتی تک‌بعدی - انتخاب $\bar{\theta}$ در یک شبکه نقاط، نسبتاً آسان است. که این امر بهینه‌سازی را از نظر محاسباتی ساده می‌نماید.

روش دیگر در تحلیل مدل‌های رگرسیون MIDAS، روش MIDAS نامقید (U-MIDAS)^۲ است که توسط فورنی و همکاران (۲۰۱۵) پیشنهاد شده است. U-MIDAS یک روش برآورد برای مدل‌های MIDAS است که برخلاف روش استاندارد MIDAS که از توابع پارامتریک (مانند چندجمله‌ای بتا و ...) برای وزن‌دهی به وقفه داده‌های با تواتر بالا استفاده می‌کند. U-MIDAS هیچ محدودیت پارامتریکی بر روی وزن‌های وقفه‌ها اعمال نمی‌کند. در عوض، این روش هر وزن وقفه را به صورت جداگانه و بدون ساختار با استفاده از حداقل مربعات معمولی (OLS) برآورد می‌کند. این روش زمانی جذاب است که اختلاف تواترها کم باشد. مثلاً در داده‌های فصلی/ماهانه (که هر فصل شامل ۳ ماه است)، U-MIDAS می‌تواند ۳، ۶ یا ۹ وقفه ماهانه را مستقیماً در مدل وارد کند و تمام پارامترها را با OLS برآورد کند. این رویکرد برای مواردی که ساختار وزن‌ها نامشخص است یا تفاوت تواترها کم است، مناسب می‌باشد. اما اگر اختلاف تواترها بزرگ باشد (مثلاً داده‌های روزانه در مدل‌های فصلی که هر فصل شامل ۹۰ روز است) U-MIDAS نیاز به تخمین تعداد زیادی پارامتر (مثلاً وزن ۹۰ وقفه) دارد که منجر به پیچیدگی محاسباتی بالا و مشکل بیش‌برازش می‌شود. نتیجه قابل توجه و غیرمنتظره این است که در مواردی که U-MIDAS جذاب به نظر می‌رسد، رویکرد جدید نیم‌رخشی برای مدل رگرسیون MIDAS از نظر محاسباتی سریع‌تر از روش OLS است فورنی و همکاران (۲۰۱۵). دلیل این امر آن است که U-MIDAS مثلاً در حالت فصلی/ماهانه یک مدل رگرسیونی با ۳، ۶ یا ۹ وقفه از داده‌های ماهانه و داده‌های فصلی با وقفه است که بر روی داده‌های فصلی جاری یا آینده تصویر می‌شوند، در حالی که روش برآورد نیم‌رخشی تنها شامل یک متغیر توضیحی (با نادیده گرفتن عرض از مبدأ در هر دو حالت) است که این متغیر توضیحی با اعمال چندجمله‌ای MIDAS (با فرض $\bar{\theta}$) بر داده‌های با تواتر بالا ساخته می‌شود. روش برآورد نیم‌رخشی از نظر محاسباتی کارآمد، از لحاظ آماری جذاب و پیچیدگی محاسباتی روش کمترین مربعات غیرخطی را برای برآورد پارامترهای مدل MIDAS را ندارد.

^۱ Profiling

^۲ Unrestricted Mixed Data Sampling

در مجموعه مدل‌های MIDAS، مطالعات متنوعی انجام شده است. کلمنتز و همکاران (۲۰۰۹) بر جنبه‌های پیش‌بینی این مدل‌ها تمرکز کردند، در حالی که بررسی‌های جامع‌تری توسط آرمستو و همکاران (۲۰۱۰) و آندرو و همکاران (۲۰۱۰) انجام شده است. گیلز و همکاران (۲۰۱۱) نخستین بار ایده نیم‌رخ را در رگرسیون چندکی MIDAS به کار گرفتند. هدف اصلی این مطالعه، ارائه روشی نوین برای برازش مدل‌های MIDAS است که مزایای محاسباتی OLS را با دقت و انعطاف‌پذیری مدل‌های پارامتری ترکیب می‌کند. نوآوری اصلی در استفاده از رویکرد نیم‌رخ برای برآورد پارامترهای چندجمله‌ای بتا نهفته است، که امکان کاهش پیچیدگی محاسباتی را فراهم می‌آورد.

یافته‌های ما نشان می‌دهد که مدل MIDAS با رویکرد نیم‌رخ، علی‌رغم دارا بودن ضریب تعیین (R^2) کمتر در نمونه، از قدرت پیش‌بینی بالاتری در خارج از نمونه برخوردار است. میانگین مربعات خطای پیش‌بینی خارج از نمونه در روش نیم‌رخ به طور قابل توجهی پایین‌تر است. این یافته نشان‌دهنده این است که U-MIDAS مشکل بیش برازشی نیز دارد. در بخش دوم مدل رگرسیون MIDAS به اختصار توضیح داده شده و رویکردهای نیم‌رخ و نامقید در این مدل‌ها بررسی می‌شوند. بخش سوم شامل شبیه‌سازی برای مقایسه این رویکرد و تحلیل نتایج حاصل از آن است.

۲ مدل رگرسیون MIDAS

مدل‌های MIDAS معمولاً بر ترکیبی از دو تواتر متفاوت، به ترتیب تواتر بالا و پایین، تمرکز دارند. فرض کنید $t = 1, \dots, T$ نشان‌دهنده واحد زمانی با تواتر پایین و m تعداد دفعاتی است که تواتر بالا در همان واحد زمانی ظاهر می‌شود. به عنوان مثال، برای رشد تولید ناخالص داخلی فصلی و شاخص‌های ماهانه به عنوان متغیرهای توضیحی، $m = 3$ خواهد بود. متغیر با فرکانس پایین با y_t^L و متغیر با فرکانس بالا نیز با $x_{t-j/m}^H$ نمایش داده می‌شود که در آن $t - j/m$ بیانگر دوره زمانی j ام (گذشته) با تواتر بالاتر است و $j = 0, 1, \dots$ برای یک ترکیب فصلی/ماهانه، $x_t^H, x_{t-1/3}^H, x_{t-2/3}^H$ به ترتیب نشان‌دهنده آخرین، دومین و اولین ماه فصل t هستند.

مدل‌های رگرسیونی MIDAS در اصل نوعی مدل‌های رگرسیونی هستند که داده‌هایی با تواترهای متفاوت (مثلاً ماهانه و فصلی) را به طور همزمان در یک مدل ترکیب می‌کنند. تأثیر متغیر با تواتر بالا با استفاده از یک تابع وزن‌دهی مدل می‌شود که به صورت خودکار وزن‌های تمام وقفه‌های گذشته را محاسبه و توزیع می‌کند. این کار باعث می‌شود تعداد پارامترهای مدل به حداقل برسد و مشکلات مربوط به انتخاب تعداد وقفه‌های مناسب نیز رفع گردد.

مدل اساسی MIDAS با یک متغیر توضیحی با تواتر بالا برای پیش‌بینی h دوره بعد در تواتر پایین، در حالتی که داده‌های با تواتر بالا تا x_t^H در دسترس هستند، به صورت زیر

$$y_{t+h}^L = a_h + b_h C(L^{1/m}; \theta_h) x_t^H + \varepsilon_{t+h}^L \quad (1.2)$$

است، که در آن $C(L^{1/m}; \theta) = \sum_{j=0}^{j_{\max}-1} c(j; \theta) L^{j/m}$ و $C(1; \theta) = \sum_{j=0}^{j_{\max}-1} c(j; \theta) = 1$ است. مؤلفه j_{\max} حداکثر تعداد وقفه‌هایی است که از متغیر با تواتر بالا (مثلاً داده‌های روزانه) در مدل استفاده می‌شود. این مدل اغلب DL-MIDAS نامیده می‌شود که DL به معنای وقفه توزیع شده است. پارامتری‌سازی ضرایب وقفه‌دار $c(j; \theta)$ به شیوه‌ای صرفه‌جویانه، یکی از ویژگی‌های کلیدی مدل‌های MIDAS است.

هرچند محدود کردن فضای پارامتری θ به یک بعد اجباری نیست، اما این کار محاسبات را تسهیل می‌کند. در این راستا، چندجمله‌ای‌های

بتا که ابتدا توسط گیلز و همکاران (۲۰۰۶) پیشنهاد شدند، گزینه مناسبی هستند. این توابع مبتنی بر توزیع بتا بوده و تنها به دو پارامتر نیاز دارند، یعنی:

$$c(j; \theta_1, \theta_2) = \frac{f(\frac{j}{j_{\max}}, \theta_1; \theta_2)}{\sum_{j=0}^{j_{\max}-1} f(\frac{j}{j_{\max}}, \theta_1; \theta_2)}, \quad (2.2)$$

که در آن

$$f(x, a, b) = \frac{x^{a-1}(1-x)^{b-1}\Gamma(a+b)}{\Gamma(a)\Gamma(b)}, \quad (3.2)$$

و $\Gamma(a) = \int_0^\infty e^{-x} x^{a-1}$ است. با ثابت نگه‌داشتن یکی از پارامترها (مثلاً $\theta_1 = 1$)، می‌توان مدل را به یک پارامتر (θ_2) تقلیل داد که برای داده‌های اقتصادی مطلوب است.

۱.۲ مدل اتورگرسیو MIDAS

روش ارائه‌شده در این مقاله به مدل DL-MIDAS در معادله (۱.۲) محدود نمی‌شود. به عنوان مثال، افزودن متغیرهای وابسته با وقفه (یعنی متغیرهای کم تواتر) منجر به دسته‌ای از مدل‌های ADL-MIDAS^۱ می‌شود. برای ساده‌سازی نمادگذاری، در ادامه مقاله (بدون از دست دادن کلیت مسئله) فرض می‌کنیم $h = 1$. با این فرض یک مدل اتورگرسیو از مرتبه p را به صورت

$$y_{t+1}^L = a + \sum_{j=1}^p \rho_j y_{t-j+1}^L + bC(L^{1/m}; \theta) x_t^H + \varepsilon_{t+1}^L \quad (4.2)$$

بازنویسی می‌کنیم. حال می‌توان روش نیم‌رخی را با ثابت در نظر گرفتن پارامترهای چند جمله‌ای MIDAS برای این مدل نیز بکار گرفت.

۲.۲ مدل MIDAS نامقید

مدل U-MIDAS که توسط فورنی و همکاران (۲۰۱۵) معرفی شد. برای توصیف این مدل فرض کنید m کوچک باشد، (مثلاً برابر ۳ برای ترکیب داده‌های فصلی/ماهانه) برخلاف مدل استاندارد MIDAS از هیچ تابع وزن‌دهی از پیش تعیین‌شده‌ای در این مدل استفاده نمی‌شود. برای هر وقفه از متغیر با تواتر بالا، یک ضریب جداگانه در نظر گرفته می‌شود. یعنی به جای برآورد $bC(L^{1/m}; \theta)$ در معادله ۴.۲ ضرایب متغیر با تواتر بالا به صورت جداگانه برآورد می‌شود. فرم کلی این مدل به صورت

$$y_{t+1} = a + \sum_{j=1}^p \rho_j y_{t-j+1}^L + \sum_{j=0}^{J_{\max}-1} c_j x_{t-j/m}^H + \varepsilon_{t+1}^L \quad (5.2)$$

است. این بدان معناست که علاوه بر پارامترهای a و ρ_j ، ما $J_{\max} - 2$ پارامتر اضافی را تخمین می‌زنیم. وقتی $m = 3$ و $J_{\max} - 1$ کوچک باشد (مثلاً حداکثر وقفه سالانه برابر ۴ باشد) و حجم نمونه به اندازه کافی بزرگ باشد که جمله خطا ε_{t+1}^L را ناهمبسته سازد، پارامترهای مدل U-MIDAS را می‌توان با حداقل مربعات معمولی (OLS) تخمین زد.

۳ شبیه‌سازی

در این بخش با استفاده از یک مطالعه شبیه‌سازی به بررسی کارایی هریک از رویکردهای ارائه شده می‌پردازیم. فرض کنید مقدار $m = 3$ برای ترکیب داده‌های فصلی/ماهانه باشد. مدل رگرسیون MIDAS را به صورت

$$y_{t+1}^L = \sum_{j=0}^{j_{\max}-1} c(j, \theta) x_{t-\frac{j}{m}}^H + \epsilon_{t+1}^L, \quad t = 1, \dots, T \quad (1.3)$$

در نظر می‌گیریم. در این مدل به دلیل اینکه ضرایب عرض از مبدأ و اتورگرسیو مربوط به متغیر کم تواتر در تحلیل ضروری نیستند، بنابراین کنار گذاشته می‌شوند. مدل U-MIDAS فرض می‌کند که ضرایب بدون هیچ محدودیتی و با استفاده از روش حداقل مربعات معمولی (OLS) برآورد می‌شوند؛ در حالی که مدل MIDAS با چندجمله‌ای بتا، وزن‌ها را به صورت بسیار مقید و پارامتریک به شکل $c(j, \theta) \propto (1 - \frac{j}{j_{\max}})^{\theta-1}$ مدل‌سازی می‌کند. روش نیم‌رخ فرض می‌کند پارامتر θ ثابت است و با روش حداقل مربعات معمولی، پارامترهای دیگر (مثل شیب) را برآورد کرده سپس مقدار θ را طوری تنظیم می‌کند که مدل بهترین برازش را داشته باشد.

در شبیه‌سازی مونته‌کارلو، داده‌های ماهانه از یک مدل $AR(1)$ با خودهمبستگی 0.8 و داده‌های فصلی از مدل MIDAS با چندجمله‌ای بتا تولید شده‌اند. مقادیر خودهمبستگی‌های $0.8, 0.1$ را نیز آزمایش شد و نتایج مشابهی حاصل گردید. مقادیر $T = \{50, 200, 2000\}$ (اندازه نمونه)، $j^{max} = \{6, 9, 12\}$ (حداکثر وقفه) و $\theta = \{2, 10\}$ (پارامترهای توزیع بتا) را در نظر می‌گیریم که ترکیبات آنها منجر به ۱۸ سناریوی مختلف می‌شود. شبیه‌سازی را به تعداد ۱۰۰۰ مرتبه تکرار کرده و در هر سناریو و تکرار، داده‌های ماهانه/فصلی را شبیه‌سازی می‌کنیم سپس مدل رگرسیون MIDAS را با استفاده از رویکرد نیم‌رخ و مدل U-MIDAS را بر داده‌ها برازش می‌دهیم. مدل‌ها را با استفاده از معیارهای ارزیابی به شرح ذیل ارزیابی می‌کنیم:

$$\begin{aligned} Bias^2 &= \sum_{j=0}^{j_{\max}-1} (\hat{c}(j, \theta) - c(j, \theta))^2, \\ R^2 &= 1 - \frac{\sum_{t=1}^T (\hat{y}_t^L - y_t^L)^2}{\sum_{t=1}^T y_t^2}, \\ MSE &= \frac{1}{r} \sum_{t=T+1}^{T+r} (\hat{y}_t^L - y_t^L)^2, \end{aligned}$$

که در آن $\hat{c}(j, \theta)$ وزن‌های برآورد شده و \hat{y}_t^L مقدار برازش شده است.

نتایج شبیه‌سازی در جدول ۱ ارائه شده است که در آن نتایج ۱۸ سناریو شبیه‌سازی برای سه شاخص عنوان شده، گزارش شده است. با توجه به جدول ۱ مشخص است که برآوردهای نیم‌رخ در تمامی سناریوها مربع اریبی کمتری نسبت به U-MIDAS دارد. به عنوان مثال، در سناریوی $T = 50, j^{max} = 6$ و $\theta = 2$ ، مربع اریبی روش نیم‌رخ 0.087 است، در حالی که این مقدار برای U-MIDAS به 0.391 می‌رسد. این نشان می‌دهد که برآورد نیم‌رخ به طور قابل توجهی دقیق‌تر از U-MIDAS است.

در بیشتر سناریوها R^2 درون‌نمونه‌ای روش نیم‌رخ کمتر از U-MIDAS است. با این حال، قدرت پیش‌بینی خارج از نمونه‌ای آن بالاتر است، به طوری که MSE پیش‌بینی آن به طور قابل توجهی کمتر از U-MIDAS است. به عنوان مثال، در سناریوی $T = 50, j^{max} = 12$ و $\theta = 2$ ، MSE پیش‌بینی روش نیم‌رخ 0.989 است، در حالی که این مقدار برای U-MIDAS برابر 1.322 (حدود ۲۶٪ بیشتر) است. این نتایج نشان می‌دهد که U-MIDAS از مشکل برازش بیش از حد رنج می‌برد. پارامترهای غیرمقید در U-MIDAS می‌توانند با R^2 درون‌نمونه‌ای بالا تنظیم شوند، اما نمی‌توانند در پیش‌بینی برون‌نمونه‌ای از یک مدل کم پارامتر پیشی

جدول ۱: مقایسه عملکرد پروفایلی و U-MIDAS در ۱۸ سناریوی شبیه‌سازی شده

T	j_{\max}	θ	$Bias^{\dagger}$	R^2	MSE Forecast	FLOPS (10^3)
پروفایلی	پروفایلی	پروفایلی	پروفایلی	پروفایلی	پروفایلی	پروفایلی
۵۰	۶	۲	۰.۰۸۷	۰.۳۹۱	۰.۶۵۸	۰.۶۹۰
۵۰	۶	۱۰	۰.۱۱۹	۰.۴۰۶	۰.۷۴۲	۰.۷۶۵
۵۰	۹	۲	۰.۰۶۶	۰.۴۹۷	۰.۶۱۰	۰.۶۸۷
۵۰	۹	۱۰	۰.۱۲۷	۰.۵۰۵	۰.۷۲۱	۰.۷۶۲
۵۰	۱۲	۲	۰.۰۵۵	۰.۶۱۹	۰.۵۷۳	۰.۶۶۹
۵۰	۱۲	۱۰	۰.۱۱۴	۰.۶۳۱	۰.۶۹۳	۰.۷۵۵
۲۰۰	۶	۲	۰.۰۴۱	۰.۱۸۱	۰.۶۶۵	۰.۶۷۷
۲۰۰	۶	۱۰	۰.۰۶۵	۰.۱۹۰	۰.۷۳۱	۰.۷۴۰
۲۰۰	۹	۲	۰.۰۲۸	۰.۲۳۴	۰.۶۲۷	۰.۶۶۴
۲۰۰	۹	۱۰	۰.۰۵۹	۰.۲۴۱	۰.۷۰۱	۰.۷۳۲
۲۰۰	۱۲	۲	۰.۰۲۲	۰.۳۱۰	۰.۵۹۸	۰.۶۴۸
۲۰۰	۱۲	۱۰	۰.۰۵۲	۰.۳۱۷	۰.۶۹۲	۰.۷۲۵
۲۰۰۰	۶	۲	۰.۰۱۳	۰.۰۶۲	۰.۶۷۰	۰.۶۷۳
۲۰۰۰	۶	۱۰	۰.۰۳۲	۰.۰۶۴	۰.۷۳۳	۰.۷۳۵
۲۰۰۰	۹	۲	۰.۰۱۰	۰.۰۸۲	۰.۶۳۶	۰.۶۴۰
۲۰۰۰	۹	۱۰	۰.۰۳۷	۰.۰۸۵	۰.۷۱۰	۰.۷۱۴
۲۰۰۰	۱۲	۲	۰.۰۰۸	۰.۱۰۲	۰.۶۰۷	۰.۶۱۶
۲۰۰۰	۱۲	۱۰	۰.۰۳۶	۰.۱۰۵	۰.۶۹۲	۰.۶۹۹

بگیرند.

با افزایش اندازه نمونه (T) و ثابت نگه داشتن j_{\max} ، هر دو روش OLS و نیم‌رخی سازگار هستند. مربع اریبی هر دو روش با افزایش اندازه نمونه کاهش می‌یابد. برای مثال، مربع اریبی روش نیم‌رخی (U-MIDAS) از ۰.۰۸۷ به ۰.۰۳۹۱ (به ۰.۰۴۱ به ۰.۰۱۸۱) و سپس به ۰.۰۱۳ (به ۰.۰۶۲) کاهش می‌یابد، همچنان که T از ۵۰ به ۲۰۰ و سپس به ۲۰۰۰ افزایش می‌یابد. نتایج شبیه‌سازی تأیید می‌کند که روش نیم‌رخی کاراتر از برآوردگر OLS است.

تعداد وقفه‌ها نیز عامل کلیدی در عملکرد پیش‌بینی است. نتایج شبیه‌سازی نشان می‌دهد که U-MIDAS در سناریوهای با j_{\max} بزرگ و T کوچک، عملکرد ضعیف‌تری نسبت به روش نیم‌رخی دارد. در حالت $T = ۵۰$ ، $j_{\max} = ۱۲$ و $\theta = ۲$ ، MSE روش نیم‌رخی ۰.۹۸۹ است، در حالی که MSE روش U-MIDAS برابر ۱.۳۲۲ است (افزایش ۲۶٪).

مقایسه تعداد عملیات محاسباتی مورد نیاز (FLOPS) نشان می‌دهد که روش پروفایلی به طور ملموسی ساده‌تر و سریع‌تر از U-MIDAS است. با رشد اندازه نمونه و تعداد وقفه‌ها، اختلاف بار محاسباتی میان دو روش افزایش می‌یابد. این مزیت محاسباتی روش نیم‌رخی را

برای کاربردهای عملی با حجم داده بزرگ و یا مدل‌های با وقفه زیاد، گزینه‌ای ایده‌آل‌تر می‌سازد.

بحث و نتیجه‌گیری

شبیه‌سازی‌های انجام‌شده نشان داد که روش نیم‌رخ‌ی در تمامی سناریوها از نظر آریبی کمتر، قدرت پیش‌بینی بالاتر و کارایی محاسباتی برتری دارد. اگرچه U-MIDAS ممکن است درون نمونه R^2 بالاتری ارائه کند اما مشکل بیش‌برازش و افزایش MSE خارج نمونه را دارد. بنابراین برای کاربردهای عملی MIDAS به ویژه در موارد نیازمند پیش‌بینی دقیق و صرفه‌جویی محاسباتی، روش نیم‌رخ‌ی توصیه می‌گردد.

مراجع

- Andreou, E., Ghysels, E., Kourtellis, A., (2010). Regression models with mixed sampling frequencies. *Journal of Econometrics* 158(2), 246–261.
- Armesto, M.T., Engemann, K.M., Owyang, M.T., (2010). Forecasting with mixed frequencies. *Federal Reserve Bank of St. Louis Review* 92(6), 521–536.
- Clements, M.P., Galvão, A.B., (2009). Forecasting US output growth using leading indicators: An appraisal using MIDAS models. *Journal of Applied Econometrics* 24(7), 1187–1206.
- Foroni, C., Marcellino, M., Schumacher, C., (2015). Unrestricted mixed data sampling (MIDAS): MIDAS regressions with unrestricted lag polynomials. *Journal of the Royal Statistical Society: Series A* 178(1), 57–82.
- Ghysels, E., Santa-Clara, P., Valkanov, R., (2004). The MIDAS touch: Mixed data sampling regression models. *UNC Working Paper*.
- Ghysels, E., Santa-Clara, P., Valkanov, R., (2006a). Predicting volatility: Getting the most out of return data sampled at different frequencies. *Journal of Econometrics* 131(1-2), 59–95.
- Ghysels, E., Qian, H., Sargent, T.J., (2011). MIDAS regressions with MIDAS shocks. *UNC Working Paper*.
- Murphy, S.A., van der Vaart, A.W., (2000). On profile likelihood. *Journal of the American Statistical Association* 95(450), 449–465.
- Patefield, W.M., (1977). On the maximized likelihood function. *Sankhyā: The Indian Journal of Statistics, Series B* 39(1), 92–96.

A Comparison of the Profiling and Unrestricted Approaches in Fitting Regression Models Based on Mixed-Frequency Data Sampling

Farzin ranjbar, Ali Aghamohammadi¹, Majid Adib²

¹Department of Statistics, University of Zanjan

²Department of Mathematics, University of Zanjan

Abstract: This article compares the profiling and unrestricted estimation methods for parameter estimation in regression models based on mixed-frequency data (MIDAS). The profile method, which combines ordinary least squares (OLS) for estimating regression coefficients with the optimization of polynomial weighting parameters, demonstrates superior computational efficiency and lower prediction error compared to the unrestricted parameter estimation approach.

Keywords: Mixed data sampling regression, Unrestricted estimation method, Profiling estimation method.

Mathematics Subject Classification (2020): 62M10, 91B84.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



برخی نتایج مجانبی در مورد اختلاف میانگین دو جامعه در داده‌های تابعی جزئی مشاهده‌شده

پویا سلیمی فر^۱، سید محمد ابراهیم حسینی نسب^۲، امید خادم نوع^۳

^۱دانشگاه شهید بهشتی

^۲دانشگاه شهید بهشتی ^۳دانشگاه زنجان

چکیده: امروزه ثبت اطلاعات از مشاهداتی مانند X_1, X_2, \dots, X_n که به صورت تابعی پیوسته از یک پدیده باشند، مرسوم است. مشاهدات تابعی اندازه‌گیری‌هایی از پدیده‌ای است که نسبت به یک متغیر مانند زمان، بسامد و . . . پیوسته است. ماهیت پیوسته این داده‌ها منبع غنی از اطلاعات است و امکان مدلسازی با دقت بالا را فراهم می‌کند، در برخی شرایط و مطالعات تجربی عواملی مانند محدودیت زمانی در مطالعه، ابزارهای اندازه‌گیری ناکارآمد یا عوامل ناشناخته سبب می‌شوند ثبت داده‌ها در سراسر دامنه‌ای معین در جریان مطالعه به صورت کامل انجام نگیرد، در این وضعیت داده‌های تابعی به صورت ناکامل یا جزئی مشاهده‌شده^۱ در دسترس هستند و روش‌های برآوردیابی معمول که برای داده‌های تابعی کامل ارائه شده است، قابل استفاده نیست. یک رویکرد این است که با در نظر گرفتن مقدار تابع در زیر مجموعه‌های مشاهده‌شده از کل دامنه به تحلیل داده‌ها پرداخته می‌شود. در این مقاله، برخی خواص مجانبی از اختلاف میانگین دو جامعه در داده‌های تابعی جزئی مشاهده‌شده ارائه می‌شود.

واژه‌های کلیدی: داده‌های تابعی جزئی مشاهده‌شده، همگرایی در احتمال، فرآیند گاوسی، عملگر قطری.

۱ مقدمه

داده‌های تابعی جزئی مشاهده شده شامل تحقق‌هایی از توابع تصادفی هستند که در کل دامنه مشاهده نشده و هر تابع در نمونه ممکن است در زیر مجموعه‌های متفاوتی از دامنه مشاهده شود و هیچ اطلاعی درباره مقادیر تابع در سایر نقاط دامنه در دسترس نباشد. با در نظر گرفتن $\tau = [0, 1]$ برای i امین متغیر تابعی $X_i \in L^2[0, 1]$ زیر مجموعه‌ای مانند $O_i \subseteq [0, 1]$ وجود دارد به طوری که $X_i(t)$ برای $t \in O_i$ مشاهده شده است و برای $t \in \tau \setminus O_i$ مشاهده نشده است، به طور کلی الگوی گمشدگی در مطالعات تابعی

^۱ سخنران، p_salimifar@sbu.ac.ir

^۱Partially observed

نیز می‌تواند در چهار حالت نمایان شود. اولین حالت که به آن اشاره می‌شود، الگوی گمشدگی کاملاً تصادفی^۲ است که در آن گمشدگی بخش‌هایی از تابع، مستقل از تابع (به این معنا که در آن ساختار گمشدگی نه به بخش‌های مشاهده شده و نیز به بخش‌های گمشده از تابع وابسته نیست) است. دومین حالت، الگوی گمشدگی تصادفی نام دارد که در آن گمشدگی تنها با بخش مشاهده شده از تابع مرتبط است و سومین الگوی گمشدگی، گمشدگی غیر تصادفی نام دارد که به هر دو بخش مشاهده شده و گمشده از تابع وابسته خواهد بود. الگوی چهارم که الگوی گمشدگی سیستماتیک^۳ است از سه وضعیت موجود متفاوت است و در آن اغلب یک عامل شناخته شده در جریان مطالعه سبب گمشدگی بخشی از تابع در دامنه می‌شود، همچنین یادآور می‌شویم که تاکنون الگوی گمشدگی تصادفی و غیر تصادفی به دلیل فرض‌های بسیار محدود کننده بر الگوی داده‌ها در پژوهش‌های مربوط به داده‌های تابعی مورد مطالعه قرار نگرفته‌اند. مسائل کاربردی مرتبط با این حالت باعث ایجاد یک مجموعه کارهای تحقیقاتی با جنبه‌های مختلف در این مسئله شده است، باگنی (۲۰۱۲)، به آزمون مشخصه‌سازی برای داده‌های تابعی ناکامل (که به ارزیابی اعتبار فرضیات احتمالی تشکیل دهنده یک مدل آماری اشاره دارد) پرداخت، کراووس (۲۰۱۵)، داده‌های پایش فشار خون سرپایی و نمایه ضربان قلب شامل الگو تغییرات ضربان قلب در حالات مختلف فعالیت، استراحت و یا شرایط خاص پزشکی را مورد مطالعه قرار داد و با استفاده از قسمت مشاهده شده از منحنی ضربان قلب برای قسمت گمشده از آن یک پیشگو ارائه کرد، سپس برای قسمت گمشده از منحنی یک باند پیشگویی ساخت. دلاگل و هال (۲۰۱۶)، داده‌های تابعی ناکامل را با استفاده از بخش‌هایی از یک مدل زنجیر مارکوف تقریب زدند. گرومنکو و همکاران (۲۰۱۷)، در یک مطالعه به بررسی روند سرمایش در لایه یونسفر با استفاده از رگرسیون تابعی با منحنی‌های ناکامل پرداختند. شیبانی و همکاران (۲۰۱۸)، به طبقه‌بندی با متغیرهای کمکی تابعی ناقص پرداختند. استفسنکی و همکاران (۲۰۱۸)، به تفکیک داده‌های تابعی جزئی مشاهده شده بر اساس تحلیل مولفه‌های اصلی پرداختند. دیسکری و همکاران (۲۰۱۹)، به بازسازی کوواریانس از مشاهدات تابعی ناکامل پرداختند. لیل و رامسدر (۲۰۱۹)، در یک مطالعه که هدف آن، برآورد منحنی میانگین قیمت مصرف انرژی در بازار ذخیره کنترلی برق در آلمان بود، تمرکز کردند و با توجه به این که در این بازار مزایده‌ای، منحنی‌های قیمت تنها به صورت جزئی در انتها دامنه مشاهده می‌شدند مکانیسم گمشدگی به صورت سیستماتیک، وابسته به استراتژی‌های معاملاتی (عامل شناخته شده) بود و این عامل سبب نقض فرض گمشدگی کاملاً تصادفی می‌شد و نتایج ناسازگاری را در برآوردگر میانگین ایجاد می‌کرد از این رو آنها از طریق قضیه اساسی حساب دیفرانسیل و انتگرال و مشتق پذیری داده‌های تابعی برآوردهای سازگاری برای تابع میانگین و عملگر کوواریانس در این وضعیت ارائه دادند. یان کیو و همکاران (۲۰۲۱)، به مطالعه‌ی آزمون فرض دو نمونه‌ای برای میانگین پرداختند. نایپ و لیل (۲۰۲۰)، یک روش بازسازی برای قسمت گمشده از منحنی‌های جزئی مشاهده شده ارائه دادند. از دیگر مطالعات می‌توان به کتب منتشر شده توسط رویین (۱۹۷۶)، هی و همکاران (۲۰۲۲)، اشاره کرد.

با توجه به این که استنباط در مورد میانگین داده‌های تابعی جزئی مشاهده‌شده از اهمیت بالایی برخوردار است و در بسیاری از مطالعات شامل آزمون فرض مورد استفاده قرار می‌گیرد در این پژوهش نیز به مطالعه‌ی برخی خواص مجانبی در مورد اختلاف میانگین یک متغیر تابعی در دو جامعه پرداخته می‌شود. نتایج این پژوهش می‌تواند در آزمون فرض برای برابری میانگین دو جامعه و در حالت تعمیم یافته آن به چند جامعه مورد استفاده قرار گیرد.

^۲Missing-Completely-At-Random

^۳Systematic Missingness

۲ برخی نتایج حدی

در مطالعه‌ی داده‌های تابعی جزئی باید دو عامل مهم در نظر گرفته شود، نخست این که وزن‌های مربوط به شانس گمشدگی در هر گروه متفاوت است و باید در ساختار اختلاف میانگین‌ها اعمال شود، دوم این که ممکن است در اکثر مطالعات واقعی برای هر گروه تابع کوواریانس متمایزی وجود داشته باشد که هر دو این عوامل باید در ساختار مسئله اعمال شود. اکنون به تعریف برخی نماد مورد استفاده در مقاله می‌پردازیم:

$O_{ij}(t)$ معرف تابع نشانگر در لحظه t برای گروه j است. از این رو $N_j(t) = \sum_{i=1}^n O_{ij}(t)$ مجموع مشاهدات در دسترس برای $j = 1, 2$ است که در آن اندیس j نشانگر تعداد گروه‌ها است و برای هر گروه مشاهدات تابعی، می‌توانند متفاوت باشند. اگر $N_j(t) = 0$ ، آنگاه $J_j(t) = I_{[N_j(t) > 0]} = 0$ ، به این معنا است که مقدار مشاهده شده‌ای از تابع برای گروه j ام در لحظه‌ی t وجود ندارد.

$$\begin{cases} \hat{\rho}(s, t) = \frac{I(s, t)}{M(s, t)} \sum_{i=1}^n U_i(s, t) \{X_i(s) - \hat{\mu}_{st}(s)\} \{X_i(t) - \hat{\mu}_{st}(t)\}, \\ I(s, t) = I_{[M(s, t) > 0]}, \\ \hat{\mu}_{st}(s) = \frac{I_{[M(s, t) > 0]}}{M(s, t)} \sum_{i=1}^n U_i(s, t) X_i(s). \end{cases}$$

$U_i(s, t) \equiv O_i(s)O_i(t)$ معرف تابع نشانگر است و به این معناست که در دو لحظه‌ی معین s و t مقدار تابع به صورت توأم در دسترس باشد، از این رو $M(s, t) \equiv \sum_{i=1}^n U_i(s, t)$ معرف مجموع مشاهدات توأم در هر گروه است و در صورتی که $I_{[M(s, t) > 0]} = 0$ (به این معنا که هیچ مقادیری از تابع در نقاط جفتی مشاهده نشود)، برآوردگر $\hat{\mu}_{st} = 0$ در نظر گرفته می‌شود، در ادامه برآورد میانگین گروهی و میانگین کل به صورت زیر است:

$$\sum_{j=1}^2 \hat{\omega}_j(t) = 1, \quad N(t) = \sum_{i=1}^n O_i(t), \quad \hat{\mu}(t) = \sum_{j=1}^2 \hat{\omega}_j(t) \hat{\mu}_j(t).$$

که در آن برآورد میانگین کل به صورت ترکیب خطی از برآورد میانگین‌های گروهی است و هر گروه با ضریب تاثیر $\hat{\omega}_i(t)$ در میانگین کل حضور دارند بعلاوه میانگین برای هر گروه و تابع وزن برای j -امین گروه به صورت زیر است:

$$\hat{\mu}_j(t) = J_j(t) N_j(t)^{-1} \sum_{i=1}^{n_j} O_{ji}(t) X_{ji}(t), \quad j \in \{1, 2\}, \quad \hat{\omega}_j(t) = \frac{\frac{N_j(t)}{\hat{r}_j^2}}{\sum_{k=1}^K \frac{N_k(t)}{\hat{r}_k^2}}.$$

بعلاوه احتمال این که در لحظه‌ی t مقدار تابع مشاهده شود و نیز احتمال این که در دو لحظه معین به صورت توأم مقادیر تابع مشاهده

شود به صورت زیر تعریف می‌شوند:

$$\pi_i(t) \equiv EO_i(t) = \Pr\{O_i(t) = 1\}, \quad \nu_i(s, t) = E[U_i(s, t)] = P\{O_i(t) = 1, O_i(s) = 1\}.$$

و ضریب استاندارد ساز در $\hat{\omega}_j(t)$ نیز به صورت $\hat{r}_j = \int_0^1 \hat{\rho}_j(t, t) dt$ تعریف می‌شود.

نخست تعریف می‌کنیم:

$$SSR_n(t) = \sum_{i=1}^2 n_i (\bar{y}_{i\cdot}(t) - \bar{y}_{\cdot\cdot}(t))^2,$$

که همان تغییرات بین گروهی است که در تحلیل واریانس مورد استفاده قرار می‌گیرد. مشابه تعریف فوق مولفه‌ی تصادفی $Z_{n_j}(t)$ نیز مشابه جملات مجموع بالا و با تعدیل آن توسط دو عامل ذکر شده برای داده‌های تابعی جزئی مشاهده‌شده به صورت زیر ارائه می‌شود:

$$Z_{n_j}(t) = \hat{r}_j^{-1} N_j(t)^{1/2} [\hat{\mu}_j(t) - \hat{\mu}(t)].$$

اکنون برای سهولت در اعمال دو مولفه‌ی مهم یعنی، ساختار گمشدگی نامتقارن و تابع کوواریانس متمایز و دوری از پیچیدگی‌های غیر ضروری با تعریف یک عملگر قطری می‌توان $Z_{n_j}(t)$ را به صورت ضرب نقطه‌ای از عناصر $\mathcal{L}^2 \rightarrow \mathcal{L}^2$: $f: \mathcal{L}^2 \rightarrow \mathcal{L}^2$ ایجاد کرد که در آن برآورد عملگر قطری یعنی $\hat{\psi}_{jl}$ بر روی هر عضو دلخواه f_l متعلق به فضای L^2 به صورت زیر عمل می‌کند:

$$(\hat{\psi}_{jl} f_l)(t) = \hat{r}_j^{-1} \left\{ \delta_{jl} - N_j(t)^{1/2} \hat{\omega}_l(t) J_l(t) N_l(t)^{-1/2} \right\} f_l(t)$$

بعلاوه تعریف می‌کنیم:

$$H_j(t) = N_j(t)^{1/2} [\hat{\mu}_j(t) - \mu(t)], \quad j = 1, 2$$

از این رو می‌توان $Z_{n_j}(t)$ را برای هر جامعه به صورت زیر بازنویسی کرد،

$$Z_{n_j}(t) = \hat{\psi}_{j1} H_1(t) + \hat{\psi}_{j2} H_2(t), \quad j = 1, 2 \quad (1.2)$$

اکنون عناصر $Z_{n_j}(t)$ به صورت زیر تعریف می‌شوند:

$$Z_{n_j}(t) = \hat{r}_j^{-1} N_j(t)^{1/2} [\hat{\mu}_j(t) - \mu(t)] = \sum_{i=1}^2 \hat{\psi}_{ji} H_i(t),$$

از این رو $Z_{n_j}(t)$ با اعمال عملگر قطری $\hat{\psi}_{jl}$ روی $H_l(t)$ به صورت زیر نمایش داده می‌شود:

$$(\hat{\psi}_{jl} H_l)(t) = \hat{r}_j^{-1} \left\{ \delta_{jl} - N_j(t)^{1/2} \hat{\omega}_l(t) J_l(t) N_l(t)^{-1/2} \right\} N_l(t)^{1/2} [\hat{\mu}_l(t) - \mu(t)], \quad l = 1, 2$$

برای حالتی که $J_j(t) = \mathbb{I}_{\{N_j(t) > 0\}} = 1$ ، به این معنا که در هر گروه مشاهده داشته باشیم، $Z_{n_j}(t)$ به صورت زیر ساده می‌شود:

$$Z_{n_j}(t) = \sum_{l=1}^2 \hat{r}_j^{-1} \left\{ \delta_{jl} N_l(t)^{1/2} - N_j(t)^{1/2} \hat{\omega}_l(t) \right\} [\hat{\mu}_l(t) - \mu(t)],$$

اما با استفاده از تعریف دلتای کرونکر، داریم:

$$\sum_{k=1}^2 \hat{r}_j^{-1} \delta_{jl} N_l(t)^{1/2} [\hat{\mu}_l(t) - \mu(t)] = \hat{r}_j^{-1} N_j(t)^{1/2} [\hat{\mu}_j(t) - \mu(t)],$$

در ادامه مولفه‌ی $Z_{n_j}(t)$ به صورت مجموع دو مولفه‌ی تصادفی تجزیه می‌شود و به همگرایی هر دو جز به صورت کلی اشاره می‌شود برای این منظور دو بخش تصادفی را بر اساس انتقال میانگین‌ها بازنویسی کرده و به همگرایی دو جز تصادفی به صورت کلی اشاره می‌کنیم.

با فرض این که میانگین‌های هر گروه به صورت $\mu_j(t) = \mu(t) + h_j(t)$ نمایش داده شوند، $Z_{n_j}(t)$ بر اساس این انتقال به اندازه‌ی هر تابع حقیقی مقداری مانند $h_j(t)$ به صورت زیر بازنویسی می‌شود:

$$Z_{n_j}(t) = \sum_{l=1}^2 \hat{r}_j^{-1} \left\{ \delta_{jl} - N_j(t)^{1/2} \hat{\omega}_l(t) N_l(t)^{-1/2} \right\} N_l(t)^{1/2} [h_l(t) + (\hat{\mu}_l(t) - \mu_l(t))]. \quad (2.2)$$

همگرایی مولفه‌های تجزیه شده $Z_{n_j}(t)$

با توجه به تعریف $Z_{n_j}(t)$ در رابطه (۲.۲) می‌توان نوشت: $Z_{n_j}(t) = \Delta_{n_j}(t) + S_{n_j}(t)$ که در آن

$$\Delta_{n_j}(t) = \sum_{l=1}^2 \hat{r}_j^{-1} \left\{ \delta_{jl} - N_j(t)^{1/2} \hat{\omega}_l(t) N_l(t)^{-1/2} \right\} N_l(t)^{1/2} h_l(t), \quad (3.2)$$

$$S_{n_j}(t) = \sum_{l=1}^2 \hat{r}_j^{-1} \left\{ \delta_{jl} - N_j(t)^{1/2} \hat{\omega}_l(t) N_l(t)^{-1/2} \right\} N_l(t)^{1/2} [\hat{\mu}_l(t) - \mu_l(t)], \quad (4.2)$$

قضیه ۱.۲. با قرار دادن $h_j(t) = \frac{d_j(t)}{n_j^{1/2}}$ ، $j = 1, 2$ که در آن $d_j(t)$ هر تابع حقیقی مقدار است، می‌توان نشان داد که همگرایی در توزیع برای $\Delta_{n_j}(t)$ و $S_{n_j}(t)$ به صورت زیر نتیجه می‌شود:

$$\Delta_{n_j}(t) = \hat{r}_j^{-1} N_j(t)^{1/2} \left\{ n_j^{-1/2} d_j(t) - \sum_{l=1}^2 \hat{\omega}_l(t) n_l^{-1/2} d_l(t) \right\} \xrightarrow{d} r_j^{-1} \pi_j(t)^{1/2} \left\{ d_j(t) - a_j^{1/2} \sum_{l=1}^2 \omega_l(t) a_l^{-1/2} d_l(t) \right\}, \quad (5.2)$$

$$S_{n_j}(t) \xrightarrow{d} r_j^{-1} \left[S_j^\infty(t) - \sum_{l=1}^2 \omega_l(t) \left[\frac{\pi_j(t) a_j}{\pi_l(t) a_l} \right]^{1/2} S_l^\infty(t) \right], \quad (6.2)$$

که در آن

$$a_j = \lim_{n \rightarrow \infty} \frac{n_j}{n}.$$

نسبت حدی اندازه نمونه هر گروه به اندازه نمونه کل گروه‌ها است، از این رو میانگین مجانبی فرآیند گاوسی با انتقال میانگین همان توزیع حدی بخش غیر تصادفی در (۵.۲) است و تابع کوواریانس آن مرتبط با بخش تصادفی در (۶.۲) است. که در آن $S_j^\infty(t) \stackrel{d}{=} GP(\circ, \mathcal{K}_j)$ به صورت فرآیند گاوسی با میانگین صفر و هسته عملگر کوواریانس زیر است:

$$\kappa_j(s, t) = [\pi_j(t)]^{-1/2} [\pi_j(s)]^{-1/2} \nu_j(s, t) \rho_j(s, t).$$

در ادامه هسته مربوط به توزیع مجانبی عملگر $Z_{n_j}(t)$ به صورت زیر ارائه می‌شود:

$$\begin{aligned}
\nu_{jk}^{\infty}(s, t) &= \text{cov} [Z_j^{\infty}(t), Z_k^{\infty}(s)] \\
&= \sum_{l=1}^p \psi_{jl}(s) \kappa_l(s, t) \psi_{kl}^*(t) \\
&= \sum_{l=1}^p r_j^{-1} \left\{ \delta_{jl} - \pi_j(s)^{1/2} a_j^{1/2} \omega_l(s) \pi_l(s)^{-1/2} a_l^{-1/2} \right\} \\
&\quad \times \kappa_l(s, t) \left\{ \delta_{kl} - \pi_k(t)^{1/2} a_k^{1/2} \omega_l(t) \pi_l(t)^{-1/2} a_l^{-1/2} \right\} r_k^{-1},
\end{aligned}$$

که در آن $\circ \xrightarrow{P} \psi - \hat{\psi}$ ، $\|\psi - \hat{\psi}\|_{\infty}$ ، $\psi_{kl}^*(t)$ به صورت عملگر الحاقی تعریف می‌شود، از این رو توزیع مجانبی $Z_{n_j}(t)$ به صورت زیر است:

$$Z_{n_j}(t) \xrightarrow{d} \mathcal{GP}(\Delta_j^{\infty}, \nu_{jj}^{\infty}),$$

که یک فرآیند گاوسی با میانگین مجانبی زیر است:

$$\Delta_j^{\infty}(t) = r_j^{-1} \pi_j(t)^{1/2} \left\{ d_j(t) - a_j^{1/2} \sum_{l=1}^p \omega_l(t) a_l^{-1/2} d_l(t) \right\}.$$

قضیه (۱.۲) به عنوان یک نتیجه مهم در تقریب توزیع آماری آزمون مورد استفاده قرار می‌گیرد.

مراجع

- [1] Bugni, F. A., Smith, J., and Lee, K. (2012). Specification test for missing functional data, *Econometric Theory*, 28(5), 959-1002.
- [2] Gromenko, O., Kokoszka, P., and Sojka, J. (2017). Evaluation of the cooling trend in the ionosphere using functional regression with incomplete curves, *Ann. Appl. Stat.*, 11(2): 898-918.
- [3] Liebl, D., Rameseder, S. (2019). Partially observed functional data: The case of systematically missing parts, *Computational Statistics Data Analysis*, 131, 104-115.
- [4] Stefanucci, M., Sangalli, L. M., and Brutti, P. (2018). PCA-based discrimination of partially observed functional data, with an application to AneuRisk65 data set, *Statistica Neerlandica*, 72(3), 246-264.
- [5] Mojirsheibani, M., Shaw, C. (2018). Classification with incomplete functional covariates. *Probability Letters*, 139, 40-46.
- [6] Delaigle, A., Hall, P. (2016). Approximating fragmented functional data by segments of Markov chains, *Biometrika*, 103(4), 779-799.
- [7] Qiu, Z., Chen, J., and Zhang, J. T. (2021). Two-sample tests for multivariate functional data with applications, *Computational Statistics Data Analysis*, 157, 107160.

- [8] Kraus, D. (2015). Components and completion of partially observed functional data, *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 77(4), 777-801.
- [9] Kneip, A., Liebl, D. (2020). On the optimal reconstruction of partially observed functional data, *The Annals of Statistics*, 1692-1717.
- [10] Descary, M. H., Panaretos, V. M. (2019). Recovering covariance from functional fragments, *Biometrika* , 106(1), 145-160.
- [11] Rubin, D. B. (1976). *Inference and missing data*. *Biometrika*, 63(3), 581-592.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30-31, 2025
University of Kurdistan



مدل خودبازگشتی صحیح مقدار مرتبه اول فصلی با توزیع حاشیه‌ای دلاپورت

مریم شالباف^۱، غلامعلی پرهام^۲

گروه آمار و ریاضی، واحد شوشتر، دانشگاه آزاد اسلامی، شوشتر، ایران

گروه آمار، دانشکده علوم ریاضی و کامپیوتر، دانشگاه شهید چمران اهواز

چکیده: در سری‌های زمانی شمارشی اغلب پدیده بیش پراکنشی مشاهده می‌شود. در این مقاله، یک فرایند خودبازگشتی صحیح مقدار مرتبه اول با ساختار فصلی معرفی می‌شود. توزیع حاشیه‌ای مدل توزیع دلاپورت می‌باشد. بر این اساس توزیع نوآوری‌ها پیچشی از توزیع پواسون با تعداد α تا توزیع هندسی تعدیل‌یافته در صفر خواهد شد. برخی از خواص مدل نمایش داده می‌شود. برای برآورد پارامترها از روش‌های یول-والکر، حداقل مربعات شرطی و حداکثر درستنمایی شرطی استفاده می‌شود. با انجام یک شبیه‌سازی عملکرد این برآوردگرها مورد ارزیابی قرار می‌گیرد. در نهایت این مدل به یک مجموعه داده واقعی برازش داده می‌شود. **واژه‌های کلیدی:** توزیع دلاپورت، مدل خودبازگشتی صحیح مقدار مرتبه اول فصلی، α تا توزیع هندسی تعدیل‌یافته، عملگر رقیق‌کننده دوجمله‌ای

۱ مقدمه

داده‌های سری زمانی با ویژگی‌های فصلی را می‌توان در زمینه‌های مختلفی مانند علوم آماری، بهداشت، اقتصاد، محیط زیست و غیره یافت. آن‌ها عمدتاً یک الگوی فصلی با دوره‌هایی را نشان می‌دهند که پس از یک بازه زمانی منظم تکرار می‌شوند. کوچکترین دوره زمانی برای این رویداد، دوره فصلی نامیده می‌شود. عوامل متعددی مانند آب‌وهوا و ویژگی‌های ذاتی می‌توانند باعث ایجاد ساختارهای فصلی شوند. مدل‌های سری زمانی گسسته مقدار (INAR) معمولاً براساس عملگر رقیق‌کننده دوجمله‌ای که توسط استیوتل و ون هارن (۱۹۷۹) به صورت $\rho \circ X = \sum_{i=1}^X Y_i$ تعریف شده است، ساخته می‌شوند، به طوری که $\{Y_t\}_{t \in \mathbb{Z}}$ یک دنباله از متغیرهای تصادفی مستقل و هم‌توزیع برنولی با احتمال موفقیت $\rho \in [0, 1]$ می‌باشد و X یک متغیر تصادفی صحیح مقدار غیرمنفی و مستقل از Y می‌باشد. **ملکنزی (۱۹۸۵)** و **آلوش و آلازید (۱۹۸۷)** مدل INAR(1) را به صورت زیر معرفی کردند:

$$X_t = \rho \circ X_{t-1} + \varepsilon_t \quad (1.1)$$

که در آن $0 \leq \rho < 1$ و ε_t دنباله‌ای از متغیرهای تصادفی با مقادیر صحیح غیرمنفی ناهمبسته با میانگین μ_ε و واریانس متناهی σ_ε^2 است. اگر $\varphi_X(s)$ و $\varphi_\varepsilon(s)$ به ترتیب نشان‌دهنده تابع مولد احتمال $\{X_t\}_{t \in \mathbb{Z}}$ و $\{\varepsilon_t\}_{t \in \mathbb{Z}}$ باشند، آنگاه توزیع حاشیه‌ای $\{X_t\}_{t \in \mathbb{Z}}$ را می‌توان به صورت $\varphi_X(s) = \varphi_X(1 - \rho + \rho s)\varphi_\varepsilon(s)$ نوشت.

متداول‌ترین فرایند برای مدل‌بندی این نوع از داده‌ها، مدل پواسون خودبازگشتی گسسته مقدار مرتبه اول (PINAR(1)) است. بیش پراکنشی (یعنی حالتی که واریانس از میانگین بزرگتر است) پدیده‌ای است که معمولاً در سری‌های شمارشی مشاهده می‌شود. چون توزیع پواسون هم پراکنش است (یعنی میانگین با واریانس برابر است) فرایند PINAR برای مدل‌سازی داده‌های با بیش پراکنشی مناسب نیست. یک رویکرد معمول برای برخورد با بیش پراکنشی در داده‌های شمارشی استفاده از توزیع‌های پواسون آمیخته است که توسط **بارتو-سوازا (۲۰۱۹)** بررسی شده است. توزیع‌های پواسون آمیخته بوسیله معرفی یک اثر تصادفی پنهان روی میانگین توزیع پواسون به دست می‌آیند. مدل خود بازگشتی صحیح مقدار با توزیع حاشیه‌ای دلاپورت توسط **شالباف و همکاران (۲۰۲۲)** بررسی شده است. در سال‌های اخیر بر روی مدل‌های خود بازگشتی صحیح مقدار مطالعات زیادی صورت گرفته است ولی فرایندهای فصلی با توزیع حاشیه‌ای گسسته به ندرت بررسی و معرفی شده‌اند. **بورگینکون و همکاران (۲۰۱۶)** فصلی بودن را براساس مدل PINAR(1) بررسی کرده‌اند. فرایند $INAR(1)_s$ فصلی مرتبه اول با دوره فصلی s ساختار ساده‌ای دارد و برای برآورد کردن پارامترهای کمتری دارد اما برای مدل‌سازی سری‌های زمانی فصلی با بیش پراکنشی مناسب نیست. **تیان و همکاران (۲۰۱۸)** یک فرایند $INAR(1)$ فصلی را براساس عملگر رقیق‌کننده دوجمله‌ای منفی و با توزیع حاشیه‌ای هندسی معرفی کردند. در این مقاله ما در جهت کمک به تحلیل سری‌های زمانی گسسته مقدار، یک فرایند خودبازگشتی صحیح مقدار مرتبه اول فصلی با توزیع حاشیه‌ای دلاپورت ($DELINAR(1)_s$) را معرفی می‌کنیم. این مدل برای مدل‌سازی سری‌های زمانی فصلی با مقدار صحیح غیرمنفی و با بیش پراکنشی مناسب است. بر این اساس مدل در بخش ۲ معرفی شده است. همچنین برخی از ویژگی‌های مدل بیان شده است. در بخش ۳، روش‌های برآورد پارامترهای مدل مورد بحث قرار می‌گیرد. در بخش ۴ برخی از نتایج پیش‌بینی ارائه شده است. بخش ۵ برخی از نتایج شبیه‌سازی را برای روش‌های برآورد نشان می‌دهد. در بخش ۶، مدل بر روی یک مجموعه داده واقعی اعمال می‌شود. در نهایت در بخش ۷ نتایج بیان شده است.

۲ فرایند $INAR(1)$ فصلی با توزیع حاشیه‌ای دلاپورت

در این بخش، ویژگی‌های ساختاری این فرایند، مانند توزیع‌های حاشیه‌ای و نوآوری، میانگین و واریانس این توزیع‌ها، تابع خودکواریانس، امید شرطی و واریانس شرطی متغیر تصادفی حاشیه‌ای و احتمالات انتقال را بررسی می‌کنیم.

۱.۲ توزیع دلاپورت

متغیر تصادفی X دارای توزیع دلاپورت با پارامترهای λ ، α و β است و آن را با نماد $DEL(\lambda, \alpha, \beta)$ نشان می‌دهند اگر تابع مولد احتمال آن به صورت $G_X(s) = e^{-\lambda(1-s)} \left(\frac{1}{1+\beta(1-s)} \right)^\alpha$ باشد، به طوری که $|s| \leq 1$ و $\lambda > 0$ و $\alpha, \beta > 0$ (**دلاپورت، ۱۹۵۹**). با دوبار مشتق گرفتن از تابع مولد احتمال توزیع دلاپورت می‌توان نشان داد که:

$$\mu = E(X_t) = \lambda + \alpha\beta \quad , \quad \sigma^2 = Var(X_t) = \lambda + \alpha\beta(1 + \beta)$$

با توجه به اینکه واریانس بزرگ‌تر از میانگین است، می‌توان گفت این توزیع بیش پراکنشی دارد.

۲.۲ فرایند خودبازگشتی صحیح مقدار مرتبه اول فصلی با توزیع حاشیه‌ای دلاپورت $DELINAR(1)_s$

فرض کنید N نشان‌دهنده مجموعه مقادیر صحیح غیرمنفی باشد. همچنین فرض کنید X یک متغیر تصادفی با مقادیر صحیح غیرمنفی و $\rho \in [0, 1]$ باشد. با استفاده از عملگر رقیق‌کننده دو جمله‌ای، مدل دلاپورت فصلی مرتبه اول به صورت

$$X_t = \rho \circ X_{t-s} + \varepsilon_t \quad (1.2)$$

تعریف می‌شود، به طوری که $\rho \in [0, 1]$ است و $\{\varepsilon_t\}$ دنباله نوآوری از متغیرهای تصادفی با مقادیر صحیح غیرمنفی و مستقل از مقادیر گذشته $\{X_t\}$ می‌باشند و $s \in \mathbb{N}$ دوره فصلی را نشان می‌دهد. متغیرهای برنولی در عملگرهای رقیق‌کننده مستقل از هم و مستقل از نوآوری $\{\varepsilon_t\}$ می‌باشند. این مدل را با نماد $DELINAR(1)_s$ نشان می‌دهیم.

قضیه ۱.۲. متغیر تصادفی ε_t را می‌توان به صورت $\varepsilon_t = Y_1 + Y_2$ نشان داد که در آن Y_1 دارای توزیع پواسون با پارامتر $\lambda(1-\rho)$ و Y_2 پیچش α تا توزیع هندسی تعدیل‌یافته در صفر است. بنابراین امید ریاضی و واریانس ε_t به صورت $\mu_\varepsilon := E(\varepsilon_t) = (1-\rho)\mu_X$ و $\sigma_\varepsilon^2 := Var(\varepsilon_t) = (1-\rho^2)\sigma_X^2 - \rho(1-\rho)\mu_X$ می‌باشند. همچنین احتمال شرطی X_t بشرط X_{t-1} به صورت زیر به دست می‌آید.

$$p_{ij} = P(X_t = j | X_{t-1} = i) = P(\rho \circ X_{t-1} + \varepsilon_t = x_t | X_{t-1} = x_{t-1}).$$

اکنون فرض کنید

$$B_{X_{t-1}}^\rho := \rho \circ X_{t-1} | X_{t-1} \sim Binomial(X_{t-1}, \rho)$$

همچنین ε_t مستقل از X_{t-1} می‌باشد و دارای ساختار پیچشی از دو متغیر تصادفی $P_{(\lambda, \rho)} \sim Po(\lambda(1-\rho))$ و پیچش α تا توزیع هندسی تعدیل‌یافته در صفر با پارامترهای $p = \frac{1}{1+\beta}$ و $p_0 = \frac{(1+\rho\beta)}{1+\beta}$ می‌باشد، که p_0 میزان جرم احتمال در صفر می‌باشد. بنابراین احتمال شرطی را می‌توان به صورت زیر نوشت:

$$\begin{aligned} P(X_t = x_t | X_{t-1} = x_{t-1}) &= P(B_{X_{t-1}}^\rho + P_{(\lambda, \rho)} + \alpha FZMG(p, p_0) = x_t) \\ &= \sum_{s=0}^{x_t} P(\alpha FZMG = s) \sum_{d=0}^{\min\{x_t-s, x_{t-1}\}} P(B_{X_{t-1}}^\rho = d) P(P_{(\lambda, \rho)} = x_t - s - d), \end{aligned} \quad (2.2)$$

۳ برآورد پارامترها

فرض کنید X_1, X_2, \dots, X_n یک نمونه از فرایند $DELINAR(1)_s$ باشند. در این فرایند ما ۴ پارامتر داریم، فرض می‌کنیم پارامتر α معلوم است پس ۳ پارامتر دیگر باید برآورد شوند. در ادامه روش‌های یول والکر، حداقل مربعات شرطی و حداکثر درستنمایی شرطی و برخی از نتایج تحلیلی و مجانبی برآوردگرها را معرفی و توضیح می‌دهیم.

۱.۳ برآورد یول والکر

فرض کنید $\gamma(k) = \sum_{t=1}^{n-k} (X_t - \bar{X})(X_{t+k} - \bar{X})$ و $0 \leq k < n$ تابع خودکواریانس نمونه‌ای باشد. چون $\rho = \frac{\gamma(1)}{\gamma(0)}$ و $E(X_t) = \lambda + \alpha\beta(1 + \beta)$ و $Var(X_t) = \lambda + \alpha\beta(1 + \beta)$ ، برآوردهای یول والکر پارامترها به صورت زیر می‌باشد:

$$\hat{\rho}_{YW} = \frac{\sum_{t=s+1}^n (X_t - \bar{X})(X_{t-s} - \bar{X})}{\sum_{t=1}^n (X_t - \bar{X})^2}, \quad \hat{\beta}_{YW} = \sqrt{\frac{S^2 - \bar{X}}{\alpha}}, \quad \hat{\lambda}_{YW} = \bar{X} - \alpha\sqrt{\frac{S^2 - \bar{X}}{\alpha}},$$

به طوری که \bar{X} و S^2 به ترتیب نشان دهنده میانگین نمونه‌ای و واریانس نمونه‌ای می‌باشد.

۲.۳ برآورد حداقل مربعات شرطی

روش حداقل مربعات شرطی را در دو مرحله اجرا می‌کنیم (کلیمکو و نلسن، ۱۹۸۷). در گام اول با نوشتن معادله حداقل مربعات شرطی برای پارامترهای ρ و μ و به کمک مشتق گیری جزئی نسبت به پارامترهای ρ و μ از معادله و کمینه کردن آن داریم:

$$\hat{\mu}_{cls} = \frac{\sum_{t=s+1}^n X_t - \hat{\rho}_{cls} \sum_{t=s+1}^n X_{t-s}}{(n-s)(1 - \hat{\rho}_{cls})}, \quad \hat{\rho}_{cls} = \frac{(n-s) \sum_{t=s+1}^n X_t X_{t-s} - \sum_{t=s+1}^n X_t \sum_{t=s+1}^n X_{t-s}}{(n-s) \sum_{t=s+1}^n X_{t-s}^2 - (\sum_{t=s+1}^n X_{t-s})^2}.$$

در گام دوم، با نوشتن معادله حداقل مربعات شرطی برای پارامتر σ^2 و کمینه کردن آن نسبت به σ^2 داریم:

$$\hat{\sigma}_{cls}^2 = \frac{\sum_{t=s+1}^n (X_t - \hat{\rho}_{cls} X_{t-s} - (1 - \hat{\rho}_{cls}) \hat{\mu}_{cls})^2 - \hat{\rho}_{cls} (1 - \hat{\rho}_{cls}) \sum_{t=s+1}^n (X_{t-s} - \hat{\mu}_{cls})}{(n-s)(1 - \hat{\rho}_{cls})}$$

سرانجام با حل دستگاه و فرض معلوم بودن α ، برآورد کمترین مربعات شرطی برای پارامترهای λ و β به صورت زیر به دست می‌آید:

$$\hat{\lambda}_{cls} = \hat{\mu}_{cls} - \alpha \hat{\beta}_{cls}, \quad \hat{\beta}_{cls} = \sqrt{\frac{\hat{\sigma}_{cls}^2 - \hat{\mu}_{cls}}{\alpha}}.$$

۳.۳ برآورد حداکثر درستنمایی شرطی

برآوردهای حداکثر درستنمایی شرطی با بیشینه کردن تابع لگاریتم درستنمایی به دست می‌آیند. تابع لگاریتم درستنمایی شرطی به صورت $CL(\rho, \lambda, \beta) = \sum_{t=s+1}^n \log P(X_t = j | X_{t-s} = i)$ می‌باشد، به طوری که $P(X_t = j | X_{t-s} = i)$ در (۲.۲) تعریف شده است. مقادیر برآورد حداکثر درستنمایی شرطی با بیشینه کردن تابع فوق بدست می‌آید. چون $CL(\rho, \lambda, \beta)$ یک تابع غیرخطی است، برآوردهای حداکثر درستنمایی شرطی باید به کمک روش‌های عددی محاسبه شوند.

۴ پیش‌بینی

در اینجا هدف، محاسبه پیش‌بینی مقدار X_{n+h} ($h \in \mathbb{N}$) براساس اطلاعات تا زمان n ام می‌باشد. توزیع X_{n+h} براساس تعریف مدل $DELINAR(1)_s$ را می‌توان به صورت زیر بیان کرد:

$$X_{n+h} \stackrel{d}{=} \rho^q \circ X_{n+h-qs} + \sum_{j=0}^{q-1} \rho^j \circ \varepsilon_{n+h-js}$$

به طوری که $q = \lceil h/s \rceil$ و $\lceil x \rceil = \min \{n \in \mathbb{N} | x \leq n\}$. در واقع $\lceil x \rceil$ نشان دهنده کوچک‌ترین عدد طبیعی بزرگ‌تر (یا مساوی) x می‌باشد. از عبارت بالا مشاهده می‌شود که توزیع پیش‌بینی h گام جلوتر شکل بسیار پیچیده‌ای دارد و بنابراین ما پیش‌بینی h گام

جلوتر را به کمک امید شرطی زیر بدست می‌آوریم:

$$\hat{X}_{n+h}|X_n = E(X_{n+h}|X_n) = \rho^q(X_{n+h-q} - \mu_x) + \mu_x.$$

به‌طوری که $\mu_x = \lambda + \alpha\beta$.

۵ شبیه‌سازی

در این قسمت عملکرد برآوردهای یول والکر، حداقل مربعات شرطی و حداکثر درستنمایی شرطی برای پارامترهای مدل برای اندازه‌های نمونه مختلف با دوره فصلی ۱۲ در جداول ۱ تا ۳ ارائه شده است. اریبی‌ها و میانگین توان دوم خطاها ۱۰۰۰ بار تکرار و محاسبه شده‌اند. این جداول نشان می‌دهد که اریبی و خطای استاندارد برآورد پارامترها با افزایش حجم نمونه برای همه موارد کاهش می‌یابد. همان‌طور که از جداول مشاهده می‌شود، روش حداقل مربعات شرطی و روش یول والکر، رفتارهای مشابهی را نشان می‌دهند. با توجه به جداول، برآوردهای حداکثر درستنمایی شرطی مقادیر کم‌تری از اریبی (مقادیر قدرمطلق) و میانگین مربعات خطا در مقایسه با برآوردهای یول والکر و حداقل مربعات شرطی را نشان می‌دهند.

جدول ۱: میزان اریبی و میانگین مربعات خطا (در داخل پرانتز) برآوردها برای $(\rho, \beta, \lambda) = (0.2, 3, 1.5)$

$\hat{\lambda}_{CML}$	$\hat{\lambda}_{cls}$	$\hat{\lambda}_{YW}$	$\hat{\beta}_{CML}$	$\hat{\beta}_{cls}$	$\hat{\beta}_{YW}$	$\hat{\rho}_{CML}$	$\hat{\rho}_{cls}$	$\hat{\rho}_{YW}$	n
۰.۱۷۶	۰.۴۹۲	۰.۰۷۵	-۰.۱۵۹	-۰.۳۲۰	-۰.۱۰۶	۰.۰۰۸	-۰.۰۲۳	-۰.۱۴۲	۲۰۰
(۰.۱۷۰۸)	(۰.۲۸۰۴)	(۰.۲۶۸۵)	(۰.۷۶۶)	(۰.۱۰۱۳)	(۰.۹۵۶)	(۰.۰۲۷)	(۰.۰۵۰)	(۰.۰۴۷)	
۰.۰۶۰	۰.۱۵۵	-۰.۰۰۳	-۰.۰۸۲	-۰.۱۳۱	-۰.۰۵۱	۰.۰۰۰	-۰.۰۱۳	-۰.۰۷۴	۴۰۰
(۰.۰۸۷۰)	(۰.۱۴۲۴)	(۰.۱۳۴۴)	(۰.۳۷۴)	(۰.۴۹۸)	(۰.۴۷۴)	(۰.۰۱۳)	(۰.۰۲۶)	(۰.۰۲۵)	
۰.۱۱۹	۰.۱۹۶	۰.۰۹۵	-۰.۰۶۷	-۰.۱۰۵	-۰.۰۵۷	-۰.۰۱۶	-۰.۰۲۹	-۰.۰۵۹	۸۰۰
(۰.۴۶۵)	(۰.۷۷۴)	(۰.۷۵۹)	(۰.۱۸۹)	(۰.۲۶۸)	(۰.۲۶۲)	(۰.۰۰۷)	(۰.۰۱۳)	(۰.۰۱۳)	

جدول ۲: میزان اریبی و میانگین مربعات خطا (در داخل پرانتز) برآوردها برای $(\rho, \beta, \lambda) = (0.5, 2, 1)$

$\hat{\lambda}_{CML}$	$\hat{\lambda}_{cls}$	$\hat{\lambda}_{YW}$	$\hat{\beta}_{CML}$	$\hat{\beta}_{cls}$	$\hat{\beta}_{YW}$	$\hat{\rho}_{CML}$	$\hat{\rho}_{cls}$	$\hat{\rho}_{YW}$	n
۰.۲۷۲	۰.۹۵۳	۰.۴۶۰	-۰.۱۱۸	-۰.۴۶۱	-۰.۲۳۵	-۰.۰۲۹	-۰.۰۸۶	-۰.۲۸۵	۲۰۰
(۰.۱۵۵۵)	(۰.۲۷۷۹)	(۰.۲۴۵۹)	(۰.۶۶۰)	(۰.۹۰۸)	(۰.۸۰۸)	(۰.۰۱۹)	(۰.۰۴۴)	(۰.۰۵۴)	
۰.۲۵۴	۰.۶۴۵	۰.۴۵۳	-۰.۱۶۳	-۰.۳۶۷	-۰.۲۷۸	-۰.۰۳۴	-۰.۰۸۱	-۰.۲۳۱	۴۰۰
(۰.۷۲۵)	(۰.۱۳۵۳)	(۰.۱۲۷۸)	(۰.۳۱۴)	(۰.۴۵۱)	(۰.۴۲۶)	(۰.۰۰۸)	(۰.۰۲۲)	(۰.۰۲۶)	
۰.۰۱۷	۰.۲۳۳	۰.۱۴۲	-۰.۰۶۸	-۰.۱۷۷	-۰.۱۳۰	-۰.۰۱۵	-۰.۰۳۷	-۰.۱۰۹	۸۰۰
(۰.۳۳۲)	(۰.۶۰۸)	(۰.۵۹۶)	(۰.۱۴۷)	(۰.۲۱۴)	(۰.۲۰۹)	(۰.۰۰۴)	(۰.۰۱۱)	(۰.۰۱۲)	

جدول ۳: میزان اریبی و میانگین مربعات خطا (در داخل پرانتز) برآوردها برای $(\rho, \beta, \lambda) = (0.8, 1.5, 4)$

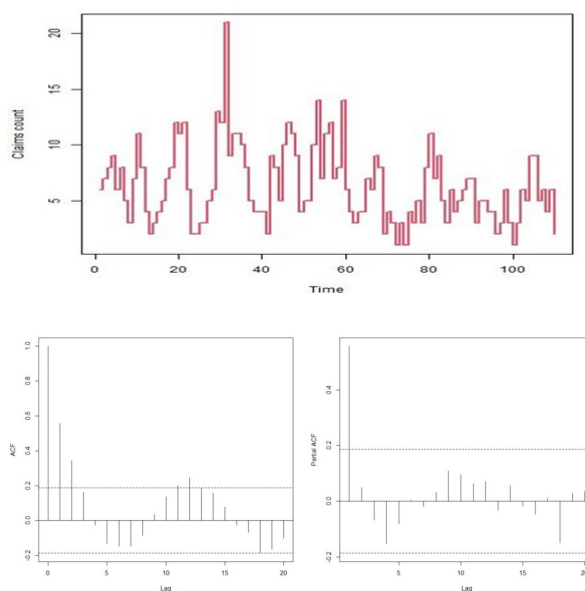
$\hat{\lambda}_{CML}$	$\hat{\lambda}_{cls}$	$\hat{\lambda}_{YW}$	$\hat{\beta}_{CML}$	$\hat{\beta}_{cls}$	$\hat{\beta}_{YW}$	$\hat{\rho}_{CML}$	$\hat{\rho}_{cls}$	$\hat{\rho}_{YW}$	n
۰.۴۶۱	۰.۲۶۰۲	۰.۱۵۳۲	-۰.۱۵۹	-۰.۱۲۴۱	-۰.۰۶۵۴	-۰.۰۳۰	-۰.۰۱۰۳	-۰.۰۵۸۵	۲۰۰
(۰.۴۴۸۸)	(۱.۲۰۰۵)	(۰.۸۱۵۱)	(۰.۱۰۸۸)	(۰.۲۶۷۶)	(۰.۱۹۳۶)	(۰.۰۰۵)	(۰.۰۲۲)	(۰.۰۵۷)	
۰.۱۸۰	۰.۱۷۶۸	۰.۱۱۲۰	-۰.۱۱۱	-۰.۰۸۸۸	-۰.۰۵۹۴	-۰.۰۲۱	-۰.۰۰۸۱	-۰.۰۳۲۶	۴۰۰
(۰.۲۴۱۶)	(۰.۵۶۴۵)	(۰.۴۹۴۰)	(۰.۵۵۵)	(۰.۱۳۳۴)	(۰.۱۱۶۵)	(۰.۰۰۲)	(۰.۰۱۰)	(۰.۰۲۱)	
۰.۳۱۸	۰.۱۱۳۰	۰.۰۹۰۷	-۰.۱۰۰	-۰.۰۵۰۲	-۰.۰۴۰۲	-۰.۰۰۹	-۰.۰۰۴۶	-۰.۱۶۹	۸۰۰
(۰.۱۱۲۰)	(۰.۲۶۷۳)	(۰.۲۴۸۶)	(۰.۲۶۹)	(۰.۶۶۶)	(۰.۶۲۵)	(۰.۰۰۱)	(۰.۰۰۵)	(۰.۰۰۸)	

۶ داده واقعی

در این بخش برای مدل سازی و پیش بینی و تجزیه و تحلیل از داده های سری زمانی شمارشی که از انجمن جبران خسارت کارگران (WCB) بریتیش کلمبیا و کانادا جمع آوری شده است، استفاده می کنیم. داده ها شامل ۱۲۰ مشاهده است که از ژانویه ۱۹۸۵ شروع شده و در دسامبر ۱۹۹۴ به پایان می رسد. میانگین نمونه ۶/۱۳ و واریانس نمونه ۱۱/۸۰ است. از این رو، به نظر می رسد داده ها بیش پراکنش داشته می باشند.

شکل ۱ نمودار سری زمانی داده ها، تابع خودهمبستگی (ACF) و تابع خودهمبستگی جزئی (PACF) را نشان می دهد. نمودار سری زمانی نشان می دهد که سری در میانگین ایستاست. الگوی کاهش هندسی با دوره فصلی ۱۲ را می توان در نمودار ACF مشاهده کرد و نشان می دهد که یک سری زمانی با رفتارهای همبستگی سریالی وجود دارد. فصلی بودن و بیش پراکنشی این مجموعه داده ها آشکار است، این ویژگی ها باعث شد که ما مدل $DELINAR(1)_{12}$ را به عنوان مدل برازش شده انتخاب کنیم. مجموعه داده ها به دو بخش تقسیم می شود. ۱۱۰ مشاهده اول برای مدل سازی سری استفاده می شوند و ۱۰ مشاهده باقیمانده برای اهداف پیش بینی در نظر گرفته شده است. همچنین برای مقایسه مدل های $INAR(1)_{12}$ ، $DELINAR(1)$ و $PINAR(1)$ را نیز در نظر می گیریم.

همان طور که قبلا اشاره شد پارامتر α را ثابت و معلوم در نظر گرفتیم و با توجه به این که α مقادیر صحیح غیرمنفی را می گیرد، بنابراین مدل را با مقادیر ۱، ۲ و ۳ برای α با رویکرد ماکزیم درستنمایی شرطی برازش کردیم تا بهترین را با توجه به معیار اطلاعاتی آکاییک (AIC) و معیار اطلاعات بیزی (BIC) انتخاب کنیم. با توجه به تجزیه و تحلیل حساسیت پارامتر α ، به نظر می رسد که معیار اطلاع آکاییک با افزایش مقدار α کمی افزایش می یابد. بنابراین ما برای مدل $\alpha = 1$ را در نظر می گیریم.



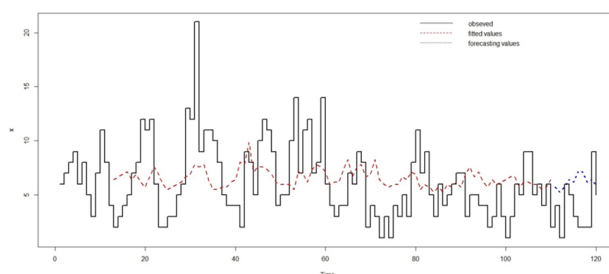
شکل ۱: نمودار سری زمانی و تابع خودهمبستگی و تابع خودهمبستگی جزئی داده های واقعی از سال ۱۹۸۵ تا ۱۹۹۴

جدول ۴ برآوردهای حداقل مربعات و حداکثر درستنمایی شرطی (با میانگین مربعات خطا در پرانتز) را برای چهار مدل برازش شده و معیار اطلاع آکاییک (AIC) و معیار اطلاع بیزی (BIC) را نشان می دهد. همان طور که می بینیم، AIC و BIC مربوط به مدل $DELINAR(1)_{12}$ در مقایسه با مدل های دیگر کمترین است. نتایج بدست آمده در جدول ۴ نشان می دهد که مدل خودبازگشتی با

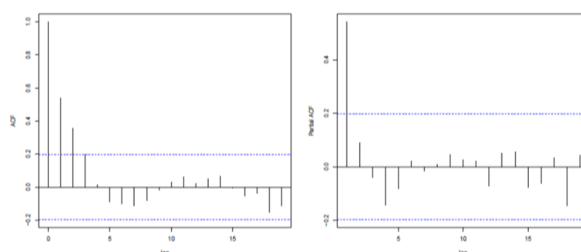
مقادیر صحیح غیرمنفی مرتبه اول دلاپورت با ساختار فصلی بهتر از سایر مدل‌های خودبازگشتی با مقادیر صحیح غیرمنفی مرتبه اول فصلی و غیرفصلی می‌باشد. جهت بررسی توانایی مدل در پیش‌بینی، از معیار ریشه میانگین مربع خطا (RMSE) استفاده شده که برای همه مدل‌ها محاسبه شده است و در جدول ۴ نوشته شده است. مقادیر جدول ۴ نشان می‌دهد که ریشه میانگین مربع خطا در هر دو مورد مدل دلاپورت فصلی و غیرفصلی خیلی مشابه است ولی در بین همه مدل‌ها مدل پواسون غیرفصلی کمترین مقدار را دارد. در شکل ۲ مقادیر مشاهده شده، مقادیر برازش و مقادیر پیش‌بینی به ترتیب با خطوط مشکی، قرمز و آبی نشان داده شده است. شکل ۳ نیز نشان دهنده تابع خودهمبستگی و تابع خودهمبستگی جزئی باقیمانده‌هاست و همان‌طور که مشاهده می‌شود همچنان همبستگی سریالی در باقیمانده‌ها دیده می‌شود.

جدول ۴: برآورد پارامترها (میانگین مربعات خطا در پرانتز) و مقادیر AIC و BIC

مدل	پارامتر	برآورد CML	برآورد CLS	AIC	BIC	RMSE
$PINAR(1)$	ρ	$0.5705(0.0005)$	$0.5651(0.0075)$	۵۴۵۸۰	۵۵۱۲۰	۳۱۰۱
	λ	$0.2705(0.0188)$	$0.2735(0.3199)$			
$INAR(1)_{12}$	ρ	$0.2459(0.0004)$	$0.2667(0.0098)$	۵۳۲۰۹	۵۳۷۴۹	۳۳۷۸
	λ	$0.6769(0.0411)$	$0.5389(0.4201)$			
$DELINAR(1)$	ρ	$0.4436(0.0016)$	$0.5651(0.0070)$	۵۲۷۰۸	۵۳۵۱۸	۳۷۷۹
	λ	$0.2768(0.0077)$	$0.1683(0.1000)$			
$DELINAR(1)_{12}$	ρ	$0.2285(0.0014)$	$0.2667(0.0100)$	۵۰۰۲۱	۵۰۸۳۱	۳۷۸۵
	λ	$0.2735(0.0032)$	$0.5089(0.3197)$			



شکل ۲: نمودار سری زمانی مقادیر مشاهده شده، برازش شده و پیش‌بینی مربوط به داده‌های واقعی از سال ۱۹۸۵ تا ۱۹۹۴



شکل ۳: تابع خودهمبستگی و تابع خودهمبستگی جزئی باقیمانده‌ها

۷ نتیجه‌گیری

در این مقاله یک فرایند خودبازگشتی مرتبه اول فصلی با توزیع حاشیه‌ای دلاپورت معرفی گردید که توزیع نوآوری مدل به صورت پیچشی از توزیع پواسون و α تا توزیع هندسی تعدیل یافته در صفر می‌باشد. در ادامه ویژگی‌های اصلی مدل مانند میانگین واریانس و تابع خودهمبستگی مورد بررسی قرار داده شد. برای برآورد پارامترها از روش‌های یول-والکر، حداقل مربعات شرطی و حداکثر درستنمایی شرطی استفاده شده است. همچنین به پیش‌بینی مدل پرداخته شده است. با انجام یک شبیه‌سازی در حجم نمونه‌های متناهی عملکرد این برآوردگرها مورد ارزیابی قرار گرفت. در نهایت این مدل به یک مجموعه داده واقعی، برازش داده شد. برای مقایسه مدل‌های $INAR(1)$ ، $PINAR(1)$ و $DELINAR(1)$ نیز در نظر گرفته شد و نتیجه نشان داد که براساس معیار AIC و BIC ، مدل پیشنهادی ما در مقایسه با سایر مدل‌های $INAR$ بهتر است. به منظور تحقیقات بیشتر، می‌توان مدل را به مرتبه‌های بیشتر گسترش داد. همچنین، طبق ACF باقیمانده‌ها که در شکل ۳ نشان داده شده است، همچنان همبستگی سریالی در باقیمانده‌ها وجود داشت بنابراین می‌توان یک مدل جدید معرفی کرد که بتواند همبستگی فصلی و سریالی را کاهش دهد.

مراجع

- Al-Osh MA, Alzaid AA. (1987) First-order interger-valued autoregressive ($INAR(1)$) process, *J Time Ser Anal.*, **8**(3), 261–275.
- Barreto-Souza W. (2019), Mixed Poisson $INAR(1)$ processes, *Stat Pap.*, **60**(3), 2119–2139.
- Bourguignon M, Vasconcellos KL, Reisen VA, Ispány M (2016), A Poisson $INAR(1)$ process with a seasonal structure, *J Stat Comput Simul*, **86**(2), 373–387.
- Delaporte P. (1959), Quelques problèmes de statistiquemathématique posés par l'assurance automobile et le bonus non sinistre, *Bulletin Trimestriel de l'Institut des Actuaire Français*, **227**, 87–102.
- Klimko, L.A. and Nelson, P.I., (1978), On conditional least squares estimation for stochastic processes, *The Annals of statistics*, 629–642.
- McKenzie E. (1985), Some simple models for discrete variate time series, *J Am Water Resour Assoc.*, **21**(4): 645–650.
- Steutel FW, Van Harn K. (1979), Discrete analogues of self-decomposability and stability, *Ann Probab*; **7**(5), 893–899.
- Shalbaf, M., Parham, G. and Chinipardaz, R., (2022), Binomial Thinning Integer-Valued AR (1) with Poisson- α Fold Zero Modified Geometric Innovations, *Journal of Sciences, Islamic Republic of Iran*, **33**(1), 55–63.
- Tian S, Wang D, Cui S. (2020), A seasonal geometric $INAR$ process based on negative binomial thinning operator, *Statistical Papers*, **61**(6), 2561–2581.

A seasonal Integer-Valued AR(1) model with Delaporte marginal distribution

Maryam Shalbaf¹, Gholam Ali Parham²

¹Department of Statistics and Mathematics, Sho.C., Islamic Azad University, Shoushtar, Iran

²Department of Statistics, Faculty of Mathematical Sciences and Computer, Shahid Chamran University of Ahvaz

Abstract: Real-count data time series often show the phenomenon of the overdispersion. In this paper, we introduce the first-order integer-valued autoregressive process with a seasonal structure. The univariate marginal distribution is derived from the Delaporte distribution and the innovations are convolution of Poisson with α -fold zero modified geometric distribution, based on binomial thinning operator. Some properties of the model are derived. The methods of Yule-Walker, conditional least squares and conditional maximum likelihood are used for estimating the parameters. The Monte Carlo experiment is conducted to evaluate the performances of these estimators in finite samples. At the end this model is illustrated using a real data set.

Keywords: Delaporte distribution, seasonal INAR(1) model, α -fold zero modified geometric, binomial thinning.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30-31, 2025
University of Kurdistan



رویکرد فازی خوشه بندی داده ها بر پایه الگوریتم FCM

مولود میاهی^۱، بهرام طارمی^۲، مینا میاهی^۳

^۱ گروه علوم پایه، دانشکده علوم پزشکی ساوه

^۲ بخش آمار، دانشکده علوم، دانشگاه شیراز

^۳ گروه ریاضی، مدعو دانشکده علوم پزشکی ساوه

چکیده: با افزایش حجم داده‌ها، تحلیل و استخراج الگوها به چالشی مهم تبدیل شده است. خوشه‌بندی یکی از روش‌های موثر برای گروه‌بندی داده‌های مشابه است. در روش‌های سنتی، هر داده تنها به یک خوشه تعلق دارد، اما در واقعیت برخی داده‌ها در مرز خوشه‌ها قرار می‌گیرند و نمی‌توان آن‌ها را به‌طور قطعی دسته‌بندی کرد. خوشه‌بندی فازی با امکان تعلق هم زمان یک داده به چند خوشه، رویکردی منعطف تر ارائه می‌دهد. این روش به دلیل مدل‌سازی بهتر، کاهش خطا و انعطاف‌پذیری مورد توجه محققان قرار گرفته است. این مقاله، به معرفی افراز فازی، خوشه بندی فازی و الگوریتم Fuzzy C-Means (FCM) می پردازد.

واژه‌های کلیدی: افراز فازی، خوشه‌بندی فازی، الگوریتم FCM

کد موضوع بندی ریاضی (۲۰۲۰): ۹۱B۸۴، ۶۲M۱۰

۱ مقدمه

خوشه‌بندی یک تکنیک اکتشافی است که با تقسیم اشیاء همگن به گروه‌های مجزا، هدف آن ساده‌سازی مجموعه داده‌های پیچیده است. الگوریتم‌های خوشه‌بندی متعددی از دیدگاه‌های مختلف توسعه داده شده‌اند که روش‌های پارتیشن‌بندی، روش‌های مبتنی بر چگالی، روش‌های مبتنی بر مدل و روش‌های سلسله مراتبی پرکاربردترین آنها در عمل هستند. (هو و همکاران، ۲۰۲۵) برخی از این الگوریتم‌ها معمولاً بر بهینه‌سازی یک تابع هدف خاص تمرکز دارند. با این حال، کمبود قابل توجهی در تفسیرپذیری و توضیح‌پذیری این نتایج خوشه‌بندی وجود دارد که چالش‌های تصمیم‌گیری را برای کاربران هنگام استفاده از این الگوریتم‌ها ایجاد می‌کند. (هو و همکاران،

^۱ M.Miahi@savehums.ac.ir

^۲ سخنران، Tarami@shirazu.ac.ir

۲۰۲۵) یکی از این روش های چالش برانگیز روش متداول الگوریتم C-Means است که خوشه ها را بر پایه کمینه سازی فاصله تا مراکز خوشه تعیین می کند. این روش به شکل و پراکندگی داده ها حساسیت ندارد و همچنین در مواجهه با داده هایی که فاصله ای مشابه از چند مرکز دارند، دارای عملکرد نامطلوبی است. (عسگریان و همکاران، ۱۳۸۶) همانطور که مشاهده می شود اصول مدل بندی کلاسیک در عدم قطعیت نمی تواند کارایی داشته باشد. سوالی در اینجا پیش می آید که با عدم قطعیت چه باید کرد؟ پاسخ این سوال را منطق فازی می دهد. (میاهی و نایینی، ۱۳۹۱) بنابراین برای رفع این محدودیت ها، الگوریتمی مانند Fuzzy C-Means (FCM) توسعه داده شده است که امکان شناسایی خوشه های فازی با اشکال متنوع را فراهم می کنند. (عسگریان و همکاران، ۱۳۸۶)

۲ دست آورد پژوهش

تعریف ۱.۲. افراز c تایی فازی: فرض کنید $X = \{x_1, x_2, \dots, x_n\}$ مجموعه ای از داده ها باشد و V_{cn} مجموعه تمام ماتریس های حقیقی $n \times c$ که $c \leq n$ است، باشد و ماتریس $\tilde{U} = [u_{ik}] \in V_{cn}$ یک افراز c تایی فازی نامیده می شود اگر شرایط زیر را داشته باشد:

$$1. \quad 0 \leq u_{ik} \leq 1 \quad 1 \leq i \leq c \quad 1 \leq k \leq n$$

$$2. \quad \sum_{i=1}^c u_{ik} = 1 \quad 1 \leq k \leq n$$

$$3. \quad 0 < \sum_{j=1}^n u_{ik} < n \quad 1 \leq i \leq c$$

مجموعه تمام ماتریس هایی که شرایط فوق را برآورده می سازند، مجموعه افرازهای فازی c تایی (M_{fc}) نامیده می شود. شرط اول نشان می دهد که در مقایسه با افراز قطعی، در افراز فازی داده ها می توانند با درجات عضویت مختلف به چندین خوشه تعلق داشته باشند. شرط دوم تضمین می کند که مجموع درجات عضویت یک داده در تمام خوشه ها برابر یک باشد، و شرط سوم بیان می کند که هیچ خوشه ای نباید کاملاً خالی یا شامل تمام داده ها باشد.

تعریف ۲.۲. تابع هدف: در افراز فازی c تایی، با استفاده از معیار پراکندگی و فاصله مینکوفسکی، تابع هدف بصورت زیر تعریف می شود:

$$minz(\tilde{U}, V) = \sum_{i=1}^c \sum_{k=1}^n (u_{ik})^m \|x_k - v_i\|^2. \quad (1.2)$$

بطوریکه

$$v_i = \frac{1}{\sum_{k=1}^n u_{ik}} \sum_{k=1}^n (u_{ik})^m . x_k \quad i = 1, 2, \dots, c. \quad (2.2)$$

است. که در آن v_i مرکز خوشه i ام و میانگین وزن دار x_k هاست که وزن هر x_k ، توان m درجه عضویتش می باشد. هرچه درجه عضویت بالاتر باشد تاثیر بیشتری در تعیین v_i دارند و همچنین با افزایش m نیز این تاثیر افزایش می یابد.

تعریف ۳.۲. نرم عمومی: با فرض اینکه G ماتریس متقارن و مثبت $p \times p$ باشد، نرم عمومی را بصورت زیر تعریف می کنیم

$$\|x_k - v_i\|_G^2 = (x_k - v_i)' . G . (x_k - v_i). \quad (3.2)$$

بدرنظر گرفتن نرم عمومی مسئله افراز فازی بصورت زیر نوشته می شود :

$$\min z(\tilde{U}, V) = \sum_{i=1}^c \sum_{k=1}^n u_{ik}^m \|x_k - v_i\|_G^2. \quad (4.2)$$

مزیت مسئله بالا این است که با کمک از معادلات دیفرانسیل می توان شرایط لازم را برای نقطه بهینه موضعی تعیین کرد. درواقع حل آن با درنظر گرفتن شرط مجموع درجات عضویت برابر با یک و مرکز خوشه بدست آمده از فرمول (۲.۲)، شرط

$$\mu_{ik} = \frac{\left(\frac{1}{\|x_k - v_i\|_G^2}\right)^{\frac{1}{m-1}}}{\sum_{j=1}^c \left(\frac{1}{\|x_k - v_{ji}\|_G^2}\right)^{\frac{1}{m-1}}}. \quad i = 1, 2, \dots, c \quad k = 1, 2, \dots, n \quad (5.2)$$

و با مشتق گیری از تابع هدف نسبت به v_i (به ازای \tilde{U} ثابت) و نسبت به μ_{ik} (به ازای v ثابت) صورت می پذیرد. در واقعیت معادلات بالا به دلیل پیچیدگی و غیرخطی بودن به صورت تحلیلی قابل حل نمی باشند. بنابراین با بهره گیری از الگوریتم های تکراری می توان مینیمم تابع هدف را بطور تقریبی بدست آورد. الگوریتمی که برای این منظور در این مقاله استفاده می شود الگوریتم FCM است.

تعریف ۴.۲: الگوریتم FCM: یک الگوریتم با یک روند تکرار بسیار ساده برای رسیدن به یک بهینه موضعی همگرا در حل معادله (۴.۲) است. برای طراحی این الگوریتم پارامترهای تعداد خوشه ها (c)، توان وزنی که فازی بودن نتایج خوشه ها را کنترل کند (m)، ماتریس مقارن و مثبت ($G_{P \times P}$)، مقدار دهی اولیه ماتریس ($\tilde{U}^{(0)}$)، معیار اختتام الگوریتم (Δ) نیاز است. این الگوریتم برای یافتن خوشه های نسبتا کروی مناسب است. با انتخاب یک مقدار ثابت برای c ($2 \leq c < n$) و رعایت شرط $\sum_{i=1}^c u_{ik} = 1$ شروع و گام های زیر را سپری می شود:

۱. مقدار m و ماتریس مقارن مثبت $G_{p \times p}$ را انتخاب نموده سپس با قرار دادن $L = 0$ ماتریس $\tilde{U}^{(L)} \in M_{fc}$ مقدار دهی اولیه می شود.

۲. مراکز c خوشه فازی اولیه $v_i^{(L)}$ را با کمک از ماتریس $\tilde{U}^{(L)}$ و رابطه v_i (فرمول (۲.۲)) بدست می آید.

۳. ماتریس درجات عضویت را به روز رسانی میشود یعنی ماتریس جدیدی به نام $\tilde{U}^{(L+1)}$ را یافته و برای یافتن هر یک از عناصر این ماتریس جدید از شرط زیر استفاده می گردد:

اگر به ازای هر i ، $x_k \neq v_i^{(L)}$ شود از μ_{ik} (فرمول (۵.۲)) کمک گرفته شود در غیر اینصورت از رابطه زیر استفاده گردد:

$$\mu_{jk}^{(L+1)} = \begin{cases} 1 & j = i \\ 0 & j \neq i \end{cases}$$

۴. محاسبه Δ و بررسی شرط توقف در این گام صورت می پذیرد. اگر $\Delta = \|\tilde{U}^{(L+1)} - \tilde{U}^{(L)}\|_G$ باشد و مقدار آن از ε بیشتر باشد آنگاه بجای L ، $L+1$ قرار داده و به گام دوم رفته و الگوریتم تکرار می گردد، درغیراینصورت الگوریتم متوقف می شود.

۳ دستورات در R

قبل از اجرای الگوریتم، داده‌ها را باید وارد R شده و در یک متغیر به نام x ذخیره گردد، سپس مراحل زیر انجام شود :

۱. نصب و بارگذاری بسته: ابتدا بسته‌ی ppclust نصب و بارگذاری می‌شود تا بتوان از الگوریتم FCM استفاده کرد.

```
install.packages("ppclust")
```

```
library(ppclust)
```

۲. اجرای الگوریتم FCM : با دستور مربوطه الگوریتم بر اساس تعداد خوشه‌ها و ضریب فازی روی داده‌ها اجرا می‌شود.

```
fcmresult <- fcm(x, centers = number of cluster, m )
```

۳. مشاهده مراکز خوشه‌ها: دستور دیگری مختصات مراکز خوشه‌ها را نشان می‌دهد.

```
fcmresult$centers
```

۴. مشاهده ماتریس درجه تعلق: براساس این ماتریس میتوان میزان تعلق هر داده به هر خوشه را نمایش داد.

```
fcmresult$u
```

۵. تولید برجسب خوشه‌ها: با استفاده از بیشترین درجه تعلق، برای هر داده یک برجسب خوشه مشخص می‌شود.

بحث و نتیجه‌گیری

تحلیل خوشه ای مبني بر افراز مجموعه اي از نقاط داده در تعدادي از خوشه ها (زیر گروه ها) استوار گردیده است، جايي که اشیاء داخل يك خوشه (زیر گروه) يك درجه معين از شبیه بودن را نشان مي دهند. خوشه بندي متداول هر نقطه (بردار ویژگی) را به يکي و فقط يکي از خوشه ها با درجه عضویت يك اختصاص مي دهد. فرضیات به خوبی مرز بین خوشه ها را تعیین می کنند ولی مرز اغلب توصیفی از داده حقیقی را منعکس نمی کند. بنابراین اگر بین زیر گروه ها ارتباط فازی باشد، جایی که توصیف گوناگونی از شباهت اشیاء در يك خوشه خاص نیاز باشد، مشکلات فراواني ایجاد مي شود. در این شرایط خوشه بندی فازی می تواند راهکار مناسبی ارائه دهد. با این حال، چالش هایی مانند تعیین تعداد بهینه خوشه ها، موقعیت اولیه مراکز، و تأثیر شکل و چگالی داده ها در نتایج نهایی، از مهم ترین مسائلی هستند که در پیاده سازی موفق این روش باید مورد توجه قرار گیرند.

مراجع

عسگریان، ا.، معین‌زاده، ح.، سریانی، م.، و حبیبی، ج. (۱۳۸۶)، رویکرد جدید برای خوشه‌بندی فازی بوسیله الگوریتم ژنتیک، سیزدهمین کنفرانس سالانه انجمن کامپیوتر ایران.

میاهی، م. و نایینی، ح. (۱۳۹۱)، خوشه‌بندی سری زمانی فازی بر اساس تابع خودهمبستگی، اولین همایش بین‌المللی اقتصادسنجی، روش‌ها و کاربردها.

Hu L., Jiang M., Liu X., and He Z. (2025), *Significance-based decision tree for interpretable categorical data clustering*, Inform. Sci., 690, 121588.

Hu L., Jiang M., Dong J., Liu X., and He Z. (2025), *Interpretable categorical data clustering via hypothesis testing*, Pattern Recogn., 111364.

A Fuzzy Approach to Data Clustering Based on the FCM Algorithm

Moloud Miah¹, Bahram Tarami², Mina Miah³

¹Department of Basic Sciences, Faculty of Medical Sciences, Saveh, Iran

²Department of Statistics, Faculty of Basic Sciences, Shiraz University, Shiraz, Iran

³Department of Mathematics, Adjunct Lecturer, Saveh University of Medical Sciences, Saveh, Iran

Abstract: With the growing volume of data, analyzing and extracting meaningful patterns has become a major challenge. Clustering is a widely used technique for grouping similar data points. In conventional clustering methods, each data point is assigned to a single cluster; however, in practice, some data points lie on the boundaries between clusters and cannot be classified definitively. Fuzzy clustering, which allows a data point to belong to multiple clusters simultaneously, provides a more flexible and realistic approach. This method has gained considerable attention from researchers due to its superior modeling capabilities, error reduction, and adaptability. This paper presents an overview of fuzzy partitioning, fuzzy clustering, and the Fuzzy C-Means (FCM) algorithm.

Keywords: Fuzzy partitioning, Fuzzy clustering, Fuzzy C-Means (FCM) algorithm

Mathematics Subject Classification (2020): 62M10, 91B84.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



بررسی رفتار مجانبی رهیافت PTE در حضور همخطی چندگانه

سید امیرحسین طباطبائی شیرازی^۱، مهدی عمادی، محمد آرشی، سولماز سیفاللهی

گروه آمار، دانشگاه فردوسی مشهد

چکیده: رگرسیون خطرات متناسب کاکس یکی از مدل‌های پرکاربرد در تحلیل داده‌های بقا با متغیرهای کمکی است. با این حال، در حضور همخطی چندگانه، کارایی این مدل کاهش می‌یابد و برآوردهای حاصل از روش استاندارد درستنمایی جزئی ممکن است غیرقابل اعتماد شوند. در این مقاله، برآوردگر پیش‌آزمون را بر اساس برآوردگر لیو بسط می‌دهیم تا دقت برآورد ضرایب افزایش یابد. ویژگی‌های نظری حدی این برآوردگر را تعیین و عملکرد آن از طریق شبیه‌سازی‌های گسترده مونت‌کارلو ارزیابی می‌گردد. همچنین، نحوه کاربرد این استراتژی برآورد بر داده‌های بقا مورد بررسی قرار می‌گیرد. نتایج نشان‌دهنده بهبودهای چشمگیر بوده و مزایای عملی این برآوردگر را برای پژوهشگران به‌خوبی آشکار می‌سازد.

واژه‌های کلیدی: مدل خطرات متناسب کاکس، برآوردگر لیو، همخطی، برآوردگر پیش‌آزمون.

کد موضوع‌بندی ریاضی (۲۰۲۰): 62F30، 62N02.

۱ مقدمه

در عصر حاضر، تحلیل داده‌های پیچیده و حجیم، به‌ویژه در حوزه‌هایی مانند پزشکی، زیست‌فناوری و مهندسی، به یکی از چالش‌های اساسی پژوهشگران تبدیل شده است. در این میان، مدل‌های رگرسیونی ابزارهای مهمی برای تفسیر روابط میان متغیرها به‌شمار می‌روند. یکی از پرکاربردترین این مدل‌ها در تحلیل داده‌های بقا، مدل خطرات متناسب کاکس (CPHM) است که با امکان بررسی هم‌زمان تأثیر چندین عامل بر نرخ وقوع یک رویداد، ابزاری قدرتمند در مطالعات طولی به‌حساب می‌آید (کاکس، ۱۹۷۲).

مدل کاکس با فرض ثبات ضرایب در طول زمان و بهره‌گیری از روش درست‌نمایی جزئی برای برآورد پارامترها، ساختاری انعطاف‌پذیر برای تحلیل زمان وقوع رویدادهای حیاتی فراهم می‌آورد. با این حال، در عمل با چالش‌هایی روبه‌روست؛ یکی از مهم‌ترین آن‌ها، همخطی چندگانه یا همبستگی شدید میان متغیرهای مستقل است. این پدیده می‌تواند منجر به ناپایداری ضرایب، افزایش واریانس برآوردها و

^۱ سخنران، tabatabaeishirazi.se@um.ac.ir

کاهش دقت و قابلیت تفسیر نتایج شود. راهکارهایی نظیر حذف یا گروه‌بندی متغیرهای همبسته، اگرچه مرسوم‌اند، اما ممکن است موجب افت کارایی مدل و بروز آریبی در نتایج گردند. از این رو، توجه به این چالش و ارائه راهکارهای مؤثر برای مقابله با آن، ضرورتی اجتناب‌ناپذیر است. در راستای مقابله با هم‌خطی در CPHM، برآوردهای متعددی پیشنهاد شده‌اند؛ از جمله برآوردهای ریح، برآوردهای نوع-لیو و برآوردهای کیبیریا-لکمان. در این میان، برآوردهای لیو به سبب ساختار ساده، برخورداری از تنها یک پارامتر تنظیم و سهولت در تعیین آن، مورد توجه ویژه قرار گرفته است. مطالعات متعدد نشان داده‌اند که این برآوردها، به‌ویژه در شرایط وجود هم‌خطی، از نظر معیار میانگین مربعات خطا، عملکرد بهتری نسبت به برآوردهای حداکثر درست‌نمایی جزئی (MPLE) دارد.

در برخی از تحلیل‌های رگرسیونی، اطلاعات پیشین درباره پارامترهای مدل ممکن است در دسترس باشد. به عنوان مثال، نتایج تحلیل‌های مقدماتی می‌توانند نشان دهند که برخی از متغیرها از لحاظ آماری معنادار نیستند و در نتیجه باید از مدل حذف شوند. این قبیل اطلاعات را می‌توان به صورت قیود خطی از نوع $C\beta = b$ در ساختار مدل اعمال کرد. برآوردهای که تحت این قیود به دست می‌آید، برآوردهای مقید نامیده می‌شود. با این وجود، در بسیاری از موارد، نسبت به صحت این قیود تردید وجود دارد. در چنین شرایطی، استفاده از برآوردهای پیش‌آزمون پیشنهاد می‌شود. این برآوردها، ترکیبی از برآوردهای معمول و برآوردهای مقید است که تصمیم‌گیری درباره استفاده از هر یک، بر اساس نتایج آزمون فرض آماری انجام می‌گیرد. مطالعات نشان داده‌اند که در صورت اعتبار اطلاعات پیشین، برآوردهای پیش‌آزمون می‌تواند نسبت به برآوردهای کلاسیک عملکرد بهتری از لحاظ میانگین مربعات خطا داشته باشد.

هدف این مقاله، معرفی برآوردهای پیش‌آزمون در (CPHM) با استفاده از برآوردهای لیو و بررسی ویژگی‌های مجانبی آن است. در ادامه، با بهره‌گیری از شبیه‌سازی مونت‌کارلو، عملکرد این برآوردها را با برآوردهای لیو در شرایط مختلف، مقایسه خواهد شد و در نهایت، کارایی آن در تحلیل داده‌های مربوط به سرطان ریه مورد ارزیابی قرار خواهد گرفت.

۲ مدل خطرات متناسب کاکس

مدل CPHM که نخستین بار توسط کاکس (۱۹۷۲) معرفی شد، به عنوان یکی از ارکان اصلی در تحلیل بقا شناخته می‌شود. این مدل، به‌ویژه به دلیل توانایی‌اش در برقراری ارتباط میان زمان وقوع یک رویداد و یک یا چند متغیر کمکی، بدون نیاز به فرض شکل مشخصی برای تابع خطر پایه، از محبوبیت فراوانی برخوردار است. فرض کنید $\tau_i = \min(t_i, c_i)$ که در آن، t_i زمان بقا در صورت مشاهده رویداد و c_i زمان سانسور در صورت عدم مشاهده رویداد است. همچنین، بردار متغیرهای کمکی را به صورت $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^T$ در نظر بگیرید. در این صورت، مدل CPHM به صورت $h(t_i | \mathbf{x}_i) = h_0(t_i) \exp(\mathbf{x}_i^T \beta)$ تعریف می‌شود که در آن، $h_0(t)$ تابع خطر پایه و $\beta = (\beta_1, \dots, \beta_p)^T$ بردار ضرایب مجهول مدل می‌باشد. برآوردهای حداکثر درست‌نمایی جزئی (MPLE) برای برآورد β به صورت زیر محاسبه شده است:

$$\hat{\beta}_{\text{MPLE}} = (\hat{Z}^T \hat{Z})^{-1} \hat{Z}^T \hat{V}, \quad (1.2)$$

بطوریکه \hat{Z} ماتریسی به ابعاد $\sum (m_i + 1) \times p$ است، m_i اندازه مجموعه‌ی ریسک آزمودنی i است و \hat{V} برداری با ابعاد $p \times 1$ می‌باشد که در $\hat{\beta}_{\text{MPLE}}$ محاسبه می‌شود تحت شرایط نظم، نشان داده شده که برآوردهای MPLE دارای توزیع مجانبی نرمال به صورت زیر می‌باشد: $\sqrt{n}(\hat{\beta}_{\text{MPLE}} - \beta) \sim N_p(0, F^{-1})$ (ماگنوس، ۱۹۹۷). فرض کنید اطلاعات پیشین درباره ضرایب مدل را بتوان در قالب q قید خطی مستقل بیان کرد. برای بررسی صحت این قیود در

مدل، می‌توان آن‌ها را بر اساس فرضیه‌های آماری $C\beta = b$ vs $C\beta \neq b$: H_0 : $C\beta = b$ vs H_1 : $C\beta \neq b$ آزمود. در این رابطه، C و b به ترتیب ماتریس و برداری معلوم با ابعاد $q \times p$ و $q \times 1$ می‌باشند. بر اساس توزیع مجانبی MPLE، آماره‌ای که برای آزمون فرض صفر H_0 به کار می‌رود، به صورت زیر تعریف می‌شود:

$$T_n = (C\hat{\beta}_{\text{MPLE}} - b)^\top (C I^{-1} C^\top)^{-1} (C\hat{\beta}_{\text{MPLE}} - b). \quad (2.2)$$

با افزایش اندازه نمونه (n) ، آماره آزمون T_n تحت H_0 دارای توزیع مجانبی کای‌دو با q درجه آزادی است. در صورتی که H_0 برقرار باشد، می‌توان برای برآورد ضرایب مدل از برآوردگر درست‌نمایی جزئی مقید (RMPLE) استفاده کرد که به صورت زیر تعریف می‌شود:

$$\hat{\beta}_{\text{RMPLE}} = \hat{\beta}_{\text{MPLE}} - A(C\hat{\beta}_{\text{MPLE}} - b), \quad (3.2)$$

$$A = F^{-1} C^\top (C F^{-1} C^\top)^{-1}$$

که در آن

۳ برآوردگرهای پیشنهادی در مدل CPHM

برآوردگر ليو یکی از روش‌های شناخته‌شده برای برآورد ضرایب در حضور همخطی در مدل‌های رگرسیونی است. در CPHM، این برآوردگر توسط احمد و همکاران (۲۰۲۳) به صورت $\hat{\beta}_{\text{CLE}} = (\hat{Z}^\top \hat{Z} + I_p)^{-1} (\hat{Z}^\top \hat{Z} + d I_p) \hat{\beta}_{\text{MPLE}}$ تعریف شده است که در آن $d \in [0, 1]$ پارامتر ليو می‌باشد. برای سادگی می‌توان برآوردگر ليو را به فرم $\hat{\beta}_{\text{CLE}} = B \hat{\beta}_{\text{MPLE}}$ نوشت که در آن $B = (\hat{Z}^\top \hat{Z} + I_p)^{-1} (\hat{Z}^\top \hat{Z} + d I_p)$ می‌باشد. برای برآورد پارامتر ليو از برآوردگری زیر در این مقاله استفاده می‌کنیم:

$$\hat{d} = 1 - \frac{\sum_{j=1}^p 1/\lambda_j (\lambda_j + 1)}{\sum_{j=1}^p \hat{\beta}_{\text{MPLE}(j)}^2 / (\lambda_j + 1)^2}. \quad (1.3)$$

با الهام گرفتن از کیریبا و صالح (۲۰۰۴) می‌توان فرم برآوردگر مقید ليو (RCLE) به صورت، $\hat{\beta}_{\text{RCLE}} = B \hat{\beta}_{\text{RMPLE}}$ تعریف کرد.

برآوردگری که هدف اصلی این مقاله می‌باشد، برآوردگر پیش آزمون (PTE) در مدل CPHM می‌باشد. این برآوردگر که از انواع برآوردگرهای انقباضی می‌باشد، بر اساس صحت و یا نادرستی فرض مطرح شده در مدل، بین دو برآوردگر مقید و نامقید تصمیم‌گیری می‌کند. با استفاده از برآوردگر MPLE و RMPLE برآوردگر پیش آزمون به فرم زیر می‌توان تعریف کرد

$$\hat{\beta}_{\text{PTE}} = \hat{\beta}_{\text{MPLE}} - (\hat{\beta}_{\text{MPLE}} - \hat{\beta}_{\text{RMPLE}}) I_{(T_n < \chi_{(q, \alpha)}^2)}, \quad (2.3)$$

و بنا به رابطه‌ای که بین $\hat{\beta}_{\text{CLE}}$ و $\hat{\beta}_{\text{RCLE}}$ به ترتیب با $\hat{\beta}_{\text{MPLE}}$ و $\hat{\beta}_{\text{RMPLE}}$ دارند، می‌توان برآوردگر پیش آزمون بر اساس برآوردگر ليو را به صورت $\hat{\beta}_{\text{CLPTE}} = B \hat{\beta}_{\text{PTE}}$ تعریف کرد. یا به عبارتی دیگر،

$$\hat{\beta}_{\text{CLPTE}} = \hat{\beta}_{\text{RCLE}} - (\hat{\beta}_{\text{CLE}} - \hat{\beta}_{\text{RCLE}})I_{(T_n < \chi_{(q, \alpha)}^2)}, \quad (3.3)$$

فرم برآوردگر پیش‌آزمون لیو حکایت از این دارد که اگر H_0 پذیرفته شود، برآوردگر پیش‌آزمون برابر RCLE و در غیر این صورت، برابر برآوردگر CLE خواهد بود. این برآوردگر توسط **ماگنوس (۱۹۹۷)** معرفی شد و به شدت در مدل‌های رگرسیونی مورد توجه قرار گرفت.

۱.۳ ویژگی برآوردگرهای پیشنهادی

در این بخش به بررسی ویژگی برآوردگرهای معرفی شده به خصوص CPTE می‌پردازیم. در ابتدا آماره آزمون T_n در رابطه (۲.۲) را تحت $H_\nu: C\beta = b + \nu, \nu \in \mathbb{R}^q$ در نظر می‌گیریم. در این حالت، زمانیکه $n \rightarrow \infty$ ، خواهیم داشت $T_n \rightarrow \infty$ به گونه‌ای که $\lim_{n \rightarrow \infty} P_\nu(T_n \geq x) = 1$. به بیان دیگر، T_n یک آزمون سازگار است. در نتیجه، تحت فرض‌های جایگزین ثابت، زمانی که $n \rightarrow \infty$ داریم

$$\begin{aligned} & n(\hat{\beta}_{\text{PTE}} - \hat{\beta}_{\text{MPLE}}) \left(\frac{1}{n} I \right) (\hat{\beta}_{\text{PTE}} - \hat{\beta}_{\text{MPLE}}) \\ &= n(C\hat{\beta}_{\text{MPLE}} - b)^\top \left(\frac{1}{n} CI^{-1}C^\top \right)^{-1} (C\hat{\beta}_{\text{MPLE}} - b) I_{(T_n < \chi_{(q, \alpha)}^2)} \\ &= T_n I_{(T_n < \chi_{(q, \alpha)}^2)} \leq \chi_{(q, \alpha)}^2 I_{(T_n < \chi_{(q, \alpha)}^2)} \rightarrow 0. \end{aligned}$$

به عبارتی دیگر $\sqrt{n}(\hat{\beta}_{\text{PTE}} - \beta) \stackrel{D}{=} \sqrt{n}(\hat{\beta}_{\text{MPLE}} - \beta)$ و در نتیجه $\sqrt{n}(\hat{\beta}_{\text{CLE}} - \beta) \stackrel{D}{=} \sqrt{n}(\hat{\beta}_{\text{CLPTE}} - \beta)$. بنابراین، به منظور تمایز مناسب میان توزیع‌های مجانبی برآوردگرها، فرض‌های جایگزین محلی $H_{\circ(n)}$ را به صورت زیر در نظر می‌گیریم:

$$H_{\circ(n)}: C\beta = b + \frac{\vartheta}{\sqrt{n}}; \quad \vartheta = (\vartheta_1, \dots, \vartheta_q)^\top \neq 0. \quad (4.3)$$

در این صورت، وقتی $n \rightarrow \infty$ ، T_n دارای توزیع غیر مرکزی کای دو با q درجه آزادی و پارامتر غیر مرکزی $\vartheta = \vartheta^\top (CI^{-1}C^\top)^{-1} \vartheta$ خواهد بود. حال می‌توان تحت $H_{\circ(n)}$ و زمانی که $n \rightarrow \infty$ اثبات کرد که

$$\begin{pmatrix} \sqrt{n}(\hat{\beta}_{\text{CLE}} - \beta) \\ \sqrt{n}(\hat{\beta}_{\text{RCLE}} - \beta) \\ \sqrt{n}(\hat{\beta}_{\text{CLE}} - \hat{\beta}_{\text{RCLE}}) \end{pmatrix} \sim N_{3p} \left(\begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}, \begin{pmatrix} BI^{-1}B^\top & B(I^{-1} - J)B^\top & BJB^\top \\ B(I^{-1} - J)B^\top & B(I^{-1} - J)B^\top & 0 \\ BJB^\top & 0 & BJB^\top \end{pmatrix} \right) \quad (5.3)$$

که در آن $\gamma = B\vartheta$ ، $J = BCI^{-1}$ ، $b_2 = b_1 - B\gamma$ و $b_3 = B\gamma$ (**کیبیریا و صالح، ۲۰۰۴**).

حال اگر اربیی مجانبی را به صورت زیر تعریف کنیم

$$\mathbb{B}(\hat{\beta}) = \lim_{n \rightarrow \infty} \mathbb{E}(\sqrt{n}(\hat{\beta} - \beta)), \quad (6.3)$$

در این صورت با استفاده از رابطه (۵.۳)، برای برآوردهای براساس برآوردگر ليو خواهیم داشت:

$$\mathbb{B}(\hat{\beta}_{\text{CLE}}) = \mathbf{b}_1, \quad \mathbb{B}(\hat{\beta}_{\text{RCLE}}) = \mathbf{b}_2, \quad \mathbb{B}(\hat{\beta}_{\text{CLPTE}}) = \mathbf{b}_1 - \mathbf{b}_2 \Upsilon_q(\Delta; \chi_{(q,\alpha)}^2).$$

که در آن $\Upsilon_q(\Delta; z)$ نمایانگر تابع توزیع تجمعی در نقطه z برای توزیع غیر مرکزی کای دو با q درجه آزادی و پارامتر غیر مرکزی Δ می باشد. همچنین، اگر واریانس برآوردگر را به فرم

$$\mathbb{V}(\hat{\beta}) = \lim_{n \rightarrow \infty} \mathbb{E}(\sqrt{n}(\hat{\beta} - \beta)\sqrt{n}(\hat{\beta} - \beta)^\top).$$

تعریف کنیم، در این صورت داریم:

$$\mathbb{V}(\hat{\beta}_{\text{CLE}}) = \mathbf{B}\mathbf{I}^{-1}\mathbf{B}^\top + \mathbf{b}_1\mathbf{b}_1^\top,$$

$$\mathbb{V}(\hat{\beta}_{\text{RCLE}}) = \mathbf{B}(\mathbf{I}^{-1} - \mathbf{J})\mathbf{B}^\top + \mathbf{b}_2\mathbf{b}_2^\top,$$

$$\mathbb{V}(\hat{\beta}_{\text{CLPTE}}) = \mathbb{V}(\hat{\beta}_{\text{CLE}}) - \left[\mathbf{b}_2(\mathbf{b}_1 - \mathbf{b}_2)\mathbf{b}_2^\top + \mathbf{B}\mathbf{J}\mathbf{B}^\top \right] \Upsilon_{q+2}(\Delta; \chi_{(q,\alpha)}^2) - \mathbf{b}_2\mathbf{b}_2^\top \Upsilon_{q+4}(\Delta; \chi_{(q,\alpha)}^2),$$

۴ شبیه سازی

در این بخش، با استفاده از شبیه سازی مونت کارلو به ارزیابی عملکرد برآوردهای پیشنهادی خواهیم پرداخت. برای بررسی عملکرد برآوردها در حضور همخطی در مدل، متغیرهای کمکی را مطابق مدل زیر تولید می کنیم:

$$x_{ij} = \sqrt{1 - \rho^2} Z_{ij} + \rho Z_{i(p+1)}; \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, p, \quad (1.4)$$

که در آن Z_{ij} نمونه های تصادفی از توزیع نرمال استاندارد می باشد. پارامتر ρ شدت همخطی بین متغیرهای کمکی را کنترل می کند و مقدار آن را $p = \{0.5, 0.7, 0.95\}$ در نظر گرفته می گیریم. همچنین $n = \{50, 100, 200\}$ ، $p = \{6, 12\}$ و احتمال سانسور را $cp = \{0.2, 0.3\}$ و زمان های بقاء پایه را توزیع نمایی با پارامتر نرخ 0.5 فرض می کنیم. مقادیر واقعی ضرایب را به صورت $\beta = (0.5, 1, -0.6, 0.8, 0, 0, \dots, 0)^\top$ تعریف می کنیم بدین معنا که، $p - 4$ ضریب آخر بردار β برابر صفرند و فرض صفر مدل را بر اساس این پارامترهای صفر در مدل در نظر می گیریم. در نهایت، داده ها برای CPHM مطابق روش توضیح داده شده توسط استین (۲۰۱۲) تولید می کنیم. علاوه بر این، برای بررسی اثر انحراف از مقدار واقعی ضرایب واقعی، بردار \mathbf{b} در فرض صفر را به صورت $\mathbf{b} = (\delta, 0_{p-5}^\top)^\top$ تعریف می کنیم که در آن $\delta \in [0, 4]$. این رابطه نشان می دهد که $\delta = 0$ متناظر با مقادیر واقعی ضرایب است و مقادیر دیگر نشان دهنده انحراف از مقادیر واقعی می باشند. برای ارزیابی عملکرد برآوردها، تولید داده ها را ۱۰۰۰ بار تکرار کرده و کارایی نسبی برآوردهای پیشنهادی را که به صورت زیر تعریف می شود، محاسبه می کنیم:

$$RE(\hat{\beta}) = \frac{\sum_{m=1}^{\infty} (\hat{\beta}_{CLE}^{(m)} - \beta)^{\top} (\hat{\beta}_{CLE}^{(m)} - \beta)}{\sum_{m=1}^{\infty} (\hat{\beta}_{*}^{(m)} - \beta)^{\top} (\hat{\beta}_{*}^{(m)} - \beta)}, \quad (2.4)$$

که در آن $\hat{\beta}_{*}^{(m)}$ نشان‌دهنده هر یک از برآوردگرهای پیشنهادی در تکرار m ام است. از نتایج شبیه‌سازی که در جدول ۱ گزارش شده است. با توجه به نتایج، زمانی که مقدار $\delta = 0$ باشد، برآوردگرهای پیشنهادی عملکرد بهتری نسبت به سایرین دارند. اما با فاصله گرفتن δ از مقدار واقعی، کارایی آن‌ها کاهش یافته و به‌ویژه برآوردگر لیو مقید به شدت افت عملکرد دارد، در حالی که کارایی برآوردگر پیش‌آزمون به مقدار یک میل می‌کند. همچنین، با افزایش اندازه نمونه، کارایی نسبی برآوردگرها کاهش می‌یابد، اما افزایش تعداد متغیرها موجب بهبود کارایی نسبی آن‌ها می‌شود. افزایش ضریب همبستگی بین متغیرها نیز منجر به ارتقای کارایی نسبی می‌گردد. در نهایت، با افزایش احتمال سانسور، هرچند مقادیر کارایی نسبی اندکی افزایش می‌یابد، اما روند کلی نتایج تغییر چندانی نمی‌کند.

جدول ۱: کارایی برآوردگرهای پیشنهادی نسبت به برآوردگر لیو در مدل CPHM.

cp = 0.2								cp = 0.3					
n	δ	ρ = 0.5		ρ = 0.7		ρ = 0.95		ρ = 0.5		ρ = 0.7		ρ = 0.95	
		RCLE	CLPTE	RCLE	CLPTE	RCLE	CLPTE	RCLE	CLPTE	RCLE	CLPTE	RCLE	CLPTE
p = 6													
50	0	1.625	1.371	1.686	1.415	1.610	1.380	1.634	1.386	1.685	1.419	1.599	1.394
	1	0.225	0.969	0.281	0.941	0.743	0.811	0.251	0.944	0.314	0.889	0.798	0.846
	2	0.064	1.000	0.082	1.000	0.309	0.872	0.072	1.000	0.094	1.000	0.350	0.835
	3	0.029	1.000	0.038	1.000	0.159	0.982	0.033	1.000	0.044	1.000	0.184	0.970
100	4	0.017	1.000	0.022	1.000	0.095	1.000	0.022	1.000	0.019	1.000	0.111	0.997
	0	1.526	1.372	1.592	1.423	1.638	1.472	1.510	1.377	1.572	1.423	1.615	1.460
	1	0.108	1.000	0.135	1.000	0.447	0.772	0.124	1.000	0.154	1.000	0.490	0.757
	2	0.029	1.000	0.036	1.000	0.147	0.993	0.033	1.000	0.042	1.000	0.167	0.987
200	3	0.013	1.000	0.016	1.000	0.070	1.000	0.015	1.000	0.019	1.000	0.080	1.000
	4	0.007	1.000	0.009	1.000	0.040	1.000	0.009	1.000	0.018	1.000	0.047	1.000
	0	1.493	1.346	1.576	1.386	1.689	1.460	1.482	1.337	1.563	1.373	1.667	1.444
	1	0.052	1.000	0.065	1.000	0.251	0.925	0.058	1.000	0.073	1.000	0.276	0.898
p = 12	2	0.013	1.000	0.017	1.000	0.072	1.000	0.015	1.000	0.019	1.000	0.081	1.000
	3	0.006	1.000	0.008	1.000	0.033	1.000	0.007	1.000	0.009	1.000	0.037	1.000
	4	0.003	1.000	0.004	1.000	0.019	1.000	0.004	1.000	0.005	1.000	0.021	1.000
	0	4.504	2.525	4.657	2.544	3.689	2.287	4.400	2.512	4.552	2.642	3.580	2.294
50	1	0.639	0.972	0.788	0.965	1.751	1.254	0.705	0.967	0.868	0.971	1.831	1.298
	2	0.185	1.000	0.236	1.000	0.754	0.936	0.208	1.000	0.267	0.999	0.829	0.937
	3	0.085	1.000	0.110	1.000	0.394	0.967	0.096	1.000	0.125	0.999	0.442	0.965
	4	0.048	1.000	0.110	1.000	0.238	1.000	0.055	1.000	0.072	1.000	0.270	0.996
100	0	3.813	2.491	4.117	2.614	4.038	2.666	3.824	2.452	4.103	2.554	3.948	2.567
	1	0.261	1.000	0.335	0.995	1.102	0.972	0.296	1.000	0.379	0.992	1.199	0.999
	2	0.069	1.000	0.090	1.000	0.363	0.976	0.079	1.000	0.103	1.000	0.410	0.958
	3	0.031	1.000	0.041	1.000	0.173	1.000	0.036	1.000	0.047	1.000	0.197	1.000
200	4	0.018	1.000	0.023	1.000	0.100	1.000	0.020	1.000	0.027	1.000	0.114	1.000
	0	3.326	2.463	3.552	2.660	3.924	2.785	3.179	2.462	3.464	2.593	3.815	2.789
	1	0.110	1.000	0.143	1.000	0.570	0.902	0.125	1.000	0.164	1.000	0.635	0.878
	2	0.028	1.000	0.037	1.000	0.164	1.000	0.032	1.000	0.043	1.000	0.109	1.000
p = 12	3	0.013	1.000	0.017	1.000	0.075	1.000	0.015	1.000	0.019	1.000	0.086	1.000
	4	0.007	1.000	0.009	1.000	0.043	1.000	0.008	1.000	0.011	1.000	0.049	1.000

۵ تحلیل داده‌های بقا

در این بخش، برای ارزیابی عملکرد برآوردگرهای پیشنهادی، داده‌های مربوط به بقای بیماران مبتلا به سرطان پیشرفته ریه از گروه درمانی سرطان North Central تحلیل خواهیم کرد. این داده‌ها، موجود در بسته آماری survival نرم‌افزار R با عنوان cancer است که شامل اطلاعاتی درباره‌ی توانایی بیماران در انجام فعالیت‌های روزمره است که با استفاده از نمرات عملکردی ارزیابی شده‌اند. پس از حذف رکوردی دارای مقادیر گمشده، مجموعه داده نهایی شامل ۱۶۷ بیمار می‌باشد. متغیرهای موجود در این مجموعه داده شامل زمان بقا بر حسب روز، وضعیت بیمار، سن، جنسیت، نمره عملکردی کارنوفسکی بر اساس ارزیابی پزشک (ph.karno)، نمره عملکردی ECOG بر اساس نظر پزشک (ph.ecog) (۰ بدون علامت، ۱ علامت‌دار اما کاملاً قادر به حرکت، ۲ بستری کمتر از ۵۰٪ از روز، ۳ بستری بیش از ۵۰٪ از روز ولی نه کاملاً بستری، ۴ کاملاً بستری) و کاهش وزن در شش ماه گذشته می‌باشد.

نتایج آزمون نسبت درست‌نمایی و آزمون والد برای برازش CPHM بر این داده‌ها، حاکی از برازش مناسب مدل با p -مقدار به ترتیب برابر با $10^{-5} \times 9$ و 10^{-4} می‌باشند. بررسی همبستگی میان متغیرهای کمکی، همبستگی قوی به میزان 0.811 بین متغیرهای ph.ecog و ph.karno را نشان می‌دهد. همچنین عدد شرطی محاسبه‌شده برابر با 223.94 بوده که وجود همخطی را تایید می‌کند. با توجه به نبود اطلاعات پیشین درباره ضرایب مدل، آزمون ضرایب انجام شد که نشان داد متغیرهای سن و کاهش وزن از نظر آماری معنادار نیستند و ضرایب آن‌ها صفر در نظر گرفته شد. برای ارزیابی عملکرد برآوردهای پیشنهادی مبتنی بر لیو، از روش بوت‌استرپ با اندازه نمونه 50 و 1000 تکرار استفاده شد. نتایج جدول ۲ نشان می‌دهد که این برآوردها دارای انحراف معیار کمتر و کارایی نسبی بالاتر از یک هستند، که به معنای برتری آن‌ها نسبت به برآوردهای لیو است؛ یافته‌ای که با نتایج شبیه‌سازی نیز هم‌خوانی دارد.

جدول ۲: برآورد، انحراف معیار sd و کارایی نسبی برآوردها برای داده‌های سرطان ریه.

CLPTE		RCLE		CLE		متغیرها
sd	برآورد	sd	برآورد	sd	برآورد	
0.274	-0.512	0.264	-0.495	0.284	-0.551	جنسیت
0.23	0.10	0.10	0.22	0.24	0.11	ph.ecog
0.10	0.352	0.01	0.345	0.13	0.356	ph.karno
1.217		1.232		-		RE

بحث و نتیجه‌گیری

در این مقاله به بررسی کارایی برآوردها که براساس برآوردهای لیو در مدل کاکس با وجود همخطی بین متغیرهای کمکی پرداخته می‌شود. ضمن معرفی برآوردها، ویژگی‌های مجانبی آن مانند اریبی و واریانس استخراج کرده و با انجام مطالعات شبیه‌سازی گسترده، عملکرد آن‌ها را در شرایط مختلف ارزیابی نمودیم. نتایج نشان داد که این برآوردها پیشنهادی به طور پیوسته کارایی نسبی بالاتری نسبت به برآوردهای اصلی لیو کسب کرده به‌ویژه در شرایطی که همخطی چندگانه شدید باشد، وجود دارد. در نهایت، مزایای عملی رویکرد پیشنهادی با استفاده از مجموعه داده سرطان ریه NCCTG نشان داده شد. این یافته‌ها بر ارزش استفاده از برآوردهای پیش‌آزمون مبتنی بر لیو در ارائه استنباط‌های قابل اطمینان‌تر در تحلیل بقا تأکید دارند.

بیانیه تأمین مالی

تحقیق محمد آرشی تحت حمایت مادی بنیاد ملی علم ایران INSF برگرفته شده از طرح شماره 4015320 انجام شده است.

مراجع

- Ahmad S., Aslam M., Ahmad S. (2023), *Extending the Liu estimator for the Cox proportional hazards regression model with multicollinearity* Commun. Stat. Simul. Comput., **53**(12), 5828–5841.
- Austin PC. (2012), *Generating survival times to simulate Cox proportional hazards models with time-varying covariates*, **31**(29), 3946–3958.

Cox DR. (1972), *Regression models and life-tables (with discussion)*, J. R. Stat. Soc., B: Methodol. ;**34**(1):187–220.

Kibria B.M.G. and Saleh A. (2004), *Performance of positive rule estimator in the ill-conditioned Gaussian regression model*, Calcutta Statistical Association Bulletin, **55**, 209–240.

Magnus J.R. (1997), *The traditional pretest estimator*, Theory of Probability and Its Applications, **44**(2), 293–308.

Investigation in PTE Asymptotic Behavior in the Present of Multicollinearity

Seyed Amirhossein Tabatabaei Shirazi, Mahdi Emadi, Mohammad Arashi, Solmaz Seifollahi

Department of Statistics, Ferdowsi University of Mashhad, Mashhad, Iran

Abstract: The Cox proportional hazards regression model is widely used for analyzing survival data with covariates. However, its performance can deteriorate in the presence of multicollinearity, leading to unreliable estimates when using the standard partial likelihood approach. In this paper, we extend the pretest estimator based on the Liu method to improve the accuracy of coefficient estimation. We derive the asymptotic theoretical properties of the proposed estimator and evaluate its performance through extensive Monte Carlo simulations. We also demonstrate how this estimation strategy can be applied to real survival data. The results show significant improvements, highlighting the practical advantages of this estimator for researchers working in survival analysis.

Keywords: Cox proportional hazards model, Liu estimator, Multicollinearity, Pretest estimator.

Mathematics Subject Classification (2020): 62N02, 62F30.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



چگونه نظریه فرآیندهای تصادفی به مدل گسترش یافته در سنتز مدرن تکامل کمک کرد؟

دکتر حمیدرضا عریضی^۱

^۱ گروه روانشناسی، دانشگاه اصفهان

چکیده: شاید مهم‌ترین نظریه‌ای که همه مفاهیم انسانی و اجتماعی را تغییر داد، ظهور تکامل داروین و ژنتیک مندلی باشد که در فاصله شش سال در قرن نوزدهم ارائه شد و تفکر بشر را در مورد انسان، محیط او، همه حیوانات و حتی گیاهان تغییر داد. جولیان هاکسلی در سال ۱۹۴۲ این دو نظریه را در یک الگوی ترکیبی با هم پیوند داد که نام آن را سنتز مدرن نهاد. این الگوی ترکیبی کمک کرد که این نظریه پا برجا باقی بماند. در مورد جهش تکامل دو دیدگاه لامارک و داروین (دانشمندان فرانسوی و انگلیسی) وجود دارد که با هم متفاوت است. در این مقاله، به کمک فرآیندهای تصادفی انواع جهش‌ها و طبقه‌بندی آن‌ها را تشریح می‌کنیم. این طبقه‌بندی به نظریه سنتز مدرن کمک کرده است تا هم‌چنان پیوند تکامل و ژنتیک استوار بماند.

واژه‌های کلیدی: فرآیندهای تصادفی - تکامل داروین - ژنتیک مندلی - جهش - طبقه‌بندی جهش‌ها.

کد موضوع‌بندی ریاضی (۲۰۲۰): 60J70.

۱ کاربرد فرآیندهای تصادفی در مدل گسترش یافته سنتز مدرن تکامل

شاید فرآیندهای تصادفی در هیچ علمی به اندازه زیست‌شناسی (تکامل) کاربرد نداشته باشد. فرآیندهای تصادفی در سه زمینه ژنوتیپ، فنوتیپ و برازش نقش اصلی را در تبیین آن داشته است. (سارکار، ۱۹۹۲).

^۱ سخنران، Dr.oreyzi@edu.ui.ac.ir

۱.۱ ارتباط نظریه فرآیندهای تصادفی به مدل گسترش یافته در سنتز مدرن تکامل

داروین در کتاب دوران ساز خود منشاء انواع که در ۱۸۵۹ منتشر شد، رویکرد تکاملی خود را با سه اصل مهم انتخاب طبیعی، انتقال صفات و فرآیندهای تصادفی شرح داد. این در هم تنیدگی با نظریه احتمال و فرآیندهای تصادفی سبب شد که آماردانان برجسته‌ای همچون **فیشر (۱۹۳۰)** و **رایت (۱۹۳۱)** به دفاع از نظریه تکامل برخیزند. این تلاش‌ها موجب شده است که احتمال و فرآیندهای تصادفی در ردیف محورهای بسیار مهمی در دفاع از نظریه تکامل قرار گیرند (**سارکار، ۱۹۹۲**). یکی از این تلاش‌ها در پوشش شرط‌بندی بیولوژیکی است. پوشش شرط‌بندی بیولوژیکی زمانی اتفاق می‌افتد که برازش اورگانیسم‌ها در شرایط معمولی کاهش یافته، هرچند در تبادل در شرایط پر فشار بایستی برازش خود را افزایش دهند. این مفهوم زمینه را برای بحث در مورد فرآیند تصادفی در آثار فیشر و رایت گشود (**کورنینگ و کافمن، ۲۰۲۳**).

موقعیتی را در نظر بگیرید که در هر سال افراد دو راهبرد تولید مثل با احتمالات a و $1-a$ را بکار ببرند. واریانس باروری (**لمان و بالو، ۲۰۰۷**) را با تابع F نشان می‌دهیم. راهبرد تولید مثل اول به کیفیت سال بستگی دارد. در سال‌هایی که کیفیت باروری بالاست، انحراف معیار باروری ($F(1(\sigma))$ و در سال‌هایی که کیفیت باروری پایین است، انحراف معیار باروری ($F(1+\sigma)$) است. راهبرد تولید مثل دوم در هر سال میزان باروری به میزان ثابتی به اندازه $F(1(c))$ کاهش می‌یابد (مقدار ثابت $c = \sigma^2$ و واریانس $\sigma^2 =$). حال یک جهش^۱ با بسامد P با راهبرد متفاوت $a + da$ را در نظر بگیرید. می‌توان تغییر در بسامد جهش در سال‌های خوب و بد را طبق فرمول (**۱.۱**) به دست آورد.

$$\begin{aligned}\Delta P_{good} &= S_{good}P(1-P) + O(da)^2 \\ \Delta P_{bad} &= S_{bad}P(1-P) + O(da)^2\end{aligned}\quad (1.1)$$

که در آن برازش خوب و بد (w_{good} و w_{bad}) برای da (راهبرد متفاوت) در هر نوع سال با فرمول‌های (**۲.۱**) تعیین می‌شود.

$$\begin{aligned}w_{good} &= 1 + S_{good} \\ w_{bad} &= 1 + S_{bad}.\end{aligned}\quad (2.1)$$

در مدل رایت و فیشر وقتی طبقات افراد در تکامل در نظر گرفته می‌شوند که بر مبنای برازش یافتن با وضعیت جهش است. O توزیع افراد در طول این طبقه‌ها^۲، d مدل زمان گسسته^۳ و S ضریب انتخاب است (**لمان و بالو، ۲۰۰۷**). **لمان و بالو (۲۰۰۷)** برای میانگین استراتژی دو راهبرد خوب و بد مقدار زیر را بدست آورده‌اند: (σ^2 واریانس و مقدار ثابت c در بالا تعریف شده‌اند).

$$da \frac{c(1-c(1-a)) - a\sigma^2}{(1-c(1-a))^2 - a^2\sigma^2}.$$

درک داروین از نظریه احتمال، متناسب با دانش زمان خودش قرن ۱۹ و پیر لاپلاس بود و هنوز اصول مهم کلموگروف پدیدار نشده بود. در فلسفه علم، پوپر بیش از همه به نقش فرآیندهای تصادفی و نظریه احتمال در نظریه تکامل داروین پرداخته است و این ارجاع به نظریه احتمال و فرآیندهای تصادفی را علاوه بر مشاهدات داروین که در کتاب منشاء انواع به آن اشاره کرده بود را ملاک علمی بودن

^۱ Mutation bias

^۲ Occupancy process

^۳ Discrete time

نظریه تکامل دانسته است. اما داروین از طریق فرانسویس گالتون (۱۸۲۲-۱۹۱۱) که آماردان بود توانست اهمیت احتمال را برای بیان نظریه تکامل پیدا کند (گالتون، ۱۸۸۹). اساس استدلال داروین بر چهار گزاره متکی است:

- ۱- در هر گونه خاص میان ویژگی تک تک موجودات تنوع هست.
 - ۲- ویژگی‌های والدین احتمالاً به فرزندان‌شان منتقل می‌شود.
 - ۳- گونه‌های جانداران به طور غریزی توانایی آن را دارند که با آهنگی هندسی تکثیر شوند.
 - ۴- به طور معمول منابع محیط برای چنین آهنگ رشدی بسنده نیست.
- گزاره دوم داروین اشاره می‌کند که ویژگی‌های والدین احتمالاً به فرزندان‌شان منتقل می‌شود. او این آگاهی را به صورتی ناهشیار دارد. تنها شش سال بعد مندل در ژنتیک مندلی ساز و کار آن را در می‌یابد. کتاب داروین در سال انتشار به پرفروش‌ترین کتاب در انگلستان تبدیل شد و با شور فراوان مورد استقبال قرار گرفت.^۴ بیشتر افرادی که کتاب او را می‌خواندند مشاهدات او که در جزایر آمریکای لاتین (و به خصوص گالاپاگوس) رخ داده بود را جالب می‌یافتند اما به مفاهیم او در نظریه احتمالات و فرآیندهای تصادفی کمتر توجه می‌کردند. اما همین اشارات متخصصان علم آمار را از همان آغاز تا امروز برمی‌انگیخت تا مانع آن شوند که بر تابوت این نظریه^۵ آخرین میخ‌ها زده شود. امروزه غیر از شواهد زیست‌شناسی، از شبیه‌سازی آزمایش رایانه‌ای برای بررسی فرآیندهای تصادفی و احتمالات (گزاره ۲) استفاده می‌شود (مود و اسلیم، ۲۰۱۲، تژادا-لاپورتا و همکاران، ۲۰۲۵).
- در برابر این نظریه به خصوص مخالفت داروین با غایت نگر^۶ و به خصوص رویکرد افلاطونی طراحی هوشمند^۷ مقاومت جامعه مذهبی در پذیرش این نظریه را برمی‌انگیخت تا اینکه یک کشیش به نام گرگور مندل^۸ از کلیسایی در شهر برنو (اینک دومین شهر بزرگ چک) مبنای استدلالی برای دو گزاره ۱ و ۲ فراهم می‌ساخت که با پرورش لوبیا در باغ حیاط کلیسا اولین اشاره‌ها به ژن را انجام داد. یکی از ایده‌های مهم تمایز مندل بین ژن‌های غالب و مغلوب بود که چرا با وجود آنکه گاهی برخی ویژگی‌های پدر بزرگ و مادر بزرگ‌ها که در فرزندان بلافاصله آن‌ها تجلی پیدا نکرده است در نوه‌های آن‌ها ظاهر می‌شود.
- نظریه مندل نظریه بسیار مهمتری نسبت به نظریه داروین بود و هنگامی که او در ۶۱ سالگی در ۱۸۸۴ از دنیا رفت فقط تعداد اندکی که بولتن خبری کلیسا در ۱۸۶۵ را خوانده بودند از کار او آگاهی داشتند. در آن زمان نظریه داروین با ژنتیک مندلی متضاد در نظر گرفته می‌شدند. اینکه گزاره ۲ داروین به خوبی با کشف مندل تبیین می‌شود کمتر مورد توجه بود. برای این تضاد بین درک داروینی و مندلی به کتاب *پیتر بولر* (۱۹۸۹) نگاه کنید که اینک خوشبختانه به فارسی توسط محی‌الدین غفرانی ترجمه شده است. با این حال توضیح این که چرا بین ژن‌های موجود درون جمعیت یک گونه واحد تفاوت و تنوع وجود دارد نه توسط داروین و نه توسط مندل قابل تشخیص نبود و همچنان به صورت یک معما باقی مانده بود. این متخصصان احتمال و فرآیندهای تصادفی بودند که علاوه بر اینکه پاسخی برای این معما یافتند، آن را در زمینه‌های گوناگون توسعه دادند (سود و میستلی، ۲۰۲۲). باید توجه کرد که مفهوم تصادف در زیست‌شناسی تکامل کاملاً خاص و مربوط به تطابق و سازگاری است (گارسایپینو، ۲۰۲۴). به طور کلی مستقل از ژنوتیپ و نیز محیط که تفاوت‌های فردی در آن رخ می‌دهد، می‌توان نظریه تکامل را در قالب فرآیندهای تصادفی مطرح کرد (لرماند و همکاران، ۲۰۰۹) به همین دلیل فرآیندهای تصادفی در نظریه مدرن تکامل جایگاه بسیار اساسی دارد.

^۴ این کتاب چهار دهه پیش توسط نورالدین فرهیخته در ایران منتشر شد.

^۵ بسیاری از دانشمندان زمان داروین آن را نظریه‌ای می‌پنداشتند که در حال مرگ است و تقریباً تا زمان سنتز مدرن درست می‌گفتند.

^۶ Teleological

^۷ Intelligent design

^۸ George Mandle

به طور کلی فرآیندهای تصادفی به سه صورت در تبیین نظریه تکامل به کار رفته است:

۱- فرآیندهای تصادفی در جهش و تغییر^۹

۲- فرآیندهای تصادفی در چرخه حیات^{۱۰}

۳- فرآیندهای تصادفی در تغییرات محیطی^{۱۱}

این سه مدل تکاملی به کمک هم می‌توانند تکامل را توضیح دهند. امروزه بیان نظریه تکامل بدون کمک فرآیندهای تصادفی بسیار دشوار است که خود کاربرد مهم این فرآیندها را نشان می‌دهد. قبل از بیان این سه نوع فرآیند تصادفی دو مدل سنتز مدرن و گسترش سنتز مدرن که در واقع به کمک مفهوم فرآیندهای تصادفی ایجاد شده‌اند را توضیح می‌دهیم. در سنتز مدرن مفهوم رانش ژنتیکی^{۱۲} که در کارهای **فیشِر** و **رایت** مورد استفاده قرار گرفته بود برای ترکیب و سنتز مفهوم تکامل در نظریه داروین و ژن در نظریه مندل به کار رفت. رانش ژنتیکی هر گونه تغییری است که در جمعیت پیش آید به نحوی که فراوانی آلل‌های^{۱۳} خاصی در جمعیت فزونی یابد. از این ترکیب داروین‌گرایی نوین^{۱۴} پدید آمد که در آن با استناد به دانسته‌های فعلی ما درباره ژن‌ها و کروموزوم‌ها به توضیح منشأ و علل تنوع ژنتیکی می‌پردازد که این خود ماده خام انتخاب طبیعی است. سنتز مدرن بعداً با مشمول کردن فرآیندهای تصادفی در تغییرات محیطی شکل گسترش یافته سنتز مدرن را پدید آورد. این شکل گسترش یافته کمک کرد تا هر دو نظریه داروین و مندل در کنار هم رشد کنند. بنابراین فرآیندهای تصادفی زبان پیوند بین دو نظریه است. در جهش و تغییر مفهوم فرآیندهای تصادفی به این صورت است که جهش‌ها از تأثیرات فنوتیپیک^{۱۵} خود به طور مستقل رخ می‌دهند و معطوف به آن نیستند. معنی این که آن‌ها تصادفی هستند این است که اعم از این که جهش برای ارگانیسم مفید یا مضر باشد، رخ می‌دهند. این مسئله که بسیاری از جهش‌ها برای ارگانیسم مخرب هستند امروزه امر شناخته شده‌ای است (نگاه کنید به **آیرواکر و کیتلی (۲۰۰۷)**). این موقعیت را می‌توان در مدل‌های هندسی تطابق و سازگاری^{۱۶} که در آن جهش‌ها اثرات فنوتیپیک نااریب دارند یافت (اور، ۲۰۰۵، مارتین و لئورماند، ۲۰۰۶). با وجود آن که میانگین اثرات برازشی جهش‌ها مخرب‌تر می‌شوند جمعیت به حالت بهینه فنوتیپی نزدیک‌تر می‌شود. **یامپولسکی و استولفتوتز (۲۰۰۱)** به حالتی در تکامل پرداخته‌اند که حجم جمعیت N و میزان جهش چندان بزرگ نباشد (محدود باشد). در این حالت جهش‌ها به سرعت تغییر می‌کنند و ترتیب ظهور جهش‌ها هم بسیار تغییر می‌کند. **یامپولسکی و استولفتوتز** به کمک فرآیندهای تصادفی، این دو عامل متغیر را سازگار با انتخاب طبیعی مدلسازی کرده‌اند. در ادامه مقاله **یامپولسکی و استولفتوتز (۲۰۰۱)**، دو نظریه احتمال تکامل موازی^{۱۷} (اور، ۲۰۰۵) و تکامل با برازش تصادفی (**جانسون و بارتون، ۲۰۰۲**) به رابطه جهش‌ها و انتخاب طبیعی پرداخته‌اند.

دومین کاربرد فرآیندهای تصادفی در تکامل مربوط به چرخه حیات است. مدل‌های تکاملی به طور ضمنی به چرخه حیات افراد مربوط می‌شوند که در آن رویدادهایی مثل تولد، تولید مثل و مرگ وجود دارد. مقیاس معمول مشاهده برای زیست‌شناسانی که جمعیت‌ها را مطالعه می‌کنند رویدادهای تصادفی مستقل برای هر فرد است. مثلاً بقای هر فرد از زمان t به زمان $t + 1$ را می‌توان به کمک احتمال

^۹ Stochasticity of Mutation and variation

^{۱۰} Stochasticity of life histories

^{۱۱} Stochasticity of environmental change

^{۱۲} Genetic drift

^{۱۳} Allele، هر نسخه از یک ژن، در جانداران از هر ژن دو آلل وجود دارد که از پدر و مادر به ارث می‌رسد.

^{۱۴} Neo-Darwinism

^{۱۵} Phenotypic effects

^{۱۶} Geometrical models of adaptation

^{۱۷} Parallel evolution

نشان داد (کورنینگ و کافمن، ۲۰۲۳). در عین حال تعداد فرزندان یک فرد و نیز تعداد گامت‌های یک ژنوتیپ را می‌توان با بازترکیب و جداسازی از توزیع احتمال (توزیع چندجمله‌ای و توزیع پواسون) بیرون کشید. این نوع فرآیندهای تصادفی می‌تواند منجر به تغییر تصادفی در اندازه جمعیت شود (فرآیندهای تصادفی جمعیت شناختی). در مدل‌های ژنتیک جمعیتی اثر این فرآیندهای تصادفی در چرخه حیات واریانس تغییر در بسامدهای ژنوتیپ از یک نسل به نسل دیگر را نشان می‌دهد که بزرگی اندازه آن به تعداد افراد بستگی دارد که به آن رانش ژنتیکی گویند. در این مورد دو مدل رایج- فیشر و مدل موران پرکاربردترین مدل‌ها هستند. در مدل استاندارد رایت-فیشر نسل‌های ناهمپوش و گسسته ارائه می‌شود که در هر نسل افراد فرزندان تولید می‌کنند و می‌میرند. برای پدیدآیی نسل بعدی، تعداد (مفروض) افراد جایگزین والدین (مستقل از همدیگر) می‌شوند. که متناظر با نمونه‌گیری چندجمله‌ای^{۱۸} از ژنوتیپ‌هاست. در مدل موران (۱۹۵۸) جمعیت در زمان‌های گسسته $t = 1, 2, 3, \dots$ نمایش داده می‌شود. فردی برای تولید مثل و فرد دیگری برای مرگ گزینش می‌شود. بنابراین این مدل جزء نمونه‌های مدل زاد و مرگ^{۱۹} است که در فرآیندهای تصادفی بسیار مطالعه شده است. در این جا نسل‌ها همپوشانی دارند و میانگین طول زندگی یک فرد در یک جمعیت به اندازه $2N$ شامل $2N$ مرحله (گام) زمانی است. می‌توان نشان داد که اگر p بسامد نوع A در یک زمان مفروض باشد واریانس بسامد $2N$ مرحله یا گام بعدی pq/N است. بنابراین در مدل موران اندازه جمعیت مؤثر $N_e = N/2$ است. در واقع بسیاری از نتایجی که در مدل فیشر- رایت بدست آمده است را می‌توان در مدل موران با جایگزینی $N/2$ به جای N بدست آورد (تژاد-لاپورتا و همکاران، ۲۰۲۵).

سومین فرآیند تصادفی مربوط به تغییرات محیطی است و در حالی که در فرآیندهای تصادفی چرخه حیات افراد مستقل از یکدیگر می‌باشند در تغییرات محیطی تأثیر محیط بر همه افراد در یک جمعیت وجود دارد. هرچند محیط ثابت نیست اما می‌توان در یک زمان آن را دارای ماهیت تناوبی (مثلاً تغییر فصل) دانست. گاهی عوامل محیطی (مثل آب و هوا) کمتر پیش‌بینی پذیر هستند، هرچند ممکن است در بدو امر فرآیندهای تصادفی غیر زیستی (مثل بروز آتشفشان) به ذهن آیند اما فرآیندهای تصادفی زیستی بسیار بیشتر هستند (در این زمینه می‌توان به مقاله لی اسمیت و همکاران (۲۰۲۵) که در آن کاربرد فرآیندهای زیستی در تغییرات محیطی دیده می‌شود مراجعه کرد). در مورد تغییرات محیطی مقیاس زمان^{۲۰} بسیار مهم است. در این حالت نوسانات محیطی نسبت به زمان نسلی^{۲۱} بسیار سریع‌تر است. اگر نوسانات فرآیند تصادفی در محاسبه میانگین افراد برای اثر فرآیندهای تصادفی بر ارگانیزم بسیار کوچک باشد، بهتر است نوسانات فرآیند تصادفی در تغییر محیطی در مقیاس‌های زمانی بزرگ مورد مطالعه قرار گیرد (سود و میستلی، ۲۰۲۲).

فرآیندهای تصادفی در فرآیند تکامل در چهار موقعیت نقش اساسی دارند که به ترتیب عبارتند از :

الف) ناسازگاری

ب) سرشت تصادفی در جهش‌های خنثی

ج) انقلاب‌هایی تکاملی

د) شکل‌دهی به انتخاب طبیعی

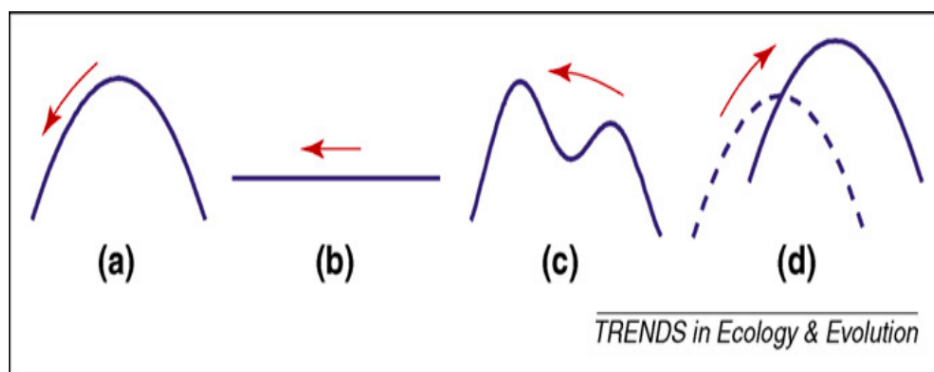
نقش فرآیندهای تصادفی در این چهار موقعیت اساسی در شکل ۱ ترسیم شده است. نخست اینکه انتخاب تصادفی همواره در جهت سازگاری و مثبت نبوده و همه شکل‌های گفته شده در فرآیندهای تصادفی حدودی برای سازگاری تعیین می‌کنند (به کسرگلو و همکاران

^{۱۸} Multinomial sampling

^{۱۹} Birth death models

^{۲۰} Timescale

^{۲۱} Generation time



شکل ۱: فرآیندهای تصادفی در فرآیند تکامل، (a) ناسازگاری، (b) سرشت تصادفی در جهش‌های خنثی، (c) انقلاب‌های تکاملی، (d) شکل‌دهی به انتخاب طبیعی.

(۲۰۲۳) نگاه کنید)، به همین دلیل در ترجمه فارسی گاهی به جای تکامل از فرگشت استفاده می‌کنند تا بر این بازگشت (عقب‌گردها) تاکید کنند. به دلیل پیشینه سازگاری یک جهش خنثی یا تغییر تصادفی در محیط ممکن است به جمعیت‌ها آسیب رسانیده و به ناسازگاری بیانجامد.

دو عامل تعیین کننده در رانش ژنتیکی μ (میزان جهش) و S (شدت انتخاب) است و تعادل بین آن‌ها با $\frac{\mu}{S}$ بدست می‌آید اما همواره چنین نیست و در موارد زیادی رانش ژنتیکی به شکل‌های دیگری بدست می‌آید (لرماند و همکاران، ۲۰۰۹). در مورد دوم یعنی سرشت تصادفی در جهش‌های خنثی باید توجه کرد که گزینش محدودیت‌های فنوتیپیکی قوی اعمال می‌شوند، مثلاً می‌توان پروتئین‌هایی که در مقیاس زمانی طولانی در ارگانیسم‌های مجزا حفظ می‌شود را مورد ملاحظه قرار داد. عدم تغییر در این مقیاس زمانی با یک انتخاب گزینش قطع می‌شود. این مربوط به نظریه‌ای است که توسط استیو می‌گولد به نام تعادل منقطع مشهور شده است (دوران، ۲۰۲۴). در یک حالت جهش‌های خنثی ممکن است نقشی مهم در تکامل بازی کنند. طول مقیاس زمانی در جهش‌های خنثی بی اهمیت است، در این حالت سرشت یک جهش را می‌توان کاملاً تابع فرآیندهای تصادفی دانست. به عبارت دیگر این جهش‌ها تنها در برخی محیط‌ها و پیشینه ژنتیکی حتمی می‌باشند و بنابراین نوعی تعادل منقطع ایجاد می‌کنند که تغییر در طول زمانی طولانی وجود ندارد اما در یک انتخاب گزینش این تعادل منقطع می‌شود. این مطالعات با واژه تعادل منقطع نام گذاری شده است که با الهام از توماس کوهن (۱۴۰۳) فلسفه علم به این‌ها اطلاق می‌شود. شرح کاملی از این مورد را در گاوریلتس (۲۰۰۳) ببینید. در مورد انقلاب‌های تکاملی منظور آن است که گاهی انتخاب طبیعی یک جمعیت را به یک اوج سازگاری موضعی می‌رساند اما از پیشرفت آن به قله‌های دیگر حتی و اگر این اوج سازگاری موضعی پائین‌تر باشد جلوگیری می‌کند. دلیل آن ایجاد اغتشاش در بسامدهای آلل‌ها می‌باشد. که رانش از یک قله به قله دیگری را پدید می‌آورد و به صورت استعاری آن را انقلاب می‌نامند. این وضعیت به‌ویژه موقعی که حجم جمعیت کوچک باشد رخ می‌دهد.

۲.۱ سوگیری جهش در یک گیاه

مفهوم فرآیندهای تصادفی در نظریه داروین اهمیت بسیار دارد، با این حال در ۲۰۲۲ یک گیاه این مفهوم را زیر سوال برد. این گیاه یک مثال نقض برای مفهوم فرآیندهای تصادفی در تکامل انواع است. مونرو و همکاران (۲۰۲۲) در شماره ۶۰۵ مجله نیچر^{۲۲} در مورد

گیاهی به نام اربیدوپسیس تالیانا^{۲۳} مقاله‌ای به چاپ رسانیدند. این گیاه اینک در پژوهش‌های گیاه‌شناسی، زیست‌شناسی مولکولی و ژنتیک ارزش بالایی یافته‌است چون بر خلاف طرفداران نظریه تکامل می‌توان به صورت مداخله‌ای جهش‌های غیرتصادفی (طراحی هوشمند) پدید آورد. تا به حال تعداد بسیار زیادی از این نوع جهش‌های غیر تصادفی در آن ایجاد کرده‌اند. این گیاه دارای ژنوم کوچکی به تعداد ۱۳۵ میلیون جفت باز است. اینک درباره این گیاه مقاله‌های زیادی چاپ شده است که حداقل در یک نمونه فرآیندهای تصادفی در تکامل را زیر سؤال می‌برد. با این حال فرآیندهای تصادفی یک پارادایم بسیار قوی در نظریه تکامل است که با یک نمونه، حداقل بنابر فلسفه توماس کوهن، زوال نخواهد یافت.

بحث و نتیجه‌گیری

فرآیندهای تصادفی امروزه در مطالعه نظریه تکامل کاربردهای بسیاری یافته است. نظریه‌های جدید به خصوص بعد از سنتز گسترش یافته تکامل همه به کمک فرآیندهای تصادفی توسعه یافته‌اند با این حال ارزشیابی این که چگونه شکل‌های مختلف فرآیندهای تصادفی بر تکامل صفات‌های خاصی تأثیر می‌گذارد، دشوار است.

از بین سه شکل فرآیند تصادفی نوع فرآیندهای تصادفی در جهش و تغییر که رانش ژنتیکی را توجیه می‌کند بیشترین کاربرد را دارد. این نظریه‌ها امروزه به این توافق رسیده‌اند که رانش ژنتیکی تطابق و سازگاری در انواع را محدود می‌کند. از زمان فیشر و رایت که دو آماردان برجسته انگلیسی و آمریکایی و هر دو مدافع نظریه داروین بودند و همه تلاش خود را به کار بردند تا آن را بر مبنای فرآیندهای تصادفی بیان کنند، تا کنون گرایش به سمت نظریه لامارک (که داروین با این شکل از بیان تکامل کاملاً مخالف بود) همه از فرآیندهای تصادفی برای توجیه نظریه تکامل استفاده می‌کنند. به عنوان مثال اپی ژنتیک که گرایش امروزه در تکامل است (نزدیک‌تر به نظریه لامارک است) در مقاله **بیور و همکاران (۲۰۲۰)** به نقش فرآیندهای تصادفی در تفاوت‌های اپی ژنتیکی پرداخته است.

مراجع

- بولر، ج. پ. (۱۹۸۹)، انقلاب مندل، ترجمه محی‌الدین زعفرانی (۱۳۷۲)، شرکت انتشارات علمی و فرهنگی، تهران.
- کوهن، ت. (۱۴۰۳)، ساختار انقلاب‌های علمی، ترجمه سعید زیبا کلام، انتشارات سمت، تهران.
- Biwer, C., Kawam, B., Chapelle, V., and Silvestre, F. (2020), The role of stochasticity in the origin of epigenetic variation in animal populations, *Integrative and Comparative Biology*, **60**(6), 1544-1557.
- Corning, P. A. , Kaufman, S. A. (2023), *Evolution on Purpose*, The Mit Press, Cambridge, MA 01239.
- Duran, N. (2024), The many ways toward punctuated evolution, *Palaeontology*, **67**(5), 12731-12736.
- Eyre-Walker, A., and Keightley, P. D. (2007), The distribution of fitness effects of new mutations, *Nature Reviews Genetics*, **8**(8), 612-618.
- Fisher, R. A. (1930), *The genetical theory of natural selection*, The Clarendon Press, Oxford, England.

- Galton, F. (1889), *Natural inheritance*, Macmillan, London.
- García-Pintos, L. P. (2024). Limits on the evolutionary rates of biological traits. *Scientific Reports*, **14**(1), 11314.
- Gavrilets S. (2003), Models of Speciation: What Have We Learned in 40 Years?, *Evolution; International Journal of Organic Evolution*, **57**, 2197–2215.
- Johnson, T., and Barton, N. H. (2002), The effect of deleterious alleles on adaptation in asexual populations. *Genetics*, **162**(1), 395-411.
- Keseroglu, K., Zinani, O. Q., Keskin, S., Seawall, H., Alpay, E. E., and Özbudak, E. M. (2023), Stochastic gene expression and environmental stressors trigger variable somite segmentation phenotypes, *Nature communications*, **14**(1) , 6497-6501.
- Lea-Smith, D. J., Hassard, F., Coulon, F., Partridge, N., Horsfall, L., Parker, K. D. and Krasnogor, N. (2025), Engineering biology applications for environmental solutions: Potential and challenges, *Nature Communications*, **16**(1), 3538-3547.
- Lehmann, L., and Balloux, F. (2007), Natural selection on fecundity variance in subdivided populations: kin selection meets bet hedging., *Genetics*, **176**(1), 361-377.
- Lenormand, T. ,Roze, D. ,Rousset. (2009) Stochasticity in Evolution, *Trends in Ecology and Evolution*, **24**(3) ,157-165.
- Monroe, J.G. ,Srikant, T., Weigle, D. (2022), Mutation Bias Reflects Natural Selection in Arabidopsis Thaliana, *Nature*, **602** ,101-105.
- Mode, C. J., and Sleeman, C. K. (2012), *Stochastic processes in genetics and evolution: computer experiments in the quantification of mutation and selection*, World Scientific, Singapore (SG).
- Moran P. (1958), Random Processes in Genetics, *Proc. Camb. Philos. Soc.*, **54**, 60–71.
- Orr, H. A. (2005), The probability of parallel evolution, *Evolution* , **59**(1), 216-220.
- Sarkar, S. (1992), The founders of Evolutionary Genetics, *Boston studies in the Philosophy and History of Science* , **142**, 89-116.
- Sood, V., and Misteli, T. (2022), The stochastic nature of genome organization and function. *Current opinion in genetics and development* , **72**, 45-52.

Tejada-Lapuerta, A., Bertin, P., Bauer, S., Aliee, H., Bengio, Y. and Theis, F. J. (2025), Causal machine learning for single-cell genomics, *Nature Genetics*, **20**, 1–12.

Wright, S. (1931), Evolution in Mendelian populations, *Genetics*, **16(2)**, 97-159.

Yampolsky, L. Y., and Stoltzfus, A. (2001), Bias in the introduction of variation as an orienting factor in evolution, *Evolution and development*, **3(2)**, 73-83.

How did the theory of stochastic processes contribute to the model developed in the modern synthesis of evolution?

Dr. Hamidreza Oreizi¹

¹Department of Psychology, University of Isfahan

Abstract: Perhaps the most important theory that changed all human and social concepts was the emergence of Darwinian evolution and Mendelian genetics, which were presented within six years of each other in the 19th century and changed human thinking about humans, their environment, all animals and even plants. In 1942, Julian Huxley linked these two theories in a combined model that he called the modern synthesis. This combined model helped this theory to survive. In this article, we explain mutation, which is viewed differently in the views of Darwin and Lamarck (English and French scientists), using random processes, and we explain the types of mutations and their classification using random processes. This classification has helped the modern synthesis theory to continue to maintain the link between evolution and genetics.

Keywords: Stochastic processes – Darwinian evolution – Mendelian genetics – Mutation – Classification of mutations.

Mathematics Subject Classification (2020): 60J70.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



B۰۶۶A"۰۶۶"۹F۰۶"۸F۰۶"۷F۰۶"۶F۰۶"۵F۰۶"۴F۰۶"۳F۰۶"۲F۰۶"۱F۰۶"۰F۰۶"

برآورد معیار $CoVaR$ بر اساس رویکرد $Copula - GARCH$ و توزیع GED

فاطمه علیزاده^۱، محمد امینی^۲، غلامرضا محتشمی برزادران^۳،

گروه آمار، دانشگاه فردوسی مشهد

چکیده: ریسک یکی از مفاهیم پایه ای در بازارهای مالی می باشد و اندازه گیری و تحلیل آن بسیار اهمیت دارد. در این مقاله یکی از مهم ترین روش های اندازه گیری ریسک یعنی ارزش در معرض خطر شرطی را در حالت دو متغیره با استفاده از دو سری زمانی بازدهی مطالعه می کنیم در ادامه ساختار وابستگی این بازدهی ها را بر اساس تابع مفصل مورد بررسی قرار داده و در سری زمانی برای پیش بینی واریانس از مدل های گارچ با توزیع حاشیه ای GED استفاده می کنیم و در پایان روش مورد نظر در داده واقعی به کار گرفته می شود. واژه های کلیدی: ریسک سیستمی، ارزش در معرض خطر شرطی، مدل گارچ، تابع مفصل، توزیع GED .
کد موضوع بندی ریاضی (۲۰۲۰): 91B84, 62H05, 62M10, 62P05.

۱ مقدمه

با گذشت زمان به دلیل گسترش حوادث نامطلوب مختلف در جهان که بخشی از آن از افزایش فعالیت های اقتصادی، اجتماعی، سیاسی و ... نشئت می گیرد، نااطمینانی درباره آینده بیشتر شده است. ریسک یکی از مفاهیم پایه ای در بازار های مالی است که تعریف واحدی از آن وجود ندارد. اقتصاددانان، دانشمندان علوم رفتاری، نظریه پردازان ریسک و آمارشناسان هر یک تصور خاصی از ریسک دارند، البته به طور سنتی ریسک به عنوان عدم اطمینان تعریف شده است و بر اساس این مفهوم، در این جا ریسک به عنوان عدم اطمینان مرتبط با وقوع یک خسارت و یا زیان تعریف می شود.

در سال های اخیر با بررسی علل نوسانات مالی مشخص می شود که عوامل مختلفی در ایجاد نوسانات در بازارهای مالی (بیمه ها، بانک ها و ...) مؤثر بوده اند از جمله بحران نفتی سال ۱۹۷۳، بحران های مالی سال های ۱۹۹۷-۱۹۹۸ در جنوب شرق آسیا و رویدادهای طبیعی نظیر زلزله و سونامی. بروز چنین حوادثی و به تبع آن نوسانات بازارهای مالی موجب ایجاد اختلال در فعالیت های سازمان های مالی، تجاری و حتی تولیدی می شود. از این رو اندازه گیری و ارزیابی ریسک اهمیت پیدا می کند و از جمله مهم ترین کاربردهای آن می توان

^۱ فاطمه علیزاده، fateme.alizade301073@gmail.com

به قیمت‌گذاری سهام‌ها، تعیین حق بیمه متناسب با ریسک و تعیین سرمایه اقتصادی^۱ شرکت که پشتوانه‌ی فعالیت‌های مالی است و خسارت‌های ناگهانی بالاتر از حد انتظار را پوشش می‌دهد، اشاره کرد.

تاکنون روش‌های متعددی برای اندازه‌گیری ریسک معرفی شده است که هر یک از این روش‌ها مزایا و معایب خاص خود را دارد. یکی از مهم‌ترین معیارهای اندازه‌گیری ریسک ارزش در معرض خطر^۲ است. این شیوه اندازه‌گیری را ابتدا تیل گلدیمن^۳ در سال ۱۹۸۰ ارائه کرد و در اواخر دهه ۱۹۸۰ توسط موسسه جی-پی مورگان^۴ گسترش یافت. این شاخص حداکثر زیان ممکن را در یک افق زمانی مشخص با توجه به یک سطح احتمال معین α بیان می‌کند و با نماد Var_{α} نشان داده می‌شود، به عنوان مثال ارزش در معرض خطر با سطح احتمال ۵٪ برای یک بازه زمانی یک روزه گویای اینست که حداکثر زیان احتمالی طی روز بعدی با احتمال ۵ درصد از مقدار ارزش در معرض خطر بالاتر می‌رود. این معیار را می‌توان با استفاده از سری بازدهی محاسبه کرد. فرض کنید $\{p_t\}$ سری زمانی بهای نوعی دارایی مالی باشد، بازدهی در زمان t ام به صورت $X_t = \log(p_t) - \log(p_{t-1})$ تعریف می‌شود.

تعریف ۱.۱. اگر $\{X_t\}$ یک سری بازدهی با تابع توزیع F_{X_t} باشد، آنگاه

$$Var_{\alpha}(X_t) = F_{X_t}^{-1}(\alpha) = \inf\{x; F_{X_t}(x) \geq \alpha\}; \alpha \in (0, 1).$$

رویدادهایی که در نهادهای مالی رخ می‌دهد ممکن است بر سیستم مالی سایر نهادها و یا حتی بر کل سیستم اقتصاد اثر بگذارد در این صورت نمی‌توان تاثیر متقابل بین ریسک‌ها را نادیده گرفت به همین دلیل از زمان شروع بحران‌های مالی و اقتصادی در سال ۲۰۰۷ ریسک سیستمی^۵ یعنی ریسکی که بحران در یک نهاد، بر سایر نهادها یا کل سیستم مالی تاثیر می‌گذارد مطرح شد. با مروری بر بحران‌های مالی در جهان ملاحظه می‌شود مهمترین خطری که ثبات مالی را تهدید می‌کند ریسک سیستمی ناشی از بحران‌های سیستم بانکی است. این بحران‌های بانکی اخیر در جهان، هزینه‌های زیادی برای آن کشورها به همراه داشته و نگرانی‌هایی را برای نظام مالی آنها به وجود آورده است. لذا به مطالعه ریسک سیستمی در بانک‌ها توجه زیادی شده و با معیارهای مختلفی ریسک سیستمی در بانک‌ها و مؤسسات مالی جهان ارزیابی گردیده است. یکی از روش‌های اندازه‌گیری ریسک سیستمی، ارزش در معرض خطر شرطی^۶ است که اولین بار توسط آدرین و برونرمریر (۲۰۱۱)^۷ معرفی شد و پس از آن جیراردی و ارگان (۲۰۱۳)^۸ تعریف دیگری از آن ارائه دادند. این معیار ارزش در معرض خطر یک سازمان با بازدهی Y را در صورتی که سازمانی دیگر با بازدهی X در بحران قرار داشته باشد بدین معنی که بازدهی آن حداکثر برابر ارزش در معرض خطر در سطح احتمال α باشد، در سطح احتمال β را محاسبه می‌کند، یعنی

$$CoVaR_{\alpha, \beta}(Y|X) = Var_{\beta}(Y | X \leq Var_{\alpha}(X)); \alpha, \beta \in (0, 1), \quad (1.1)$$

که به آن ارزش در معرض خطر شرطی ($CoVaR$) گویند. در صورتی که بازدهی‌های Y وابسته باشند می‌توان از توابع مفصل برای بررسی ساختار وابستگی بین بازدهی‌ها استفاده کرد. از جمله مطالعاتی که در زمینه برآورد این معیارها با استفاده از تابع مفصل انجام شده است می‌توان به موارد زیر اشاره کرد.

¹Economic Capital

²Value at Risk

³Till Guldman

⁴J.P. Morgan

⁵Systemic Risk

⁶Conditional Value-at-risk

⁷Adrin and Brunnermeier

⁸Girardi and Ergun

ماینیک و اسکینینگ (۲۰۱۴) معیار $CoVaR$ را براساس توابع مفصل گاوسی و استودنت، ربرودو و یوگولینی (۲۰۱۵) در مفصل‌های تغییر شکل یافته گامبل و برناردی و همکاران (۲۰۱۷) نیز برای توابع مفصل کران‌های فرشه، مارشال اولکین، FGM و برخی مفصل‌های ارشمیدسی یعنی گامبل، فرانک، جو و کلیتون تحلیل کرده‌اند. کاریمالیس و نومیکس (۲۰۱۸)، $CoVaR$ را با توجه به نتایج برناردی و همکاران (۲۰۱۷) و مدل‌های سری زمانی گارچ بررسی و سپس ریسک سیستمی در سیستم بانکی اروپا را با توابع مفصل فرانک و جو محاسبه کرده‌اند. جی و همکاران (۲۰۱۹) $CoVaR$ را با استفاده از مفصل‌های زمان متغیر برای قیمت نفت و دلار، هوانگ و همکاران (۲۰۲۴) به روش شبیه‌سازی مونت کارلو و مفصل، کیلمن و همکاران (۲۰۲۲) و وانگ و همکاران (۲۰۲۵) با مفصل‌های واین مطالعه کرده‌اند.

۱.۱ تابع مفصل

در سال‌های اخیر استفاده از توابع مفصل^۹ به عنوان ابزاری مفید برای مدل‌سازی ساختار وابستگی بین دو یا چند متغیر گسترش یافته است و کاربرد بسیاری در علوم بیمه و اقتصاد دارد. مفصل‌ها توابعی هستند که بین توابع توزیع چند متغیره و توزیع حاشیه ای آن‌ها ارتباط برقرار می‌کنند. از طرفی مفصل خود تابع توزیعی است که توزیع حاشیه‌ای آن از توزیع یکنواخت استاندارد پیروی می‌کند. قضیه ۲.۱ که با نام قضیه اسکالر معروف است رابطه بین توزیع‌های حاشیه ای، تابع مفصل و توزیع توام دو متغیر تصادفی وابسته را بیان می‌کند.

قضیه ۲.۱. اگر (X, Y) یک بردار تصادفی با تابع توزیع توام $F(\cdot, \cdot)$ و توابع توزیع حاشیه ای F_X و F_Y باشند آنگاه مفصل $C: I^2 \rightarrow I$ وجود دارد به طوریکه:

$$F(x, y) = C(F_X(x), F_Y(y)) ; \forall x, y \in R.$$

اگر F_X و F_Y پیوسته باشد آنگاه C یکتاست در غیر این صورت به صورت یکتا روی $Range F_X \times Range F_Y$ تعیین می‌شود و با داشتن توزیع‌های حاشیه‌ای و توام برای یافتن مفصل متناظر قرار می‌دهیم: $u = F_X(x)$ ، $v = F_Y(y)$ و داریم:

$$C(u, v) = F(F_X^{-1}(u), F_Y^{-1}(v)) ; \forall u, v \in [0, 1].$$

خانواده‌ای از مفصل‌ها به نام مفصل‌های ارشمیدسی، وجود دارند که با استفاده از یک تابع به نام تابع مولد $\phi: [0, 1] \rightarrow [0, \infty)$ که در آن $\phi(1) = 0$ و $\phi(0) = \infty$ می‌باشد به شکل زیر ساخته می‌شوند:

$$C(u, v) = \phi^{-1}(\phi(u) + \phi(v)) \quad (2.1)$$

۲.۱ مدل بندی تغییر پذیری

مدل بندی تغییر پذیری که اغلب به صورت واریانس تعریف می‌شود برای پیش بینی نوسانات آینده یکی از مباحث مهم در مطالعات مالی می‌باشد. اغلب مشاهده می‌شود که بازدهی روزانه سهام بعد از یک دوره تغییرات شدید بها، واریانس شرطی بزرگتری نسبت به یک دوره نسبتاً پایدار دارند. مدل‌های $ARCH$ و $GARCH$ ^{۱۰} برای فرایند واریانس شرطی که با آن بتوانیم تغییر پذیری مقادیر آینده را بر

^۹Copula Function

^{۱۰}Generalized Auto regressive Conditional Heteroscedasticity

اساس داده های جاری و گذشته پیشگویی کنیم، در مدل ARCH(q) فرض می شود که سری بازدهی $\{X_t\}$ به صورت

$$X_t = \mu + a_t, \quad a_t = \sigma_{t|t-1}^{1/2} \varepsilon_t, \quad (3.1)$$

$$\sigma_{t|t-1}^{1/2} = \omega + \alpha_1 a_{t-1}^{1/2} + \alpha_2 a_{t-2}^{1/2} + \dots + \alpha_q a_{t-q}^{1/2}. \quad (4.1)$$

تولید می شود که در آن پارامترهای α_i و ω مجهولند، $\{\varepsilon_t\}$ دنباله ای از متغیرهای تصادفی مستقل و هم توزیع با میانگین صفر و واریانس واحد است و ε_t برای $j = 1, 2, \dots$ مستقل از a_{t-j} است.

در رویکردی دیگر علاوه بر q تاخیر از توان دوم a_t ، p تاخیر از واریانس شرطی را هم وارد مدل می کند. یعنی،

$$\sigma_{t|t-1}^{1/2} = \omega + \alpha_1 a_{t-1}^{1/2} + \dots + \alpha_q a_{t-q}^{1/2} + \beta_1 \sigma_{t-1|t-2}^{1/2} + \dots + \beta_p \sigma_{t-p|t-p-1}^{1/2}. \quad (5.1)$$

۲ ارزش در معرض خطر شرطی

ماینینگ و اسکینگ (۲۰۱۴) برای دو متغیر بازدهی X و Y با توابع توزیع حاشیه ای $F_X(x)$ و $F_Y(y)$ ساختار وابستگی تابع مفصل $C(u, v)$ و توزیع توأم $F(x, y)$ ارزش در معرض خطر شرطی به ازای هر $\alpha, \beta \in [0, 1]$ به شرح زیر به دست آوردند:

$$CoVaR_{1-\alpha, 1-\beta}(Y_t | X_t) = F_{Y_t}^{-1}(F_{V|U \leq \alpha}^{-1}(\beta)) = VaR_{F_{V|U \leq \alpha}^{-1}(\beta)}(Y_t) \quad (1.2)$$

که در آن (U, V) دارای توزیع توأم $C(u, v)$ با توابع توزیع حاشیه ای یکنواخت استاندارد می باشند و

$$F_{V|U \leq \alpha}(v) = \frac{C(\alpha, v)}{\alpha}. \quad (2.2)$$

علاوه بر این برناردی (۲۰۱۷) نشان داد با در نظر گرفتن $w^* = F_{V|U \leq \alpha}^{-1}(\beta)$ ، و استفاده از (۲.۲)، w^* از طریق حل معادله زیر به دست می آید:

$$C(\alpha, w^*) = \alpha\beta. \quad (3.2)$$

سپس w^* را در مفصل های کران های فرچت، مارشال اولکین، FGM و مفصل های ارشمیدسی فرانک، کلیتون و گامبل محاسبه کرده است.

حال در حالت خاص اگر $C(u, v)$ یک مفصل ارشمیدسی با تابع مولد $\varphi(t)$ باشد با محاسبات جبری ساده داریم:

$$w^* = \varphi^{-1}(\varphi(\alpha\beta) - \varphi(\alpha)), \quad (4.2)$$

در نتیجه

$$CoVaR_{1-\alpha, 1-\beta}(Y_t | X_t) = VaR_{w^*}(Y_t). \quad (5.2)$$

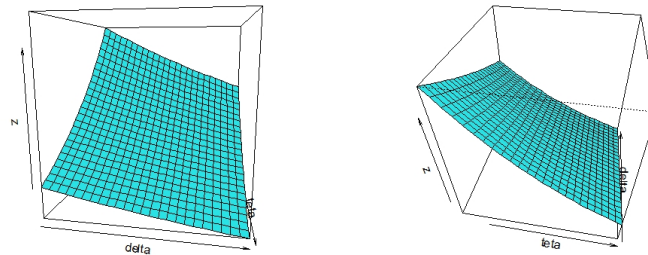
:

حال اگر X و Y دارای ساختار وابستگی مفصل BB_1 با تابع مولد $\varphi(t) = (t^{-\theta} - 1)^\delta$ باشد، آنگاه معکوس تابع مولد به صورت $\varphi^{-1}(s) = (s^{\frac{1}{\delta}} + 1)^{-\frac{1}{\theta}}$ و با استفاده از رابطه (۴.۲) داریم:

$$w^* = [((\alpha\beta)^{-\theta} - 1)^\delta - (\alpha^{-\theta} - 1)^\delta]^{\frac{1}{\delta}} + 1]^{-\frac{1}{\theta}}$$

(۶.۲)

در شکل ۱ تغییرات w^* براساس تغییرات θ و δ برای $\alpha = \beta = 0.5$ رسم شده است که در آن مشاهده می‌شود w^* با ثابت نگاه داشتن یکی از پارامترهای θ و δ نسبت به پارامتر دیگر نزولی می‌باشد و از آنجایی که معکوس تابع توزیع تابعی صعودی نسبت به سطح احتمال (w^*) است از این رو با افزایش هر یک از مقادیر θ و δ مقدار $CoVaR$ کاهش پیدا می‌کند. اکنون فرض کنید سری زمانی



شکل ۱: نمودار w^* براساس تغییرات θ و δ مفصل $BB\backslash$

بازدهی به صورت $Y_t = \mu_t + \sigma_{t|t-1}\varepsilon_t$ باشد. ابتدا مناسب‌ترین مدل سری زمانی بازدهی Y_t را به دست آورده، پارامترهای مدل سری زمانی، پارامتر توزیع خطای تعمیم یافته و میانگین (μ_t) و انحراف معیار ($\sigma_{t|t-1}$) زمان آتی را برآورد کرده و برای برآورد ارزش در معرض خطر شرطی در سطح احتمال w^* که به روش زیر بدست می‌آید عمل می‌کنیم:

$$P(Y_t < CoVaR) = P(\mu + \sigma_{t|t-1}\varepsilon_t < CoVaR) = w^* \Rightarrow$$

$$P(\varepsilon_t < \frac{CoVaR - \mu}{\sigma_{t|t-1}}) = w^* \Rightarrow$$

$$CoVaR = \mu + \sigma_{t|t-1}F_{\varepsilon_t}^{-1}(w^*).$$

که در آن وقتی خطاها از توزیع GED تبعیت کنند تابع چگالی آن به صورت $f_{\varepsilon_t}(x) = \frac{\beta}{2\Gamma(\frac{1}{\beta})}e^{-|x|^\beta}$ است. برای $w^* < \frac{1}{2}$ با تغییر متغیرهای $y = (-x)^\beta$ و $z = x^\beta$ داریم:

$$F_{\varepsilon_t}(\eta) = \frac{\beta}{2\Gamma(\frac{1}{\beta})} \left(\int_0^\infty \frac{1}{\beta} y^{\frac{1}{\beta}-1} e^{-y} dy + \int_0^{\eta^\beta} \frac{1}{\beta} z^{\frac{1}{\beta}-1} e^{-z} dz \right),$$

و با تعاریف تابع گاما و گامای ناقص پایینی می‌توان نوشت:

$$F_{\varepsilon_t}(\eta) = \frac{1}{2\Gamma(\frac{1}{\beta})} (2\Gamma(\frac{1}{\beta}) - \Gamma(\frac{1}{\beta}, \eta^\beta)). \quad (7.2)$$

و $CoVaR$ از طریق حل معادله زیر بدست می‌آید.

$$1 - \frac{\Gamma(\frac{1}{\beta}, CoVaR^\beta)}{2\Gamma(\frac{1}{\beta})} = w^*. \quad (8.2)$$

و برای $w^* > \frac{1}{2}$ با استفاده از رابطه

$$F_{\varepsilon_t}(\eta) = \frac{\Gamma(\frac{1}{\beta}, (-\eta)^\beta)}{2\Gamma(\frac{1}{\beta})}. \quad (9.2)$$

از طریق حل معادله

$$\frac{\Gamma(\frac{1}{\beta}, (-CoVaR)^\beta)}{2\Gamma(\frac{1}{\beta})} = w^*, \quad (10.2)$$

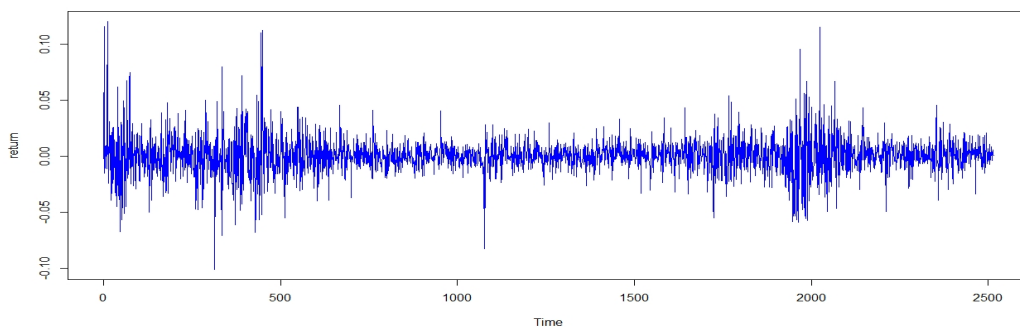
برای حل این معادلات نیازمند روش‌های عددی می‌باشیم.

۳ کاربرد در داده واقعی

در این بخش داده‌های بازدهی روزانه شرکت‌های *SP* و *IBM* از سال ۲۰۰۲ تا ۲۰۱۰ مورد بررسی قرار گرفته است.

ابتدا برای محاسبه مقدار ساختار وابستگی بین بازدهی‌ها را مورد بررسی قرار می‌دهیم. با مقایسه تاو کندال، وابستگی دمی بالا، پایین و معیارهای *AIC* و *BIC* توابع مفصل مختلف برازش داده شده مفصل *BB1* مفصل مناسبی برای بیان ساختار وابستگی می‌باشد، پی-مقدار آزمون نیکویی برازش کرامر ون مایسز و کولموگروف اسمیرنوف بیشتر از پنج صدم بنابراین فرض مناسب بودن مفصل *BB1* را رد نمی‌کنیم.

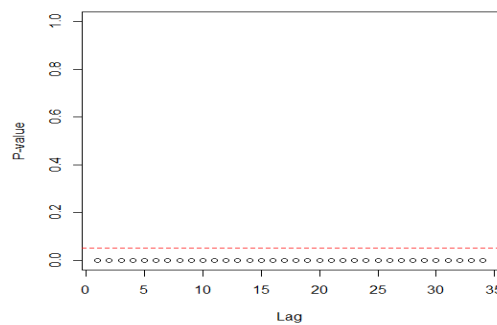
برآورد پارامترهای این مفصل به روش ماکزیم درستنمایی $\theta = 0.76$ و $\delta = 0.32$ بدست آمد بنابراین مقدار w^* برای $\alpha = \beta$ ۰/۵ بر اساس رابطه ۶.۲ تقریباً ۰/۷۲ محاسبه شد.



شکل ۲: بازدهی روزانه سهام IBM

شکل ۲ از سری بازدهی IBM نشان می‌دهد که در برخی از دوره‌ها بسیار ناپایدار شدند به عبارت دیگر واریانس شرطی سری زمانی در طی زمان تغییر می‌کند، به همین منظور در ادامه به بررسی مدل‌های آرچ و گارچ می‌پردازیم. در مورد داده‌های بازدهی سهام IBM برای بررسی وجود اثرهای ARCH با رسم پی-مقدارهای آزمون مک‌لیود-لی در شکل ۳ مشاهده می‌شود که همگی در سطح ۵ درصد معنادارند و به طور رسمی شواهدی قوی به نفع ARCH در داده‌ها را نشان می‌دهد.

با بررسی سری توان دوم و قدرمطلق داده‌های زیان سهام IBM، نمودارهای خودهمبستگی، (acf) خودهمبستگی جزئی، $(pacf)$ و خصوصاً خودهمبستگی گسترده $(eacf)$ یک مدل $ARMA(1, 1)$ را گزارش می‌کند، بنابراین می‌توان یک مدل $GARCH(1, 1)$ را برای داده‌های اصلی در نظر گرفت.



شکل ۳: پی-مقدارهای آزمون مکلیود-لی برای زیان روزانه سهام IBM

اگر چه نشانه ای از یک مدل $GARCH(2, 2)$ نیز وجود دارد، بنابراین با استفاده از روش‌های مشخص سازی مثل AIC و BIC $GARCH(1, 1)$ را انتخاب می‌کنیم. پیش از آن که مدل برازنده شده‌ای را بپذیریم کار اساسی این است که واریسی کنیم آیا مدل به درستی مشخص شده یا نه. اگر مدل به درستی مشخص شده باشد مانده های استاندارد شده یک سری نوفه سفیدند. در این مورد در نمودارهای ACF و $PACF$ مانده های استاندارد شده و توان دومشان هیچ کدام از مقادیر خارج از بازه نیستند که این پیشنهاد می‌کند سری باقی مانده ها نوفه سفید است، علاوه بر این نتایج آزمون های کلی نگر هم حاکی از مدل نوفه سفید برای باقیمانده های استاندارد شده است. برای باقی مانده ها توزیع های خطای تعمیم یافته، نرمال و استودنت مورد بررسی قرار گرفت که توزیع خطای تعمیم یافته با مقایسه معیارهای AIC و MSE عملکرد بهتری داشت.

پس از برآورد ضرایب مدل سری زمانی مقدار $CoVaR$ شرکت IBM در صورت بحران در شرکت SP و محاسبه مقدار w^* برای $\alpha = \beta = 0.05$ که بر اساس رابطه ۶.۲ تقریباً ۰.۷۲٪ محاسبه شد بر اساس مدل سری زمانی $GARCH(1, 1)$ و توزیع حاشیه‌ای GED با پارامتر $3/174$ مقدار $CoVaR$ برابر با $2/641$ بدست آمد.

۴ نتیجه‌گیری

در این مقاله یکی از روش‌های مهم بررسی ریسک سیستمی یعنی ارزش در معرض خطر شرطی بر اساس تابع مفصل ارشمیدی $BB1$ و مدل‌های سری زمانی گارچ در صورتی که مانده‌های آن از توزیع GED پیروی کند، مورد مطالعه قرار گرفت و در پایان این معیار برای شرکت IMB در صورتی که شرکت SP در شرایط بحرانی قرار گیرد محاسبه شد که در آن توابع مفصل مختلف برای بررسی ساختار وابستگی بازدهی مدل‌های سری زمانی گارچ برای پیش بینی نوسان روز آینده شرکت IMB بررسی شد، نتایج حاکی از آن بود که مفصل $BB1$ و سری زمانی $GARCH(1, 1)$ بهترین مدل‌بندی را ارائه می‌دهند.

مراجع

- [1] Adrian, T., Brunnermeier, M.K., (2016). CoVaR, *Amer. Econ. Rev.*, **106**, 1705–1741.
- [2] Bernardi, M., Durante, F., Jaworski, P., (2017). CoVaR of families of copulas. *Statist. Probab. Lett.* **120**, 8–17.

- [3] Girardi, G., Ergün, A.T., (2013). Systemic risk measurement: multivariate GARCH estimation of Co-VaR. *J. Bank. Finance*, **37**, 3169–3180.
- [4] Huang, W., Lin, N., Hong, L. J. (2024). Monte Carlo Estimation of CoVaR, *Operations Research*, **72(6)**, 2337-2357.
- [5] Ji, Q., Bouri, E., Roubaud, D., Shahzad, S. J. H. (2018). Risk spillover between energy and agricultural commodity markets: A dependence-switching CoVaR-copula model, *Energy Economics*, **75**, 14-27.
- [6] Karimalis, E. N., Nomikos, N. K., (2018). Measuring systemic risk in the European banking sector: a copula CoVaR approach, *The European Journal of Finance*, **24(11)**, 944-975.
- MReboredo, Juan C., Ugolini, A., (2015). Systemic risk in European sovereign debt markets: A CoVaR-copula approach, *Journal of International Money and Finance*, **51**, 214–244.
- [7] Nelsen, R.B. (2006). *An Introduction to Copulas*, vol. 139. Springer Science & Business Media.
- [8] Wang, J., Yan, X., Cao, Y., Wang, X. (2025), Multi-scale dependence and risk contagion among international financial markets based on VMD-Vine copula-CoVaR, *Applied Economics*, **57(6)**, 658-677.

Estimating CoVaR Based on Copula-GARCH method and GED Distribution

Fatemeh Alizadeh¹, Mohammad Amini³, Gholam Reza Mohtashami Borzadaran²

Department of Statistics, Ferdowsi University of Mashhad

Abstract: Risk is one of the basic concepts in financial markets and its measurement and analysis is very important. In this article, we study one of the most important methods of measuring risk, namely conditional value at risk, in the bivariate case using two time series of returns. Next, we examine the dependence structure of these returns based on the detailed function and use GARCH models with marginal distributions (GED) in the time series to predict the variance. Finally, the method is applied to real data.

Keywords: Systemic Risk, Conditional Value at Risk, GARCH Model, Copula Function, GED Distribution.

Mathematics Subject Classification (2020): 62P05, 62M10, 62H05, 91B84.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



مروری بر استفاده از خانواده توزیع‌های فاز-نوع در مدل شوک‌های تعمیم یافته

سیروس فتحی منش^۱، محی‌الدین ایزدی^۲

^۱ گروه آمار- دانشگاه کردستان

^۲ گروه آمار- دانشگاه رازی

چکیده: استفاده از خانواده توزیع‌های فاز-نوع که بر اساس فرایندهای مارکف زمان گسسته و پیوسته تعریف می‌شوند در مدل‌بندی متغیرهای تصادفی نامنفی گسسته و پیوسته رو به گسترش است. در دهه اخیر، استفاده از این خانواده توزیع‌ها در برخی مدل‌های شوک مورد توجه پژوهشگران زمینه قابلیت اعتماد قرار گرفته است. در این مقاله هدف ما مروری بر نحوه استفاده از خانواده توزیع‌های فاز-نوع در سه مدل شوک کرانگین تعمیم‌یافته، تجمعی تعمیم‌یافته و دلتای تعمیم یافته است. استفاده از این خانواده از توزیع‌ها در مدل شوک‌های ذکرشده توجیه می‌شود و سپس توزیع تعداد شوک‌ها و طول عمر سیستم بر اساس روابط ماتریسی موجود در خانواده توزیع‌های فاز-نوعی بیان می‌شوند.

کلمات کلیدی: خانواده توزیع‌های فاز-نوع، فرایندهای مارکف، مدل شوک تجمعی تعمیم یافته، مدل شوک دلتای تعمیم یافته، مدل شوک کرانگین تعمیم یافته.

۱ مقدمه

در مدل‌های آماری استفاده از خانواده توزیع‌های منعطف که با تغییر پارامتر، دامنه وسیعی از توزیع‌ها را شامل شود بسیار حائز اهمیت است، زیرا اشتباه در انتخاب توزیع برای برازش به یک مجموعه داده متناسب با وسعت دامنه آن خانواده، کاهش می‌یابد. خانواده توزیع‌های فاز-نوعی^۱ یک خانواده گسترده از توزیع‌های آماری است که برای هر دو حالت گسسته و پیوسته قابل تعریف می‌باشند. این خانواده به‌طور گسترده‌ای در مدل‌بندی متغیرهای تصادفی نامنفی در زمینه‌های مختلفی مانند قابلیت اعتماد، نظریه صف، ریسک‌های بیمه‌ای و مالی، زیست‌شناسی و سایر علوم مورد استفاده قرار می‌گیرد. ایده اولیه این توزیع در **ارلنگ** (۱۹۰۹) مطرح گردید و سپس

^۱ سخنران، s.fathimanesh@uok.ac.ir

^۱ phase-type distribution

نیوتس^۲ در سال ۱۹۷۵ آن را در یک چارچوب ریاضی مدون بسط داد و به عنوان یک خانواده از توزیع‌ها معرفی نمود (نیوتس ، ۱۹۷۵).

خاصیت بسته بودن این خانواده از توزیع‌ها نسبت به مجموع متناهی، توزیع‌های آمیخته و سیستم‌های منسجم قابلیت اعتمادی در اسف و لویکسون (۱۹۸۲) مورد بررسی قرار گرفت. در دهه ۹۰ میلادی و بعد از آن، پژوهشگران مختلفی به تطبیق این توزیع‌ها بر پدیده‌های تصادفی مختلف و همچنین به کنکاش بیشتری از خواص این خانواده پرداختند. در سینی (۱۹۹۰) ویژگی‌ها و نحوه مشخص‌سازی این توزیع‌ها مورد مطالعه قرار گرفت. همچنین، نحوه محاسبه برآورد درست‌نمایی ماکسیمم پارامترهای این خانواده از توزیع‌ها در بابیو و کومنی (۱۹۹۲) شرح داده شد. در آلن (۱۹۹۵) و فدی (۱۹۹۵) از توزیع‌های فاز-نوع به ترتیب برای مدل‌بندی داده‌های تحلیل بقا و قابلیت اعتماد استفاده شد. در لی (۲۰۰۳) و کای و لی (۲۰۰۵) به کاربرد حالت چند متغیره این توزیع‌ها در مدل‌های شوک و نظریه ریسک پرداخته شد. همچنین، بلاد (۲۰۰۵) یک مقاله مروری در زمینه مدل‌بندی ریسک با استفاده از توزیع‌های فاز-نوع نوشت. از جمله مقالات پیشرو در زمینه کاربرد این توزیع‌ها در مدل‌های شوک می‌توان به منتورو و همکاران (۲۰۰۹) و سگویا و لابیو (۲۰۱۳) اشاره نمود. در چند سال اخیر نیز موجی جدید از توجه به توزیع‌های فاز-نوع در انواع مدل شوک ایجاد شده است. از جمله این مقالات می‌توان به ایرلماز (۲۰۱۷)، ژائو و همکاران (۲۰۱۸)، ازکوت و ایرلماز (۲۰۱۹)، وانگ و همکاران (۲۰۲۲)، ایرلماز و انلو (۲۰۲۳)، اوزکوت و همکاران (۲۰۲۴)، دانگ و همکاران (۲۰۲۵) و منش و همکاران (۲۰۲۵) اشاره نمود. در این مقاله، هدف ما مروری بر استفاده از خانواده توزیع‌های فاز-نوع در چند مدل شوک جدید و مختلف معرفی شده در چند سال اخیر است. برای این منظور به‌طور ویژه، سه مدل شوک: کرانگین تعمیم یافته^۳، دلتای تعمیم یافته^۴ و تجمعی تعمیم یافته^۵ را مورد بررسی قرار می‌دهیم و ویژگی‌های توزیعی و قابلیت اعتمادی طول عمر سیستم‌های تحت این گونه شوک‌ها را بر مبنای توزیع‌های فاز-نوع بیان می‌کنیم. در ادامه این مقاله نخست، به معرفی اجمالی خانواده توزیع‌های فاز-نوع گسسته و پیوسته می‌پردازیم. سپس مدل شوک‌های پایه‌ای را بصورت مختصر توضیح داده و در نهایت نحوه استفاده از خانواده توزیع‌های فاز-نوع را برای مدل شوک‌های بیان شده در بالا را شرح می‌دهیم.

۲ توزیع‌های فاز-نوع

توزیع‌های فاز-نوع به هر دو صورت گسسته و پیوسته و بر مبنای فرآیندهای مارکف گسسته و پیوسته قابل تعریف هستند. در این بخش نخست، حالت گسسته و سپس حالت پیوسته را بیان می‌کنیم.

تعریف ۱.۲. فرایند مارکف زمان-گسسته $\{X_n, n = 0, 1, 2, \dots\}$ با فضای حالت $L = \{1, 2, \dots, m, m+1\}$ ، که در آن m حالت اول گذرا و حالت $(m+1)$ ام جاذب است، ماتریس احتمال انتقال

$$P = \begin{pmatrix} T & t \\ 0 & 1 \end{pmatrix}$$

که در آن T یک زیر ماتریس $m \times m$ مربوط به احتمال انتقال حالت‌های گذرا است و توزیع اولیه $\pi_0 = (\alpha', \alpha_{m+1})'$ را در نظر

²Neuts

³generalized extreme shock model

⁴generalized δ shock model

⁵generalized cumulative shock model

بگیرید. اگر متغیر تصادفی X را زمان تا جذب فرایند تعریف کنیم در اینصورت X دارای توزیع فاز-نوع گسسته با پارامترهای α و \mathbf{T} است و آن را با نماد $X \sim PH_d(\alpha, \mathbf{T})$ نشان می‌دهیم.

برای این توزیع می‌توان روابط زیر را بدست آورد.

• تابع جرم احتمال:

$$p_X(k) = \alpha' \mathbf{T}^{k-1} \mathbf{t}, \quad k = 1, 2, 3, \dots$$

• تابع توزیع تجمعی:

$$F_X(k) = 1 - \alpha' \mathbf{T}^k \mathbf{1}, \quad k = 0, 1, 2, 3, \dots$$

• تابع مولد احتمال:

$$H_X(z) = E(z^X) = \alpha_{m+1} + z \alpha' (\mathbf{I} - z \mathbf{T})^{-1} \mathbf{t},$$

•

$$\frac{d^k H(z)}{dz^k} \Big|_{z=1} = E(X(X-1) \dots (X-(k-1))) = k! \alpha' (\mathbf{I} - \mathbf{T})^{-k} \mathbf{T}^k \mathbf{1}.$$

تعریف ۲.۲. فرایند مارکف زمان-پیوسته $\{X_t, t \geq 0\}$ با فضای حالت $L = \{1, 2, \dots, m, m+1\}$ که در آن m حالت اول گذرا و حالت $m+1$ مام جاذب است، ماتریس پارامترهای بینهایت کوچک^۶

$$\mathbf{Q} = \begin{pmatrix} \mathbf{T}_{m \times m} & \mathbf{t}_{m \times 1} \\ \mathbf{o}_{1 \times m} & \mathbf{o} \end{pmatrix}$$

که در آن $\mathbf{T}_{m \times m}$ یک زیر ماتریس $m \times m$ مربوط به پارامترهای بی‌نهایت کوچک حالت‌های گذرا است و توزیع اولیه $\pi_0 = (\alpha', \alpha_{m+1})'$ را در نظر بگیرید. اگر متغیر تصادفی پیوسته X را زمان تا جذب فرایند تعریف کنیم در اینصورت X دارای توزیع فاز-نوع پیوسته با پارامترهای α و \mathbf{T} است و آن را با نماد $X \sim PH_c(\alpha, \mathbf{T})$ نشان می‌دهیم.

برای حالت پیوسته نیز می‌توان روابط زیر را بدست آورد.

• تابع چگالی احتمال:

$$f_X(x) = \alpha' \exp(\mathbf{T}x) \mathbf{t}, \quad x \geq 0$$

که در آن

$$\exp(\mathbf{T}x) = \sum_{i=0}^{\infty} \frac{(\mathbf{T}x)^i}{i!}$$

به ماتریس نمایی معروف است و $\mathbf{T}^0 = \mathbf{I}$ تعریف می‌شود.

• تابع توزیع تجمعی:

$$F_X(x) = 1 - \alpha' \exp(\mathbf{T}x) \mathbf{1}, \quad x \geq 0$$

که در آن $\mathbf{1}' = (1, \dots, 1)$.

^۶Infinitesimal

- تابع تبدیل لاپلاس:

$$L_X(s) = E(\exp(-sX)) = \alpha_{m+1} + \alpha'(s\mathbf{I} - z\mathbf{T})^{-1}\mathbf{t},$$

این خانواده از توزیع‌ها دارای ویژگی‌های بستاری جالبی در هر دو حالت گسسته و پیوسته می‌باشند که در پایین به بیان مختصری از آنها می‌پردازیم.

- نسبت به ترکیبات خطی بسته است. به بیان دیگر اگر X_1, \dots, X_n متغیرهای تصادفی مستقل که همگی عضو خانواده توزیع‌های

فاز-نوع گسسته (پیوسته) باشند، آنگاه برای مقادیر ثابت و نامنفی a_1, \dots, a_n

$\sum_{i=1}^n a_i X_i$ نیز دارای توزیع فاز-نوع گسسته (پیوسته) است.

- نسبت به آمیختگی متناهی بسته است. به بیان دیگر اگر X_1, \dots, X_n متغیرهای تصادفی مستقل که همگی عضو خانواده

توزیع‌های فاز-نوع گسسته (پیوسته) باشند، آنگاه برای متغیرهای تصادفی $I_i \sim Ber(p_i)$ با شرط $\sum_{i=1}^n I_i = 1$ نتیجه

می‌شود که توزیع آمیخته متناهی آنها یعنی $\sum_{i=1}^n I_i X_i$ نیز دارای توزیع‌های فاز-نوع گسسته (پیوسته) است.

- نسبت به مجموع تصادفی بسته است. به بیان دیگر اگر X_1, \dots, X_n, \dots دنباله‌ای از متغیرهای تصادفی مستقل و هم‌توزیع

$(i.i.d)$ فاز-نوع پیوسته باشند و N دارای توزیع فاز-نوع گسسته باشد، آنگاه $\sum_{i=1}^N X_i$ دارای توزیع فاز-نوع پیوسته

است.

- نسبت به آماره‌های مرتب بسته است. به بیان دیگر اگر X_1, \dots, X_n متغیرهای تصادفی مستقل که همگی عضو خانواده

توزیع‌های فاز-نوع گسسته (پیوسته) باشند، آنگاه آماره‌های مرتب متناظر با آنها، یعنی $X_{(1)}, \dots, X_{(n)}$ نیز دارای توزیع

فاز-نوع گسسته (پیوسته) می‌باشند.

- خانواده توزیع‌های فاز-نوع در بین خانواده توزیع‌های نامنفی چگال است یعنی برای هر متغیر تصادفی نامنفی دلخواه می‌توان

دنباله‌ای از متغیرهای تصادفی عضو خانواده توزیع‌های فاز-نوع را به گونه‌ای یافت که به متغیر تصادفی مورد نظر در توزیع همگرا

باشد.

- ممکن است برای یک متغیر تصادفی عضو خانواده توزیع‌های فاز-نوعی بیش از یک نمایش از لحاظ شیوه تعریف پارامترها

وجود داشته باشد و لزوماً یکتا نمی‌باشند اما در هر صورت بدیهی است که ویژگی‌های توزیعی در نهایت یکسان خواهند بود.

برای توضیحات بیشتر در مورد این خانواده از توزیع‌ها و اثبات روابط بالا می‌توان به [بلاد و نیلسون \(۲۰۱۷\)](#) و [هی \(۲۰۱۴\)](#) مراجعه

نمود.

۳ مدل‌های شوک

یکی از مدل‌های رایج در قابلیت اعتماد در مدل‌بندی طول عمر سیستم‌ها یا قطعاتی با عامل خرابی خارجی (خارج از سیستم) مدل‌های

شوکی می‌باشند. در این مدل‌ها، خرابی سیستم بر اساس الگوی معینی از نحوه وارد شدن شوک‌ها از لحاظ شدت و فاصله زمانی بین

⁷independent and identical distribution

شوکه‌ها تعریف می‌شود. شوک‌های خارجی، بسته به سیستم مورد مطالعه، می‌توانند نوسانات الکتریکی و دمایی، ضربات و تکانه‌ها، خرابی بخشی از سیستم و غیره باشد. همچنین، این شوک‌ها در مسائل مالی و بیمه نیز بصورت اتفاقات نامطلوب قابل تعبیر است. به عنوان مثال، در یک مدل ریسک، هر ادعای خسارت می‌تواند به عنوان یک شوک تلقی گردد. در قابلیت اعتماد، انواع مختلفی از مدل‌های شوک تعریف شده‌اند و از آنها جهت مدل‌بندی طول عمر سیستم‌های تصادفی مختلف به منظور مطالعه شاخص‌های قابلیت اعتمادی و یا سیاست نگهداری بهینه استفاده شده است.

فرض کنید یک سیستم از لحظه صفر شروع به کار می‌کند و شوک‌ها نیز طبق یک فرایند تصادفی معین، به این سیستم وارد می‌شوند. فاصله زمانی بین شوک $(i-1)$ ام و i ام را با X_i ، شدت شوک i ام را با Y_i ، تعداد شوک‌های وارد شده به سیستم تا لحظه خرابی را با N ، مقدار ثابت و معین حد آستانه را با d و طول عمر سیستم را با $T = \sum_{i=1}^N X_i$ نشان می‌دهیم. مدل‌های شوک مختلفی بر اساس نحوه تعریف N قابل تعریف هستند که از بین آنها سه مدل شوک زیر پایه‌ای و بنیادی هستند که بقیه مدل شوک‌ها نیز بصورت تعمیمی از این مدل شوک‌ها تعریف می‌شوند.

• **مدل شوک کرانگین:** در این مدل تا زمانیکه شدت شوک‌های وارده از یک حد آستانه مانند d بیشتر نشود سیستم فعال است و در نتیجه $N = \min\{n, Y_n > d\}$ (گات و سلر، ۱۹۹۹).

• **مدل شوک تجمعی:** در این مدل تا زمانیکه مجموع شدت شوک‌های وارده از یک حد آستانه مانند d بیشتر نشود سیستم فعال است و در نتیجه $N = \min\{n, Y_1 + \dots, Y_n > d\}$ (گات، ۱۹۹۰).

• **مدل شوک دلتا:** در این مدل، بدون توجه به میزان شدت شوک‌ها، اگر فاصله زمانی بین دو شوک پیاپی از حد آستانه δ کمتر باشد سیستم فعال است و در نتیجه $N = \min\{n, X_1 > \delta, \dots, X_{n-1} > \delta, X_n \leq \delta\}$ (لی و ژائو، ۲۰۰۷).

نویسندگان بسیاری تعمیم‌هایی از مدل‌های بالا و همچنین مدل‌های شوک آمیخته را که ترکیبی از دو یا چند مدل شوک‌های قبلی هستند را ارائه کرده‌اند. برای توضیحات بیشتر به عنوان نمونه می‌توان به گات (۲۰۰۱) و مالور و همکاران (۲۰۰۶) مراجعه کرد.

۴ مدل‌های شوک بر اساس خانواده توزیع‌های فاز-نوع

با توجه به نحوه تعریف N در بسیاری از موارد و تحت شرایط نه‌چندان محدود کننده، می‌توان از توزیع‌های فاز-نوع گسسته برای مدل‌بندی N و همچنین توزیع‌های فاز-نوع پیوسته برای مدل‌بندی T استفاده نمود. در این بخش، نحوه استفاده از توزیع‌های فازنوع را برای چند مدل شوک تعمیم یافته مورد بحث قرار می‌دهیم.

۱.۴ مدل شوک کرانگین تعمیم یافته

این مدل شوک بوسیله بوزبولت و ایرلماز (۲۰۲۰) و برخی خواص آن نیز بوسیله منش و همکاران (۲۰۲۳) مورد بررسی قرار گرفت. در این مدل مانند مدل شوک کرانگین معمولی، اگر شوک وارد شده به سیستم از حد آستانه d عبور کند سیستم از کار می‌افتد با این تفاوت که در حالت تعمیم یافته به جای یک منبع شوک، m منبع شوک مستقل وجود دارند. اگر $Y_{i,k}$ برای $k = 1, \dots, m$ نشان‌دهنده میزان شدت i امین شوک وارد به سیستم از منبع شوک k ام باشد و d_k حد آستانه آن باشد در اینصورت $P_k = P(Y_i^k > d_k)$ احتمال خرابی

سیستم توسط این شوک است. همچنین احتمال وارد شدن شوک از طریق منبع k ام برابر π_k است. تحت این مفروضات و نمادگذاری‌ها، دو نوع مدل کرانگین تعمیم یافته در **بوزبولت و ایرلماز (۲۰۲۰)** تعریف شده است

- **مدل ۱:** در این مدل، در طول فعال بودن سیستم، شوک‌ها می‌توانند از منابع مختلف به سیستم وارد شوند.
- **مدل ۲:** در این مدل، اولین شوک از هر منبعی وارد شود در ادامه بقیه شوک‌ها نیز تا زمان از کار افتادن سیستم فقط از آن منبع وارد می‌شوند.

فرض کنید تعداد شوک‌های وارد شده به سیستم تا زمان خرابی مطابق این دو مدل به ترتیب N_1 و N_2 و طول عمر سیستم تحت این دو مدل به ترتیب $T_1 = \sum_{i=1}^{N_1} X_i$ و $T_2 = \sum_{i=1}^{N_2} X_i$ باشند. در اینصورت در مدل ۱، N_1 زمان جذب یک فرآیند مارکف گسسته با فضای وضعیت $L = \{0, 1\}$ است که در آن ۱ فعال بودن سیستم (وضعیت گذرا) و ۰ از کار افتادن سیستم (وضعیت جاذب) را نشان می‌دهد. بنابراین،

$$N_1 \sim Ph_d(1, 1 - \sum_{i=1}^m p_i \pi_i) \sim Geo(1 - \sum_{i=1}^m p_i \pi_i).$$

در مدل ۲، N_2 زمان جذب یک فرآیند مارکف گسسته با فضای وضعیت $L = \{1, 2, \dots, m, m+1\}$ است. وضعیت i نشان دهنده این است که شوک وارد شده از منبع i باعث از کار افتادن سیستم نشده است و وضعیت $m+1$ از کار افتادن سیستم را نشان می‌دهد. در این حالت، توزیع اولیه فرایند $\mathbf{a} = (\pi_1, \dots, \pi_m)'$ و زیر ماتریس احتمال انتقال حالت‌های گذرا به ترتیب بصورت

$$\mathbf{Q} = \begin{pmatrix} 1 - p_1 & 0 & \dots & 0 \\ 0 & 1 - p_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 - p_m \end{pmatrix}$$

می‌باشند. در نتیجه $N_2 \sim PH_d(\mathbf{a}, \mathbf{Q})$. اگر X_i ها را نیز متغیرهای تصادفی $(i.i.d)$ از توزیع $PH_c(\alpha, \mathbf{A})$ باشند آنگاه T_1 و T_2 نیز عضو خانواده توزیع‌های فازنوع بصورت زیر هستند.

$$\begin{aligned} T_1 &\sim PH_c(\alpha, \mathbf{A} \otimes \mathbf{I} + (-\mathbf{A}'\alpha) \otimes (1 - \sum_{i=1}^n \pi_i p_i)) \\ T_2 &\sim PH_c(\alpha \otimes \mathbf{a}, \mathbf{A} \otimes \mathbf{I} + (-\mathbf{A}'\alpha) \otimes \mathbf{Q}). \end{aligned}$$

در این روابط، نماد \otimes نشان دهنده ضرب کرونکر^۸ بین دو ماتریس است. برای توضیحات بیشتر در مورد این نوع ضرب ماتریسی می‌توان به **گراهام (۲۰۱۸)** مراجعه کرد. حال، با توجه به روابط موجود در توزیع‌های فاز-نوعی می‌توان ویژگی‌های قابلیت اعتمادی مختلف مانند تابع قابلیت اعتماد، تابع نرخ خطر و میانگین طول عمر باقی مانده را برای T_1 و T_2 و بصورت روابط ماتریسی یافت.

۲.۴ مدل شوک تجمعی تعمیم یافته

در مدل شوک تجمعی تعمیم یافته همانند مدل شوک تجمعی، اگر مجموع شوک‌های وارد شده به سیستم از یک حد آستانه عبور کند سیستم از کار می‌افتد اما با این تفاوت که در حالت تعمیم یافته به جای یک منبع شوک، m منبع شوک مستقل وجود دارد که هر شوک

^۸kroncker product

با احتمال π_i از منبع i صادر می‌شود (گانگ و همکاران، ۲۰۲۰). اگر برای $i = 1, \dots, m$ شدت شوک وارد شده از منبع i را با $Y^{(i)}$ نشان دهیم و شوک وارد شده به سیستم را با Z نشان دهیم در اینصورت

$$P(Z > t) = \sum_{i=1}^m \pi_i P(Y^{(i)} > t).$$

در این مدل

$$N = \min\{n, Z_1 + \dots + Z_n > d\}, \quad T = \sum_{i=1}^N X_i.$$

حال اگر فرض کنیم $Y^{(i)} \sim PH_c(\alpha_i, \mathbf{A}_i)$ در اینصورت از بسته بودن خانواده توزیع‌های فازنوعی نسبت به توزیع‌های آمیخته نتیجه می‌شود که

$$Z = \sum_{i=1}^m I_i Y^{(i)} \sim PH_c(\alpha, \mathbf{A}), \quad \alpha = (\pi_1 \alpha_1, \dots, \pi_m \alpha_m), \quad \mathbf{A} = \text{diag}\{\mathbf{A}_1, \dots, \mathbf{A}_m\}.$$

همچنین با توجه به بسته بودن این خانواده نسبت به مجموع، داریم

$$\sum_{i=1}^n Z_i \sim PH_c(\gamma, \Gamma)$$

که در آن

$$\gamma = (\alpha', c\alpha', \dots, c^{n-1}\alpha')', \quad c = 1 - \alpha'1 \quad \mathbf{A}^\circ = -\mathbf{A}1,$$

و

$$\Gamma = \begin{pmatrix} \mathbf{A} & \mathbf{A}^\circ \alpha' & \dots & c^{n-2} \mathbf{A}^\circ \alpha' \\ \circ & \mathbf{A} & \dots & c^{n-3} \mathbf{A}^\circ \alpha' \\ \vdots & \vdots & \ddots & \vdots \\ \circ & \circ & \dots & \mathbf{A} \end{pmatrix}.$$

اکنون، با داشتن توزیع $\sum_{i=1}^n Z_i$ می‌توانیم توزیع N را از رابطه زیر بیابیم.

$$\begin{aligned} P(N = n) &= P\left(\sum_{i=1}^n Z_i > d, \sum_{i=1}^{n-1} Z_i \leq d\right) \\ &= P\left(\sum_{i=1}^n Z_i > d\right) - P\left(\sum_{i=1}^{n-1} Z_i > d\right). \end{aligned}$$

با در نظر گرفتن فرض استقلال فاصله زمانی بین شوک‌ها و شدت آنها، تابع قابلیت اعتماد سیستم نیز بصورت یک توزیع آمیخته به شکل زیر قابل محاسبه است.

$$P(T > t) = P\left(\sum_{i=1}^N X_i > t\right) = \sum_{n=1}^{\infty} P\left(\sum_{i=1}^n X_i > t\right) p(N = n).$$

علاوه براین، اگر فرض کنیم فاصله زمانی بین ورود شوک‌ها iid از توزیع فاز-نوعی هستند آنگاه $\sum_{i=1}^n X_i$ نیز عضو این خانواده باقی می‌مانند و توزیع T قابل بیان بصورت آمیخته‌ای از توزیع‌های فاز-نوعی هستند.

۳.۴ مدل شوک دلتای تعمیم یافته

مدل شوک دلتای تعمیم یافته بوسیله ایرلماز و انلو (۲۰۲۳) و به عنوان تعمیمی از مدل شوک دلتا ارائه گردید. بر اساس این مدل، سیستم در لحظه صفر فعال است و فاصله بین زمان‌های ورود شوک متغیرهای تصادفی iid می‌باشند. این سیستم تا زمانیکه m شوک پیاپی در یک فاصله زمانی کمتر از δ به سیستم وارد نشود، سیستم فعال است. اگر تعداد شوک‌های وارد شده به این سیستم تا لحظه خرابی را با $N_{\delta,m}$ نشان دهیم در اینصورت برای $n \geq m$ داریم

$$\{N_{\delta,m} = n\} = \{X_1 + \dots + X_m \geq \delta, X_1 + \dots + X_{m+1} \geq \delta, \dots, X_{n-m} + \dots + X_{n-1} \geq \delta, X_{n-m+1} + \dots + X_n < \delta\}$$

و طول عمر سیستم نیز بصورت $T = \sum_{i=1}^{N_{\delta,m}} X_i$ خواهد بود. اگر $m = 1$ باشد این مدل، به مدل شوک δ معمولی تبدیل می‌شود. به جهت وابستگی مجموع زمان‌های تعریف شده در تعریف $N_{\delta,m}$ ، محاسبه دقیق توزیع آن مشکل است ایرلماز و انلو (۲۰۲۳) با تعریف $p_1 = P(X_1 + \dots + X_m \geq \delta, X_1 + \dots + X_{m+1} \geq \delta)$ و $p_2 = P(X_1 + \dots + X_m \geq \delta)$ آن را با متغیر تصادفی $N_{\delta,m}^*$ بصورت زیر تقریب زدند.

$$P(N_{\delta,m}^* = n) = \begin{cases} 0 & n < m \\ 1 - p_2, & n = m \\ p_2 \left(\frac{p_1}{p_2}\right)^{n-m-1} \left(1 - \frac{p_1}{p_2}\right), & n > m. \end{cases}$$

تحت فرض فاز-نوعی بودن توزیع X_i ها و با فرض $1 = P(X_i > 0)$ توزیع تقریبی $N_{\delta,m}$ برای هر دو حالت گسسته و پیوسته بودن X_i ها قابل محاسبه است.

حالت گسسته

در این حالت، فرض کنید $X \sim PH_d(\mathbf{a}, \mathbf{Q}_{k \times k})$ باشد در اینصورت با توجه به خواص توزیع‌های فاز-نوعی $\sum_{i=1}^m X_i \sim PH_d(\pi_m, \mathbf{Q}_m)$ که در آن $\pi_m = (1, 0, \dots, 0)'$ یک بردار mk بعدی، $\mathbf{Q}_1 = \mathbf{Q}$ و

$$\mathbf{Q}_m = \begin{pmatrix} \mathbf{Q}_{m-1} & (\mathbf{I}_{(m-1)k} - \mathbf{Q}_{m-1})\mathbf{a}' \\ 0 & \mathbf{Q} \end{pmatrix}.$$

با در نظر گرفتن این مفروضات p_1 و p_2 از روابط زیر قابل محاسبه هستند.

$$\begin{aligned} p_1 &= \pi'_{m-1} \mathbf{Q}_{m-1}^{\delta-1} \mathbf{1} + \sum_{s=1}^{\delta-1} \mathbf{a}' (\mathbf{Q}^{\delta-s-1} \mathbf{1})' \pi'_{m-1} \mathbf{Q}_{m-1}^{s-1} (\mathbf{I}_{(m-1)k} - \mathbf{Q}_{m-1}) \mathbf{1} \\ p_2 &= \pi'_m \mathbf{Q}_{m-1}^{\delta-1} \mathbf{1}. \end{aligned}$$

حالت پیوسته

اگر فاصله زمانی بین ورود شوک‌ها متغیرهای تصادفی پیوسته باشند مقادیر p_1 و p_2 در حالت کلی بصورت زیر قابل محاسبه هستند.

$$p_1 = \bar{F}_{m-1}(\delta) + \int_0^\delta (\bar{F}(\delta-s))' f_{m-1}(s) ds \quad (1.4)$$

$$p_2 = \bar{F}_{m-1}(\delta) + \int_0^\delta \bar{F}(\delta-s) f_{m-1}(s) ds \quad (2.4)$$

در روابط (۱.۴) و (۲.۴) عبارت \bar{F}_{m-1} و f_{m-1} به ترتیب تابع قابلیت اعتماد و تابع چگالی احتمال $X_1 + \dots + X_{m-1}$ می‌باشند. اگر فرض کنیم $X \sim PH_c(\alpha, A)$ ، در اینصورت با توجه به خواص توزیع‌های فاز-نوعی روابط (۱.۴) و (۲.۴) بصورت زیر قابل محاسبه هستند.

$$p_1 = \alpha' \exp(As) 1 + \int_0^\delta (\alpha' \exp(A(\delta - s)) 1)' \alpha' \exp(As) a^\circ ds$$

$$p_2 = (\alpha', c\alpha') \exp\left(\begin{pmatrix} A & A^\circ \alpha' \\ 0 & A \end{pmatrix} \delta\right) 1.$$

مراجع

- AALEN, O.O. (1995), Phase type distributions in survival analysis, *Scand. J. Statist* , **22**, 447-463.
- Assaf, D. and Levikson, B. (1982), Closure of phase-type distributions under operations arising in reliability theory, *The Annals of Probability* , **10**, 265-269.
- Bladt, M. (2005), A review on phase-type distributions and their use in risk theory, *ASTIN Bulletin: The Journal of the IAA*, **35(1)**, 145-161.
- Bladt, M., & Nielsen, B. F. (2017), *Matrix-exponential distributions in applied probability (Vol. 81, pp. 1-736)*, New York: Springer.
- Bobbio A, Cumani A (1992), ML estimation of the parameters of a PH distribution in triangular canonical form, in: Balbo G, Serazzi G (eds) Computer performance evaluation, Elsevier, Amsterdam, 173-206.
- Bozbulut, A.A. & Eryilmaz, S. (2020), extreme shock models and their applications, *Communications in Statistics – Simulation and Computation*, **49(1)** 110-120.
- Cai, J., and Li, H. (2005), Multivariate risk model of phase type, with applications to shock models., *Insurance: Mathematics and Economics*, **36(2)**, 137-152.
- Dong, Q., Bai, M., Yan, Z., & Wu, B. (2025), An Optimal load adjustment policy for multi-state k-out-of-n balanced systems with self-healing mechanisms, *Reliability Engineering & System Safety*, 111091.
- ERLANG, A. (1909), Sandsynlighedsregning og telefonsamtaler, *Nyt tidsskrift for Matematik* , **20**, 33-39.
- Eryilmaz, S. (2017), Computing optimal replacement time and mean residual life in reliability shock models, *Computers and industrial engineering*, **103**, 40-45.
- Eryilmaz, S., & Unlu, K. D. (2023), A new generalized δ -shock model and its application to 1-out-of-(m+1): G cold standby system, *Reliability Engineering & System Safety*, **234**, 109203.

- Faddy MJ. (1995), Phase-type distributions for failure times, *Math Comput Model* , **22**, 63-70. Phase-type distributions for failure times.
- Gong, M., Eryilmaz, S., & Xie, M. (2020), Reliability assessment of system under a generalized cumulative shock model, *Proceedings of the Institution of Mechanical Engineers, Part O: Journal of Risk and Reliability*, **234(1)**, 129-137.
- Graham, A. (2018), *Kronecker products and matrix calculus with applications*, Courier Dover Publications.
- Gut, A. (1990), Cumulative shock models, *Advances in Applied Probability*, **22 (2)** 504-507.
- Gut, A. (2001), Mixed shock models, *Bernoulli*, **7(3)** 541-555.
- Gut, A., & Hušler, J. R. (1999), Extreme shock models, *Extremes*, **2** 295-307.
- He, Q. M. (2014), *Fundamentals of matrix-analytic methods (Vol. 365)*, New York: Springer.
- Li, H., (2003), Association of multivariate phase-type distributions, with applications to shock models., *Statist. Probab. Lett* , **64**, 281-392.
- Li, Z. and Zhao, P. (2007), Reliability Analysis on the δ -Shock Model of Complex Systems, *IEEE Transactions on Reliability*, **56(2)** 340-348.
- Mallor, F., Omei, E. and Santos, J. (2006), Asymptotic results for a run and cumulative mixed shock model, *Journal of Mathematical Sciences*, **138(1)** 5410-5414.
- Manesh, S. F., Izadi, M., & Khaledi, B. E. (2023), On stochastic ordering among extreme shock models, *Probability in the Engineering and Informational Sciences*, **37(4)** 961-972.
- Manesh, S. F., Izadi, M., & Khaledi, B. E. (2025), A new mixed generalized δ -shock model, *statistics & probability letters*, submitted.
- Montoro-Cazorla, D., Pérez-Ocón, R., and Segovia, M. C. (2009), Shock and wear models under policy N using phase-type distributions, *Applied mathematical modelling*, **33(1)**, 543-554.
- Neuts, M.F., (1975), Probability distributions of phase type, In: *Liver Amicorum Prof. Emeritus H. Florin*, University of Louvain, Belgium, **20**, 173-206.
- O'CINNEIDE, C.A. (1990), Characterization of phase-type distributions, *Comm. Statist Stochastic Models* , **6**, 1-57.

- Ozkut, M., and Eryilmaz, S. (2019), Reliability analysis under Marshall–Olkin run shock model, *Journal of Computational and Applied Mathematics*, **349**,52-59.
- Ozkut, M., Kan, C., & Franko, C. (2024), Analyzing the multi-state system under a run shock mode, *Probability in the Engineering and Informational Sciences*, **38(4)**,619-931.
- Segovia, M. C., and Labeau, P. E. (2013), Reliability of a multi-state system subject to shocks using phase-type distributions, *Applied mathematical modelling*, **37(7)**,488-554.
- Wang, X., Zhao, X., Wu, C., & Wang, S. (2022), Mixed shock model for multi-state weighted k-out-of-n: F systems with degraded resistance against shocks, *Reliability Engineering & System Safety*, **217**,108098.
- Zhao, X., Guo, X., and Wang, X. (2018), Reliability and maintenance policies for a two-stage shock model with self-healing mechanism, *Reliability Engineering and System Safety*, **172**,185-194.

A review of the use of the family of phase-type distributions in the generalized shock model.

Sirous Fathi Manesh ¹, Muhyiddin Izadi ²

¹ Department of Statistics, University of Kurdistan, Sanandaj, Iran

²Department of Statistics, Razi University, Kermanshah, Iran

Abstract: The utilization of phase-type (PH) distributions - defined through discrete-time and continuous-time Markov processes - is increasingly prevalent in modelling non-negative discrete and continuous random variables. In recent years, these distributions have gained considerable attention among reliability researchers for their application in various shock models. This paper presents a comprehensive review of PH distribution applications in three principal shock models: the generalized shock model, generalized cumulative shock model, and generalized delta shock model. We provide theoretical justification for employing PH distributions in these shock models, followed by derivation of both shock count distribution and system lifetime distribution using the inherent matrix-based relationships characteristic of PH distributions.

Keywords: Generalized cumulative shock model, Generalized delta shock model, Generalized extreme shock model, Markov Processes, Phase-type distribution family.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



بررسی عملکرد مجانبی برآوردگر ترکیبی در مدل‌های خطی با خطا در اندازه‌گیری

فاطمه قپانی^۱،

گروه آمار و ریاضی، واحد شوشتر، دانشگاه آزاد اسلامی، شوشتر، ایران

چکیده: در این مقاله یک برآوردگر ترکیبی برای بهبود دقت برآوردها در مدل‌های با خطی در اندازه‌گیری خطی تحت محدودیت‌های خطی تصادفی معرفی و ویژگی‌های مجانبی برآوردگر معرفی شده تعیین می‌شود. عملکرد برآوردگر معرفی شده با استفاده از معیار ماتریس میانگین مربعات خطا نسبت به بعضی برآوردهای موجود مورد بررسی قرار می‌گیرد. در پایان عملکرد برآوردگر معرفی شده با استفاده از یک مطالعه‌ی شبیه سازی ارزیابی می‌شود.

واژه‌های کلیدی: برآوردگر ترکیبی، مدل خطی با خطا در اندازه‌گیری، محدودیت‌های خطی تصادفی، هم خطی .
کد موضوع بندی ریاضی (۲۰۲۰): 62J07، 62J05.

۱ مقدمه

وجود وابستگی خطی بین ستون‌های ماتریس متغیرهای توضیحی باعث کاهش دقت برآورد و حساسیت برآورد نسبت به تغییر کوچکی در مجموعه داده‌ها می‌شود. یک روش برخورد با مسئله هم خطی به دست آوردن برآورد پارامترهای مدل با اضافه کردن اطلاعات پیشین به صورت محدودیت‌های خطی دقیق یا تصادفی روی پارامترهای نامعلوم مدل است [رائو و همکاران \(۲۰۰۸\)](#). در این مقاله یک برآوردگر ترکیبی با استفاده از روش کبریا - لکمن و در مدل‌های خطی با خطا در اندازه‌گیری تحت محدودیت‌های خطی تصادفی معرفی می‌شود. هدف اصلی برآورد ترکیبی، ترکیب برآوردهای مختلف برای کاهش انحراف و خطا در برآورد ضرایب مدل است. یک مدل با خطا در اندازه‌گیری را به صورت

$$y = Z\beta + \varepsilon, \varepsilon \sim N(0, \sigma^2 I_n)$$

$$X = Z + U, U \sim MN(0, I_n \otimes \Sigma)$$

^۱ سخنران، Fatemeh.Ghapani@iau.ac.ir

در این مدل $y = (y_1, y_2, \dots, y_n)'$ بردار $n \times 1$ از مشاهدات، Z ماتریس $n \times p$ از متغیرهای توضیحی، β بردار $p \times 1$ از پارامترهای ناشناخته مدل و ε بردار $n \times 1$ از خطای تصادفی مدل است. X ماتریس مشاهده شده Z با خطای اندازه گیری U است. فرض کنید ε و U از هم مستقل و U یک ماتریس با عناصر معلوم است. فرض کنید بردار پارامتر β تحت محدودیت‌های خطی تصادفی به صورت زیر باشد:

$$r = R\beta + e, e \sim N(0, \sigma^2 W)$$

در این رابطه r یک بردار معلوم $m \times 1$ از مشاهدات، R یک ماتریس معلوم $m \times p$ پرتبه سطری $m < p$ و e بردار $m \times 1$ از خطای تصادفی است. برآوردگر آمیخته در مدل‌های با خطا در اندازه‌گیری توسط قیانی و همکاران (۲۰۱۵) به صورت زیر معرفی شد:

$$\hat{\beta}_{ME} = S_r^{-1}(X'y + R'W^{-1}r),$$

$$S_r = S + R'W^{-1}R.$$

۲ برآوردگر ترکیبی

برای تعیین برآوردگر جدید در مدل‌های با خطا در اندازه‌گیری براساس کبریا و لکمن (۲۰۲۰) رابطه‌ی زیر را در نظر می‌گیریم:

$$\Phi_1 = (y - X\beta)'(y - X\beta) - n\beta'\Sigma\beta + k[(\beta + \hat{\beta})'(\beta + \hat{\beta}) - c]$$

در این رابطه c ضریب ثابت و k ضریب لاکرانژ در نظر گرفته می‌شود. با مشتق گرفتن از Φ_1 برآوردگر $\hat{\beta}_{KL}$ به صورت زیر تعیین می‌شود:

$$\hat{\beta}_{KL} = S_k^{-1}(X'y - k\hat{\beta})$$

به منظور در نظر گرفتن محدودیت‌های تصادفی در برآوردگر $\hat{\beta}_{KL}$ رابطه‌ی زیر را در نظر می‌گیریم:

$$\Phi_2 = (y - X\beta)'(y - X\beta) - n\beta'\Sigma\beta + k[(\beta + \hat{\beta}_{ME})'(\beta + \hat{\beta}_{ME}) - c] + (r - R\beta)'W^{-1}(r - R\beta)$$

با مشتق گرفتن از Φ_2 نسبت به β برآوردگر ترکیبی به صورت زیر معرفی می‌شود:

$$\hat{\beta}_{MKL} = (S + R'W^{-1}R + kI_p)^{-1}(X'y - k\hat{\beta}_{ME} + R'W^{-1}r)$$

۳ خواص مجانبی

در این بخش خواص مجانبی برآوردگر $\hat{\beta}_{MKL}$ مورد بررسی قرار می‌گیرد. فرض می‌کنیم وقتی n به بینهایت میل می‌کند حدهای $n^{-1}(Z'Z + R'W^{-1}R - kI_p)$ و $n^{-1}(Z'Z + R'W^{-1}R + kI_p)$ ، $n^{-1}(Z'Z + R'W^{-1}R)$ وجود دارند.

قضیه ۱.۳. تحت فرض‌های بیان شده $\hat{\beta}_{MKL}$ به طور مجانبی دارای توزیع نرمال با میانگین و کواریانس ماتریس به صورت زیر است

$$E(\hat{\beta}_{MKL}) = G_{rk}\beta$$

$$AVar(\hat{\beta}_{MKL}) = A_r^{-1}(G_{rk}BG_{rk} + \sigma^2 G_{rk}A_r G_{rk})A_r^{-1}$$

که در آن $G_{rk} = (Z'Z + R'W^{-1}R + kI_p)^{-1}(Z'Z + R'W^{-1}R - kI_p)$ ، $B = (n\sigma^2 + \beta'Z'Z\beta)\Sigma$ و $A_r = Z'Z + R'W^{-1}R$

اثبات. از آنجا که $E(X'X) = Z'Z + n\Sigma$ فانگ و همکاران (۲۰۰۳)، بنابراین داریم:

$$n^{-1}(X'X + R'W^{-1}R) = n^{-1}(Z'Z + R'W^{-1}R) + \Sigma + O_p(n^{-\frac{1}{2}})$$

همچنین داریم

$$\begin{aligned}\sqrt{n}\hat{\beta}_{MKL} &= \left[n^{-1}(Z'Z + R'W^{-1}R) + O_p(n^{-\frac{1}{2}}) \right]^{-1} n^{-\frac{1}{2}} \left[G_{rk}(X'y + R'W^{-1}r) + O_p(n^{\frac{1}{2}}) \right] \\ &= \left[I_p + O_p(n^{-\frac{1}{2}}) \right]^{-1} C^{-1} n^{-\frac{1}{2}} \left[G_{rk}(X'y + R'W^{-1}r) + O_p(n^{\frac{1}{2}}) \right]\end{aligned}$$

زمانی که حد $C = n^{-1}(Z'Z + R'W^{-1}R)$ وجود دارد رابطه‌ی زیر نتیجه می‌شود:

$$\sqrt{n}\hat{\beta}_{MKL} = C^{-1}\xi + O_p(n^{-\frac{1}{2}}),$$

که $\xi = n^{-\frac{1}{2}} [G_{rk}(X'y + R'W^{-1}r)]$ به طور مجانبی نرمال است (فانگ و همکاران (۲۰۰۳)).

از آنجا که $E[G_{rk}(X'y + R'W^{-1}r)] = G_{rk}A_r\beta$ ، نتیجه گرفته می‌شود که $E(\xi) = n^{-\frac{1}{2}}G_{rk}A_r\beta$. همچنین

$$\sqrt{n}(\hat{\beta}_{MKL} - G_{rk}\beta) = C^{-1}[\xi - E(\xi)] + O_p(n^{-\frac{1}{2}})$$

بنابراین $\sqrt{n}(\hat{\beta}_{MKL} - G_{rk}\beta)$ دارای توزیع مجانبی نرمال با میانگین صفر است. همچنین داریم. $AVar(\sqrt{n}\hat{\beta}_{MKL}) =$

$C^{-1}Var(\xi)C^{-1}$ برای محاسبه واریانس ξ روش زیر را بکار می‌بریم:

$$\begin{aligned}Var(\xi) &= E_{y_r}[Var(\xi|y_r)] + Var_{y_r}[E(\xi|y_r)] \\ &= n^{-1}E_{y_r}(G_{rk}y'y\Sigma G_{rk}) + n^{-1}Var_{y_r}[G_{rk}(Z'y + R'W^{-1}r)],\end{aligned}$$

در این رابطه E_{y_r} و Var_{y_r} امید و واریانس نسبت به بردار $y_r' = (y', r')$ می‌باشند. از آنجا که $E(y'y) = n\sigma^2 + \beta'Z'Z\beta$ و

$Var(\xi) = n^{-1}(G_{rk}BG_{rk} + \sigma^2 G_{rk}A_r G_{rk})$ می‌توان نوشت $Var_{y_r}[G_{rk}(Z'y + R'W^{-1}r)] = \sigma^2 G_{rk}A_r G_{rk}$

بنابراین خواهیم نوشت:

$$AVar(\hat{\beta}_{MKL}) = A_r^{-1}(G_{rk}BG_{rk} + \sigma^2 G_{rk}A_r G_{rk})A_r^{-1}.$$

□

۴ ماتریس میانگین مربعات خطا

برای بررسی کارایی برآوردگر معرفی شده نسبت به برآوردهای موجود دیگر از معیار ماتریس میانگین مربعات خطا استفاده می‌شود. ماتریس میانگین مربعات خطای برآوردگر $\hat{\beta}$ از پارامتر β به صورت زیر تعریف می‌شود:

$$MSEM(\hat{\beta}) = E(\hat{\beta} - \beta)(\hat{\beta} - \beta)' = Var(\hat{\beta}) + Bias(\hat{\beta})Bias(\hat{\beta})'$$

که $Bias(\hat{\beta}) = E(\hat{\beta}) - \beta$ مقدار اریبی برآوردگر $\hat{\beta}$ است. مقدار میانگین مربعات خطا معیار دیگری برای بررسی کارایی یک برآوردگر است که به صورت زیر تعریف می‌شود:

$$MSE(\hat{\beta}) = tr [MSEM(\hat{\beta})] = tr [Var(\hat{\beta})] + Bias(\hat{\beta})'Bias(\hat{\beta})$$

ماتریس میانگین مربعات برآوردگرهای ذکر شده به طور مجانبی به صورت زیر نوشته می‌شود:

$$AMSEM(\hat{\beta}) = A^{-1}(B + \sigma^2 A)A^{-1}$$

$$AMSEM(\hat{\beta}_{KL}) = A^{-1}(G_k B G_k + \sigma^2 G_k A G_k)A^{-1} + b_1 b_1'$$

$$AMSEM(\hat{\beta}_{ME}) = A_r^{-1}(B + \sigma^2 A_r)A_r^{-1}$$

$$AMSEM(\hat{\beta}_{MKL}) = A_r^{-1}(G_{rk} B G_{rk} + \sigma^2 G_{rk} A_r G_{rk})A_r^{-1} + b_2 b_2'$$

$$b_2 = -\sigma^2 A_{rk}^{-1} \beta \text{ و } A_k = Z'Z + kI_p, \quad b_1 = -\sigma^2 A_k^{-1} \beta \text{ که}$$

۱.۴ مقایسه ماتریس میانگین مربعات خطای برآوردهای $\hat{\beta}$ و $\hat{\beta}_{KL}$

به منظور مقایسه کارایی دو برآوردگر $\hat{\beta}_{KL}$ و $\hat{\beta}$ ماتریس زیر را در نظر می‌گیریم:

$$\Delta_1 = AMSEM(\hat{\beta}) - AMSEM(\hat{\beta}_{KL}) = D_1 - b_1 b_1'$$

$$= \sigma^2 A_k^{-1}(A^{-1}B + BA^{-1} + \sigma^2 I_p)A_k^{-1} - b_1 b_1',$$

از آنجا که $D_1 = \sigma^2 A_k^{-1}(A^{-1}B + BA^{-1} + \sigma^2 I_p)A_k^{-1}$ یک ماتریس معین مثبت است. بنابراین Δ_1 معین مثبت است هرگاه $b_1' D_1^{-1} b_1 \leq 1$.

۲.۴ مقایسه ماتریس میانگین مربعات خطای برآوردهای $\hat{\beta}_{ME}$ و $\hat{\beta}_{MKL}$

به منظور مقایسه ماتریس میانگین مربعات خطای دو برآوردگر $\hat{\beta}_{MKL}$ و $\hat{\beta}_{ME}$ ماتریس زیر را در نظر می‌گیریم:

$$\Delta_2 = AMSEM(\hat{\beta}_{ME}) - AMSEM(\hat{\beta}_{MKL}) = D_2 - b_2 b_2'$$

$$= \sigma^2 A_{rk}^{-1}(A_r^{-1}B + BA_r^{-1} + \sigma^2 I_p)A_{rk}^{-1} - b_2 b_2',$$

از آنجا که $D_2 = \sigma^2 A_{rk}^{-1}(A_r^{-1}B + BA_r^{-1} + \sigma^2 I_p)A_{rk}^{-1}$ یک ماتریس معین مثبت است بنابراین Δ_2 معین مثبت است هرگاه $b_2' D_2^{-1} b_2 \leq 1$.

۳.۴ مقایسه ماتریس میانگین مربعات برآوردگرهای $\hat{\beta}_{KL}$ و $\hat{\beta}_{MKL}$

به منظور بررسی دو برآوردگر در معیار ماتریس میانگین مربعات خطا ماتریس زیر در نظر گرفته می‌شود:

$$\Delta_3 = AMSEM(\hat{\beta}_{KL}) - AMSEM(\hat{\beta}_{MKL}) = D_3 + b_1 b_1' - b_2 b_2',$$

در این رابطه $D_3 = AVar(\hat{\beta}_{KL}) - AVar(\hat{\beta}_{MKL})$. واضح است که $AVar(\hat{\beta}_{KL}) > 0$ و $AVar(\hat{\beta}_{MKL}) > 0$. بنابراین هرگاه داشته باشیم $1 < \lambda_{\max} [AVar(\hat{\beta}_{MKL}) AVar(\hat{\beta}_{KL})^{-1}] < 1$ آنگاه $D_3 > 0$ بنابراین $\Delta_3 \geq 0$ اگر فقط اگر داشته باشیم $b_2' (D_3 + b_1 b_1')^{-1} b_2 \leq 1$.

۵ برآورد پارامتر k

در این بخش پارامتر اریب k را طوری برآورد می‌کنیم که ماتریس Δ_1 یک ماتریس معین مثبت شود. ماتریس Δ_1 را می‌توان به صورت زیر نوشت:

$$\Delta_1 = 2k A_k^{-1} (A^{-1} B + B A^{-1} + 2\sigma^2 I_p - 2k \beta' \beta) A_k^{-1},$$

از آنجا که $A_k^{-1} > 0$ ، $A^{-1} B > 0$ و $B A^{-1} > 0$. بنابراین Δ_1 یک ماتریس معین مثبت است هرگاه $2\sigma^2 I_p - 2k \beta' \beta$ نیمه معین مثبت باشد. بنابراین یک برآورد مناسب برای پارامتر اریب k به صورت $\hat{k} < \frac{\hat{\sigma}^2}{\hat{\beta}' \hat{\beta}}$ تعیین می‌شود.

۶ مطالعه‌ی شبیه‌سازی

به منظور ارزیابی نتایج تئوری یک مطالعه‌ی شبیه‌سازی انجام شده است. برای در نظر گرفتن درجات مختلف هم‌خطی با توجه به مقاله [مک‌دونالد و گالارنو \(۱۹۷۵\)](#) متغیرهای توضیحی به صورت زیر تولید می‌شوند:

$$z_{il} = (1 - \rho^2)^{\frac{1}{2}} w_{il} + \rho w_{i,p+1}, i = 1, \dots, n, l = 1, \dots, p,$$

که در آن w_{ij} متغیرهای شبه تصادفی مستقل نرمال استاندارد و ρ^2 همبستگی بین دو متغیر توضیحی را نشان می‌دهد. داده‌ها همچنین استاندارد می‌شوند تا $Z'Z$ به فرم ماتریس همبستگی شود. j -امین مجموعه از داده‌های شبیه‌سازی از مدل با خطا در اندازه‌گیری تحت محدودیت‌های خطی تصادفی به صورت زیر تولید می‌شود:

$$y_j = Z\beta + \varepsilon_j,$$

$$X_j = Z + U_j, j = 1, \dots, 1000,$$

$$r_j = R\beta + e_j,$$

در این رابطه $y_j = (y_{1j}, \dots, y_{nj})'$ ، $Z = (z^{(1)}, z^{(1)}, z^{(3)})$ ، به طوری که $z^{(l)} = (z_{1l}, \dots, z_{nl})'$ و $l = 1, 2, 3$ و همچنین R یک ماتریس معلوم و عناصر آن از توزیع نرمال استاندارد تولید شده‌اند و $\varepsilon_j \sim N(0, \sigma^2 I_n)$ تولید می‌شود. در این مطالعه $p = 3$ ، $\sigma^2 = 0.5$ ، $m = 1$ و $\Sigma = \text{diag}(0.05, 0.05, 0.05)$ و $e_j \sim N(0, \sigma^2 I_m)$ و $r_j = (r_{1j}, \dots, r_{mj})'$.

در نظر گرفته شده است. برای انجام شبیه‌سازی از نرم افزار R استفاده شده است. برای هر ترکیب از پارامترها ۱۰۰۰ تکرار انجام می‌شود. عملکرد برآوردگرها با استفاده از مقدار میانگین مربعات خطا از رابطه‌ی زیر ارزیابی می‌شود:

$$MSE(\tilde{\beta}) = \frac{1}{1000} \sum_{j=1}^{1000} \sum_{l=1}^2 (\tilde{\beta}_{lj} - \beta_l)^2$$

نتایج شبیه‌سازی در جداول ۱ و ۲ ارائه شده است.

جدول ۱: برآورد مقدار میانگین مربعات خطا برآوردهای مختلف با $n = 50$

$\hat{\beta}_{MKL}$	$\hat{\beta}_{ME}$	$\hat{\beta}_{KL}$	$\hat{\beta}$	ρ
۰/۰۵۷۱	۰/۰۶۱۱۰	۰/۰۵۷۵	۰/۰۶۱۵	۰/۸۰
۰/۱۱۲۲	۰/۱۲۵۱	۰/۱۱۳۵	۰/۱۲۶۷	۰/۹۰
۰/۲۴۹۱	۰/۲۹۵۱	۰/۲۵۴۴	۰/۳۰۲۳	۰/۹۵

جدول ۲: برآورد مقدار میانگین مربعات خطا برآوردهای مختلف با $n = 100$

$\hat{\beta}_{MKL}$	$\hat{\beta}_{ME}$	$\hat{\beta}_{KL}$	$\hat{\beta}$	ρ
۰/۰۴۵۵	۰/۰۴۷۷۰	۰/۰۴۵۷	۰/۰۴۷۹	۰/۸۰
۰/۱۰۵۸	۰/۱۱۴۷	۰/۱۰۶۷	۰/۱۱۵۸	۰/۹۰
۰/۲۳۷۰	۰/۲۶۶۹	۰/۲۴۰۵	۰/۲۷۱۳	۰/۹۵

نتایج شبیه‌سازی ارائه شده در جداول ۱ و ۲ بیانگر آن است که با افزایش تعداد نمونه مقدار میانگین مربعات خطای برآوردها کاهش می‌یابد. با افزایش سطح هم‌خطی مقدار میانگین مربعات خطای همه برآوردها افزایش پیدا می‌کند و در تمام موارد مقدار میانگین مربعات خطای برآوردگر $\hat{\beta}_{MKL}$ نسبت به سایر برآوردهای ذکر شده کمتر است. بنابراین در حضور هم‌خطی استفاده از برآوردگر جدید پیشنهاد می‌شود.

بحث و نتیجه‌گیری

در این مقاله یک برآوردگر ترکیبی در مدل‌های خطی با خطا در اندازه‌گیری تحت محدودیتهای خطی تصادفی معرفی گردید. ویژگی‌های مجانبی برآوردگر معرفی شده بررسی و با استفاده از معیار ماتریس میانگین توان‌های دوم خطا کارایی برآوردها مورد ارزیابی قرار گرفت. با استفاده از مطالعه‌ی شبیه‌سازی نشان داده شد که برای مقدار مشخصی از پارامتر اریب برآوردگر معرفی شده کارایی بیشتری نسبت به برآوردگر محدود شده دارد.

- Fung, W. K., Zhong, X. P. and Wei, B. C. (2003). On estimation and influence diagnostics in linear mixed measurement error models. *American Journal of Mathematical and Management Sciences*, **23**, 37–59.
- Ghapani, F., A. R. Rasekh, and B. Babadi. (2015). Detection of outliers and influential observations in linear ridge measurement error models with stochastic linear restrictions. *Journal of Sciences Islamic Republic of Iran*. **26**(4), 355–66.
- Kibria, B. M. G. and Lukman, A. F. (2020). A new ridge-type estimator for the linear regression model. *Simulations and applications. Hindawi Scientifica*. <https://doi.org/10.1155/2020/9758378>.
- McDonald, C., Galarneau, D. A. (1975). A Monte Carlo Evaluation of some Ridge-Type Estimators. *Journal of the American Statistical Association*. **70**, 407-416.
- Rao, C. R., Toutenburg, H., Shalabh and Heumann, C. (2008), *Linear Models and Generalizations*, Springer. Berlin.

Investigating the asymptotic performance of the hybrid estimator in the linear measurement error models

Fatemeh Ghapani

Department of Mathematics and Statistics, Sho.C., Islamic Azad University, Shoushtar, Iran

Abstract: In this paper, a hybrid estimator is introduced for improving the accuracy of estimators in linear measurement error models with stochastic linear restrictions, and the asymptotic properties of the introduced estimator are investigated. The performance of the introduced estimator is examined using the mean square error matrix criterion compared to some existing estimators. Finally, the performance of the introduced estimator is evaluated using a simulation study.

Keywords: Hybrid estimator, Linear measurement error model, Stochastic linear restrictions, Multicollinearities.

Mathematics Subject Classification (2020): 62J05, 62J07.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



مقایسه روش‌های کلاسیک و بیزی در برآورد پارامتر قابلیت اعتماد تنش-مقاومت برای توزیع گمپرتز تعمیم‌یافته یک‌ه

اعظم کاراندیش مروستی^۱، احسان ارمز^۲، مریم بصیرت^۳

گروه ریاضی و آمار، واحد مشهد، دانشگاه آزاد اسلامی، مشهد، ایران

گروه ریاضی و آمار، واحد مشهد، دانشگاه آزاد اسلامی، مشهد، ایران

گروه ریاضی و آمار، واحد مشهد، دانشگاه آزاد اسلامی، مشهد، ایران

چکیده: در این مقاله، برآورد پارامتر قابلیت اعتماد تنش-مقاومت $R = P(Y < X)$ مورد بررسی قرار گرفته است؛ حالتی که در آن X (مقاومت) و Y (تنش) متغیرهای تصادفی مستقلی هستند که از توزیع گمپرتز تعمیم‌یافته یک‌ه (UGG) پیروی می‌کنند. به‌منظور برآورد پارامتر R ، از سه رویکرد اصلی شامل روش ماکسیمم درستنمایی (MLE)، روش خودگردان (Bootstrap)، و روش بیزی (Bayesian) با بهره‌گیری از الگوریتم‌های گیبز و متروپلیس-هستینگس استفاده شد. مطالعه حاضر علاوه بر شبیه‌سازی داده‌ها، از داده‌های واقعی بارندگی دو منطقه اقلیمی ایران (یاسوج و ساری) طی سال‌های ۱۳۹۰ تا ۱۴۰۲ بهره برد تا اعتبار برآوردها در شرایط واقعی سنجیده شود. نتایج نشان داد که روش بیزی، به‌ویژه با در نظر گرفتن وزن‌دهی متفاوت داده‌ها (نسبت‌های $m = 0.3$ و $n = 0.7$)، عملکرد پایدارتر و فاصله‌های اعتباری دقیق‌تری نسبت به سایر روش‌ها ارائه می‌دهد. همچنین روش خودگردان توانست در داده‌های محدود، پوشش مناسبی برای برآورد R فراهم کند، اگرچه میانگین برآورد آن نسبت به MLE پایین‌تر بود. در مجموع، برآوردهای بیزی در شرایط داده‌های نامتقارن یا کوچک، دقت بالاتری نشان دادند، که این امر در مقایسه با مطالعات پیشین نیز تأیید می‌شود. یافته‌های تحقیق نشان می‌دهند که ترکیب توزیع‌های انعطاف‌پذیر مانند گمپرتز تعمیم‌یافته یک‌ه با رویکردهای بیزی می‌تواند چارچوب قابل اعتمادتری برای تحلیل قابلیت اطمینان سامانه‌ها در حوزه‌های مهندسی و اقلیم‌شناسی فراهم آورد.

واژه‌های کلیدی: توزیع گمپرتز تعمیم‌یافته یک‌ه، تنش-مقاومت، ماکسیمم درستنمایی، روش بیزی.

کد موضوع‌بندی ریاضی (۲۰۲۰): 62Mxx

۱ مقدمه

توزیع‌های آماری مهم‌ترین ابزار در مدل‌بندی داده‌های واقعی هستند. با توجه به رشد روزافزون آمار در تمامی زمینه‌های علمی، معرفی توزیع‌های جدید با انعطاف بیشتر به منظور تبیین و تفسیر بهتر پدیده‌ها، یک ضرورت انکارناپذیر است. در دهه‌های اخیر تلاش‌های زیادی برای تحقق این موضوع انجام شده است. بیشتر روش‌هایی که قبل از سال ۱۹۸۰ میلادی برای ساخت توزیع‌ها استفاده شده است مبتنی بر معادلات دیفرانسیل بود که به پیدایش دستگاه پیرسنی و معادلات بور انجامید. بسیاری از توزیع‌های کلاسیک در این قالب می‌گنجد. بعد از سال ۱۹۸۰ میلادی تا ۲۰۱۴ میلادی بیشتر روش‌ها مبتنی بر افزودن پارامترها به توزیع‌های کلاسیک به منظور افزایش انعطاف‌پذیری، انجام شده است. امروزه استفاده از تجهیزات اندازه‌گیری مدرن مانند حسگرها در چندین شاخه از علوم کاربردی از جمله قابلیت اطمینان، مهندسی کیفیت و علوم زیست‌پزشکی و اجتماعی به سرعت در حال رشد است. این موضوع منجر به ایجاد مجموعه داده‌های جدید با رفتارهای مختلف در چولگی و کشیدگی می‌شود که نمی‌توانند به خوبی توسط توزیع‌های آماری کلاسیک مانند توزیع‌های نرمال و وایبول توصیف شوند. بر این اساس، خواسته‌های جدیدی برای توسعه توزیع‌های انعطاف‌پذیرتر برای پوشش این پدیده‌ها به وجود آمده است. در بیشتر موقعیت‌ها، توزیع‌های جدید با گسترش توزیع‌های شناخته شده موجود ساخته می‌شوند. یکی از توزیع‌های معروف در قابلیت اطمینان، زیست‌شناسی و کهرولت شناسی توزیع گمپرتز است که **گمپرتز (۱۸۲۵)** آن را معرفی کرد. **ویلکنز (۲۰۰۱)** نشان داد که توزیع گمپرتز روابط دقیق یا حدی با توزیع‌هایی مانند وایبول، لوژستیک تعمیم یافته، نمایی، نمایی دوگانه و گامبل دارد. قابلیت اعتماد تنش-مقاومت که به صورت $R = P(Y < X)$ تعریف می‌شود، یکی از معیارهای مهم در تحلیل قابلیت سامانه‌هاست. در این مقاله، دو متغیر تصادفی X و Y به ترتیب معرف مقاومت و تنش هستند و فرض می‌شود که هر دو از توزیع گمپرتز تعمیم‌یافته یکه (UGG) تبعیت می‌کنند. تابع چگالی احتمال این توزیع در رابطه (۱.۱) ارائه شده است:

$$f(x) = \theta \lambda x^{-(c+1)} e^{-\frac{\lambda}{c}(x^{-c}-1)} [1 - e^{-\frac{\lambda}{c}(x^{-c}-1)}]^{\theta-1}; \quad 0 < x < 1, \theta > 0, \lambda > 0, c > 0 \quad (1.1)$$

لذا مقدار تحلیلی R به صورت زیر به دست می‌آید:

$$R = P(Y < X) = \int_0^1 P(Y < X | X = x) f_X(x) dx = \frac{\theta_2}{\theta_1 + \theta_2}.$$

بنابراین قابلیت اعتماد تنش-مقاومت فقط تابعی از θ_1 و θ_2 است، که صرفاً به پارامترهای شکل توزیع وابسته است.

۲ روش‌های برآورد

۱.۲ روش ماکسیمم درستنمایی (MLE)

تابع لگاریتم درستنمایی داده‌ها از معادله (۱.۲) استخراج شده است:

$$\begin{aligned} l(\theta_1, \theta_2, \lambda, c) &= \log L(\theta_1, \theta_2, \lambda, c) \\ &= m \log(\theta_1) + n \log(\theta_2) + (m+n) \log \lambda - (c+1) \left[\sum_{i=1}^m \log x_i + \sum_{i=1}^n \log y_i \right] \\ &\quad - \frac{\lambda}{c} \left[\sum_{i=1}^m (x_i^{-c} - 1) + \sum_{i=1}^n (y_i^{-c} - 1) \right] + (\theta_1 - 1) \sum_{i=1}^m \log(1 - e^{-\frac{\lambda}{c}(x_i^{-c}-1)}) \\ &\quad + (\theta_2 - 1) \sum_{i=1}^n \log(1 - e^{-\frac{\lambda}{c}(y_i^{-c}-1)}). \end{aligned} \quad (1.2)$$

با محاسبه مشتق‌های جزئی نسبت به پارامترها، معادلات به‌دست‌آمده برای θ_1 ، θ_2 ، λ و c حل می‌شوند و سپس مقدار \hat{R} از فرمول زیر محاسبه می‌شود:

$$\begin{aligned} \frac{\partial}{\partial \theta_1} l(\theta_1, \theta_2, \lambda, c) &= \frac{m}{\theta_1} + \sum_{i=1}^m \ln(C(x_i)), \\ \frac{\partial}{\partial \theta_2} l(\theta_1, \theta_2, \lambda, c) &= \frac{n}{\theta_2} + \sum_{i=1}^n \ln(C(y_i)), \\ \frac{\partial}{\partial \lambda} l(\theta_1, \theta_2, \lambda, c) &= \frac{m+n}{\lambda} - \left(\sum_{i=1}^m A(x_i) + \sum_{i=1}^n A(y_i) \right) \\ &\quad + (\theta_1 - 1) \left(\sum_{i=1}^m \frac{A(x_i) e^{-\lambda A(x_i)}}{C(x_i)} \right) + (\theta_2 - 1) \left(\sum_{i=1}^n \frac{A(y_i) e^{-\lambda A(y_i)}}{C(y_i)} \right), \\ \frac{\partial}{\partial c} l(\theta_1, \theta_2, \lambda, c) &= - \left(\sum_{i=1}^m \ln(x_i) + \sum_{i=1}^n \ln(y_i) \right) \\ &\quad + \frac{\lambda}{c} \left(\sum_{i=1}^m (A(x_i) - B(x_i)) + \sum_{i=1}^n (A(y_i) - B(y_i)) \right) \\ &\quad - (\theta_1 - 1) \frac{\lambda}{c} \sum_{i=1}^m \frac{A(x_i) + B(x_i)}{C(x_i)} e^{-\lambda A(x_i)} \\ &\quad - (\theta_2 - 1) \frac{\lambda}{c} \sum_{i=1}^n \frac{A(y_i) + B(y_i)}{C(y_i)} e^{-\lambda A(y_i)}, \end{aligned}$$

که در آن $C(x_i) = 1 - e^{-\lambda A(x_i)}$ و $B(x_i) = x_i^{-c} \ln(x_i)$ ، $A(x_i) = (x_i^{-c} - 1)/c$ در آن ماکسیمم درستنمایی $\hat{\theta}_1$ ، $\hat{\theta}_2$ ، $\hat{\lambda}$ و \hat{c} را با قرار دادن مشتق جزئی برابر با صفر و حل همزمان این معادلات می‌توان محاسبه کرد. از این معادلات بلافاصله نتیجه می‌شود که برآوردگر ماکسیمم درستنمایی θ_1 و θ_2 به صورت زیر است:

$$\begin{aligned} \hat{\theta}_1 &= - \frac{m}{\sum_{i=1}^m \ln(1 - e^{-\frac{\hat{\lambda}}{c}(x_i^{-\hat{c}}-1))}}, \\ \hat{\theta}_2 &= - \frac{n}{\sum_{i=1}^n \ln(1 - e^{-\frac{\hat{\lambda}}{c}(y_i^{-\hat{c}}-1))}}. \end{aligned}$$

شایان ذکر است که برآوردهای λ و c را می‌توان به صورت عددی بدست آورد. با استفاده از ویژگی پایایی برآوردهای ماکسیم درستنمایی، داریم:

$$\hat{R} = \frac{\hat{\theta}_2}{\hat{\theta}_1 + \hat{\theta}_2}.$$

۲.۲ روش خودگردان (Bootstrap)

دو نوع فاصله اطمینان به روش خودگردان بررسی شد: فاصله اطمینان خودگردان پایه و فاصله اطمینان خودگردان صدکی

فاصله‌های اطمینان خودگردان پایه

یک فاصله‌ی اطمینان خودگردان پایه $100(1 - \gamma)\%$ برای R به صورت زیر می باشد.

$$\left(\hat{R} - \hat{R}_{((B+1)(1-\frac{\gamma}{100}))}^*, \hat{R} - \hat{R}_{((B+1)(\frac{\gamma}{100}))}^* \right)$$

که در آن \hat{R} ، برآوردهای ماکسیم درستنمایی برای R است و $\hat{R}_{((B+1)\eta)}^*$ ، $\eta = (B+1)^{-1}$ -مین مقدار مرتب شده از برآوردهای خودگردان $\{\hat{R}_s^*, s = 1, 2, \dots, B\}$ است.

فاصله‌های اطمینان خودگردان صدکی

یک فاصله‌ی اطمینان خودگردان صدکی $100(1 - \gamma)\%$ برای R به صورت زیر است.

$$\left(\hat{R}_{((B+1)(\frac{\gamma}{100}))}^*, \hat{R}_{((B+1)(1-\frac{\gamma}{100}))}^* \right).$$

۳.۲ روش بیزی

با فرض توزیع‌های پیشین گاما برای پارامترها، توزیع پسین پارامترها طبق روابط زیر حاصل می‌گردد.

$$\begin{aligned} \pi(\theta_1, \theta_2, \lambda, c | \mathbf{x}, \mathbf{y}) &= \frac{L(\mathbf{x}, \mathbf{y} | \theta_1, \theta_2, \lambda, c) \pi(\theta_1, \theta_2, \lambda, c)}{\int_{\theta_1} \int_{\theta_2} \int_{\lambda} \int_c L(\mathbf{x}, \mathbf{y} | \theta_1, \theta_2, \lambda, c) \pi(\theta_1) \pi(\theta_2) \pi(\lambda) \pi(c) dc d\lambda d\theta_2 d\theta_1} \\ &\propto L(\mathbf{x}, \mathbf{y} | \theta_1, \theta_2, \lambda, c) \pi(\theta_1) \pi(\theta_2) \pi(\lambda) \pi(c) \\ &\propto \left(\prod_{i=1}^m x_i \right)^{-(c+1)} e^{-\frac{\lambda}{c} \sum_{i=1}^m (x_i^{-c} - 1)} \times \prod_{i=1}^m \left(1 - e^{-\frac{\lambda}{c} (x_i^{-c} - 1)} \right)^{\theta_1 - 1} \\ &\times \left(\prod_{i=1}^n y_i \right)^{-(c+1)} e^{-\frac{\lambda}{c} \sum_{i=1}^n (y_i^{-c} - 1)} \times \prod_{i=1}^n \left(1 - e^{-\frac{\lambda}{c} (y_i^{-c} - 1)} \right)^{\theta_2 - 1} \\ &\times \theta_1^{m+a_1-1} e^{-\theta_1 b_1} \times \theta_2^{n+a_2-1} e^{-\theta_2 b_2} \times \lambda^{m+n+d-1} e^{-\lambda e} \times c^{f-1} e^{-cg} \end{aligned}$$

توابع چگالی پسین $\theta_1, \theta_2, \lambda$ و c عبارتند از:

$$\begin{aligned}\pi(\theta_1 | \lambda, c, \mathbf{x}) &\propto \theta_1^{m+a_1-1} e^{-\theta_1 b_1} \prod_{i=1}^m \left[1 - e^{-\frac{\lambda}{c}(x_i^{-c}-1)} \right]^{\theta_1-1}, \\ \pi(\theta_2 | \lambda, c, \mathbf{y}) &\propto \theta_2^{n+a_2-1} e^{-\theta_2 b_2} \prod_{i=1}^n \left[1 - e^{-\frac{\lambda}{c}(y_i^{-c}-1)} \right]^{\theta_2-1}, \\ \pi(\lambda | \theta_1, \theta_2, c, \mathbf{x}, \mathbf{y}) &\propto e^{-\frac{\lambda}{c}[\sum_{i=1}^m (x_i^{-c}-1) + \sum_{i=1}^n (y_i^{-c}-1)]} \lambda^{m+n+d-1} e^{-\lambda e} \\ &\quad \times \prod_{i=1}^m \left(1 - e^{-\frac{\lambda}{c}(x_i^{-c}-1)} \right)^{\theta_1-1} \prod_{i=1}^n \left(1 - e^{-\frac{\lambda}{c}(y_i^{-c}-1)} \right)^{\theta_2-1}, \\ \pi(c | \lambda, \theta_1, \theta_2, c, \underline{x}, \underline{y}) &\propto \left(\prod_{i=1}^m x_i \prod_{i=1}^n y_i \right)^{-(c+1)} e^{-\frac{\lambda}{c}[\sum_{i=1}^m (x_i^{-c}-1) + \sum_{i=1}^n (y_i^{-c}-1)]} \\ &\quad e^{-cg} c^{f-1} \prod_{i=1}^m \left(1 - e^{-\frac{\lambda}{c}(x_i^{-c}-1)} \right)^{\theta_1-1} \prod_{i=1}^n \left(1 - e^{-\frac{\lambda}{c}(y_i^{-c}-1)} \right)^{\theta_2-1}.\end{aligned}$$

از الگوریتم گیز و متروپلیس-هستینگس برای نمونه‌گیری از توزیع‌های پسین استفاده شد و برآورد بیزی R بر اساس تابع زیان مربعات میانگین و برای تابع زیان لینکس به صورت زیر محاسبه گردید.

تابع زیان مربع خطا باشد، برآورد R عبارت است از:

$$\hat{R}_{BS} = \frac{1}{B-M} \sum_{i=M+1}^N R^{(i)}.$$

تابع زیان لینکس باشد برابر است با

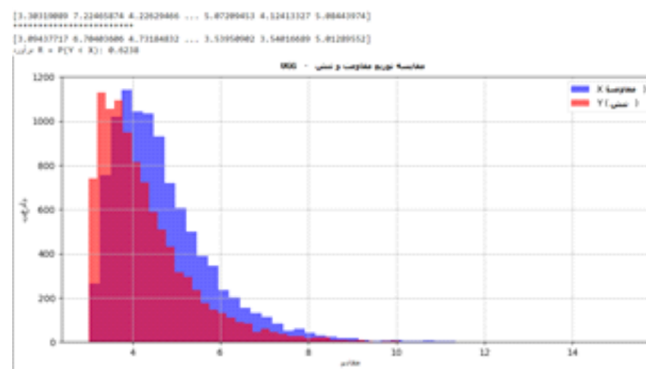
$$\hat{R}_{BL} = -\frac{1}{\phi} \log \left(\frac{1}{B-M} \sum_{i=1+M}^N e^{-\phi R^{(i)}} \right)$$

که در آن M دوره سوخت (یعنی تعداد تکرارها قبل از رسیدن به توزیع مانا) است.

۳ تحلیل شبیه‌سازی

با استفاده از داده‌های شبیه‌سازی شده برای مقادیر مختلف θ_1 و θ_2 ، مقادیر مختلفی از R تولید و با روش‌های بیان شده در بخش دوم برآورد شدند. با استفاده از برنامه پایتون تعداد ۱۰ هزار عدد تصادفی با توزیع گمپرتز تعمیم‌یافته یک‌تولد و متناسب با داده‌های مذکور مقدار قابلیت اعتماد برآورد شده است. نمودار مقایسه‌ای توزیع مشاهدات در مقاومت و تنش در نمودار زیر برای داده‌های شبیه‌سازی شده ارائه شده است.

برآورد $R = P(Y < X) : ۰/۶۲۳۸$



شکل ۱: مقایسه توزیع تنش-مقاومت در داده‌های شبه سازی شده

نتایج تحلیل مشاهدات مذکور نشان داد در برآورد ماکسیمم درستنمایی مقدار θ_1 ، $8/696$ و مقدار θ_2 ، $5/914$ و مقدار تنش مقاومت برابر $0/4048$ ($R = 0/4048$) می‌باشد.

ضمناً در روش خودگردان مقدار تنش-مقاومت برابر با $0/6348$ ($R = 0/6348$) با فاصله اطمینان ۹۵ درصدی ($0/6650$ و $0/6060$) به دست آمده است. برآورد بیزی با میانگین توزیع پسین $0/4050$ و تابع لینکس $0/7011$ که فاصله اطمینان ۹۵ درصدی ($0/4263$ و $0/3850$) حاصل شده است.

البته متناسب با پژوهش کاراندیش مروستی و همکاران (۲۰۲۴) می‌توان نتایج بارندگی دو استان کهگیلویه و بویراحمد (یاسوج) و مازندران (ساری) را از سال ۱۳۹۰ تا ۱۴۰۲ با یکدیگر مقایسه و به ترتیب به عنوان تنش-مقاومت در نظر گرفت.

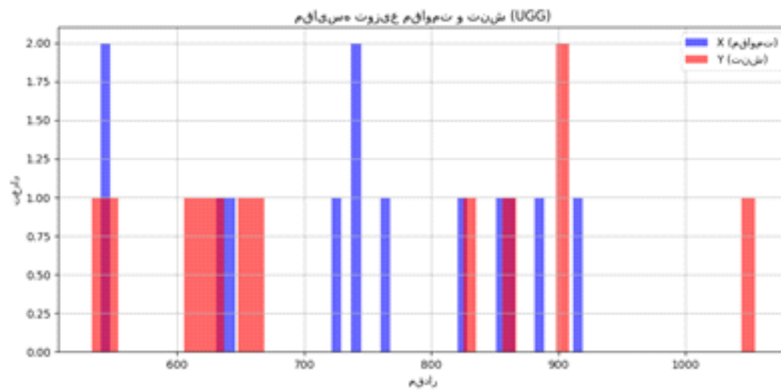
جدول ۱: بارندگی استان‌های کهگیلویه و بویراحمد (یاسوج) و مازندران (ساری) در سال‌های ۱۳۹۰ تا ۱۴۰۲

سال	سال	سال	سال	سال	سال	سال	سال	سال	سال	سال	سال	سال
۱۳۹۰	۱۳۹۱	۱۳۹۲	۱۳۹۳	۱۳۹۴	۱۳۹۵	۱۳۹۶	۱۳۹۷	۱۳۹۸	۱۳۹۹	۱۴۰۰	۱۴۰۱	۱۴۰۲
۹۱۹	۸۵۵	۷۴۱	۶۴۱	۷۲۵	۴۸۲۴	۹۵۳۹	۹۸۸۷	۷۸۶۱	۶۳۱	۴۷۴۴	۲۵۴۵	۴۷۶۶
۶۶۱	۸۲۶	۶۴۸	۶۲۴	۶۱۱	۵۵۳۲	۲۳۵۶	۷۱۰۵۴	۲۸۶۵	۲۹۰۷	۹۶۳۱	۹۹۰۱	۹۵۴۴

هر چند هر یک از استان‌ها در برخی از سال‌ها نسبت به دیگری برتری داشته است؛ برای داده‌های فوق مقدار R به عنوان تنش-مقاومت یک می‌شود.

مهم‌ترین نتایج در جداول فوق عبارتند از:

- افزایش حجم نمونه موجب کاهش خطای میانگین مربعات (MSE) می‌شود.
 - روش بیزی با تابع زیان لینکس کمترین اریبی و خطا را ارائه داد.
 - روش صدکی خودگردان و روش گیبز، فاصله‌های اطمینان پوشاننده دقیق‌تری نسبت به MLE فراهم کردند.
- تفاوت برآوردهای R در روش کلاسیک و بیزی نمایانگر حساسیت روش‌ها به انتخاب پیشین‌ها و نوع داده‌هاست. نتایج حاصل از روش‌های مختلف در ادامه به صورت خلاصه شده است.



شکل ۲: مقایسه توزیع تنش-مقاومت در داده‌های واقعی شهرهای یاسوج و ساری

جدول ۲: نتایج برآورد برای داده های واقعی

روش	برآورد θ_1 (ياسوج)	برآورد θ_2 (ساری)	مقدار تحلیلی R
روش MLE	۱/۶۴	۱/۳۱	۰/۶۳۶
روش خودگردان	میانگین R		۰/۴۸۷
	برآورد فاصله اطمینان ۹۵%	(۰/۳۴۶, ۰/۶۱۵)	
	میانگین پسین R (MSE)		۰/۷۰۰۸
	تابع زیان لینکس ($a = 2$)		۰/۷۰۰۸
	برآورد فاصله اطمینان ۹۵%	(۰/۶۹۴, ۰/۷۰۶)	

بحث و نتیجه گیری

نتایج این مطالعه نشان داد که در برآورد پارامتر قابلیت اعتماد تنش-مقاومت $R = p(Y < X)$ ، عملکرد روش‌های مختلف به‌طور معناداری به ساختار داده‌ها و روش مورد استفاده بستگی دارد. در داده‌های واقعی مورد بررسی، روش ماکسیمم درست‌نمایی (MLE) مقدار بالاتری از R (حدود ۰/۶۳۶) را نسبت به روش‌های خودگردان و بیزی ارائه داد، که می‌تواند ناشی از تأکید این روش بر ساختار مدل و نادیده‌گیری پراکندگی‌های جزئی در داده‌های محدود باشد. در مقابل، روش خودگردان (Bootstrap) مقدار میانگین پایین‌تری برای R (حدود ۰/۴۸۷) با فاصله اطمینان گسترده‌تری ارائه کرد، که نشان‌دهنده حساسیت این روش به پراکندگی‌های نمونه‌ای و قدرت پوشش خوب آن برای داده‌های کوچک است. روش بیزی (Bayesian) با استفاده از توزیع‌های پیشین نمایی ساده و الگوریتم متروپلیس، برآوردی میانی (۰/۷۰۰۸) با فاصله اعتباری متوازن ارائه داد که با نتایج مطالعات مشابه مانند تاراج و کوتز (۲۰۰۷)^۱ و ویو و همکاران (۲۰۱۴)^۲ در شرایط داده‌های کوچک و ناهمگن هم‌راستاست. این یافته‌ها بر اهمیت انعطاف‌پذیری روش‌های بیزی به‌ویژه در شرایطی با حجم داده محدود یا توزیع‌های نامتقارن تأکید می‌کند. به‌کارگیری توزیع‌های پیشین مناسب، به‌ویژه پیشین‌های آگاهانه یا بر اساس دانش زمینه‌ای، موجب کاهش اربیبی و افزایش دقت برآورد می‌شود. همچنین، تطابق نتایج این تحقیق با مطالعاتی نظیر احمدی

¹Nadarajah & Kotz²Wu et al.

و زارعی (۲۰۲۰) و کاراندیش مروستی و همکاران (۲۰۲۴) که عملکرد روش‌های بیزی را در تحلیل قابلیت اعتماد برای توزیع‌های تعمیم‌یافته تأیید می‌کنند، بر اعتبار یافته‌های حاضر می‌افزاید.

مراجع

- Ahmadi, J., and Zarei, H. (2020), Bayesian inference for the stress-strength reliability parameter based on generalized distributions. *Journal of Statistical Research of Iran (JSRI)*, **17(1)**, 45–62.
- Chen, M. H., and Shao, Q. M. (1999), Monte Carlo estimation of bayesian credible and HPD intervals. *Journal of Computational and Graphical Statistics*, **8(1)**, 69–92.
- Davison, A. C., and D. V. Hinkley. (1997), Bootstrap methods and their application. *California: Cambridge University Press*.
- Downton F. (1973), The estimation of $Pr(Y < X)$ in the normal case. *Technometrics*, **15 (3)**, 551–8. doi: 10.1080/00401706.1973.10489081.
- Harris. (1970), The estimation of reliability from stress-strength relationships. *Technometrics*, **12 (1)**, 49–54. doi: 10.1080/00401706.1970.10488633.
- Gompertz, B. (1825), On the nature of the function expressive of the law of human mortality, and on a new mode of determining the value of life contingencies. *Esq. FRS & c. Philosophical Transactions of the Royal Society of London*, **115**, 513–583. In a letter to Francis Baily, 1825.
- Jha, M. K., S. Dey, and Y. M. Tripathi. (2020), Reliability estimation in a multicomponent stress-strength based on unit-Gompertz distribution. *International Journal of Quality & Reliability Management*, **37 (3)**, 428–50. doi: 10.1108/IJQRM-04-2019-0136.
- Karandish Marvasti, A., E. Ormoz, and M. Basirat. (2022), Unit generalized Gompertz distribution and its application. *Andishe Ye Amari*, **53 (1)**, 81-92. (In Persian).
- Karandish Marvasti, A., E. Ormoz, and M. Basirat. (2024), Estimation of stress-strength reliability based on unit generalized Gompertz distribution. *Communications in Statistics - Theory and Methods*. Doi: 10.1080/03610926.2024.2428988.
- Nadarajah, S., and Kotz, S. (2007), Reliability modeling using the generalized Gompertz distribution. *IEEE Transactions on Reliability*, **56(4)**, 570–575.

- Willekens, F. (2001), Gompertz in context: the Gompertz and related distributions. *In Forecasting Mortality in Developed Countries, European studies of population*, ed.E.Tabeau,A.vandenBerg Jeths, C. Heathcote, vol. 9, 105–26. Dordrecht: Springer. doi: 10.1007/0-306-47562-6-5.
- Wu, J., Zhang, Z., and Zheng, G. (2014), Bayesian estimation of stress-strength reliability using progressive type-II censored data from the generalized distributions. *Computational Statistics & Data Analysis*, **71**, 128–139.

Comparison of Classical and Bayesian Methods in Estimation of Stress-Strength Reliability Parameter for the Unit Generalized Gompertz Distribution

Azam Karandishmarvasti¹, Ehsan Ormoz², Maryam Basirat³

¹Department of Mathematics and statistics, Ma.C., Islamic Azad University, Mashhad, Iran

²Department of Mathematics and statistics, Ma.C., Islamic Azad University, Mashhad, Iran

³Department of Mathematics and statistics, Ma.C., Islamic Azad University, Mashhad, Iran

Abstract: In this paper, the estimation of the stress-resistance reliability parameter $R = P(Y < X)$ is investigated; A case where (X resistance) and (Y tension) are independent random variables that follow the unit generalized Gompertz (UGG) distribution. In order to estimate the parameter R , three main approaches were used, including the maximum likelihood method (MLE), the bootstrap method, and the Bayesian method using the Gibbs and Metropolis-Hastings algorithms. In addition to data simulation, the present study used real rainfall data from two climatic regions of Iran (Yasuj and Sari) during the years 2011 to 2023 to measure the validity of the estimates in real conditions. The results showed that the Bayesian method, especially considering different data weighting (ratio $m = 0.3$ and $n = 0.7$), provides more stable performance and more accurate confidence intervals than other methods. The bootstrap method was also able to provide adequate coverage for estimating R in limited data. Although its average estimate was lower than that of MLE. Overall, Bayesian estimators showed higher accuracy in conditions of asymmetric or small data, which is also confirmed by comparison with previous studies. The research findings indicate that combining flexible distributions such as the generalized Gompertz unit with Bayesian approaches can provide a more reliable framework for analyzing the reliability of systems in the fields of engineering and climatology.

Keywords: Unit generalized Gompertz distribution; stress-strength; maximum likelihood estimation; bayesian estimation.

Mathematics Subject Classification (2020): 62Mxx.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



کاربرد روش استوار مونت کارلوی بهینه در استنباط بیزی بدون درست‌نمایی برای داده‌های وابسته

امید کریمی^۱، فاطمه حسینی

گروه آمار دانشگاه سمنان

چکیده: اغلب موارد در تحلیل فرایندهای سری‌زمانی و فضایی با مدل‌های آماری پیچیده‌ای مواجه می‌شویم که تابع درست‌نمایی آن‌ها براحتی قابل محاسبه نیستند. به همین دلیل، بسیاری از روش‌های استنباط بدون درست‌نمایی با محدودیت‌هایی در دقت و کارایی مواجه هستند. در این مقاله، یک روش استنباط بدون درست‌نمایی با استفاده از روش استوار مونت کارلوی بهینه را برای تحلیل داده‌های وابسته فضایی ارائه می‌کنیم. این رهیافت یک چارچوب نوین و کارآمد برای استنباط بدون درست‌نمایی است که نمونه‌های وزنی دقیقی از توزیع پسین ارائه می‌دهد. در نهایت، روش پیشنهادی روی داده‌های شبیه‌سازی شده پیاده‌سازی و ارزیابی می‌شود. **واژه‌های کلیدی:** استنباط بیزی، داده‌های وابسته، نمونه‌های مونت کارلویی، استنباط بدون درست‌نمایی. **کد موضوع‌بندی ریاضی (۲۰۲۰):** 62M05, 62H11, 62F15.

۱ مقدمه

مدل‌های مبتنی بر شبیه‌ساز به دلیل انعطاف‌پذیری بالایی که در مدل‌سازی ارائه می‌دهند، بسیار جذاب هستند. در واقع، هر مکانیزم تولید داده که بتوان آن را به صورت مجموعه‌ای محدود از گام‌های الگوریتمی تعریف کرد، می‌تواند به عنوان یک مدل مبتنی بر شبیه‌ساز برنامه‌نویسی شود. از این رو، این مدل‌ها اغلب برای شبیه‌سازی پدیده‌های فیزیکی در علوم طبیعی، مانند ژنتیک، اپیدمیولوژی، یا علوم اعصاب، به کار می‌روند (**گاتمان و کورندر، ۲۰۱۶**). در این مدل‌ها، تولید نمونه‌ها با استفاده از روش‌های مونت کارلویی امکان پذیر است، اما محاسبه تابع درست‌نمایی غیرممکن است. ناتوانی در محاسبه تابع درست‌نمایی، استنباط بدون درست‌نمایی (LFI)^۱، یعنی تقریب توزیع پسین را به چالشی بزرگ تبدیل می‌کند. روش مونت کارلوی بهینه (OMC)^۲ که توسط **میدز و ولینگ (۲۰۱۵)** پیشنهاد شد،

^۱ سخنران، omid.karimi@semnan.ac.ir

^۱ Likelihood-Free Inference

^۲ Optimization Monte Carlo

یک روش نوین برای LFI است. ایده اصلی این روش، تبدیل مکانیزم تصادفی تولید داده به مجموعه‌ای از فرآیندهای بهینه‌سازی است. گولمیس و همکاران (۲۰۱۶) برخی محدودیت‌های کلیدی OMC را شناسایی و نسخه بهبودیافته‌ای به نام مونت کارلوی بهینه استوار (ROMC)^۳ را پیشنهاد کردند. استواری در این روش به معنای کاهش اتکای صرف به اطلاعات موضعی (نقطه‌ای) در فضای پارامتر و استفاده از نواحی پذیرش برای تخمین دقیق‌تر توزیع پسین است. برخلاف روش OMC که تنها بر مقدار بهینه تکیه دارد، ROMC از مجموعه‌ای از نقاط در ناحیه پذیرش استفاده می‌کند و در نتیجه دقت بالاتری دارد. از منظر نظری، روش ROMC در چارچوب کلی فرآیندهای تصادفی تحلیل می‌شود. در واقع، ساختار شبیه‌ساز به‌عنوان نگاشتی از متغیرهای تصادفی کمکی به داده‌ها عمل می‌کند و الگوریتم ROMC با استفاده از نواحی پذیرش در فضای پارامتر، به تخمین توزیع پسین در شرایطی می‌پردازد که تابع درست‌نمایی غیرقابل محاسبه است. از این حیث، ROMC را می‌توان رویکردی مبتنی بر تحلیل هندسی نواحی با احتمال بالا در فرآیندهای تصادفی شرطی دانست.

در این مقاله، نحوه پیاده‌سازی ROMC را برای تقریب بهینه توزیع پسین روی داده‌های وابسته ارائه می‌کنیم. ROMC یک چارچوب استنباط بدون درست‌نمایی است که مجموعه‌ای از گام‌های الگوریتمی را برای تقریب توزیع پسین تعریف می‌کند، بدون آنکه الگوریتم خاصی را برای هر گام الزامی کند. بنابراین، ROMC را می‌توان به عنوان یک روش پایه استفاده کرد. در نهایت با یک مطالعه شبیه‌سازی روش پیشنهادی را ارزیابی و تحلیل می‌کنیم.

۲ مدل‌های مبتنی بر شبیه‌ساز

در این بخش، ابتدا مقدمه‌ای کوتاه درباره مدل‌های مبتنی بر شبیه‌ساز ارائه می‌دهیم، سپس به روش OMC و نسخه استوار آن ROMC می‌پردازیم. مدل مبتنی بر شبیه‌ساز، مکانیزمی تصادفی و پارامتری برای تولید داده است که امکان نمونه‌گیری از داده‌ها را فراهم می‌کند، اما محاسبه تابع درست‌نمایی در آن امکان‌پذیر نیست. در واقع یک مدل مبتنی بر شبیه‌ساز، مجموعه‌ای پارامتری از توابع چگالی احتمال $\{p(\mathbf{y} | \boldsymbol{\theta})\}_{\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^D}$ بر فضای پارامتر $\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^D$ را ارائه می‌دهد که شکل بسته ندارد یا از نظر محاسباتی غیرقابل محاسبه است. در چنین مواردی، تنها دسترسی به شبیه‌ساز $m_T(\boldsymbol{\theta})$ ممکن است، یعنی یک مکانیزم مبتنی بر الگوریتم برنامه‌نویسی که به صورت تصادفی داده‌های \mathbf{y} را از مجموعه‌ای از پارامترها $\boldsymbol{\theta}$ تولید می‌کند. فرآیند نمونه‌گیری از شبیه‌ساز را با $\mathbf{y} \rightarrow m_T(\boldsymbol{\theta})$ نشان می‌دهیم. همان‌طور که میدز و ولینگ (۲۰۱۵) نشان داده‌اند، می‌توان تصادفی بودن شبیه‌ساز را با معرفی مجموعه‌ای از متغیرهای تصادفی $\mathbf{u} \sim p(\mathbf{u})$ جدا کرد. بنابراین، برای یک جفت مشخص $(\boldsymbol{\theta}, \mathbf{u})$ ، شبیه‌ساز به یک نگاشت قطعی g تبدیل می‌شود، به طوری که $\mathbf{y} = g(\boldsymbol{\theta}, \mathbf{u})$ است. فرض کنید \mathbf{y}_0 داده‌ی مشاهده‌شده باشد، مشکل اصلی ناتوانی در محاسبه تابع درست‌نمایی $L(\boldsymbol{\theta}) = p(\mathbf{y}_0 | \boldsymbol{\theta})$ است. در واقع درست‌نمایی، احتمال تولید داده \mathbf{y} مشابه داده مشاهده‌شده \mathbf{y}_0 با استفاده از پارامترهای $\boldsymbol{\theta}$ است. بنابراین، تابع درست‌نمایی $L(\boldsymbol{\theta})$ به صورت

$$L(\boldsymbol{\theta}) = \lim_{\epsilon \rightarrow 0} c_\epsilon \int_{\mathbf{y} \in B_{d,\epsilon}(\mathbf{y}_0)} p(\mathbf{y} | \boldsymbol{\theta}) d\mathbf{y} = \lim_{\epsilon \rightarrow 0} c_\epsilon \Pr(g(\boldsymbol{\theta}, \mathbf{u}) \in B_{d,\epsilon}(\mathbf{y}_0) | \boldsymbol{\theta})$$

می‌توان نوشت. که در آن c_ϵ یک ضریب تناسب وابسته به ϵ است و $B_{d,\epsilon}(\mathbf{y}_0)$ ناحیه‌ای به شعاع ϵ اطراف \mathbf{y}_0 است که با یک تابع فاصله d به صورت $B_{d,\epsilon}(\mathbf{y}_0) := \{\mathbf{y} : d(\mathbf{y}, \mathbf{y}_0) \leq \epsilon\}$ تعریف می‌شود.

در مواردی که \mathbf{y} بعد بالا باشد از آماره خلاصه T استفاده و سپس فاصله d اعمال می‌شود. در این حالت، $B_{d,\epsilon}(\mathbf{y}_0) := \{\mathbf{y} : d(T(\mathbf{y}), T(\mathbf{y}_0)) \leq \epsilon\}$

³Robust OMC

$\{d(T(\mathbf{y}), T(\mathbf{y}_*)) \leq \epsilon\}$ خواهد شد. محاسبه $\Pr(g(\boldsymbol{\theta}, \mathbf{u}) \in B_{d,\epsilon}(\mathbf{y}_*) | \boldsymbol{\theta})$ به عنوان کسری از نمونه‌هایی که در ناحیه ϵ اطراف \mathbf{y}_* قرار دارند، در حالت حدی $\epsilon \rightarrow 0$ از نظر محاسباتی غیرممکن است. بنابراین، محدودیت به $\epsilon > 0$ کاهش می‌یابد که منجر به درست‌نمایی تقریبی $L_{d,\epsilon}(\boldsymbol{\theta}) = \Pr(\mathbf{y} \in B_{d,\epsilon}(\mathbf{y}_*) | \boldsymbol{\theta})$ می‌شود. در نتیجه، توزیع پسین تقریبی را می‌توان به صورت $p_{d,\epsilon}(\boldsymbol{\theta} | \mathbf{y}_*) \propto L_{d,\epsilon}(\boldsymbol{\theta})p(\boldsymbol{\theta})$ نوشت. این روش تنها راه حل برای مواجهه با مشکل محاسباتی تابع درست‌نمایی نیست، بلکه رویکردهای دیگر شامل مدل‌سازی رابطه (تصادفی) بین $\boldsymbol{\theta}$ و \mathbf{y} ، معکوس آن، یا تعریف استنباط بدون درست‌نمایی به عنوان یک مسئله برآورد نسبت هستند (توماس و همکاران، ۲۰۲۲؛ هرمانس و همکاران، ۲۰۲۰).

۱.۲ مونت کارلوی بهینه (OMC)

روش OMC برای اولین بار توسط میدز و ولینگ (۲۰۱۵) ارایه شده است. برای این منظور ابتدا تابع نشانگر $\mu_{B_{d,\epsilon}(\mathbf{y})}(\mathbf{x})$ را به صورت زیر نوشت:

$$\mu_{B_{d,\epsilon}(\mathbf{y})}(\mathbf{x}) = \begin{cases} 1 & \text{اگر } \mathbf{x} \in B_{d,\epsilon}(\mathbf{y}) \\ 0 & \text{در غیر این صورت} \end{cases}$$

$$L_{d,\epsilon}(\boldsymbol{\theta}) = \Pr(\mathbf{y} \in B_{d,\epsilon}(\mathbf{y}_*) | \boldsymbol{\theta}) = \int_{\mathbf{u}} \mu_{B_{d,\epsilon}(\mathbf{y}_*)}(g(\boldsymbol{\theta}, \mathbf{u})) d\mathbf{u}.$$

در آن صورت با استفاده از نمونه‌های مونت کارلویی برگرفته از شبیه‌ساز به صورت زیر تقریب زده می‌شود:

$$L_{d,\epsilon}(\boldsymbol{\theta}) \approx \frac{1}{N} \sum_{i=1}^N \mu_{B_{d,\epsilon}(\mathbf{y}_*)}(g(\boldsymbol{\theta}, \mathbf{u}_i)), \quad \mathbf{u}_i \sim p(\mathbf{u}). \quad (1.2)$$

در معادله (۱.۲)، برای هر \mathbf{u}_i ، ناحیه‌ای به صورت $C_\epsilon^i = \{\boldsymbol{\theta} : g(\boldsymbol{\theta}, \mathbf{u}_i) \in B_{d,\epsilon}(\mathbf{y}_*)\}$ در فضای پارامتر $\boldsymbol{\theta}$ وجود دارد که تابع نشانگر مقدار یک را برمی‌گرداند. بنابراین، می‌توان درست‌نمایی تقریبی و توزیع پسین را به صورت زیر نوشت:

$$L_{d,\epsilon}(\boldsymbol{\theta}) \approx \frac{1}{N} \sum_{i=1}^N \mu_{C_\epsilon^i}(\boldsymbol{\theta}), \quad p_{d,\epsilon}(\boldsymbol{\theta} | \mathbf{y}_*) \propto p(\boldsymbol{\theta}) \sum_{i=1}^N \mu_{C_\epsilon^i}(\boldsymbol{\theta}). \quad (2.2)$$

در واقع، تقریب درست‌نمایی و توزیع پسین به توصیف مجموعه‌های C_ϵ^i تبدیل می‌شود. در روش OMC، فرض می‌کند که فاصله d ، فاصله اقلیدسی $\|\cdot\|_2$ بین آمار خلاصه T داده‌های مشاهده‌شده و تولیدشده است و نواحی C_ϵ^i می‌توانند به خوبی با بیضی‌های بسیار کوچک تقریبی شوند. این فرضیات منجر به تقریب توزیع پسین با استفاده از نمونه‌های وزنی $\boldsymbol{\theta}_i^*$ می‌شوند که کمترین فاصله بین داده‌های مشاهده‌شده و شبیه‌سازی‌شده را برای هر تحقق $\mathbf{u}_i \sim p(\mathbf{u})$ به صورت $\boldsymbol{\theta}_i^* = \argmin_{\boldsymbol{\theta}} \|\Phi(\mathbf{y}_*) - \Phi(g(\boldsymbol{\theta}, \mathbf{u}_i))\|_2$ به دست می‌آورند.

۲.۲ روش استوار مونت کارلوی بهینه (ROMC)

گولمیس و همکاران (۲۰۱۶) نشان دادند که فرض بیضی‌های بسیار کوچک می‌تواند به توزیع‌های پسین بیش از حد مطمئن منجر شود. به طور شهودی، این مشکل به این دلیل رخ می‌دهد که وزن‌ها در OMC تنها بر اساس اطلاعات محلی در $\boldsymbol{\theta}_i^*$ محاسبه می‌شوند و استفاده صرف از اطلاعات محلی می‌تواند گمراه‌کننده باشد. برای مثال، اگر انحنای $\|T(\mathbf{y}_*) - T(g(\boldsymbol{\theta}, \mathbf{u}_i))\|_2$ در $\boldsymbol{\theta}_i^*$ تقریباً صاف باشد، ممکن است به اشتباه نشان دهد که C_ϵ^i بسیار بزرگ‌تر از اندازه واقعی آن است. روش استوار OMC پیشنهادی این مشکل را برطرف می‌کند. ROMC نواحی پذیرش C_ϵ^i را با جعبه‌های مرزی D -بعدی \hat{C}_ϵ^i تقریبی می‌کند. یک توزیع یکنواخت $q_i(\boldsymbol{\theta})$ روی هر

جعبه مرزی تعریف می‌شود و به عنوان توزیع پیشنهادی برای تولید نمونه‌های پسین $\theta_{ij} \sim q_i$ عمل می‌کند. نمونه‌ها وزن به صورت $w_{ij} = \mu_{C_\epsilon^i}(\theta_{ij}) \frac{p(\theta_{ij})}{q_i(\theta_{ij})}$ دریافت می‌کنند که استفاده از توزیع‌های پیشنهادی $q_i(\theta)$ به جای پیشین $p(\theta)$ را تعدیل می‌کند. با توجه به نمونه‌های وزنی، $\mathbb{E}_{p(\theta|y_\circ)}[h(\theta)]$ از یک تابع $h(\theta)$ می‌تواند به صورت $\mathbb{E}_{p(\theta|y_\circ)}[h(\theta)] \approx \frac{\sum_{ij} w_{ij} h(\theta_{ij})}{\sum_{ij} w_{ij}}$ تقریبی شود. برای محاسبه وزن w_{ij} باید بررسی کنیم که آیا نمونه‌های θ_{ij} ، که از جعبه‌های مرزی گرفته شده‌اند، در ناحیه پذیرش C_ϵ^i قرار دارند. این بررسی به عنوان یک مکانیزم ایمنی عمل می‌کند که نادرستی‌های احتمالی در ساخت \hat{C}_ϵ^i را اصلاح می‌کند. اما این بررسی نیازمند ارزیابی تابع فاصله $d_i(\theta_{ij})$ است که در مدل‌های پیچیده ممکن است پرهزینه باشد. به طور کلی، ROMC به دو بخش آموزش و استنباط تقسیم می‌شود. بخش آموزش شامل گام‌هایی برای تخمین نواحی پیشنهادی است و بخش استنباط نمونه‌های وزنی را محاسبه می‌کند. در واقع ابتدا، نواحی پیشنهادی \hat{C}_ϵ^i تعیین و برآورد می‌شوند و سپس توزیع پسین تقریب زده می‌شود. در ادامه یک مطالعه شبیه‌سازی برای پیاده‌سازی روش پیشنهادی روی داده‌های وابسته ارائه می‌شود.

۳ مطالعه شبیه‌سازی

در این بخش، روش پیشنهادی را ابتدا روی مدل میانگین متحرک مرتبه دوم (MA2)، که یکی از مدل‌های فرایندهای سری زمانی است، پیاده‌سازی و ارزیابی می‌کنیم. سپس روی یک مدل رگرسیون فضایی ساده اجرا و تحلیل می‌شود. استنباط با سه رویکرد مختلف روش ROMC انجام می‌شود: (۱) با بهینه‌ساز مبتنی بر گرادین، (۲) با بهینه‌سازی بیزی و (۳) با برازش یک شبکه عصبی به عنوان مدل جایگزین. مورد آخر نشان‌دهنده قابلیت توسعه روش پیشنهادی است، که در آن بخشی از ROMC با یک روش مشخص توسط کاربر جایگزین می‌شود.

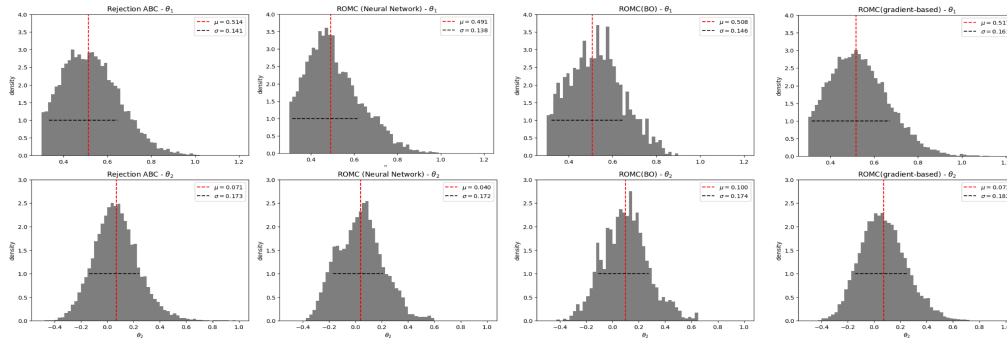
۱.۳ مدل میانگین متحرک مرتبه دوم

مدل MA2 به صورت $y_t = w_t + \theta_1 w_{t-1} + \theta_2 w_{t-2}$ ، $t = 1, \dots, K$ تعریف می‌شود، که در آن $\theta_1, \theta_2 \in \mathbb{R}$ و $w_t \sim \mathcal{N}(0, 1)$ است. متغیر تصادفی w_k نوفه سفید است و دو پارامتر مورد نظر، θ_1 و θ_2 ، وابستگی به مشاهدات قبلی را مدل می‌کنند. پارامتر K تعداد مشاهدات متوالی است که $K = 100$ در نظر گرفته شده است. برای اطمینان از قابل‌شناسایی بودن از پیشین پیشنهادی مارین و همکاران (۲۰۱۲) به صورت $p(\theta) = p(\theta_1)p(\theta_2 | \theta_1) = \mathcal{U}(\theta_1; -2, 2)\mathcal{U}(\theta_2; \theta_1 - 1, \theta_1 + 1)$ استفاده می‌کنیم. بردار مشاهده $y_\circ = (y_1, \dots, y_{100})$ با $y_\circ = (0.6, 0.2)$ تولید شده است. به دلیل بالا بودن بُعد خروجی y ، از آماره خلاصه استفاده می‌کنیم. با توجه به اینکه بردار خروجی یک تحقق سری‌زمانی را نشان می‌دهد، از اتوکواریانس‌ها با تأخیر ۱ و ۲ به صورت زیر استفاده شده است:

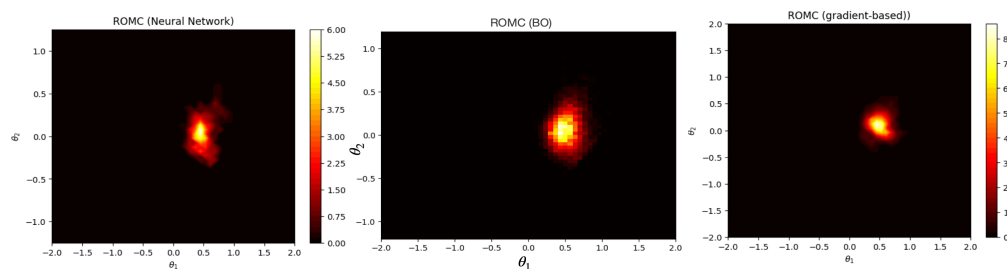
$$s_1(y) = \frac{1}{T-1} \sum_{t=2}^T y_t y_{t-1}, \quad s_2(y) = \frac{1}{T-2} \sum_{t=3}^T y_t y_{t-2},$$

$$s(y) = (s_1(y), s_2(y)), \quad d = \|s(y) - s(y_\circ)\|_2^2, \quad (1.3)$$

فاصله بین مشاهده و خروجی شبیه‌ساز با فاصله اقلیدسی، طبق معادله ۱.۳، اندازه‌گیری می‌شود. همچنین برای شبکه عصبی، از کلاس MLPRegressor از بسته scikit-learn استفاده می‌کنیم. که در آن جایگزین تابع فاصله واقعی d_i در نواحی پیشنهادی می‌شود. بنابراین، تمام اقدامات استنباط، یعنی نمونه‌گیری، محاسبه امید، و ارزیابی پسین، بر اساس \hat{d}_i انجام می‌شود. در این جا از یک



شکل ۱: هیستوگرام‌های توزیع‌های پسین حاشیه‌ای برای هر روش



شکل ۲: توزیع‌های پسین حاشیه‌ای پارامترهای مدل فضایی برای سه نسخه ROMC

شبکه عصبی با دو لایه مخفی، هر یک با ۱۰ نورون، استفاده می‌کنیم و آن را با ۵۰۰ نمونه از هر ناحیه پیشنهادی آموزش می‌دهیم. برای مقایسه نتایج استنباط ROMC با الگوریتم سنتی بیز تقریبی ABC، نیز در تحلیل آورده شده است. در جدول ۱، میانگین و انحراف معیار نمونه‌های پسینی پارامترها برای هر روش ارائه شده است. مشاهده می‌شود که نمونه‌های به‌دست‌آمده از توزیع پسین تقریبی در تمام روش‌ها معادل و نزدیک به مقدار واقعی هستند. در شکل ۱، هیستوگرام‌های توزیع‌های پسین حاشیه‌ای پارامترها برای همه روش‌ها

جدول ۱: مقایسه نتایج پسین برای پارامترهای مدل $MA(2)$ با استفاده از نسخه‌های مختلف روش ROMC و الگوریتم ABC

	$\hat{\mu}_{\theta_1}$	$\hat{\sigma}_{\theta_1}$	$\hat{\mu}_{\theta_2}$	$\hat{\sigma}_{\theta_2}$
نمونه‌گیری رد ABC	۰/۵۱۶	۰/۱۴۲	۰/۰۷۰	۰/۱۷۲
ROMC (گرادیانی)	۰/۵۰۱	۰/۱۴۲	۰/۰۳۳	۰/۱۶۹
ROMC (بیزی)	۰/۵۱۳	۰/۱۶۹	۰/۰۹۰	۰/۱۷۴
ROMC (شبکه عصبی)	۰/۴۹۱	۰/۱۳۸	۰/۰۴۰	۰/۱۷۲

رسم شده‌اند که نشان می‌دهد توزیع پارامترها در همه روش‌ها نسبتاً مشابه به هم هستند. در شکل ۲، توزیع پسین برای سه نسخه مختلف روش ROMC ارائه شده است. نتایج نشان می‌دهند که تمام نسخه‌های ROMC نتایج نسبتاً مناسب و یکسانی ارائه می‌دهند که با الگوریتم نمونه‌گیری رد ABC هم‌خوانی دارند. با موازی کردن این روش می‌توان سرعت محاسبات را تا ۵ برابر افزایش داد.

۲.۳ مدل رگرسیون فضایی

در این بخش، روش پیشنهادی را روی مدل رگرسیون فضایی در یک شبکه منظم 10×10 با مختصات $s_i = (i, j)$ و $i, j = 1, \dots, 10$ پیاده‌سازی و ارزیابی می‌کنیم. مشابه بخش قبل استنباط با سه رویکرد مختلف روش ROMC انجام می‌شود. مدل

رگرسیون فضایی به صورت $y(s_i) = \beta_0 + \beta_1 x(s_i) + w(s_i) + \varepsilon_i$ که در آن $i = 1, \dots, 100$ یک متغیر توضیحی در مکان s_i است که از توزیع یکنواخت $\mathcal{U}(0, 1)$ تولید شده است، $\beta = (\beta_0, \beta_1)$ بردار ضرایب رگرسیون، $w(s_i)$ فرآیند گاوسی با کوواریانس نمایی $\text{Cov}(w(s_i), w(s_j)) = \sigma^2 \exp(-\phi \|s_i - s_j\|)$ با $\sigma^2 = 1$ و $\varepsilon_i \sim \mathcal{N}(0, \tau^2)$ نوفه سفید است. فاصله $\|s_i - s_j\|$ فاصله اقلیدسی بین مکان‌های s_i و s_j روی شبکه منظم است. پارامترهای مورد نظر برای استنباط $\theta = (\beta_0, \beta_1, \phi, \tau^2)$ هستند. از پیشین‌های معمول زیر استفاده می‌کنیم:

$$p(\theta) = p(\beta_0)p(\beta_1)p(\phi)p(\tau^2) = \mathcal{N}(\beta_0; 0, 10)\mathcal{N}(\beta_1; 0, 10)\mathcal{U}(\phi; 0.1, 5)\mathcal{IG}(\tau^2; 2, 1). \quad (2.3)$$

که در آن \mathcal{IG} توزیع گامای معکوس است. بردار مشاهده $\mathbf{y}_0 = (y(s_1), \dots, y(s_{100}))$ با $\mathbf{y}_0 = (1, 0.5, 1, 0.1)$ θ^* تولید شده است. انتخاب آماره‌های خلاصه $T(y)$ برای داده‌های فضایی نقشی کلیدی در عملکرد الگوریتم ROMC ایفا می‌کند. در این مطالعه، از میانگین کلی مشاهدات و واریوگرام تجربی به صورت زیر استفاده شده است تا به ترتیب ویژگی‌های سطح کلی و ساختار همبستگی فضایی داده‌ها منعکس گردد. این آماره‌ها نقش تعیین‌کننده‌ای در تعریف تابع فاصله و شکل‌گیری نواحی پذیرش دارند:

$$T_1(\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n y(s_i), \quad T_2(\mathbf{y}) = \frac{1}{|\mathcal{N}|} \sum_{(i,j) \in \mathcal{N}} (y(s_i) - y(s_j))^2, \\ T(\mathbf{y}) = (T_1(\mathbf{y}), T_2(\mathbf{y})), \quad d = \|T(\mathbf{y}) - T(\mathbf{y}_0)\|_2^2,$$

که در آن \mathcal{N} مجموعه جفت‌های مکان‌هایی است که فاصله اقلیدسی آن‌ها کمتر از ۲ واحد است. در واقع از همسایگی‌های نزدیک مشابه استفاده کرده‌ایم. در این‌جا از یک شبکه عصبی با دو لایه مخفی، هر یک با ۱۲ نورون، استفاده می‌کنیم و آن را با ۵۰۰ نمونه از هر ناحیه پیشنهادی آموزش می‌دهیم. برای مقایسه نتایج استنباط ROMC با الگوریتم سنتی نمونه‌گیری رد ABC، نیز در این تحلیل آورده شده است. در جدول ۲، میانگین و انحراف معیار نمونه‌های پسینی پارامترهای β_0, β_1, ϕ و τ^2 برای همه روش‌ها ارائه شده است. مشاهده می‌شود که نمونه‌های به دست آمده از توزیع پسین تقریبی در تمام روش‌ها معادل و تقریباً به مقدار واقعی نزدیک هستند. نتایج

جدول ۲: خلاصه نتایج توزیع‌های پسین حاشیه‌ای پارامترهای مدل فضایی برای سه نسخه ROMC و نمونه‌گیری رد ABC

$\hat{\sigma}_{\tau^2}$	$\hat{\mu}_{\tau^2}$	$\hat{\sigma}_{\phi}$	$\hat{\mu}_{\phi}$	$\hat{\sigma}_{\beta_1}$	$\hat{\mu}_{\beta_1}$	$\hat{\sigma}_{\beta_0}$	$\hat{\mu}_{\beta_0}$	
۰/۰۲۵	۰/۰۹۸	۰/۲۱۲	۱/۰۱۰	۰/۰۸۵	۰/۴۹۵	۰/۱۴۸	۰/۹۹۲	نمونه‌گیری رد ABC
۰/۰۲۴	۰/۰۹۹	۰/۲۰۸	۱/۰۰۷	۰/۰۸۳	۰/۴۹۹	۰/۱۴۵	۰/۹۹۸	ROMC (مبتنی بر گردایان)
۰/۰۲۶	۰/۰۹۷	۰/۲۱۵	۱/۰۱۳	۰/۰۸۷	۰/۴۹۴	۰/۱۵۰	۰/۹۸۹	ROMC (بهینه‌سازی بیزی)
۰/۰۲۳	۰/۰۹۶	۰/۲۱۰	۱/۰۰۵	۰/۰۸۲	۰/۴۹۰	۰/۱۴۳	۰/۹۸۷	ROMC (شبکه عصبی)

نشان می‌دهند که تمام نسخه‌های ROMC نتایج منسجمی ارائه می‌دهند که با الگوریتم نمونه‌گیری رد ABC هم‌خوانی دارند. با وجود قابلیت بالای موازی‌سازی در روش ROMC، این روش همچنان دارای محدودیت‌هایی در کاربردهای واقعی است. حل تعداد زیادی مسئله بهینه‌سازی برای تعیین نقاط θ^* و ساخت نواحی پذیرش، زمان‌بر بوده و نیازمند منابع محاسباتی قابل توجهی است، به‌ویژه در مدل‌هایی با ابعاد بالا یا شبیه‌سازهای پیچیده. همچنین در محیط‌هایی با محدودیت در دسترسی به پردازنده‌های چند هسته‌ای، بهره‌گیری کامل از ظرفیت موازی‌سازی با چالش‌هایی مواجه خواهد بود.

۴ بحث و نتیجه‌گیری

در این مقاله، روش استنباط بدون درست‌نمایی ROMC برای تحلیل داده‌های وابسته، از جمله سری‌های زمانی و داده‌های فضایی، ارائه شد. این روش با استفاده از سه رویکرد مختلف پیاده‌سازی و ارزیابی شد. نتایج مطالعه شبیه‌سازی روی مدل میانگین متحرک مرتبه دوم و مدل رگرسیون فضایی نشان داد که هر سه نسخه ROMC نمونه‌های پسینی با میانگین و انحراف معیار نزدیک به مقادیر واقعی تولید می‌کنند و با الگوریتم نمونه‌گیری رد ABC نتایج مشابهی دارند. همچنین، موازی‌سازی روش ROMC می‌تواند سرعت محاسبات را تا ۵ برابر افزایش دهد، که نشان‌دهنده کارایی بالای این روش در مسائل پیچیده است.

جنبه اصلی این روش این است که ROMC یک چارچوب انعطاف‌پذیر و قابل گسترش ارائه می‌دهد که به کاربران اجازه می‌دهد تا بخش‌های مختلف روش، مانند تابع فاصله یا الگوریتم بهینه‌سازی، را با رویکردهای دلخواه جایگزین کنند و این امکان استفاده گسترده‌تر و توسعه بیشتر توسط پژوهشگران را فراهم می‌کند. با این حال، چالش‌هایی برای بهبود کارایی ROMC، به‌ویژه در مسائل با ابعاد بالا، همچنان وجود دارد. مهم‌ترین چالش، امکان اجرای ROMC در محیط‌های توزیع‌شده، مانند خوشه‌های محاسباتی، است که می‌تواند زمان محاسبات را به‌طور قابل توجهی کاهش دهد. علاوه بر این، تنظیم تعداد نوروها و لایه‌های شبکه عصبی برای مسائل پیچیده‌تر یا بهینه‌سازی انتخاب آماره‌های خلاصه برای داده‌های فضایی می‌تواند دقت تقریب توزیع پسین را بهبود بخشد. در آینده، تمرکز بر توسعه الگوریتم‌های موازی و بهینه‌سازی‌های عددی می‌تواند قابلیت‌های ROMC را برای کاربردهای گسترده‌تر تقویت کند.

مراجع

- Gutmann, M. U. and Corander, J. (2016), *Bayesian Optimization for Likelihood-Free Inference of Simulator-Based Statistical Models*, *Journal of Machine Learning Research*, **17**, 1–47
- Gkolemis, V., Gutmann, M. U. and Pesonen, H. (2024), *An Extendable Python Implementation of Robust Optimization Monte Carlo*, *Journal of Statistical Software*, **110**, 1–26.
- Hermans, J., Begy, V. and Louppe, G. (2020), *Likelihood-Free MCMC with Amortized Approximate Ratio Estimators*, *International Conference on Machine Learning*, 4239–4248, PMLR.
- Marin, J. M., Pudlo, P., Robert, C. P. and Ryder, R. J. (2012), *Approximate Bayesian Computational Methods*, *Statistics and Computing*, **22**, 1167–1180.
- Meeds, T. and Welling, M. (2015), *Optimization Monte Carlo: Efficient and Embarrassingly Parallel Likelihood-Free Inference*, *Advances in Neural Information Processing Systems*, **28**.
- Thomas, O., Dutta, R., Corander, J., Kaski, S. and Gutmann, M. U. (2022), *Likelihood-Free Inference by Ratio Estimation*, *Bayesian Analysis*, **17**, 1–31.

Robust Optimal Monte Carlo Method

Omid Karimi, Fatemeh Hosseini

Department of Statistics, Semnan University, Semnan, Iran

Abstract: In most cases in the analysis of time-series and spatial processes, we encounter complex statistical models whose likelihood functions are not easily computable. For this reason, many likelihood-free inference methods face limitations in terms of accuracy and efficiency. In this paper, we present a likelihood-free inference method using an optimal robust Monte Carlo approach for analyzing spatially dependent data. This approach provides a novel and efficient framework for likelihood-free inference, offering accurate weighted samples from the posterior distribution. Finally, the proposed method is implemented and evaluated on simulated data.

Keywords: Bayesian inference, dependent data, Monte Carlo methods, likelihood-free inference.

Mathematics Subject Classification (2020): 65C05, 90C26, 62H11.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30-31, 2025
University of Kurdistan



مفصل‌های دایره‌ای-خطی: ویژگی‌ها و کاربردها

نجیب‌الله کریمی^۱، علی دست‌برآورده^۲

^۱ دانشجوی دکتری آمار، دانشگاه یزد

^۲ عضو هیأت علمی گروه آمار، دانشگاه یزد

چکیده: این مقاله مطالعه مروری از تحقیقات جدید در زمینه مفصل‌های دایره‌ای-خطی است. برای تحلیل ساختار وابستگی داده‌های دایره‌ای و خطی در حالت‌های مختلف، چند مفصل دایره‌ای-خطی نیز مورد مطالعه قرار گرفته‌اند. علاوه بر این، به معرفی و بکارگیری بسته cylcop در نرم‌افزار R نیز پرداخته شده است. در نهایت، از طریق شبیه‌سازی، مفصل‌های ارائه شده مورد بررسی قرار گرفته و پارامترهای آنها برآورد شده‌اند. نتایج نشان می‌دهند که مفصل‌های با بخش‌های ویژه برای داده‌های حرکتی مناسب هستند.

واژه‌های کلیدی: داده‌های حرکتی، داده‌های دایره‌ای-خطی، مفصل با بخش‌های ویژه، مفصل دایره‌ای-خطی.

کد موضوع‌بندی ریاضی (۲۰۲۰): 62F99, 62G30, 62H11.

۱ مقدمه

در بسیاری از حوزه‌های علمی مانند هواشناسی، زیست‌شناسی، علوم زمین و روان‌شناسی با داده‌هایی مواجه هستیم که ماهیتی زاویه‌ای دارند. این داده‌ها که به داده‌های دایره‌ای معروف هستند، بر روی محیط یک دایره واحد تعریف می‌شوند و ویژگی اصلی آن‌ها تناوب و برابری در انتها و ابتدای مقیاس اندازه‌گیری است. نمونه‌هایی از داده‌های دایره‌ای عبارت‌اند از: جهت باد، چرخه‌های زمان در طی شبانه‌روز، زاویه‌های فاز در پردازش سیگنال و جهت حرکت حیوانات. برای مثال، زوایای 0° و 360° درجه در یک فضای دایره‌ای معادل هستند. به همین دلیل، روش‌های آماری متداول که بر پایه مفروضات داده‌های خطی بنا شده‌اند، برای تحلیل این نوع داده‌ها مناسب نیستند (ماردیا و جاپ، ۲۰۰۰).

در بسیاری از مطالعات، متغیر دایره‌ای با یک یا چند متغیر خطی همراه است. برای مثال، در بررسی رفتار حیوانات، ممکن است جهت حرکت (دایره‌ای) با سرعت حرکت (خطی) ترکیب شوند یا در تحلیل باد، سرعت و جهت به‌طور توأم بررسی شوند. به داده‌هایی که

از رفتار حیوانات یا وزیدن باد به‌دست می‌آیند، داده‌های حرکتی نیز می‌گویند. در رفتار حیوانات، داده‌های مربوط به جهت حرکت را معمولاً تحت عنوان زاویه‌های چرخش و داده‌های مربوط به سرعت حرکت را تحت عنوان طول‌های گام بررسی می‌کنند (هودل و فبرینگ، ۲۰۲۱).

در چنین مواردی، درک دقیق وابستگی بین متغیرهای دایره‌ای و خطی برای تفسیر درست پدیده‌ها ضروری است. یکی از روش‌های منعطف و قدرتمند برای مدل‌سازی وابستگی بین متغیرها، استفاده از مفصل‌ها است. مفصل‌ها این امکان را فراهم می‌کنند که ساختار وابستگی میان متغیرها جدا از توزیع‌های حاشیه‌ای مدل‌سازی شود (نلسون، ۲۰۰۶). با این وجود، توسعه مفصل‌هایی که به‌طور خاص برای مدل‌سازی وابستگی بین متغیرهای دایره‌ای و خطی طراحی شده باشند، نسبتاً جدید و در حال گسترش است. این مفصل‌ها، قادرند با حفظ ماهیت تناوبی متغیرهای دایره‌ای و خواص متغیر خطی، ساختار وابستگی را به دقت مدل کنند (دورنت و سمپی، ۲۰۱۵). از سوی دیگر، پژوهش‌هایی مانند نگار و همکاران (۲۰۲۳) با تمرکز بر کاربردهای عملی مدل‌سازی وابستگی دایره‌ای-خطی در داده‌های زیستی و پزشکی، مسیرهای جدیدی را در این داده‌ها گشوده‌اند.

این مقاله، علاوه بر مقدمه، سه بخش دیگر دارد. در بخش دوم به مفاهیم پایه مرتبط با متغیرهای دایره‌ای و مفصل‌های خطی پرداخته شده و محدودیت‌های این مفصل‌ها در مدل‌سازی وابستگی بین داده‌های دایره‌ای-خطی، بیان شده است. در بخش سوم، مفصل‌های دایره‌ای-خطی معرفی شده و ویژگی‌های آن‌ها مورد بررسی قرار گرفته است. همچنین، به کاربرد و بکارگیری بسته cylcop نیز اشاره شده است. در بخش چهارم، از طریق شبیه‌سازی، کاربردهای مفصل‌های دایره‌ای-خطی در مدل‌سازی وابستگی بین متغیرهای دایره‌ای و خطی، بیان شده است.

۲ مفاهیم پایه

۱.۲ متغیرهای دایره‌ای

متغیرهای دایره‌ای به متغیرهای اطلاق می‌شوند که مقادیر آنها در فضای دایره‌ای تعریف شده و خاصیت تناوبی دارند. این متغیرها که معمولاً با $\theta \in [0, 2\pi)$ نمایش داده می‌شوند، در تحلیل‌های جهت‌دار و زاویه‌ای کاربرد فراوانی دارند. از جمله مهم‌ترین کاربردهای عملی می‌توان به جهت باد در هواشناسی، زوایای حرکت در ناوبری و زمان وقوع پدیده‌های تناوبی در زیست‌شناسی اشاره کرد. مدل‌سازی توزیع‌های احتمالی متغیرهای دایره‌ای موسوم به آمار جهتی است (ماریا و جاپ، ۲۰۰۰). تابع چگالی احتمال یک متغیر دایره‌ای، Θ ، تابعی با دوره تناوب 2π است که در رابطه زیر صدق می‌کند:

$$f_{\Theta}(\theta) = f_{\Theta}(\theta + 2\pi), \quad \theta \in [0, 2\pi), \quad k \in \mathbb{Z}. \quad (1.2)$$

مدل‌های توزیع خطی ممکن است برای متغیرهای دایره‌ای نامعتبر باشند. بنابراین، استفاده از مدل‌های توزیع خاص برای متغیرهای دایره‌ای ضروری است. به‌عنوان مثال می‌توان به توزیع فون‌میزس، توزیع کوشی پوشانده‌شده و توزیع نرمال پوشانده‌شده اشاره کرد. در اینجا فقط به توزیع فون‌میزس که به عنوان معادل دایره‌ای توزیع نرمال شناخته می‌شود، اشاره می‌کنیم. تابع چگالی این توزیع به‌صورت زیر تعریف می‌شود:

$$f(\theta; \mu, k) = \frac{e^{k \cos(\theta - \mu)}}{2\pi I_0(k)}, \quad (2.2)$$

که در آن $\mu \in [0, 2\pi)$ پارامتر مکان (میانگین جهت)، $k \geq 0$ پارامتر تمرکز و $I_0(k)$ تابع بسل اصلاح‌شده از نوع اول و مرتبه صفر است (ابرامویتز و استیگون، ۱۹۶۵). این توزیع در مدل‌سازی بسیاری از پدیده‌های طبیعی با ماهیت تناوبی کاربرد گسترده‌ای دارد.

۲.۲ مفصل و ساختار وابستگی

نظریه مفصل بر پایه قضیه اسکالر (اسکلر، ۱۹۵۹) بنا شده است و روش منعطف‌تر و مناسب‌تر برای ساخت توزیع توأم چندمتغیره (که در اینجا ما دومتغیره آن را بررسی می‌کنیم) فراهم می‌آورد. قضیه اسکالر بیان می‌کند که تابع توزیع توأم چندمتغیره را می‌توان به صورت ترکیبی از یک تابع مفصل که ساختار وابستگی بین متغیرها را توصیف می‌کند و توابع توزیع تجمعی حاشیه‌ای برای تمام متغیرهای حاشیه‌ای نوشت. در مفصل‌ها، ساختار وابستگی میان متغیرها جدا از توزیع هر متغیر در نظر گرفته می‌شود. بنابراین، مدل مناسب برای توزیع حاشیه‌ای هر متغیر را می‌توان به طور مستقل و آزادانه انتخاب کرد. فرض کنید $F_{X_1, X_2}(x_1, x_2)$ و $f_{X_1, X_2}(x_1, x_2)$ به ترتیب تابع توزیع تجمعی توأم و تابع چگالی احتمال توأم برای متغیرهای تصادفی (X_1, X_2) باشند. در این صورت، برای هر $(x_1, x_2) \in \mathbb{R}^2$ ، با استفاده از یک تابع مفصل C داریم (وانگ و همکاران، ۲۰۲۱):

$$F_{X_1, X_2}(x_1, x_2) = C(u_1, u_2), \quad (3.2)$$

$$f_{X_1, X_2}(x_1, x_2) = c(u_1, u_2) f_{X_1}(x_1) f_{X_2}(x_2), \quad (4.2)$$

که در آن $u_i = F_{X_i}(x_i)$ ، $i = 1, 2$ ، مقدار تابع توزیع حاشیه‌ای X_i در نقطه x_i ، $f_{X_i}(x_i)$ مقدار تابع چگالی متناظر با آن و c تابع چگالی مفصل C هستند. تابع چگالی c به صورت زیر محاسبه می‌گردد:

$$c(u_1, u_2) = \frac{\partial^2 C(u_1, u_2)}{\partial u_1 \partial u_2}. \quad (5.2)$$

تعریف ۱.۲. مفصل تجربی: با استفاده از یک نمونه شامل n مشاهدی دوبعدی مستقل و هم‌توزیع به صورت (x_i, y_i) که از یک توزیع توأم استخراج شده‌اند، می‌توان مفصل تجربی را به صورت زیر محاسبه کرد (جینیست و فوری، ۲۰۰۷):

$$C_n(u, v) = \frac{1}{n} \sum_{i=1}^n I\left(\frac{r_i}{n+1} \leq u, \frac{s_i}{n+1} \leq v\right), \quad (6.2)$$

که در آن $I(\cdot)$ تابع نشانگر، r_i و s_i به ترتیب رتبه مشاهدات x_i و y_i هستند. هنگامی که یکی از متغیرهای تصادفی ماهیت دایره‌ای داشته باشد، می‌توان همان رویه را دنبال کرد، اما برای تعیین رتبه‌های متغیر دایره‌ای، نیاز به تعریف یک مرجع خواهیم داشت. در این مقاله، مقدار مرجع را $-\pi$ در نظر می‌گیریم.

۳.۲ محدودیت‌های مفصل‌های کلاسیک برای داده‌های دایره‌ای-خطی

هنگامی که یکی از متغیرها دایره‌ای و دیگری خطی باشد، مفصل‌های کلاسیک با چالش‌های اساسی مواجه می‌شوند. مشکل اصلی ناشی از ناسازگاری ذاتی بین ساختار تناوبی متغیرهای دایره‌ای و مفروضات اساسی مفصل‌های کلاسیک است (جانسون و وهرلی، ۱۹۷۸). این ناسازگاری می‌تواند منجر به برآوردهای نادرست و تفسیرهای اشتباه از روابط بین متغیرها شود. علاوه بر این، پیچیدگی‌های محاسباتی ناشی از نیاز به تبدیل‌های خاص برای متغیرهای دایره‌ای، استفاده از روش‌های استاندارد را با دشواری مواجه می‌سازد. در مفصل‌های کلاسیک، فرض بر این است که داده‌ها از یک ترتیب ویژه پیروی می‌کنند، در حالی که داده‌های دایره‌ای دارای ماهیت تناوبی‌اند و آغاز و

پایان مشخص ندارند. همچنین، استفاده از توزیع‌های حاشیه‌ای نامناسب برای داده‌های دایره‌ای در چارچوب مفصل‌های کلاسیک می‌تواند به نتایج آماری گمراه‌کننده منجر شود. برای رفع این چالش‌ها، مفصل‌های سازگار با داده‌های دایره‌ای-خطی ضروری به نظر می‌رسد.

۳ مفصل دایره‌ای-خطی

۱.۳ معرفی و ویژگی‌ها

مفصل‌های دایره‌ای-خطی از دهه ۱۹۷۰ در ادبیات آماری از جمله در آثار جانسون و وهرلی (۱۹۷۷، ۱۹۷۸) مطرح بوده‌اند. هرچند در آن زمان هنوز واژه «مفصل» به‌کار نمی‌رفت. از آن زمان تا کنون، این حوزه موضوع پژوهش‌های فعال بوده است. فرض کنید Θ یک متغیر تصادفی دایره‌ای پیوسته و X یک متغیر تصادفی خطی پیوسته باشند. در حالی که X روی کل خط حقیقی ($x \in \mathbb{R}$) و متغیر زاویه‌ای Θ روی دایره واحد ($\theta \in S$) تعریف می‌شوند، در اینجا ما آن را به‌صورت بازه‌ای از طول 2π ، یعنی $[a, a + 2\pi)$ در نظر می‌گیریم.

هر تابع چگالی احتمال دایره‌ای-خطی پیوسته روی زیرمجموعه‌ای از فضای استوانه‌ای ($S \times \mathbb{R}$) تعریف شده است. بنابراین، دامنه مفصل دایره‌ای-خطی نیز به جای مربع واحد $([0, 1]^2)$ ، سطح یک استوانه با ارتفاع واحد و محیط واحد خواهد بود و با این فرض، در ادامه u را به‌عنوان بعد دایره‌ای و v را به‌عنوان بُعد خطی در نظر می‌گیریم. برای اختصار، مفصل‌هایی که چگالی آن‌ها دارای ویژگی زیر باشند را «دوره‌ای»^۱ در جهت « u » یا به اختصار «دوره‌ای» می‌نامیم (ویژگی اصلی یک مفصل دایره‌ای-خطی):

$$c(0, v) = c(1, v), \forall v \in [0, 1]. \quad (1.3)$$

بررسی رفتار حیوانات یکی از واضح‌ترین کاربردهای مفصل‌های دایره‌ای-خطی است. البته نتایج آن را می‌توان برای بسیاری از داده‌های دایره‌ای-خطی دیگر نیز تعمیم داد. در رفتار حیوانات، برای نمایش زاویه‌های چرخش معمولاً بازه $\Theta \in [-\pi, \pi)$ را در نظر می‌گیرند. در این بازه، زوایای منفی به چرخش به چپ و زوایای مثبت به چرخش به راست اشاره دارند. تکیه‌گاه متغیر تصادفی X (طول گام)، به اعداد حقیقی نامنفی محدود شده است. زاویه‌های چرخش تا حدودی با سایر متغیرهای دایره‌ای متفاوت‌اند، زیرا منطقی است فرض کنیم که یک حیوان تمایل ذاتی به چرخش به چپ یا راست ندارد. این فرض را می‌توان با استفاده از یک تابع چگالی احتمال حاشیه‌ای که حول $\theta = 0$ متقارن است، به همراه مفصلی که نه تنها چگالی آن دوره‌ای است (به رابطه ۱.۳ توجه کنید)، بلکه در بُعد u نیز حول $u = 0.5$ متقارن است، محقق ساخت (هودل و فبریگ، ۲۰۲۱). منظور ما از مفصل متقارن در u یا به‌طور خلاصه مفصل متقارن، مفصلی است که تابع چگالی آن در رابطه زیر صدق می‌کند:

$$c(u, v) = c(1 - u, v) \quad \forall u, v \in [0, 1]. \quad (2.3)$$

هرچند انواع دیگر تقارن نیز برای مفصل‌ها قابل تعریف است، اما در اینجا به همین تقارن بسنده می‌کنیم. (به نلسون (۲۰۰۶) مراجعه شود). هر مفصل متقارن، الزاماً دوره‌ای نیز است.

¹Periodic

۲.۳ چند مفصل دایره‌ای-خطی

در این بخش، چند مفصل دایره‌ای-خطی که ویژگی‌های مورد نیاز در رابطه (۱.۳) را دارند، معرفی می‌شوند. این مفصل‌ها برای مدل‌سازی وابستگی بین متغیرهای دایره‌ای و خطی طراحی شده‌اند.

• مفصل جانسون-وهرلی: مفصل جانسون-وهرلی به صورت زیر تعریف می‌شود (وانگ و همکاران، ۲۰۲۱):

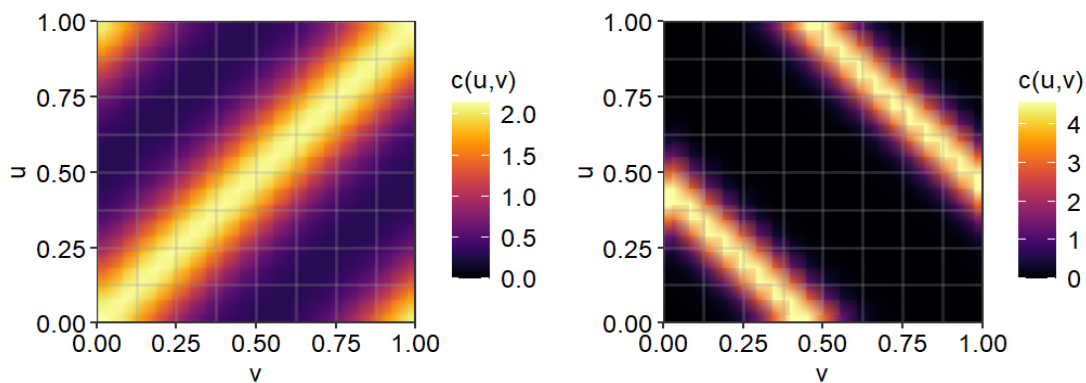
$$c(u, v) = \sqrt{2} \pi g(\eta), \quad \eta = \begin{cases} \sqrt{2} \pi(u - v), & u \geq v \\ \sqrt{2} \pi(u - v) + \sqrt{2} \pi, & u < v, \end{cases} \quad (۳.۳)$$

$$C(u, v) = \int_0^u dw_1 \int_0^v \sqrt{2} \pi g[\sqrt{2} \pi(w_1 - w_2)] dw_2 = \int_0^u [-G[\sqrt{2} \pi(w_1 - w_2)]^v] dw_1, \quad (۴.۳)$$

که در آن $\eta \in [0, \sqrt{2} \pi]$ یک متغیر دایره‌ای، $g(\cdot)$ تابع چگالی احتمال η و $G(\cdot)$ تابع توزیع متناظر با $g(\cdot)$ هستند. در رابطه (۳.۳)، اگر به جای $g(\cdot)$ توزیع فون‌میزس را قرار دهیم، به مفصل به دست آمده مفصل فون‌میزس استوانه‌ای می‌گویند. این مفصل به صورت زیر است:

$$C(u, v) = uv + \frac{1}{\sqrt{2} \pi^2 I_0(k)} \sum_{j=1}^{\infty} I_j(k) \left(\frac{\cos[\sqrt{2} \pi u - \sqrt{2} \pi v - \mu]}{j^2} + \frac{-\cos[j(-\sqrt{2} \pi v - \mu)] + \cos[-j\mu]}{j^2} \right), \quad (۵.۳)$$

که در آن $I_j(\cdot)$ تابع بسل اصلاح‌شده مرتبه j است (ابرامویتز و استیگون، ۱۹۶۵). پارامتر μ تابع چگالی احتمال را در جهت v جابجا می‌کند و k تمرکز را تعیین می‌نماید. تابع چگالی این مفصل در دو حالت مختلف ($\mu = \pi, k = 4$ و $\mu = 0, k = 1$) در شکل ۱ نشان داده شده است. شکل‌های سمت چپ و راست به ترتیب نشان می‌دهند که مفصل همبستگی مثبت دارد یا منفی. برای $k = 0$ مفصل استقلال را داریم و برای $k = \infty$ با شرط همبستگی مثبت به کران بالایی (M) و با شرط همبستگی منفی به کران پایینی (W) فرشه-هافدینگ میل می‌کند. در حالی که مفصل فون‌میزس استوانه‌ای یک مفصل دایره‌ای-خطی و دارای ویژگی تناوبی است، اما متقارن نیست؛ موضوعی که کارایی آن را برای مدل‌سازی داده‌های حرکتی محدود می‌کند. در ادامه، مفصل‌هایی معرفی می‌شوند که چگالی‌های آن‌ها هم تناوبی و هم متقارن در بُعد u ($u = 0/5$) هستند.



شکل ۱: چگالی مفصل فون‌میزس استوانه‌ای در دو حالت (چپ: $\mu = 0, k = 1$; راست: $\mu = \pi, k = 4$).

• مفصل چندجمله‌ای مرتبه دو^۲: مفصل چندجمله‌ای مرتبه دو به صورت زیر تعریف می‌شود:

$$c(u, v) = 1 + 2\pi\alpha \cos(2\pi u)(1 - 2v), \quad (6.3)$$

$$C(u, v) = uv + \alpha \sin(2\pi u)v(1 - v), \quad (7.3)$$

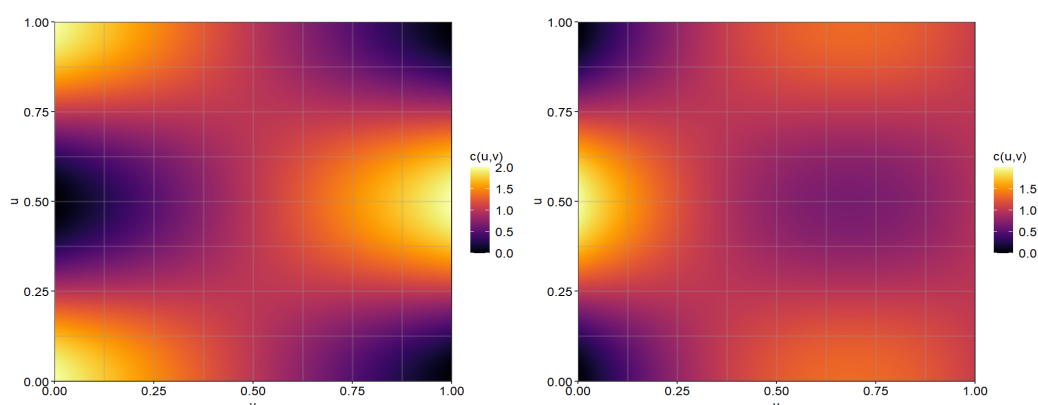
که در آن α پارامتر مفصل است و $|\alpha| \leq (2\pi)^{-1}$.

• مفصل چندجمله‌ای مرتبه سه^۳: به طور مشابه، مفصل چندجمله‌ای مرتبه سه به صورت زیر تعریف می‌گردد:

$$c(u, v) = 1 + 2\pi\alpha \cos(2\pi u)(1 - 4v + 3v^2) + 2\pi\beta \cos(2\pi u)(2v - 3v^2), \quad (8.3)$$

$$C(u, v) = v(1 - v)[\alpha \sin(2\pi u)(1 - v) - \beta \sin(2\pi u)v], \quad (9.3)$$

که در آن $|\alpha| \leq (2\pi)^{-1}$ و $|\beta| \leq (2\pi)^{-1}$. برای آسانی کار، از مفصل چندجمله‌ای مرتبه دو به عنوان مفصل مرتبه دو و از مفصل چندجمله‌ای مرتبه سه به عنوان مفصل مرتبه سه یاد می‌کنیم. آشکار است که مفصل‌های دایره‌ای-خطی مرتبه دو و مرتبه سه در بُعد u دایره‌ای متناوب و متقارن و در بُعد v خطی هستند. مفصل‌های مرتبه دو و مرتبه سه در مقاله [هودل و فیریگ \(۲۰۲۱\)](#) با جزئیات ارائه شده‌اند. همچنین، [هودل \(۲۰۲۲\)](#) بسته cylcop را در نرم‌افزار R برای مفصل‌های دایره‌ای-خطی با تقارن زاویه‌ای، ارائه کرد. توابع چگالی این دو مفصل در شکل ۲ رسم شده است. در این مفصل‌ها برای $\alpha = 0$ (و $\beta = 0$) مفصل استقلال به دست می‌آید، اما کران‌های فرشه-هافدینگ را به دست آورده نمی‌توانیم. برای $\alpha > 0$ (و $\beta < 0$)، مفصل همبستگی مثبت را بین u و v نشان می‌دهد هرگاه $u \in [0, 0.5]$ و همبستگی منفی را نشان می‌دهد هرگاه $u \in [0.5, 1]$. برای حالت مقابل یعنی $\alpha < 0$ (و $\beta > 0$) نیز این جهت‌های همبستگی عکس حالت قبل هستند.



شکل ۲: چگالی‌های مفصل دایره‌ای-خطی با بخش‌های ویژه نسبت به v . چپ: چگالی مفصل دایره‌ای-خطی مرتبه دو ($\alpha = \frac{1}{3\pi}$), راست: چگالی مفصل دایره‌ای-خطی مرتبه سه ($\alpha = -\frac{1}{3\pi}$, $\beta = 0.1$).

²Copula with Quadratic Sections

³Copula with Cubic Sections

۳.۳ برآورد پارامترها

با در اختیار داشتن مجموعه‌ای شامل n مشاهده‌ی مستقل و هم‌توزیع ($i.i.d$) از یک جفت متغیر دایره‌ای و خطی (θ_i, x_i) ، یک خانواده مفصل مفروض $C(u, v; \alpha)$ و توزیع‌های حاشیه‌ای مفروض $F_\Theta(\theta; \beta)$ و $F_X(x; \gamma)$ ، به‌طور نظری می‌توانیم بردار پارامترهای (α, β, γ) را از طریق ماکسیم‌سازی لگاریتم تابع درستنمایی توأم برآورد کنیم. با این حال، این رویکرد مستلزم بهینه‌سازی هم‌زمان تعداد زیادی از پارامترها است و مشخص‌سازی نادرست توزیع‌های حاشیه‌ای می‌تواند منجر به برآوردهای اریب برای پارامترهای مفصل شود و برعکس. اگر توزیع‌های حاشیه‌ای معلوم باشند، لگاریتم تابع درستنمایی به‌صورت زیر خواهد بود:

$$\ell(\alpha) = \sum_{i=1}^n \log [c(F_\Theta(\theta_i), F_X(x_i); \alpha)]. \quad (۱۰.۳)$$

این رویکرد به دو شیوه جایگزین برای برآورد پارامترهای مفصل منتهی می‌شود که در آنها برآوردهایی از F_Θ و F_X در رابطه (۱۰.۳) جای‌گذاری می‌شوند. این برآوردها، یعنی $\hat{F}_\Theta(\theta)$ و $\hat{F}_X(x)$ می‌توانند به‌صورت پارامتری با استفاده از روش ماکسیم درستنمایی (MLE)^۴ و ناپارامتری به‌صورت شبه‌مشاهده، یعنی نمونه‌گیری از مفصل تجربی حاصل شوند (هودل و فبریگ، ۲۰۲۱). برآورد پارامترهای مفصل با استفاده از حاشیه‌های ناپارامتری برآوردشده، به‌عنوان برآورد ماکسیم شبه‌درستنمایی (MPLE)^۵ شناخته می‌شود. برای بسیاری از مفصل‌های خطی-خطی، برآورد پارامترها را می‌توان به‌صورت تحلیلی از روی برآورد یک معیار وابستگی مانند ضریب اسپیرمن یا ضریب کندال به‌دست آورد؛ اما برای بیش‌تر مفصل‌های معرفی‌شده در این مقاله، راه‌حل تحلیلی در دسترس نیست و بنابراین نیاز به بهینه‌سازی عددی داریم.

برای سایر مفصل‌های دایره‌ای-خطی، می‌توان از جستجوی شبکه‌ای^۶ بهره‌گرفت تا پارامترهایی را یافت که اختلاف یک معیار وابستگی بین مفصل تجربی و مفصل پارامتری را مینم کند (هودل و فبریگ، ۲۰۲۱). در نهایت، انتخاب مدل را می‌توان بر پایه معیار اطلاع آکاییک (AIC)^۷ انجام داد؛ ولی در زمینه MPLE مورد مناقشه نیز بوده است؛ یعنی تعدادی از پژوهش‌گران باورمند هستند انتخاب مدلی‌که پارامترهای آن به روش MPLE برآورد شده‌اند، بر مبنای معیار اطلاع آکاییک درست نیست (گرونبرگ و هجورت، ۲۰۱۴).

۴.۳ معیارهای همبستگی

معیارهای پارامتری همبستگی دایره‌ای-خطی مستلزم آن هستند که متغیر خطی دارای توزیع نرمال باشد (جانسون و وهرلی، ۱۹۷۷؛ ماریا، ۱۹۷۶) و بنابراین، به‌طور کلی برای تمام داده‌های دایره‌ای-خطی قابل استفاده نیستند. از این‌رو، هودل و فبریگ (۲۰۲۱) از یک برآورد ناپارامتری برای ضریب همبستگی دایره‌ای-خطی، با نماد D ، استفاده کردند که برای نخستین‌بار توسط ماریا (۱۹۷۶) معرفی شد و در مطالعاتی چون سلو و همکاران (۱۹۸۸) و اخیراً توسط تاو (۲۰۱۵) به کار رفته است. این ضریب با رتبه‌بندی داده‌ها برآورد می‌شود و به‌صورت زیر است:

$$D = k \left(\left(\sum_i s_i \cos \frac{\gamma \pi r_i}{n} \right)^2 + \left(\sum_i s_i \sin \frac{\gamma \pi r_i}{n} \right)^2 \right), \quad (۱۱.۳)$$

^۴Maximum Likelihood Estimator

^۵Maximum Pseudo-likelihood Estimator

^۶Grid-Search

^۷Akaike Information Criterion

که در آن k یک ثابت نرمال‌سازی، r_i رتبه زاویه‌ای (دایره‌ای) مشاهده i (با نقطه مرجع $-\pi$) و s_i رتبه خطی همان مشاهده هستند. با این حال، ضریب D فقط می‌تواند مقادیری بین 0 و 1 را بگیرد و بنابراین نمی‌تواند همبستگی مثبت یا منفی را از هم تمایز دهد. همچنین، لیگوی و همکاران (۲۰۱۹) یک معیار همبستگی را بر مبنای اطلاع متقابل معیار واگرایی کولگ-لیبلر بین توزیع توأم و حاصل ضرب توزیع‌های حاشیه‌ای پیاده‌سازی کرده‌اند.

$$\begin{aligned} I(\Theta, X) &= \int_{-\pi}^{\pi} \int_0^{\infty} f_{\Theta, X}(\theta, x) \log \left(\frac{f_{\Theta, X}(\theta, x)}{f_{\Theta}(\theta) f_X(x)} \right) d\theta dx \\ &= \int_{-\pi}^{\pi} \int_0^{\infty} c(F_{\Theta}(\theta), F_X(x)) f_{\Theta}(x) f_X(x) \log \left(c(F_{\Theta}(\theta), F_X(x)) \right) d\theta dx \\ &= \int_0^1 \int_0^1 c(u, v) \log(c(u, v)) du dv. \end{aligned} \quad (12.3)$$

برابری آخر با تغییر متغیر $u = F_{\Theta}(\theta)$ و $v = F_X(x)$ به دست آمده است. از این عبارت، بلافاصله مشخص می‌شود که اطلاع متقابل برابر منفی آنترپوی مرتبط با تابع چگالی مفصل $c(u, v)$ است (ما و سون، ۲۰۱۱). از آنجایی که محاسبه انتگرال در خط آخر رابطه (۱۲.۳) دشوار است، نمی‌توانیم پارامترهای خانواده مفصل انتخاب شده را به طور مستقیم تعیین کنیم. با این حال، نشان داده شده است که اطلاع متقابل نسبت به پارامتر وابستگی در دامنه وسیعی از مفصل‌های دومتغیره، رفتار یکنوا دارد (تنزیر و ایلیدن، ۲۰۱۶). بنابراین، انتظار می‌رود که چارچوب بهینه‌سازی تقریبی توضیح داده شده در بالا، در استفاده از اطلاع متقابل به عنوان معیار وابستگی، به خوبی عمل کند.

می‌توان برآوردی از اطلاع متقابل را با دسته‌بندی داده‌ها (مشاهدات متغیر تصادفی دایره‌ای و خطی) و گرفتن میانگین اطلاع متقابل نقطه‌ای-که براساس فراوانی نقاط داده در هر دسته محاسبه می‌شود- بر روی تمام ترکیب‌های دسته‌ها به دست آورد. هم چنین توجه داشته باشید که اطلاع متقابل داده‌ها برابر با اطلاع متقابل مفصل تجربی است. در صورت نیاز، می‌توان اطلاع متقابل را با استفاده از آنترپوی H به بازه $[0, 1]$ طوری نرمال‌سازی کرد که مقدار صفر مفصل استقلال و مقدار یک مفصل بالایی فرشه-هافدینگ را معرفی کنند (تنزیر و ایلیدن، ۲۰۱۶):

$$I_{norm}(u, v) = \frac{I(u, v)}{\sqrt{H(u)H(v)}}. \quad (13.3)$$

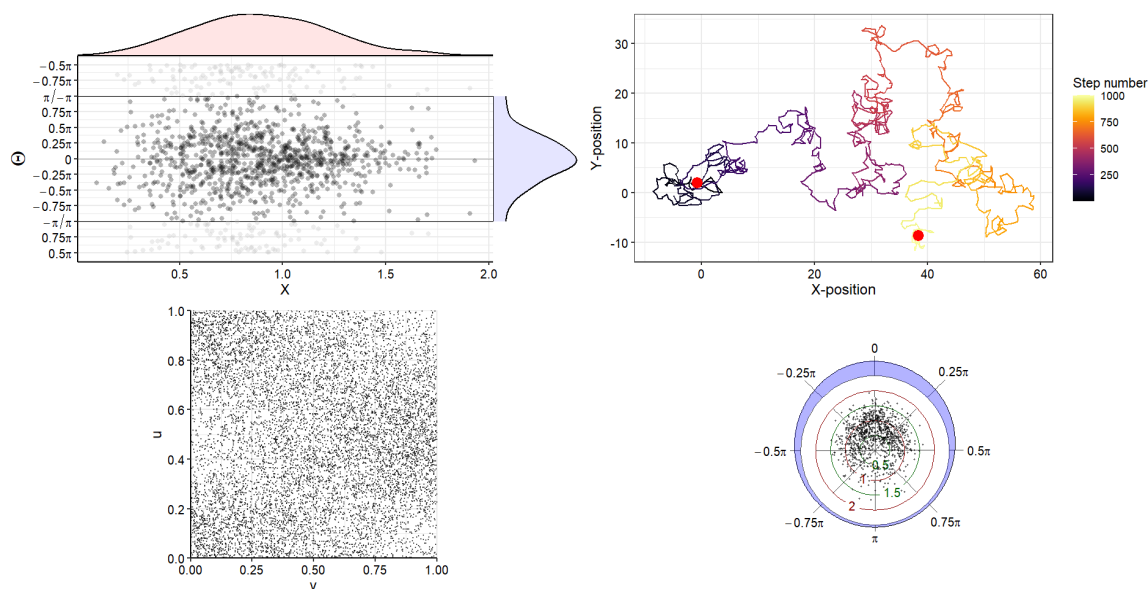
در نهایت، اگر بپذیریم که مفصل واقعی نسبت به متغیر u متقارن است، می‌توانیم مفصل تجربی را نیز متقارن‌سازی کنیم (به این صورت که مقادیر u بزرگ‌تر از 0.5 را با $1 - u$ جایگزین کنیم).

۴ کاربردها

همان‌طور که قبلاً اشاره شد، مفصل دایره‌ای-خطی فون‌میزس متقارن نیست و استفاده‌ی آن برای داده‌های حرکتی مانند داده‌های باد یا حرکت حیوانات که معمولاً متقارن فرض می‌شوند، برآوردهای مناسب به دست نمی‌دهد. اما مفصل‌های مرتبه دو و مرتبه سه این مشکل را برطرف می‌کنند. در این مقاله از مفصل مرتبه دو با $\alpha = 0.8$ ، تعداد 1000 داده تولید شده که توزیع حاشیه‌ای دایره‌ای (زاویه‌های چرخش) آن را توزیع فون‌میزس با پارامترهای $(\mu, k) = (0, 1)$ و توزیع حاشیه‌ای خطی (طول‌های گام) آن را توزیع وایبل با پارامترهای $(a, b) = (3, 1)$ تشکیل می‌دهند. چون مفصل‌های معرفی شده بیش تر به موضوعات حرکتی می‌پردازند، داده‌های تولید شده از آن‌ها را داده‌های مسیر^۸ نیز می‌گویند. نمودارهای این داده‌ها در شکل ۳ آورده شده است. در این شکل، نمودار بالا (سمت چپ)، نمودار پراکنش

^۸Trajectory Data

زاویه‌های چرخش و طول‌های گام به‌ترتیب با برآوردهای چگالی ناپارامتری به رنگ‌های آبی و سرخ را نشان می‌دهد؛ نمودار بالا (سمت راست)، نمودار مسیر است که از دنباله زاویه‌های چرخش و طول‌های گام مشتق شده است؛ نمودار پایین (سمت چپ)، نمودار پراکنش مفصل به‌کاربرده شده را نشان می‌دهد و نمودار پایین (سمت راست)، نمودار پراکنش دایره‌ای زاویه‌های چرخش و طول‌های گام است.



شکل ۳: نمایش داده‌های دایره‌ای-خطی که با استفاده از مفصل مرتبه دو با حاشیه‌های فون‌میزس و وایبل شبیه‌سازی شده است.

ضریب همبستگی بین متغیرهای دایره‌ای و خطی در دو حالت بررسی شده است. مقدار ضریب همبستگی رابطه (۱۱.۳) $D = 0.436$ و مقدار ضریب همبستگی اطلاع متقابل $I_{norm}(u, v) = 0.321$ هستند. هردو معیار، همبستگی بین زاویه‌های چرخش و طول‌های گام را کم نشان می‌دهند.

در ادامه، از مفصل مرتبه سه با پارامترهای $(\alpha, \beta) = (0.08, -\frac{1}{\pi})$ و توزیع‌های حاشیه‌ای مفصل قبلی داده تولید کردیم. تعداد داده‌های تولیدشده را در اینجا نیز ۱۰۰۰ در نظر گرفتیم و از رسم نمودارهای آن خودداری کردیم. ضریب همبستگی رابطه (۱۱.۳) و اطلاع متقابل به‌ترتیب به‌صورت $D = 0.955$ و $I_{norm}(u, v) = 0.422$ به‌دست آمدند. بر مبنای این مفصل نیز ضریب همبستگی‌ها کوچک هستند. در بخش بعدی، پارامترهای مفصل‌ها به‌روش‌های ماکسیمم درستنمایی و ماکسیمم شبه‌درستنمایی برآورد شده و اریبی و جذر میانگین مربعات خطای آن‌ها (RMSE)^۹ برای نمونه‌های مختلف در جدول ۱ آورده شده است. هردو روش برآوردهای نزدیک به پارامترهای اصلی را به‌دست می‌دهند. در جدول مشاهده می‌شود که با افزایش اندازه نمونه، اریبی و جذر میانگین مربعات خطا کاهش پیدا می‌کنند و هردو خطای کمی دارند. علاوه‌براین، هردو روش برآوردیابی برای هردو مفصل با پارامترهای مختلف، نتایج قابل قبولی ارائه می‌دهند و بنابراین، استفاده‌ی آن‌ها برای داده‌های واقعی مناسب به‌نظر می‌رسد.

^۹Root-Mean-Square Error

جدول ۱: بررسی عملکرد روش‌های MLE و MPLE برای برآورد پارامترهای مفصل‌ها

مفصل‌های ویژه		روش برآورد	اندازه نمونه (n)	$a = -0.1592$		$b = 0.1592$		$a = 0.08$		$b = -0.1592$	
				اریبی	RMSE	اریبی	RMSE	اریبی	RMSE	اریبی	RMSE
مرتبه سه	MLE	۱۰۰	۰/۰۱۸۴	۰/۰۳۶۸	-۰/۰۲۰۵	۰/۰۳۷۵	۰/۰۰۱۹	۰/۰۵۶۶	۰/۰۲۲۶	۰/۰۴۲۶	۰/۰۱۵۲
		۵۰۰	۰/۰۰۶۳	۰/۰۱۲۷	-۰/۰۰۶۸	۰/۰۱۳۴	۰/۰۰۱۲	۰/۰۲۶۳	۰/۰۰۷۹	۰/۰۱۵۲	۰/۰۱۱۱
		۱۰۰۰	۰/۰۰۴۴	۰/۰۰۸۸	-۰/۰۰۴۴	۰/۰۰۸۷	۰/۰۰۰۴	۰/۰۱۷۷	۰/۰۰۵۹	۰/۰۱۱۱	۰/۰۴۰۷
	MPLE	۱۰۰	۰/۰۱۶۰	۰/۰۳۳۵	-۰/۰۱۸۳	۰/۰۲۵۳	۰/۰۰۶۲	۰/۰۵۴۸	۰/۰۲۱۶	۰/۰۴۰۷	۰/۰۱۴۹
		۵۰۰	۰/۰۰۶۱	۰/۰۱۲۹	-۰/۰۰۶۶	۰/۰۱۳۱	۰/۰۰۱۳	۰/۰۲۵۵	۰/۰۰۷۷	۰/۰۱۴۹	۰/۰۱۰۷
		۱۰۰۰	۰/۰۰۴۶	۰/۰۰۹۰	-۰/۰۰۴۶	۰/۰۰۸۹	۰/۰۰۱۸	۰/۰۱۸۵	۰/۰۰۵۹	۰/۰۱۰۷	
مرتبه دو	MLE	۱۰۰	۰/۰۱۳۷	۰/۰۲۴۴			-۰/۰۰۳۰	۰/۰۳۷۱			
		۵۰۰	۰/۰۰۴۷	۰/۰۰۸۸			۰/۰۰۰۰	۰/۰۱۶۴			
		۱۰۰۰	۰/۰۰۳۱	۰/۰۰۵۸			-۰/۰۰۰۶	۰/۰۱۱۴			
	MPLE	۱۰۰	۰/۰۱۲۷	۰/۰۲۳۶			-۰/۰۰۱۴	۰/۰۳۹۰			
		۵۰۰	۰/۰۰۴۷	۰/۰۰۹۰			۰/۰۰۰۴	۰/۰۱۶۲			
		۱۰۰۰	۰/۰۰۳۲	۰/۰۰۶۰			-۰/۰۰۰۲	۰/۰۱۲۰			

بحث و نتیجه‌گیری

در این مقاله، ضمن معرفی داده‌های دایره‌ای و مفصل‌ها، مفصل‌های دایره‌ای-خطی و ویژگی‌های آن‌ها مورد بررسی قرار گرفت و به چند مفصل مهم در این حوزه پرداخته شد. نتایج شبیه‌سازی نشان دادند که مفصل‌های با بخش‌های ویژه از جمله مفصل مرتبه دو و مفصل مرتبه سه، می‌توانند برازش‌های مناسبی برای داده‌های حرکتی ارائه دهند و ساختار وابستگی در این نوع داده‌ها را به‌طور مناسب مدل می‌کنند. داده‌های حرکتی شامل داده‌های دایره‌ای و خطی می‌شوند و ویژگی‌های تناوب و تقارن را هم‌زمان دارا می‌باشند. نظر به این ویژگی‌ها، هر مفصل دایره‌ای-خطی نمی‌تواند وابستگی بین آن‌ها را به‌درستی مدل کند.

قدردانی و تشکر

نویسندگان مقاله از پیشنهادات داوران گرامی که موجب بهبود و ارائه بهتر آن شده‌اند، کمال سپاس و امتنان دارند.

- Abramowitz, M. and Stegun, I. A. (1965), *Handbook of Mathematical Functions*, New York, Dover.
- Durante, F. and Sempi, C. (2015), *Principles of Copula Theory*, Second Edition, Chapman and Hall/CRC.
- Genest, C. and Favre, A. C. (2007), Everything You Always Wanted to Know about Copula Modeling but Were Afraid to Ask, *Journal of Hydrologic Engineering*, **24**(4), 347–368.
- Gronneberg, S. and Hjort, N. L. (2014), The Copula Information Criteria, *Scandinavian Journal of Statistics*, **41**(2), 436–459.
- Hodel, F.H. and Fieberg, J. R. (2021), Cylcop: A Package for Circular-Linear Copulae with Angular Symmetry, *BioRxiv* (2021): 2021-07.
- Hodel, F. H. (2022), Package 'cylcop', URL: <https://cran.r-project.org/package=Cylcops>.
- Johnson, R. A and Wehrly, T. E. (1977), Measures and Models for Angular Correlation and Angular-Linear Correlation, *Journal of the Royal Statistical Society, Series B (Methodological)*, **39**(2), 222–229.
- Johnson, R. A and Wehrly, T. E. (1978), Some Angular-Linear Distributions and Related Regression Models, *Journal of the American Statistical Association*, **73**(363), 602–606.
- Leguey, I., Larranaga, P., Bielza, C. and Kato, K. (2019), A Circular-Linear Dependence Measure under Johnson-Wehrly Distributions and its Application in Bayesian Networks, *Information Sciences*, **486**, 240–253.
- Ma, J. and Sun, Z. (2011), Mutual Information is Copula Entropy, *Tsinghua Science and Technology*, **16**(1), 51–54.
- Mardia, K. V. (1976), Linear-Circular Correlation Coefficient and Rhythmometry, *Biometrika*, 403–405.
- Mardia, K. V. and Jupp, P. E. (2000), *Directional Statistics*, Second Edition, New York, John Wiley and Sons, Ltd, Chichester.
- Nagar, P., Bekker, A., Arashi, M., Kat C. J. and Barnard, A. C. (2024), A Dependent Circular-Linear Model For Multivariate Biomechanical Data: ILIZARV RING FIXATOR STUDY, *Statistical Methods in Medical Research*, **33**(9), 1660-1672.
- Nelsen, R. B. (2006), *An Introduction to Copulas*, Second Edition, New York, Springer.
- Sklar, M. (1959), Fonctions de repartition a n dimensions et leurs marges, *Annales de l'ISUP*, **8**(3), 229–231.

Solow, A. R., Bullister, J. L. and Nevison, C. (1988), An Application of Circular-Linear Correlation Analysis to the Relationship between Freon Concentration and Wind Direction in Woods Hole, Massachusetts, *Environmental Monitoring and Assessment*, **10**(3), 219–228.

Tenzer, Y. and Elidan, G. (2016), On the Monotonicity of the Copula Entropy, *arXiv preprint arXiv:1611.06714*.

Tu, R. (2015), A Study of the Parametric and Nonparametric Linear-Circular Correlation Coefficient, *California Polytechnic State University, San Luis Obispo*, 1–24, URL: <http://digitalcommons.calpoly.edu/statsp/51/>.

Wang, Z. W., Zhang, W. M., Zhang, Y. F. and Liu, Z. (2021), Circular-Linear-Linear Probabilistic Model Based on Vine Copulas: An Application to the Joint Distribution of Wind Direction, Wind Speed, and Air Temperature, *Journal of Wind Engineering and Industrial Aerodynamics*, **215**, 104704.

Circular-Linear Copulas: Properties and Applications

Najibullah Karimi¹, Ali Dastbaravarde²

¹PhD Student of Statistics, Yazd University

²Department of Statistics, Faculty of Mathematics, Yazd University, 89195-741, Yazd, Iran

Abstract: This review paper presents recent studies in the field of circular-linear copulas. To analyze the dependence structure of circular and linear data, various forms of circular-linear copulas are examined. The use and implementation of the cylcop package in R software have been discussed as well. Furthermore, through simulation, copula parameters are estimated. The results indicate that the copulas with specific sections are suitable for modeling movement data.

Keywords: Movement Data, Circular-Linear Data, Copula with Special Sections, Circular-Linear Copula.

Mathematics Subject Classification (2020): 62H11, 62G30, 62F99.



پانزدهمین سمینار احتمال
و فرایندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



مدل سازی فرایندهای سری زمانی شبکه عصبی بیزی برای شناسایی پیش زلزله با ناهنجاری های یونوسفری زلزله خاش

محدثه کیکاوسی^۱، فاطمه حسینی، امید کریمی

گروه آمار، دانشگاه سمنان

چکیده: فرایندهای سری زمانی حافظه بلند و کوتاه مدت برای مدل کردن داده های وابسته به گذشته به کار می روند که قادرند اطلاعات مربوط به زمان های گذشته را برای مدت طولانی حفظ کنند. با بهره گیری از شبکه های عصبی و روش های بیزی می توان روند زمانی ناهنجاری های یونوسفری را شناسایی و تحلیل کرد. تحلیل ناهنجاری های یونوسفری به شناسایی وقوع زمین لرزه ها با استفاده از پس لرزه ها کمک می کند. در این مقاله، مدل های حافظه بلند مدت با روش های کلاسیک و بیزی برای تحلیل و پیش بینی روندهای زمانی ناهنجاری های مربوط به زلزله خاش سال ۱۳۹۲ به کار گرفته می شوند. نتایج حاصل نشان می دهند که مدل پیشنهادی توانایی مناسبی در شناسایی الگوهای پنهان یونوسفری پیش از زمین لرزه دارند و می توان از آن ها به عنوان ابزاری کارآمد در تحلیل پیش نشانگرهای لرزه ای بهره گرفت.

واژه های کلیدی: فرایندهای زمانی، شبکه های عصبی بیزی، حافظه طولانی کوتاه مدت، زمین لرزه خاش.

کد موضوع بندی ریاضی (۲۰۲۰): 65K10, 62H11, 62J12.

۱ مقدمه

پیش بینی زلزله همواره یکی از چالش های مهم در علوم زمین بوده است. در سال های اخیر، استفاده از تحلیل تغییرات محتوای الکترونی کل^۱ (TEC) به عنوان یکی از روش های نوین برای شناسایی پیش نشانگرهای زلزله مورد توجه قرار گرفته است. پژوهش های متعددی نشان داده اند که بلایای طبیعی می توانند از طریق سازوکارها و شدت های مختلف، لایه یونوسفر^۲ را دچار اختلال کنند (چن و همکاران، ۲۰۲۰؛ فری شه و همکاران، ۲۰۲۴). با توجه به پیچیدگی رفتار زمین، استفاده از شبکه های عصبی با مدل سازی عدم قطعیت می تواند بینش عمیق تری در تحلیل این ناهنجاری ها ارائه دهد. در این مقاله، با استفاده از داده های TEC استخراج شده از فایل های یونکس^۳ و

^۱ سخنران، mohadeseh-keykavosi@semnan.ac.ir

^۱Total Electron Content

^۲Ionosphere

^۳IONEX

به کارگیری مدل‌های حافظه کوتاه و بلند مدت^۴ (LSTM) و LSTM بیزی^۵ (سقیب و همکاران، ۲۰۲۴) به تحلیل ناهنجاری‌های یونسفری پیش از وقوع زلزله شهر خاش در استان سیستان و بلوچستان پرداخته‌ایم. این زلزله ۲۷ فروردین ۱۳۹۲ با بزرگای حدود ۷/۵ رخ داده است. یکی از اهداف مهم، بررسی پیش‌بینی روند زمانی داده‌های TEC در بازه زمانی پیش از زلزله و ارزیابی هم‌زمانی ناهنجاری‌های شناسایی‌شده با زمان وقوع این زلزله می‌باشد.

۲ مدل‌های LSTM کلاسیک و بیزی

برای مدل‌کردن روند زمانی داده‌های TEC، از دو مدل شبکه عصبی بازگشتی LSTM کلاسیک و بیزی برای شناسایی ناهنجاری‌های احتمالی پیش از زلزله استفاده می‌شود. طبق آگاروال (۲۰۱۸) مدل LSTM کلاسیک یکی از انواع شبکه‌های عصبی بازگشتی است که برای مدل‌سازی داده‌های سری زمانی و وابسته به گذشته طراحی شده است. این مدل با استفاده از ساختار حافظه‌ای ویژه، قادر است اطلاعات مربوط به زمان‌های گذشته را برای مدت طولانی‌تری حفظ کند و از آن‌ها در یادگیری وابستگی‌های زمانی پیچیده بهره‌گیرد. در مدل‌های LSTM، چهار مؤلفه‌ی وضعیت سلول و دروازه ورودی^۶ و دروازه فراموشی^۷ و دروازه خروجی^۸ وجود دارد که نقش کلیدی در پردازش و حفظ اطلاعات ایفا می‌کنند. هر یک از این اجزا وظیفه‌ای خاص را بر عهده دارند و با همکاری یکدیگر، توانایی مدل در یادگیری و حفظ اطلاعات بلندمدت را بهبود می‌بخشند. برای درک بهتر عملکرد این مدل، فرض کنید در زمان t ، ورودی X_t ، بردار وزن‌ها $W_t = [W_f, W_i, W_c, W_o]$ و بردار اربیبی‌ها $b_t = [b_f, b_i, b_c, b_o]$ وجود دارد. خروجی مدل در این زمان نیز O_t خواهد بود. دروازه فراموشی به صورت $f_t = \sigma(W_f \cdot [O_{t-1}, X_t] + b_f)$ است که در آن، f_t خروجی دروازه فراموشی در زمان t است. تابع σ نشان‌دهنده تابع فعال‌سازی سیگموئید^۹ است که خروجی را در محدوده بین ۰ تا ۱ قرار می‌دهد. O_{t-1} خروجی مدل در زمان قبلی ($t-1$) است. همچنین، W_f وزن‌های مرتبط با دروازه فراموشی و b_f اربیبی این دروازه است. دروازه فراموشی با استفاده از این مکانیزم، تصمیم می‌گیرد که چه میزان از اطلاعات قبلی در وضعیت سلول حفظ شود. اگر f_t بزرگتر از ۰/۵ باشد، اطلاعات قبلی در حالت سلولی حفظ می‌شود؛ در غیر این صورت، اطلاعات فراموش می‌شود. این فرایند به مدل LSTM اجازه می‌دهد تا به صورت انتخابی و هوشمندانه، اطلاعات مهم را نگه داشته و اطلاعات غیرضروری را حذف کند، که این امر در پردازش توالی‌های زمانی و یادگیری الگوها بسیار حیاتی است.

دروازه ورودی در مدل LSTM به صورت $i_t = \sigma(W_i \cdot [O_{t-1}, X_t] + b_i)$ محاسبه می‌شود. که در آن، i_t خروجی دروازه ورودی در زمان t است و نشان می‌دهد چه اطلاعاتی از ورودی جدید باید وارد وضعیت سلول شود. وضعیت سلول پیشنهادی در مدل LSTM به صورت $\tilde{C}_t = \tanh(W_c \cdot [O_{t-1}, X_t] + b_c)$ تعریف می‌شود که در آن، \tilde{C}_t مقدار پیشنهادی برای به‌روزرسانی وضعیت سلول در زمان t است. تابع \tanh به عنوان تابع فعال‌سازی استفاده می‌شود و خروجی آن در محدوده بین -۱ و ۱ قرار می‌گیرد. W_c وزن‌های مرتبط با وضعیت سلول پیشنهادی و b_c اربیبی آن است. سپس وضعیت سلول به صورت $C_t = f_t \cdot C_{t-1} + i_t \cdot \tilde{C}_t$ به‌روزرسانی می‌شود. که در آن، C_t وضعیت سلول در زمان t است که با ترکیب اطلاعات قبلی و اطلاعات جدید به‌دست می‌آید. دروازه خروجی مدل LSTM به صورت $h_t = \sigma(W_o \cdot [O_{t-1}, X_t] + b_o)$ به دست می‌آید که در آن، h_t سیگنال خروجی دروازه خروجی

^۴Long Short-Term Memory

^۵Bayesian LSTM

^۶Input Gate

^۷Forget Gate

^۸Output Gate

^۹Sigmoid

در زمان t است که مقداری بین 0 و 1 اختیار می‌کند. پس از محاسبه h_t ، خروجی نهایی مدل در زمان t (O_t) با استفاده از وضعیت سلول به‌روزرسانی شده (C_t) و سیگنال خروجی دروازه خروجی به صورت $O_t = h_t \cdot \tanh(C_t)$ به دست می‌آید. در این مقاله، مدل LSTM پایه‌ای به‌کار گرفته شده است که شامل یک لایه LSTM با تعداد مشخصی واحد حافظه و یک لایه خروجی چگالی است. این مدل به‌صورت کلاسیک و قطعی آموزش داده می‌شود؛ به این معنا که در مرحله‌ی آموزش، از تکنیک دراپ آوت^{۱۰} برای جلوگیری از بیش‌برازش استفاده می‌شود، اما در مرحله آزمون، دراپ آوت غیرفعال می‌گردد. در نتیجه، مدل تنها یک مقدار مشخص و قطعی برای TEC در هر زمان پیش‌بینی می‌کند. این مدل توانایی مناسبی در یادگیری الگوهای زمانی TEC دارد، اما به دلیل ماهیت قطعی خود، امکان سنجش عدم قطعیت پیش‌بینی‌ها را فراهم نمی‌کند. به همین دلیل، در ادامه از رهیافت بیزی برای بهبود مدل LSTM بهره‌گرفته می‌شود تا علاوه بر پیش‌بینی، ارزیابی میزان اطمینان مدل نیز امکان‌پذیر باشد.

در ادامه مؤلفه‌های W_t و b_t را با استفاده از رهیافت بیزی بهینه می‌کنیم. در مدل‌های LSTM، وزن‌ها (W_t) و اربیی‌ها (b_t) با استفاده از توزیع‌های پیشین مدل‌سازی می‌شوند. این رهیافت عدم قطعیت‌ها را در پارامترها در نظر می‌گیرد و برازش‌ها را بهینه می‌کند. وزن‌ها و اربیی‌ها دارای توزیع نرمال به صورت $W_t \sim \mathcal{N}(\mu_0, \sigma^2)$ و $b_t \sim \mathcal{N}(\mu_0, \sigma^2)$ است که در آن، μ_0 میانگین اولیه‌ی پیشین است که انتظار اولیه برای مقدار وزن‌ها و اربیی‌ها را نشان می‌دهد و σ^2 واریانس پیشین است که میزان عدم قطعیت در پارامترها را مشخص می‌کند. در مدل LSTM، تنظیم اربیی برای بهینه‌سازی عملکرد مدل بسیار مهم است. همچنین، برای بهبود توانایی تعمیم‌پذیری مدل و جلوگیری از بیش‌برازش، از تکنیک دراپ آوت استفاده می‌شود. برای افزایش اعتمادپذیری مدل و لحاظ‌کردن عدم قطعیت در پیش‌بینی‌ها، از ساختار بیزی در قالب شبکه عصبی LSTM استفاده می‌شود. در این ساختار، با بهره‌گیری از روش‌های مونت کارلویی^{۱۱}، لایه‌های دراپ آوت نه‌تنها در زمان آموزش بلکه در مرحله‌ی آزمون نیز فعال باقی می‌مانند. نمونه‌های تصادفی مونت کارلویی امکان برآورد توزیع پسین را فراهم می‌سازد. با نمونه‌گیری چندباره از مدل، می‌توان میانگین پیش‌بینی‌ها و انحراف معیار آن‌ها را محاسبه کرد. تابع زیان در مدل LSTM به‌صورت $\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, O_t) + \lambda \sum_{j=1}^L (\|W_t\|^2 + \|b_t\|^2)$ تعریف می‌شود، که شامل دو بخش است. بخش اول، میانگین خطای پیش‌بینی را محاسبه می‌کند، که با استفاده از تابع \mathcal{L} اختلاف بین مقادیر واقعی (y_i) و پیش‌بینی‌های مدل (O_t) را اندازه‌گیری می‌کند. بخش دوم، یک عبارت تنظیم‌گر^{۱۲} است که با مجموع مربعات وزن‌ها (W) و اربیی‌ها (b) و ضرب آن در پارامتر λ از بیش‌برازش جلوگیری می‌کند.

طبق هستی و همکاران (۲۰۰۹)، برآوردگر ماکسیمم درست‌نمایی برای وزن‌ها (W_t) با بهینه‌کردن تابع زیان به دست می‌آید. این رویکرد تعادلی مناسب بین دقت پیش‌بینی و جلوگیری از بیش‌برازش برقرار می‌کند و به بهبود عملکرد و تعمیم‌پذیری مدل کمک می‌نماید. در رهیافت بیزی امکان به دست آوردن توزیع احتمالی مشاهدات پیش‌بینی‌شده را فراهم می‌کند. تابع درست‌نمایی به‌صورت زیر تعریف می‌شود:

$$p(y | X, W, \sigma^2) \propto \frac{1}{n\sqrt{\sigma^2}} e^{-\frac{1}{2\sigma^2}((y-XW)^T(y-XW))}.$$

این فرمول احتمال اینکه بردار خروجی y از داده‌های ورودی X و پارامترهای مدل W تولید شده باشد را نشان می‌دهد. که در آن، $(y - XW)^T(y - XW)$ مجموع مربعات خطاها (SSE)^{۱۳} است که اختلاف بین مقادیر واقعی y و مقادیر پیش‌بینی‌شده XW را اندازه‌گیری می‌کند. در مدل بیزی، از توزیع پیشین مزدوج برای پارامترها استفاده می‌شود. توزیع پیشین به‌صورت $p(W, \sigma^2) =$

¹⁰ Dropout

¹¹ Monte Carlo Dropout

¹² Regularization

¹³ sum squares errors

$p(W, \sigma^2 | y, X) \propto p(y | X, W, \sigma^2) p(W | \sigma^2) p(\sigma^2)$ می‌باشد. در آن صورت توزیع پسین به صورت تعریف می‌شود. توزیع پسین پارامترهای مدل به صورت

$$p(W, \sigma^2 | y, X) \propto \frac{1}{n\sqrt{\sigma^2}} e^{-\frac{1}{2\sigma^2}(y-XW)^T(y-XW)} \\ \times (\sigma^2)^{-(n-v)/2} e^{-\frac{1}{2\sigma^2}(W-\mu_0)^T\Lambda_0(W-\mu_0)} \\ \times (\sigma^2)^{-(\alpha_0-1)} e^{-\frac{b_0}{\sigma^2}}$$

به دست می‌آید. که در آن $(\sigma^2)^{-(n-v)/2} e^{-\frac{1}{2\sigma^2}(W-\mu_0)^T\Lambda_0(W-\mu_0)}$ توزیع پیشین وزن‌ها را با میانگین μ_0 و ماتریس دقت $\Lambda_0 = cI$ مشخص می‌کند و $(\sigma^2)^{-(\alpha_0-1)} e^{-\frac{b_0}{\sigma^2}}$ توزیع پیشین واریانس را با پارامترهای α_0 و β_0 مدل‌سازی می‌کند. ماتریس رنج cI نقشی کلیدی در تنظیم مدل و حل مشکل چندخطی بودن دارد. با اضافه کردن این ماتریس به $X^T X$ ، ماتریس معکوس پذیر شده و از بیش‌برازش جلوگیری می‌کند. پارامتر c شدت تنظیم را کنترل می‌کند؛ هرچه c بزرگتر باشد، تأثیر پیشین قوی‌تر است. میانگین پسین μ_n با فرمول $\mu_n = (X^T X + \Lambda_0)^{-1}(X^T y \hat{w} + \Lambda_0 \mu_0)$ محاسبه می‌شود که ترکیبی از اطلاعات داده‌ها و باورهای پیشین است. به‌روزرسانی پارامترهای بیزی شامل محاسبه $\Lambda_n = X^T X + \Lambda_0$ و $a_n = a_0 + \frac{n}{\sigma^2}$ است. همچنین، W_t با فرمول $W_t = W_0 + \frac{1}{\sigma^2}(y^T y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n)$ به‌روزرسانی می‌شود. بنابراین توزیع پیشگویی نهایی به صورت

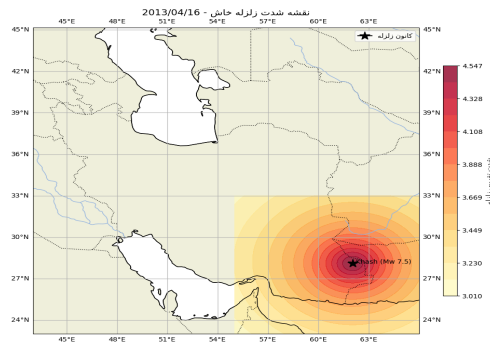
$$p(y | m) = \frac{p(m) \times p(y | X, W, \sigma, m)}{p(W, \sigma | y, X, m)}$$

محاسبه می‌شود. این فرمول بر اساس قانون بیز، توزیع پیشگویی $p(y | m)$ را به عنوان نسبت درست‌نمایی $p(y | X, W, \sigma, m)$ و توزیع پیشین $p(m)$ به توزیع پسین پارامترها $p(W, \sigma | y, X, m)$ محاسبه می‌کند. پیش‌بینی نهایی دارای توزیع نرمال است که نه تنها یک مقدار نقطه‌ای، بلکه یک بازه اطمینان برای مقدار خروجی ارائه می‌دهد. این بازه اطمینان، عدم قطعیت‌های موجود در برآورد پارامترها و نویز مدل را در نظر می‌گیرد و به کاربر امکان می‌دهد تا درک بهتری از محدوده احتمالی مقادیر پیش‌بینی شده داشته باشد. در نتیجه، LSTM بیزی نه تنها پیش‌بینی عددی ارائه می‌دهد، بلکه بازه اطمینان برای آن پیش‌بینی را نیز مشخص می‌کند. این ویژگی برای تحلیل پیش‌نشانگرهای زلزله که با عدم قطعیت‌های طبیعی همراه‌اند، اهمیت ویژه‌ای دارد.

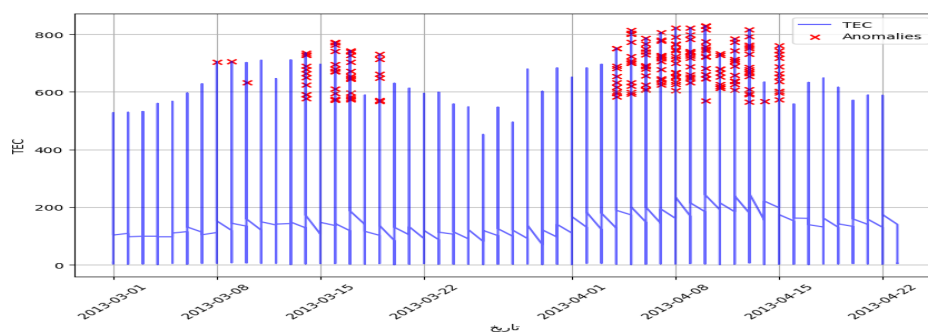
۳ تحلیل داده‌های زلزله خاش

برای مطالعه موردی، زلزله خاش در تاریخ ۲۷ فروردین ۱۳۹۲ انتخاب شد. پس از استخراج داده‌ها، به منظور تحلیل دقیق‌تر و شناسایی ناهنجاری‌ها، شاخص نمره‌گذاری استاندارد^{۱۴} به صورت $Z_i = \frac{x_i - \mu}{\sigma}$ برای هر روز محاسبه شد، که در آن x_i مقدار TEC در زمان مشخص، μ میانگین و σ انحراف معیار مقادیر TEC در یک بازه زمانی مرجع است. شکل ۱ موقعیت و شدت زلزله خاش را روی نقشه کشور ایران نشان می‌دهد، این زلزله با بزرگای ۷/۵ در عمق حدود ۹۲ کیلومتری زمین به وقوع پیوست. دلیل انتخاب این زلزله، قدرت بالا، گستره تأثیر قابل توجه و ثبت دقیق آن در منابع بین‌المللی بوده است که امکان بررسی ناهنجاری‌های یونوسفری پیش از رخداد را با دقت و وضوح بیشتری فراهم می‌سازد. شکل ۲ روند تغییرات TEC را در بازه زمانی بین ابتدای مارس تا اواخر آوریل ۲۰۱۳ برای ایستگاه خاش نمایش می‌دهد. در این بازه زمانی، دو دوره‌ی پرتراکم از ناهنجاری‌ها به وضوح قابل مشاهده هستند: دوره اول: از حدود ۸ مارس تا ۱۶ مارس، که در آن تعداد قابل توجهی از مقادیر TEC از روند معمول خود فراتر رفته‌اند. دوره دوم: از حدود ۵

¹⁴Z-score



شکل ۱: نقشه شدت زلزله خاش



شکل ۲: نمودار TEC با نقاط ناهنجاری: مقادیر TEC با استفاده از خطوط نمایش داده شده‌اند و نقاط ناهنجار استاندارد شده با علامت ستاره مشخص شده‌اند.

آوریل تا ۱۵ آوریل، که در این بازه نیز مقادیر TEC دارای پراکندگی زیاد و رفتار غیرعادی بوده‌اند. این رفتار ناهنجار ممکن است با پیش‌نشانگرهای ژئوفیزیکی یا ژئومغناطیسی زلزله‌ها مرتبط باشد و می‌تواند به عنوان هشدار اولیه مورد استفاده قرار می‌گیرد.

در ادامه، مدل‌های LSTM کلاسیک و بیزی روی داده‌ها اجرا و مورد ارزیابی قرار گرفت. هدف اصلی، بررسی دقت پیش‌بینی این دو مدل، ارزیابی ناهنجاری‌های احتمالی در داده‌ها پیش از رخداد زلزله و تحلیل عدم قطعیت در پیش‌بینی‌ها است. برای این منظور داده‌ها به دو بخش آموزش^{۱۵} (۸۰٪) و آزمون^{۱۶} (۲۰٪) تقسیم شدند. هر دو مدل روی بازه زمانی پیش از زلزله آموزش داده شدند و عملکرد آن‌ها در بازه نزدیک به وقوع زلزله (تقریباً ۴۶ روز قبل تا ۷ روز بعد از زلزله) مورد ارزیابی قرار گرفتند. معیارهای میانگین قدرمطلق خطا (MAE)^{۱۷} و ریشه میانگین مربعات خطا (RMSE)^{۱۸} و نرمال میانگین قدرمطلق خطا (NMAE)^{۱۹} و نرمال ریشه میانگین مربعات خطا (NRMSE)^{۲۰} برای ارزیابی دقت این مدل‌ها به صورت $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2}$ ، $NMAE = \frac{MAE}{\max(y_i) - \min(y_i)}$ ، $NRMSE = \frac{RMSE}{\max(y_i) - \min(y_i)}$ محاسبه شده‌اند، که در آن y_i داده‌های واقعی TEC و \hat{y}_i مقدار پیشگویی آن است. نتایج ارزیابی مدل‌های LSTM کلاسیک و بیزی روی داده‌های خاش در جدول ۱ خلاصه شده است. نتایج جدول ۱ نشان می‌دهد که مدل LSTM قادر است تغییرات TEC را با خطای نسبتاً کنترل‌شده‌ای

¹⁵Train

¹⁶Test

¹⁷Mean Absolute Error

¹⁸Root Mean Squared Error

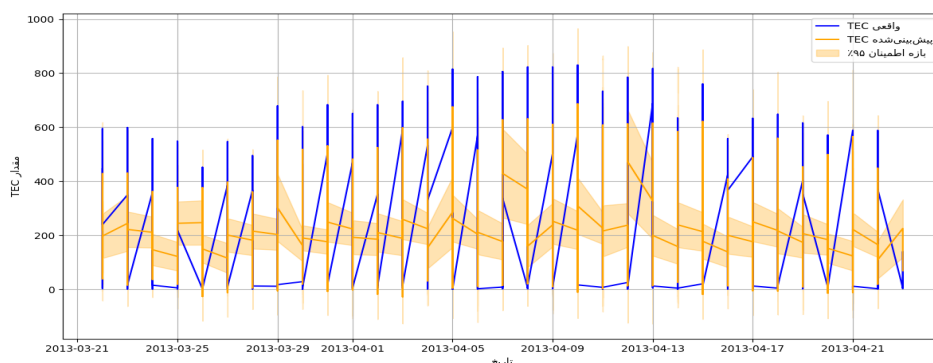
¹⁹Normal Mean Absolute Error

²⁰Normal Root Mean Squared Error

جدول ۱: خلاصه نتایج دو مدل LSTM کلاسیک و بیزی

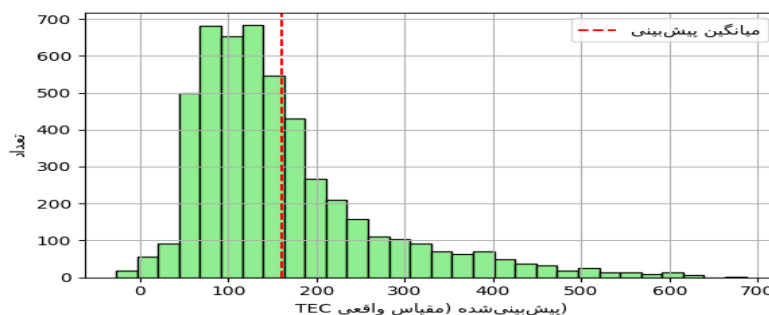
معیار	LSTM کلاسیک	LSTM بیزی
MAE	۱۳۵/۷۰	۱۳۵/۱۶
RMSE	۱۸۷/۵۰	۱۷۵/۱۷
NMAE	۰/۱۶۳۷	۰/۱۶۳۰
NRMSE	۰/۲۲۶۲	۰/۲۱۱۳

پیش‌بینی نماید، به‌ویژه آن‌که مقدار NMAE کمتر از ۰/۲ بوده که نشان‌دهنده پیش‌بینی قابل قبول در مقیاس تغییرات طبیعی داده‌هاست. همان‌گونه که مشاهده می‌شود، هرچند دقت کلی مدل LSTM بیزی نسبت به LSTM تفاوت زیادی ندارد، اما مقدار RMSE کاهش یافته است که بیانگر کاهش اثر نمونه‌های دارای خطای بالا در مدل بیزی است. همه تحلیل داده‌ها با پایتون انجام شده است. در شکل ۳، علاوه بر مقایسه مقدار واقعی با مقدار پیش‌بینی شده، بازه اطمینان ۹۵٪ نیز نمایش داده شده است.



شکل ۳: نمودار پیش‌بینی سری زمانی TEC با مدل LSTM بیزی

مدل LSTM بیزی نه تنها پیش‌بینی را انجام می‌دهد، بلکه عدم قطعیت آن را نیز نشان می‌دهد. عرض بازه‌های اطمینان در حوالی روزهای پیش از زلزله افزایش می‌یابد، که می‌تواند نشانه‌ای از افزایش عدم اطمینان مدل و احتمال بروز ناهنجاری باشد. این ویژگی روش‌های بیزی برای تحلیل پیش‌نشانگرهای زلزله بسیار کاربردی است، زیرا به جای پیش‌بینی قطعی، بازه‌ای از پیش‌بینی‌ها ارائه می‌دهد. شکل ۴ توزیع میانگین پیش‌بینی‌ها توسط مدل بیزی را نمایش می‌دهد. توزیع به‌صورت نامتقارن و حول مقدار میانگین



شکل ۴: نمودار هیستوگرام پیش‌بینی مدل LSTM بیزی

متمرکز است، اما نسبت به هیستوگرام، LSTM پراکندگی کمتری دارد که باعث دقت پیشگویی در مدل بیزی است. این نامتقارنی می‌تواند بازتاب‌دهنده ناپایداری یا تنش در لایه یونوسفر باشد. مدل LSTM بیزی دقت نسبتاً بالاتری در پیش‌بینی دارد و مهم‌تر از آن، با ارائه بازه اطمینان، قدرت تحلیل ناهنجاری‌های ساختاری در داده را فراهم می‌آورد. به‌ویژه در نواحی زمانی پیش از زلزله، افزایش عرض بازه اطمینان می‌تواند به عنوان شاخصی برای شناسایی ناهنجاری در لایه یونوسفر در نظر گرفته شود.

بحث و نتیجه‌گیری

در این مقاله با استفاده از مدل سری‌زمانی LSTM کلاسیک و بیزی با شاخص نمره‌گذاری استاندارد، ناهنجاری‌های احتمالی در داده‌های TEC استخراج‌شده از فایل‌های یونکس شناسایی گردید. این ناهنجاری‌ها عمدتاً به‌صورت خوشه‌ای و در بازه‌های زمانی مشخصی پیش از وقوع برخی زلزله‌ها ظاهر شدند که می‌توانند به‌عنوان نشانه‌هایی از پیش‌نشانگرهای یونوسفری تفسیر شوند. در ادامه، به‌منظور بررسی و مدل‌کردن رفتار زمانی این ناهنجاری‌ها، از دو مدل یادگیری عمیق شامل LSTM کلاسیک و LSTM بیزی با مرزهای عدم‌قطعیت استفاده شد. نتایج این مدل‌ها نشان داد که مدل بیزی علاوه بر ارائه دقت قابل‌قبول، قابلیت نمایش عدم‌قطعیت پیش‌بینی‌ها را نیز داراست که در تحلیل داده‌های حساس و پرنوسان مانند TEC بسیار ارزشمند است. با توجه به این یافته‌ها، پیشنهاد می‌شود در مطالعات آتی مدل‌های پیشرفته‌تر بیزی، گراف‌پایه به‌کار گرفته شود تا همبستگی دقیق‌تری بین ناهنجاری‌های TEC مشخص گردد و گامی مؤثر در راستای توسعه سامانه‌های هشدار زلزله برداشته شود.

مراجع

- Chen J., Zhang X., Ren X., Zhang J., Freeshah M., Zhao Z. (2020), Ionospheric disturbances detected during a typhoon based on GNSS phase observations: A case study for typhoon Mangkhut over Hong Kong, *Adv. Space Res.*, **66**, 1743–1753.
- Freeshah M., Adil M.A., Sentürk E., Zhang X., Ren X., Liu H., Osama N. (2024a), A cyclone formation, eastward plume drag, ionhydration process, and the consequent ionospheric changes following the 2022 Hunga Tonga-Hunga Ha'apai volcanic eruption, *Advances in Space Research*, **73** (5), 2457–2470.
- Hastie, T., Tibshirani, R., Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.
- Saqib, M., Sentürk, E., Adil, M.A., Freeshah, M. (2024). Seismo-ionospheric precursory detection using hybrid Bayesian-LSTM network model with uncertainty-boundaries and anomaly-intensity. *Advances in Space Research*, **74**, 1828–1842.
- Aggarwal, C. C. (2018). *Neural networks and deep learning: A textbook*. Springer International Publishing.
- <https://doi.org/10.1007/978-3-319-94463-0>

Bayesian Neural Network Temporal Processes Modeling for Earthquake Precursors Identification Using Ionospheric Anomalies of Khash Earthquake”

Mohadeseh Keykavosi, Fatemeh Hoseini, Omid Karimi

Department of statistics, Semnan University

Abstract: Long-short term memory time series processes are employed to model past-dependent data, capable of retaining information from past times for extended periods. By leveraging neural networks and Bayesian methods, the temporal trends of ionospheric anomalies can be identified and analyzed. The analysis of ionospheric anomalies aids in detecting earthquake occurrences using aftershocks. In this paper, long-term memory models, both classical and Bayesian, are utilized to analyze and predict the temporal trends of anomalies related to the 2013 Khash earthquake. The results demonstrate that the proposed model has a suitable capability in identifying hidden ionospheric patterns before earthquakes and can be used as an effective tool in analyzing seismic precursors.

Keywords: Temporal Processes, Bayesian Neural Networks, Long- Short term Memory, Khash Earthquake.

Mathematics Subject Classification (2020): 62J12, 62H11, 65K10.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30- 31, 2025
University of Kurdistan



پیش‌بینی زمان حوادث ترافیکی با استفاده از شبکه عصبی بقا

فاطمه محمدی^۱، محمد آرشی^۱ آرزو حبیبی‌راد^۱ ابوالفضل محمدزاده مقدم^۲

^۱گروه آمار دانشگاه فردوسی مشهد

^۲گروه عمران دانشگاه فردوسی مشهد

چکیده: حوادث ترافیکی یکی از چالش‌های اصلی در مدیریت ترافیک شهری هستند که می‌توانند تأثیرات منفی بر ایمنی و کارایی سیستم‌های حمل و نقل داشته باشند. آژانس‌های حمل و نقل برای مدیریت حوادث ترافیکی و به کارگیری استراتژی‌های مناسب نیاز به پیش‌بینی مدت زمان حادثه دارند. داده‌های زمان حوادث ترافیکی داده‌های زمان تا رویداد هستند. یکی از ویژگی‌های این داده‌ها وابستگی به مدت زمان است. یکی از مدل‌های رایج برای در نظر گرفتن این ویژگی، مدل‌های نرخ خطر مبنا است. عملکرد مدل‌های نرخ خطر مبنا در پیش‌بینی زمان حوادث ترافیکی، با توجه به فرضیاتی که دارند ممکن است دارای محدودیت‌هایی باشد. از جمله این فرضیات می‌توان به فرض خطرات متناسب، در نظر گرفتن توزیع خاص برای داده‌ها، اثر ثابت با زمان عوامل تأثیرگذار و رابطه خطی بین متغیرهای مستقل و وابسته اشاره کرد و ممکن است برخی از این فرضیات در داده‌های حوادث ترافیکی وجود نداشته باشد. یکی دیگر از روش‌های رایج برای مدل‌بندی زمان حوادث ترافیکی شبکه‌های عصبی مصنوعی می‌باشد. نقطه قوت این مدل‌ها این است که برخی از فرضیات مشکل‌ساز مدل‌های نرخ خطر مبنا را ندارند. اگر هدف این مدل‌ها پیش‌بینی مدت زمان پاکسازی حادثه باشد، نمی‌توانند وابستگی به مدت زمان را در داده‌های حوادث ترافیکی در نظر بگیرند. لذا از مدل‌های شبکه‌های عصبی با هدف برآورد تابع بقا استفاده می‌شود. اما این مدل‌ها هم نمی‌توانند مستقیم مدت زمان یک حادثه ترافیکی را پیش‌بینی کنند. بنابراین در این مطالعه مدل شبکه عصبی مبتنی بر تحلیل بقا معرفی می‌شود که بر اساس ساختار یادگیری چند وظیفه‌ای، دو وظیفه برآورد تابع بقا و پیش‌بینی زمان حوادث ترافیکی را همزمان انجام می‌دهد. مدل پیشنهادی می‌تواند وابستگی به مدت زمان را هنگام پیش‌بینی زمان حادثه با برآورد همزمان تابع بقا در نظر بگیرد. نتایج، عملکرد برتر مدل مورد مطالعه را نسبت به مدل‌های آماری و یادگیری ماشین نشان می‌دهد.

واژه‌های کلیدی: تحلیل بقای عمیق، مدل‌های نرخ خطر مبنا، زمان پاکسازی حادثه، تابع بقا، یادگیری چند وظیفه‌ای

کد موضوع بندی ریاضی (۲۰۲۰): xxCxx, xxBxx, xxAxx.

۱ مقدمه

حوادث ترافیکی تأثیرات منفی بر ایمنی و کارایی سیستم‌های حمل و نقل دارند. از جمله این تأثیرات می‌توان به ایجاد ترافیک و تأخیر در عبور و مرور، آسیب‌ها و تلفات جانی و بروز حوادث ثانویه اشاره کرد. بنابراین به عنوان یکی از چالش‌های اصلی در مدیریت ترافیک محسوب می‌شوند. سازمان‌های مدیریت ترافیک برای به‌کارگیری استراتژی‌های مناسب و کاهش این اثرات منفی نیاز به پیش‌بینی مدت زمان حادثه دارند. مدت زمان حادثه را می‌توان به چهار مرحله تقسیم کرد: زمان تشخیص، زمان پاسخ، زمان پاکسازی و زمان بازیابی. زمان پاکسازی که زمان ورود بین یک تیم پاسخ و پاکسازی کامل حادثه است طولانی‌ترین مرحله است. یکی از اهداف مدل‌بندی داده‌های حوادث ترافیکی، پیش‌بینی زمان پاکسازی حادثه است. داده‌های زمان پاکسازی حادثه داده‌های مدت زمان هستند. یکی از ویژگی‌های این داده‌ها، وابستگی به مدت زمان است. برای در نظر گرفتن این ویژگی عمده‌تاً از مدل‌های آماری نرخ خطر مبنا استفاده شده است. **نام و منرینگ (۲۰۰۰)؛ جونز و همکاران (۱۹۹۱)؛ چانگ (۲۰۱۰). نام و منرینگ (۲۰۰۰)** مدت زمان حادثه را با استفاده از یک مدل خطر متناسب پارامتری که متغیر وابسته آن نرخ خطر است، تحلیل کردند. آنها گزارش داده‌اند که پارامترهای برآورد شده و عملکرد پیش‌بینی مدل زمان شکست شتابیده (AFT) در طول سال‌های مختلف پایدار است. مدل خطرات متناسب کاکس دارای فرض خطرات متناسب هستند که این فرض ممکن است برای همه متغیرهای مستقل برقرار نباشد. در مطالعات قبلی برای پیش‌بینی مدت زمان حادثه از مدل‌های زمان شکست شتابیده (AFT) استفاده شده است. **جونز و همکاران (۱۹۹۱)** مدل زمان شکست شتابیده (AFT) که متغیر وابسته آن میانگین مدت زمان حادثه بود برای مدل‌بندی داده‌های حوادث ترافیکی به کار بردند. **چانگ (۲۰۱۰)** مدل زمان شکست شتابیده را برای مجموعه داده‌های حادثه در مقیاس بزرگ در بزرگراه اعمال کرد **توسلی حجتی و همکاران (۲۰۱۳)** ناهمگنی مشاهده نشده را با اعمال پارامترهای تصادفی در مدل زمان شکست شتابیده در نظر گرفت. مدل زمان شکست شتابیده (AFT) یک مدل پرکاربرد است که اثرات ثابت با زمان عوامل تأثیرگذار را فرض می‌کند. همچنین یکی دیگر از فرضیات این مدل‌ها رابطه خطی متغیرهای مستقل با لگاریتم میانگین زمان پاکسازی است. این دو فرض ممکن است برای داده‌های حوادث ترافیکی و در یک آزادراه برقرار نباشد. بنابراین استفاده از این مدل‌ها دارای محدودیت‌هایی می‌باشد. یکی دیگر از روش‌های رایج برای مدل‌سازی زمان پاکسازی حادثه مدل‌های یادگیری ماشین از جمله مدل جنگل تصادفی، ماشین‌گرادیان تقویتی و شبکه‌های عصبی مصنوعی است. با این حال مدل‌های یادگیری ماشین میانگین زمان پاکسازی را به عنوان متغیر وابسته در نظر می‌گیرند و ویژگی وابستگی به مدت زمان را نادیده می‌گیرند. مطالعه‌ای اخیر گزارش داده است که عملکرد مدل‌های یادگیری ماشین قابل مقایسه یا بدتر از مدل‌های زمان شکست شتابیده است. **وی و همکاران (۲۰۰۷)** برای پیش‌بینی مدت زمان حادثه از پیش‌بینی متوالی استفاده کردند که عملکرد مدل‌های شبکه‌های عصبی مصنوعی را بهبود بخشید. **حمد و همکاران (۲۰۲۰)** مدل جنگل تصادفی (RF) را برای پیش‌بینی زمان پاکسازی حادثه توسعه دادند. همچنین **ما و همکاران (۲۰۱۷)** یک مدل GBM را برای پیش‌بینی زمان پاکسازی حادثه و همچنین تحلیل عوامل تأثیرگذار توسعه دادند. نتایج آن‌ها عملکرد پیش‌بینی بهتر GBM را نسبت به RF و ANN نشان داد. **لی و همکاران (۲۰۲۰)** پیش‌بینی مدت زمان حادثه را با استفاده از شبکه عصبی عمیق انجام داد. نتایج آن‌ها نشان داد که DNN بهتر از مدل‌های یادگیری ماشین معمولی با هدف پیش‌بینی میانگین زمان حادثه، عمل می‌کند. مطالعه اخیر **تانگ و همکاران (۲۰۲۰)** گزارش داد که عملکرد پیش‌بینی مدل‌های یادگیری ماشین قابل مقایسه یا بدتر از مدل‌های زمان شکست شتابیده (AFT) با توزیع‌های فرضی مختلف بوده است. مطالعات قبلی از مدل‌های یادگیری ماشین و شبکه‌های عصبی عمیق به عنوان رگرسیون برای زمان پاکسازی حادثه استفاده می‌کردند. اگر هدف این مدل‌ها پیش‌بینی زمان پاکسازی حادثه باشد نمی‌توانند وابستگی به مدت زمان را که یکی از ویژگی‌های مهم داده‌های زمان تا رویداد

است به طور مستقیم در نظر بگیرند. این مدل‌ها تنها یک وظیفه انجام می‌دهند و آن هم پیش‌بینی میانگین زمان پاکسازی حادثه است. اگر مدل‌های یادگیری ماشین و شبکه‌های عصبی عمیق^۱ با هدف برآورد تابع بقا طراحی شوند می‌توانند مدل‌های مناسبی برای تحلیل داده‌های زمان تا رویداد باشند زیرا می‌توانند ویژگی وابستگی به مدت زمان در این داده‌ها را در نظر بگیرند. **ایشواران و همکاران (۲۰۰۸)** با توسعه جنگل بقای تصادفی که مدل جنگل تصادفی را برای داده‌های زمان تا رویداد گسترش می‌دهد، تابع بقا را برآورد کند. **هوئورن و همکاران (۲۰۰۶)** مدل ماشین تقویتی گرایان اصلاح شده^۲ را پیشنهاد کردند. **کاترمن و همکاران (۲۰۱۸)** DeepSurv را ابداع کردند که با استفاده از شبکه‌های عصبی عمیق روابط غیرخطی و تعاملی را در نظر بگیرد و تابع بقا را برآورد کند. این مدل از مدل‌های خطرات متناسب کاکس و جنگل بقای تصادفی بهتر عمل کرد. **لی و همکاران (۲۰۱۸)** DeepHit را پیشنهاد کردند که به طور مستقیم تابع بقا را با استفاده از شبکه‌های عصبی عمیق برآورد می‌کند. با اینکه مدل‌های یادگیری ماشین با هدف برآورد تابع بقا می‌توانند وابستگی به مدت زمان را در نظر بگیرند اما به صراحت نمی‌توانند زمان پاکسازی حادثه را پیش‌بینی کنند. در ادامه مدلی ارائه می‌شود که می‌تواند به طور همزمان هم تابع بقا را برآورد کند و هم مدت زمان حادثه را با استفاده از ساختار یادگیری چند وظیفه‌ای پیش‌بینی کند.

۲ تعاریف اولیه

داده‌های زمان پاکسازی حادثه داده‌های مدت زمان هستند که زمانی است تا وقوع رویداد. منظور از رویداد پاکسازی کامل حادثه است. برخی از مفاهیم مورد نیاز در تحلیل داده‌های زمان تا رویداد شامل تابع بقا، تابع نرخ خطر، داده سانسور شده، ویژگی وابستگی به مدت زمان و میانگین طول عمر می‌باشد. در ابتدا لازم است این مفاهیم توضیح داده شود.

تعریف ۱.۲. تابع بقا^۳ احتمال این است که زمان بقا بیش از یک زمان مشخص باشد. این تابع، که به نام تابع قابلیت اطمینان نیز مشهور است احتمال زنده ماندن در طول زمان t را محاسبه می‌کند. عموماً این تابع را با استفاده از روش‌های ناپارامتری مانند کاپلان-مایر یا نلسون آلن و یا روش‌های پارامتریک و از روی تابع توزیع به دست می‌آورند. متغیر تصادفی T که همان زمان بقا است می‌تواند گسسته یا پیوسته باشد اما همواره باید $T \geq 0$ باشد.

$$S(t) = Pr(T > t) = 1 - F(t)$$

تعریف ۲.۲. تابع نرخ خطر^۴ که با $h(t)$ نشان داده می‌شود، میزان مرگ و میر لحظه‌ای را نشان می‌دهد. این تابع شانس رخداد پیشامد را در لحظه t با شرط آن که واحد آزمایشی تا این لحظه زنده باشد، بیان می‌کند.

• اگر T متغیر تصادفی گسسته باشد:

$$h(t) = Pr(T = t | T \geq t)$$

• اگر T متغیر تصادفی پیوسته باشد:

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{Pr(t < T < t + \Delta t)}{\Delta t}$$

¹Deep Neural Network

²modified Gradient Boosting Machine

³Survival Function

⁴Hazard Function

با استفاده از تابع نرخ خطر می‌توان بقا را برآورد کرد. s متغیر انتگرال‌گیری است و مقدار $[0, t]$ را می‌گیرد.

$$S(t) = \exp\left(-\int_0^t h(s) ds\right)$$

با استفاده از تابع نرخ خطر می‌توان تابع نرخ خطر تجمعی را به دست آورد که با نماد $H(t)$ نشان داده می‌شود.

$$H(t) = \int_0^t h(s) ds$$

بنابراین:

$$S(t) = \exp(-H(t))$$

تعریف ۳.۲. در کنار داده‌های زمان تا رویداد، وضعیت آن داده هم ثبت می‌شود. منظور از وضعیت داده این است که آیا داده سانسور شده است یا داده کامل است. منظور از داده کامل این است که اطلاعات داده از اول مطالعه تا زمانی که رویداد برایش اتفاق می‌افتد موجود باشد. منظور از داده سانسور شده هم این است که اطلاعات کاملی از زمان دقیق وقوع رویداد برای برخی نمونه‌ها وجود ندارد. سانسورها می‌توانند انواع مختلفی داشته باشند:

- سانسور راست

زمانی که رویداد بعد از پایان مطالعه رخ دهد یا به عبارتی دیگر رویداد تا انتهای مطالعه رخ نمی‌دهد.

- سانسور چپ

زمانی که رویداد قبل از شروع مطالعه رخ داده باشد.

- سانسور فاصله‌ای

زمانی که رویداد بین دو بازه زمانی مشخص رخ داده باشد.

تعریف ۴.۲. ویژگی وابستگی به مدت زمان به این معنی است که احتمال پایان پاکسازی حادثه به طول مدت زمانی که پاکسازی به طول انجامیده است بستگی دارد. بدین منظور که در این نوع داده‌ها متغیر وابسته، زمان تا رویداد است. این زمان همواره مثبت است و در مطالعات تحلیل بقا می‌تواند به صورت پیوسته یا گسسته در نظر گرفته شود.

تعریف ۵.۲. مساحت زیر منحنی تابع بقا، میانگین طول عمر گفته می‌شود.

$$E(T) = \int_0^{\infty} S(t) dt$$

۳ مدل‌های تحلیل بقا

در آمار برای مدل‌بندی داده‌های زمان تا رویداد از مدل‌های تحلیل بقا استفاده می‌شود این مدل‌ها شامل مدل‌های نرخ خطر مینا هستند. مدل‌های نرخ خطر مینا قادر هستند که ویژگی وابستگی به مدت زمان را در نظر بگیرند. این مدل‌ها می‌توانند دو نوع متغیر وابسته

از جمله میانگین زمان رویداد و تابع نرخ خطر را در نظر بگیرند. مدل‌های نرخ خطر مبنا که متغیر وابسته آن‌ها تابع نرخ خطر است شامل مدل‌های خطرات متناسب کاکس^۵ است. همچنین مدل‌هایی که متغیر وابسته آن‌ها میانگین زمان بقا است شامل مدل‌های زمان شکست شتابیده (AFT) می‌باشد. در حالی که هدف مدل خطرات متناسب کاکس برآورد اثرات متغیرهای مستقل بر توابع نرخ خطر یا بقا است، هدف مدل زمان شکست شتابیده برآورد مستقیم اثرات بر روی مدت زمان (یا لگاریتم زمان) است.

• مدل خطرات متناسب کاکس

$$h(t, x) = h_0(t) \times \exp(\beta_1 x_1 + \dots + \beta_p x_p)$$

در مدل کاکس، $h_0(t)$ تابع نرخ خطر پایه است و نشان می‌دهد زمانی که هیچ‌یک از متغیرهای مستقل بر پاسخ تأثیر نداشته باشند، میزان خطر وقوع رویداد چقدر است. همچنین، برای $x_i = 1, \dots, p$ ویژگی‌هایی هستند که بر نرخ خطر تأثیر می‌گذارند و β_i نیز میزان تأثیر هر یک از آن‌ها را مشخص می‌کند. بردار β شامل ضرایب مربوط به متغیرهای مستقل، و \mathbf{X} برداری از این متغیرهای مستقل یا همان ویژگی‌ها است.

این مدل‌ها یک خطر متناسب را فرض می‌کنند که نشان می‌دهد اثر متغیرهای مستقل ثابت از زمان است. منظور از فرض خطرات متناسب این است که نسبت مخاطره برای واحد امⁱ به واحد j ام مستقل از زمان باشد. نرخ خطر پایه مقدار خطر را زمانی که همه متغیرهای مستقل صفر هستند، نشان می‌دهد. خطر به عنوان یک تابع ضربی از یک خطر پایه و یک تابعی از متغیرهای مستقل است. مدل کاکس یک مدل نیمه پارامتری است ولی می‌تواند پارامتری هم باشد.

مدل زمان شکست شتابیده این مدل یک رابطه خطی بین لگاریتم زمان بقا و متغیرهای مستقل را فرض می‌کند.

$$\ln t = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$$

متغیر وابسته این مدل را می‌توان تابع نرخ خطر نیز در نظر گرفت:

$$h(t, x) = \exp(-\beta \mathbf{X}) \cdot h_0(\exp(-\beta \mathbf{X}) \times t)$$

می‌توان توزیع‌های فرضی وایبول، لگ نرمال و لگ لجستیک را در نظر گرفت که در نرخ خطر پایه تأثیرگذار هستند.

مدل‌های شکنندگی

در برخی مطالعات تحلیل بقا، ممکن است زیر گروه‌های جامعه دارای یک متغیر تصادفی مشترک مشاهده نشده باشند که قادر به وارد کردن آن‌ها در مطالعه نیستیم. اگر اثر این متغیرها در نظر گرفته نشود بخشی از تغییرات تابع نرخ خطر که می‌توانست با وجود این عامل توضیح داده شود کاهش می‌یابد و با جمله خطا جمع شده و در نتیجه باعث افزایش تغییرات تابع نرخ خطر نسبت به حالتی که عامل مورد نظر در مدل وجود داشته باشد، می‌گردد که باعث به دست آمدن برآوردهای اریب برای پارامترهای مدل رگرسیونی می‌گردد. رویکرد پارامتر تصادفی (عامل شکنندگی) روشی است که به طور گسترده برای پرداختن به ناهمگنی مشاهده نشده پذیرفته شده است به طوری که به برخی از پارامترها اجازه می‌دهد تا در مشاهدات با توزیع فرضی، تغییر کنند.

$$\beta_n = \beta + \delta_n$$

⁵Cox proportional hazard Model

این عامل شکنندگی را می‌توان هم در مدل‌های خطرات متناسب به کار برد و هم در مدل‌های زمان شکست شتابیده از آن استفاده کرد. اگر این عامل تصادفی در مدل‌های زمان شکست شتابیده به کار رود به نام مدل‌های (Random parameter AFT) شناخته شده و پیش‌بینی و تحلیل دقیق‌تری از عوامل تاثیرگذار مرتبط با ناهمگنی مشاهده نشده ارائه می‌کنند.

۴ مدل‌های یادگیری ماشین

یکی دیگر از روش‌های رایج برای مدل‌سازی داده‌های زمان پاکسازی حادثه، مدل‌های یادگیری ماشین می‌باشد. از جمله می‌توان به مدل‌های شبکه‌های عصبی مصنوعی (ANN)، جنگل تصادفی (RF) و ماشین‌های گرادیان تقویتی (GBM) اشاره کرد. مدل‌های شبکه‌های عصبی عمیق (DNN) یک روند اجتناب ناپذیر در زمینه‌های مختلف تحقیقاتی می‌باشد.

۱.۴ جنگل بقای تصادفی

جنگل تصادفی^۶ یکی از الگوریتم‌های یادگیری ماشین مبتنی بر درخت تصمیم است که برای مسائل رگرسیون و طبقه‌بندی استفاده می‌شود. همچنین اگر مدل جنگل تصادفی با هدف برآورد تابع بقا طراحی شود^۷ می‌تواند برای داده‌های زمان تا رویداد مناسب باشد زیرا ویژگی وابستگی به مدت زمان را در این داده‌ها در نظر می‌گیرد. در این نوع الگوریتم‌ها از چندین مدل ضعیف استفاده می‌شود تا عملکرد مدل کلی را بهبود ببخشد. در مدل جنگل تصادفی از چندین درخت تصمیم استفاده می‌شود. همان‌طور که در جنگل تصادفی با هدف رگرسیون مدل‌های ضعیف درخت تصمیم به صورت تصادفی و مستقل آموزش می‌بینند در مدل جنگل بقای تصادفی (RSF) هم درخت‌ها به طور جداگانه با هدف برآورد تابع بقا آموزش می‌بینند. در مدل جنگل بقای تصادفی هر درخت تصمیم به جای پیش‌بینی یک مقدار عددی در گره آخر، تابع بقا را به روش‌های ناپارامتری مثل کاپلان مه‌یر برآورد می‌کند. پس از آموزش درخت‌ها، برای هر نمونه جدید پیش‌بینی تابع بقا به صورت میانگین از پیش‌بینی تمام درخت‌ها محاسبه می‌شود.

الگوریتم ۱. فرض کنید داده‌ها به صورت زیر باشند:

$$\mathcal{D} = \{(X_i, T_i, \delta_i)\}_{i=1}^n$$

که در آن: X_i : بردار ویژگی‌ها، T_i : زمان مشاهده‌شده (کمترین مقدار بین زمان وقوع و سانسور)، و $\delta_i \in \{0, 1\}$: نشان‌گر وقوع رویداد است. (۱=رخ داده، ۰=سانسور شده)

• ساخت درخت‌ها

۱. یک نمونه بوت‌استرپ از داده‌ها بگیرید.
۲. در هر گره، m ویژگی به‌طور تصادفی انتخاب کنید.
۳. برای هر تقسیم، آماره آزمون لاگ-رتبه را محاسبه کنید.
۴. تقسیم با بیشترین مقدار آماره لاگ-رتبه را انتخاب کنید.

^۶Random Forest

^۷Random Survival Forest

۵. تقسیم گره‌ها تا رسیدن به شرط توقف ادامه می‌یابد.

- برآورد در گره نهایی

در هر گره نهایی، تابع بقا با استفاده از برآوردگر کاپلان-مایر محاسبه می‌شود:

$$\hat{S}(t) = \prod_{j: t_j \leq t} \left(1 - \frac{d_j}{r_j}\right)$$

که: d_j تعداد رخدادها در زمان t_j و r_j تعداد افراد در معرض خطر در زمان t_j ، می‌باشد.

- پیش‌بینی نهایی جنگل

میانگین برآوردهای بقا از تمام درخت‌ها به صورت زیر محاسبه می‌شود:

$$\hat{S}(t | X) = \frac{1}{B} \sum_{b=1}^B \hat{S}^{(b)}(t | X)$$

۲.۴ ماشین‌گرادیان تقویت‌شده

ماشین‌گرادیان تقویتی^۸ یکی از الگوریتم‌های یادگیری جمعی است. این مدل هم بر مبنا چندین درخت تصمیم است. تفاوت بین مدل‌های Adaboost و ماشین‌گرادیان تقویتی در نحوه به‌روزرسانی مدل و نحوه تعامل مدل‌های ضعیف است. Adaboost مدل‌های ضعیف درخت تصمیم را به صورت متوالی تقویت می‌کند و در هر مرحله به مواردی که اشتباه پیش‌بینی شده‌اند وزن‌هایی اختصاص می‌دهد تا در مدل بعدی بیشتر به آن‌ها توجه شود. در مدل ماشین‌گرادیان تقویتی به همه نمونه‌ها وزن‌های یکسانی اختصاص داده می‌شود و با استفاده از روش عددی گرادیان نزولی پیش‌بینی‌ها بهبود می‌بخشد. به عبارت دیگر در هر مرحله در مدل ماشین‌گرادیان تقویتی هدف بهبود عملکرد کلی مدل در پیش‌بینی است.

۳.۴ شبکه عصبی مصنوعی

شبکه‌های عصبی^۹ شامل سه نوع لایه از جمله لایه ورودی، لایه خروجی و لایه‌های پنهان هستند. تعداد نورون‌ها در لایه ورودی برابر با تعداد ویژگی‌های داده‌های ورودی است و تعداد نورون‌ها در لایه خروجی بستگی به هدف پیش‌بینی دارد. تعداد نورون‌های لایه پنهان بین تعداد نورون‌های لایه ورودی و تعداد نورون‌های لایه خروجی قرار می‌گیرد. به عبارتی، تعداد نورون‌ها در لایه پنهان ابرپارامتری است که با آزمایش مقادیر مختلف تعیین می‌شود تا بهترین عملکرد آموزشی شبکه عصبی به دست آید.

آموزش شبکه عصبی شامل دو مرحله از پیش‌انتشار و پس‌انتشار است. پیش‌انتشار فرآیندی است که در آن داده‌های ورودی به لایه ورودی ارسال می‌شود، توسط لایه پنهان پردازش می‌شود و در نهایت به لایه خروجی منتقل می‌گردد. جزئیات پیش‌انتشار به صورت زیر است:

⁸Gradient Boosting Machine

⁹Neural Network

اگر شبکه عصبی با یک لایه پنهان و یک لایه خروجی داشته باشیم، رابطه بین لایه ورودی و لایه پنهان به صورت ریاضی به صورت زیر بیان می شود:

$$neth_1 = \sum w_i x_i + b_1$$

که در آن x_i بردار ورودی است، w_j بردارهای وزن بین لایه ورودی و لایه پنهان هستند، و b_1 بایاس بین لایه های محاسبه شده است که به طور انعطاف پذیری برای تنظیم عملکرد شبکه عصبی طراحی شده است. سپس، $neth_1$ از توابع فعال سازی از جمله ReLu عبور می کند. تابع فعال سازی ReLu به صورت $\max(0, neth_1)$ است.

$$h_i = g(neth_1) = \max(0, neth_1)$$

h_i مقادیر گره های لایه پنهان هستند که دوباره در وزن های اتصالات بین لایه پنهان و لایه خروجی w_j ضرب شده و با مقدار بایاس هر گره لایه خروجی جمع می شود. در نهایت از تابع فعال سازی لایه آخر $g(\cdot)$ عبور کرده و مقدار خروجی هر گره در لایه خروجی که همان Z است، به دست می آید. Z به صورت زیر محاسبه می شود:

$$Net_z = \sum w_j h_i + b_2$$

و

$$Z = g(Net_z) = \max(0, Net_z)$$

که در آن w_j بردارهای وزن بین لایه پنهان و لایه خروجی هستند و b_2 بایاس بین لایه پنهان و لایه خروجی است. اگر Z با خروجی مورد انتظار مطابقت نداشته باشد، خطای پیش بینی با استفاده از تابع ضرر انتخاب شده محاسبه خواهد شد. سپس، پس انتشار خطای پیش بینی را به نورون های شبکه باز می گرداند و از الگوریتم نزول گرادینت برای اصلاح وزن های نورون ها به منظور کاهش خطاهای پیش بینی استفاده می کند. هنگامی که خطای پیش بینی با انتظارات مطابقت داشت، پیش بینی شبکه عصبی انجام می شود.

۵ مدل تحلیل بقای عمیق چند وظیفه ای

در بحث شبکه های عصبی اگر چند وظیفه با هم مرتبط بودند می توانیم یک مدل شبکه عصبی طراحی کنیم که همزمان قادر به انجام دو وظیفه باشد. مدل مورد استفاده در این مطالعه یک شبکه عصبی عمیق (منظور تعداد زیاد لایه های پنهان) دو وظیفه ای است که می تواند ویژگی وابستگی به مدت زمان را با استفاده از اطلاعات مشترکی که از برآورد تابع بقا به دست می آورد در پیش بینی زمان پاک سازی حادثه در نظر بگیرد. در یادگیری چند وظیفه ای، چندین کار مرتبط به طور همزمان یاد می گیرند و از اطلاعات مشترک بین آن ها برای بهبود عملکرد کلی استفاده می شود. یادگیری چند وظیفه ای به جای آموزش یک مدل جداگانه برای هر کار، یک مدل واحد را برای انجام چندین کار آموزش می دهد. با ورود یک بردار متغیرهای مستقل به مدل، چندین خروجی (هدف یا متغیرهای وابسته) به دست می آوریم. در

یادگیری چند وظیفه‌ای، برخی از لایه‌ها یا پارامترها در بین وظایف به اشتراک گذاشته می‌شوند و به مدل اجازه می‌دهند ویژگی‌های مشترکی را که برای همه وظایف مفید هستند، یاد بگیرد. این مدل به طور همزمان روی وظایف مختلف آموزش داده می‌شود و لایه‌های مشترک بر اساس زیان ترکیبی از همه وظایف به روز می‌شوند. علاوه بر لایه‌های مشترک، مدل‌های یادگیری چند وظیفه‌ای معمولاً دارای لایه‌های اختصاصی مربوط به وظایف هستند که جنبه‌های منحصر به فرد هر وظیفه را مدیریت می‌کنند. از جمله مزایای یادگیری چند وظیفه‌ای این است که در صورت مرتبط بودن وظایف، عملکرد آنها با اطلاعات مشترک هم‌دیگر بهبود می‌بخشد و مدل یادگیری چند وظیفه‌ای می‌تواند به عنوان یک تنظیم‌کننده عمل کند و از توجه بیش از حد مدل در یک وظیفه جلوگیری کند. همچنین از معایب ساختار یادگیری چند وظیفه‌ای می‌توان گفت با افزایش تعداد وظایف پیچیدگی و هزینه محاسباتی یادگیری چند وظیفه‌ای می‌تواند به طور قابل توجهی افزایش یابد. نتایج نشان می‌دهند که انجام وظایف چندگانه به طور مکمل عملکرد وظایف مختلف را افزایش می‌دهد اگر اطلاعات نهفته‌ای که از مولفه‌های مشترک استخراج می‌شود به دو وظیفه مشابه کمک کند ساختار یادگیری چند وظیفه‌ای می‌تواند به طور موثر اطلاعات پنهان مشترک را جذب کند. وظیفه پیش‌بینی میانه زمان پاسخگویی بخشی از برآورد تابع بقا است بنابراین دو وظیفه مرتبط هستند. در این مطالعه به دنبال انجام دو وظیفه یعنی برآورد تابع بقا و پیش‌بینی میانگین زمان پاکسازی حادثه به طور همزمان هستیم. وظیفه برآورد تابع بقا مربوط به مدل‌های خطرات متناسب کاکس است. همچنین پیش‌بینی میانگین زمان پاکسازی حادثه وظیفه مدل‌های زمان شکست شتابیده است.

برای لایه اختصاصی مربوط به وظیفه برآورد تابع بقا، این مطالعه ۱۰ دقیقه را به عنوان فاصله زمانی برای برآورد تابع بقا با توجه به کاربرد آن در مدیریت حادثه در نظر گرفته است. تابع بقا یک احتمال است، احتمال بقا در هر کدام از بازه‌های ۱۰ دقیقه‌ای محاسبه می‌شود. برای این منظور تابع فعالسازی لایه آخر سافت مکس^{۱۰} در نظر گرفته می‌شود. برای وظیفه برآورد تابع بقا دو تابع زیان در نظر گرفته شده است:

• تابع زیان لگاریتم درستنمایی

در داده‌های زمان تا رویداد داده‌های سانسور شده وجود دارد. این داده‌ها اطلاعات کاملی درباره چگالی به ما نمی‌دهند و بنابراین در تابع درستنمایی باید به همان میزانی که به ما اطلاعات می‌دهند اثر داشته باشند بنابراین باید تابع ماسکی تعریف شود تا داده‌های سانسور شده را از بدون سانسور تشخیص دهد.

$$L_1 = - \sum_{t=1}^{n_{out}} y_t \cdot \log \hat{y}_t$$

n_{out} تعداد بازه‌هایی است که تابع بقا در آن بازه‌ها برآورد می‌شود. y_t مقدار پاسخ واقعی است که از داده‌های مشاهده شده به دست می‌آید و میزان وقوع رویداد برای یک حادثه در زمان t ام است و \hat{y}_t احتمال برآورد شده رخ دادن پاکسازی حادثه در بازه زمانی t ام است که خروجی مدلی است که وظیفه برآورد تابع بقا را انجام می‌دهد.

• تابع زیان رتبه بندی^{۱۱}

این تابع زیان با استفاده از نرخ خطر تجمعی ترتیب وقوع رویدادها را بررسی می‌کند به طوری که حادثه‌هایی که رویداد برای آن‌ها زودتر اتفاق می‌افتد احتمال بقا کمتر و نرخ خطر بیشتری نسبت به حادثه‌هایی دارند که همان رویداد را دیرتر تجربه کردند. اگر مدلی بتواند این ترتیب‌ها را درست پیش‌بینی کند مثل این است که توانسته تابع بقا را دقیق برآورد کند. این تابع زیان برای

¹⁰Softmax

¹¹Rank Loss

بررسی این است که آیا مدل ترتیب درست وقوع رویدادها را پیش‌بینی کرده است یا خیر. هر چه این زیان کمتر باشد یعنی برآورد تابع بقا دقیق‌تر است. با توجه به عملکرد این تابع زیان باید وقتی که می‌خواهیم این زیان را محاسبه کنیم تابع ماسکی تعریف کنیم تا برای داده‌هایی که در طی زمان اندازه‌گیری شدند هم بتواند ترتیب درست رویدادها را به درستی اعمال کند. این تابع زیان مانند شاخص تطابق وابسته به زمان (C-Index) عمل می‌کند.

$$L_2 = \sum_{i \neq j} A_{i,j} \eta(\hat{H}(t_i|x_{t_i}), \hat{H}(t_j|x_{t_j}))$$

که در آن:

$$\eta(x, y) = \exp\left(\frac{y - x}{\sigma}\right)$$

و

$$A_{i,j} = \begin{cases} 1 & \text{اگر } t_i < t_j \\ 0 & \text{در غیر این صورت} \end{cases}$$

در لایه آخر شبکه عصبی مربوط به وظیفه پیش‌بینی میانگین زمان پاکسازی حادثه، از تابع فعالسازی خطی استفاده شده است. تابع زیان مربوط به آموزش این وظیفه مجموع میانگین مربعات خطا (MSE) است.

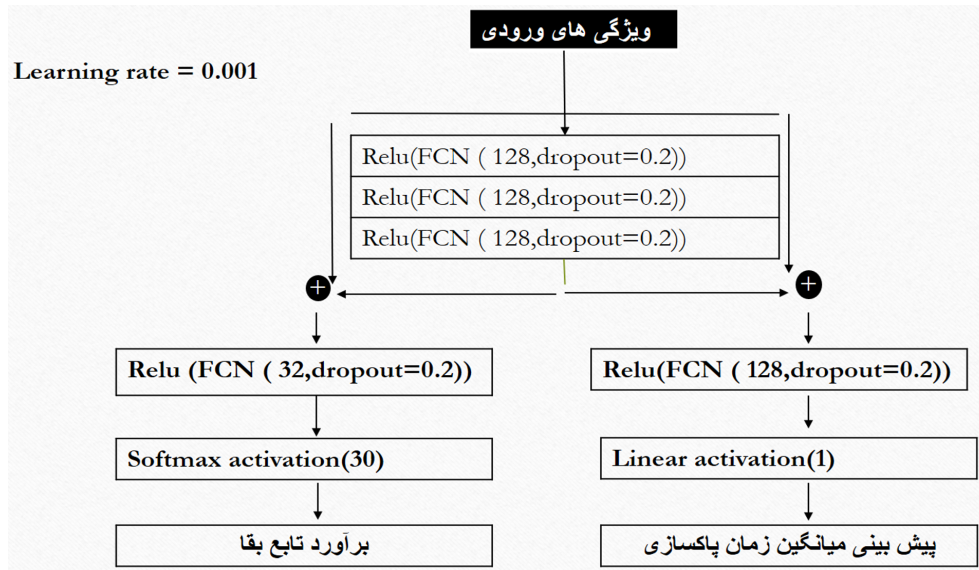
$$MSE = \frac{1}{K} \sum_{k=1}^K (T_k - \hat{T}_k)^2$$

در این فرمول K تعداد حوادث مشاهده شده است. T_k و \hat{T}_k به ترتیب مدت زمان واقعی پاکسازی حادثه و مدت زمان پاکسازی تخمین‌زده شده توسط مدل می‌باشد.

با تابع زیان ترکیبی زیر لایه‌های مشترک آموزش داده می‌شود:

$$L_{total} = L_1 + w_1 \cdot L_2 + w_2 \cdot L_3$$

وزن‌های هر تابع زیان جزو ابر پارامترهایی است که با استفاده از روش‌های جستجوی شبکه‌ای کالیبره می‌شوند. از الگوریتم بهینه سازی Adam برای بهینه کردن وزن‌های شبکه عصبی استفاده شده است. نرخ یادگیری در این مدل 0.001 در نظر گرفته شده است. ساختار شبیه‌سازی مدل در شکل ۱ نشان داده شده است. در شکل ۱ بین دو وظیفه، لایه‌های مشترک وجود دارند که ویژگی‌های پنهان مشترک بین دو وظیفه را استخراج می‌کنند. تابع فعالسازی لایه‌های مشترک ReLu می‌باشد و برای هر لایه پنهان در لایه‌های مشترک لایه Dropout در نظر گرفته شده است که روشی برای جلوگیری از بیش‌برازش مدل است. در این روش به طور تصادفی ۲۰ درصد از نورون‌های هر لایه در محاسبات خروجی غیرفعال می‌شود. لایه‌های مشترک با تابع زیان ترکیبی L_{total} آموزش می‌بیند و ویژگی‌های مشترک بین دو وظیفه را استخراج می‌کند. تعداد گره‌ها در هر لایه ۱۲۸ نورون است و لایه‌ها کاملاً به هم متصل هستند. در یادگیری چند وظیفه‌ای باید وظایف باهم مرتبط باشند تا از ویژگی‌های مشترک‌شان برای بهبود عملکرد هر وظیفه بهره برد. بعد از لایه‌های مشترک بین وظایف به وظایف اختصاصی می‌رسیم. شبکه عصبی اول وظیفه برآورد تابع بقا را بر عهده دارد و پارامترهای این شبکه با تابع زیان L_1



شکل ۱: ساختار مدل تحلیل بقای عمیق چند وظیفه‌ای

و L_2 آموزش می‌بیند و تابع فعالسازی لایه آخر این مدل تابع سافت مکس است که برای محاسبه احتمال وقوع پاکسازی حادثه در هر بازه t است. شبکه عصبی دوم با وظیفه پیش‌بینی زمان پاکسازی است که پارامترهای این مدل نیز با تابع زیان L_3 آموزش می‌بیند و تابع فعالسازی لایه آخر این شبکه تابع خطی است.

۶ مثال کاربردی

داده‌های حوادث ترافیکی شامل ۳۹۷۶۰ حادثه با ویژگی‌های زمانی، جغرافیایی و حادثه است. داده‌ها مربوط به تصادف آزادراه‌های کشور کره جنوبی از سال ۲۰۱۴ تا ۲۰۱۹ می‌باشد. به طور تصادفی ۷۰ درصد از داده‌ها برای آموزش و اعتبارسنجی مدل و ۳۰ درصد از داده‌ها برای آزمون مدل استفاده شده است. اغلب برخی ویژگی‌های گمشده یا غیر واقعی در مجموعه داده‌ها وجود دارد که با روش‌های پیش‌پردازش داده‌ها برطرف می‌شوند. فقط از داده‌های کامل حادثه استفاده شده است که شامل تمام ویژگی‌های زمانی، جغرافیایی حادثه و جریان ترافیک است با حذف نقاط پرت در نهایت ۱۶۹۲۹ داده برای توسعه مدل به کار گرفته شده است. در این پژوهش همه داده‌ها کامل هستند و سانسور وجود ندارد. تمام متغیرهای طبقه بندی در طول آموزش مدل به متغیرهای باینری تبدیل می‌شوند و متغیرهای مربوط به ترافیک پیوسته هستند و استاندارد می‌شوند. در این مثال، مدل زمان شکست شتابیده را با در نظر گرفتن توزیع‌های فرضی مختلف برای نرخ خطر پایه، به کار می‌گیرد.

مدل زمان شکست شتابیده (AFT) (توزیع وایبول)
مدل زمان شکست شتابیده (AFT) (توزیع لگ نرمال)
مدل زمان شکست شتابیده (AFT) (توزیع لگ لجستیک)

در این مطالعه مدل زمان شکست شتابیده با یک پارامتر تصادفی (RPAFT) برای در نظر گرفتن ناهمگنی مشاهده نشده به کار گرفته شد. از مدل (GBM) و (RF) با دو هدف پیش‌بینی میانگین زمان پاکسازی و برآورد تابع بقا برای بررسی عملکرد مدل هدف

استفاده شده است. ابرپارامترهای اصلی با استفاده از تکنیک (K-fold) تنظیم شده اند. در این روش هر ترکیب از ابر پارامترها ارزیابی می‌شوند و خطای هر ترکیب به دست می‌آید. ترکیبی که کمترین خطا را در معیارهای میانگین خطای قدر مطلق (MAE) و میانگین درصد خطای قدر مطلق (MAPE) و شاخص تطابق وابسته به زمان (c-td) داشته باشد، انتخاب می‌شود.

برای جنگل تصادفی (RF) چهار ابرپارامتر از جمله تعداد درخت، حداکثر عمق هر درخت، نرخ نمونه‌گیری برای آموزش هر درخت و تعداد متغیرهای مورد استفاده در هر گره تقسیم وجود دارد. برای ماشین‌گردان تقویتی (GBM) نیز چهار ابر پارامتر از جمله تعداد درخت، حداکثر عمق هر درخت، نرخ یادگیری و تعداد متغیرهای در نظر گرفته‌شده در هر درخت وجود دارد.

مدل هدف، جداگانه با هر کدام از مدل‌های تک وظیفه‌ای شبکه عصبی عمیق (DNN) مقایسه می‌شود تا نشان‌دهنده بهبود عملکرد یادگیری چند وظیفه‌ای باشد. ترکیب ابر پارامترها از جمله تعداد لایه‌ها و نورون‌های هر شبکه، نرخ یادگیری، نرخ dropout و وزن‌ها برای هر زیان به طور گسترده تنظیم شده‌اند تا معماری مدل نهایی را انتخاب کنند. عملکرد مدل برای پیش‌بینی میانگین مدت زمان رویداد با میانگین خطای قدرمطلق (MAE) و میانگین درصد خطای قدرمطلق (MAPE) ارزیابی می‌شود.

$$MAE = \frac{1}{K} \sum_{k=1}^K |T_k - \hat{T}_k|$$

$$MAPE = \frac{1}{K} \sum_{k=1}^K \left| \frac{T_k - \hat{T}_k}{T_k} \right| \times 100$$

عملکرد برای وظیفه برآورد توزیع زمان بقا (تابع بقا) با شاخص C-td ارزیابی می‌شود:

$$C - td = \frac{\sum_{i \neq j} A_{i,j} \cdot B_{i,j}}{\sum_{i \neq j} A_{i,j}}$$

که در آن:

$$A_{i,j} = \begin{cases} 1 & \text{اگر } t_i < t_j \\ 0 & \text{در غیر این صورت} \end{cases}$$

و

$$B_{i,j} = \begin{cases} 1 & \text{اگر } \hat{H}(t_i | (X_{I,t_i}, X_{T,t_i})) > \hat{H}(t_j | (X_{I,t_j}, X_{T,t_j})) \\ 0 & \text{در غیر این صورت} \end{cases}$$

این شاخص توانایی مدل را در پیش‌بینی صحیح ترتیب وقوع رویدادها، بر اساس این فرض اندازه‌گیری می‌کند که رویدادهایی که مدت بیشتری طول کشیده‌اند، باید مخاطره کمتری نسبت به رویدادهایی داشته باشند که مدت کوتاه‌تری دوام آورده‌اند. شاخص تطابق وابسته به زمان تمام بازه‌های زمانی مخاطره تجمعی را ارزیابی می‌کند، بنابراین می‌تواند مواردی را که اثرات عوامل تاثیرگذار در طول زمان تغییر می‌کنند (وابسته به زمان هستند) را نیز ارزیابی کند. این ویژگی به مدل اجازه می‌دهد تا در شرایط پویا و پیچیده‌تر، مانند شبکه‌های ترافیکی که متغیرهای مستقل ممکن است در زمان‌های مختلف تأثیرات متفاوتی داشته باشند، عملکرد بهتری داشته باشد.

بحث و نتیجه‌گیری

جداول زیر عملکرد پیش‌بینی مدل هدف را در مقایسه با سایر مدل‌ها نشان می‌دهد.

• مدل‌های آماری

مدل‌های آماری	MAE	MAPE	AIC
مدل زمان شکست شتابیده (AFT) (توزیع وایبول)	۱۸.۸۶۵	۵۴.۴۶۵	۲۲۵۵۲.۸۴۴
مدل زمان شکست شتابیده (AFT) (توزیع لگ نرمال)	۱۸.۳۳۵	۵۱.۳۹۱	۱۸۶۱۳.۵۳۶
مدل زمان شکست شتابیده (AFT) (توزیع لگ لجستیک)	۱۸.۴۶۱	۵۲.۱۹۸	۱۸۷۰۰.۴۲۰
مدل شکنندگی (RPAFT) (توزیع لگ لجستیک)	۱۷.۱۱۴	۴۲.۰۶۶	۱۸۲۷۳.۵۸۳

از بین مدل‌های آماری فوق مدل زمان شکست شتابیده با توزیع لگ لجستیک از مدل‌های دیگر AFT از نظر MAE و MAPE و AIC بهتر عمل می‌کند. همچنین مدل شکنندگی با یک پارامتر تصادفی RPAFT بهتر از مدل زمان شکست شتابیده با توزیع فرضی لگ لجستیک عمل می‌کند. این نشان‌دهنده وجود ناهمگنی مشاهده نشده در بین داده‌ها می‌باشد.

• مدل‌های یادگیری ماشین

مدل‌های یادگیری ماشین	MAE	MAPE	ctd
RF-Survival	۱۷.۱۱۱	۴۰.۴۹۱	۰.۶۰۲
RF-Regression	۱۸.۱۷۸	۵۳.۵۵۸	
GBM-Survival	۱۶.۷۶۸	۴۰.۵۷۸	۰.۶۰۷
GBM-Regression	۱۶.۷۹۶	۴۲.۰۰۵	

همه مدل‌های یادگیری ماشین و شبکه‌های عصبی عمیق عملکرد بهتری نسبت به مدل زمان شکست شتابیده با توزیع لگ لجستیک دارند و به این معنی که در نظر گرفتن یک رابطه تعاملی و غیرخطی عملکرد پیش‌بینی را بهبود می‌بخشد. مدل GBM-Regression و DNN-Regression و مدل هدف بهتر از مدل شکنندگی RPAFT عمل می‌کنند. عملکرد هردو مدل جنگل تصادفی (RF) و مدل ماشین‌گرادیان تقویتی (GBM) با تغییر وظایف‌شان از پیش‌بینی میانگین زمان پاکسازی به برآورد تابع بقا، بهبود یافته است.

• مدل‌های شبکه‌های عصبی عمیق

مدل‌های شبکه‌های عصبی عمیق	MAE	MAPE	ctd
DNN-Survival	۱۷.۵۵۳	۴۳.۱۹۵	۰.۶۲۴
DNN-Regression	۱۶.۸۰۸	۴۱.۴۵۲	

برخلاف مدل‌های یادگیری ماشین، شبکه عصبی عمیق با هدف برآورد تابع بقا عملکرد پیش‌بینی کننده بهتری نسبت به شبکه عصبی عمیق با هدف پیش‌بینی میانگین زمان پاکسازی حادثه، ندارد.

• مدل هدف

مدل هدف	MAE	MAPE	ctd
Multi task Deep Survival Analysis Model	۱۶.۵۵۱	۳۹.۶۰۹	۰.۶۳۷

برخلاف مدل‌های دیگر که به طور جداگانه وظایف پیش‌بینی مدت زمان حادثه و برآورد تابع بقا را انجام می‌دهند مدل هدف این وظایف را به طور همزمان انجام می‌دهد. مدل هدف بهترین عملکرد را هم برای پیش‌بینی میانگین زمان پاکسازی (از نظر MAE و MAPE) و هم برای برآورد تابع بقا (از نظر c-td) نشان می‌دهد. برای پیاده‌سازی مدل روی داده‌ها از نرم افزار پایتون نسخه ۸.۳ استفاده شده است.

قدردانی و تشکر

نویسندگان از داوران محترم به خاطر دقت نظر و صرف وقت در مطالعه مقاله و ارائه پیشنهادات برای بهبود آن، تقدیر و تشکر می‌کنند.

مراجع

- Nam, D., Mannering, F. (2000). *An exploratory hazard based analysis of highway incident duration*. Transportation Research Part A: Policy and Practice, 34(2), 85-102.
- Jones, B., Janssen, L., and Mannering, F. (1991). *Analysis of the frequency and duration of freeway accidents in Seattle*. Accident, Analysis and Prevention, 23(4), 239–255. [https://doi.org/10.1016/0001-4575\(91\)90003-N](https://doi.org/10.1016/0001-4575(91)90003-N)
- Chung, Y. (2010). *Development of an accident duration prediction model on the Korean Freeway Systems*. Accident; Analysis and Prevention, 42(1), 282–289. <https://doi.org/10.1016/j.aap.2009.08.005>
- Tavassoli Hojati, A., Ferreira, L., Washington, S., and Charles, P. (2013). *Hazard based models for freeway traffic incident duration*. Accident; Analysis and Prevention, 52, 171–181. <https://doi.org/10.1016/j.aap.2012.12.037>
- Wei, C., Lee, Y. (2007). *Sequential forecast of incident duration using artificial neural networks models*. Accident; Analysis and Prevention, 39(5), 944–954. <https://doi.org/10.1016/j.aap.2006.12.017>
- Hamad, K., Al-Ruzouq, R., Zeiada, W., Abu Dabous, S., and Khalil, M. A. (2020). *Predicting incident duration using random forests*. Transportmetrica A: Transport Science, 16(3), 1269–1293. <https://doi.org/10.1080/23249935.2020.1733132>
- Ma, X., Ding, C., Luan, S., Wang, Y., Wang, Y. (2017). *Prioritizing influential factors for freeway incident clearance time prediction using the gradient boosting decision trees method*. IEEE Transactions on Intelligent Transportation Systems, 18(9), 2303–2310. <https://doi.org/10.1109/TITS.2016.2635719>

- Li, L., Sheng, X., Du, B., Wang, Y., and Ran, B. (2020). *A deep fusion model based on restricted Boltzmann machines for traffic accident duration prediction. Engineering Applications of Artificial Intelligence*, 93(April), 103686. <https://doi.org/10.1016/j.engappai.2020.103686>
- Tang, J., Zheng, L., Han, C., Yin, W., Zhang, Y., Zou, Y., and Huang, H. (2020). *Statistical and machine-learning methods for clearance time prediction of road incidents: A methodology review. Analytic Methods in Accident Research*, 27, 100123. <https://doi.org/10.1016/j.amar.2020.100123>
- Ishwaran, H., Kogalur, U. B., Blackstone, E. H., and Lauer, M. S. (2008). *Random survival forests. The Annals of Applied Statistics*, 2(3), 841–860. <https://doi.org/10.1214/08-AOAS169>
- Hothorn, T., Buhlmann, P., Dudoit, S., Molinaro, A., and Van Der Laan, M. J. (2006). *Survival ensembles. Biostatistics (Oxford, England)*, 7(3), 355–373. <https://doi.org/10.1093/biostatistics/kxj011>
- Katzman, J. L., Shaham, U., Cloninger, A., Bates, J., Jiang, T., and Kluger, Y. (2018). *DeepSurv: Personalized treatment recommender system using a Cox proportional hazards deep neural network. BMC Medical Research Methodology*, 18(1), 24. <https://doi.org/10.1186/s12874-018-0482-1>
- Lee, C., Zame, W. R., Yoon, J., and Van Der Schaar, M. (2018). *DeepHit: A deep learning approach to survival analysis with competing risks. 32nd AAAI Conference on Artificial Intelligence, AAAI 2018*, 2314–2321. <https://doi.org/10.1609/aaai.v32i1.11842>



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30-31, 2025
University of Kurdistan



مدل‌بندی تصادفات جاده‌ای استان خراسان جنوبی با استفاده از مدل‌های اتورگرسیو صحیح مقدار

زهره نخعی زاده^۱، سارا جمهوری^۲

^{۱،۲} گروه آمار، دانشگاه بیرجند

چکیده: در علم آمار، به دنباله‌ای از مشاهدات که به طور معمول در فواصل زمانی منظم ثبت می‌شوند، "سری زمانی" گفته می‌شود. بسیاری از متغیرها، مانند تعداد تصادفات، قربانیان جرایم و سایر موارد مشابه، ماهیتی گسسته دارند. بنابراین، تحلیل و پیش‌بینی این نوع داده‌ها با استفاده از مدل‌های سری زمانی شمارشی، مناسب‌تر از دیگر روش‌ها است. در این مقاله، از مدل‌های سری زمانی شمارشی برای بررسی داده‌های تصادفات روزانه برون‌شهری در شهرستان‌های استان خراسان جنوبی استفاده شده است. هدف از این پژوهش، شناسایی مدل مناسبی است که بهترین برازش را به داده‌ها داشته باشد. پس از مقایسه مدل‌های مختلف، مدلی که بهترین تطابق را با داده‌های موجود نشان می‌دهد، انتخاب شده و بر اساس آن پیش‌بینی‌های لازم انجام می‌شود.

واژه‌های کلیدی: پیش‌بینی، عملگر نازک‌کننده، سری زمانی شمارشی.

کد موضوع‌بندی ریاضی (۲۰۲۰): 62M10.

۱ مقدمه

تصادفات جاده‌ای یکی از مهم‌ترین و چالش‌برانگیزترین مسائل در حوزه حمل و نقل مدرن محسوب می‌شوند. این تصادفات بر اثر عوامل مختلفی از جمله افزایش جمعیت، رشد تقاضای سفر و بنابراین استفاده بیشتر از وسیله‌های نقلیه موتوری رخ می‌دهند. در نتیجه، ایمنی جاده‌ها و خیابان‌ها اهمیت ویژه‌ای پیدا می‌کند، چرا که کاهش تصادفات و جلوگیری از وقوع آن‌ها هدف اصلی در ارتقاء ایمنی حمل و نقل است.

در بسیاری از کشورها، از جمله ایران، تصادفات جاده‌ای به عنوان یکی از اصلی‌ترین دلایل مرگ‌ومیر و جراحات تشخیص داده شده است. در نتیجه توجه به این امر یکی از اندیشه‌های اصلی در برنامه‌ریزی‌های حمل و نقل و سیاست‌های ایمنی است. بر اساس آمارهای

^۱ سخنران، znakhaezadeh72@gmail.com

موجود، ایران یکی از کشورهایی است که با بالاترین نرخ مرگ‌ومیر ناشی از تصادفات جاده‌ای مواجه است. تحلیل آمارهای تصادفات می‌تواند به مدیران و کارشناسان این حوزه کمک کند تا روندهای تصادفات را بهتر درک کرده و اقدامات موثری برای کاهش آن‌ها انجام دهند.

برای پیش‌بینی تعداد تصادفات در مناطق خاص و در بازه‌های زمانی مختلف، از مدل‌های پیش‌بینی و سری‌های زمانی استفاده می‌شود. در سال‌های اخیر، تکنیک‌های رگرسیونی که بر فرض استقلال داده‌ها تکیه دارند، کمتر مورد استفاده قرار می‌گیرند، چرا که محدودیت‌هایی در مدل‌سازی سری‌های زمانی پیچیده دارند. در این راستا، محققان در قالب مدل‌های سری زمانی پیوسته مقدار، همچون مدل‌های $ARMA$ ، $SARIMA$ و $ARIMA$ ، داده‌های مربوط به تصادفات را مورد مطالعه قرار داده‌اند و نتایج مفیدی در حوزه پیش‌بینی و تحلیل روندهای تصادفات ارائه داده‌اند. مطالعات متعددی در قلمرو مدل‌های سری زمانی انجام شده است که از این میان می‌توان به یوسف‌زاده چابک و همکاران (۲۰۱۶)، آگینگ و همکاران (۲۰۲۳)، نصیری و همکاران (۲۰۲۳) و ولوگدین و همکاران (۲۰۲۴) و از میان پژوهش‌های داخلی به توکلی و رحیم‌اف (۱۳۹۰)، عسگری و همکاران (۱۳۹۴) و امید و همکاران (۱۳۹۶) اشاره کرد.

به دلیل ماهیت گسسته و غیرمنفی داده‌های تصادفات، مدل‌های سنتی و پیوسته مقدار سری‌های زمانی برای تحلیل این داده‌ها مناسب نیستند. از این‌رو، توسعه مدل‌هایی که به طور خاص برای این نوع داده‌ها طراحی شده‌اند، ضروری است. یکی از ساده‌ترین و در عین حال پرکاربردترین مدل‌ها برای داده‌های شمارشی، مدل $INAR(1)$ معرفی شده توسط مک‌کنزی (۱۹۸۵) است. این مدل نسخه‌ای صحیح مقدار از مدل خودرگرسیو مرتبه اول $AR(1)$ است و با وجود ساختار ساده، از انعطاف‌پذیری بالایی برخوردار است. مدل $INAR(1)$ نه تنها امکان تحلیل فرآیندهای شمارشی ایستا را فراهم می‌کند، بلکه مبنایی برای توسعه مدل‌های پیشرفته‌تر نیز محسوب می‌شود. همچنین، مفاهیم مهمی مانند برآورد پارامترها، تحلیل مانده‌ها و استنباط آماری را می‌توان به سادگی در قالب این مدل پیاده‌سازی کرد. اساس ساخت مدل‌های $INAR(1)$ مبتنی بر عملگر نازک‌کننده^۱ است. با در نظر گرفتن عملگرهای نازک‌کننده متفاوت، مدل‌های سری زمانی متنوعی به دست می‌آیند. در پژوهش حاضر از سری‌های زمانی شمارشی برای مدل‌سازی داده‌های تصادفات جاده‌ای در سطح استان خراسان جنوبی بهره گرفته شده است.

۲ مدل‌های شمارشی

در این بخش، خانواده مدل‌های سری زمانی شمارشی $INAR(1)$ و خواص آن مورد بررسی قرار گرفته است. ابتدا مفهوم عملگر نازک‌کننده دوجمله‌ای^۲ بیان می‌شود.

تعریف ۱.۲. فرض کنید X یک متغیر تصادفی گسسته و $\{Z_i, i = 1, 2, \dots\}$ دنباله‌ای از متغیرهای تصادفی مستقل و هم‌توزیع با توزیع مشترک برنولی و احتمال موفقیت $(0, 1)$ باشند. در این صورت توزیع شرطی متغیر تصادفی Z_i به $\alpha \circ X = \sum_{i=1}^X Z_i$ به شرط X دوجمله‌ای $Bin(X, \alpha)$ است. عملگر \circ عملگر نازک‌کننده دوجمله‌ای نامیده می‌شود.

ال‌اوش و الزید (۱۹۸۷) مدل $INAR(1)$ را با استفاده از عملگر دوجمله‌ای به صورت زیر معرفی کردند.

^۱Thinning operator

^۲Binomial thinning

تعریف ۲.۲. فرآیند تصادفی شمارشی $\{X_t, t \in \mathcal{T}\}$ یک فرآیند $INAR(1)$ نامیده می شود هرگاه

$$X_t = \alpha \circ X_{t-1} + \varepsilon_t, \quad t \in \mathbb{Z}, \quad \alpha \in (0, 1), \quad (1.2)$$

که در آن جملات نوآوری^۳ $\{\varepsilon_t; t \in \mathcal{T}\}$ دنباله ای از متغیرهای تصادفی نامنفی، صحیح مقدار، مستقل و هم توزیع با میانگین μ_ε و واریانس σ_ε^2 است. همچنین فرض می شود که متغیرهای تصادفی موجود در عملگر نازک کننده α و ε_t از $\{X_t; t \in \mathcal{T}\}$ مستقل هستند.

ال اوش و الزید (۱۹۸۷) نشان دادند که فرآیند $INAR(1)$ ، یک زنجیر مارکوف همگن با احتمال انتقال یک مرحله ای به صورت

زیر است

$$\begin{aligned} p_{k|l} &:= P(X_t = k | X_{t-1} = l) \\ &= \sum_{j=0}^{\min\{k, l\}} \binom{l}{j} \alpha^j (1 - \alpha)^{l-j} P(\varepsilon_t = k - j). \end{aligned} \quad (2.2)$$

توجه ۳.۲. با انتخاب توزیع های مختلفی برای جملات نوآوری، مدل های $INAR(1)$ مختلفی حاصل می شود. یکی از مدل های پرکاربرد، حالتی است که در آن توزیع حاشیه ای X_t پواسون با میانگین $\mu = \frac{\lambda}{(1-\alpha)}$ است و $\{\varepsilon_t; t \in \mathcal{T}\}$ دنباله ای از متغیرهای تصادفی پواسون با میانگین و واریانس λ می باشد.

توجه ۴.۲. با وجود اینکه توزیع پواسون یکی از مهم ترین توزیع های آماری مرتبط با فرآیندهای شمارشی است، اما استفاده از آن همیشه در مدل های سری زمانی شمارشی مناسب نیست. زیرا در این توزیع میانگین و واریانس با یکدیگر برابرند و این خاصیت عملاً در سری های زمانی مشاهده نمی شود. به همین علت، مدل های دیگری مورد مطالعه قرار گرفته است. اگر در (۱.۲) جملات نوآوری از توزیع دو جمله ای منفی $NB(n, \pi)$ پیروی کنند، مدل $NBINAR(1)$ حاصل می شود.

توجه ۵.۲. در برخی مدل های $INAR(1)$ فرض می شود که ضریب نازک کننده متغیری تصادفی است. در این صورت مدل به صورت $X_t = \alpha_\phi \circ X_{t-1} + \varepsilon_t$ نوشته می شود. **جو (۱۹۹۶)** با در نظر گرفتن $\phi = \frac{1}{(n+1)}$ و توزیع بتا $Beta(\frac{1-\phi}{\phi}\alpha, \frac{1-\phi}{\phi}(1-\alpha))$ برای α_ϕ وقتی جملات نوآوری دارای توزیع $NB(n(1-\alpha), \pi)$ باشند، مدل $RNBINAR(1)$ را معرفی کرد. در این مدل، توزیع حاشیه ای X_t دو جمله ای منفی $NB(n, \pi)$ خواهد بود.

۱.۲ برآورد پارامترهای مدل

روش های متعددی مانند روش گشتاورها، برآورد حداقل مربعات شرطی و برآورد ماکزیمم درستنمایی برای تخمین پارامترهای مدل وجود دارد. برای انجام برآورد پارامترها بر پایه روش درستنمایی، می توان از احتمال های شرطی و خاصیت مارکوفی سری زمانی $INAR(1)$ بهره برد. تابع لگاریتم درستنمایی در مدل $INAR(1)$ با بردار پارامتر θ عبارت است از

$$\ell(\theta) = \log p_{x_1}(\theta) + \sum_{t=2}^T \log p_{x_t|x_{t-1}}(\theta), \quad (3.2)$$

که در آن احتمالات انتقال با توجه به رابطه (۲.۲) محاسبه می شوند. پارامتر θ در مدل های مختلف $INAR(1)$ متفاوت است. به عنوان مثال در مدل پواسون $INAR(1)$ که در ملاحظه ۳.۲ معرفی شد $\theta = (\alpha, \lambda)$ و در مدل $NBINAR(1)$ معرفی شده در ملاحظه ۴.۲، $\theta = (\alpha, n, \pi)$ است. در این صورت برآوردگر ماکسیمم درستنمایی پارامتر θ به صورت $\hat{\theta} = \operatorname{argmax}_{\theta} \ell(\theta)$ می باشد.

³Innovations

۲.۲ تشخیص و کفایت مدل

در این بخش، ابزارهای مورد نیاز برای تشخیص مناسب‌ترین مدل برای برازش به داده‌ها معرفی می‌شوند. برای تشخیص ساختار وابستگی سری، نمودار خودهمبستگی جزئی نمونه‌ای ($SPACF$) رسم می‌شود. اگر $\hat{\rho}_{par}(k)$ در تاخیر اول مخالف صفر و در بقیه تاخیرها به سمت صفر میل کند، مدل $INAR(1)$ برای برازش به داده‌ها مناسب خواهد بود. برای تعیین مدل $INAR(1)$ مناسب برای داده‌ها، می‌توان از شاخص $\hat{I} = \frac{S^2}{\bar{X}}$ استفاده کرد که در آن $\bar{X} = \frac{1}{T} \sum_{t=1}^T X_t$ و $S^2 = \frac{1}{T} \sum_{t=1}^T (X_t - \bar{X})^2$. اگر مدل پواسون برای برازش به داده‌ها مناسب باشد، در اینصورت باید مقدار این آماره تقریباً برابر یک باشد. همچنین این آماره به طور مجانبی از توزیع نرمال با

$$E(\hat{I}) = 1 - \frac{1}{T} \frac{1 + \alpha}{1 - \alpha}, \quad V[\hat{I}] \approx \frac{2}{T} \frac{1 + \alpha^2}{1 - \alpha^2},$$

پیروی می‌کند (شوئر و ویب، ۲۰۱۴). اگر چندین مدل برای برازش به داده‌ها مناسب باشند، از معیارهایی مانند آکائیک و معیار اطلاع بیزی که به صورت زیر هستند

$$AIC = -2\ell(\theta) + 2m, \quad BIC = -2\ell(\theta) + 2 \log T,$$

برای شناسایی مدل مناسب‌تر استفاده می‌شود (m تعداد پارامترهای مجهول مدل و T تعداد مشاهدات است). بعد از شناسایی بهترین مدل، لازم است که کفایت مدل بررسی شود. محاسبه باقی‌مانده‌های استاندارد شده پیرسون اولین رویکرد در کفایت مدل است. این مقادیر طبق رابطه زیر حاصل می‌شود.

$$e_t = \frac{x_t - E(X_t | x_{t-1})}{\sqrt{V(X_t | x_{t-1})}}, \quad t = 2, 3, \dots, T.$$

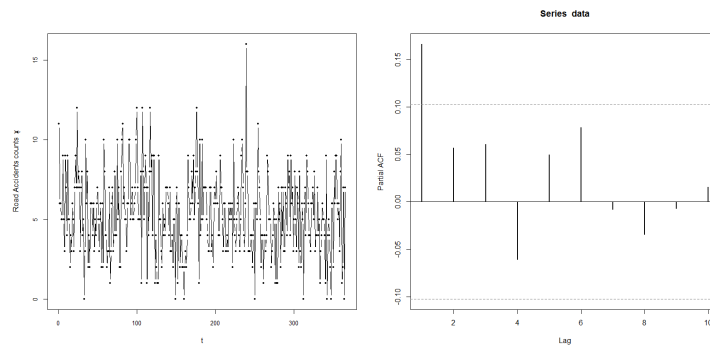
برای یک مدل مناسب انتظار می‌رود که این باقی‌مانده‌ها دارای میانگین و واریانس به ترتیب برابر صفر و یک بوده و با یکدیگر همبستگی نداشته باشد. واریانس بزرگتر (کوچکتر) از یک نشان دهنده وجود بیش‌پراکندگی (کم‌پراکندگی) نسبت به مدل در نظر گرفته شده است (هاروی و فرناندز، ۱۹۸۹).

۳ پیش‌بینی تصادفات جاده‌ای

در این پژوهش، ابتدا آمار و ارقام مربوط به تصادفات جاده‌ای^۴ در سه سال متوالی ۱۴۰۱، ۱۴۰۲ و ۱۴۰۳ گزارش شده است. این گزارش نشان‌دهنده وضعیت تصادفات، میزان تلفات، محل وقوع و زمان‌بندی آن‌ها است که به منظور شناسایی روندها و ارائه راهکارهای بهبود امنیت جاده‌ای اهمیت فراوان دارد.

در سال ۱۴۰۱، تعداد کل تصادفات ثبت‌شده به ۲۰۵۲ مورد رسید که از این تعداد، ۲۵۸ مورد مرگ و میر ناشی از تصادف گزارش شده است. همچنین، حدود ۳۴ درصد از این تصادفات در جاده‌های روستایی، ۶۰ درصد در جاده‌های اصلی و ۶ درصد در جاده‌های فرعی رخ داده است. در این سال، ۶۴ درصد تصادفات در طول روز و ۳۶ درصد در شب اتفاق افتاده است، که نشان‌دهنده اهمیت زمان‌بندی و دقت در رانندگی در ساعات مختلف شبانه‌روز است. در سال ۱۴۰۲، تعداد تصادفات کمی کاهش یافته و به ۱۹۸۶ مورد رسید. در این سال، تعداد ۲۴۸ فوتی به ثبت رسیده که نسبت به سال قبل کمی کاهش یافته است. همچنین ۵۷ درصد تصادفات در جاده‌های اصلی، ۳۶ درصد در جاده‌های روستایی و ۷ درصد در جاده‌های فرعی رخ داده است. الگوی زمانی همچنان برتری تصادفات

^۴ داده‌ها از پلیس راه استان خراسان جنوبی اخذ شده است.



شکل ۱: تعداد تصادفات روزانه جاده‌ای استان خراسان جنوبی و نمودار ACF داده‌ها

در روز را نشان می‌دهد (۶۳ درصد در روز و ۳۷ درصد در شب). در سال ۱۴۰۳، تعداد تصادفات دوباره کمی کاهش یافته و به ۱۹۵۳ مورد رسیده است. تعداد فوتی‌ها نیز به ۲۰۰ تن کاهش یافته است، که نشان‌دهنده نسبتی پایدار است. در این سال، توزیع جغرافیایی تصادفات نسبت به سال‌های قبل تقریباً ثابت است (۵۷ درصد در جاده‌های اصلی، ۳۵ درصد در جاده‌های روستایی و ۸ درصد در جاده‌های فرعی). همچنین، الگوهای زمانی همچنان بر تصادفات روز تأکید دارند (۶۴ درصد در طول روز و ۳۶ درصد در شب).

بر اساس این آمار، می‌توان نتیجه گرفت که تعداد تصادفات جاده‌ای در طی سه سال مورد بررسی کاهش نسبی یافته است؛ اما با توجه به میزان تلفات، نیاز به تمرکز بیشتر و اتخاذ اقدامات مؤثر در حوزه ایمنی جاده‌ها و آموزش رانندگان احساس می‌شود. همچنین، سهم قابل توجه تصادفات در جاده‌های روستایی و در ساعات روز نشان می‌دهد که برنامه‌ریزی برای ارتقاء زیرساخت‌ها و اقدامات مراقبتی در این فضاها اهمیت دارد. به طور کلی، تمرکز بر کاهش تصادفات در جاده‌های روستایی و کنترل سرعت و رعایت قوانین در زمان‌های مختلف شبانه‌روز می‌تواند نقش مؤثری در بهبود وضعیت ایمنی جاده‌ای کشور ایفا کند. در این بخش مجموعه داده‌های تصادفات روزانه جاده‌ای در استان خراسان جنوبی در سال ۱۴۰۳ مورد بررسی قرار می‌گیرد ($T = ۳۶۵$). کد نویسی‌ها به کمک نرم افزار R انجام شده است. شکل (a) ۱، نمودار سری زمانی تعداد تصادفات را نشان می‌دهد. مشاهدات در فاصله صفر تا ۱۶ تصادف در روز قرار دارند. روند یا تغییرات فصلی قابل توجهی در داده‌ها مشاهده نمی‌شود. نوسانات بالا و پایین موجود در این نمودار بیان‌کننده سطح متوسطی از همبستگی است. این موضوع از نمودار ACF نمونه‌ای نیز مشهود است. با توجه به مقدار $\hat{\rho}_{par}(1) = 0.1662$ ملاحظه می‌شود که این مقدار به‌طور قابل توجهی از صفر بزرگتر است. بنابراین، مدلی شبیه به $AR(1)$ برای توصیف مشاهدات مناسب است. به وضوح، مقدار میانگین $\bar{x} \approx 5.3205$ کمتر از واریانس $s^2 \approx 6.6579$ است. لذا، در داده‌ها پراکندگی زیادی وجود دارد. از طرف دیگر، مقدار احتمال برآورد شده در نقطه صفر برابر $p_0 \approx 0.2191$ است که بزرگتر از مقدار متناظر آن از توزیع پواسون، یعنی $exp(-\bar{x}) \approx 0.0048$ می‌باشد. بنابراین، برازش یک مدل $INAR(1)$ به داده‌ها، قابل قبول به نظر می‌رسد. برای انجام آزمون تشخیص بیش‌پراکندگی در مدل $INAR(1)$ ، فرضیه صفر اینکه داده‌ها از مدل $INAR(1)$ با نوآوری‌های پواسون ^۵ پیروی می‌کنند، آزمون می‌شود. تحت فرضیه صفر، $E(\hat{I}) = 0.9961$ و $V(\hat{I}) = 0.0700$ است. با توجه به مقدار $p - value = 0.0004$ مشخص می‌شود که پراکندگی بیش از حد و احتمالاً قابل توجه است. بنابراین مدل $RCNBINAR(1)$ برای برازش به داده‌ها مناسب‌تر است. این مدل با مدل پواسون $INAR(1)$ از نظر معیارهای AIC و BIC مقایسه می‌شود. برآورد میانگین و شاخص پراکندگی عوامل نوآوری‌ها به ترتیب $\bar{\mu}_\varepsilon = \bar{x}(1 - \hat{\alpha}) \approx 4.4358$ و $\hat{I} = \hat{I}(1 + \hat{\alpha}) - \hat{\alpha} \approx 1.2891$ است. از بررسی مقادیر AIC و BIC موجود در جدول

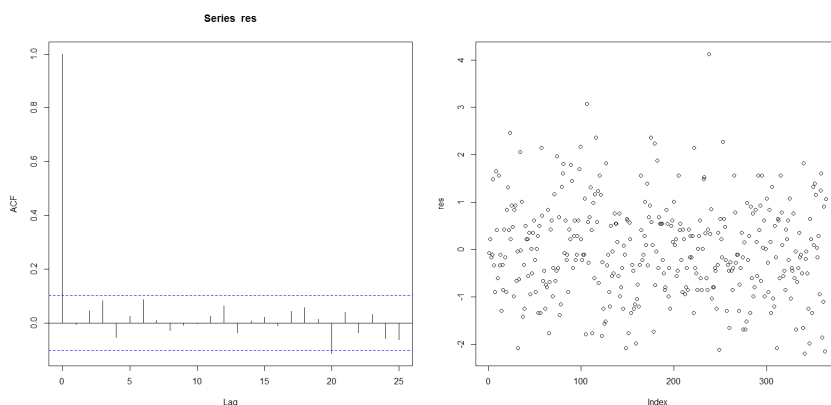
جدول ۱: آمار توصیفی داده‌های تصادفات در شهرهای استان خراسان جنوبی در سال ۱۴۰۳

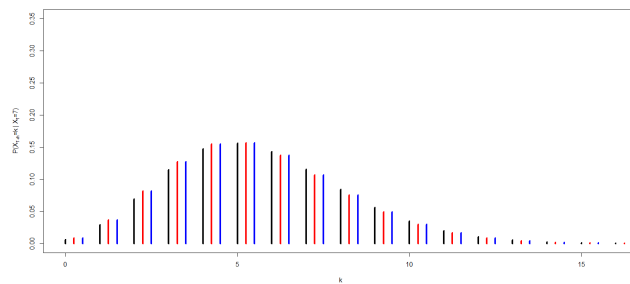
متغیر	کشیگی	چولگی	ماکزیم	مینم	انحراف معیار	میانه	میانگین
تصادفات سال ۱۴۰۳	۳۳۱۹	۰/۳۴۹	۱۶	۰	۲/۵۸۰	۵	۵/۳۲۱

جدول ۲: برآوردهای درستنمایی ماکسیم پارامترها و AIC و BIC در مدل‌های مختلف

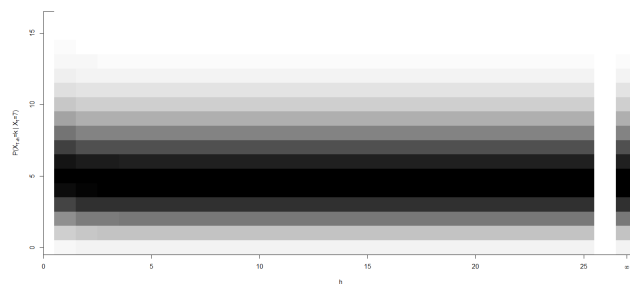
Model	Parameter			AIC	BIC
	۱	۲	۳		
$i.i.d.Poisson$	۵/۳۲۰			۱۷۳۵/۱۰۷	۱۷۳۹/۰۰۷
μ	۰/۱۲۰				
$Poisson.INAR(1)$	۴/۶۰۴	۰/۱۳۵		۱۷۲۶/۹۱۶	۱۷۳۴/۷۱۶
(λ, α)	(۰/۲۴۵)	(۰/۰۴۱)			
$i.i.d.NB$	۱۹/۷۲۶	۰/۷۸۵		۱۷۲۶/۸۵۱	۱۷۳۴/۶۵۱
(n, π)	(۷/۲۴۵)	(۰/۰۶۱)			
$NBINAR(1)$	۱۵/۳۴۰	۰/۷۷۳	۰/۱۵۵	۱۷۱۹/۳۲۹	۱۷۳۱/۰۲۸
(n, π, α)	(۶/۳۷۸)	(۰/۰۶۹)	(۰/۰۴۹)		
$RCNBINAR(1)$	۲۱/۴۵۹	۰/۰۸۰۱	۰/۱۵۸	۱۷۱۹/۲۳۰	۱۷۳۱/۰۳۰
(n, π, α)	(۹/۳۱۶)	(۰/۰۸۰۱)	(۰/۱۵۸)		

۲ مشخص می‌شود که مدل‌های $INAR(1)$ بهتر از مدل $i.i.d$ عمل می‌کنند. مقدار برآورد شده α به طور قابل توجهی از صفر بزرگتر است. از مقایسه مدل‌های $INAR(1)$ برحسب معیارهای AIC و BIC ، ملاحظه می‌شود که مدل $RCNBINAR(1)$ برای برازش به داده‌ها مناسب‌تر است. برای بررسی کفایت این مدل، نمودارهای پراکنش و ACF نمونه‌ای باقیمانده‌های پیرسون در شکل ۲ رسم شده‌اند. با توجه به این شکل، هیچگونه همبستگی معنی‌داری بین باقیمانده‌ها مشاهده نمی‌شود. برای این مدل، توزیع‌های پیش‌بینی h -گام روبه جلو به شرط مقدار آخرین مشاهده ($x_{۳۶۵} = ۷$)، با استفاده از رابطه (۲.۲) به دست می‌آیند. نمودار این توزیع‌ها در شکل ۳ و ۴ رسم شده است. با توجه به شکل ۳ پیش‌بینی مبتنی بر میانه توزیع به ازای سه مقدار $h = ۱$ (مشکی)، $h = ۵$ (قرمز) و $h = ۱۰$ (آبی) در نقطه $k = ۵$ اتفاق می‌افتد. شکل ۴ نیز مقدار احتمال بیشتر را روی مقادیر ۴، ۵ و ۶ نشان می‌دهد. لذا انتظار می‌رود که در روزهای آتی، به طور متوسط ۵ تصادف در روز اتفاق بیاقتد.

شکل ۲: نمودارهای پراکنش و ACF نمونه‌ای باقیمانده‌های پیرسون



شکل ۳: هیستوگرام توزیع پیش‌بینی h گام به جلو برای $h = 1, 5, 10$



شکل ۴: توزیع پیش‌بینی h گام به جلو برای $h = 1, \dots, 25, \infty$

بحث و نتیجه‌گیری

در این مقاله، با استفاده از مدل‌های سری زمانی شمارشی، تعداد تصادفات رانندگی جاده‌ای در سطح استان خراسان جنوبی، مورد مطالعه قرار گرفت. پس از تحلیل، مشخص گردید که مدل $RCNBINAR(1)$ برای برازش به داده‌ها، مناسب‌تر از سایر مدل‌های رقیب است.

مراجع

- امیدی، نبی، خلیلی، کرم، امید، محمد رضا. و جعفری اسکندری، میثم. (۱۳۹۶). ارائه مدل اقتصادسنجی پیش‌بینی کشته شدگان و مصدومان تصادفات جاده‌ای، مطالعات مدیریت ترافیک، شماره ۴۴، ۲۱-۴۴.
- توکلی شیرازی، نیما و رحیم اف، کامران. (۱۳۹۰). مدلسازی سری زمانی مجروحین و متوفیات حوادث تصادفات جاده‌ای ایران با رویکردهای $ARMA$ ، AR ، MA ، دومین کنفرانس ملی تصادفات جاده‌ای، سوانح ریلی و هوایی، زنجان.
- عسگری، حشمت اله، امید، محمد رضا، امید، نبی و مفتاحی، هادی. (۱۳۹۴). پیش‌بینی بروز تخلفات از نوع سرعت غیرمجاز در آزادراه‌ها، مطالعات مدیریت ترافیک. شماره ۳۸، ۲۵-۳۸.

Agyemang, E. F., Mensah, J. A., Ocran, E., Opoku, E. Nortey, E. N. N. (2023). Time series based road traffic

- accidents forecasting via SARIMA and Facebook Prophet model with potential changepoints. *Heliyon*, 9, e22544.
- Al-Osh, M. A., Alzaid, A. A. (1987). First-order integer-valued autoregressive (INAR(1)) process. *Journal of Time Series Analysis*, 8(3), 261-275.
- Freeland, R. K. (1998). Statistical analysis of discrete time series with application to the analysis of workers' compensation claims data (Doctoral dissertation, University of British Columbia).
- Harvey, A.C., Fernandes, C. (1989). Time series models for count or qualitative observations. *Journal of Business & Economic Statistics*, 7(4), 407-417.
- Joe, H. (1996). Time series models with univariate margins in the convolution-closed infinitely divisible class. *Journal of Applied Probability*, 33(3), 664-677.
- McKenzie, E. (1985). Some simple models for discrete variate time series. *Water Resources Bulletin*, 21(4), 645-650.
- Nasiri, H., Mohammadpour, S. I., Dahaghin, M. (2023). Forecasting time trend of road traffic crashes in Iran using the macro-scale traffic flow characteristics. *Heliyon*, 9(3), e14481.
- Schweer, S., Weiß, C.H. (2014). Compound Poisson INAR(1) processes: stochastic properties and testing for overdispersion. *Computational Statistics and Data Analysis*, 77, 267-284.
- Vologdin, S., Kasatkina, E., Kasatkin, A. (2024). Analyzing and forecasting road traffic accidents and their consequences: a case study of the Udmurt republic. *E3S Web of Conferences*, 471, 06006.
- Yousefzadeh-Chabok, S., Ranjbar-Taklimie, F., Malekpouri, R., Razzaghi, A. (2016). A Time Series Model for Assessing the Trend and Forecasting the Road Traffic Accident Mortality. *Arch Trauma Res*, 5(3), e36570.

Modeling Road Traffic Accidents in South Khorasan Province Using Integer-Valued Autoregressive Models

Nakhaeezadeh, Z.¹, Jomhoori, S.²

^{1,2}Department of Statistics, University of Birjand, Iran.

Abstract: In statistics, a sequence of observations that are typically recorded at regular time intervals is referred to as a "time series". Many variables, such as the number of accidents, crime victims, and other similar events, exhibit a discrete nature. Therefore, the analysis and forecasting of these types of data are more appropriately conducted using count time series models. In this article, count time series models are applied to analyze daily intercity accident data from the counties of South Khorasan province. These data span a three-year period, from 2021 to 2024. The objective of this research is to identify the most suitable model that provides the best fit to the data. After comparing various models, the one that demonstrates the best match with the existing data is selected, and based on it, the necessary forecasts are made.

Keywords: Counting timeseries , Forecasting, Thinning operator.

Mathematics Subject Classification (2020): 62M10.



پانزدهمین سمینار احتمال
و فرآیندهای تصادفی

۸ و ۹ شهریور ۱۴۰۴
دانشگاه کردستان



Seminar On Probability
and Stochastic Processes

August 30-31, 2025
University of Kurdistan



فرایند چوب‌شکنی رگسیون-بتا با کوواریانس نامانا وابسته به متغیر کمکی

اسماعیل یارعلی^۱، فیروزه ریواز^۲، مجید جعفری خالدی^۳

^۱ دانش‌آموخته دکتری آمار، دانشگاه شهید بهشتی، تهران

^۲ گروه آمار، دانشگاه شهید بهشتی، تهران

^۳ گروه آمار، دانشگاه تربیت مدرس، تهران

چکیده: در این مقاله رویکردی جدیدی برای ساخت فرایند چوب‌شکنی^۱ مورد مطالعه قرار می‌گیرد به‌طوری‌که توزیع تصادفی وابسته به متغیر کمکی است. این شیوه مدل‌بندی علاوه بر آن‌که نامانایی متأثر از متغیر کمکی را در ساختار وابستگی داده‌ها را پوشش می‌دهد، امکان افزایش ناهمبستگی ناحیه مورد مطالعه را نیز فراهم می‌کند. با در نظر گرفتن این فرایند به عنوان توزیع اثرهای تصادفی در یک چارچوب مدل خطی آمیخته، داده‌ها مدل شده و استنباط بیزی آن ارائه می‌شود. سپس با استفاده از یک مجموعه داده‌ی واقعی مربوط به میزان بارش عملکرد مدل معرفی شده در مقایسه با مدل‌های رقیب مورد ارزیابی قرار می‌گیرد.

واژه‌های کلیدی: فرایند چوب‌شکنی، رگسیون-بتا، نامانایی، متغیر کمکی، ساختار همبستگی.

کد موضوع‌بندی ریاضی (۲۰۲۰): 62F15, 62M30, 62H11.

۱ مقدمه

یکی از مؤلفه‌های اصلی تحلیل داده‌های فضایی، تعیین وابستگی فضایی است. این کار عموماً از طریق تابع کوواریانس که متعلق به رده‌ای پارامتری از مدل‌های مانا است، انجام می‌شود. این درحالی است که در بسیاری از کاربردها، فرض مانایی واقع‌بینانه نیست. در واقع، در بسیاری از کاربردها عوامل موضعی مانند توپوگرافیکی منطقه مورد مطالعه بر همبستگی فضایی تأثیرگذار هستند و باعث ایجاد نامانایی در ساختار وابستگی فضایی می‌شوند. لذا ساخت مدل‌های کوواریانس نامانای متأثر از متغیرهای کمکی موضعی بسیار حائز اهمیت است. تاکنون رهیافت‌های مختلفی برای الحاق متغیر کمکی در ساختار مرتبه دوم فرایندهای فضایی ارائه شده است.

^۱ سخنران، e_yarali@sbu.ac.ir

^۱ Stick-breaking proces

از جمله می‌توان به روش هموارسازی هسته (ریچ و همکاران، ۲۰۱۱)، روش پیچش فرایند (رایزر و کالدر، ۲۰۱۵)، روش دگرپسی (اشمیت و همکاران، ۲۰۱۱)، روش بسط توابع پایه (یارعلی و ریواز، ۲۰۲۰) و استفاده از معادلات مشتقات جزئی تصادفی (اینگریستن و همکاران، ۲۰۱۴) اشاره کرد. گرچه الحاق متغیر تصادفی در ساختار وابستگی باعث بهبود در مدل‌های فضایی هم به لحاظ کاهش بعد تعداد پارامتر و به دنبال آن محاسبات ساده و هم تفسیر راحت‌تر نامانایی شده است، اما در تمامی این روش‌ها برای الحاق متغیر کمکی در ساختار وابستگی فضایی فرض شده است میدان تصادفی مورد مطالعه گاوسی است، در حالی‌که چنین فرضی در عمل به سختی برقرار بوده و داده‌ها شواهدی از تخطی مفروضات توزیع گاوسی را از خود نشان می‌دهند. برای رفع این محدودیت‌ها، در سال‌های اخیر مدل‌های فضایی ناگاوسی متعددی از جمله مدل‌های چوله گاوسی، رهیافت پیچش فرایند گاوسی-لگ گاوسی و فرایندهای گاوسی تبدیل یافته g و h توکی معرفی شده است. علی‌رغم ویژگی‌های مناسب این مدل‌ها، ساختار همبستگی القاء شده توسط آنها مبتنی بر فرض محدودکننده مانایی و همسانگردی میدان تصادفی مورد مطالعه است. درحالی‌که در برخی از کاربردها ارتباط ساختار فضایی با متغیر / متغیرهای کمکی موضعی باعث ایجاد نامانایی در ساختار وابستگی داده‌ها می‌شود.

یک رویکرد مناسب استفاده از رهیافت بیزی ناپارامتری است که در آن توزیع داده‌ها نامعلوم و تصادفی در نظر گرفته می‌شود و برای آن یک پیشینی مناسب انتخاب می‌شود. در این رویکرد، معمولاً از پیشینی فرایند چوب شکنی استفاده می‌شود از جمله می‌توان به فرایندهای چوب شکنی هسته (ریچ و فونتنیس، ۲۰۰۷؛ برزگر و ریواز، ۲۰۲۰؛ دحدوح و خالدی، ۲۰۲۰)، فرایندهای چوب شکنی پروبیت (رودریگوئز و دانسون، ۲۰۱۱)، فرایندهای چوب شکنی لوحیت (رن و همکاران، ۲۰۱۱)، فرایندهای چوب شکنی تعمیم یافته (بارسلا و همکاران، ۲۰۱۸) و فرایندهای چوب شکنی پنهان (رودریگوئز و همکاران، ۲۰۱۰) اشاره کرد. علی‌رغم ایجاد انعطاف بالا در مدل‌های معرفی شده، ساختار همبستگی القا شده از این مدل‌ها بسیار پیچیده و وابستگی فضایی وابسته به متغیر کمکی نیست. برای رفع این مشکل یارعلی و همکاران (۲۰۲۲) با استفاده از رهیافت بیزی ناپارامتری، رده‌ی جدیدی از فرایند چوب شکنی را معرفی کردند به‌طوری‌که نامانایی از طریق متغیر کمکی در همبستگی فضایی القا می‌شود. این شیوه مدل‌بندی علاوه بر آن‌که رفتارهای توزیعی شامل عدم تقارن و چند مدی بودن را مدل‌بندی می‌کند، نامانایی متأثر از متغیر کمکی را در ساختار کوواریانس پوشش می‌دهد. همچنین افزایش برداری القا شده توسط این پیشینی پیشنهادی وابسته به متغیر کمکی بوده که امکان افزایش برداری ناحیه‌ی مورد مطالعه را بر اساس اطلاعات متغیر کمکی فراهم می‌کند براین اساس در بخش ۲ فرایند چوب شکنی رگرسیون-بتا که توسط یارعلی و همکاران (۲۰۲۲) معرفی شد، بیان می‌شود. در ادامه، با بکارگیری یک مدل بیزی مبتنی بر فرایند چوب شکنی رگرسیون-بتا، استنباط‌های پسینی و پیشگویی فضایی ارائه می‌شوند. همچنین با استفاده از یک مجموعه داده‌ی واقعی عملکرد مدل معرفی شده با مدل‌های رقیب از نظر پیشگویی فضایی و افزایش برداری ناحیه مورد مطالعه مورد ارزیابی قرار می‌گیرد.

۲ فرایند چوب شکنی رگرسیون-بتا

در این بخش با استفاده از رهیافت بیزی ناپارامتری، فرایند فضایی چوب شکنی رگرسیون-بتا که توسط یارعلی و همکاران (۲۰۲۲) بیان شد، معرفی می‌شود. بر این اساس فرض کنید $Y(\cdot) = \{Y(s), s \in D \subset \mathcal{R}^d, d \neq 1\}$ یک میدان تصادفی باشد که روی ناحیه D تعریف شده است. به‌علاوه فرض کنید $Y(s)$ مقدار مشاهد شده در موقعیت s است که از مدل

$$Y(s) = w'(s)\beta + \eta_x(s) + \epsilon(s), \quad (1.2)$$

پیروی می‌کند که در آن $w(s)$ بردار متغیرهای کمکی در مکان s و $\beta \in \mathcal{R}^p$ بردار ضرایب رگرسیونی نامعلوم است. همچنین $\epsilon(s)$ نشان‌دهنده خطای اندازه‌گیری است که فرض می‌شود دارای توزیع نرمال با میانگین صفر و واریانس τ^2 است. اکنون برای الحاق متغیر کمکی در ساختار همبستگی فضایی، فرض می‌شود فرایند فضایی $\eta_x(s)$ به صورت

$$\eta_x(s) = G^{-1}(V_x(s)), \quad (2.2)$$

است که در آن $G(\cdot)$ اندازه احتمال تصادفی گسسته است که به صورت $G(\cdot) = \sum_{l=1}^L p_l \delta_{\theta_l}(\cdot)$ تعریف می‌شود که در آن $\delta_{\theta_l}(\cdot)$ تابع دیراک^۲، $G^{-1}(\cdot)$ وارون تعمیم‌یافته تابع $G(\cdot)$ ، و اتم‌های $\theta_1, \dots, \theta_L$ تحقق‌هایی مستقل از توزیع نرمال $N(0, \sigma^2)$ هستند. همچنین $V_x(\cdot) = \{V_x(s), s \in \mathcal{R}^d\}$ یک فرایند فضایی رگرسیون-بتا است که به صورت زیر تعریف می‌شود.

تعریف ۱.۲. فرض کنید $\{Z(s), s \in \mathcal{R}^d\}$ یک فرایند گاوسی مانا با میانگین صفر، واریانس یک و تابع همبستگی $\rho(s_i, s_j; \lambda)$ باشد. اگر $V_x(s)$ به شرط $Z(s)$ ، یک فرایند مستقل و دارای توزیع $Beta(\mu_x(s)\psi, \psi(1 - \mu_x(s)))$ با تابع چگالی

$$f(v_x) = \frac{\Gamma[\psi]}{\Gamma[\mu_x(s)\psi]\Gamma[\psi(1 - \mu_x(s))]} v_x^{\mu_x(s)\psi-1} (1 - v_x)^{\psi(1 - \mu_x(s))-1}, \quad 0 < \mu_x(s) < 1, \quad \psi > 0, \quad (3.2)$$

باشد که در آن $\mu_x(s) = \mathbb{E}[V_x(s)|Z(s)]$ به صورت

$$\mu_x(s) = \frac{\exp(x'(s)\gamma + Z(s))}{1 + \exp(x'(s)\gamma + Z(s))}, \quad (4.2)$$

تعریف می‌شود به طوری که $x(s) = (x_1(s), \dots, x_q(s))'$ برداری از متغیرهای کمکی در مکان s ، $\gamma = (\gamma_1, \dots, \gamma_q)'$ برداری از ضرایب نامعلوم که اثر متغیرهای کمکی را در همبستگی فضایی کنترل می‌کند، آنگاه $V_x(\cdot) = \{V_x(s), s \in \mathcal{R}^d\}$ را یک فرایند فضایی رگرسیون-بتا گویند.

حال، برای هر مجموعه متناهی از موقعیت‌های s_1, \dots, s_n روی ناحیه \mathcal{D} ، احتمال شرطی

$$P_x \left(\eta_x(s_1) = \theta_{l_1}, \dots, \eta_x(s_n) = \theta_{l_n} | \{v_l\}_{l=1}^L, \{\theta_l\}_{l=1}^L \right)$$

به صورت

$$\begin{aligned} & P_x \left(\eta_x(s_1) = \theta_{l_1}, \dots, \eta_x(s_n) = \theta_{l_n} | \{v_l\}_{l=1}^L, \{\theta_l\}_{l=1}^L \right) \\ &= P \left(\sum_{j=1}^{l_1-1} p_j \leq V_x(s_1) < \sum_{j=1}^{l_1} p_j; \dots; \sum_{j=1}^{l_n-1} p_j \leq V_x(s_n) < \sum_{j=1}^{l_n} p_j \right), \end{aligned} \quad (5.2)$$

تعریف می‌شود به طوری که l_1, \dots, l_n یکی از مقادیر $\{1, \dots, L\}$ را می‌گیرد. در این صورت توزیع پیشینی القاشده توسط مدل سلسله‌مراتبی بالا تحت اندازه احتمال P_x را یک فرایند چوب‌شکنی رگرسیون-بتا نامیده و با نماد $P_x \sim BRSBP(\alpha, \psi, \gamma, \lambda, \sigma^2)$ نمایش می‌دهیم.

یک رهیافت دیگر برای ساخت فرایند چوب‌شکنی رگرسیون-بتا استفاده از فرایندهای برچسب‌گذاری است. برای این منظور، فرض

کنید $\xi_x(\cdot) = \{\xi_x(s), s \in \mathcal{D}\}$ یک فرایند برچسب‌گذاری وابسته به متغیرهای کمکی است که مقدار l را می‌گیرد اگر

$$\sum_{j=1}^{l-1} p_j \leq V_x(s) < \sum_{j=1}^l p_j. \quad (6.2)$$

²Dirac function

حال برای هر مجموعه متناهی s_1, \dots, s_n وزن‌های توأم وابسته فضایی برای فرایند برچسب‌گذاری به صورت

$$Q_{x(s_1), \dots, x(s_n)}(l_1, \dots, l_n) = P\left(\xi_x(s_1) = l_1, \dots, \xi_x(s_n) = l_n \mid \{v_l\}_{l=1}^L\right) \\ = P\left(\sum_{j=1}^{l_1-1} p_j \leq V_x(s_1) < \sum_{j=1}^{l_1} p_j; \dots; \sum_{j=1}^{l_n-1} p_j \leq V_x(s_n) < \sum_{j=1}^{l_n} p_j\right), \quad (7.2)$$

تعریف می‌شود. در این صورت به شرط فرایند برچسب‌گذاری، اثر تصادفی فضایی (۲.۲) را می‌توان به صورت $\eta_x(s) = \theta_{\xi_x(s)}$ نوشت. بنابراین، با حاشیه‌سازی روی ξ_x داریم $\eta_x | P_x \sim P_x$ و $\eta_x \sim BRSBP(\alpha, \psi, \gamma, \lambda, \sigma^2)$ از این‌رو با استفاده از فرایند برچسب‌گذاری (۷.۲) می‌توان یک تعریف معادل و یکسان برای $BRSBP$ ارائه کرد.

۱.۲ ویژگی‌های فرایند چوب‌شکنی رگرسیون-بتا

- فرایند چوب‌شکنی رگرسیون-بتا با توزیع متناهی بعد (۵.۲) در اصول سازگاری کولموگوروف صدق می‌کند.
- فرایند چوب‌شکنی رگرسیون-بتا رده‌ی بزرگتری از فرایندهای چوب‌شکنی است که در حالت خاص فرایند چوب‌شکنی پنهان (رودریگوز و همکاران، ۲۰۱۰) را در خود دارد. به بیان واضح‌تر، اگر در مدل $BRSBP$ به جای فرایند فضایی رگرسیون-بتا، فرایند فضایی یکنواخت $U(s)$ در نظر گرفته شود آنگاه فرایند چوب‌شکنی پنهان ($LaSBP$) حاصل می‌شود. علاوه بر این اگر در مدل پیشنهادی، ψ به بینهایت میل کند آنگاه مدل $BRSBP$ به یک مدل رگرسیون لجستیک تبدیل می‌شود. همچنین اگر $\mu_x(s)$ به صورت مدل رگرسیون پروبیت مدل‌بندی شود و ψ به بینهایت میل کند آنگاه مدل $BRSBP$ به مدل رگرسیون پروبیت تبدیل می‌شود. بنابراین مدل پیشنهادی رده‌ی بزرگتری است که در حالت خاص مدل‌های رگرسیون لجستیک و پروبیت را در خود دارد. همچنین توزیع بتا نسبت به توزیع یکنواخت انعطاف بیشتری دارد و پارامترهای این توزیع می‌تواند رفتارهای مختلفی از داده‌ها را مدل‌بندی کند.
- افزایش‌دهی القا شده توسط مدل $BRSBP$ روی اثر فضایی مبتنی بر متغیر کمکی است که با گسسته‌سازی فرایند فضایی رگرسیون-بتا حاصل می‌شود. همچنین در تعریف مدل $BRSBP$ از یک فرایند گاوسی برای القای وابستگی فضایی استفاده شده است که این کار موجب سادگی در محاسبات و برازش مدل نسبت به فرایندهای چوب‌شکنی دیگر می‌شود.

۳ استنباط بیزی و پیشگویی فضایی

از آن‌جا که هر یک از پارامترها ویژگی خاصی از میدان تصادفی را بیان می‌کنند، فرض می‌شود کلیه‌ی پارامترها مستقل از هم هستند. برای ضرایب رگرسیونی β و γ توزیع‌های پیشینی ناآگاهی‌بخش نرمال چندمتغیره در نظر گرفته می‌شود. برای پارامترهای واریانس σ^2 و τ^2 به ترتیب پیشینی مزدوج گامای و ارون ناآگاهی‌بخش $IG(s_{\sigma^2}, r_{\sigma^2})$ و $IG(s_{\tau^2}, r_{\tau^2})$ استفاده می‌شود. همچنین پیشینی گاما $G(s_\psi, r_\psi)$ و $G(s_\alpha, r_\alpha)$ به ترتیب برای پارامترهای α و ψ در نظر گرفته می‌شود. همچنین برای برازش مدل پیشنهادی، تابع همبستگی نمایی با پارامتر λ به صورت $\rho(s, s'; \lambda) = \exp\{-\frac{\|s-s'\|}{\lambda}\}$ فرض می‌شود. برای پارامتر همبستگی λ نیز توزیع گاما $G(a_\lambda, b_\lambda)$ در نظر گرفته می‌شود.

فرض کنید $\Theta = \{\beta, \gamma, \sigma^2, \tau^2, \lambda, \alpha, \psi\}$ بردار پارامترهای نامعلوم باشد. با توجه به این‌که توزیع پسینی دارای فرم بسته نیست، با اتخاذ روش نمونه‌گیری گیبز و داده‌افزایی توزیع‌های تمام‌شرطی تعیین و از آن‌ها نمونه‌گیری می‌شود. با توجه به این‌که برخی از

توزیع‌های تمام‌شرطی دارای فرم استاندارد نیستند از الگوریتم متروپولیس-هستینگس برای نمونه‌گیری از آن‌ها استفاده می‌شود. به منظور خلاصه‌سازی جزئیات محاسبات بیان نشده است.

یکی از اهداف اصلی تحلیل داده‌های فضایی پیشگویی فرایند مورد مطالعه در مکان‌های فاقد نمونه است. برای این منظور فرض کنید $y^* = y(s_0)$ و $y = (y(s_1), \dots, y(s_n))'$ به‌ترتیب نشان‌دهنده‌ی متغیر پاسخ در هر موقعیت دلخواه جدید s_0 و برداری از مشاهدات از فرایند Y در مکان‌های s_1, \dots, s_n باشند. در این صورت تحت مدل معرفی شده، توزیع پیشگوی فضایی برای مکان s_0 به صورت

$$p(y^*|y) = \int p(y^*|\xi_x(s_0), \theta_{\xi_x(s_0)}, \beta, \tau^2) p(\xi_x(s_0), \theta_{\xi_x(s_0)}, \beta, \tau^2|y) d\beta d\tau^2 d\xi_x(s_0) d\theta_{\xi_x(s_0)} \quad (1.3)$$

است که در آن

$$p(y^*|\xi_x(s_0), \theta_{\xi_x(s_0)}, \beta, \tau^2) \sim N(w'\beta + \theta_{\xi_x(s_0)}, \tau^2). \quad (2.3)$$

در رابطه (۱.۳) محاسبه‌ی انتگرال به صورت تحلیلی ناممکن است و با استفاده از روش‌های مونت‌کارلو می‌توان آن را تقریب زد. در این راستا، به ازای هر مشاهده پسینی از $\xi_x(s_0)$ ، $\theta_{\xi_x(s_0)}$ و β و τ^2 یک مشاهده از توزیع (۲.۳) تولید می‌شود. سپس با تکرار مراحل فوق به اندازه مورد نیاز، امکان دستیابی به مشاهدات از توزیع پیشگو به صورت $\{y(s_0)^{(m)}, j = 1, \dots, M\}$ مهیا می‌شود. در نتیجه پیشگویی فضایی بیزی در مکان s_0 به صورت زیر بدست می‌آید:

$$\hat{y}(s_0) \equiv \hat{\mathbb{E}}(y(s_0)|y) = \frac{1}{M} \sum_{m=1}^M y^{(m)}(s_0).$$

۴ مثال کاربردی

در این بخش از مقاله به عملکرد و ارزیابی مدل معرفی شده از دیدگاه پیشگویی فضایی و افرازبندی ناحیه مورد مطالعه در مقایسه با مدل‌های معرفی شده توسط ریچ و همکاران (۲۰۱۱)، رابیز و کالدِر (۲۰۱۵)، گرمسی و لی (۲۰۰۸)، پیچ و کوئینانا (۲۰۱۶) و رودریگوئز و همکاران (۲۰۱۰) پرداخته می‌شود. برای این منظور، این مدل‌ها به اختصار به‌ترتیب به صورت $TGP, M2, M1$ و $LaSBP, PPM$ و $BRSBP$ نمایش داده می‌شوند. داده‌های مورد بررسی متوسط میزان بارش برحسب میلی‌متر مربوط به ۴۱۳ ایستگاه هواشناسی در سال ۲۰۱۶ برای کشور ایران است. این داده‌ها به وسیله‌ی سازمان هواشناسی و وزارت نیرو اندازه‌گیری شده است. در ادامه برای بررسی نامانایی داده‌ها در ساختار دوم فرایند از آزمون نامانایی که توسط سویترباندیوپی و سوبارائو (۲۰۱۷) بیان گردید، استفاده شده است که در سطح ۵ درصد مانایی داده‌ها رد می‌شود. همچنین با توجه به این که میزان بارش به متغیرهای هواشناسی از جمله دما وابسته است، دما به عنوان متغیر کمکی برای برازش مدل‌ها انتخاب می‌کنیم. از آن‌جا که هدف این مقاله بررسی تاثیر متغیر کمکی در همبستگی فضایی است، ابتدا روند از داده‌ها حذف و با داده‌های بدون روند برازش مدل‌ها را انجام می‌دهیم.

در ادامه از معیار میانگین توان دوم خطا^۳ ($MSPE$) برای ارزیابی قدرت پیشگویی مدل‌های بیان شده، استفاده می‌شود. برای این منظور ۱۰ مجموعه داده با حذف تصادفی ۱۰ درصد از داده‌ها تولید شده و سپس با استفاده از رهیافت بیزی مدل‌های $M1, TGP, M2, PPM, LaSBP$ و $BRSBP$ به داده‌ها برازش داده شده است. برای برازش مدل، سه زنجیر با نقاط شروع اولیه

³Mean Squared Prediction Error

جدول ۱: مقایسه‌ی عملکرد پیشگویی مدل‌های: M_1 , M_2 , $LaSBP$, TGP , PPM و $BRSBP$.

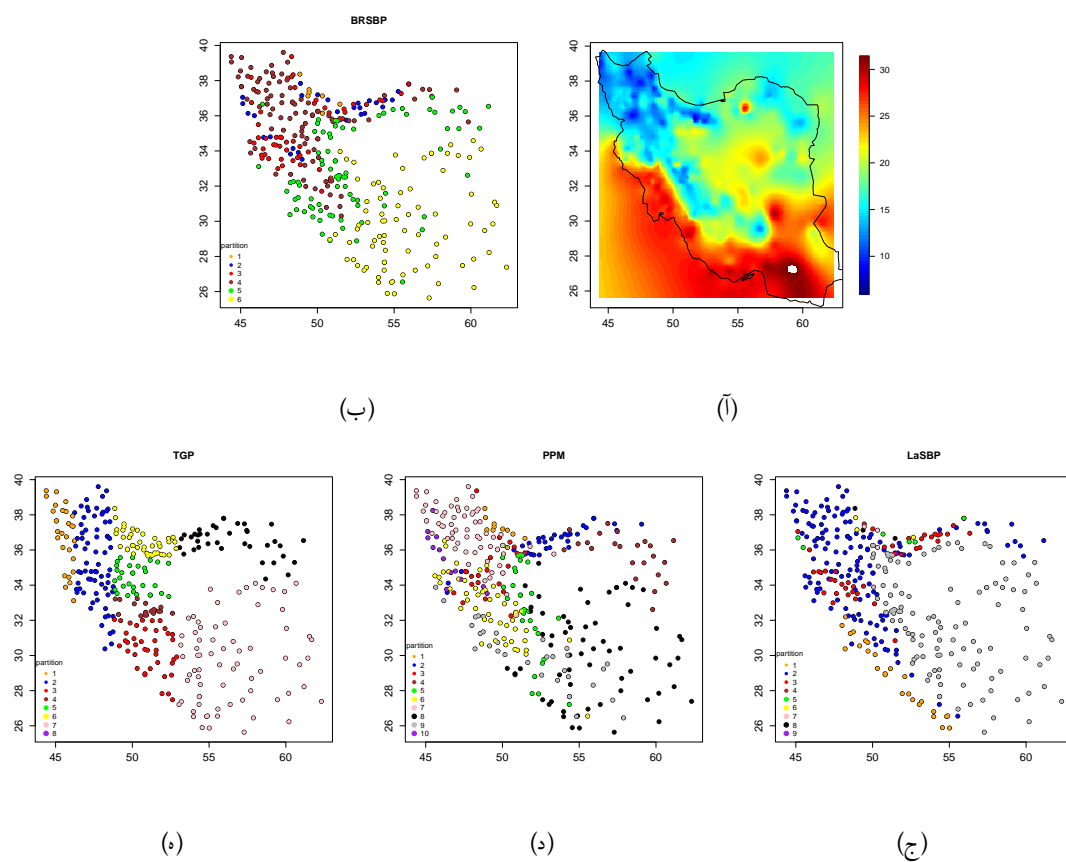
مدل	M_1	M_2	$LaSBP$	$BRSBP$	PPM	TGP
MSPE	۱/۹۲	۱/۳۶	۱/۰۸	۰/۶۴	۰/۸۸	۰/۷۶

متفاوت و ۲۰۰۰۰ تکرار الگوریتم MCMC اجرا شد که تعداد ۵۰۰۰ تکرار اولیه به عنوان دوره‌ی داغیدن در نظر گرفته شد. به منظور خلاصه‌سازی، جزئیات انتخاب پیشینی و نحوه‌ی استنباط بیزی این مدل‌ها بیان نشده است. در جدول ۱ نتایج هر یک از مدل‌ها براساس معیار $MSPE$ آمده است. همان‌طور که مشاهده می‌شود مدل $BRSBP$ عملکرد بهتری نسبت به دیگر مدل‌ها دارد. به‌خصوص، مدل‌های M_1 و M_2 عملکرد ضعیف‌تری نسبت به مدل پیشنهادی دارند. در واقع مدل $BRSBP$ مقدار $MSPE$ کمتری نسبت به مدل‌های M_1 , M_2 , PPM و TGP داراست. مقدار $MSPE$ برای مدل پیشنهادی ۰/۶۴ و برای مدل $LaSBP$ برابر با ۱/۰۸ است که تقریباً ۴۱٪ کوچکتر است.

در شکل ۱ نمودار افرازبندی ایستگاه‌های هواشناسی حاصل از برازش مدل $BRSBP$ آمده است. همان‌طور که مشاهده می‌شود ایستگاه‌هایی که رفتار دمایی یکسانی دارند در خوشه‌های یکسان قرار گرفته‌اند (شکل ب ۱). به عنوان مثال، مناطق شرقی و جنوب شرقی ایران که دارای متوسط دمای بالا و میزان بارش کم هستند، در خوشه‌های یکسان قرار گرفته‌اند. همچنین، در شکل‌های ج ۱، د ۱ و ه ۱ به‌ترتیب افرازبندی ایستگاه‌های هواشناسی حاصل از برازش مدل‌های $LaSBP$, PPM و TGP آمده است. همان‌طور که مشاهده می‌شود، افرازبندی ایجاد شده توسط مدل‌های $LaSBP$ و PPM تقریباً شبیه مدل $BRSBP$ است اما تعداد خوشه‌های ایجاد شده براساس مدل $BRSBP$ کمتر است. اما در مدل TGP ایستگاه‌های هواشناسی نزدیک به هم، در خوشه‌های یکسان قرار گرفته‌اند. با این حال، با توجه به تنوع آب و هوایی که در ایستگاه‌های نزدیک به هم وجود دارد، افرازبندی ایجاد شده توسط مدل TGP منطقی و معقول به نظر نمی‌رسد.

بحث و نتیجه‌گیری

در این مقاله با استفاده از فرایند رگرسیون-بتا، رویکرد جدیدی برای ساخت فرایندهای فضایی مورد مطالعه قرار گرفت به‌طوری‌که توزیع تصادفی وابسته به متغیر کمکی است. همانند فرایندهای چوب‌شکنی تعمیم‌یافته تابع کوواریانس مدل معرفی شده نامانا و وابسته به متغیر کمکی است. از ویژگی‌های مهم مدل معرفی شده می‌توان به خوشه‌بندی بر اساس متغیر کمکی اشاره کرد که باعث بهبود در افرازبندی ناحیه مورد مطالعه می‌شود. همچنین با استفاده از یک مجموعه داده واقعی عملکرد مدل از نظر پیشگویی فضایی و خوشه‌بندی مورد ارزیابی قرار گرفت.



شکل ۱: (آ) نمودار رویه دما و (ب، ج، د، ه) نمودارهای افزاینده‌های ایستگاه‌های هواشناسی از برازش مدل‌های $BRSBP$ ، $LaSBP$ ، TGP و PPM .

- Bandyopadhyay, S., and Rao, S. S. (2017), *A test for stationarity for irregularly spaced spatial data*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 79, 95–123.
- Barcella, W., M. De Iorio, S. Favaro, and G. L. Rosner (2018), *Dependent generalized dirichlet process priors for the analysis of acute lymphoblastic leukemia*, Biostatistics, 19(3), 342–358.
- Barzegar, Z., and Rivaz, F. (2020), *A scalable Bayesian nonparametric model for large spatio-temporal data*, Computational Statistics, 35, 153–173.
- Dahdouh, O., and Khaledi, M. J., (2020), *Generalized spatial stick-breaking processes*, Communications in Statistics Simulation and Computation, 1-20.
- Dunson, D. B. and J.-H. Park (2008), *Kernel stick-breaking processes*, Biometrika 95(2), 307–323.
- Gramacy, R. B. and H. K. H. Lee (2008), *Bayesian treed gaussian process models with an application to computer modeling*. Journal of the American Statistical Association, 103(483), 1119–1130.
- Ingebrigtsen, R., F. Lindgren, and I. Steinsland (2014), *Spatial models with explanatory variables in the dependence structure*, Spatial Statistics 8, 20–38.
- Page, G. L. and F. A. Quintana (2016), *Spatial product partition models*, Bayesian Analysis 11(1), 265–298. 149.
- Reich, B. J. and M. Fuentes (2007), *A multivariate semiparametric Bayesian spatial modeling framework for hurricane surface wind fields*, The Annals of Applied Statistics 1(1), 249–264.
- Reich, B.J., Eidsvik, J., Guindani, M., Nail, A.J., Schmidt, A.M., (2011), *A class of covariate-dependent spatiotemporal covariance functions for the analysis of daily ozone concentration*, The annals of applied statistics. 5 (4), 2425—2447.
- Ren, L., L. Du, L. Carin, and D. B. Dunson (2011), *Logistic stick-breaking process*, Journal of Machine Learning Research 12(1).
- Risser, M.D., Calder, C.A., (2015), *Regression-based covariance functions for nonstationary spatial modeling*, Environmetrics 26 (4), 284—297.
- Rodriguez, A. and D. B. Dunson (2011), *Nonparametric Bayesian models through probit stick-breaking processes*, Bayesian Analysis 6(1).

- Rodríguez, A., Dunson, D.B. and Gelfand, A.E., (2010), *Latent stick-breaking processes*, Journal of the American Statistical Association, 105(490), pp.647-659.
- Schmidt, A. M., P. Guttorp, and A. O'Hagan (2011), *Considering covariates in the covariance structure of spatial processes*, Environmetrics 22(4), 487–500.
- Yarali, E., Rivaz, F., (2020), *Incorporating covariate information in the covariance structure of misaligned spatial data*, Environmetrics 31 (6), e2623.
- Yarali, E., Rivaz, F., and Khaledi, M. J. (2022), *A Bayesian nonparametric spatial model with covariate-dependent joint weights*, Spatial Statistics, 100662.

Beta regression stick-breaking process with covariate dependent nonstationary covariance

Esmail Yarali, Firoozeh Rivaz¹, and Majid Jafari Khaledi²

¹Department of Statistics, Shahid Beheshti University, Tehran.

²Department of Statistics, Tarbiat Modares University, Tehran.

Abstract: This paper introduces a new class of stick-breaking processes that incorporate covariate information into the random joint distributions of spatial processes. The proposed modeling approach allows for a nonstationary covariance function that is driven by covariates, and it also induces a covariate-dependent random partitioning scheme. Using this process as a random effect distribution within a linear mixed model framework, the data can be effectively modeled, and a Bayesian inference approach is developed. The performance of the proposed model is evaluated and compared with that of competing spatial models using a real-world example.

Keywords: Stick-breaking process, Beta-Regression, Covariate information, Non-stationarity, Covariance structure.

Mathematics Subject Classification (2020): 62H11, 62M30, 62F15.
