

مدل‌بندی رگرسیونی داده‌های کروی با استفاده از تابع مخاطره هاورسین و روش کم‌ترین توان‌های دوم خطا

میثم مقیم‌بیگی، موسی گلعلی‌زاده*؛ دانشگاه تربیت مدرس، گروه آمار

پذیرش ۹۷/۰۶/۲۶

دریافت ۹۷/۰۲/۲۱

چکیده

از دیرباز محققان به تحلیل آماری داده‌ها روی سطح کروی زمین توجه داشته‌اند. داده‌هایی از این دست می‌تواند مربوط به مهاجرت برخی از حیوانات از منطقه‌ای به منطقه‌ای دیگر باشد. آن‌گاه مدل‌بندی آماری مسیر حرکت آن‌ها به محققان علوم‌زیستی کمک می‌کند تا بتوانند برای حرکت آن‌ها پیش‌گویی داشته و هم‌چنین محدوده‌ای را برآورد کنند که حضور حیوانات در آن منطقه محتمل‌تر باشد. برای بررسی چنین پدیده‌هایی در این مقاله، مدل‌بندی آماری مسیر حرکت اشیاء روی کره به روش‌های ناپارامتری و کم‌ترین توان‌های دوم خطا مدنظر قرار گرفته است. این مدل‌بندی براساس دو مدل جدا از هم برای زوایای شکل گرفته روی کره انجام می‌گیرد. مدل‌های ارائه شده با استفاده از داده‌های شبیه‌سازی شده و داده واقعی ارزیابی خواهند شد.

واژه‌های کلیدی: تابع مخاطره، داده‌های کروی، رگرسیون ناپارامتری، مدل طولی.

مقدمه

از نقطه نظر هندسی داده‌های کروی داده‌هایی هستند که مقادیرشان را روی کره واحد اختیار می‌کنند. روش‌های زیادی برای برازش یک خم و به‌ویژه منحنی‌های رگرسیونی به‌روی داده‌های کروی وجود دارد. برای معرفی مدل‌های رگرسیون گولد^۱ [۱] از زوایای نقاط موجود روی کره استفاده کرد. او توزیع فیشر را به‌عنوان توزیع خطا در تحلیل‌های خود به‌کار گرفت. نسخه ناپارامتری از مدل او را تامپسون و کلارک^۲ [۲] پیشنهاد دادند. داده‌هایی که نزدیک به قطب شمال یا جنوب باشند رفتار متفاوتی نسبت به داده‌های دور از قطب دارند. این مسئله مشکل اصلی مدل معرفی شده آن‌ها بود. از این‌رو، ایشان سعی کردند داده‌ها را به‌طریقی دور از قطب نگه دارند. آن‌ها بعداً توانستند با استفاده از صفحه مماسی بر این مسئله غلبه کنند و پیشنهاد استفاده از اسپلاین‌ها را در آن صفحه ارائه دادند [۳]. پس از آن‌ها فیشر^۳ و همکاران [۴] دو خانواده از اسپلاین‌های کروی را برای داده‌های کروی پیشنهاد کردند. آن‌ها دو خانواده از خم‌ها را با استفاده از هندسه دیفرانسیل که برای برازش اسپلاین مناسب و به‌علاوه تحت انتخاب مختصات پایا بودند را معرفی کردند.

یکی از شیوه‌های پیش‌بینی در آمار، استفاده از روش‌های رگرسیون ناپارامتری است. بررسی مدل‌های هموار نیز در بسیاری از نوشته‌ها مورد توجه بوده است. هردوی این رویکردها همراه با شیوه‌های دیگر در مورد آمار ناقلیدسی (آمار

نویسنده مسئول golalizadeh@modares.ac.ir

1. Gould
2. Thompson and Clark
3. Fisher

مربوط به تحلیل داده‌های موجود در فضای ناقلیدسی) شامل آمار جهتی نیز مد نظر قرار گرفته است. مسیر اسپلاینی با استفاده از پارامترهای دورانی اولین بار در [۵] برای داده‌های دایره‌ای یا همان داده‌های زاویه‌ای پیشنهاد شد. ساختار مدل رگرسیونی ناپارامتری با استفاده از کمینه‌سازی تابع مخاطره اقلیدسی نیز در [۶] بررسی شد. این تابع مخاطره از این نظر اقلیدسی خوانده می‌شود که بر اساس فواصل اقلیدسی شکل گرفته است.

هدف اصلی مقاله حاضر این است که با استفاده از ایده معرفی شده در [۶] برای داده‌های روی دایره و همچنین ارائه شده در [۱] به معرفی مدل رگرسیونی ناپارامتری و کم‌ترین توان‌های دوم خطا برای داده‌های کروی بپردازد. اگرچه در مدل ناپارامتری معرفی شده در [۱] مدل‌های ارائه شده بر اساس زوایا مستقل از هم هستند، اما با انتخاب یک تابع مخاطره مناسب می‌توان هم‌بستگی میان زوایا را روی کره لحاظ کرد. این موضوع در مقاله حاضر دنبال شده است.

مدل طولی ناپارامتری برای داده‌های کروی

از اصول مبانی هندسه می‌دانیم که نقاط روی کره می‌توانند به وسیله دو زاویه مشخص شوند. زوایای θ و ϕ از معروف‌ترین این زوایا هستند. حال اگر فرض شود n مشاهده روی کره در اختیارند و هدف تحقیق، پایه‌ریزی یک مدل‌بندی آماری مناسب باشد یکی از فرض‌های اساسی می‌تواند استقلال مشاهدات روی کره باشد. گولد [۱] با این فرض مدل‌های نوع خطی را بنا نهاد. تامپسون و کلارک [۲] نیز برای مدل پارامتری‌شان از ویژگی‌های خاص زوایا کمک گرفتند. با اختیار داشتن متغیرهای (θ_i, ϕ_i) به ازای $i = 1, 2, \dots, n$ ، مدل پارامتری آن‌ها بر اساس توابع،

$$\begin{cases} \phi_i = f_1(t_i) + \varepsilon_{i1} \\ \theta_i = f_2(t_i) + \varepsilon_{i2} \end{cases} \quad (1)$$

به طوری که ε_{i1} و ε_{i2} خطاهای مدل و t_i متغیر تبیینی است که معمولاً متغیر زمان در نظر گرفته می‌شود. تامپسون و کلارک [۲] فرض کردند که θ_i ها دارای توزیع نرمال استاندارد با میانگین $f_2(t_i)$ و واریانس $1/\kappa$ و ϕ_i ها دارای توزیع نرمال با میانگین $f_1(t_i)$ و واریانس $(1/\kappa) \sin^2 f_2(t_i)$ هستند. دی‌مارزیو^۱ و همکاران [۶] یک مدل ناپارامتری برای مشاهدات زاویه‌ای در زمان t_i به صورت $\phi_i = g(t_i) + \varepsilon_i$ معرفی کردند که در آن ε_i یک متغیر تصادفی زاویه‌ای با میانگین معلوم و پارامتر تمرکز متناهی است. به منظور برآورد پارامترها، آن‌ها از تابع مخاطره اقلیدسی (تابع مخاطره‌ای که در فضای اقلیدسی استفاده می‌شود) استفاده کردند. به منظور ارائه مدلی رگرسیونی ناپارامتری برای داده‌های کروی، تعمیم روش آن‌ها از دایره روی کره مفید خواهد بود. از این رو، در این مقاله با استفاده از مدل و ایده معرفی شده در [۲]، [۶] و یک تابع مخاطره مناسب برای داده‌های کروی، مدل رگرسیونی به روش ناپارامتری برای داده‌های کروی ارائه می‌شود.

یکی از روابط مهم در حوزه جهت‌یابی^۲ فاصله هاورسین^۱ است. این فرمول بر اساس فاصله دایره بزرگ بین دو نقطه روی کره از طریق طول و عرض جغرافیایی به دست می‌آید. برای هر دو نقطه (θ_1, ϕ_1) و (θ_2, ϕ_2) روی کره با شعاع یک، فاصله دایره بزرگ بین آن‌ها به صورت:

1. Di Marzio
2. Navigation

$$d = r \sin^{-1} \left(\sqrt{\sin^2 \left(\frac{\theta_r - \theta_l}{r} \right) + \cos(\theta_l) \cos(\theta_r) \sin^2 \left(\frac{\phi_r - \phi_l}{r} \right)} \right)$$

و فاصله هاورسین بین آن‌ها بدین صورت نوشته می‌شود [۷]:

$$D = \text{hav}(d) = \text{hav}(\theta_r - \theta_l) + \cos \theta_l \cos \theta_r \text{hav}(\phi_r - \phi_l)$$

که در آن $\text{hav}(d) = \sin^2(d/2)$. بر اساس این رابطه می‌توان تابع مخاطره هاورسین را روی کره بدین صورت تعریف کرد:

$$E[D | t] = E[\sin \theta \sin f_r(t) - \cos \theta \cos f_r(t) \cos(f_l(t) - \phi) | t] \quad (2)$$

اگر برای دامنه عمومی زمانی t ، فرض شود

$$m_l(t) = E[\cos \theta \sin \phi | t], \quad m_r(t) = E[\cos \theta \cos \phi | t],$$

در این صورت تابع مخاطره (۲) زمانی که $f_l(t) = \tan^{-1}(m_l(t) / m_r(t))$ کمینه می‌شود. هم‌چنین با تعریف

$$m_r(t) = E[\sin \theta | t], \quad m_f(t) = E[\cos \theta \cos(f_l(t) - \phi) | t]$$

تابع مخاطره (۲) به‌ازای $f_r(t) = \tan^{-1}(m_r(t) / m_f(t))$ کمینه می‌شود. به‌منظور اثبات این موضوع کافی است از تابع مخاطره (۲) نسبت به $f_l(t)$ و $f_r(t)$ مشتق گرفته شود. با مشتق‌گیری از تابع مخاطره (۲) و برابر صفر قرار دادن نتایج داریم:

$$\begin{aligned} \frac{\partial E[D | t]}{\partial f_l(t)} &= \cos f_r(t) E[\cos \theta \sin(f_l(t) - \phi) | t] \\ &= \cos f_r(t) E[\cos \theta \cos \phi \sin f_l(t) - \cos \theta \sin \phi \cos f_l(t)] \\ &= \cos f_r(t) \cos f_l(t) m_r(t) - \cos f_r(t) \cos f_l(t) m_l(t) = 0, \\ \frac{\partial E[D | t]}{\partial f_r(t)} &= -\cos f_r(t) E[\sin \theta | t] + \sin f_r(t) E[\cos \theta \cos(f_l(t) - \phi) | t] \\ &= -\cos f_r(t) m_r(t) + \sin f_r(t) m_f(t) = 0, \end{aligned}$$

ملاحظه می‌شود که حل دو معادله اخیر منجر به روابط

$$\widehat{f}_l(t) = \tan^{-1}(m_l(t) / m_r(t)) \quad \text{و} \quad \widehat{f}_r(t) = \tan^{-1}(m_r(t) / m_f(t))$$

می‌شوند. حال به‌منظور این که نشان دهیم توابع حاصل کمینه‌کننده تابع مخاطره (۲) هستند، مشتق مرتبه دوم را

نیز محاسبه می‌کنیم. از این‌رو داریم:

$$\begin{aligned} \xi_{l1}^2 &= \frac{\partial^2 E[D | t]}{\partial f_l^2(t)} = \cos f_r(t) E[\cos \theta \cos(f_l(t) - \phi) | t], \\ \xi_{r1}^2 &= \frac{\partial^2 E[D | t]}{\partial f_r^2(t)} = \sin f_r(t) E[\sin \theta | t] + \cos f_r(t) E[\cos \theta \cos(f_l(t) - \phi) | t], \\ \xi_{lr}^2 &= \frac{\partial^2 E[D | t]}{\partial f_l(t) \partial f_r(t)} = -\sin f_r(t) E[\cos \theta \sin(f_l(t) - \phi) | t]. \end{aligned}$$

با استفاده از آزمون مشتق دوم در نقاط بحرانی $\widehat{f}_l(t)$ و $\widehat{f}_r(t)$ بدین صورت داریم:

$$\begin{aligned} \xi_1^2 \xi_2^2 - (\xi_1^2)^2 \left| \begin{array}{c} \widehat{f}_1(t) \\ \widehat{f}_2(t) \end{array} \right| &= \cos f_2(t) m_2(t) (\sin f_1(t) m_1(t) + \cos f_1(t) m_2(t)) \\ &\quad - [\sin f_2(t) E[\cos \theta \sin(f_1(t) - \phi) | t]]^2 \\ &= \cos f_2(t) m_2(t) \sin f_1(t) m_1(t) + \cos^2 f_2(t) m_2^2(t) - 0 > 0. \end{aligned}$$

$$\xi_1^2 \left| \begin{array}{c} \widehat{f}_1(t) \\ \widehat{f}_2(t) \end{array} \right| = \cos f_2(t) \cos f_1(t) m_2(t) + \cos f_2(t) \sin f_1(t) m_1(t) > 0, \quad \forall; -\frac{\pi}{2} < \theta, \phi < \frac{\pi}{2},$$

که در آن $\sin \tan^{-1}(x) = x / \sqrt{1+x^2}$ و $\cos \tan^{-1}(x) = 1 / \sqrt{1+x^2}$ است. بنابراین کمینه شدن تابع مخاطره بر بازه $-\pi/2 < \theta, \phi < \pi/2$ تضمین می‌شود. از طرفی به این نکته توجه شود که در صورتی که داده‌های روی نیمکره بالایی قرار نداشته باشند با استفاده از توابع دوران می‌توان آن‌ها را به نیم‌کره بالا انتقال و پس از تحلیل همراه با نتایج حاصل به مکان اولیه منتقل شوند. از آن‌جاکه محققان در مسائل واقعی با مقادیر مشاهده شده (گسسته) از (θ_i, ϕ_i) روبرو هستند، از این‌رو، روش مونت کارلویی با وزن‌های مناسب می‌تواند برای محاسبه برآوردگرها استفاده شود. برای معرفی یک برآورد از $f_1(t)$ و $f_2(t)$ بر این اساس، فرض کنید $n, (\theta_i, \phi_i)$ مشاهده وابسته روی کره در n زمان مجزا t_1, t_2, \dots, t_n باشند. برای هر $j = 1, 2, 3, 4$ برآوردگرهای گشتاوری از $m_j(t)$ بدین‌صورت قابل معرفی هستند:

$$\widehat{m}_1(t) = \frac{1}{n} \sum_{i=1}^n \cos \theta_i \sin \phi_i W_1(t_i - t),$$

$$\widehat{m}_2(t) = \frac{1}{n} \sum_{i=1}^n \cos \theta_i \cos \phi_i W_2(t_i - t),$$

$$\widehat{m}_3(t) = \frac{1}{n} \sum_{i=1}^n \sin \theta_i W_3(t_i - t),$$

$$\widehat{m}_4(t) = \frac{1}{n} \sum_{i=1}^n \cos \theta_i \cos(\phi_i - f_1(t)) W_4(t_i - t),$$

به‌طوری‌که W_1 و W_2 توابع وزنی هستند و

$$\widehat{f}_1(t) = \tan^{-1}(\widehat{m}_1(t) / \widehat{m}_2(t)).$$

می‌توان ملاحظه کرد که برآورد $f_2(t)$ بدین‌صورت است:

$$\widehat{f}_2(t) = \tan^{-1}(\widehat{m}_3(t) / \widehat{m}_4(t))$$

واضح است که برآورد $f_1(t)$ مستقل از $f_2(t)$ بوده است، در صورتی که برآورد $f_2(t)$ وابسته به مقدار $\widehat{f}_1(t)$ است. به‌منظور برآورد توابع $m_j(t)$ نیاز به معرفی توابع وزن مناسب است. توابع وزن بسیاری در نوشتگان در زمینه‌های مختلف آماری استفاده شده است که از معروف‌ترین آن‌ها می‌توان به توابع کرنل اشاره کرد [۸]. یکی از توابع کرنل که در نوشتگان بسیار استفاده می‌شود تابع چگالی نرمال است. در این مقاله با توجه به آسانی کار با تابع چگالی نرمال و ویژگی متقارن بودن آن از تابع کرنل نرمال استفاده می‌شود. به‌همین منظور فرض می‌شود که توابع وزن W_1 و W_2 توابع چگالی نرمال با میانگین صفر و واریانس σ_1^2 و σ_2^2 باشند. از آن‌جاکه پارامترهای واریانس نقش پهنای باند را دارند

از این‌رو، نیاز به برآورد این پارامترها است. روش‌های مختلفی برای برآورد پهنای باند (h) در تابع کرنل وجود دارد که برای مثال می‌توان به روش پهنای باند بهینه، روش ارزیابی متقابل بیشینه درست‌نمایی، ارزیابی متقابل نااریب و اریب اشاره کرد. در این مقاله به منظور برآورد پارامتر پهنای باند از روش ارزیابی متقابل بیشینه درست‌نمایی استفاده می‌کنیم. این روش در [۹] و [۱۰] معرفی شد که در آن پارامتر پهنای باند براساس مشاهدات X_1, \dots, X_n با تابع چگالی $l_h(x)$ از طریق تابع درست‌نمایی نما $\prod_{i=1}^n \hat{l}_h(x_i)$ برآورد می‌شود. یادآوری می‌شود که تابع $\hat{l}_h(x)$ برآورد تابع چگالی $l_h(x)$ به صورت

$$\hat{l}_h(x) = \frac{1}{(n-1)h} \sum_{i=1}^n g\left(\frac{x_i - x}{h}\right)$$

که در آن g تابع کرنل است. به‌طور بدیهی یک برآورد بیشینه درست‌نمایی برابر $h = 0$ است. به‌منظور حذف این برآورد ناکارآمد کافی است تابع اعتبارسنج $\hat{l}_h(x)$ را با تابع اعتبارسنج $\hat{l}_{h,i}(x)$ جای‌گزین کنیم که برای $k = 1, 2$ و زمان T_j و T_i به صورت

$$\hat{l}_{h,i,k}(x_i) = \frac{1}{(n-1)h_k} \sum_{j \neq i} g_k\left(\frac{T_j - T_i}{h_k}\right) \quad (3)$$

است. بنابراین یک برآورد اعتبارسنجی بیشینه درست‌نمایی به صورت

$$h_{mlcv,k} = \arg \max_{h_k > 0} MLCV(h_k)$$

است که در آن

$$MLCV(h_k) = \left(n^{-1} \sum_{i=1}^n \log \left[\sum_{j \neq i} g_k\left(\frac{T_j - T_i}{h_k}\right) \right] - \log[(n-1)h_k] \right).$$

در ادامه روش دیگری برای مدل‌بندی آماری مسیر حرکت اشیاء روی کره معرفی می‌شود. این روش مبتنی بر استفاده از تابع ربط میان زوایا و متغیر تبیینی و براساس روش کم‌ترین توان دوم خطا است.

روش کم‌ترین توان دوم خطا مبتنی بر تابع ربط

یکی از توابع ربط برای ایجاد ارتباط در داده‌های زاویه‌ای تابع $\tan^{-1}(\cdot)$ است. این تابع ربط به‌دلیل ایجاد ارتباط میان اعداد حقیقی و زوایا در مدل‌بندی داده‌های زاویه‌ای محبوبیت زیادی دارد [۱۱]. با استفاده از این تابع می‌توان مدل (۱) را به صورت (۴) نوشت:

$$\begin{cases} \phi_i = a_1 + f_1(\tan^{-1}(b_1 + c_1 t_i)) + \varepsilon_{i_1} \\ \theta_i = a_2 + f_2(\tan^{-1}(b_2 + c_2 t_i)) + \varepsilon_{i_2} \end{cases} \quad (4)$$

که در آن برای $k = 1, 2$ ، a_k ، b_k و c_k پارامترهای مدل و f_k توابعی معلوم هستند. به‌منظور برآورد پارامترها کافی است کمیت $\sum_{i=1}^n \varepsilon_{ik}^2$ براساس پارامترهای مدل کمینه شوند. به‌همین منظور کافی است از $\sum_{i=1}^n \varepsilon_{ik}^2$ نسبت به پارامترها مشتق گرفته شود. با تعریف $\omega_{i_1} = \phi_i$ و $\omega_{i_2} = \theta_i$ برآورد کم‌ترین توان‌های دوم خطای پارامترها از حل معادلات (۵)، (۶)، (۷)

(۷)، (۶)، (۵)

$$\sum_{i=1}^n (\omega_{ik} - a_k - f_k(\tan^{-1}(b_k + c_k t_i))) = 0, \tag{۵}$$

$$\sum_{i=1}^n \frac{f_k'(\tan^{-1}(b_k + c_k t_i))}{1 + (b_k + c_k t_i)^2} (\omega_{ik} - a_k - f_k(\tan^{-1}(b_k + c_k t_i))) = 0, \tag{۶}$$

$$\sum_{i=1}^n \frac{f_k'(\tan^{-1}(b_k + c_k t_i)) t_i}{1 + (b_k + c_k t_i)^2} (\omega_{ik} - a_k - f_k(\tan^{-1}(b_k + c_k t_i))) = 0. \tag{۷}$$

و با استفاده از الگوریتم بازگشتی زیر به دست می‌آیند:

الف) ابتدا با فرض معلوم بودن تابع f_k و مقدار اولیه $a_k = \bar{\omega}_k$ پارامترهای b_k و c_k به صورت هم‌زمان و با استفاده

از روش کم‌ترین توان دوم خطا در رابطه (۵) برآورد می‌شوند.

ب) با استفاده از رابطه (۶) و با فرض معلوم بودن a_k و c_k پارامتر b_k برآورد می‌شود.

ج) با داشتن مقادیر برآورد شده a_k و b_k پارامتر c_k با استفاده از رابطه (۷) برآورد می‌شود.

د) پس از برآورد b_k و c_k ، با استفاده از رابطه (۵) پارامتر a_k به روز رسانی می‌شود.

ه) با زگشت به مرحله (ب).

بررسی شبیه‌سازی برای مدل طولی

به منظور بررسی عملکرد مدل (۱) از یک ایده ساده در شبیه‌سازی نقاط در صفحه استفاده می‌شود. در مفهوم آمار شکل، شکل برابر با همه اطلاعات باقی‌مانده از یک شیء پس از حذف اثرات مکان، مقیاس و دوران از یک شیء است [۱۲]. بر اساس این تعریف، شکل هر مثلث می‌تواند به وسیله نقطه‌ای روی کره توصیف شود مختصات کره‌ی چنین مثلث‌هایی بر سطح کره در [۱۳] ارائه شده است. به عبارت دیگر، هر شکل هر مثلث، به دور از اثر دوران، مکان و مقیاس را به طور ساده می‌توان با دو زاویه درونی آن مشخص کرد. این دو زاویه می‌تواند مختصات یک نقطه روی کره باشد. از این رو، برای شبیه‌سازی نقاط روی کره کافی است با فرض ثابت بودن دو رأس مثلث، رأس سوم آن شبیه‌سازی شود. این روش کمک می‌کند که نقاط روی صفحه شبیه‌سازی شود که به مراتب آسان‌تر از شبیه‌سازی روی کره است. به همین منظور، ابتدا فرض کنید رأس سوم مثلث‌ها دارای خاصیت مارکف از مرتبه اول در صفحه باشند. سپس بردار تصادفی $X = (X_1, X_2)^T$ از توزیع نرمال ۴ متغیره با میانگین $\mu = (\mu_1, \mu_2)^T$ و ماتریس کوواریانس

$$\begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

را در نظر بگیرید که در آن $|\Sigma_{11}| > 0$ و $|\Sigma_{22}| > 0$. از ادبیات آمار چند متغیره می‌دانیم که X_2 به شرط $X_1 = x_1$ دارای توزیع نرمال دو متغیره با میانگین $(X_1 - \mu_1) + \Sigma_{21} \Sigma_{11}^{-1} (x_1 - \mu_1)$ و ماتریس کوواریانس $\Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$ است. توجه داریم که بردار دو متغیره X_1 مربوط به طول و عرض نقطه اول و بردار دو متغیره X_2 مربوط به طول و عرض نقطه دوم است. در نتیجه با استفاده از توزیع شرطی و با یک الگوریتم تکراری می‌توان مثلث‌ها را تولید کرد. به همین منظور ابتدا فرض کنید $\mu_1 = (0/25, 0/77)^T$ و $\mu_2 = (\mu_{21}, \mu_{22})^T$ تابعی از زمان باشد که در آن

$$\begin{cases} \mu_{21} = 0/0017t^2 - 0/07t + 0/39 \\ \mu_{22} = 0/00012t^3 - 0/005t^2 + 0/06t + 0/67. \end{cases}$$

هم‌چنین به‌منظور سادگی مسئله فرض کنید ماتریس‌های کوواریانس مدل به‌صورت مؤلفه‌هایی از مقادیر ثابت بدین‌صورت باشند:

$$\Sigma_{11} = \begin{bmatrix} 0/15 & 0/10 \\ 0/10 & 0/15 \end{bmatrix}, \quad \Sigma_{12} = \begin{bmatrix} 0/14 & 0/13 \\ 0/11 & 0/10 \end{bmatrix},$$

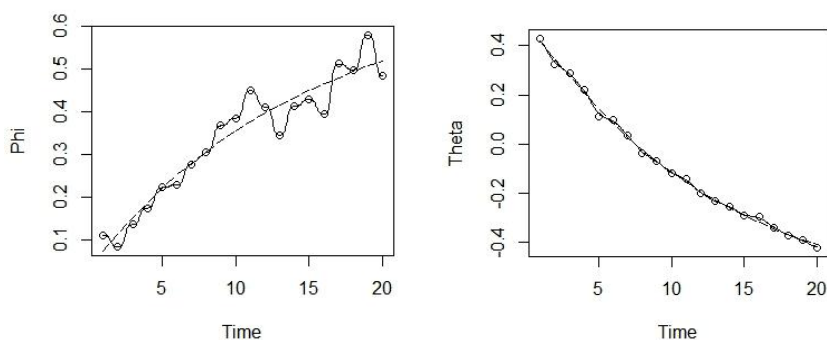
$$\Sigma_{21} = \begin{bmatrix} 0/14 & 0/11 \\ 0/13 & 0/10 \end{bmatrix}, \quad \Sigma_{22} = \begin{bmatrix} 0/15 & 0/13 \\ 0/13 & 0/15 \end{bmatrix},$$

الگوریتم ذیل برای شبیه‌سازی مثلث‌ها در ۲۰ زمان متوالی در نظر گرفته شده است:

۱. به‌عنوان اولین گام، مقدار X_1 از توزیع نرمال دو متغیره با میانگین μ_1 و کوواریانس Σ_{11} تولید شود.
۲. با فرض $t = 1$ و با استفاده از توزیع شرطی X_2 به شرط $X_1 = X_1$ مقدار X_2 تولید شود.
۳. با تغییر t به $t + 1$ و هم‌چنین جابه‌جایی X_1 به X_2 دو گام (۱) و (۲)، به اندازه ۲۰ بار تکرار تا داده‌ها شبیه‌سازی شوند.

پس از شبیه‌سازی مثلث‌ها که در واقع نقاطی روی کره هستند مدل (۱) برای آن‌ها به‌کار گرفته می‌شود. به همین منظور، ابتدا فرض کنید برای $i = 1, 2, \dots, 20$ ، دوتایی (θ_i, ϕ_i) نمایش قطبی از نقاط در زمان t باشند. در مدل‌بندی ناپارامتری فرض شد برای $j = 1, 2, 3, 4$ وزن‌های W_j به‌صورت تابع چگالی نرمال با میانگین صفر و واریانس σ_j^2 باشند. مقادیر σ_j^2 با استفاده از تابع اعتبارسنج (۳) برآورد شد که این مقادیر برابر $\hat{\sigma}_1^2 = \hat{\sigma}_2^2 = 0/3$ و $\hat{\sigma}_3^2 = \hat{\sigma}_4^2 = 0/5$ به‌دست آمدند. هم‌چنین مقدار ضریب تعیین رگرسیونی R^2 برای میانگین عرض‌های جغرافیایی (θ_i) برابر $0/9996$ و میانگین طول‌های جغرافیایی (ϕ_i) برابر $0/9998$ حاصل شد. لازم به اشاره است که میانگین توان دوم خطا در روش ناپارامتری برابر $2/27 \times 10^{-5}$ به‌دست آمد.

با استفاده از روش کم‌ترین توان‌های دوم خطا و با تعریف $f_1(x) = -f_2(x) = x$ پارامترها برآورد شدند که مقدار ضریب تعیین تعمیم یافته برای دو مدل (۴) به‌ترتیب برابر $0/9146$ و $0/9972$ به‌دست آمدند. مقدار میانگین توان دوم خطا با استفاده از روش کم‌ترین توان‌های دوم خطا برابر $0/0019$ شد که این مقدار بسیار بزرگ‌تر از مقدار میانگین توان دوم خطا به‌روش ناپارامتری است. شکل ۱ نمایشی از پیش‌بینی مسیر با استفاده از روش ناپارامتری (خطوط ممتد) و هم‌چنین کم‌ترین توان‌های دوم خطا (خطوط بریده) برحسب زمان را نشان می‌دهد.



شکل ۱. مسیر حرکت زوایا (بر حسب رادیان) در طول زمان همراه با پیش‌گویی آن‌ها با استفاده از روش ناپارامتری (خطوط ممتد) و روش کم‌ترین توان‌های دوم خطا (خطوط بریده)

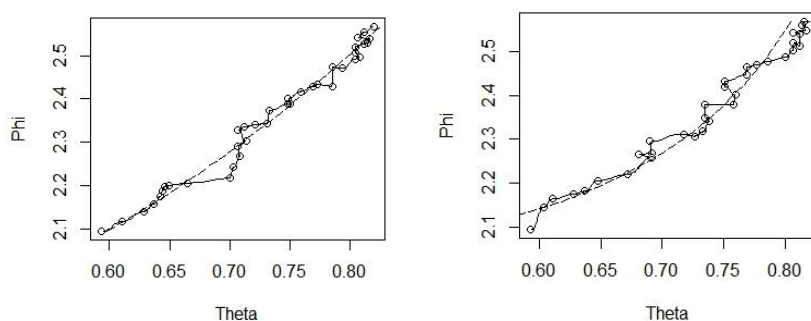
کاربرد در مثال واقعی

به‌منظور بررسی مدل ارائه شده، داده‌های مربوط به مهاجرت فیل‌های دریایی واقع در ناحیه شرق اقیانوس آرام شمالی در نزدیکی سواحل کالیفرنیا جنوبی در نظر گرفته شده است. فیل‌های دریایی ۸ تا ۱۰ ماه از سال را به مهاجرت برای یافتن غذا می‌پردازند. دانشمندان به‌منظور بررسی الگوی مهاجرت آن‌ها با استفاده از سنسورهایی که به بدن آن‌ها متصل شده بودند، موقعیت فیل‌ها را در هر روز از شروع مهاجرت در ۷۵ روز ثبت کرده‌اند. این مهاجرت مربوط به ۴۱ روز مسیر رفت و ۳۴ روز مسیر برگشت است. این داده‌ها قبلاً توسط چند آماردان بررسی شده است که برای مثال می‌توان به [۱۴] اشاره کرد.

از آنجایی که مسیر مهاجرت فیل‌های دریایی دو مسیر رفت و برگشت است؛ مدل‌بندی حرکت‌شان به‌صورت جداگانه برای دو مسیر رفت و برگشت انجام می‌گیرد. به‌علاوه، این مدل‌بندی با استفاده از دو روش ناپارامتری و کم‌ترین توان‌های دوم خطا انجام می‌شود. در روش ناپارامتری ابتدا با استفاده از تابع اعتبارسنج پارامترهای واریانس برآورد شدند. این مقادیر $\hat{\sigma}_1^2 = \hat{\sigma}_2^2 = 0.31$ و $\hat{\sigma}_3^2 = \hat{\sigma}_4^2 = 0.44$ پس از برآورد پارامترهای واریانس، مدل ناپارامتری (۱) که در بخش ۱ تشریح شد برای داده‌های مختصات حرکت فیل‌های دریایی استفاده شد. مقدار میانگین توان دوم خطا در روش ناپارامتری برای مسیر رفت و برگشت به‌ترتیب برابر $4/84 \times 10^{-7}$ و $6/57 \times 10^{-7}$ به‌دست آمد.

مدل‌بندی داده‌ها با استفاده از روش کم‌ترین توان‌های دوم خطا نیز انجام شد که برای مسیر رفت با فرض $f_1(x) = f_2(x) = x$ مقدار میانگین مربعات خطا برابر 0.002 به‌دست آمد. هم‌چنین ضریب تعیین تعمیم یافته برای دو مدل (۴) به‌ترتیب برابر 0.9944 و 0.9818 به‌دست آمد. برای مسیر برگشت با فرض $f_1(x) = x$ و $f_2(x) = x + x^2$ مقدار میانگین دوم خطا برابر 0.1745 به‌دست آمدند. در این مسیر ضریب تعیین هریک از مدل‌های (۴) به‌ترتیب برابر 0.9956 و 0.9598 بود.

شکل ۲ مسیر رفت و برگشت فیل‌های دریایی همراه با برآورد آن‌ها را با استفاده از زوایا نمایش می‌دهد. در این شکل دایره‌ها مختصات حرکت فیل‌ها در هر زمان، خطوط ممتد پیش‌گویی مسیر حرکت آن‌ها با استفاده از روش ناپارامتری و خطوط بریده پیش‌گویی مسیر حرکت با استفاده از روش کم‌ترین توان‌های دوم خطا است.



شکل ۲. سمت راست مسیر رفت و سمت چپ مسیر برگشت (زوایا برحسب رادیان) از حرکت فیل‌های دریایی همراه با پیش‌گویی آن‌ها با استفاده از روش ناپارامتری (خطوط ممتد) و کم‌ترین توان‌های دوم خطا (خطوط بریده)

بحث و نتیجه‌گیری

در سال‌های اخیر محققان با داده‌هایی روبه‌رو می‌شوند که طبیعت ناقطری دارند. برای مثال داده‌هایی که روی کره زمین قرار دارند، داده‌هایی هستند که به محدودیت نرم مشخص بودن آن‌ها باید توجه کرد. مدل‌بندی این داده‌ها با رویکردهای آماری یکی از دغدغه‌های اصلی محققان بوده و هست. در مقاله حاضر، سعی شد به بخشی از موضوع تحلیل داده‌های کروی که با زوایا مشخص می‌شوند پاسخ داده شود. از این‌رو، رویکرد مدل‌بندی رگرسیونی براساس زوایای مستقل طوری تطبیق داده شد تا دو روش مدل‌بندی رگرسیونی طولی ناپارامتری و کمترین توان‌های دوم خطا برای تحلیل داده‌های کروی ارائه شود.

در دیدگاه روش ناپارامتری به مدل‌بندی رگرسیونی داده‌های کروی، تابع مخاطره هاورسین معرفی و کمیته‌سازی آن مدنظر قرار گرفت. به‌علاوه، ارزیابی آن با استفاده از داده‌های شبیه‌سازی شده و مثال واقعی بررسی شد. همچنین مدل‌بندی رگرسیونی به‌روش کمترین توان‌های دوم خطا با استفاده از تابع ربط مناسب ارائه شد. اگرچه این دیدگاه کارایی روش ناپارامتری را ندارد؛ اما قادر است روش مناسبی در پیش‌گویی مسیر حرکت روی کره ارائه کند. روش‌های ارائه شده به‌دلیل سادگی می‌توانند در مسائل کاربردی مورد اقبال محققان قرار گیرند. استفاده از فاصله هاورسین نیز به‌عنوان فاصله‌ای مناسب داده‌های کروی می‌تواند بر مقبولیت روش پیشنهادی بیافزاید. به‌منظور افزایش کارایی مدل‌های معرفی شده پیشنهاد می‌شود از دیگر توابع ربط مبتنی بر روش کمترین توان‌های دوم خطا و توابع مخاطره مناسب برای داده‌های کروی استفاده شود.

منابع

1. Gould, A. L. "A regression technique for angular variates", *Biometrics*, 25 (1969) 683-700.
2. Thompson R., Clark R. M. "Fitting polar wander paths. *Physics of the Earth and Planetary Interiors*" (1981) 27, 1-7.
3. Thompson R., Clark R. M. "A robust least-squares Gondwanan apparent polar wander path and the question of palaeomagnetic assessment of Gondwanan reconstruction", *Earth and Planetary Science Letters*, 57 (1982) 152-158.
4. Fisher, N. I. Lewis T., Embleton B, J, J., "Statistical Analysis of Spherical Data", New York: Cambridge University Press (1987).
5. Parker R. L., Denham C. R., "Interpolation of unit vectors", *Geophysical Journal International*, 58 (1979) 685-687.
6. Di Marzio M., Panzera A., Taylor C. C., "Non-parametric regression for circular responses", *Scandinavian Journal of Statistics*, 40 (2013) 238-255.
7. Abramowitz M., Stegun I. A., "Handbook of Mathematical Functions: With Formulas, Graphs, and Mathematical Tables", New York: Dover Publications (1965).
8. Ramsay J. O., Silverman B. W., "Applied Functional Data Analysis Methods and Case Studies", New York: Springer-Verlag (2002).

9. Habbema J. D. F., Hermans J., Van den Broek K., "A stepwise discriminant analysis program using density estimation", *Compstat 1974, Proceedings in computational statistics*, PhysicaVerlag, Wien. (1974) 101-110.
10. Duin R. P. W., "On the choice of smoothing parameters for parzen estimators of probability density functions", *IEEE Transactions on Computers C-25* (1976) 1175-1179.
11. Fisher N. I. and Lee A. J., "Regression models for an angular response", *Biometrics*, 48 (1992) 665-677.
12. Kendall D.G., "The diffusion of shape", *Advances in Applied Probability*, 9 (1977) 428-430.
13. Dryden I. L., Mardia K. V., "Statistical Shape Analysis", Chichester: John Wiley & Sons (1998).
14. Brillinger D. R., Stewart B. S., "Elephant-seal movements: Modelling migration", *Canadian Journal of Statistics*, 26 (1998) 431-443.