

مدل‌بندی داده‌های فازی با رگرسیون اسپلاین تطبیقی چندگانه

جلال چاچی^۱، غلامرضا حسامیان^۲

^۱ گروه آمار، دانشگاه سمنان

^۲ گروه آمار، دانشگاه پیام‌نور

تاریخ دریافت: ۱۳۹۲/۷/۳۰ تاریخ آخرین بازنگری: ۱۳۹۳/۳/۱۸

چکیده: در این مقاله به مدل‌بندی داده‌های ورودی دقیق-خروجی فازی پرداخته می‌شود و رویکرد رگرسیون مارس فازی با پارامترهای دقیق و جملات خطای فازی معرفی می‌گردد. روش پیشنهادی شامل دو مرحله است: در مرحله اول با استفاده از رگرسیون اسپلاین تطبیقی چندگانه (مارس) مراکز متغیر وابسته برآورد می‌شوند، و در مرحله دوم کمترین مقادیر خطاهای فازی بر اساس یک مسأله بهینه‌سازی غیر خطی به دست می‌آیند. در انتها کاربرد مدل پیشنهاد شده در مدل‌بندی داده‌های واقعی در مهندسی آب نشان داده می‌شود. نتایج تجربی این مثال برتری روش پیشنهادی را در مقایسه با برخی از روش‌های متداول رگرسیون فازی کمترین توان‌های دوم خطا نشان می‌دهد.

واژه‌های کلیدی: رگرسیون اسپلاین تطبیقی چندگانه (مارس)، داده‌های فازی، سامانه استنتاج فازی، دبی رودخانه.

آدرس الکترونیک مسئول مقاله: جلال چاچی، jchachi@profs.semnan.ac.ir

کلمه موضوعی: بندی ریاضی (۲۰۱۰): ۶۲۷۸۶، ۶۲۷۰۸، ۹۴D۰۵

۱ مقدمه

مدل‌های رگرسیونی (آماری و فازی) الگوهایی فراهم می‌آورند که می‌توان بر پایه آنها ارتباط بین مجموعه‌ای از متغیرهای تبیینی و پاسخ را بررسی کرد. بر خلاف رگرسیون آماری، رگرسیون فازی را می‌توان در موارد زیر به کار برد:

- داده‌ها نادقیق باشند؛
- رابطه بین متغیرها از نوع امکانی (و نه احتمالی) باشد؛
- فرضیات زیربنایی در مدل‌های رگرسیون آماری مورد تردید باشند.

در رگرسیون فازی، همانند رگرسیون آماری، یک فرم تابعی بین مشاهدات متغیرهای تبیینی و پاسخ و به منظور کنترل و پیش‌بینی مقادیر متغیر پاسخ در نظر گرفته می‌شود. بر این اساس و به منظور کنترل ابهام مقدار پیش‌بینی یا فاصله آن با مقدار واقعی متغیر پاسخ، شیوه‌های رگرسیون فازی را می‌توان از لحاظ رویکرد به دو نوع کلی طبقه‌بندی کرد:

(۱) **رویکردهای موسوم به رگرسیون امکانی:** در رگرسیون امکانی خطای مدل در قالب توزیع‌های امکانی ضرایب مدل منظور می‌شود. در این مدل‌ها خطای پیش‌بینی متغیر وابسته تعبیر امکانی دارد. به سخن دیگر هنگامی که بر پایه مدل رگرسیونی و به‌ازای مقادیر متغیرهای تبیینی مقدار متغیر پاسخ را پیش‌بینی می‌کنیم، مقادیر پیش‌بینی تعبیر امکانی دارد و نه تعبیر احتمالی. یافتن ضرایب این گونه مدل‌ها اغلب مستلزم استفاده از روش‌های برنامه‌ریزی خطی یا غیرخطی است. در این رویکردها ابهام کل مدل که برابر با مجموع پهنای‌های مقادیر برآورد شده متغیر پاسخ است، تحت شرایط و قیودی کمینه می‌شود.

(۲) **شیوه‌های مبتنی بر روش کمترین توان‌های دوم:** در این رویکردها بر پایه تعریف‌هایی برای فاصله بین اعداد فازی روش کمترین توان‌های دوم کلاسیک تعمیم می‌یابد و بر اساس آن ضرایب مدل رگرسیون فازی برآورد می‌شوند. به عبارتی در این شیوه، مجموع فاصله‌ها بین مقادیر متغیر پاسخ

مشاهده شده و مقادیر متغیر پاسخ برآورد شده، بر اساس تعریف یک فاصله بین اعداد فازی، کمینه می‌شود. این رویکردها نیز تنوع بسیاری دارند. این تنوع عوامل مختلفی دارد، از جمله: دقیق یا نادقیق بودن مشاهدات متغیرها، تنوع در تعریفی که برای فاصله بین مجموعه‌های فازی در نظر گرفته می‌شود، دقیق یا نادقیق منظور نمودن عرض از مبدا و برخی عوامل دیگر است.

از آنجا که رگرسیون فازی از جمله موضوعاتی است که بسیار مورد توجه و مطالعه محققان قرار گرفته است و تنوع گسترده‌ای در روش‌ها و رویکردهای معرفی شده وجود دارد، می‌توان برای مطالعه و مرور روش‌های معرفی شده در این زمینه می‌توان به ارقامی (۱۳۸۱)، میرزایی و ارقامی (۱۳۸۶)، چاچی (۱۳۹۱) و دورسو و همکاران (۲۰۱۰، ۲۰۱۱) مراجعه نمود. برخی از این رویکردها از جهات مختلف مورد بررسی و نقد قرار گرفته‌اند. در جهت رفع مشکلات بالا تحقیقات بسیاری صورت گرفته که می‌توان به چن و دنگ (۲۰۰۸) و لو و ونگ (۲۰۰۹) رجوع کرد.

در ادامه برای رفع برخی از انتقادات مطرح شده، از روش رگرسیون اسپلاین تطبیقی چندگانه (مارس^۱) به عنوان جایگزینی برای روش کمترین توان‌های دوم در محیط فازی استفاده می‌شود. روش مارس یک شیوه مدل‌بندی رگرسیون قطعه‌ای ناپارامتری است که هیچ فرضیه زیربنایی درباره رابطه تابعی بین متغیرهای تبیینی و پاسخ در نظر نمی‌گیرد. این روش برای مدل کردن داده‌هایی که دامنه وسیعی دارند، یا رابطه بین متغیرها خطی نمی‌باشد، بسیار مناسب است (فریدمن، ۱۹۹۱؛ هستی و همکاران، ۲۰۰۹). در طول چند دهه اخیر روش مارس با بسیاری از روش‌های پارامتری و ناپارامتری متداول (از جمله روش کمترین توان‌های دوم خطی) از لحاظ دقت برازش، کارایی، استواری، و سادگی محاسبات مورد مقایسه قرار گرفته است (دی-آندرس و همکاران، ۲۰۱۱؛ کرینر، ۲۰۰۷؛ لی و چن، ۲۰۰۵؛ لی و همکاران، ۲۰۰۶).

در این مقاله با روش مارس داده‌های فازی مدل‌بندی می‌شوند و روش پیشنهادی با روش‌های مشابه در رگرسیون فازی با رویکرد کمترین توان‌های دوم مورد مقایسه

^۱ Multivariate Adaptive Regression Splines (MARS)

قرار می‌گیرد. از آنجا که روش مارس در مقایسه با روش کمترین توان‌های دوم نسبت به داده‌های پرت استوار است، به نظر می‌رسد که روش پیشنهادی برتری قابل توجهی به روش کمترین توان‌های دوم در مدل‌بندی داده‌هایی که در آنها مشاهدات پرت وجود داشته باشد، دارد. روش پیشنهادی یک روش دو مرحله‌ای است. در مرحله اول، با استفاده از مارس مراکز متغیر پاسخ برآورد می‌شوند. سپس در مرحله دوم و با استفاده از روش‌های بهینه‌سازی، جملات خطای فازی متناظر با داده‌ها به دست می‌آیند. در این روش ضرایب مدل دقیق هستند و جملات خطای فازی برای در نظر گرفتن ابهام به مدل افزوده می‌شوند.

در بخش ۲، انگیزه‌ها و اهداف ارائه مدل جدید بیان می‌شود. در بخش ۳، مدل رگرسیون مارس فازی با پهنای متغیر بیان و تشریح می‌شود. نحوه پیش‌بینی یک مقدار جدید با استفاده از سامانه استنتاج فازی در بخش ۴ بیان و تشریح می‌شود. دو ملاک نیکویی برازش برای مدل‌های رگرسیون فازی در بخش ۵ ذکر می‌شوند. در بخش ۶، با استفاده از داده‌های واقعی در مهندسی آب یک مطالعه مقایسه‌ای بین روش پیشنهادی و چند روش مطرح در رگرسیون فازی انجام می‌شود. در انتها به بحث و نتیجه‌گیری پرداخته می‌شود.

۲ اهداف

تا کنون رویکردهای مختلف و متنوعی در زمینه رگرسیون فازی ارائه شده است. بسیاری از این رویکردها از جهات مختلف مورد بررسی و نقد قرار گرفته‌اند (چاچی، ۱۳۹۱). مشکل اساسی بسیاری از روش‌های موجود عبارتند از:

- مدل‌های رگرسیون امکانی و مدل‌های مبتنی بر کمترین توان‌های دوم فازی نسبت به داده‌(های) پرت حساس هستند؛
- در برخی از روش‌های موجود با افزایش مقدار متغیر مستقل پهنای برآورد شده متغیر پاسخ افزایش می‌یابد، که به آن مشکل افزایش پهنای گفته می‌شود؛
- در برخی از روش‌های مدل‌بندی مشاهدات فازی متغیر پاسخی که پهنای مشاهده شده آن روند صعودی، نزولی یا متغیر دارند، ممکن است برآوردهایی

با خطای زیاد حاصل شوند.

اگرچه برای حل انتقادات فوق برخی روش‌ها در رگرسیون فازی ارائه شده است، اما با وجود این هنوز روشی که به‌طور کامل این مشکلات را برطرف کند، ارائه نشده است. برای مثال، روش‌های زیر در برطرف نمودن حساسیت مدل نسبت به داده پرت معرفی شدند:

- رویکردهای رگرسیون فازی مبتنی بر روش کمترین قدرمطلق خطا (چاچی، ۱۳۹۱؛ کلکین نما و طاهری، ۲۰۱۲)؛

- رویکردهایی که در آنها از روش‌های رگرسیون آماری استوار، از جمله M -برآوردگرها، S -برآوردگرها و غیره، استفاده می‌شود (دورسو و همکاران، ۲۰۱۱).

اما در رویکردهای بالا مسأله افزایش پهناها همچنان حل نشده است. از سوی دیگر، در زمینه برطرف کردن مسأله افزایش پهناها نیز می‌توان به روش‌های زیر اشاره کرد:

- روش‌هایی که در آنها اغلب ضرایب مدل دقیق هستند، یا
- یک جمله فازی به‌منظور در نظر گرفتن ابهام فازی به مدل افزوده می‌شود (چن و دنگ، ۲۰۰۸؛ لو و ونگ، ۲۰۰۹).

این روش‌ها نیز در چارچوب رگرسیون کمترین توان‌های دوم فازی قرار دارند و از جهت حساسیت نسبت به داده پرت مورد انتقاد هستند. در میان روش‌های فوق باید به روش چن و دنگ (۲۰۰۸) اشاره نمود که یک روش رگرسیون فازی با پهناهای متغیر است که مسأله افزایش پهناها را برطرف می‌کند. این روش همچنین، برآوردهای خوبی را برای مشاهدات فازی متغیر پاسخی که پهناهای آن روند صعودی، نزولی، ثابت و حتی متغیر دارند ارائه می‌دهد. اما همچنان که گفته شد، این روش نیز با مشکل تأثیرپذیری زیاد نسبت به داده پرت مواجه است.

در ادامه یک مدل رگرسیون فازی ارائه می‌شود که نسبت به داده‌های (های) پرت استوار است و دیگر این که مسأله افزایش پهناها را برحسب افزایش مقادیر متغیر تبیینی و پهناهای متغیر پاسخ همزمان برطرف می‌کند.

۳ بیان و تشریح مدل

فرض کنید مجموعه‌ای از مشاهدات مربوط به یک یا چند متغیر تبیینی و یک متغیر پاسخ در اختیار داریم که مشاهدات متغیر(های) تبیینی دقیق هستند ولی مشاهدات متغیر پاسخ نادقیق (فازی) هستند. برای مدل‌بندی چنین داده‌هایی یک مدل رگرسیون فازی با پهناهای متغیر ارائه می‌شود. مدل مربوطه به صورت

$$\tilde{y}_i = \beta_0 + \sum_{m=1}^M \beta_m B_m(x_i) \oplus \tilde{\varepsilon}_i, \quad i = 1, \dots, n$$

در نظر گرفته می‌شود، که در آن $\tilde{y}_i = (y_i, y_i^l, y_i^r)_{LR}$ i امین مشاهده فازی متغیر پاسخ، $x_i = [x_{0i}, x_{1i}, \dots, x_{ki}] \in \mathbb{R}^{k+1}$ ، $(i = 1, \dots, n; k < n; x_{0i} = 1)$ ، i امین مشاهده حقیقی مقدار متغیرهای تبیینی، β_m ، $m = 0, 1, \dots, M$ ، ضرایب دقیق مدل مارس برای توابع پایه^۲ $B_m(x_i)$ و $\tilde{\varepsilon}_i = (\varepsilon_i, \ell_i, r_i)_{LR}$ خطای فازی مدل برای مشاهده i ام است.

برای برآورد پارامترها $(\beta_m, m = 0, 1, \dots, M)$ و جملات خطای فازی مدل $(\tilde{\varepsilon}_i = (\varepsilon_i, \ell_i, r_i)_{LR}, i = 1, \dots, n)$ یک روش دو مرحله‌ای معرفی می‌شود. در مرحله اول ضرایب دقیق $\beta = [\beta_0, \beta_1, \dots, \beta_M]^t$ برآورد می‌شوند، و در مرحله دوم جملات خطای فازی بر اساس یک مسأله بهینه‌سازی به دست می‌آیند.

مرحله ۱. ضرایب رگرسیون $\beta = [\beta_0, \beta_1, \dots, \beta_M]^t$ با مدل‌بندی مراکز متغیر پاسخ فازی بر روی مقادیر متغیر تبیینی، از طریق مدل مارس

$$y_i = \beta_0 + \sum_{m=1}^M \beta_m B_m(x_i) + \varepsilon_i, \quad i = 1, \dots, n$$

^۲ Basis functions

به دست آورده می شود. ضرایب مدل در روش مارس با اجرای برنامه های earth یا mars در نرم افزار R برآورد می شوند (فوکس و ویزبرگ، ۲۰۱۱). سپس مراکز متغیر پاسخ فازی به صورت زیر به دست می آیند.

$$\hat{y}_i = \hat{\beta}_0 + \sum_{m=1}^M \hat{\beta}_m B_m(x_i), \quad i = 1, \dots, n$$

مرحله ۲. عبارت خطای فازی $\tilde{e}_i = (0, \ell_i, r_i)_{LR}$ به ازای $i = 1, \dots, n$ ، به گونه ای برآورد می شود که:

الف- پهنای برآورد شده هر \tilde{e}_i یعنی $\ell_i + r_i$ ، مساوی پهنای مشاهده شده \tilde{y}_i ، یعنی $y_i^l + y_i^r$ باشد،

ب- برآورد مشاهده \tilde{y}_i ، یعنی \hat{y}_i ، کمترین خطای برآورد را داشته باشد.

برای این منظور خطاهای فازی را از طریق حل یک مسأله برنامه ریزی غیرخطی و با در نظر گرفتن شرط اول به عنوان قیود آن و شرط دوم به عنوان تابع هدف آن به دست آورده می شود که در آن تابع هدف مجموع اندازه های تشابه بین توابع عضویت مشاهده شده و برآورد شده متغیر پاسخ به صورت (لو و ونگ، ۲۰۰۹)

$$\sum_{i=1}^n \frac{\int \min\{\tilde{y}_i(x), \hat{y}_i(x)\} dx}{\int \max\{\tilde{y}_i(x), \hat{y}_i(x)\} dx}$$

است و قیود به گونه ای اختیار می شوند که \hat{y}_i پهنای برآورد شده متغیر پاسخ مساوی \hat{y}_i پهنای مشاهده شده آن باشد. یعنی شرط $\ell_i + r_i = y_i^l + y_i^r$ برای هر $i = 1, \dots, n$ ، برقرار باشد. این قید به این دلیل در نظر گرفته می شود که پهنای برآورد شده کنترل شود، یا به عبارتی مقادیر پهنای خیلی کم یا خیلی زیاد برای متغیر پاسخ برآورد نشود (چن و دنگ، ۲۰۰۸). برای در نظر گرفتن این قیود فرض می شود y_m^l و y_m^r به ترتیب کمترین مقدار پهنای چپ و راست مشاهدات متغیر وابسته باشند، یعنی

$$y_m^l = \min\{y_1^l, \dots, y_n^l\}, \quad y_m^r = \min\{y_1^r, \dots, y_n^r\}. \quad \text{www.SID.ir}$$

۸ مدل‌بندی داده‌های فازی با رگرسیون اسپلاین تطبیقی چندگانه

قسمتی از پهنای برآورد شده متغیر پاسخ مقدار ثابت $y_m^l + y_m^r$ اختیار می‌شود. بدیهی است که

$$D_i = (y_i^l + y_i^r) - (y_m^l + y_m^r) \geq 0, \quad i = 1, \dots, n.$$

اما این مقدار ثابت باید از دو طرف گسترش یابد تا در نهایت مساوی با i امین پهنای مشاهده شده متغیر وابسته شود. برای این منظور، مقدار D_i به دو قسمت d_i و $D_i - d_i$ تقسیم می‌شود، که در آن $0 \leq d_i \leq D_i$. بنابراین، قیود مسأله بهینه‌سازی غیرخطی به صورت

$$\tilde{\varepsilon}_i = (0, y_m^l + d_i, y_m^r + D_i - d_i)_{LR}, \quad 0 \leq d_i \leq D_i, \quad i = 1, \dots, n,$$

در نظر گرفته می‌شود. سرانجام، با توجه به تابع هدف و قیود بیان شده مسأله بهینه‌سازی غیرخطی که از طریق آن مقادیر d_1, \dots, d_n برآورد می‌شوند به صورت

$$\max_{d_1, \dots, d_n} \sum_{i=1}^n \int \min\{\tilde{y}_i(x), \hat{y}_i(x)\} dx, \quad \hat{y}_i = \hat{\beta}_0 + \sum_{m=1}^M \hat{\beta}_m B_m(x_i) \oplus \tilde{\varepsilon}_i,$$

$$\tilde{\varepsilon}_i = (0, y_m^l + d_i, y_m^r + D_i - d_i)_{LR}, \quad 0 \leq d_i \leq D_i, \quad i = 1, \dots, n.$$

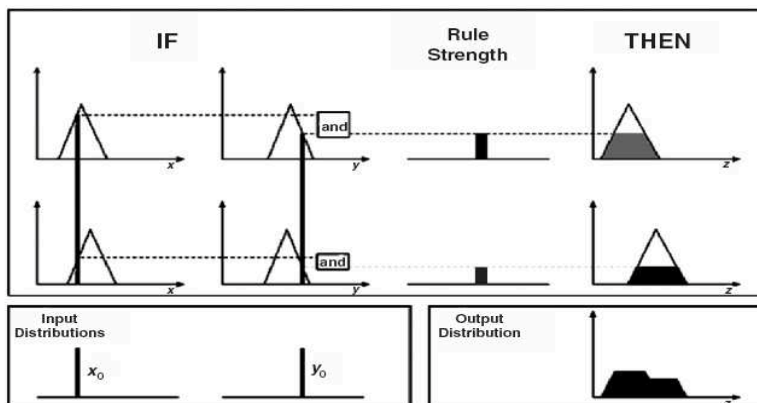
در نظر گرفته می‌شود. پس از تعیین مقادیر $\hat{d}_1, \dots, \hat{d}_n$ مدل رگرسیون مارس فازی به صورت زیر حاصل می‌شود

$$\begin{aligned} \hat{y}_i &= \hat{\beta}_0 + \sum_{m=1}^M \hat{\beta}_m B_m(x_i) \oplus \hat{\varepsilon}_i \\ &= \hat{\beta}_0 + \sum_{m=1}^M \hat{\beta}_m B_m(x_i) \oplus (0, y_m^l + \hat{d}_i, y_m^r + D_i - \hat{d}_i)_{LR}, \quad i = 1, \dots, n. \end{aligned}$$

۴ پیش‌بینی متغیر پاسخ با سامانه استنتاج فازی

پس از تعیین مقادیر $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_M$ و $\hat{d}_1, \dots, \hat{d}_n$ مدل رگرسیون مارس فازی با پهنای متغیر به صورت

$$\hat{y}_i = \hat{\beta}_0 + \sum_{m=1}^M \hat{\beta}_m B_m(x_i) \oplus (0, y_m^l + \hat{d}_i, y_m^r + D_i - \hat{d}_i)_{LR}, \quad i = 1, \dots, n$$



شکل ۱: سامانه استنتاج فازی با دو متغیر ورودی دقیق و دو قانون اگر-آنگاه فازی

نتیجه می‌شود. اکنون با استفاده از این مدل می‌توان به پیش‌بینی مقادیر متغیر پاسخ بر اساس مشاهدات جدیدی از متغیر تبیینی پرداخت. برای پیش‌بینی، چن و دنگ (۲۰۰۸) از سامانه استنتاج فازی استفاده کردند که در ادامه توضیح داده می‌شود و برای پیش‌بینی متغیر پاسخ بر اساس مدل پیشنهاد شده در این مقاله از این شیوه استفاده می‌شود. ساختار اصلی یک سامانه استنتاج فازی از سه بخش اصلی تشکیل می‌شود (زیمرمن، ۲۰۰۱):

- الف- بخش قوانین: شامل مجموعه‌ای از قوانین اگر-آنگاه فازی است؛
- ب- بخش مجموعه داده‌ها یا اطلاعات فازی داده‌ها: شامل توابع عضویت داده‌های فازی به کار رفته در قوانین اگر-آنگاه فازی است؛
- ج- فرآیند استنتاج: شامل قواعد تصمیم‌گیری بر اساس اطلاعات به دست آمده و روش الحاق نتایج حاصله از قوانین اگر-آنگاه فازی است.

در شکل ۱ سامانه استنتاج فازی ممدانی با دو متغیر ورودی دقیق و دو قانون اگر-آنگاه فازی نشان داده شده است.

فرض کنید $x^* = [1, x_1^*, \dots, x_k^*]$ مقدار جدید متغیر تبیینی و $Y = \hat{\beta}_0 + \sum_{m=1}^M \hat{\beta}_m B_m(x_k^*)$ پیش‌بینی مرکز متغیر پاسخ بر اساس مدل مارس

۱۰ مدل‌بندی داده‌های فازی با رگرسیون اسپلاین تطبیقی چندگانه

باشد. اکنون $\hat{\epsilon}$ جمله خطای فازی متناظر با این مشاهده برای مدل

$$\tilde{Y} = \hat{\beta}_0 + \sum_{m=1}^M \hat{\beta}_m B_m(x_i^*) \oplus \hat{\epsilon},$$

از طریق یک سامانه استنتاج فازی ممدانی با یک متغیر ورودی-خروجی به دست آورده می‌شود. فرض کنید

$$(\tilde{y}_a, \hat{\epsilon}_a) \quad a \in A = \{i : \tilde{y}_i(Y) > 0, i = 1, \dots, n\},$$

مشاهداتی از متغیر پاسخ به همراه جملات خطای فازی متناظر با آنها باشند که توسط مقدار جدید x^* فعال می‌شوند. حال a امین قانون اگر-آنگاه فازی به صورت

$$R^a: \text{ اگر } Y \text{ به صورت } \tilde{y}_a \text{ باشد، آنگاه } \hat{\epsilon} \text{ به صورت } \hat{\epsilon}_a \text{ است،}$$

نوشته می‌شود، که در آن $a \in A$ و R^a قانون a ام است. در سامانه استنتاج فازی مقدم و تالی هر قانون اگر-آنگاه فازی به صورت گزاره‌های فازی هستند. به علاوه، بر اساس قوانین در نظر گرفته شده خروجی این سامانه نیز به صورت یک مجموعه فازی است. اکنون برای مقدار Y ، هر تابع عضویت \tilde{y}_a در خروجی به میزان درجه عضویت $\tilde{y}_a(Y)$ فعال می‌شود، که از آن در به دست آوردن جمله خطای فازی استفاده می‌شود. در نهایت، با الحاق خروجی‌های تمام قوانینی که فعال شده‌اند، جمله خطای فازی به دست می‌آید.

توجه شود که حتی اگر تمام مشاهدات متغیر پاسخ اعداد فازی LR باشند، تابع عضویت به دست آمده برای خطای فازی هر شکل نامنظمی می‌تواند داشته باشد (شکل ۱). اما مناسب است که خطای برآورد شده نیز به صورت یک عدد فازی LR به دست آید، زیرا معمولاً ترجیح داده می‌شود که مقدار پیش‌بینی نیز به صورت یک عدد فازی LR باشد. برای این منظور، خطای فازی تبدیل به یک عدد فازی LR می‌شود تا پیش‌بینی متغیر پاسخ نیز یک عدد فازی LR شود. بنابراین، خطای فازی به صورت عدد فازی LR $\hat{\epsilon} = (\hat{\epsilon}, \hat{\ell}, \hat{r})$ در نظر گرفته می‌شود که در آن $\hat{\epsilon}$ مقدار غیرفازی شده عبارت خطای فازی بر اساس روش مرکز ثقل است و $\hat{\ell}$ و \hat{r} کمترین و بیشترین مقادیر ممکن برای پهناهای چپ و راست عبارت خطای فازی هستند.

سرانجام، برای مقدار جدید x^* مقدار پیش‌بینی متغیر پاسخ به صورت

$$\hat{Y}_{x^*} = Y \oplus \hat{\varepsilon} = (Y + \hat{\varepsilon}, \hat{\ell}, \hat{r})_{LR},$$

به دست می‌آید.

۵ ملاک‌های ارزیابی مدل

برای ارزیابی مدل‌های رگرسیون فازی چندین ملاک توسط مولفان مورد استفاده قرار گرفته است (چن و دنگ، ۲۰۰۸؛ لو و ونگ، ۲۰۰۹؛ کلکین‌نما و طاهری، ۲۰۱۲). در ادامه دو ملاک متداول برای بررسی نیکویی برازش مدل‌های رگرسیون فازی یادآوری می‌شوند. از این دو ملاک برای مقایسه مدل رگرسیونی پیشنهاد شده در این مقاله و چند رویکرد دیگر در زمینه رگرسیون فازی استفاده می‌شود. این دو ملاک میانگین اندازه‌های تشابه^۳ و میانگین قدرمطلق خطاهای^۴ برآورد هستند که به ترتیب به صورت زیر تعریف می‌شوند

$$MSM = \frac{1}{n} \sum_{i=1}^n \frac{\int \min\{\tilde{y}_i(x), \hat{y}_i(x)\} dx}{\int \max\{\tilde{y}_i(x), \hat{y}_i(x)\} dx},$$

$$MAE = \frac{1}{n} \sum_{i=1}^n \int |\tilde{y}_i(x) - \hat{y}_i(x)| dx.$$

شایان ذکر است که در یک مدل رگرسیون فازی MAE می‌تواند هر مقدار بزرگتر از صفر باشد، ولی MSM همواره عددی بین صفر و یک است. اگر $MAE = 0$ یا $MSM = 1$ باشد، در این صورت مدل رگرسیون فازی مربوطه برازش کامل به داده‌ها دارد. بنابراین در مقایسه بین دو مدل رگرسیون فازی هر مدلی که مقدار MAE کوچکتری (نزدیک‌تر به صفر) داشته باشد، یا مقدار MSM بزرگتری (نزدیک‌تر به یک) داشته باشد، آن مدل برازش بهتری به داده‌ها دارد.

^۳ Mean of similarity measures

^۴ Mean of absolute errors

۶ مثال کاربردی

یکی از مسایل مهم در مهندسی آب اندازه‌گیری دبی یا بار معلق^۵ و سرعت آب^۶ در حوزه‌های آبریز است. در مطالعات کاربردی برآورد دبی بر اساس سرعت آب اهمیت زیادی دارد. بر اساس مطالعه‌ای در دریند واقع در خراسان شمالی با استفاده از ابزارهای استاندارد، برخی خصوصیات آب‌شناسی حوزه‌های آبریز منطقه ثبت گردید. در این خصوص از ایستگاههای مختلف تعداد ۱۸۴ جفت مشاهده مربوط به دبی و سرعت آب حوزه‌های آبریز منطقه به‌دست آمد. بنا به محدودیت‌های آزمایشگاهی و ابزارهای اندازه‌گیری و ماهیت متغیر پاسخ، این مشاهدات به‌صورت اعداد فازی مثلثی در جدول ۱ گزارش شده‌اند. نمودار پراکنش مراکز متغیر پاسخ بر حسب متغیر مستقل به همراه برآورد مارس آنها در شکل ۲ نشان داده شده است. با توجه به این نمودار واضح است که یک خط راست دقت کافی برای مدل‌بندی این داده‌ها را ندارد. در ادامه با رگرسیون مارس فازی به برآورد متغیر پاسخ فازی $(\tilde{y} = (y, s)_T)$ بر اساس مقادیر متغیر تبیینی (x) پرداخته می‌شود. اما بر خلاف شیوه‌های مدل‌سازی پارامتری متداول، در روش مارس ضرایب مدل برای تمام دامنه مشاهدات متغیرها یکسان در نظر گرفته نمی‌شود و به‌جای آن از خطوط شکسته‌ای^۷ به منظور برازش یک تابع پیوسته چند ضابطه‌ای بهینه به داده‌ها استفاده می‌شود. این رویکرد برای مدل‌بندی مسائلی که در آنها تعداد متغیرها بسیار زیاد است یا با حجم وسیعی از داده‌ها روبرو هستیم مفید است.

با به‌کار بردن روش پیشنهاد شده در این مقاله، مدل

$$\begin{aligned} \hat{y}_i = & 2/82 + 2/27 \max\{0, x_i + 0/24\} - 0/76 \max\{0, -0/24 - x_i\} \\ & - 11/98 \max\{0, x_i - 2/17\} + 15/98 \max\{0, x_i - 2/32\} \\ & \oplus (0, 0/03 + \hat{d}_i, 0/03 + D_i - \hat{d}_i)_T, \quad i = 1, \dots, 184, \end{aligned}$$

^۵ Suspended load

^۶ Discharge

^۷ Splines

جدول ۱: داده‌های مهندسی آب

\hat{d}_i	$\tilde{y}_i = (y_i, s_i)_T$	x_i	i
۰/۳۳	$(۲/۷۵, ۰/۵۵)_T$	-۰/۵۱	۱
۰/۴۰	$(۲/۴۳, ۰/۲۳)_T$	-۰/۲۱	۲
۰/۵۰	$(۳/۲۵, ۱/۱۳)_T$	-۰/۲۴	۳
۰/۳۸	$(۳/۵۳, ۱/۴۹)_T$	-۰/۳۳	۴
۱/۳۹	$(۳/۶۵, ۱/۶۹)_T$	۰/۰۴	۵
\vdots	\vdots	\vdots	\vdots
۲/۲۶	$(۶/۹۰, ۴/۳۸)_T$	۰/۹۰	۱۸۰
۰/۴۸	$(۷/۶۶, ۵/۰۰)_T$	۰/۴۶	۱۸۱
۳/۱۸	$(۴/۹۰, ۲/۷۸)_T$	۰/۸۱	۱۸۲
۳/۸۳	$(۶/۹۱, ۴/۲۵)_T$	۱/۴۴	۱۸۳
۴/۵۳	$(۷/۸۸, ۵/۷۶)_T$	۲/۳۰	۱۸۴

برای این داده‌ها به دست می‌آید. برآوردهای \hat{d}_i ، $i = 1, \dots, 184$ در جدول ۱ آورده شده‌اند.

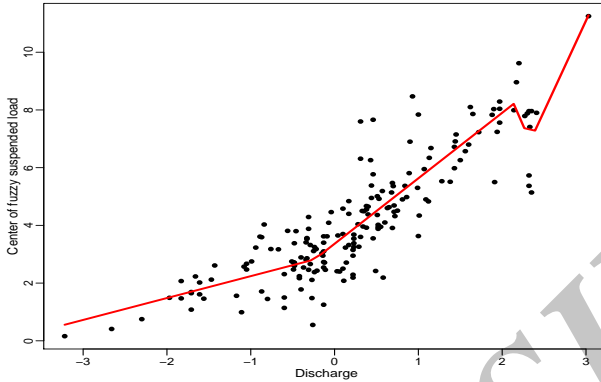
در ادامه روش پیشنهاد شده در این مقاله با رویکردهای ژو و لی (۲۰۰۱) (XL) ، چن و دنگ (۲۰۰۸) (CD) ، و فرارو و همکاران (۲۰۱۰) (F) در رگرسیون فازی مقایسه می‌شود. توجه شود که در روش ژو و لی (۲۰۰۱) مقادیر متغیر تبیینی نامنفی در نظر گرفته می‌شوند و در اینجا بدون اینکه از کلیت مساله کاسته شود، داده‌های متغیر تبیینی به $x'_i = x_i - \min_i\{x_i\} = x_i - (-۳/۲۲) \geq ۰$ تبدیل می‌شود. این سه مدل به صورت

$$\hat{y}_i^{XL} = (-۱/۵۳ + ۱/۶۴ x'_i, \max\{۰, -۲/۲۰ + ۱/۲۲ x'_i\})_T,$$

$$\hat{y}_i^{CD} = ۳/۷۶ + ۱/۶۴ x_i \oplus \hat{\varepsilon}_i,$$

$$\hat{y}_i^F = (۳/۷۶ + ۱/۶۴ x_i, \exp\{۰/۱۰ + ۰/۷۴ x_i\})_T,$$

به دست می‌آیند. نتایج نیکویی برآزش سه مدل بالا به همراه مدل رگرسیون مارتس فازی در جدول ۲ بیانگر برتری مدل رگرسیون مارتس فازی بر سه روش دیگر است.



شکل ۲: نمودار پراکنش مراکز متغیر پاسخ بر حسب مقادیر متغیر تبیینی به همراه برآورد مارس آنها

جدول ۲: مقادیر نیکویی برازش مدل‌های رگرسیون فازی

<i>MAE</i>	<i>MSM</i>	مدل‌های پیشنهاد شده توسط
۱/۳۰۵۵	۰/۴۳۸۲	ژو و لی (۲۰۰۱)
۰/۵۴۳۱	۰/۶۲۲۹	چن و دنگ (۲۰۰۸)
۱/۳۳۴۵	۰/۳۹۲۷	فرارو و همکاران (۲۰۱۰)
۰/۴۶۸۳	۰/۶۷۴۶	رگرسیون مارس فازی

برای پیش‌بینی مقدار فازی متغیر پاسخ بر اساس مدل رگرسیون مارس فازی به‌ازای $x = ۳/۴۰$ ابتدا مرکز متغیر پاسخ به‌صورت

$$\begin{aligned}
 Y &= ۲/۸۲ + ۲/۲۷ \max\{۰, ۳/۴۰ + ۰/۲۴\} \\
 &\quad - ۰/۷۶ \max\{۰, -۰/۲۴ - ۳/۴۰\} - ۱۱/۹۸ \max\{۰, ۳/۴۰ - ۲/۱۷\} \\
 &\quad + ۱۵/۹۸ \max\{۰, ۳/۴۰ - ۲/۳۲\} \\
 &= ۱۳/۶۰,
 \end{aligned}$$

برآورد می‌شود. سپس طبق روش معرفی شده در بخش ۴، به منظور فراهم آوردن یک سامانه استنتاج فازی، داریم

$$A = \{i \mid \tilde{y}_i(۱۳/۶۰) \geq ۰, i = ۱, \dots, ۱۸۴\} \quad \text{www.SID.ir}$$

$$= \{82, 97, 98, 141, 146, 184\}.$$

لذا خطاهای فازی

$$\hat{\varepsilon}_i = (0, 0/03 + \hat{d}_i, 0/03 + D_i - \hat{d}_i)_T, \quad i \in A,$$

در این سامانه فعال هستند. حال فرض کنید Y مرکز برآورد شده و $\hat{\varepsilon}$ خطای فازی برآورد شده متناظر با آن باشد. بنابراین در سامانه استنتاج فازی قوانین اگر-آنگاه زیر را داریم:

(۱) اگر Y به صورت $(8/47, 5/47)_T$ باشد، آنگاه $\hat{\varepsilon}$ به صورت $(0, 1/31, 9/63)_T$ است؛

(۲) اگر Y به صورت $(8/29, 5/63)_T$ باشد، آنگاه $\hat{\varepsilon}$ به صورت $(0, 4/99, 6/27)_T$ است؛

(۳) اگر Y به صورت $(8/96, 5/45)_T$ باشد، آنگاه $\hat{\varepsilon}$ به صورت $(0, 4/50, 6/40)_T$ است؛

(۴) اگر Y به صورت $(9/62, 5/02)_T$ باشد، آنگاه $\hat{\varepsilon}$ به صورت $(0, 2/74, 7/30)_T$ است؛

(۵) اگر Y به صورت $(11/25, 6/65)_T$ باشد، آنگاه $\hat{\varepsilon}$ به صورت $(0, 6/70, 6/60)_T$ است؛

(۶) اگر Y به صورت $(7/88, 5/76)_T$ باشد، آنگاه $\hat{\varepsilon}$ به صورت $(0, 4/56, 6/96)_T$ است.

چون $Y = 13/60$ ، با اجرای قوانین بالا در قسمت سامانه استنتاج فازی در نرم افزار MATLAB (مطلب، ۲۰۰۷)، خطای فازی $\hat{\varepsilon} = (0/29, 6/99, 9/34)_T$ به دست می آید. اکنون با جایگزین کردن خطای فازی به دست آمده در مدل رگرسیون مارس فازی داریم

$$\hat{Y} = 13/60 \oplus (0/29, 6/99, 9/34)_T = (13/89, 6/99, 9/34)_T.$$

بحث و نتیجه‌گیری

در این مقاله روشی دو مرحله‌ای برای معرفی مدل رگرسیون فازی با پهناهای متغیر بیان شد. در این رویکرد افزایش مقادیر متغیر تبیینی تأثیری بر افزایش یا کاهش پهناهای برآورد شده برای متغیر پاسخ ندارد. همچنین این روش می‌تواند مشاهدات فازی را که در آنها پهناهای متغیر پاسخ، روند نزولی، صعودی، ثابت یا متغیر دارند، به خوبی مدل‌بندی کند.

چون در روش پیشنهادی، مراکز متغیر پاسخ فازی بر اساس روش مارس که یک روش استوار نسبت به داده‌های پرت است، برآورد شد، اثرات منفی و نامناسب داده‌های پرت در برآورد مراکز متغیر پاسخ وارد نمی‌شوند. از طرف دیگر، در برآورد جملات خطای فازی از مراکز برآورد شده استفاده می‌شود، بنابراین این برآوردها نیز تحت تأثیر اثرات منفی داده‌(های) پرت واقع نخواهند شد.

اگرچه رویکرد پیشنهادی با الهام از رویکرد چن و دنگ (۲۰۰۸) تدوین شده است، اما این دو رویکرد در شیوه برآورد ضرایب رگرسیونی، قیود به‌کار رفته در مسأله بهینه‌سازی برای به‌دست آوردن جملات خطا و تعداد پارامترهای برآورد شده، متفاوت هستند. از طرفی می‌توان گفت که مدل پیشنهادی در مقایسه با روش چن و دنگ کمتر تحت تأثیر داده (یا داده‌های) پرت قرار می‌گیرد، که این موضوع دلیل برتری روش پیشنهاد شده در این مقاله بر روش آنها است. روش چن و دنگ یک روش رگرسیون کمترین توان‌های دوم فازی با پهناهای متغیر است که در برازش، برتری قابل توجهی بر بسیاری از روش‌های رگرسیون فازی دارد.

تقدیر و تشکر

نویسندگان مقاله از جناب آقای مهندس رضایی پژند کمال تشکر و قدردانی را بابت انتشار و در اختیار گذاشتن داده‌های واقعی به‌کار رفته در این مقاله را دارند. همچنین از داوران محترم که نظرات ارزشمند ایشان باعث بهبود مطالب ارائه شده در این مقاله گردید، کمال تشکر و قدردانی را داریم.

مراجع

- ارقامی، ن. ر. (۱۳۸۱)، مروری بر رگرسیون فازی، گزارش نخستین سمینار مجموعه‌های مشکک و کاربردهای آن، دانشگاه شهید باهنر کرمان، ۱-۱۸.
- چاچی، ج. (۱۳۹۱)، روش‌های آماری بر اساس اطلاعات نادقیق، رساله دکترای آمار، دانشگاه صنعتی اصفهان، دانشکده علوم ریاضی.
- میرزایی یگانه، ش.، ارقامی، ن. ر. (۱۳۸۶)، رگرسیون فازی: مروری بر چند رویکرد، اندیشه آماری، ۱۲، ۳۵-۴۷.
- De Andrés, J., Lorca, P., De Cos Juez, F. J. and Sánchez-Lasheras, F. (2011), Bankruptcy Forecasting: A Hybrid Approach Using Fuzzy C-Means Clustering and Multivariate Adaptive Regression Splines (Mars), *Expert Systems with Applications*, **38**, 1866-1875.
- Chen, S. P. and Dang, J. F. (2008), A Variable Spread Fuzzy Linear Regression Model with Higher Explanatory Power and Forecasting Accuracy, *Information Sciences*, **178**, 3973-3988.
- D'Urso, P., Massari, R. and Santoro, A. (2010), A Class of Fuzzy Cluster-wise Regression Models, *Information Sciences*, **180**, 4737-4762.
- D'Urso, P., Massari, R. and Santoro, A. (2011), Robust Fuzzy Regression Analysis, *Information Sciences*, **181**, 4154-4174.
- Ferraro, M. B., Coppi, R., González-Rodríguez, G. and Colubi, A. (2010), A Linear Regression Model for Imprecise Response, *International Journal of Approximate Reasoning*, **51**, 759-770.
- Fox, J. and Weisberg, S. (2011), *An R Companion to Applied Regression*, 2nd Edt., Sage Publications, Thousand Oaks, CA.

- Friedman, J. (1991), Multivariate Adaptive Regression Splines, *The Annals of Statistics*, **19**, 1-141.
- Hastie, T., Tibshirani, R. and Friedman, J. H. (2009), *The Elements of Statistical Learning*, 2nd Edt., Springer.
- Kelkinnama, M. and Taheri, S. M. (2012), Fuzzy Least-Absolutes Regression Using Shape Preserving Operations, *Information Sciences*, **214**, 105-120.
- Kriner, M. (2007), *Survival Analysis with Multivariate Adaptive Regression Splines*, Ph.D. Dissertation, Fakultät für Mathematik, Informatik und Statistik, Ludwig-Maximilians-Universität München.
- Lee, T. S. and Chen, I. F. (2005), A Two Stage Hybrid Credit Scoring Model Using Artificial Neural Networks and Multivariate Adaptive Regression Splines, *Expert Systems with Applications*, **28**, 743-752.
- Lee, T. S., Chiu, C. C., Chou, Y. C. and Lu, C. J. (2006), Mining The Customer Credit Using Classification and Regression Tree and Multivariate Adaptive Regression Splines, *Computational Statistics and Data Analysis*, **50**, 1113-1130.
- Lu, J. and Wang, R. (2009), An Enhanced Fuzzy Linear Regression Model with More Flexible Spreads, *Fuzzy Sets and Systems*, **160**, 2505-2523.
- Xu, R. and Li, C. (2001), Multidimensional Least-Squares Fitting with a Fuzzy Model, *Fuzzy Sets and Systems*, **119**, 215-223.
- Zimmermann, H. J. (2001), *Fuzzy Set Theory and Its Applications*, 4th Edt., Kluwer Nihoff, Boston.