

# روش اسپلاین‌های رگرسیونی تطبیقی چندگانه (MARS) با پاسخ‌های فازی و کاربرد آن در پزشکی اجتماعی<sup>۱</sup>

مینا رضایی، روشنگ علی محمدی، مهشید نامداری

دانشکده علوم ریاضی، دانشگاه الزهرا

دانشکده دندانپزشکی، دانشگاه علوم پزشکی شهید بهشتی

تاریخ پذیرش: ۱۴۰۰/۲/۶

تاریخ دریافت: ۱۳۹۹/۱۱/۱۹

نوع مقاله: علمی-پژوهشی

## چکیده

نخست مدل رگرسیون اسپلاین تطبیقی چندگانه (MARS) را در حالتی که مقادیر متغیرهای پاسخ، فازی هستند مرور می‌کنیم. سپس این مدل را بر داده‌های مربوط به یک بررسی در زمینه پزشکی اجتماعی به کار می‌گیریم. این بررسی، درباره ارتباط بین آگاهی افراد درباره بیماری‌های سرطان و سطح اجتماعی و اقتصادی افراد است. دانش افراد درباره سرطان به صورت مقادیر زبانی (/ کلامی) بیان شده است و در قالب مجموعه‌های فازی صورت بندی شده‌اند. نتایج مدلسازی را با نتایج به دست آمده با روش کمترین توانهای دوم فازی و کمترین مربعات خطای فازی مقایسه می‌کنیم که حاکی از برتری مدل MARS است. این مقایسه بر پایه دو ملاک رایج نیکویی برازش انجام گرفته است.

<sup>۱</sup> این مقاله نسخه اصلاح و تکمیل شده مقاله ای است که در هفتمین کنگره مشترک سیستم‌های فازی و هوشمند

ایران ارائه شده است.

Mathematics Subject Classification (2010): 99X99; 99X99, Email: Mi.Rezaei@student.alzahra.ac.ir.

عبارات و کلمات کلیدی: رگرسیون اسپلاین تطبیقی چندگانه (مارس)، متغیر زبانی، داده فازی، سطح اجتماعی و اقتصادی، رگرسیون کمترین توانهای دوم فازی

## ۱ سرآغاز

مدل های رگرسیونی الگوهایی ریاضی را با هدف بررسی ارتباط بین متغیر پاسخ ( / وابسته / خروجی) و متغیرهای تبیینی ( / مستقل / توضیحی / ورودی/ پیش بین) فراهم میکنند. رویکرد های متفاوتی در رگرسیون فازی وجود دارد. در این مقاله به معرفی روشی جدید می پردازیم و با مثال کاربردی این روش را توضیح می دهیم. *MARS* یک روش ناپارامتری است که هیچ فرضیه زیر بنایی درباره رابطه تابعی بین متغیرهای تبیینی و در نظر نمی گیرد. این روش، روشی کارآمد برای داده هایی است که دامنه وسیعی دارند و همچنین داده پرت نیز در این مدل موثر است. در بخش ۲ به مرور و بررسی کوتاه روش *MARS* می پردازیم. در بخش ۳ اهداف و انگیزه استفاده از روش *MARS* فازی را توضیح می دهیم. بخش ۴ روش *MARS* برای مدل سازی داده هایی با پاسخ فازی را شرح می دهیم و روش ساخت مدل را در این حالت تشریح می کنیم. در بخش ۵ دو معیار نیکویی برازش را معرفی می کنیم که در ادامه از آن استفاده می کنیم. بخش ۶ مثالی درباره سطح اجتماعی و اقتصادی افراد و آگاهی آنان درباره بیماری سرطان است که با استفاده از روش *MARS* ارتباط بین آن ها را مدل بندی کرده ایم و سپس این روش را با روش رگرسیون چن و دنگ و کمترین مربعات خطا مقایسه می کنیم و در نهایت در بخش ۸ به نتیجه گیری می پردازیم.

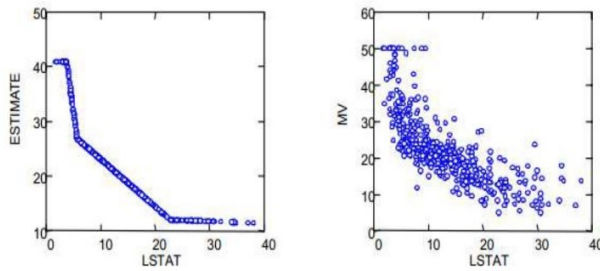
## ۲ رگرسیون تطبیقی چندگانه اسپلاین *MARS*

روش رگرسیون تطبیقی چندگانه اسپلاین (*MARS*) اولین بار توسط فریدمن [۷] معرفی شده و برای مسائلی با بعدهای بالا هنگامی که تعداد متغیر های پیش بین زیاد است یا حالتی که ارتباط بین متغیرها در بازه های مختلف از الگوهای مختلفی پیروی می کند، خوب عمل می کند. این روش می تواند به عنوان تعمیمی از رگرسیون خطی گام به گام یا اصلاحی از روش درخت رگرسیونی (*CART*) در نظر گرفته شود [۸]. به طوری که در روش *MARS* رویه پاسخ پیوسته است، اما در روش درخت رگرسیونی پیوسته نبوده و ناپیوستگی های آن در مرز نواحی افراز شده شکل می گیرد. در ادامه *MARS* را از دید اول یعنی تعمیمی از رگرسیون خطی گام به گام مرور می کنیم. مدل سازی در این روش بر اساس برازش رگرسیون های خطی قطعه ای که ساده ترین نوع اسپلاین ها

است، صورت می‌گیرد. در روش MARS، متغیرهای موثر و نقاط مرزی بازه‌ها (گره‌ها) برای هر متغیر از طریق روشی سریع تشخیص داده می‌شوند. علاوه بر جستجوی یک به یک متغیرها، MARS به جستجوی اثرات متقابل بین متغیرها تا هر مرتبه‌ای که مورد نظر باشد نیز می‌پردازد. مدل بهینه MARS در یک فرآیند دو مرحله‌ای انتخاب می‌شود. در مرحله اول، MARS یک مدل بیش از حد بزرگ را با روشی خاص می‌سازد، و در مرحله دوم توابع پایه‌ای که کمترین سهم را در مدل دارند، حذف می‌شوند [۴].

ابتدا با شکل زیر مفهوم گره‌ها را معرفی می‌کنیم، در شکل ۱ قاب سمت راست نشان دهنده مجموعه‌ای از داده‌ها و قاب سمت چپ یک برآورد MARS را با سه گره برای این داده‌ها نشان می‌دهد.

در MARS گره‌ها به وسیله یک روش جستجو تعیین می‌گردد. در این راستا، چنانچه معیار



شکل ۱: مجموعه‌ای از داده‌ها و یک برآورد MARS

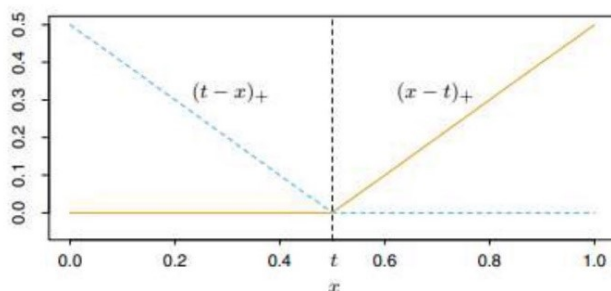
بهینگی برازش مدل رگرسیون، ضریب تعیین ( $R^2$ ) باشد در این روش MARS با بررسی تعداد زیادی از گره‌های بالقوه، گره‌ای که بیشترین مقدار ضریب تعیین را دارد، انتخاب می‌شود [۱۲]. گاه پیدا کردن بهترین جفت گره‌ها، نیاز به محاسبات زیادی دارد و همچنین یافتن بهترین مجموعه گره‌ها، وقتی که تعداد گره‌های مورد نیاز مجهول باشد، دشوار است. MARS مکان و تعداد مورد نیاز مورد نیاز گره‌ها را با یک روش پیش‌رو-پس‌رو می‌یابد. در مرحله پیش‌رو، مدلی که بیشترین برآورد دارد با گره‌های زیادی تولید می‌گردد. سپس گره‌هایی که کمترین سهم را در برازش کلی دارند، حذف می‌شوند. بنابر این انتخاب گره در مرحله پیش‌رو شامل مکان‌های نادرست زیادی برای گره‌ها خواهد بود، اما این گره‌های نادرست، احتمالاً در مرحله پس‌رو از مدل

حذف می‌شوند. از این رو می‌توان روش *MARS* را یک روش تطبیقی نامید [۶].  
 اساس روش *MARS* متکی بر توابعی قطعه ای موسوم به توابع اسپلاین است که به صورت زیر تعریف می‌شوند

$$h_1(x) = (t - x)_+ = \begin{cases} t - x & t > x \\ 0 & \text{در غیر این صورت} \end{cases}$$

$$h_2(x) = (x - t)_+ = \begin{cases} x - t & x > t \\ 0 & \text{در غیر این صورت} \end{cases}$$

که در آن  $t$  "گره" نامیده می‌شود و نماد "+"، نشان دهنده قسمت مثبت است. در این روش فرض بر آن است که گره  $t$  می‌تواند مقادیر مشاهدات متغیر  $x$  را اختیار کند، یعنی  $t \in \{x_1, x_2, \dots, x_n\}$ . در شکل ۲، نمودار توابع  $(x - 0.5)_+$  و  $(0.5 - x)_+$  به عنوان مثال نشان داده شده است. در این شکل نمودار توابع  $(x - t)_+$ ، خط توپر و  $(t - x)_+$ ، خط نقطه چین می‌باشد.



شکل ۲: نمودار توابع پایه

## ۱.۲ روش ساختن MARS

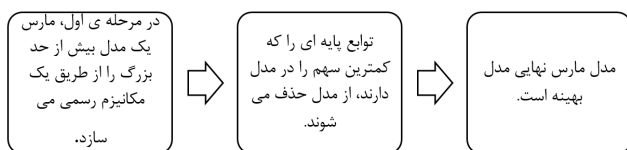
راهبرد ساختن مدل شبیه رگرسیون خطی گام به گام است. اما به جای استفاده از متغیرهای ورودی، از توابع در مجموعه  $C$  و حاصل ضرب هایشان (اثرهای متقابل) استفاده می‌شود.

$$C = \{(X_j - t)_+, (t - X_j)_+\}, \quad t \in \{x_{1j}, x_{2j}, \dots, x_{nj}\}, \quad j = 1, \dots, p.$$

اگر داده تکراری در مجموعه داده‌ها وجود نداشته باشد، مجموعه  $C$  شامل  $2np$  تابع پایه خواهد بود. مدل مارس مورد نظر دارای ساختار تابعی زیر است

$$f(X) = \beta_0 + \sum_{m=1}^M \beta_m h_m(X), \quad (1)$$

که در آن هر  $h_m(X)$  یک تابع در  $C$  یا حاصلضرب دو یا چند تا از این توابع است همچنین  $M$  تعداد توابع موجود در مدل است که پس از اعمال مرحلهٔ پس‌رو مشخص می‌گردد. با داشتن توابع  $h_m$ ، ضرایب  $\beta_m$  بوسیله مینیم کردن مجموع توان‌های دوم خطا، برآورد می‌شوند. در نرم افزار R، پکیج earth برای محاسبهٔ روش MARS تعریف شده است، به صورت کلی می‌توان گفت فرآیند کلی ساخت مدل مارس به شکل زیر است:



شکل ۳: فرآیند کلی ساخت مدل مارس

### ۳ اهداف و انگیزه استفاده از روش *MARS* فازی

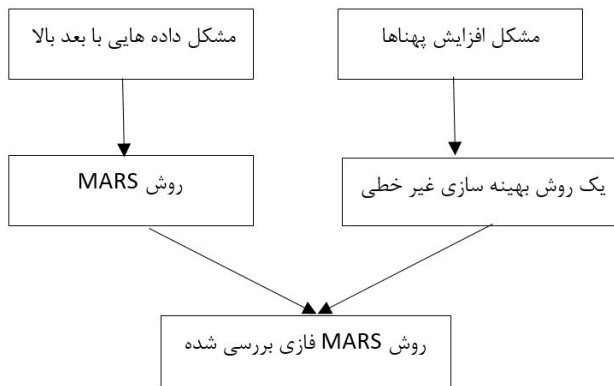
مدل های رگرسیونی (فازی یا آماری) برای بررسی ارتباط بین متغیر پاسخ و متغیرهای تبیینی الگویی را فراهم می‌کند. رگرسیون فازی را (برخلاف رگرسیون آماری) در موارد زیر می‌توان به کار برد

۱ داده ها نادقیق باشند.

۲ رابطه بین متغیر ها امکانی باشد.

۳ فرضیات زیربنایی در مدل های رگرسیون آماری مورد تردید باشد.

۴ حجم نمونه کم باشد.



شکل ۴: اهداف و انگیزه استفاده از روش *MARS* فازی

رگرسیون فازی را میتوان بسته به این که هر کدام از متغیر های تبیینی، پاسخ و ضرایب مدل دقیق یا فازی در نظر گرفته شود به انواع مختلفی تقسیم بندی کرد. علاوه بر تقسیم بندی ذکر شده رگرسیون فازی رویکرد های مختلفی وجود دارد که عمده ترین آن ها عبارتند از رویکرد کمترین مربعات فازی و رویکرد رگرسیون امکانی.

رویکردهای رگرسیون فازی بسیار متنوع است. برای مروری بر آن ها می‌توان به [۱۰] [۹] [۳]

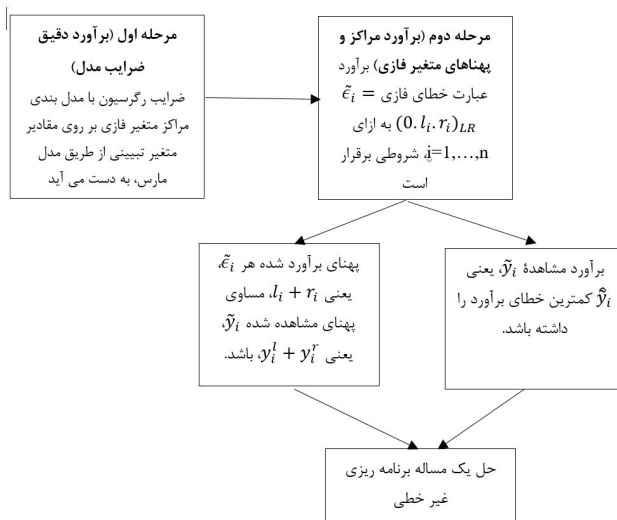
[۲] [۱] مراجعه کرد. هر یک از رویکردها برتری‌ها و ایراداتی دارند. حساسیت نسبت به داده‌های پرت و ابهام زیاد در مقادیر پیش‌بینی از جمله ایرادات برخی از روش‌های رگرسیون فازی است. بسیاری از این رویکردها از جهات مختلف مورد بررسی قرار گرفته‌اند، مشکل بسیاری از روش‌های موجود عبارتند از، بیشتر مدل‌های رگرسیون امکانی و مدل‌های مبتنی بر کمترین مربعات فازی نسبت به داده‌های (های) پرت حساس هستند؛ در برخی از روش‌ها با افزایش مقدار متغیر مستقل، پهنای برآورد شده‌ی متغیر پاسخ افزایش می‌یابد، که به آن مشکل افزایش پهنای گفته می‌شود؛ در برخی از روش‌ها که پهنای مشاهده شده متغیر پاسخ آن روند صعودی، نزولی و یا متغیر دارند، ممکن است برآورد‌هایی با خطای زیاد حاصل شود. اخیراً چاچی و همکاران [۵] یک مدل رگرسیون فازی ارائه داده‌اند که نسبت به داده‌های (های) پرت استوار است، و دیگر اینکه مساله افزایش پهنای را برحسب افزایش مقادیر متغیر تبیینی و پهنای متغیر پاسخ همزمان برطرف می‌کند. روش ارائه شده مبتنی بر اسپلاین‌های رگرسیونی است، این روش تکنیک Multivariate Adaptive Regression Spline (MARS) نام دارد، آن‌ها با استفاده از روش MARS داده‌هایی که متغیر پاسخ آن‌ها فازی است را تحلیل کردند. در شکل ۴ که با الهام از [۵] کشیده شده است، دلایل استفاده از روش مارس فازی شرح داده شده است. در ادامه ضمن تشریح و بررسی این مدل، کاربردی از آن را در مطالعات پزشکی معرفی و بررسی می‌کنیم.

## ۴ روش *MARS* با پاسخ‌های فازی

فرض کنید مجموعه‌ای از مشاهدات مربوط به یک یا چند متغیر پیش‌بین و یک متغیر پاسخ در اختیار داریم که مشاهدات متغیر (های) پیش‌بین دقیق هستند ولی مشاهدات متغیر پاسخ نادقیق (/ فازی) هستند. برای مدل بندی چنین داده‌هایی یک مدل رگرسیون فازی با پهنای متغیر ارائه شده است [۵]، مدل مربوطه به صورت زیر است

$$\tilde{y} = \beta_0 + \sum_{m=1}^M \beta_m B_m(x_i) \oplus \tilde{\epsilon}_i \quad i = 1, \dots, n \quad (2)$$

که در آن  $\tilde{y} = (y_i, y_i^l, y_i^r)$  ،  $i$  امین مشاهده فازی متغیر پاسخ،  
 $i$  امین مشاهده حقیقی مقدار متغیرهای تبیینی،  $x_i = [x_{0i}, x_{1i}, \dots, x_{ki}]$   
 $B_m(x_i)$  پایه<sup>۲</sup> ، ضرایب دقیق مدل ماریس برای توابع پایه  $m = 0, 1, \dots, M$  ،  $\beta_m$   
 و  $\tilde{\epsilon}_i = (0, l_i, r_i)$  خطای فازی مدل برای مشاهده  $i$  ام است.  
 برای برآورد پارامترها و جملات خطای فازی مدل یک روش دو مرحله ای معرفی



شکل ۵: فرایند برآورد یک مدل ماریس فازی

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 مساله بهینه سازی به دست می‌آیند.

مرحله ۱. (برآورد ضرایب دقیق مدل) ضرایب رگرسیون  $\beta = [\beta_0 + \beta_1 + \dots + \beta_M]^t$

<sup>2</sup>basic function



فازی بر روی مقادیر متغیر تبیینی، از طریق مدل مارس

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

به دست آورده می‌شود، ضرایب مدل در روش مارس با اجرای برنامه ی *earth* و یا *MARS* در نرم افزار *R* برآورد می‌شوند. سپس مراکز متغیر پاسخ فازی به صورت زیر به دست می‌آید.

مرحله ۲. (برآورد مراکز و پهناهای متغیر پاسخ فازی) عبارت خطای فازی

$\tilde{\epsilon}_i = (^\circ, l_i, r_i)_{LR}$  به ازای  $i = 1, \dots, n$  به گونه ای برآورد می‌شود که:

الف) پهناهای برآورد شده هر  $\epsilon_i$  یعنی  $l_i + r_i$ ، مساوی پهناهای 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 مساله 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 می‌شود که در آن تابع هدف مجموع اندازه های تشابه بین توابع عضویت مشاهده شده و برآورد شده متغیر پاسخ به صورت

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

است و قیود به گونه ای اختیار می‌شوند که  $i$  امین پهناهای مشاهده شده آن باشد. یعنی شرط  $l_i + r_i = y_i^l + y_i^r$ ،  $i = 1, \dots, n$ ، برقرار باشد. این قید به این دلیل در نظر گرفته می‌شود که پهناهای برآورد شده کنترل شود یا به عبارتی پهناهای خیلی کم یا خیلی زیاد برای متغیر پاسخ برآورد نشود. برای در نظر گرفتن این قیود فرض می‌شود  $y_m^l$  و  $y_m^r$  به ترتیب کمترین مقدار پهناهای چپ و راست مشاهدات متغیر وابسته باشند، یعنی

$$y_m^l = \min\{y_1^l, \dots, y_n^l\}, y_m^r = \min\{y_1^r, \dots, y_n^r\}$$

قسمتی از پهنای برآورد شده متغیر پاسخ مقدار ثابت  $y_m^l + y_m^r$  اختیار می‌شود. واضح 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  $0 \leq d_i \leq D_i, i = 1, \dots, n$  بنا بر این، قیود مساله بهینه سازی غیرخطی به صورت

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

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

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

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

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

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

به صورت کلی در شکل ۵ میتوان فرایند کلی برآورد مدل ماریس فازی را مشاهده کرد.

## ۵ نیکویی برازش

برای ارزیابی مدل های رگرسیون فازی چند ملاک توسط مولفان معرفی شده است. در این بررسی، از دو ملاک متداول برای نیکویی برازش رگرسیون فازی، به شرح زیر، استفاده می شود [۱۱] [۶].

$$MSM = \frac{1}{n} \sum_{i=1}^n \frac{\int \min\{\tilde{y}_i(x), \hat{y}_i(x)\} dx}{\int \max\{\tilde{y}_i(x), \hat{y}_i(x)\} dx}, \quad (3)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n \int |\tilde{y}_i - \hat{y}_i| dx \quad (4)$$

توجه کنید که MSM بین صفر و یک است ولی MAE می تواند هر مقدار بزرگتر از صفر باشد. ضمن اینکه، MSM به بیانی ملاک شباهت است و MAE ملاکی از فاصله بین دو عدد فازی است. در مقایسه بین دو مدل رگرسیون فازی هر مدلی که مقدار MAE کوچکتری (نزدیک به صفر) و یا مقدار MSM بزرگتری (نزدیک به یک داشته باشد)، برازش بهتری به داده ها دارد.

## ۶ کاربرد در پزشکی اجتماعی

رشته پزشکی اجتماعی یک شاخه تخصصی از پزشکی عمومی است که در آن به بررسی ارتباط بین سطح اجتماعی، اقتصادی، شرایط زندگی افراد و سلامت آن ها و سطح آگاهی درباره بیماری ها و عملکرد پزشکی آن ها می پردازد. در این علم به اثر نابرابری های اجتماعی بر روی بیماری هایی که فرد ممکن است درگیر آن شود می پردازد. در این مثال به بررسی آگاهی افراد از علائم و ریسک فاکتورهای سرطان پرداخته ایم و اثر متغیری به نام  $SES^3$  را بر سطح آگاهی مردم بررسی می کنیم.

<sup>3</sup> Social Economic Status

جدول ۱: نمونه ای از اطلاعات جمع آوری شده توسط پرسشنامه

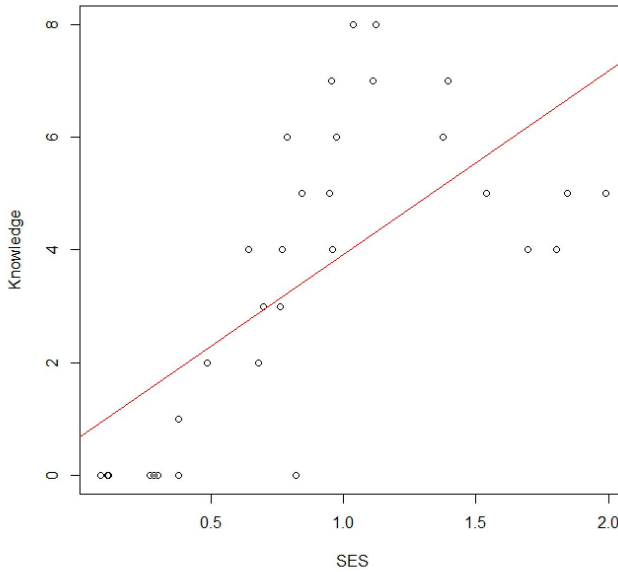
$i$	$x_i$	$\tilde{y}_i = (y_i, s_i)_T$	$d_i$
۱	۰/۰۸۱۹۱۳۴	$(۰, ۰/۸۵)_T$	۰/۳۲
۲	۰/۰۸۱۹۱۳۴	$(۰, ۰/۳۴)_T$	۰/۱۹
۳	۰/۱۰۷۴۴۶۴	$(۰, ۰/۰۶)_T$	۰/۰۱
⋮	⋮	⋮	⋮
۳۱	۱/۸۴۵۰۸۳	$(۵, ۰/۲۲)_T$	۰/۱۱
۳۲	۱/۹۸۹۳۵	$(۵, ۰/۰۷)_T$	۰/۰۲

*SES* یک معیار اقتصادی و اجتماعی است که تمام مقادیر اقتصادی و اجتماعی فرد و خانواده آن فرد را مورد اندازه‌گیری قرار می‌دهد. این معیار بر اساس درآمد، سطح تحصیلات و بهره‌مندی از مجموعه‌ای از امکانات رفاهی و شغل فرد می‌باشد. همچنین محل زندگی، داشتن انواع لوازم زندگی، مساحت ملک مسکونی، استیجاری یا صاحب ملک بودن و غیره نیز مورد بررسی قرار می‌گیرد. *SES* معمولاً برای بیان تفاوت سطح اقتصادی افراد مورد استفاده قرار می‌گیرد. حال به بررسی ارتباط بین *SES* افراد و سطح آگاهی آن‌ها می‌پردازیم.

یک پرسشنامه از ۳۲ نفر (والدین یک مهد کودک در تهران) پرسیده شده است، این پرسشنامه در سال ۱۳۹۵-۱۳۹۶ از افراد بالغ (۲۵-۸۰ سال) تکمیل شده است، در این پرسشنامه ۸ سوال درباره نشانه‌های سرطان دهان از افراد پرسیده شده است. اگر امتیاز فردی ۸ باشد، یعنی اطلاع آن شخص از نشانه‌های سرطان کامل است. برای اینکه اطلاعات به دست آمده دقیق باشد از افراد دو بار نظر سنجی صورت گرفته است که در صورت تفاوت در جواب‌ها این متغیر فازی در نظر گرفته شود و این تفاوت به عنوان پهنای عدد فازی فرض می‌شود. برای مثال برای مشاهده پنجم داریم:  $x = ۰/۱۱$  و  $y = (۰, ۰/۰۹)_T$  یعنی *SES* فرد مورد بررسی برابر ۰/۱۱ است و سطح دانش فرد حدود ۰ است. در اینجا میزان آگاهی افراد عدد فازی مثلثی در نظر گرفته می‌شود. مشاهدات در جدول ۱ آورده شده است. مقادیر  $\hat{d}_i$  محاسبه شده و مقادیر متغیرها را در ۱ مشاهده می‌کنید.

در شکل ۶ نمودار پراکنش *SES* افراد و میزان آگاهی آن‌ها درباره سرطان دهان مشاهده می‌کنیم. با توجه به این نمودار واضح است که یک خط راست دقت کافی برای مدل‌بندی به این داده‌ها را

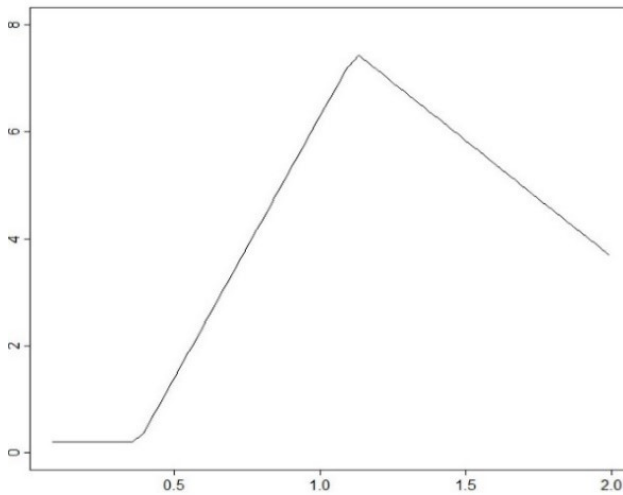
ندارد زیرا داده‌ها در بازه‌های مختلف رفتار متفاوتی از خود نشان می‌دهد. بنابراین روش مارس می‌تواند روش مناسبی برای برآورد داده‌ها باشد، می‌بینیم که مارس دارای دو گره است (در دو نقطه رفتار تابع عوض می‌شود). در ادامه با مارس فازی به برآورد متغیر پاسخ فازی  $\tilde{y} = (y, s)_T$  بر اساس متغیر تبیینی  $x$  می‌پردازیم، اما برخلاف شیوه‌های متداول مانند کمترین مربعات از خطوط شکسته ای به منظور برازش یک تابع پیوسته چند ضابطه ای بهینه به داده‌ها استفاده می‌شود.



شکل ۶: نمودار پراکنش داده‌ها و برازش رگرسیون کمترین مربعات خطا

$$\hat{y}_i^{MARS} = 0/19 + 9/78 \max\{0, x_i - 0/37\} - 14/14 \max\{0, x_i - 1/12\}$$

$$\oplus (0, 0/02 + \hat{d}_i, 0/02 + D_i - \hat{d}_i)_T \quad (5)$$



شکل ۷: مدل *MARS* برآورد شده به داده ها

## ۷ مقایسه با روش های دیگر

در بخش قبل برای داده های جمع آوری شده در یک مطالعه درباره آگاهی افراد درباره بیماری افراد و سطح اجتماعی آن ها با استفاده از روش ماریس مدلی را برآورد کردیم، ادامه روش ارائه شده با روش چن و دنگ [۶] و روش کمترین مربعات خطا [۱۳] مقایسه می شود. مدل برآورد شده با شیوه چن و دنگ به داده های پزشکی اجتماعی به صورت ۶ و با استفاده از روش کمترین قدر مطلق خطا به صورت ۷ برآورد شده است.

$$\hat{y}_i^{CD} = 0/68 + 3/31x_i \oplus \hat{\varepsilon}_i \quad (6)$$

$$\hat{y}_i^{LA} = (1/14, 0/41)_T \oplus (0/107, 0/00)_T x_i \quad (7)$$

در جدول ۲ مقادیر MAE و MSM مدل پیشنهادی، مدل چنگ و دنگ و روش کمترین مربعات خطا مقایسه شده است.

جدول ۲: مقایسه روش *MARS* و روش چن و دنگ برحسب دو معیار نیکویی برازش

<i>MAE</i>	<i>MSM</i>	
۰/۶۵	۰/۶۱	روش <i>MARS</i>
۱/۵۸۴	۰/۵۳	روش چن و دنگ
۰/۹۸	۰/۵۸	روش کمترین مربعات خطا

طبق جدول ۲ مشاهده می‌شود مدل *MARS*، مدل مناسب تری است زیرا مقدار *MSM* آن بیشتر و *MAE* آن کمتر است در نتیجه برای داده های تحت بررسی روش ماریس فازی روش کارآمد تری نسبت به روش های معمول دیگر است. در روش ماریس فازی مقدار *MAE* ۰/۶۵ است که با مقایسه با مقدار *MAE* به دست آمده در مدل های چنگ و دنگ (۱/۵۸۴) و روش کمترین مقدر مطلق خطا (۰/۹۸) مقدار کمتری است و طبق بخش ۵ مدل *MARS* برای داده های موجود عملکرد بهتری دارد.

مقدار *MSM* در روش های ماریس، چن و دنگ و کمترین قدر مطلق خطا به ترتیب ۰/۶۱، ۰/۵۳ و ۰/۵۸ به دست آمده که حاکی از برتری مدل ماریس فازی نسبت به مدل های دیگر است. با وجود اینکه روش کمترین مربعات خطا برای داده هایی که ارتباط خطی ندارند روشی کارا می باشد؟؟، اما در داده های موجود در پزشکی اجتماعی روش ماریس عملکرد بهتری داشته است.

## ۸ نتیجه گیری

در این مقاله روشی دو مرحله ای برای معرفی مدل رگرسیون فازی با پهنای متغیر توضیح داده شد. در روش *MARS* فازی، افزایش متغیر تبیینی تأثیری بر افزایش یا کاهش پهنای برآورد شده برای متغیر پاسخ ندارد. هم چنین روش *MARS* فازی می‌تواند مشاهدات فازی را که در آن ها پهنای متغیر پاسخ، روند صعودی، نزولی، ثابت یا متغیر دارند به خوبی برآورد کند. چون در این روش، مراکز متغیر پاسخ فازی بر اساس روشی که نسبت به داده های پرت استوار است، برآورد میشود، اثرات منفی و نامناسب داده های پرت در برآورد مراکز متغیر پاسخ وارد نمی‌شوند. از طرفی می‌توان گفت که مدل پیشنهادی در مقایسه با روش چن و دنگ کمتر تحت تأثیر داده (یا

داده های) پرت قرار میگیرد، که این موضوع دلیل برتری روش پیشنهاد شده در مقاله در برابر روش آن ها است.

کاربرد این روش بر پایه دادههایی واقعی مبتنی بر یک بررسی در حوزه پزشکی اجتماعی بررسی شد. نتایج مدل سازی حاکی از برتری مدل رگرسیون *MARS* فازی بود. گفتنی است که روش چن و دنگ یک روش رگرسیون کمترین توان دوم فازی با پهنای متغیر است که در برازش، برتری قابل توجهی نسبت به دیگر روش ها دارد. لذا مقایسه با روش چن و دنگ در واقع مقایسه با بسیاری از روش های رایج در رگرسیون فازی است.

## مراجع

- [۱] چاچی. ج. (۱۳۹۱) روش های آماری بر اساس اطلاعات نادقیق، رساله دکترای آمار، دانشگاه صنعتی اصفهان، دانشکده علوم ریاضی.
- [۲] طاهری، س.م. (۱۳۹۶) رویکردهای ابتکاری در رگرسیون فازی، اندیشه آماری، ۲۲ (شماره ۲)، ۴۳-۵۲.
- [۳] میرزایی یگانه، ش.، ارقامی، ن. ر. (۱۳۸۶) رگرسیون فازی: مروری بر چند رویکرد، اندیشه آماری، ۱۲، ۳۵-۴۷.
- [۴] ولی زاده، ت. (۱۳۹۱) اسپلاین های تطبیقی و غیر تطبیقی در مدل های رگرسیونی نیمه پارامتری، رساله کارشناسی ارشد آمار، دانشگاه صنعتی شاهرود، دانشکده ریاضی.
- [5] Chachi J., Taheri, S. M., Pazhand, H. R., (2016) Suspended Load Estimating Using L1-Fuzzy Regression, L2-Fuzzy Regression and MARS-Fuzzy Regression Model., Hydrological Science Journal, 6, 1489-1502.
- [6] Chen, S. P. and Dang, J. F. A (2008) Variable Spread Fuzzy Linear Regression Model With Higher Explanatory power and Forecasting Accuracy, Information Science, 38, 3973-3988.



- [7] Friedman, J. H. (1990) Multivariate adaptive regression splines. *Ann. Statist.*, 19, 11-41.
- [8] Hastie, T., Tibshirani, R. and Friedman, J. H. (2009) *The elements of statistical learning Data mining, Inference and prediction*, 2nd Ed., Springer, New York.
- [9] Khammar A., Arefi, M., Akbari, M. (2021) A general approach to fuzzy regression models based on different loss functions, *Soft Computing*, 25(2), 1-15.
- [10] Khammar A., Arefi M., Akbari M. (2020) Robust fuzzy varying coefficient regression model based on Huber loss function, 2020 8th Iranian Joint Congress on Fuzzy and intelligent Systems (CFIS), September 2-4, 2020, Mashhad, Iran.
- [11] Kelkinnama, M. and Taheri, S.M. (2012) Fuzzy Least-Absolute Regression Using Shape Preserving Operation, *Information Sciences*, 214, 105-120.
- [12] Salford System. (2001) *MARS user guide*.
- [13] Zeng W., Feng Q. and Li J. (2017) Fuzzy least absolute linear regression, *Applied Soft Computing*, 52, 1009-1019.
- [14] Jiang, L. and Liao, H. (2020) Mixed fuzzy least absolute regression analysis with quantitative and probabilistic linguistic information, *Fuzzy Sets and Systems*, 387, 35-48.