



رگرسیون ریج با متغیر پاسخ و پیشگوی فازی

محمدرضا ربیعی^۱، منصوره مقدس^۲

^۱ استادیار گروه آمار، دانشکده علوم ریاضی، دانشگاه صنعتی شاهرود

^۲ کارشناس ارشد آمار اقتصادی و اجتماعی، دانشگاه فردوسی مشهد

چکیده: پارامترهای مدل‌های رگرسیون فازی به طوری معمول به کمک مسائل بهینه‌سازی برآورد می‌گردند. در این مسائل عموماً فاصله بین پاسخ مشاهده‌شده و برآورد متناظر آن‌ها با شرایط خاصی مینیمم می‌گردد. در رگرسیون فازی نیز همانند رگرسیون معمولی وجود همخطی در بین متغیرهای پیشگو می‌تواند از کارایی مدل به دست آمده کم کند. از این رو لازم است به روش‌هایی مانند رگرسیون ریج روی آورد. در این مقاله یک مدل رگرسیون ریج فازی با متغیرهای پیشگو و پاسخ فازی و ضرایب دقیق پیشنهاد می‌کنیم. واژه‌های کلیدی: رگرسیون چندگانه فازی، رگرسیون ریج، همخطی، مسئله بهینه‌سازی.
کد موضوع بندی ریاضی (۲۰۱۰): 62J07, 62A86.

۱ مقدمه

بررسی رابطه بین یک متغیر وابسته و یک یا چند متغیر مستقل مورد توجه بسیاری از محققین می‌باشد تا بتوانند با استفاده از این رابطه به توصیف و پیش بینی متغیر وابسته بپردازند، از این رابطه در آمار با عنوان رگرسیون نام برده می‌شود. حال اگر متغیرها و یا رابطه بین آن‌ها اعداد مبهمی باشند محیطی که در آن به دنبال رابطه رگرسیونی بین آن‌ها هستیم یک محیط فازی بوده و در نتیجه یک تحلیل رگرسیون فازی خواهیم داشت، در حقیقت رگرسیون فازی یک روش قدرتمند برای تحلیل پدیده‌ها در محیط فازی است. رگرسیون فازی اولین بار توسط **تاناکا و همکاران (۱۹۸۲)** معرفی شد و پس از آن توسط بسیاری از محققان اصلاح و گسترش یافت.

مدل‌های رگرسیونی در حوزه کاربردی وسیعی مورد استفاده قرار می‌گیرد. یک مسئله جدی که می‌تواند استفاده از مدل رگرسیونی را با اشکال مواجه کند، همخطی یا وابستگی خطی نزدیک بین متغیرهای رگرسیونی است. وجود همخطی توانایی برآورد ضرایب مدل رگرسیون را با مشکل روبرو می‌کند. بنابراین لازم است روش‌هایی معرفی گردد که با وجود همخطی در بین متغیرهای پیشگو، بتواند همچنان مدل رگرسیون کارایی لازم را داشته باشد. برای این منظور، روش‌هایی مانند رگرسیون ریج معرفی شده است.

^۱محمدرضا ربیعی : Rabie1354@yahoo.com

وجود همخطی در رگرسیون فازی نیز می‌تواند در دسر ساز باشد از این رو محققان به معرفی روش‌های رگرسیون ریح فازی روی آوردند. **دونوسو و همکاران (۲۰۰۷)** یک روش رگرسیون فازی بر پایه رگرسیون ریح پیشنهاد کردند که به مقابله با مشکل همخطی در محیط فازی پرداخت. در مدل رگرسیون آن‌ها متغیرهای پیشگو (ورودی) دقیق و متغیرهای پاسخ (خروجی) و پارامترهای مدل فازی در نظر گرفته شده است، اما باید توجه داشت، گاهی متغیرهای پیشگو مبهم بوده و لازم است از مدل رگرسیون با متغیر پیشگوی فازی استفاده شود.

ما در این مقاله، یک روش رگرسیون ریح فازی برای حالت متغیرهای پیشگو و پاسخ فازی متقارن و پارامترهای دقیق ارائه می‌کنیم. این مقاله متشکل از بخش‌های زیر است: در بخش بعد به بیان یک سری تعاریف و قضایای به کار برده در مقاله می‌پردازیم. بخش ۳ مقدماتی درباره مدل رگرسیون ریح بیان می‌گردد. بخش ۴ مدل رگرسیون ریح فازی با متغیرهای پیشگو و پاسخ فازی متقارن و ضرایب دقیق معرفی شده و یک مثال عددی ارائه می‌شود. بخش آخر نیز به نتیجه‌گیری اختصاص دارد.

۲ تعاریف و قضایا

در این بخش به ارائه تعاریف و قضایای مورد نیاز می‌پردازیم که از مرجع **دوبوا و پراد (۱۹۸۰)** گرفته شده است.

تعریف ۱.۲. تابع عضویت عدد فازی مثلثی $\tilde{A} = (a, c_L, c_R)$ به صورت زیر است:

$$\mu_{\tilde{A}}(x) = \begin{cases} 1 - \left(\frac{a-x}{c_L}\right) & a - c_L \leq x \leq a \\ 1 - \left(\frac{x-a}{c_R}\right) & a < x \leq a + c_R \\ 0 & o.w \end{cases}$$

که در آن a مرکز عدد فازی و c_L و c_R به ترتیب پهنای چپ و راست عدد فازی است. مجموعه اعداد فازی مثلثی را با $T(R)$ نشان می‌دهیم. در صورتی که $c_L = c_R$ ، عدد فازی مثلثی متقارن نامیده می‌شود و با نماد $\tilde{A} = (a, c)_T$ نشان داده می‌شود.

قضیه ۲.۲. فرض کنید $\tilde{A}_1 = (a_1, c_{L1}, c_{R1})_T$ و $\tilde{A}_2 = (a_2, c_{L2}, c_{R2})_T$ اعداد فازی مثلثی و k یک عدد حقیقی باشد. در این صورت داریم:

$$\tilde{A}_1 + \tilde{A}_2 = (a_1 + a_2, c_{L1} + c_{L2}, c_{R1} + c_{R2})_T$$

$$k\tilde{A}_1 = \begin{cases} (ka_1, kc_{L1}, kc_{R1})_T & k \geq 0 \\ (ka_1, -kc_{R1}, -kc_{L1})_T & k < 0 \end{cases}$$

۳ رگرسیون ریح

در فضای رگرسیون احتمالی با حضور همخطی، درمینان ماتریس $X'X$ به سمت صفر میل می‌کند و این امر باعث نامطلوب شدن پارامترهای برآورد شده می‌گردد. وجود چندهمخطی آثار جدی متعددی بر برآوردهای حداقل مربعات ضرایب رگرسیونی دارد. به عنوان مثال فرض کنید فقط دو متغیر رگرسیونی X_1 و X_2 وجود داشته باشند. با این فرض که X_1 و X_2 و Y به طول واحد مقیاس-سازی

شده باشند، مدل را به صورت

$$Y = \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

در نظر می-گیریم و معادلات حداقل مربعات نرمال عبارتند از:

$$(X'X)\hat{\beta} = X'Y$$

یا

$$\begin{pmatrix} 1 & r_{12} \\ r_{12} & 1 \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} r_{1y} \\ r_{2y} \end{pmatrix}$$

که در آن r_{12} ضریب همبستگی پیرسن بین X_1 و X_2 و برای هر $j = 1, 2$ ، r_{jy} ضریب همبستگی پیرسن بین X_j و Y می باشد.

معکوس ماتریس $X'X$ عبارتست از:

$$C = (X'X)^{-1} = \begin{pmatrix} \frac{1}{1-r_{12}^2} & \frac{-r_{12}}{1-r_{12}^2} \\ \frac{-r_{12}}{1-r_{12}^2} & \frac{1}{1-r_{12}^2} r_{2y} \end{pmatrix} \quad (1.3)$$

و برآوردگرهای ضرایب رگرسیونی عبارتند از:

$$\hat{\beta}_1 = \frac{r_{1y} - r_{12}r_{2y}}{1 - r_{12}^2}, \quad \hat{\beta}_2 = \frac{r_{2y} - r_{12}r_{1y}}{1 - r_{12}^2}$$

اگر بین X_1 و X_2 همخطی شدید موجود باشد، در این صورت ضریب همبستگی بزرگ خواهد بود و با توجه به رابطه (۱.۳) اگر

$1 \rightarrow |r_{12}|$ آنگاه داریم:

$$V(\hat{\beta}_j) = \frac{1}{1 - r_{12}^2} \times \sigma^2 \rightarrow \infty$$

بسته به اینکه $1 \rightarrow r_{12}$ یا $-1 \rightarrow r_{12}$ خواهیم داشت:

$$Cov(\hat{\beta}_1, \hat{\beta}_2) = \frac{-r_{12}}{1 - r_{12}^2} \times \sigma^2 \rightarrow \pm\infty$$

بنابراین همخطی زیاد بین X_1 و X_2 منجر به واریانس‌ها و کواریانس‌های بزرگ برای برآوردگرهای حداقل مربعات ضرایب رگرسیونی خواهد شد. وقتی که بیش از دو متغیر رگرسیونی وجود داشته باشد، همخطی اثرات مشابهی ایجاد می-کند. همخطی همچنین به ایجاد برآوردهای حداقل مربعات $\hat{\beta}_j$ که از نظر قدرمطلق خیلی بزرگ می‌باشند، منجر خواهد شد. (رجوع شود به **دراپر و اسمیت (۱۹۸۱)**) یکی از روش‌های مقابله با مشکل همخطی، رگرسیون ریج است. در رگرسیون ریج به ماتریس $X'X$ یک پارامتر λ اضافه می‌شود، به عبارت دیگر ماتریس $(X'X - \lambda I)$ جانشین ماتریس $X'X$ می‌گردد، در نتیجه به منظور برآورد پارامترها رابطه زیر به دست می‌آید:

$$\hat{\beta} = (X'X - \lambda I)^{-1} X'Y$$

هر چه λ بیشتر باشد واریانس پارامترها کمتر اما اریبی آن‌ها بیشتر می‌گردد. هدف مقدار کوچک قابل قبول از λ است که در آن برآورد ریج پایدار شود. مقدار λ معمولاً توسط محقق تعیین می‌گردد و معمولاً عددی بین صفر و یک است. اگرچه در روش رگرسیون ریج، ضرایب برآوردشده نسبتاً اریب می‌باشند، ولی در مقایسه با برآوردهای کمترین مربعات، خطای کوچکتر و پایداری بیشتری دارند. در واقع برآورد ریج یک تبدیل خطی از برآورد حداقل مربعات است.

۴ مدل رگرسیون ریح فازی پیشنهادی

دونوسو و همکاران (۲۰۰۷) با تعمیم روش هاستی و همکاران (۲۰۰۵) در محیط فازی توانستند یک مسئله مینیمم‌سازی برای برآورد پارامترهای مدل رگرسیون ریح ارائه دهند. روش آن‌ها زمانی که متغیرهای پیشگو، دقیق و متغیر پاسخ، فازی فرض هستند قابل استفاده است. اما زمانی که متغیرهای پیشگو دارای ابهام باشند، روش آن‌ها کارایی خود را از دست می‌دهد. از این رو باید به مدل‌هایی روی آورد که علاوه بر فازی بودن متغیر پاسخ، متغیرهای پیشگو نیز فازی در نظر گرفته شوند. در ادامه سعی کردیم با ارائه روشی جدید مشکل این گونه مدل‌ها را برطرف کنیم.

مدل رگرسیون زیر را در نظر بگیرید:

$$\tilde{Y} = A_0 + A_1 \tilde{X}_1 + \dots + A_p \tilde{X}_p \quad (1.4)$$

به طوری که برای هر $i = 1, 2, \dots, n$ و $j = 1, 2, \dots, p$ مقادیر $\tilde{X}_{ij} = (x_{ij}, c_{x_{ij}})_T$ ، متغیرهای پیشگو، $\tilde{Y}_i = (y_i, c_{y_i})_T$ متغیر پاسخ فازی متقارن و A_0 و A_j پارامترهای مدل رگرسیونی هستند. ضرب بین پارامترهای دقیق مدل و متغیرهای پیشگوی فازی زمانی که اعداد مثلثی فازی باشند با استفاده از قضیه ۲.۲ قابل تعریف است. بنابراین مدل رگرسیون (۱.۴) به صورت زیر بازنویسی می‌شود:

$$\begin{aligned} \tilde{Y}_i &= A_0 + A_1 \tilde{X}_1 + \dots + A_p \tilde{X}_p \\ &= \left(A_0 + \sum_{j=1}^p A_j x_{ij}, \sum_{j=1}^p |A_j| c_{x_{ij}} \right) \end{aligned} \quad (2.4)$$

حال با توجه به مدل رگرسیون (۲.۴) و تعمیم روش دونوسو و همکاران (۲۰۰۷) برای حالت متغیرهای پیشگو و پاسخ فازی و ضرایب دقیق، تابع هدف به صورت زیر خواهد بود:

$$\begin{aligned} &k_1 \sum_{i=1}^n (y_i - (A_0 + \sum_{j=1}^p A_j x_{ij}))^2 + k_2 \left(\sum_{i=1}^n (y_i - c_{y_i} - (A_0 + \sum_{j=1}^p A_j x_{ij} - \sum_{j=1}^p |A_j| c_{x_{ij}}))^2 \right) \\ &+ \sum_{i=1}^n (y_i + c_{y_i} - (A_0 + \sum_{j=1}^p A_j x_{ij} + \sum_{j=1}^p |A_j| c_{x_{ij}}))^2 + \lambda (A_0^2 + \sum_{j=1}^p A_j^2) \end{aligned} \quad (3.4)$$

که در آن k_1 و k_2 اعداد حقیقی بین صفر و یک هستند که توسط محقق تعیین می‌شوند. در صورتی که $k_1 > k_2$ اهمیت مرکز بیشتر از پهنای پاسخ‌ها در نظر گرفته می‌شود و اگر $k_1 < k_2$ پهنای پاسخ‌ها از اهمیت بیشتری برخوردار می‌گردد. همچنین در صورتی که $k_1 = k_2$ ، اهمیت مرکز و پهنای یکسان فرض شده است. پارامتر λ نیز همان پارامتر ریح می‌باشد. با حل مسئله مینیمم‌سازی (۳.۴)، پارامترهای مدل رگرسیون (۱.۴) برآورد می‌شوند. در ادامه برای توضیح بیشتر، مدل رگرسیونی فوق را بر روی مشاهدات واقعی اجراء می‌کنیم.

مثال ۱.۴. در این مثال از داده‌های تاناکا و همکاران (۱۹۸۲) استفاده کردیم. این داده‌ها در جدول ۱ آمده است، که در آن متغیر پاسخ، قیمت فروش خانه (ven ۱۰۰۰)، و متغیرهای پیشگو به ترتیب، کیفیت مصالح ساختمانی (به ترتیب کیفیت پایین، متوسط، بالا)، مساحت طبقه اول (m^2)، مساحت طبقه دوم (m^2)، تعداد اتاق‌ها و تعداد اتاق‌های به سبک ژاپنی می‌باشد. در ابتدا به بررسی همخطی بین متغیرهای پیشگو می‌پردازیم. مقدار آماره VIF برای بررسی همخطی بین مشاهدات پیشگو در جدول ۲ گزارش شده است. با توجه به مقدار بیشتر از ۵ آماره VIF برای متغیرهای X_1 ، X_2 و X_5 مشکل همخطی وجود دارد. از این رو بهتر است از روش رگرسیون

جدول ۱: داده‌های مقاله تاناکا و همکاران

ردیف	\tilde{Y}_i	X_1	X_2	X_3	X_4	X_5
۱	$(6060, 550)_T$	۱	۳۸,۰۹	۳۶,۴۳	۵	۱
۲	$(7100, 50)_T$	۱	۶۲,۱۰	۲۶,۵۰	۶	۱
۳	$(8080, 400)_T$	۱	۶۳,۷۶	۴۴,۷۱	۷	۱
۴	$(8280, 150)_T$	۱	۷۴,۵۲	۳۸,۰۹	۸	۱
۵	$(8650, 750)_T$	۱	۷۵,۳۸	۴۱,۴۰	۷	۲
۶	$(8520, 450)_T$	۲	۵۲,۹۹	۲۶,۴۹	۴	۲
۷	$(9170, 700)_T$	۲	۶۲,۹۳	۲۶,۴۹	۵	۲
۸	$(10310, 200)_T$	۲	۷۲,۰۴	۳۳,۱۲	۶	۳
۹	$(10920, 600)_T$	۲	۷۶,۱۲	۴۳,۰۶	۷	۲
۱۰	$(12030, 100)_T$	۲	۹۰,۲۶	۴۲,۶۴	۷	۲
۱۱	$(13940, 350)_T$	۳	۸۵,۷۰	۳۱,۳۳	۶	۳
۱۲	$(14200, 250)_T$	۳	۹۵,۲۷	۲۷,۶۴	۶	۳
۱۳	$(16010, 300)_T$	۳	۱۰۵,۹۸	۲۷,۶۴	۶	۳
۱۴	$(16320, 500)_T$	۳	۷۹,۲۵	۶۶,۸۱	۶	۳
۱۵	$(16990, 650)_T$	۳	۱۲۰,۵۰	۳۲,۲۵	۶	۳

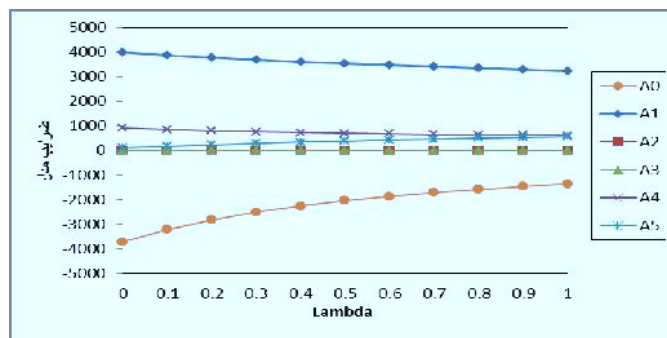
ریج برای برآورد پارامترهای مدل رگرسیون فازی استفاده شود. با توجه به این که متغیرهای پیشگو، دقیق هستند، برای دو متغیر X_2 و X_3 (این دو متغیر پیوسته هستند، در نتیجه امکان خطا در اندازه گیری وجود دارد. بنابراین آن‌ها را به صورت فازی در نظر می‌گیریم) پهنایی متناسب با مقدار مشاهده شده در نظر گرفته می‌شود. برای این منظور، با مشورت متخصص، به ازای هر $i = 1, 2, \dots, n$ $S_{ij} X_{ij}$ ، $j = 2, 3$ برای پهنای مشاهده مورد نظر انتخاب شدند. به طوری که $0 \leq S_{ij} \leq 1$ و مقادیر تصادفی باشند (تمامی مشاهدات اعداد فازی مثلثی متقارن هستند). مسئله بهینه‌سازی برای $k_1 = k_2 = 1$ و λ های متفاوت از صفر تا یک با گام 0.1 انجام می‌گیرد، پارامترهای برآورد شده در شکل ۱ نشان داده شده است. با توجه به شکل ۱، برآوردها به طور تقریبی از $\lambda = 0.8$ هموار شده‌اند. بنابراین مدل رگرسیون نهایی به صورت زیر است:

$$\tilde{Y} = -1580.86 + 3343.37\tilde{X}_1 + 8.85\tilde{X}_2 + 10.80\tilde{X}_3 + 633.51\tilde{X}_4 + 500.97\tilde{X}_5$$

برای پیش‌بینی از مدل فوق، فرض کنید بردار $\{(1, 0)_T, (80/21, 0/50)_T, (54/84, 0/25)_T, (3, 0)_T, (1, 0)_T\}$ بردار مشاهدات متغیرهای پیشگو باشد، با جایگذاری در مدل، مقدار $\hat{\tilde{Y}} = (5466/14, 7/12)_T$ برای قیمت خانه پیش‌بینی می‌شود.

جدول ۲: داده‌های مقاله تاناکا و همکاران

نام متغیر	VIF
X_1	۹,۵۰
X_2	۶,۰۵
X_3	۱,۶۷
X_4	۳,۸۷
X_5	۶,۱۳

شکل ۱: ضرایب برای λ های متفاوت

بحث و نتیجه‌گیری

در این مقاله تلاش کردیم برای مشاهدات فازی که با مشکل همخطی روبرو هستند، مدل رگرسیون ریج فازی را برازش دهیم. برای به انجام رساندن این هدف، با تعمیم یک مدل رگرسیون ریج فازی، یک مسئله برنامه‌ریزی درجه دوم پیشنهاد گردید، به طوری که علاوه بر غلبه بر مشکل همخطی، برای متغیرهای پیشگوی فازی نیز قابل استفاده باشد.

مراجع

- Asai H., Tanaka S. and Uegima K. (1982), Linear Regression Analysis With Fuzzy Model, *IEEE Trans. Systems Man Cybern* 12, 6, 903-907.
- Donoso, S., Nicolás M., and Vila M. A., (2007), Fuzzy ridge regression with non symmetric membership functions and quadratic models. *Intelligent Data Engineering and Automated Learning-IDEAL 2007*. Springer Berlin Heidelberg, 135-144.
- Draper N. R. and Smith H. (1981), *Applied Regression Analysis*. John Wiley & Sons, Inc, New York.

Dubois, D., and Prade, H. (1980), Fuzzy Sets and Systems: Theory and Applications. *Mathematics in Science and Engineering*. Vol. 144.

Hastie T., Tibshirani, R. Friedman, J. and Franklin, J. (2005), The elements of statistical learning: data mining, inference and prediction, *The Mathematical Intelligencer*, 27, 2, 83-85.