

رگرسیون خطی فازی و کاربردهای آن در پژوهش‌های علوم اجتماعی

ابراهیم خدایی *

(تاریخ دریافت: ۱۳۸۸/۶/۳۰ ، تاریخ تصویب: ۱۳۸۸/۱۲/۱۸)

چکیده

در این مقاله، تئوری رگرسیون کلاسیک و رگرسیون فازی مقایسه شده‌اند و مثالی کاربردی از کاربرد رگرسیون فازی در علوم اجتماعی عرضه می‌شود. رگرسیون کلاسیک برپایه فرض دقیق بودن متغیرهای مورد مطالعه و مشاهدات مرتبط با آنها استوار بوده و در نهایت، روابط بین متغیرها را نیز دقیقاً مشخص می‌کند. در مدل‌سازی مسائل اجتماعی، عموماً با مشاهدات نادقیق یا روابط مبهم روبرویم. بنابراین استفاده از روش‌های برازش توابع که به تبیین ساختار مبهم داده‌ها و روابط بین آنها می‌پردازند، ضروری است. در بررسی حاضر، رگرسیون خطی فازی معدل کتبی دیپلم روی نمرات ریاضی، فیزیک و جنسیت داوطلبان آزمون سراسری سال ۱۳۸۷ در قالب مثالی کاربردی با فرض فازی بودن ضرایب و غیر فازی بودن داده‌ها بررسی می‌شود. برای آزمون رگرسیون فازی از شاخص اطمینان و مقدار ابهام مدل استفاده خواهد شد.

واژگان کلیدی: رگرسیون خطی کلاسیک، مجموعه فازی، رگرسیون فازی خطی، آزمون سراسری.

* استادیار آمار سازمان سنجش آموزش کشور، khodaie@sanjesh.org

مقدمه

پیش از گسترش علوم، خصوصاً علوم انسانی، بررسی و تحلیل جوامع انسانی عموماً در فلسفه صورت می‌گرفت. در این دوره، مسائل اجتماعات انسانی با دیدگاه‌هایی فلسفی که اساساً ذهنی و انتزاعی بودند تفسیر می‌شدند و اندیشمندان به دنبال بررسی تجربی روابط و پدیده‌های انسانی و اجتماعی نبودند. روایی این‌گونه حکم‌ها و نظرات نیز با معیارهای ذهنی و عقلی ارزشیابی می‌شد، نه با ملاک‌های عینی و تجربی. توسعه دیدگاه‌های علمی در اوایل قرن بیستم و ظهور اثبات‌گرایی دانشمندان حوزه علوم انسانی را به پژوهش‌های عینی و تجربی رهنمون ساخت و صاحب‌نظران را به عرضه معرفتی جدید از واقعیت‌های اجتماعی و تحقیق در خصوص مسائل جامعه با روش عینی هدایت کرد. در این میان، نظریه‌پردازی و روش‌شناسی در علوم اجتماعی به استفاده از علم آمار به مثابه ابزاری دقیق برای محک تجربی پدیده‌ها گرایش پیدا کرد. با نگاهی جزئی‌تر می‌توان جامعه‌شناسی را به دو حوزه نظری و تجربی تقسیم کرد. در بحث‌های نظری، با اندیشه‌ها و نظریه‌ها و مفاهیم انتزاعی و روابط بین آنها سرو کار داریم، اما در حوزه تجربی با جهان تجربی و داده‌ها روبرویم و در هر دو حوزه، برای نتیجه‌گیری و تولید اندیشه و حل مسائل به آمار و ریاضیات به مثابه یک ابزاری مناسب برای تلخیص و استنباط داده‌ها نیاز است. یکی از ابزارهای بسیار مهم در تحلیل روابط بین متغیرها رگرسیون است. کاربردهای روش‌های رگرسیونی به تفصیل در حوزه علوم اجتماعی به بحث و بررسی گذاشته شده است و در خصوص لزوم این ابزار برای تحقیق‌های علوم اجتماعی اجماع نظر وجود دارد. با توجه به این‌که به‌کارگیری شیوه‌ها و ابزارهای نظریه مجموعه‌های فازی در گسترش و تعمیق روش‌های آماری چندی است مورد توجه محققان حوزه‌های علوم، به خصوص علوم اجتماعی و جامعه‌شناسی، قرار گرفته است، هدف این مقاله نگاهی اجمالی به روش رگرسیون فازی از نظر آموزش علاقه‌مندان و محققان حوزه علوم اجتماعی است. با این دیدگاه، در این مقاله کلیات روش‌ها، اهداف و فرض‌های لازم در رگرسیون کلاسیک و نظریه مجموعه‌ها، آمار، تئوری و آزمون‌های رگرسیون فازی به بحث و بررسی گذاشته می‌شود. سپس، رابطه بین نمرات، معدل و جنسیت داوطلبان در آزمون سراسری سال ۱۳۸۷ از منظر رگرسیون کلاسیک و فازی تجزیه و تحلیل و نتایج به‌دست‌آمده با همدیگر مقایسه می‌شوند.

رگرسیون خطی

نظریه رگرسیون را اولین بار فرانسیس گالتون در ۱۸۸۵ مطرح کرد. معنی کلمه رگرسیون بازگشت است و گالتون نشان داد که میانگین قد فرزندان افراد بلندقد و کوتاه‌قد به میانگین جامعه اصلی بازگشت پیدا می‌کند. بعد از عرضه تئوری رگرسیون، مثل تئوری احتمال که با پیش‌بینی برد و باخت قماربازان شکل گرفت، محققان حوزه‌های گوناگون علوم دریافتند که کاربردهای رگرسیون بسیار فراتر از بحث اولیه آن بوده است، به طوری که در حال حاضر اگر بخواهیم تعریفی برای این تئوری به دست دهیم، می‌توان به شکل خلاصه گفت که تحلیل رگرسیون روشی آماری است که تغییرات متغیر وابسته از طریق متغیر یا متغیرهای مستقل تبیین و پیش‌بینی می‌شود. روش‌های رگرسیونی از نظر تئوری به سه نوع رگرسیون خطی، غیر خطی و ناپارامتری تقسیم می‌شوند. در این مقاله، فقط رگرسیون خطی مورد بحث و بررسی قرار خواهد گرفت. رگرسیون خطی را برحسب تعداد متغیرهای وابسته و مستقل می‌توان در قالب سه نوع تقسیم‌بندی زیر لحاظ کرد:

(الف) رگرسیون ساده خطی: در این حالت، یک متغیر وابسته و یک متغیر مستقل وجود دارد.

(ب) رگرسیون چندگانه: در این حالت، یک متغیر وابسته و بیش از یک متغیر مستقل فرض می‌شود.

(ج) رگرسیون چندمتغیره: در این حالت، چند متغیر وابسته و چند متغیر مستقل برای تجزیه و تحلیل فرض می‌شوند.

در رگرسیون خطی ساده، اگر فرض کنیم دو متغیر (y_i, x_i) ، $i = 1, 2, \dots, n$ وجود داشته باشد، رابطه به صورت زیر است:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad \varepsilon_i = 1, 2, \dots, n \quad (1)$$

در رگرسیون چندگانه، بیش از یک متغیر مستقل وجود دارد، به طوری که اگر فرض کنیم دو متغیر مستقل Z و X و متغیر وابسته Y موجود باشد، رابطه مورد نظر به صورت زیر خواهد بود:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 z + \varepsilon_i$$

در رابطه (۱) β_0 و β_1 را پارامترهای مدل نامیده و هدف برآورد β_0 و β_1 است. برای برآورد پارامترهای فوق روش‌های متفاوتی وجود دارد که می‌توان به روش‌های حداقل مربعات و حداکثر درست‌نمایی اشاره کرد. در روش حداقل مربعات، β_0 و β_1 طوری برآورد می‌شوند که مقدار باقیمانده یا خطا در رابطه (۱)، یعنی $\varepsilon_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$ در رابطه $\sum_{i=1}^n \varepsilon_i^2$ حداقل مقدار خود را داشته باشد.

با استفاده از مشتق‌گیری و سایر روش‌های ریاضی، برآورد پارامترها به صورت زیر خواهد بود:

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad (2)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (3)$$

که در رابطه (۲) و (۳)، \bar{x} و \bar{y} میانگین‌های نمونه X و Y اند.

در ساختن مدل رگرسیونی (۱) فرض می‌شود که ε_i ها مستقل و دارای توزیع نرمال $N(0, \sigma^2)$ اند (σ^2 نامعلوم است). در رگرسیون کلاسیک، پس از پیدا کردن معادله رگرسیونی از آزمون‌های کلاسیک آماری از قبیل آزمون F ، آزمون t به منظور اعتباربخشی به مدل استفاده می‌شود و علاوه بر این، همان طوری که در مقدمه ذکر شد، هدف تئوری رگرسیون پیش‌بینی و تبیین واریانس متغیر وابسته به وسیله متغیرهای مستقل است و آماره‌ای که برای این امر استفاده می‌شود عبارت است از ضریب تبیین:

$$r^2 = \frac{\text{واریانس معلوم}}{\text{کل واریانس}} = \frac{\sum (\hat{y} - \bar{y})^2}{\sum (y - \bar{y})^2} = \frac{n-2}{n-2} \quad (4)$$

Γ^2 یا توان دوم ضرب همبستگی نشان دهنده واریانس مشترک X و Y است یا به عبارت دیگر، مقدار Γ^2 نشان دهنده میزان تغییرات متغیر وابسته است که به وسیله مدل رگرسیونی تبیین می‌شود. همچنین، Γ^2 را گاهی معیار کاهش نسبی خطا نیز می‌نامند.

مجموعه‌ها و آمار فازی

ریاضیات یک فرامجموعه از منطق بولی است که بر مفهوم درستی نسبی دلالت می‌کند. منطق کلاسیک هر چیزی را بر اساس یک سیستم دوتایی نشان می‌دهد (درست یا غلط، \circ یا 1 ، سفید یا سیاه)، ولی منطق فازی درستی هر چیزی را با یک عدد که مقدار آن بین صفر و یک است نشان می‌دهد. مثلاً اگر رنگ سیاه را با عدد صفر و رنگ سفید را با عدد 1 نشان دهیم، آنگاه رنگ خاکستری عددی نزدیک به صفر خواهد بود. در ۱۹۶۵، دکتر لطفی زاده نظریه سیستم‌های فازی را معرفی کرد. منطق فازی، معتقد است که ابهام در ماهیت علم است. لطفی زاده معتقد است که باید به دنبال ساختن مدل‌هایی بود که ابهام را به منزله بخشی از سیستم برخلاف دیگران که معتقدند باید تقریب‌ها را دقیق‌تر کرد تا بهره‌وری افزایش یابد مدل کند. در منطق ارسطویی، یک دسته‌بندی درست و نادرست وجود دارد. همه گزاره‌ها درست یا نادرست‌اند. اما در منطق فازی، گزاره‌ها مقداری درست و مقداری نادرست‌اند. مثلاً گزاره "هوا سرد است" یک گزاره منطقی فازی است که درستی آن بسته به شرایط جغرافیایی و فرد قضاوت‌کننده گاهی کم و گاهی زیاد است.

قبل از بحث درباره مجموعه‌های فازی بهتر است نگاهی اجمالی به نظریه مجموعه‌های کلاسیک بیان‌داریم که با تعریف مجموعه شروع می‌شود. گردایه معینی از اشیاء را مجموعه می‌نامند. اشیاء این گردایه اعضا یا عناصر مجموعه نامیده می‌شوند. مثلاً اعداد طبیعی کمتر یا مساوی ۸ یک مجموعه را تشکیل می‌دهند. اگر این مجموعه را با A نشان دهیم، خواهیم داشت:

$$A: \{1, 2, \dots, 8\}$$

مجموعه A را می‌توان به صورت‌های $A = \{x \in \mathbb{N} \mid x \leq 8\}$ یا $A = \{x \in X \mid P(x)\}$ نیز نشان داد که P ویژگی x را نشان می‌دهد. برای متعلق بودن عناصر یک مجموعه می‌توان از توابع نشانگر استفاده کرد. مثلاً برای مجموعه فوق‌الذکر تابع نشانگر به صورت زیر تعریف می‌شود:

$$I_A(x) = \begin{cases} 1 & x \in A \\ \circ & x \notin A \end{cases}$$

چنان‌که ملاحظه می‌شود، دامنه تابع نشانگر مجموعه مرجع و برد آن مجموعه دو عضوی $\{\circ, 1\}$ است. یعنی $I_A(x): X \rightarrow \{\circ, 1\}$. واضح است که هر مجموعه یک تابع نشانگر دارد و بالعکس.

در مثال بالا، ویژگی مجموعه A دقیقاً تعریف شده است و عضویت یک عنصر در مجموعه کاملاً مشخص است، اما در زندگی روزمره با مفاهیم و استدلال‌هایی سروکار داریم که نادقیق‌اند، از قبیل قد بلند، سردی هوا، مسافت طولانی، فشار خون بالا. نظریه مجموعه‌های فازی نظریه‌ای برای مدل‌بندی و تجزیه و تحلیل مفاهیم نادقیق فوق‌الذکر است. در نظریه مجموعه‌های کلاسیک، برد تابع نشانگر مجموعه $\{\circ, 1\}$ است، اما در منطق فازی برد مجموعه به جای مجموعه $\{\circ, 1\}$ بازه $[\circ, 1]$ در نظر گرفته می‌شود. به عبارت دیگر، در این حالت به جای تعلق قطعی، یا عدم تعلق قطعی میزان عضویت در بازه $[\circ, 1]$ تعریف می‌شود و هر قدر میزان عضویت به عدد یک نزدیک‌تر باشد، یعنی تعلق بیشتر عضو به مجموعه، و هر قدر به صفر نزدیک‌تر باشد، یعنی عدم تعلق به مجموعه.

در این تئوری، به جای مجموعه نشانگر از تابع عضویت $[0,1] : X \rightarrow \mu_A(x)$ استفاده می‌شود، به عبارت دیگر، $\mu_A(x)$ میزان عضویت x در مجموعه فازی A را نشان می‌دهد. با توجه به این که در تئوری مجموعه‌ها خصوصیات ریاضی گوناگونی از اجتماع، اشتراک متمم، مکمل و افراز وجود دارد، طبیعتاً این خصوصیات در تئوری مجموعه‌های فازی نیز قابل تعریف‌اند.^۲

با توجه به گسترش کاربردی منطق فازی در علوم گوناگون از قبیل فنی‌مهندسی، علوم انسانی و علوم پزشکی مجال بحث و بررسی درخصوص این منطق در این مقاله فراهم نیست و فقط نکته مهمی که باید در اینجا بدان اشاره شود این است که آیا با توجه به حوزه مورد بحث، یعنی احتمالات، اشتراکی بین تئوری منطق فازی و علم آمار وجود دارد یا خیر؟ یا آیا بین منطق فازی و علم آمار مرزبندی دقیقی وجود دارد یا خیر؟ در پاسخ می‌توان گفت نظریه آمار و نظریه مجموعه‌های فازی هر دو برای مطالعه الگوها و سیستم‌هایی شامل عدم قطعیت و پویا وضع شده‌اند، اما با این تفاوت که نظریه آمار برای مطالعه الگوهای مبتنی بر عدم قطعیت آماری (منسوب به پیشامدهای تصادفی) و نظریه مجموعه‌های فازی برای مطالعه الگوهای مبتنی بر عدم قطعیت امکانی (ناشی از ابهام و نادقیق بودن) گسترش یافته‌اند.

این دو نظریه متناقض یکدیگر نبوده و شامل همدیگر هم نیستند، اما می‌توان برای حل یک مسئله از هر دو روش‌های کلاسیک آماری و روش‌های فازی استفاده کرد. در این زمینه مثلاً می‌توان به آزمون فرض‌های فازی رگرسیون فازی اشاره داشت.

رگرسیون خطی فازی

چنان که در بخش قبلی ملاحظه شد، رگرسیون کلاسیک فرضیاتی را در زمینه توزیع احتمالی خطاها در نظر می‌گیرد. اگرچه مدل رگرسیون خطی کلاسیک کاربردهای بسیار دارد، اما در بعضی مواقع ساختن مدل با مشکلاتی مواجه است که عبارت‌ند از تعداد کم یا نامناسب بودن مشاهدات، مشکلات تعریف تابع توزیع مناسب، ابهام در رابطه بین متغیرهای وابسته و مستقل، ابهام در وقوع یا درجه وقوع رویدادها، بی‌دقتی و خطا. مثلاً به کاربردن تحلیل رگرسیون آماری ممکن است باعث نتیجه‌گیری اشتباه شود. برای حل این مسئله می‌توان از روش‌های دیگر از قبیل رگرسیون استوار و رگرسیون خطی فازی استفاده کرد که در اینجا رگرسیون خطی فازی مورد بحث قرار می‌گیرد. رگرسیون فازی در حالت کلی به سه نوع تقسیم می‌شود:

(الف) رگرسیون فازی در حالتی که روابط بین متغیرها (ضرایب مدل رگرسیونی) فازی فرض شوند.

(ب) رگرسیون فازی در حالتی که مشاهدات در متغیر وابسته و متغیرهای مستقل نادقیق و فازی باشند.

(ج) رگرسیون فازی در حالتی که هم روابط بین متغیرها و هم مشاهدات فازی در نظر گرفته شوند.

رگرسیون خطی با ضرایب فازی را اولین بار تاناکا و همکارانش در ۱۹۸۲ معرفی کردند. بعد از اولین مقاله ایشان، مقالات متعددی در خصوص تئوری فوق‌الذکر به دست ایشان و افراد دیگری با تکیه بر گسترش تئوری و مثال‌های کاربردی منتشر شد که حشمتی و کاندل (۱۹۸۵) برای پیش‌بینی هوا از آن جمله است.

در رگرسیون خطی با ضرایب فازی، فرض می‌شود که مشاهدات و متغیرها دقیق و ابهام در مدل و ضرایب رگرسیون است. فرض کنیم Y متغیر وابسته و X_1, X_2, \dots, X_p متغیرهای مستقل و تعداد مشاهدات n باشد، صورت کلی مدل رگرسیون خطی فازی به شکل زیر خواهد بود:

$$\tilde{Y} = f(x, A) = \tilde{A}_0 + \tilde{A}_1 x_1 + \tilde{A}_2 x_2 + \dots + \tilde{A}_p x_p \quad (5)$$

هدف برآورد پارامترهای مدل (۵) یعنی $\tilde{A}_0, \tilde{A}_1, \dots, \tilde{A}_p$ است، به قسمتی که مدل (۵) بهترین برازش را برای داده‌ها به دست دهد. برای پیدا کردن پارامترهای فوق از تابع عضویت مثلثی زیر استفاده خواهیم کرد. توجه شود که می‌توان از توابع عضویت دیگر، از قبیل نرمال، استفاده کرد، اما در این مقاله فقط تابع عضویت مثلثی مورد بحث و بررسی قرار می‌گیرد. تابع عضویت مثلثی به صورت زیر تعریف می‌شود:

^۲. نک منابع گوناگون فارسی، به ویژه عادل آذر (۱۳۸۶)

$$\tilde{A}_{(x)} = \begin{cases} 1 - \frac{a-x}{s^L} & a - s^L \leq x \leq a \\ 1 - \frac{x-a}{s^R} & a < x \leq a + s^R \end{cases} \quad (6)$$

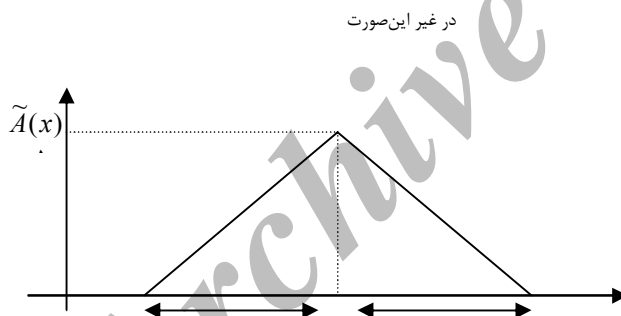
که a مقدار میانه، s^L کران پایین و s^R کران بالای \tilde{A} را نشان می‌دهند، به قسمیت که $\tilde{A} = (a, s^L, s^R)_T$. اگر $s^L = s^R = s$ باشد، در این صورت \tilde{A} عدد مثلثی متقارن و در غیر این صورت نامتقارن خواهد بود. در صورت متقارن بودن، \tilde{A} را با $(a, s)_T$ نشان خواهیم داد. به عبارت دیگر، در حالت $\tilde{A} = (a, s)_T$ ، تابع عضویت $\tilde{A}(x)$ به صورت زیر خواهد بود:

$$\tilde{A}_{(x)} = \begin{cases} 1 - \frac{a-x}{s} & a - s \leq x \leq a \\ 1 - \frac{x-a}{s} & a < x \leq a + s \end{cases} \quad (7)$$

تابع عضویت ضریب فازی \tilde{A} که یک عدد فازی مثلثی فرض می‌شود به صورت زیر است که در شکل (۱) به تصویر کشده شده است.

$$\tilde{A}_{(x)} = \begin{cases} 1 - \frac{|a_i - x|}{s_i} & a_i - s_i \leq x \leq a_i + s_i \\ 0 & \text{در غیر این صورت} \end{cases}$$

شکل (۱)



تابع عضویت عدد فازی \tilde{A} تابع دو پارامتر a, s است، به طوری که پارامتر a ، پارامتر میانه و پارامتر s پارامتر گستره عددی فازی است که نشان‌دهنده میزان فازی بودن عدد است. در واقع، عدد فازی \tilde{A} با دو پارامتر a, s نشان‌دهنده تقریباً a است. پارامترهای فازی $\tilde{A} = \{\tilde{A}_0, \tilde{A}_1, \dots, \tilde{A}_p\}$ می‌توانند به صورت بردار $\tilde{A} = \{a, s\}$ نشان داده شوند، به قسمی که داریم:

$$a = (a_0, a_1, \dots, a_p)$$

$$s = (s_0, s_1, \dots, s_p)$$

بنابراین، خروجی رگرسیون می‌تواند به صورت ذیر نشان داده شود:

$$\tilde{y} = (a_0, s_0) + (a_1, s_1)x_1 + (a_2, s_2)x_2 + \dots + (a_p, s_p)x_p \quad (8)$$

در نتیجه، تابع عضویت متغیر خروجی رگرسیون به صورت ذیر بدست می‌آید:

$$\mu_{\tilde{y}}(y) = \begin{cases} \max(\min\{\tilde{A}(x)\}) & \{x \mid y = f(x, a) = \phi\} \\ 0 & \text{در غیر این صورت} \end{cases} \quad (9)$$

با جایگزینی رابطه (۹) در رابطه (۸) خواهیم داشت:

$$\mu_{\tilde{y}}(y) = \begin{cases} 1 - \frac{|y - \sum_{i=1}^n a_i x_i|}{\sum_{i=1}^n s_i |x_i|} & x_i \neq 0 \\ 1 & x_i = 0, y = 0 \\ 0 & x_i = 0, y \neq 0 \end{cases} \quad (10)$$

مثال: فرض کنید یک مدل رگرسیون فازی به شکل زیر باشد:

$$\tilde{Y} = \tilde{A}_1 x_1 + \tilde{A}_2 x_2$$

پارامترهای مدل اعداد فازی مثلثی فرض شده و با پارامترهای \mathbf{a} و \mathbf{s} به شرح زیر تعریف می‌شوند:

$$\tilde{A}_1 = (4, 2) \quad , \quad \tilde{A}_2 = (5, 3)$$

تابع رگرسیون به صورت زیر است:

$$\tilde{Y} = (4, 2)x_1 + (5, 3)x_2$$

در نتیجه، خروجی مدل \tilde{Y} نیز یک عدد فازی مثلثی خواهد بود که تابع عضویت آن به صورت زیر به دست می‌آید:

$$\mu_{\tilde{y}}(y) = \begin{cases} 1 - \frac{|y - 4x_1 - 5x_2|}{2|x_1| + 3|x_2|} & x_1, x_2 \neq 0 \\ 1 & x_1 = x_2 = y = 0 \\ 0 & x_1 = x_2 = 0, y \neq 0 \end{cases}$$

در این صورت، پارامترهای a_y و s_y برای خروجی مدل رگرسیون فازی به شرح زیر است:

$$a_y = 4x_1 + 5x_2$$

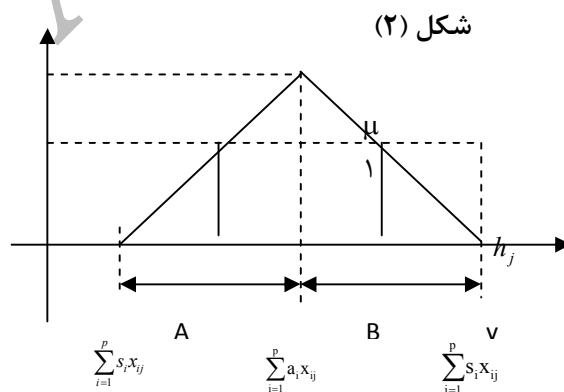
$$s_y = 2|x_1| + 3|x_2|$$

حالت داده‌های غیر فازی در رگرسیون می‌تواند تبدیل به یک مدل برنامه‌ریزی خطی شود. در این حالت، هدف مدل رگرسیون تعیین مقادیر بهینه پارامترهای \tilde{A} است، به قسمی که مجموعه فازی خروجی مدل رگرسیون شامل (y_i) دارای درجه عضویت بزرگ‌تر یا مساوی h باشد. یعنی:

$$\mu_{\tilde{y}_j}(y_j) \geq h_j \quad j=1, 2, \dots, n \quad (11)$$

که n تعداد مشاهدات است و درجه h را کاربر تعیین می‌کند. با افزایش h ، میزان فازی بودن خروجی کاهش می‌یابد. شکل (۲)

تابع عضویت خروجی مدل رگرسیون را نشان می‌دهد:



معادله (۱۱) نشان‌دهنده این است که متغیر وابسته فازای می‌بایست بین دو مقدار A و B قرار گیرد (شکل ۲). در این مدل رگرسیونی، هدف یافتن ضرایب فازای است که گستره خروجی فازای را برای همه مقادیر مجموعه داده‌ها حداقل کند:

$$O = \min \left\{ \sum_{i=1}^p \sum_{j=1}^n s_i x_{ij} \right\} \quad (12)$$

تابع هدف مدل در مواجهه با دو محدودیت زیر حداقل می‌شود:

$$y_j \geq \sum_{i=1}^p a_i x_{ij} - (1-h) \sum_{i=1}^p s_i x_{ij} \quad (13)$$

$$y_j \leq \sum_{i=1}^p a_i x_{ij} + (1-h) \sum_{i=1}^p s_i x_{ij} \quad (14)$$

از آنجا که برای هر مشاهده دو نوع محدودیت وجود دارد، تعداد کل محدودیت‌ها $2n$ خواهد بود.

مثال: فرض کنید برای یک مدل رگرسیون با دو متغیر مستقل مقادیر زیر را داریم:

x_1	x_2	y
3	5	12
2	4	8
1	6	11

داده‌های مشاهده‌شده غیر فازای‌اند، اما فرض بر این است که ساختار مدل فازای است. با فرض این که پارامترهای مدل (a, s) مثلثی بوده پارامترهای مدل با حل مسئله برنامه‌ریزی خطی زیر به دست می‌آیند.
تابع هدف:

$$O = \min \{s_1(3+2+1) + s_2(5+4+6)\}$$

محدودیت‌ها:

$$12 \geq 3a_1 + 5a_2 - (1-h)(3s_1 + 5s_2)$$

$$12 \leq 3a_1 + 5a_2 + (1-h)(3s_1 + 5s_2)$$

$$8 \geq 2a_1 + 4a_2 - (1-h)(2s_1 + 4s_2)$$

$$8 \leq 2a_1 + 4a_2 + (1-h)(2s_1 + 4s_2)$$

$$11 \geq a_1 + 6a_2 - (1-h)(1s_1 + 6s_2)$$

$$11 \leq a_1 + 6a_2 + (1-h)(1s_1 + 6s_2)$$

با حل مسئله برنامه‌ریزی خطی فوق فرض می‌کنیم جواب‌های زیر به دست می‌آید (با فرض $h=0.5$)

$$a_1 = 1.04, \quad a_2 = 1.62, \quad s_1 = 0.53, \quad s_2 = 0.00$$

و همچنین، مقدار مینیمم شده تابع هدف یعنی O برابر 3.23 و

$$\tilde{A}_1 = (1.04, 0.53), \quad \tilde{A}_2 = (1.62, 0.00)$$

در این حالت، علاقه‌مندیم به‌ازای مقادیر متغیرهای مستقل، یعنی x_1, x_2 مقدار، y را پیش‌بینی کنیم. در این صورت،

$$x_2 = 6, \quad x_1 = 1$$

$$\tilde{Y}^a(z) = 1.04 \times 1 + 1.62 \times 6 = 10.76$$

$$\tilde{Y}^s(x) = 0.53 \times 1 + 0.00 \times 6 = 0.53$$

مقدار مشاهده‌شده، یعنی ۱۱، قابل مقایسه است. با تغییر h مقادیر a_i ها تغییری پیدا نمی‌کنند، اما مقادیر s_i ها و O تغییر پیدا می‌کنند.

ارزیابی مدل برازش شده

مدل‌های رگرسیونی فازی با شاخص‌ها و آماره‌های گوناگونی ارزیابی می‌شوند. یکی از شاخص‌های ارزیابی رگرسیون فازی شاخص اطمینان (Index of Confidence) است که با فرمول زیر بیان می‌شود:

$$IC = 1 - \frac{SSE}{SST} \quad (15)$$

که SSE و SST عبارتند از:

$$SSE = 2 \sum_{i=1}^n \left(Y_i - \tilde{Y}_i^a \right)^2$$
$$SST = \sum_{i=1}^n \left(Y_i - \tilde{Y}_i^L \right)^2 + \sum_{i=1}^n \left(\tilde{Y}_i^R - Y_i \right)^2$$

معمولاً در انتخاب مدل‌های گوناگونی را در نظر گرفته و مدلی را انتخاب می‌کنیم که دارای IC کوچک باشد. البته میزان ابهام مدل را نیز باید در نظر گرفت.

مثال کاربردی: بررسی رابطه بین معدل کتبی دیپلم و نمرات ریاضی و فیزیک و جنسیت داوطلبان

از ۱۳۵۷ و بعد از انقلاب اسلامی، تغییراتی اساسی در سیستم پذیرش دانشجو در کشور رخ داده است. هدف از این تغییرات و اصلاحات ایجاد شانس برابر برای همه داوطلبان و تسهیل ورود متقاضیان به دوره‌های آموزش عالی در کنار توجه به افزایش ظرفیت‌های این دوره‌ها بوده است. ظرفیت‌های مطلوب آموزش عالی طی ده سال اخیر تقریباً ثابت مانده است و افزایش ظرفیت‌ها بیشتر در دوره‌های پیام نور و غیرانتفاعی اتفاق افتاده است که از نظر داوطلبان دوره‌هایی چندین مورد نبوده‌اند. از طرف دیگر به‌رغم افزایش ظرفیت‌های دوره‌های مورد اشاره، تعداد متقاضیان ورود به آموزش عالی کاهش نیافته است. از جمله تعداد شرکت‌کنندگان در آزمون سراسری ۱۳۷۹ برابر ۱.۳۳۹.۱۳۴ و در ۱۳۸۷ برابر ۱.۵۰۶.۰۰۰، اما و ظرفیت‌های دوره روزانه دانشگاه‌ها طی ده سال گذشته تقریباً ثابت و به‌طور متوسط حدود ۸۰۰۰۰ نفر بوده است. بنابراین، داوطلبان برای ورود به دانشگاه‌ها و موسسات آموزش عالی به‌منظور ورود به مقاطع بالاتر، مخصوصاً دوره‌های روزانه، در رقابت بسیار فشرده‌ای درگیرند. رقابت فشرده داوطلبان ورود به دانشگاه‌ها و مراکز آموزش عالی در یک آزمون چهار ساعته باعث ایجاد حساسیت، استرس زیاد داوطلبان و مشکلات اجتماعی متعددی در کشور شده است. این امر مراجع قانون‌گذاری، از قبیل مجلس شورای اسلامی و شورای انقلاب فرهنگی، را واداشته است که قوانین و ضوابط متعددی را به تصویب برسانند. آخرین قانون مربوط به آبان ۱۳۸۶ است. در این قانون مجلس شورای اسلامی مصوب کرد که از ۱۳۹۰، سوابق تحصیلی استاندارد دانش‌آموزان جایگزین آزمون سراسری شود. صرف‌نظر از امکان اجرای این قانون، در این بخش علاوه بر بحث و بررسی رگرسیون کلاسیک و فازی تلاش بر این است که به این سوال پاسخ داده شود:

آیا بین سوابق تحصیلی نهایی موجود دانش‌آموزان (معدل کتبی دیپلم یا سال سوم متوسطه) و نمرات آزمون سراسری رابطه‌ای وجود دارد یا خیر؟

در این بررسی، اطلاعات همه شرکت‌کنندگان گروه آزمایشی علوم ریاضی و فنی در سهمیه منطقه یک در آزمون سراسری ۱۳۸۷، شامل نمرات هفت درس امتحانی (ادبیات فارسی، عربی، فرهنگ و معارف اسلامی، زبان خارجی، ریاضیات، فیزیک و مکانیک و شیمی)، معدل کتبی دیپلم و جنس داوطلبان به منزله جامعه مورد بررسی تجزیه و تحلیل می‌شود. جامعه آماری عبارت است از اطلاعات همه شرکت‌کنندگان آزمون سراسری سال ۱۳۸۷ در گروه آزمایشی علوم ریاضی و فنی در منطقه یک به تعداد ۱۰۳۹۳۶ نفر که از بین این افراد ۴۳۴۳ نفر داده گمشده و بقیه، در منطقه یک حدود ۹۹.۵۹۳ نفر برای تجزیه و تحلیل استفاده شده است. در ادامه، با توجه به این که بحث و بررسی رگرسیون خطی فازی هدف اصلی این مقاله است، نمونه‌ای به حجم ۷۴ را با روش نمونه‌گیری تصادفی ساده و بدون جایگذاری با خطای ۸ درصد استخراج و با دو روش رگرسیون کلاسیک و فازی بررسی و

تجزیه و تحلیل می‌کنیم. نرم افزار مورد استفاده برای رگرسیون کلاسیک SPSS و برای رگرسیون فازی Lingo11.0 و Excel است. این بخش به سه قسمت زیر تقسیم می‌شود:

الف) رگرسیون کلاسیک چندگانه معدل روی نمرات و جنس در جامعه آماری

ب) رگرسیون کلاسیک چندگانه معدل روی نمرات و جنس در نمونه آماری

ج) رگرسیون فازی چندگانه معدل روی نمرات و جنس در نمونه آماری

الف) رگرسیون کلاسیک چندگانه معدل روی نمرات و جنس در جامعه آماری

پس از اجرای رگرسیون چندگانه گام به گام معدل کتبی دیپلم (متغیر وابسته) روی نمرات هفت درس امتحانی (ادبیات فارسی، عربی، فرهنگ و معارف اسلامی، زبان خارجی، ریاضیات، فیزیک و مکانیک و شیمی) و جنس همه داوطلبان منطقه یک گروه آزمایشی علوم ریاضی و فنی (جامعه آماری شامل ۹۹۵۹۳ نفر، ۴۴.۵ درصد زن و ۵۵.۵ درصد مرد) مدل زیر به دست می‌آید.

Sig.	t	Standardized Coefficients	Unstandardized Coefficients		Model
			Std. Error	B	
.۰۰۰۰	۱۲۲۵.۶		۱.۰۱۲	۱۲۴۰.۰	1 ضریب ثابت a
	۷۱			۱۷	
.۰۰۰۰	۱۰۴.۲۱	.۲۴۴	۱.۳۰۳	۱۳۵.۷۴	جنس jnc
	۲			۵	
.۰۰۰۰	۸۸.۱۱۱	.۰۳۷۰	.۰۰۰۸	.۰۰۶۷۴	ریاضی ng
				5	
.۰۰۰۰	۷۱.۰۶۵	.۰۲۹۸	.۰۰۰۷	.۰۰۴۸۲	فیزیک ng
				6	

چنان که ملاحظه می‌شود، سایر متغیرهای مستقل به علت عدم معنی دار بودن از مدل حذف شده‌اند. آزمون‌های آماری از قبیل F و نرمال بودن باقیمانده‌ها با آزمون کلموگروف و اسمیرینف در سطح ۰.۰۵ تایید شده‌اند. ضریب تبیین یا R^2 برابر ۰.۴۵۳ است. به عبارت دیگر، ۴۵.۳ تغییرات متغیر معدل کتبی دیپلم با متغیرهای جنسیت و نمرات ریاضی و فیزیک داوطلبان تبیین می‌شود.

ب) رگرسیون کلاسیک چندگانه معدل روی نمرات و جنس در نمونه آماری

در این مرحله، یک نمونه به حجم ۷۵ به روش نمونه‌گیری تصادفی ساده و بدون جایگذاری با خطای ۸ درصد استخراج می‌شود. پس از اجرای رگرسیون چندگانه گام به گام معدل کتبی دیپلم (متغیر وابسته) روی نمرات دروس امتحانی ریاضیات، فیزیک و جنس روی نمونه استخراج شده (نمونه آماری شامل ۱۰۰ نفر، ۴۴ درصد زن و ۵۶ درصد مرد) مدل زیر به دست می‌آید.

Sig.	t	Standardized Coefficients	Unstandardized Coefficients		Model
	B	Beta	Std. Error	B	
۰.۰۰۰	۳۳.۳۰ ۲		۳۷.۸۲۰	۱۲۵۹.۴ ۶۴	ضریب ثابت
۰.۰۵۲	۱.۸۹۹	۰.۱۸۳	۴۹.۴۳۵	۹۳.۸۵۷	n جنس
۰.۰۰۴	۳.۰۰۲	۰.۴۳۷	۰.۲۳۳	۰.۶۹۹	n ریاضی
۰.۰۴۰	۱.۸۲۸	۰.۲۶۳	۰.۲۰۲	۰.۳۷۰	n فیزیک

ضریب تبیین یا R^2 برابر ۰.۴۱ است. به عبارت دیگر، ۴۱ تغییرات متغیر معدل کتبی دیپلم با متغیرهای جنسیت و نمرات ریاضی و فیزیک داوطلبان در نمونه استخراج شده در رگرسیون تبیین می‌شود. همچنین، آزمون F یعنی معنی‌داری کل رگرسیون برازش شده نیز معنی‌دار است. نکته مهم عدم نرمال بودن باقیمانده‌ها با آزمون کلموگروف و اسمیرینف در سطح 0.05 است که با توجه به کم بودن حجم نمونه است و برای حل این مشکل رگرسیون فازی پیشنهاد می‌شود که در بخش بعد به آن خواهیم پرداخت. با توجه به جداول بالا نتایج به دست آمده در نمونه‌ای با خطای ۸ درصد قابل تعمیم به جامعه است.

ج) رگرسیون فازی چندگانه معدل روی نمرات و جنس در نمونه آماری

تئوری مجموعه‌ها و رگرسیون خطی فازی در بخش سوم و چهارم عرضه شده‌اند. در این بخش می‌خواهیم رگرسیون خطی فازی را با فرض ابهام در ضرایب در خصوص معدل کتبی دیپلم (متغیر وابسته) روی متغیرهای مستقل نمرات ریاضی، فیزیک و جنسیت داوطلبان به دست بیاوریم. تعداد نمونه‌ها برابر ۷۴ و متغیرهای جنس، نمره ریاضی و فیزیک را به ترتیب با JNC، NG5 و NG6 نشان می‌دهیم. برای پیدا کردن تابع هدف (۱۲) داریم: $\sum JNC = 27$ ، $\sum NG5 = 9212$ و $\sum NG6 = 9295$. بنابراین:

$$O = 148 + 27 \times S_1 + 9212 \times S_2 + 9295 \times S_3$$

می‌خواهیم تابع O را با فرض قیدهای معادله (۱۳) و (۱۴) که در مجموع ۱۴۸ قید است حداقل کنیم. پس از تشکیل قیدها و حل معادله خطی بالا با توجه به h‌های گوناگون، مقادیر $a_0, a_1, a_2, a_3, s_0, s_1, s_2, s_3$ براساس جدول شماره ۱ با استفاده از نرم-افزار Lingo 11.0 چنین به دست می‌آید:

جدول ۱. نتایج حاصل از برازش توابع معدل، نمرات و جنس

h	a_0	s_0	a_1	s_1	a_2	s_2	a_3	s_3
۰.۲	۱۱۸۳	۳۰.۶	۲۵۴	۱۵۷	۰.۵۴	۱.۳	۰.۳۳	۰
۰.۴	۱۱۸۳	۴۰.۸	۲۵۴	۲۱۰	۰.۵۴	۱.۷	۰.۳۳	۰
۰.۵	۱۱۸۳	۴۸.۹	۲۵۴	۲۵۲	۰.۵۴	۲.۰	۰.۳۳	۰
۰.۶	۱۱۸۳	۶۱.۲	۲۵۴	۳۱۵	۰.۵۴	۲.۵	۰.۳۳	۰
۰.۸	۱۱۸۳	۱۲۲.۴	۲۵۴	۶۳۰	۰.۵۴	۵	۰.۳۳	۰

پس از پیدا کردن مقادیر a_i و S_i ها، بایستی مقادیر مرکزی (Y^a) و پراکندگی داده‌ها ($Y_{s0.8}, Y_{s0.5}, Y_{s0.2}, Y_{s0.4}, Y_{s0.6}$) را مانند مثال بخش قبلی به دست بیاوریم که معادلات زیر خواهند بود:

$$Y^a = 1183.517 + 254.1639 * jnc + 0.5359716 * ng5 + 0.334554 * ng6.$$

$$Y_{s0.2} = 306.2444 + 157.5180 * jnc + 1.269992 * ng5 + 0.000000 * ng6.$$

$$Y_{s0.4} = 408.3258 + 210.0240 * jnc + 1.693323 * ng5 + 0.000000 * ng6.$$

$$Y_{s0.5} = 489.9910 + 252.0288 * jnc + 2.031987 * ng5 + 0.000000 * ng6.$$

$$Y_{s0.6} = 612.4887 + 315.0360 * jnc + 2.539984 * ng5 + 0.000000 * ng6.$$

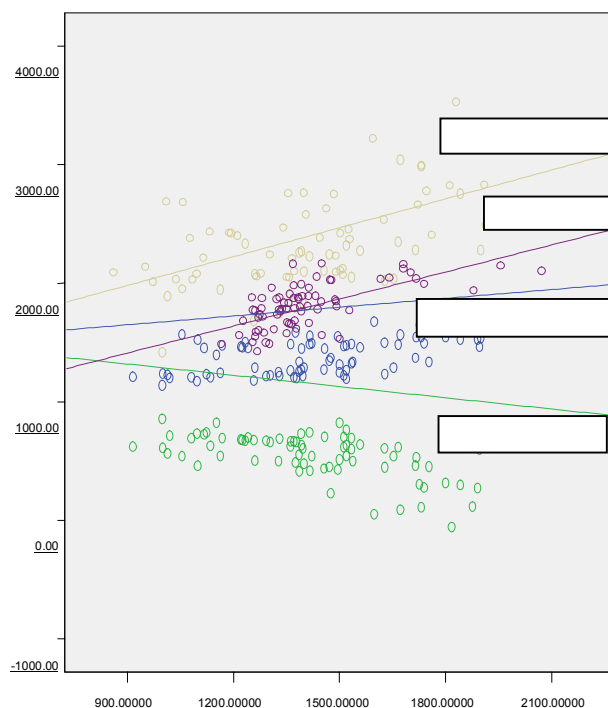
$$Y_{s0.8} = 1224.977 + 630.0720 * jnc + 5.079968 * ng5 + 0.000000 * ng6.$$

پس از پیدا کردن کران‌های بالا و پایین برای h های متفاوت، شاخص اطمینان IC با استفاده از فرمول ۱۵ به شکل جدول شماره ۲ بدست می‌آید:

جدول ۲. مقادیر شاخص اطمینان و ابهام کل برای h های متفاوت

h	$\frac{SSE}{SST}$	IC	ابهام کل
۰.۲	۰.۱۳	۰.۸۷	۶۱۲۷۶
۰.۴	۰.۰۷	۰.۹۳	۸۱۷۰۱
۰.۵	۰.۰۵	۰.۹۵	۹۸۰۴۲
۰.۶	۰.۰۴	۰.۹۶	۱۲۲۵۵۲
۰.۸	۰.۰۱	۰.۹۹	۲۴۵۱۰۵

نمودار ۳. نمودار رگرسیون فازی پیش‌بینی کننده (شامل تابع کران بالا و پایین و مرکزی) برای $h=0.5$



به منظور ارزیابی مدل‌های حاصل از رگرسیون فازی، شاخص اطمینان و مقدار ابهام کل مدل برای h ‌های متفاوت محاسبه و در جدول شماره ۲ آمده است. چنان‌که ملاحظه می‌شود، با بزرگ شدن مقدار h سطح اعتبار مدل بالا می‌رود و در عین حال، این امر باعث افزایش ابهام کل مدل نیز می‌شود. لذا برای انتخاب h صحیح به میزان ابهام کل نیز باید توجه داشت. البته می‌توان سطح اعتبار $h=0.5$ را به منزله سطح اعتباری معقول و متداول در نظر گرفت. در داده‌های مورد بررسی سطح اعتبار ۰.۵ به سبب IC بالا و افزایش اندک ابهام مدل معقول به نظر می‌رسد. ضریب همبستگی بین پیش‌بینی‌کننده‌های رگرسیون کلاسیک و فازی مرکزی ۰.۹ است که نشان‌دهنده نزدیک بودن نتایج دو روش است. اما رگرسیون فازی این امتیاز را دارد که فاقد فرض‌های محدودکننده رگرسیون کلاسیک است، مخصوصاً در این نمونه که توزیع باقیمانده‌ها نیز نرمال نیست. علاوه بر این، چنان‌که در نمودار ۳ ملاحظه می‌شود، رگرسیون کلاسیک برازش‌شده در داخل فاصله اطمینان رگرسیون فازی برازش‌شده با $h=0.5$ قرار می‌گیرد. بنابراین، به نظر می‌رسد در نمونه‌های کوچک، به علت عدم پایایی فرض‌های رگرسیون کلاسیک، به کار بردن رگرسیون خطی فازی معقول‌تر است.

خلاصه و نتیجه گیری

در این مقاله، روش‌های رگرسیون کلاسیک و فازی با تأکید بر بعد آموزشی برای علاقه‌مندان رشته‌های علوم اجتماعی به بحث و بررسی گذاشته شد و یک مثال کاربردی در این حوزه به تفصیل تجزیه و تحلیل شد. بایستی توجه داشت که در کاربرد روش‌های فازی در علوم اجتماعی و حتی سایر علوم با توجه به ماهیت علم فازی، مرزبندی‌های تصمیم‌گیری به روشنی روش‌های کلاسیک نیست و در اکثر موارد محقق است که باید تصمیم بگیرد. در این خصوص می‌توان به انتخاب سطح فازی بودن یا به کار بردن اعداد فازی غیر مثلثی، مانند اعداد فازی نرمال، اشاره کرد.

منابع

- آذر، عادل و فرجی، حجت (۱۳۸۶) علم مدیریت فازی، تهران: دانشگاه تربیت مدرس.
- امینی، امیرمظفر و مهدی خیاطی (۱۳۸۵) "عوامل مؤثر بر عدم موفقیت طرح تشکیل تعاونی‌های آب ایران (استفاده از رگرسیون فازی)" *اقتصاد کشاورزی و توسعه*، شماره ۵۳، ۱۴ (۱): ۶۹-۹۱.
- ایمانی، محمدتقی (۱۳۷۶) "ابهام روش‌شناسی و تنگناهای پژوهش در ایران"، در مجموعه مقالات نخستین سمینار آموزش عالی ایران، تهران: انتشارات دانشگاه علامه طباطبایی.
- بارت کاسکو (۱۳۸۰) *تفکر فازی*، ترجمه علی غفاری، عادل مقصودپور، علیرضا پورممتاز و جمشید قسیم، تهران: دانشگاه خواجه نصیرالدین طوسی.
- تیموری، حبیب‌الله (۱۳۸۶) *آمار در علوم اجتماعی*، تهران: کیان مهر.
- طاهری، محمود (۱۳۸۷) *مقدمه‌ای بر احتمال و آمار فازی*، کرمان: دانشگاه شهید باهنر کرمان.
- لی وانگ (۱۳۷۸) *سیستم‌های فازی و کنترل فازی*، ترجمه محمد تشنه لب، نیما صفارپور، داریوش افیونی، تهران: دانشگاه خواجه نصیرالدین طوسی.
- محمدی، جهانگرد؛ طاهری، محمود (۱۳۸۴) "برازش توابع انتقالی خاک با استفاده از رگرسیون فازی"، در *علوم و فنون کشاورزی و منابع طبیعی*، ۹ (۲): ۵۱-۶۱.
- Arnold F. Shapiro (2004) *Fuzzy Regression Models*, Pennsylvania: Pen State University.
- Chang Ping-Teng Kan (1997) "Fuzzy Regression Analysis" in *Fuzzy Sets and Systems*, volume 90.
- E. Pasha, T. Razzaghnia1, T. Allahviranloo, Gh. Yari, H. R. Mostafaei, (2007) "Fuzzy Linear Regression Models with Fuzzy Entropy", *Applied Mathematical Sciences*, 1 (35): 1715 – 1724.
- Galton, F. (1885) "Regression toward Mediocrity in Heredity Stature" in *Journal of the Anthropological Institute*, : 15: 246-263.
- Modarres. M, E. Nasrabdi and M. M Nasrabadi (2004). "Fuzzy linear Regression Analysis from the Risk Point of view", *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 12 (5): 635-649.
- Tanaka, H., S. Uejima and K. Asai (1982) "Linear Regression Analysis with Fuzzy Model" in *IEEE Transactions on Systems, Man, and Cybernet*, 12: 903-907.
- Liang. W and Chen. Y (2004) "Estimation of Weibull Parameters using a Fuzzy Least Squares", *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 12 (5): 701-711
- Zadeh, L. A. (1965) *Fuzzy sets, Inform and control*, 8: 338-353.