

حاشا
الرحمن الرحيم



دانشکده علوم ریاضی

رشته آمار، گرایش آمار ریاضی

پایان نامه کارشناسی ارشد

رگرسیون لجستیک در محیط فازی و کاربردهای آن در علوم پزشکی

نگارنده: میترا خراباتی

استاد راهنما

دکتر محمدرضا ربیعی

استاد مشاور

دکتر فاطمه سلمانی

شهریور ۱۳۹۸

شماره:

تاریخ:

باسمه تعالی



دانشگاه تبریز

مدیریت تحصیلات تکمیلی

فرم شماره (۳) صورتجلسه نهایی دفاع از پایان نامه دوره کارشناسی ارشد

با نام و یاد خداوند متعال، ارزیابی جلسه دفاع از پایان نامه کارشناسی ارشد خانم میترا خراباتی با شماره دانشجویی ۹۵۰۵۳۷۴ رشته آمار گرایش ریاضی تحت عنوان رگرسیون لجستیک در محیط فازی و کاربردها آن در علوم پزشکی که در تاریخ ۹۸/۶/۱۱ با حضور هیأت داوران در دانشگاه صنعتی شاهرود برگزار گردید به شرح ذیل اعلام می گردد:

<input type="checkbox"/> مردود <input checked="" type="checkbox"/> قبول (با درجه: عالی.....)			
<input type="checkbox"/> عملی <input checked="" type="checkbox"/> نظری			
امضاء	مرتبه علمی	نام و نام خانوادگی	عضو هیأت داوران
		دکتر محمدرضا ربیعی	۱- استاد راهنمای اول
			۲- استاد راهنمای دوم
		دکتر فاطمه سلمانی	۳- استاد مشاور
		دکتر حسین باغیشی	۴- نماینده تحصیلات تکمیلی
		دکتر محمد آرشی	۵- استاد ممتحن اول
		دکتر احمد نزاکتی	۶- استاد ممتحن دوم



نام و نام خانوادگی رئیس دانشکده: دکتر ابراهیم هاشمی

تاریخ و امضاء و مهر دانشکده:

تصوه: در صورتی که کسی مردود شود حداکثر یکبار دیگر (در مدت مجاز تحصیل) می تواند از پایان نامه خود دفاع نماید (دفاع مجدد نباید زودتر از ۴ ماه برگزار شود).

تقدیم به روح آسمانی پدرم
و وجود پربرکت مادرم که همواره مشوق من
در کسب علم و دانش بوده اند.

سپاس بی‌کران خدای راست، که دستانم را توان نوشتن از اوست تا به انجام برسانم.

در این مجال، فرصت را غنیمت شمرده، از زحمات تمامی بزرگوارانی که در این راه پشتیبانم بودند، قدردانی می‌نمایم.

وظیفه خود می‌دانم از جناب آقای دکتر محمدرضا ربیعی، استاد راهنمای ارجمندم، به خاطر تمامی محبت‌ها و راهنمایی‌های بی‌شائبه ایشا در تمام مراحل پایان‌نامه کمال تشکر را بنمایم.

از سرکار خانم دکتر فاطمه سلمانی، استاد مشاور گرامی که همواره با صبر و حوصله پاسخگوی سؤالات بنده بودند، متشکرم.

از داوران محترم، جناب آقای دکتر محمد آرشی و جناب آقای دکتر احمد نزاکتی و نماینده محترم تحصیلات تکمیلی، جناب آقای دکتر حسین باغیشنی بسیار متشکرم.

در انتها بر خود لازم می‌دانم از خانواده خود به‌ویژه مادر بزرگوارم و دوست عزیزم خانم کتایون مسلمی، که با محبت‌های بی‌دریغشان در طی این دوران، هر چند از راه دور، بنده را مورد حمایت خودشان قرار دادند، قدردانی می‌نمایم.

از پروردگار منان خواستارم که به من توان جبران گوشه‌ای از آن همه لطف و محبت تمامی بزرگواران را ارزانی فرماید.

میترا خراباتی

شهریور ۱۳۹۸

تعهد نامه

اینجانب **میترا خراباتی** دانشجوی کارشناسی ارشد رشته **آمار علوم ریاضی** دانشگاه صنعتی شاهرود، نویسنده پایان نامه با عنوان **رگرسیون لجستیک در محیط فازی و کاربردهای آن در علوم پزشکی**، تحت راهنمایی **محمدرضا ربیعی** متعهد می شوم:

- تحقیقات در این پایان نامه توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است.
- در استفاده از نتایج پژوهش های دیگر پژوهش گران، به مرجع مورد استفاده استناد شده است.
- مطالب این پایان نامه، تا کنون توسط خود، یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارایه نشده است.
- حقوق معنوی این اثر، به دانشگاه صنعتی شاهرود تعلق دارد، و مقالات مستخرج با نام “دانشگاه صنعتی شاهرود” یا “Shahrood University of Technology” به چاپ خواهد رسید.
- حقوق معنوی تمام افرادی که در به دست آوردن نتایج اصلی پایان نامه تاثیرگذار بوده اند، در مقالات مستخرج از پایان نامه رعایت می گردد.
- در تمام مراحل انجام این پایان نامه، در مواردی که از موجود زنده (یا بافت های آنها) استفاده شده است، ضوابط و اصول اخلاقی رعایت شده است.
- در تمام مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته (یا استفاده شده است)، اصل رازداری و اصول اخلاق انسانی رعایت شده است.

میترا خراباتی

شهریور ۱۳۹۸

مالکیت نتایج و حق نشر

- تمام حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه های رایانه ای، نرم افزارها و تجهیزات ساخته شده) متعلق به دانشگاه صنعتی شاهرود می باشد. این مطلب باید به نحو مقتضی، در تولیدات علمی مربوطه ذکر شود.
- استفاده از اطلاعات و نتایج موجود در این پایان نامه بدون ذکر منبع مجاز نمی باشد.

چکیده

رگرسیون لجستیک، در مدل‌سازی داده‌های آماری گسسته کمک می‌کند. مدل‌های آماری با داده‌های گسسته در علوم پزشکی (مانند مرگ/حیات، وجود بیماری/عدم بیماری و...) کاربرد فراوان دارند. در تشخیص بیماری، پزشکان از منابع اطلاعاتی مختلفی، مانند تاریخچه بالینی بیماری، معاینات بدنی، تست‌های آزمایشگاهی و غیره استفاده می‌کنند. اما قرار دادن افراد در دو گروه بیمار و سالم اصولاً در هاله‌ای از ابهام است و مشاهدات مبهم در تشخیص بالینی، فراوان دیده می‌شود. در چنین مواردی به دلیل نادقیق بودن متغیر پاسخ، رگرسیون لجستیک مناسب نیست. برای مدل‌ساز مشاهدات پاسخ دودویی فازی، مدل رگرسیون لجستیک فازی در این پایان‌نامه ارائه شده است. برای محاسبه امکان موفقیت در مدل لجستیک فازی، اصطلاح ترم‌های زبانی مانند «کم»، «متوسط»، «زیاد»، «...» تعریف می‌شود. سپس، با استفاده از اصل گسترش، تبدیل لجیت «نسبت بخت امکانی» براساس مجموعه‌ای از مشاهدات متغیرهای تبیینی دقیق مدل‌سازی می‌شود. همچنین برای برآورد ضرایب مدل پیشنهادی از روش امکانی و روش کمترین توان دوم فازی استفاده می‌کنیم. همچنین برای ارزیابی مدل، دو روش انتخاب مدل ارائه می‌دهیم.

کلمات کلیدی: رگرسیون فازی، رگرسیون لجستیک، رگرسیون لجستیک فازی، کمترین توان دوم خطا، انتخاب مدل

لیست مقالات مستخرج از پایان نامه

۱. خراباتی، میترا. ربیعی، محمدرضا، (۱۳۹۸). شیوه ای جدید در تحلیل رگرسیون لجستیک فازی، نهمین سمینار آمار و احتمال فازی، بابلسر
۲. ربیعی، محمدرضا. خراباتی، میترا (۱۳۹۸). معرفی بسته fuzzyreg برای تحلیل برخی از روشهای رگرسیون فازی در R، نهمین سمینار آمار و احتمال فازی، بابلسر

فهرست مطالب

فهرست تصاویر

فهرست جداول

۱	مقدمات و مفاهیم رگرسیون لجستیک	۱
۱	مقدمه	۱.۱
۲	رگرسیون	۲.۱
۳	مدل رگرسیون خطی	۱.۲.۱
۳	مدل رگرسیون خطی تعمیم یافته	۲.۲.۱
۴	رگرسیون لجستیک	۳.۱
۵	میانگین تابع پاسخ زمانی در حال دودویی	۱.۳.۱
۷	تابع پاسخ لجستیک	۲.۳.۱
۸	رگرسیون لجستیک ساده	۴.۱
۸	مدل رگرسیون لجستیک ساده	۱.۴.۱
۹	تابع درست‌نمایی	۲.۴.۱
۱۰	برآورد پارامتر	۳.۴.۱
۱۳	تعبیر ضریب b_1	۴.۴.۱
۱۴	رگرسیون لجستیک چندگانه	۵.۱
۱۵	برآورد پارامترهای لجستیک چندگانه	۱.۵.۱
۱۷	مجموعه‌ها و رگرسیون فازی	۲
۱۷	مجموعه‌های فازی	۱.۲
۲۱	معرفی تابع f_{apply}	۱.۱.۲
۲۱	اعداد فازی	۲.۲
۲۳	حساب اعداد فازی	۳.۲
۲۷	فاصله بین دو عدد فازی	۱.۳.۲

۲۸	رگرسیون فازی	۴.۲
۳۰	رگرسیون امکانی	۵.۲
۳۱	رگرسیون کمترین توان‌های دوم فازی	۶.۲
۳۲	مدل رگرسیون خطی با ضرایب فازی	۱.۶.۲
۳۵		رگرسیون لجستیک فازی با رویکرد امکانی	۳
۳۵	معرفی مدل رگرسیون لجستیک فازی	۱.۳
۳۷	برازش مدل رگرسیون لجستیک با رویکرد امکانی	۱.۱.۳
۳۹	برآورد ضرایب فازی	۲.۱.۳
۴۰	معیار نیکویی برازش	۳.۱.۳
۴۱	انتخاب متغیر در رگرسیون لجستیک فازی	۲.۳
۴۲	روش پیشرو برای انتخاب متغیر به شیوه سلمانی و همکاران	۱.۲.۳
۴۲	معیار LE	۲.۲.۳
۴۳	الگوریتم انتخاب متغیر به روش پیشرو	۳.۲.۳
۴۳	کاربرد رگرسیون لجستیک فازی در بیماری لوپوس	۴.۲.۳
۴۹		رگرسیون لجستیک فازی با رویکرد کمترین توان دوم	۴
۴۹	روش مجموع توان دوم برای داده‌های فازی	۱.۴
۵۰	برآورد پارامترها با متر ژو	۱.۱.۴
۵۲	مثال کاربردی	۲.۱.۴
۵۵	روش پیشنهادی برای برآورد پارامترها	۲.۴
۵۷	مثال کاربردی	۱.۲.۴
۵۸	معیار نیکویی برازش	۲.۲.۴
۵۹	روش پیشرو برای انتخاب متغیر به شیوه کیی و همکاران	۳.۲.۴
۶۵		مراجع	
۶۹		آ معرفی بسته FuzzyNumbers	
۷۷		ب کدهای استفاده شده در پایان نامه	
۷۷	کدهای انتخاب مدل با معیار LE	۱.ب
۸۳	کدهای متر ژو	۲.ب
۸۴	کدهای متر کیی و همکاران	۳.ب

فهرست تصاویر

۷ نمودار چگالی نرمال (-) و چگالی لجستیک (-)	۱.۱
۴۵ نمودار پیشنهادی μ_i	۱.۳
۴۸ الگوریتم انتخاب متغیر سلمانی در مدل رگرسیون لجستیک فازی	۲.۳
۶۱ الگوریتم انتخاب متغیر کپی و همکاران در مدل رگرسیون لجستیک فازی	۱.۴

فهرست جداول

۴۵	موارد مشکوک به لوپوس و عوامل خطر مربوط به آن	۱.۳
۴۶	مقادیر مشاهده شده \tilde{y}_i	۲.۳
۴۶	مرحله اول انتخاب بهترین مدل	۳.۳
۴۷	مرحله دوم انتخاب بهترین مدل	۴.۳
۴۷	مرحله سوم انتخاب بهترین مدل	۵.۳
۴۸	مرحله چهارم انتخاب بهترین مدل	۶.۳
۶۰	مراحل انتخاب مدل	۱.۴

فصل ۱

مقدمات و مفاهیم رگرسیون لجستیک

۱.۱ مقدمه

تشخیص بیماری غالباً یک فرآیند پیچیده است که تبحر و تجربه‌ی خاصی را می‌طلبد. پزشکان برای تشخیص بیماری از منابع اطلاعاتی گوناگون استفاده می‌کنند. تاریخچه‌ی بالینی بیمار، معاینات بدنی، تست‌های آزمایشگاهی، نتایج دستگاه‌های اندازه‌گیری و شرایط و علائم از پیش تعیین شده همگی معیار تشخیص بیماری قرار می‌گیرند. اما عملاً در بسیاری از مواقع، طبقه‌بندی افراد در دو گروه بیمار و سالم هم‌چنان در هاله‌ای از ابهام است. صحبت‌های ذهنی و تا حدی غلوآمیز بیمار درباره‌ی سابقه‌ی بیماری خود، اشتباه پزشک در معاینات بدنی به‌ویژه در مورد بیماری‌هایی که علائم بالینی مشترکی ندارند، نادقیق بودن مرز بین افراد سالم و بیمار در تست‌های آزمایشگاهی، ابهام در تجزیه و تحلیل نتایج حاصل از دستگاه‌های اندازه‌گیری مانند سونوگرافی و بالاخره عدم توافق کلی در تعریف علائم بیماری‌هایی که برای تشخیص آن‌ها آزمایشات کلینیکی خاصی وجود ندارند، از جمله مواردی هستند که منجر به مشاهدات نادقیق در مطالعات بالینی می‌شوند. مدل‌سازی و تحلیل این گونه نتایج با استفاده از روش‌های معمول آماری امکان‌پذیر نمی‌باشد. برای مدل‌سازی و تحلیل چنین مشاهداتی، به‌خصوص در مواردی که حجم نمونه کم است، مدل‌سازی‌های مبتنی بر منطق و مجموعه‌های فازی راه‌گشا هست.

در فصل اول، به معرفی رگرسیون لجستیک می‌پردازیم. ابتدا رگرسیون خطی تعمیم یافته

را معرفی می‌کنیم، سپس به معرفی رگرسیون لجستیک ساده و رگرسیون لجستیک چندگانه می‌پردازیم و در نهایت فاصله اطمینان و آزمون‌های مرتبط با رگرسیون لجستیک را بررسی می‌کنیم. مطالب این بخش به صورت عمده از جان نتر و همکاران [۲۶] گرفته شده‌اند.

۲.۱ رگرسیون

تحلیل رگرسیونی یک روش آماری قدرتمند برای مدل‌سازی رابطه‌ی بین یک متغیر وابسته و یک یا چندین متغیر مستقل است. از اهداف اصلی تحلیل رگرسیونی کشف رابطه‌ی تابعی بین متغیر وابسته و متغیرهای مستقل است که به منظور کنترل مقادیر متغیر وابسته و یا پیش‌بینی آن در آینده صورت می‌گیرد. مدل ارتباطی حاکم بر متغیرهای پاسخ و مجموعه متغیرهای تبیینی، توسط افراد متخصص در زمینه مورد مطالعه، بر مبنای دانش یا قضاوت‌های عینی و ذهنی آن‌ها تعیین می‌شود. این مدل که می‌تواند پارامترهای نامعلوم زیادی را در بر داشته باشد، یک مدل پارامتری نامیده می‌شود. اگر در معادله رگرسیونی تنها یک متغیر تبیینی وجود داشته باشد، آن را مدل رگرسیونی خطی ساده می‌نامند. در یک مدل رگرسیونی، هدف بررسی اثر متغیرهای تبیینی بر روی متغیر پاسخ است. معمولاً، متغیرهای تبیینی را با x و متغیرهای پاسخ را با y نمایش می‌دهند. همچنین این ضابطه می‌تواند خطی یا غیرخطی باشد. باید توجه داشته باشیم که منظور از رابطه خطی یا غیرخطی، رابطه بین y و x ‌ها نیست، بلکه در واقع پارامترهای رگرسیونی به‌طور خطی یا غیرخطی وارد معادله رگرسیون می‌شوند. به‌عنوان مثال، هریک از روابط زیر خطی‌اند؛ هر چند بین x و y ممکن است خطی نباشد:

$$y = \beta_0 + \beta_1 x + \epsilon$$

$$y = \beta_0 + \beta_1 e^x + \epsilon$$

$$y = \beta_0 + \beta_1 \log x + \epsilon$$

همچنین ϵ جمله خطا می‌باشد که متغیری تصادفی و غیر قابل مشاهده است.

مثال‌هایی از یک مدل رگرسیونی عبارتند از:

۱. کشف رابطه بین سن (متغیر تبیینی) و وجود بیماری کرونر قلبی (متغیر پاسخ)
۲. بررسی تاثیر میزان قند و کلسیم خون (متغیر تبیینی) بر مبتلا بودن فرد به بیماری دیابت (متغیر پاسخ)
۳. بررسی تاثیر سن (متغیر تبیینی) بر تعدا ضربان قلب (متغیر تبیینی)

در مدل رگرسیون خطی، متغیر پاسخ Y یک متغیر تصادفی پیوسته است، در حالی که x تصادفی نبوده و توسط تحلیل گر کنترل و با خطای قابل اغماضی اندازه‌گیری می‌شود. بنابراین به‌ازای هر مقدار ممکن x برای y یک توزیع احتمال وجود دارد. به‌طور کلی، منظور از مدل

رگرسیونی این است که $E(y|x)$ را بر حسب تابعی از x ، یا y را بر حسب تابعی از x ، به همراه یک جمله خطا، ϵ ، که امید ریاضی آن صفر می‌شود، بنویسیم. بنابراین داریم:

$$y = \beta_0 + \beta_1 x + \epsilon, \quad E(Y|X) = \beta_0 + \beta_1 x$$

هر دو معادل‌اند.

۱.۲.۱ مدل رگرسیون خطی

مدل ارتباطی (ضابطه ریاضی) حاکم بر متغیر پاسخ و مجموعه متغیرهای تبیینی، توسط افراد متخصص در زمینه مورد مطالعه، بر مبنای دانش یا قضاوت‌های عینی و ذهنی آن‌ها تعیین می‌شود. فرض کنید که خط رگرسیون y نسبت به x به صورت $y = \beta_0 + \beta_1 x$ باشد، آن‌گاه مدل خطی مرتبه اول را به صورت

$$y = \beta_0 + \beta_1 x + \epsilon$$

می‌توانیم بنویسیم. یعنی برای x داده شده، مقدار مشاهده شده y مربوط به آن عبارت است از: $\beta_0 + \beta_1 x$ به اضافه مقدار خطای ϵ که ممکن است هر y مشخصی خارج از خط رگرسیون قرار گیرد. به بیان دیگر y یک متغیر تصادفی غیرقابل مشاهده و x یک متغیر ریاضی (غیرتصادفی) است. معادله بالا مدلی است که بر اساس گمان و فرض خطی بودن ارتباط بین x و y بیان شده‌است.

برای مدل رگرسیون خطی چهار پذیره در نظر گرفته می‌شود:

۱. خطی بودن ضابطه تابع رگرسیون $E(y|x)$

۲. ثابت بودن واریانس جمله خطا

۳. ناهمبسته بودن مولفه‌های خطا

۴. (بر حسب نیاز) نرمال بودن خطا

تخطی از هر کدام از این پذیره‌ها (به ویژه پذیره اول)، می‌تواند منجر به نامناسب بودن مدل مورد نظر شود.

۲.۲.۱ مدل رگرسیون خطی تعمیم یافته

در سال‌های اخیر، مدل‌های خطی تعمیم یافته (GLM) به طور چشم‌گیری مورد استفاده قرار گرفته‌اند و متون آماری متعددی در این زمینه به رشته تحریر درآمده‌اند. قسمتی از روند رشد این دسته از مدل‌ها را می‌توان بدین شرح نوشت: مدل‌های خطی تعمیم یافته، به عنوان تعمیم مدل‌های خطی، توسط نلدر و ودربرن [۲۵] معرفی شد. بررسی نظری دقیق این دسته

از مدل‌ها در مک کالا و نلدر [۲۳] آمده است. میر و همکاران [۲۲] نیز منبع خوبی برای دیدگاه شهودی این مدل‌هاست.

در مدل‌های خطی تعمیم یافته، توزیع‌هایی به جز نرمال نیز می‌توانند به عنوان توزیع متغیر پاسخ در نظر گرفته شوند. البته توزیع در نظر گرفته شده باید متعلق به خانواده نمایی یا شبیه به این خانواده باشد. [۲۳] میانگین پاسخ به طور مستقیم مدل بندی نمی‌شود بلکه تبدیلی از آن با استفاده از تابع پیوند مدل بندی می‌گردد.

۳.۱ رگرسیون لجستیک

در روش مدل سازی، زمانی که مقادیر متغیر پاسخ، گسسته است دیگر نمی‌توان از روش‌های رگرسیون خطی استفاده کرد. روش رگرسیون لجستیک توانایی مدل سازی این گونه مسائل را برای متغیرهای تبیینی پیوسته و گسسته داراست. مدل رگرسیون لجستیکی توسط کاکس برای توصیف وابستگی یک متغیر دوتایی به مجموعه‌ای از متغیرهای پیوسته معرفی شد. [۷] در کلاس مدل‌های غیرخطی دسته‌ای از مدل‌ها وجود دارند که ذاتاً خطی هستند؛ یعنی با تبدیلی بر روی متغیرها، رابطه بین آن‌ها خطی می‌شود. مدل رگرسیون لجیت یکی از این مدل‌هاست که با تبدیل لجیت، خطی می‌شود.

رگرسیون لجستیک یکی از تکنیک‌های کاربردی برای تحلیل داده‌های رده بندی شده است که از آن برای بیان رابطه بین متغیرهای تبیینی با یک متغیر پاسخ از نوع دو سطحی استفاده می‌شود. به عنوان نمونه، اگر نتیجه آزمایشی را به صورت موفقیت و شکست معرفی کنیم، در این حالت متغیر پاسخ دیگر پیوسته نبوده، و به صورت رده بندی شده خواهد بود. در این حالت برای رده بندی مشاهدات جدید، می‌توان از مدل رگرسیون لجستیک استفاده کرد که عضوی از رده مدل‌های خطی تعمیم یافته محسوب می‌شوند.

مدل رگرسیون خطی تابعی پارامتریک است و به صورت زیر بیان می‌شود:

$$y = \beta_0 + \beta_1 x + \epsilon$$

$$E(Y|X) = \beta_0 + \beta_1 x$$

$$y = E(Y|X) + \epsilon$$

که در آن β_0 و β_1 پارامترهایی هستند که به روش کمترین مجموع مربعات خطا برآورد شده و X متغیر تبیینی مدل با مقدار مشاهده شده x و ϵ متغیر تصادفی دارای توزیع نرمال با میانگین صفر و واریانس σ^2 می‌باشد. گفته شد در رده بندی مدل‌های خطی، مدل‌هایی وجود دارند که ذاتاً خطی هستند. به عبارت دیگر، با بعضی تبدیل‌های مناسب روی متغیرهای مدل، رابطه بین آنها خطی می‌شود. مدل رگرسیون لجستیک یکی از این مدل‌هاست. از این روش برای مدل سازی رابطه بین متغیر پاسخ رسته‌ای و یا یک یا چند متغیر تبیینی (کمی یا کیفی) استفاده می‌شود. حال مجموعه داده‌های پاسخ را در نظر بگیرید که در آن پاسخ به صورت دودویی است

$(0, 1)$ ، یعنی متغیر پاسخ گسسته است. در این صورت به جای استفاده از مدل بندی رگرسیون ساده، از مدل بندی رگرسیون لجستیک استفاده می کنیم. اگر رگرسیون لجستیک احتمال وقوع موفقیت $y = 1$ را $\pi(x)$ و احتمال وقوع شکست $y = 0$ را $1 - \pi(x)$ در نظر بگیریم تابع لجیت (تبدیل لجیت) که نسبت موفقیت به شکست است به صورت زیر است:

$$\ln\left(\frac{\pi(x)}{1 - \pi(x)}\right) = \beta_0 + \beta_1 x$$

استفاده می کنیم. هدف به دست آوردن $E(Y = 1|X)$ است.

$$E(Y = 1|X = x) = \pi(x) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$$

برای یک مشاهده جدید مانند x^* ، پیشگویی حاصل از مدل برازش شده، $\hat{\pi}(x^*)$ مقداری بین صفر و یک است. اگر هدف رده بندی متغیر پاسخ متناظر با این مشاهده باشد، معمولاً با تعریف مقدار آستانه، رده مشاهده پیش گویی می شود. مثلاً اگر $\hat{\pi}(x^*)$ آنگاه رده برابر یک و در اگر برابر با $1 - \hat{\pi}(x^*)$ رده برابر با صفر است. برای انجام این مرحله، ابتدا باید پارامترهای مدل یعنی ضرایب β_0 و β_1 برآورد شوند. برای به دست آوردن این برآوردها یکی از شیوه کم ترین توان های دوم خطا است.

۱.۳.۱ میانگین تابع پاسخ زمانی در حال دودویی

فرض کنید که مدل رگرسیون خطی ساده به صورت زیر تعریف شود:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad y_i = 0, 1 \quad (1.1)$$

وقتی خروجی متغیر پاسخ Y_i دودویی است و مقادیر 0 یا 1 می گیرد. زمانی که $E[\epsilon_i] = 0$ باشد داریم:

$$E[\epsilon_i] = \beta_0 + \beta_1 x_i \quad (2.1)$$

چون که Y_i یک متغیر تصادفی برنولی است می توان توزیع احتمال را به صورت زیر نوشت

Y_i	احتمال
1	$P(Y_i = 1) = \pi_i$
0	$P(Y_i = 0) = 1 - \pi_i$

بنابراین، π_i احتمال $Y_i = 1$ و $1 - \pi_i$ احتمال $Y_i = 0$ است. با استفاده از تعریف امید ریاضی داریم.

$$E[Y_i] = 1(\pi_i) + 0(1 - \pi_i) = \pi_i = P(Y_i = 1) \quad (3.1)$$

از (۳.۱) می‌توان نتیجه گرفت

$$E[Y_i] = \beta_0 + \beta_1 X_i = \pi_i$$

زمانی که متغیر دودویی باشد و از رگرسیون ساده استفاده کنیم ۳ مشکل داریم:

۱. نرمال نبودن عناصر خطا. برای هر متغیر پاسخ دودویی ۰ و ۱، هر عنصر خطا به صورت $\epsilon_i = Y_i - (\beta_0 + \beta_1 X_i)$ لذا که می‌تواند دو مقدار زیر را بگیرد.

$$\epsilon_i = 1 - \beta_0 - \beta_1 X_i \quad : Y_i = 1 \text{ زمانی که}$$

$$\epsilon_i = -\beta_0 - \beta_1 X_i \quad : Y_i = 0 \text{ زمانی که}$$

واضح است که، مدل رگرسیون خطی ساده، زمانی که فرض می‌شود ϵ_i به‌طور نرمال توزیع شده است، مناسب نیست.

۲. واریانس خطای غیر ثابت. مشکل دیگر این است که واریانس خطا ϵ_i زمانی که متغیر دودویی داریم نمی‌تواند ثابت باشد. با توجه به فرمول (۲.۱)، $V^2[Y_i]$ برای مدل رگرسیون خطی ساده می‌بایست به روش زیر به‌دست آید.

$$V^2[Y_i] = E[(Y_i - E[Y_i])^2] = (1 - \pi_i)^2 \pi_i + (0 - \pi_i)^2 (1 - \pi_i)$$

یا

$$V^2[Y_i] = \pi_i(1 - \pi_i) = (E[Y_i])(1 - E[Y_i]) \quad (۴.۱)$$

واریانس ϵ_i مشابه واریانس Y_i است زیرا $\epsilon_i = Y_i - \pi_i$ و π_i ثابت است.

$$V^2[\epsilon_i] = \pi_i(1 - \pi_i) = E[Y_i](1 - E[Y_i]) \quad (۵.۱)$$

یا

$$V^2[\epsilon_i] = (\beta_0 + \beta_1 X_i)(1 - \beta_0 - \beta_1 X_i) \quad (۶.۱)$$

در فرمول ۶.۱، $V^2[\epsilon_i]$ وابسته به X است بنابراین واریانس خطا در سطوح مختلف X ، متفاوت خواهد بود و از مجموع مربعات معمولی قابل محاسبه نیست.

۳. محدودیت بر روی تابع پاسخ. چون تابع پاسخ به‌صورت احتمالی بیان می‌شود، زمانی که Y ها دودویی هستند در این صورت متغیر پاسخ می‌بایست به‌صورت زیر محدود شود:

$$0 \leq E[Y] = \pi \leq 1$$

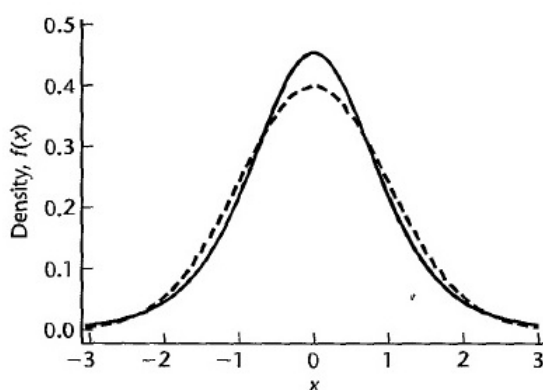
توابع پاسخ زیادی به‌طور خودکار این محدودیت را دارا نیستند. یک تابع پاسخ خطی، به‌طور مثال، ممکن است محدودیت‌های این ویژگی را در محدوده متغیر پیشگو که هدف مدل است، برطرف کند. کمترین توان دوم معمولی بهینه نمی‌شود.

۲.۳.۱ تابع پاسخ لجستیک

فرض کنید که مدل دارای خطای نرمال است، برای متغیر پاسخ اصلی در فرمول زیر

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \quad (7.1)$$

که فرض می‌شود ϵ_i توزیع نرمال با میانگین ۰ و واریانس σ^2 دارد، فرض بر این است که ϕ ، دارای تابع تجمعی نرمال استاندارد برای مدل π_i است. یک تابع خطا جایگزین، که به جای توزیع نرمال، توزیع لجستیک است. شکل (۱.۱) چگالی نرمال استاندارد و تابع چگالی لجستیک را نمایش می‌دهد، هر دو دارای میانگین ۰ و واریانس ۱ هستند. طبق شکل دُم تابع لجستیک سنگین‌تر به نظر می‌رسد.



شکل ۱.۱: نمودار چگالی نرمال (-) و چگالی لجستیک (-)

تابع لجستیک با متغیر تصادفی ϵ_L ، دارای میانگین ۰ و واریانس $\sigma = \frac{\pi}{\sqrt{3}}$ به صورت زیر است.

$$f_L(\epsilon_L) = \frac{\exp(\epsilon_L)}{[1 + \exp(\epsilon_L)]^2}$$

تابع توزیع تجمعی نیز به صورت زیر است.

$$F_L(\epsilon_L) = \frac{\exp(\epsilon_L)}{1 + \exp(\epsilon_L)}$$

فرض کنید که ϵ_i در معادله (۷.۱) توزیع لجستیک با میانگین ۰ و انحراف معیار σ باشد، داریم.

$$P(Y_i = 1) = P\left(\frac{\epsilon_i}{\sigma} \leq \beta_0^* + \beta_1^* X_i\right)$$

که $\frac{\epsilon_i}{\sigma}$ دارای توزیع لجستیک با میانگین ۰ و انحراف معیار ۱ است. با ضرب طرفین نامساوی در $\frac{\pi}{\sqrt{3}}$ داریم:

$$\begin{aligned} P(Y_i = 1) &= \pi_i = P\left(\frac{\pi}{\sqrt{3}} \frac{\epsilon_i}{\sigma} \leq \frac{\pi}{\sqrt{3}} \beta_0^* + \frac{\pi}{\sqrt{3}} \beta_1^* X_i\right) \\ &= P(\epsilon_L \leq \beta_0 + \beta_1 X_i) \\ &= F_L(\beta_0 + \beta_1 X_i) \\ &= \frac{\exp(\beta_0 + \beta_1 X_i)}{1 + \exp(\beta_0 + \beta_1 X_i)} \end{aligned}$$

زمانی که $\beta_0 = (\pi/\sqrt{3})\beta_0^*$ و $\beta_1 = (\pi/\sqrt{3})\beta_1^*$ باشد نشان می‌دهد که پارامترهای رگرسیون لجستیک است. به‌طور خلاصه، تابع توزیع میانگین لجستیک به صورت زیر تعریف می‌شود.

$$E[Y_i] = \pi_i = F_L(\beta_0 + \beta_1 X_i) = \frac{\exp(\beta_0 + \beta_1 X_i)}{1 + \exp(\beta_0 + \beta_1 X_i)}$$

در یک نمایش ساده جبری از تساوی بالا می‌توان نوشت

$$E[Y_i] = \pi_i = [1 + \exp(-\beta_0 - \beta_1 X_i)]^{-1} \quad (۸.۱)$$

معکوس تابع توزیع تجمعی F_L به صورت زیر تعریف می‌شود

$$F_L^{-1}(\pi_i) = \beta_0 + \beta_1 X_i = \pi_i' \quad (۹.۱)$$

را $F_L^{-1}(\pi_i)$ تبدیل لجیت می‌گویند.

$$F_L^{-1}(\pi_i) = \ln\left(\frac{\pi_i}{1 - \pi_i}\right) \quad (۱۰.۱)$$

به نسبت بالا نسبت بخت^۱ گفته می‌شود.

۴.۱ رگرسیون لجستیک ساده

برای برآورد پارامترهای رگرسیون لجستیک، باید از روش درست‌نمایی ماکزیمم استفاده کنیم. این روش، یک روش مناسب برای مقابله با مشکلات مربوط به متغیر پاسخ Y_i که دودویی است، می‌باشد. در ابتدا احتیاج داریم که کارکرد احتمالی مشترک نمونه‌های مشاهده شده را تعمیم دهیم. به‌جای استفاده از توزیع نرمال برای مشاهدات Y همانطور که قبلاً انجام شد، اکنون باید از توزیع برنولی برای متغیر تصادفی دودویی استفاده کنیم.

۱.۴.۱ مدل رگرسیون لجستیک ساده

ابتدا، به بیان یک فرمول برای مدل رگرسیون لجستیک ساده نیاز داریم. گفته شد که وقتی متغیر پاسخ دودویی باشد، احتمالات π و $1 - \pi$ به ترتیب مقادیر ۱ و ۰ را می‌گیرند. Y متغیر

^۱ odds ratio

تصادفی برنولی با پارامتر π $E[Y] = \pi$ است، که می‌توانیم مدل ساده رگرسیون لجستیک را به صورت زیر بیان کرد.

$$Y_i = E[Y_i | X_i] + \epsilon_i$$

از آنجایی که تابع ϵ_i وابسته به تابع برنولی پاسخ Y_i است.

$$E[Y_i | X_i] = \pi_i = \frac{\exp(\beta_0 + \beta_1 X_i)}{1 + \exp(\beta_0 + \beta_1 X_i)} \quad (11.1)$$

فرض کنید که مشاهدات x ثابت هستند. از سوی دیگر، اگر x تصادفی باشد، $E[Y_i]$ به عنوان یک معادله شرطی با توجه به مقدار x مشاهده می‌شود.

۲.۴.۱ تابع درست‌نمایی

از آنجایی که مشاهده Y_i یک متغیر تصادفی برنولی عادی است داریم

$$P(Y_i = 1) = \pi_i \quad (12.1)$$

$$P(Y_i = 0) = 1 - \pi_i \quad (13.1)$$

و یا به طور خلاصه به صورت تابع احتمال بنویسیم

$$f_i(Y_i) = \pi_i^{Y_i} (1 - \pi_i)^{1 - Y_i} \quad Y_i = 0, 1; \quad i = 1, \dots, n$$

باید توجه داشت که $f_i(1) = \pi_i$ و $f_i(0) = 1 - \pi_i$ است، از این رو $f(y_i)$ نشان دهنده احتمال $y_i = 1$ یا $y_i = 0$ است.

از آنجایی که Y_i مستقل است، تابع احتمال توام به صورت زیر نوشته می‌شود:

$$g(Y_1, \dots, Y_n) = \prod_{i=1}^n f_i(Y_i) = \prod_{i=1}^n \pi_i^{Y_i} (1 - \pi_i)^{1 - Y_i}$$

با لگاریتم گرفتن از تابع احتمال به راحتی می‌توان برآورد درست‌نمایی ماکزیمم را محاسبه کرد.

$$\begin{aligned} \ln g(Y_1, \dots, Y_n) &= \ln \prod_{i=1}^n \pi_i^{Y_i} (1 - \pi_i)^{1 - Y_i} \\ &= \sum_{i=1}^n [Y_i \ln \pi_i + (1 - Y_i) \ln(1 - \pi_i)] \quad (14.1) \\ &= \sum_{i=1}^n \left[Y_i \ln \left(\frac{\pi_i}{1 - \pi_i} \right) \right] + \sum_{i=1}^n \ln(1 - \pi_i). \end{aligned}$$

از رابطه (11.1) داریم:

$$1 - \pi_i = [1 + \exp(\beta_0 + \beta_1 X_i)]^{-1} \quad (15.1)$$

علاوه بر این، از (۱۰.۱) می‌توان نتیجه گرفت.

$$\text{Ln} \left(\frac{\pi_i}{1 - \pi_i} \right) = \beta_0 + \beta_1 X_i \quad (16.1)$$

بنابراین می‌توان با توجه به (۱۴.۱) بیان کرد

$$\text{Ln} L(\beta_0, \beta_1) = \sum_{i=1}^n Y_i (\beta_0 + \beta_1 X_i) - \sum_{i=1}^n \text{Ln} [1 + \exp(\beta_0 + \beta_1 X_i)] \quad (17.1)$$

از آن جایی که $L(\beta_0, \beta_1)$ جایگزین $g(Y_1, \dots, Y_n)$ است می‌توانیم از آن برای برآورد پارامترهای تابع درستنمایی استفاده کنیم.

۳.۴.۱ برآورد پارامتر

در رگرسیون لجستیک، پارامترها معمولاً به روش درستنمایی ماکسیمم برآورد می‌شوند. اما چون صورت بسته‌ای برای برآوردهای درستنمایی ماکسیمم این مدل وجود ندارد، از روش‌های عددی برای برآورد آن‌ها استفاده می‌شود.

همانند روش رگرسیون خطی می‌توان مدل رگرسیون لجستیک را به صورت زیر بیان کرد

$$Y = E(Y|X) + \epsilon \quad \Rightarrow \quad Y = \pi(x) + \epsilon$$

با توجه به توزیع $Y|X$ ، از روش درستنمایی ماکزیمم برای برآورد پارامترها استفاده می‌کنیم. از آن جایی که $Y|X = x$ دارای توزیع دوجمله‌ای $Y|X = x \sim B(1, \pi(x))$ لذا داریم:

$$P(Y = y_i) = \pi(x_i)^{y_i} (1 - \pi(x_i))^{1-y_i}$$

که در آن:

$$\begin{aligned} \pi(x_i) &= E(Y_i | X_i = x_i) = P(Y_i = 1, X_i = x_i) \\ &= \frac{1}{1 + \exp(\beta_0 + \sum_{i=1}^n \beta_1 x_i)} \end{aligned} \quad (18.1)$$

در این صورت تابع درستنمایی براساس یک نمونه n تایی مستقل عبارت است از:

$$L(\beta|x) = \prod_{i=1}^n P(Y = y_i)$$

که در آن $\beta = (\beta_0, \beta_1)^T$ می‌باشد. با لگاریتم گرفتن از طرفین داریم:

$$\ln L(\beta|x) = \sum_{i=1}^n y_i \{ \ln(\pi(x_i) + (1 - y_i)) \}$$

با جایگذاری رابطه (۱۸.۱) داریم:

$$\begin{aligned}
 \ln L &= \sum_{i=1}^n \left\{ y_i \ln \left(\frac{1}{1 + \exp(\beta_0 + \sum_{i=1}^n \beta_i x_i)} \right) \right. \\
 &\quad \left. + (1 - y_i) \ln \left(\frac{\exp(\beta_0 + \sum_{i=1}^n \beta_i x_i)}{1 + \exp(\beta_0 + \sum_{i=1}^n \beta_i x_i)} \right) \right\} \\
 &= \sum_{i=1}^n \left\{ \ln \left(\frac{\exp(\beta_0 + \sum_{i=1}^n \beta_i x_i)}{1 + \exp(\beta_0 + \sum_{i=1}^n \beta_i x_i)} \right) \right. \\
 &\quad \left. - y_i \ln(\exp(\beta_0 + \sum_{i=1}^n \beta_i x_i)) \right\} \\
 &= \sum_{i=1}^n \left\{ \ln(\exp(\beta_0 + \sum_{i=1}^n \beta_i x_i)) \right. \\
 &\quad \left. - y_i \ln(\exp(\beta_0 + \sum_{i=1}^n \beta_i x_i)) \right. \\
 &\quad \left. - \ln(1 + \exp(\beta_0 + \sum_{i=1}^n \beta_i x_i)) \right\} \\
 &= \sum_{i=1}^n \left\{ (\beta_0 + \sum_{i=1}^n \beta_i x_i) \right. \\
 &\quad \left. - y_i (\beta_0 + \sum_{i=1}^n \beta_i x_i) - \ln(1 + \exp(\beta_0 + \sum_{i=1}^n \beta_i x_i)) \right\}
 \end{aligned}$$

حال با مشتق‌گیری از رابطه بالا نسبت به پارامترها داریم:

$$\frac{d \ln L}{d \beta_0} = \sum (y_i - \pi(x_i))$$

$$\frac{d \ln L}{d \beta_1} = \sum x_i (y_i - \pi(x_i))$$

باید به این نکته توجه داشت که در مدل رگرسیون لجستیک ارتباط Y و X غیر خطی است و حل معادلات فوق مانند رگرسیون خطی نمی‌باشد. لذا باید برای حل معادلات فوق از روش‌های حل عددی استفاده کرد. پس از برآورد پارامترهای مورد نظر، حال می‌توان برای اختصاص دادن هر مشاهده به یکی از مقادیر ۱ یا ۰، بیشترین احتمال $P(Y_i = y_i | X_i = x_i)$ را در نظر گرفت که در آن‌ها دارای مقادیر ۰ یا ۱ می‌باشند. بنابراین Y را به صفر نسبت می‌دهیم هرگاه:

$$P(Y = 0 | x) > P(Y = 1 | x)$$

$$\frac{P(Y = 0 | x)}{P(Y = 1 | x)} > 1$$

در نتیجه خواهیم داشت:

$$\exp(\beta_0 + \sum_{i=1}^n \beta_i x_i) > 1$$

که اگر از طرفین نامساوی فوق لگاریتم بگیریم خواهد شد:

$$\beta_0 + \sum_{i=1}^n \beta_i x_i > 0$$

$$\ln \frac{1-\pi}{\pi} = \beta_0 + \sum_{i=1}^n \beta_i x_i > 0$$

بنابراین قاعده تصمیم برای پیش‌بینی کلاس مورد نظر براساس روش رگرسیون لجستیک به صورت زیر می‌باشد:

$$\hat{\beta}_0 + \sum \hat{\beta}_i x_i < 0 \quad \text{اگر } Y = 1 \text{ آنگاه}$$

$$\hat{\beta}_0 + \sum \hat{\beta}_i x_i > 0 \quad \text{اگر } Y = 0 \text{ آنگاه}$$

همان‌طور که گفته شد روش‌های عددی برای حل این مسائل و به دست آوردن مقادیر b_0 و b_1 مناسب است. برای محاسبه، چندین بار روش‌های عددی را تکرار می‌کنیم، یکی از روش‌ها به طور منظم وزن دادن به کم‌ترین مربعات است. با تکیه بر برنامه‌های آماری معتبر که به طور مشخص برای رگرسیون لجستیک طراحی شده‌اند می‌توان مقادیر b_0 و b_1 را برآورد کرد. همان‌طور که گفته شد روش‌های عددی برای حل این مسائل و به دست آوردن مقادیر b_0 و b_1 مناسب است. برای محاسبه، چندین بار روش‌های عددی را تکرار می‌کنیم، یکی از روش‌ها به طور منظم وزن دادن به کمترین توان‌های دوم است. با تکیه بر برنامه‌های آماری معتبر که به طور مشخص برای رگرسیون لجستیک طراحی شده‌اند می‌توان مقادیر b_0 و b_1 را برآورد کرد. یک راه دیگر برای به دست آوردن برآورد درست‌نمایی ماکزیمم، باید از π_1 برای نشان دادن مقدار برازش شده در i امین مقدار استفاده کنیم.

$$\hat{\pi}_i = \frac{\exp(b_0 + 1X_i)}{1 + \exp(b_0 + 1X_i)} \quad (19.1)$$

تابع پاسخ لجستیک برازش شده به صورت زیر است.

$$\hat{\pi} = \frac{\exp(b_0 + 1X)}{1 + \exp(b_0 + 1X)} \quad (20.1)$$

اگر از تعریف رگرسیون لجستیک در معادله (۹.۱) استفاده کنیم می‌توانیم با توجه به (۲۰.۱) تابع پاسخ برازش شده را به دست آوریم.

$$\hat{\pi}' = b_0 + b_1 X \quad (21.1)$$

در حالی که

$$\hat{\pi}' = \text{Ln} \left(\frac{\hat{\pi}}{1 - \hat{\pi}} \right) \quad (22.1)$$

معادله (۲۱.۱) را تابع پاسخ لجیت برازش شده می‌نامند. یک تابع پاسخ لجستیک برازش شده به دست می‌آید که معمولاً از گام‌های بعدی و امتحان کردن مناسب بودن تابع پاسخ برازش شده خوب باشد، می‌توان نتیجه‌گیری‌ها و پیشگویی‌های متنوعی داشته باشیم.

۴.۴.۱ تعبیر ضریب b_1

تعبیر ضریب b_1 در برآورد رگرسیونی، در تابع پاسخ لجستیک برازش شده در معادله (۲۳.۱)، تعبیری مستقیم از شیب در مدل رگرسیون خطی نیست. دلیلش این است که تاثیر افزایش یک واحد در X ‌های مختلف برای مدل رگرسیون لجستیک با توجه به موقعیت شروع نقطه X است. یک تعبیر برای b_1 را می‌توان این‌گونه گفت، به ازای هر واحد افزایش X ، تابع لجستیک برازش شده که با برآورد تابع بخت، $\hat{\pi}/(1 - \hat{\pi})$ شده است در $\exp(b_1)$ ضرب می‌شود. فرض کنید که مقدار تابع لجیت برازش شده (۲۱.۱) اگر $X = X_j$:

$$\hat{\pi}'(X_j) = b_0 + b_1 X_j \quad (23.1)$$

نمادگذاری $\hat{\pi}(X_j)$ نشان می‌دهد که به‌طور مشخص سطح X مرتبط با مقدار برازش شده است. به هر حال می‌توانیم فرض کنیم که مقدار تابع پاسخ لجیت برازش شده برای $X = X_j + 1$ به صورت زیر تعریف می‌شود:

$$\hat{\pi}'(X_j + 1) = b_0 + b_1(X_j + 1) \quad (24.1)$$

تفاضل بین دو مقدار برازش شده به صورت زیر است:

$$\hat{\pi}'(X_j + 1) - \hat{\pi}'(X_j) = b_1 \quad (25.1)$$

حال با توجه به معادله (۲۲.۱)، $\hat{\pi}(X_j)$ ، زمانی که $X = X_j$ باشد، لگاریتم برآورد $odds$ است که با $\text{Ln}(odds_1)$ مشخص می‌شود. به صورت مشابه، $\hat{\pi}(X_j + 1)$ ، زمانی که $X = X_j + 1$ باشد لگاریتم برآورد $odds$ است که با $\text{Ln}(odds_2)$ نشان داده می‌شود. بنابراین داریم:

$$\ln(odds_2) - \ln(odds_1) = \ln\left(\frac{odds_2}{odds_1}\right) = b_1 \quad (26.1)$$

با گرفتن آنتی لگ از طرفین، می توان مقدار نسبت $odds$ را به دست آورد که آن را نسبت بخت می نامند و با \widehat{OR} نشان می دهند.

$$\widehat{OR} = \frac{odds_2}{odds_1} = \exp(b_1) \quad (27.1)$$

۵.۱ رگرسیون لجستیک چندگانه

مدل رگرسیون لجستیک ساده که قبلاً معرفی شد، یک بیان ساده برای پیش گویی متغیر است. در حقیقت، چندین متغیر پیش گو معمولاً به رگرسیون لجستیک برای به دست آوردن توصیف کافی و پیش بینی های مفید، نیاز دارند.

با تعمیم دادن مدل رگرسیون لجستیک ساده، به سادگی می توانیم $\beta_0 + \beta_1 X$ در (۸.۱) را با $\beta_0 + \beta_1 X_1 + \dots + \beta_{p-1} X_{p-1}$ جایگزینی کنیم. با ساده سازی فرمول ها، می توان آن را به صورت ماتریس درآورد. بردارهای زیر را در نظر بگیرید:

$$\beta_{p \times 1} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix} \quad \mathbf{X}_{p \times 1} = \begin{bmatrix} 1 \\ X_1 \\ X_2 \\ \vdots \\ X_{p-1} \end{bmatrix} \quad \mathbf{X}_i_{p \times 1} = \begin{bmatrix} 1 \\ X_{i1} \\ X_{i2} \\ \vdots \\ X_{i,p-1} \end{bmatrix} \quad (28.1)$$

در این صورت داریم

$$X^T \beta = \beta_0 + \beta_1 X_1 + \dots + \beta_{p-1} X_{p-1} \quad (29.1)$$

$$X_i^T \beta = \beta_0 + \beta_1 X_{i1} + \dots + \beta_{p-1} X_{i,p-1}$$

با تعمیم رگرسیون لجستیک ساده به رگرسیون لجستیک چندگانه داریم:

$$E[Y] = \frac{\exp(X^T \beta)}{1 + \exp(X^T \beta)} \quad (30.1)$$

و معادله پاسخ لجستیک (۹.۱) به صورت زیر تعمیم داده می شود

$$E[Y] = \left[1 + \exp(-X^T \beta)\right]^{-1} \quad (31.1)$$

به صورت مشابه، تبدیل لجیت (۱۰.۱) به صورت

$$\pi' = \ln\left(\frac{\pi}{1-\pi}\right)$$

است. بنابراین منجر به تابع پاسخ لجیت یا پیشگوی خطی می شود.

$$\pi' = X^T \beta \quad (۳۲.۱)$$

می توان مدل رگرسیون لجستیک چندگانه را به صورت زیر بیان کرد:
 Y_i ها متغیرهای تصادفی برنولی تصادفی با مقدار مورد انتظار $E[Y_i] = \pi_i$ هستند که

$$E[Y_i] = \pi_i = \frac{\exp(X_i^T \beta)}{1 + \exp(X_i^T \beta)} \quad (۳۳.۱)$$

بنابراین، مشاهدات X ثابت فرض می شوند. با توجه به این که، متغیر X تصادفی است $E[Y_i]$ یک میانگین که از مقادیر $X_{i,1}, \dots, X_{i,p-1}$ گرفته شده است.

مانند رگرسیون لجستیک ساده در (۸.۱)، رگرسیون لجستیک چندگانه در (۳۳.۱) یکنواخت و S شکل است. با توجه به این که $X\beta$ تقریباً خطی است زمانی که π بین $0/2$ و $0/8$ باشد. متغیر X ممکن است متفاوت از متغیر پیش گو باشد یا بعضاً ممکن است منحنی و یا اثر متقابل را نشان دهد. به هر حال، متغیرهای پیش گو ممکن است کمی یا کیفی باشند.

۱.۵.۱ برآورد پارامترهای لجستیک چندگانه

اکنون می خواهیم از روش درست‌نمایی ماکزیمم برای برآورد پارامترها در تابع لجستیک چندگانه (۳۳.۱) استفاده کنیم. تابع \log -درست‌نمایی برای رگرسیون لجستیک ساده در رابطه (۱۷.۱) را می توان برای رگرسیون لجستیک چندگانه تعمیم داد.

$$\text{Ln}L(\beta) = \sum_{i=1}^n Y_i (X_i^T \beta) - \sum_{i=1}^n \text{Ln}[1 + \exp(X_i^T \beta)] \quad (۳۴.۱)$$

در روش‌های عددی به جای مقدار $\beta_0, \dots, \beta_{p-1}$ از $\text{Ln}L(\beta)$ استفاده شده است و به روش درست‌نمایی ماکزیمم b_0, b_1, \dots, b_{p-1} را برآورد می کند. فرض کنید b برداری از برآوردهای درست‌نمایی ماکزیمم باشد.

$$\mathbf{b}_{p \times 1} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_{p-1} \end{bmatrix} \quad (۳۵.۱)$$

برازش تابع پاسخ لجستیک و برازش مقادیر را می توان به صورت زیر بیان کرد.

$$\hat{\pi} = \frac{\exp(\mathbf{X}^T \mathbf{b})}{1 + \exp(\mathbf{X}^T \mathbf{b})} = |1 + \exp(-\mathbf{X}^T \mathbf{b})|^{-1} \quad (۳۶.۱)$$

$$\hat{\pi}_i = \frac{\exp(\mathbf{X}_i^\top \mathbf{b})}{1 + \exp(\mathbf{X}_i^\top \mathbf{b})} = |1 + \exp(-\mathbf{X}_i^\top \mathbf{b})|^{-1} \quad (۳۷.۱)$$

زمانی که

$$\mathbf{X}^\top \mathbf{b} = b_0 + b_1 X_1 + \cdots + b_{p-1} X_{p-1} \quad (۳۸.۱)$$

$$\mathbf{X}_i^\top \mathbf{b} = b_0 + b_1 X_{i1} + \cdots + b_{p-1} X_{i,p-1} \quad (۳۹.۱)$$

فصل ۲

مجموعه‌ها و رگرسیون فازی

در این فصل تعاریف و قضایای نظریه مجموعه‌های فازی و رگرسیون فازی به منظور آشنایی خواننده به اختصار آورده شده است.

۱.۲ مجموعه‌های فازی

نظریه مجموعه فازی در سال ۱۹۶۵ توسط پروفیسور لطفی عسکرزاده^۱ دانشمند ایرانی تبار و استاد دانشگاه برکلی آمریکا، ارائه شده است. این نظریه از زمان ارائه آن تا کنون، گسترش و تعمیق زیادی یافته و کاربردهای گوناگونی در زمینه‌های مختلف پیدا کرده است.

در نظریه مجموعه‌های قطعی، یک مجموعه معمولی با یک ویژگی دقیق و معین مشخص می‌شود. به‌دین ترتیب که اگر عضو مجموعه مرجع، آن خصوصیت را داشته باشد، عضو مجموعه و اگر فاقد آن خصوصیت باشد، عضو مجموعه نیست. مثلاً، اول بودن اعداد حقیقی یک خصوصیت خوش تعریف است یعنی برای هر عدد حقیقی می‌توان گفت که یا عدد اول است یا عدد اول نیست. با فرض اینکه X یک مجموعه مرجع باشد، برای هر زیر مجموعه معمولی از A از مجموعه مرجع، که با یک خصوصیت خوش تعریف مشخص می‌شود، می‌توان یک تابع

^۱L. A. Zadeh

نشانگر تعریف کرد که با هر عضو متعلق به A ، مقدار صفر را اختیار کند، یعنی

$$I_A(x) = \begin{cases} 1 & x \in A \\ 0 & x \notin A \end{cases}$$

اما اکثر مفاهیمی که با آنها برخورد داریم و براساس آن زندگی روزمره خود را پیش می‌بریم، مفاهیمی دقیق نیستند، مانند هوای خیلی سرد، ویژگی زیبا بودن، قد خیلی بلند، مجموعه اعداد بزرگ و برای این خصوصیات مرز مشخص و دقیقی را نمی‌توان تعیین کرد. برای مثال از نظر یک فرد، هوای ۵ درجه سانتی‌گراد و از نظر فرد دیگری هوای ۲ درجه سانتی‌گراد، هوای خیلی سرد است. نظریه مجموعه‌های قطعی از عهده صورت‌بندی این مفاهیم و خصوصیات برنمی‌آید.

نظریه مجموعه‌های فازی، نظریه‌ای برای صورت‌بندی و تجزیه و تحلیل این گونه مفاهیم است. این نظریه تعمیم طبیعی از نظریه مجموعه‌ای قطعی است. به‌طور مثال، فرض کنید می‌خواهیم مجموعه اعداد بزرگ را تعریف کنیم. طبیعی است که «بزرگ بودن» یک خصوصیت خوش تعریفی نیست. بنا به پیشنهاد زاده، اجازه می‌دهیم که هر x در مجموعه اعداد حقیقی، به اندازه یک مقدار بین صفر و یک، عضو مجموعه اعداد بزرگ باشد، به طوری که هر چه آن x بزرگ‌تر باشد، عدد مربوط به عضویت آن در مجموعه اعداد بزرگ به یک نزدیک‌تر، و هر چه x کوچک‌تر باشد، عدد مربوط به عضویت آن در مجموعه اعداد بزرگ به صفر نزدیک‌تر است. در این صورت می‌گوییم این عدد به اندازه مثلاً $0/8$ عضو مجموعه اعداد بزرگ است. اساس کار نظریه مجموعه‌های فازی، گسترش مفهوم تابع نشانگر یک مجموعه، با برد $[0, 1]$ به جای $\{0, 1\}$ است.

تعریف ۱.۱.۲. یک زیر مجموعه فازی A (از این پس کوتاه‌ی: مجموعه فازی) از مجموعه مرجع X ، توسط یک تابع عضویت $[0, 1] : X \rightarrow \mu_A(x)$ مشخص می‌شود که در آن برای هر $x \in X$ ، مقدار $\mu_A(x)$ میزان عضویت x در مجموعه فازی A است.

توجه کنید که اگر x کاملاً در A عضو باشد، $\mu_A(x) = 1$ و اگر اصلاً در A عضو نباشد، $\mu_A(x) = 0$ است. لذا یک مجموعه معمولی و تابع نشانگر آن حالت خاصی از مجموعه فازی و تابع عضویت آن است.

برای نشان دادن یک مجموعه فازی، روش‌های مختلفی رایج است. یک روش متداول، به کار بردن مستقیم تابع عضویت مجموعه فازی است، که از آن وقتی که مجموعه مرجع پیوسته باشد، استفاده می‌شود.

مثال ۱.۱.۲. فرض کنید $X = \{1, \dots, 7\}$. می‌خواهیم یک مجموعه فازی تعریف کنیم که اعضای آن، ویژگی نادقیق «کوچک» را داشته باشند. برای مدل‌سازی این مجموعه کافی است تابع عضویت این مجموعه فازی را مشخص کنیم. تعیین این تابع بستگی به نظر تصمیم‌گیرنده

دارد. مثلاً یک تابع می‌تواند به صورت زیر تعریف شود

$$A(x) = \begin{cases} 1 & x = 1 \\ 0.9 & x = 2 \\ 0.7 & x = 3 \\ 0.5 & x = 4 \\ 0.3 & x = 5 \\ 0.1 & x = 6 \\ 0 & x = 7 \end{cases}$$

برای مثال $A(3) = 0.7$ بدین معنی است که از نظر تصمیم‌گیر، عدد ۳ به اندازه ۰/۷ به مجموعه فازی اعداد کوچک تعلق دارد. به سخن دیگر، از نظر تصمیم‌گیر گزاره « ۳ عددی کوچک است » به اندازه ۰/۷ درست است.

تعریف تابع عضویت یک مجموعه فازی بستگی به نظر فرد تصمیم‌گیرنده دارد. در واقع این تعریف جنبه ذهنی و شخصی دارد. لذا می‌توان توابع عضویت مختلفی را برای یک مجموعه فازی که بیانگر یک ویژگی فازی باشد، تصور کرد.

تعریف ۲.۱.۲. فرض کنید X یک مجموعه مرجع و A یک زیر مجموعه فازی از آن باشد،

$$\sup(A) = \{x \in X | A(x) > 0\}.$$

تعریف ۳.۱.۲. ارتفاع یک مجموعه فازی نامیده می‌شود. اگر این عدد برابر یک باشد آن‌گاه مجموعه فازی A را نرمال و در غیر این صورت آن را زیرنرمال می‌نامیم. البته هر مجموعه فازی زیرنرمال را می‌توان با تقسیم درجات عضویت آن بر ارتفاع A نرمال کرد.

مثال ۲.۱.۲. در مثال قبل، $\sup(A) = \{1, 2, \dots, 7\}$. همچنین مجموعه فازی A مجموعه فازی نرمال است.

تعریف ۴.۱.۲. زیر مجموعه معمولی از عناصر X را که درجه عضویت آنها در مجموعه A ، حداقل به بزرگی α باشد، $-\alpha$ برش A (مجموعه تراز α ام وابسته به A) گوییم و با A_α نشان می‌دهیم.

$$A_\alpha = \{x \in X | A(x) \geq \alpha\}, \quad 0 < \alpha \leq 1$$

مثال ۳.۱.۲. فرض کنید $X = \{0, 1, 2, \dots\}$ مجموعه مقادیر یک متغیر تصادفی پواسن مربوط به تعداد تصادفات شبانه‌روزی در یک ناحیه باشد و فرض کنید مجموعه فازی A بیانگر «تعداد

تصادفات کم» با تابع عضویت زیر تعریف شود.

$$A(x) = \begin{cases} 1 & x = 0 \\ 0/8 & x = 1 \\ 0/6 & x = 2 \\ 0/4 & x = 3 \\ 0/2 & x = 4 \end{cases}$$

در این صورت چند α -برش A عبارتند از

$$\begin{aligned} A_{0/8} &= \{0, 1, 2, 3, 4\} & A_{0/5} &= \{0, 1, 2\} \\ A_{0/6} &= \{0, 1, 2\} & A_{0/85} &= \{0\} \end{aligned}$$

به‌طور خلاصه می‌توان نوشت

$$A_\alpha = \begin{cases} \{0, 1, 2, 3, 4\} & 0 < \alpha \leq 0/2 \\ \{0, 1, 2, 3\} & 0/2 < \alpha \leq 0/4 \\ \{0, 1, 2\} & 0/4 < \alpha \leq 0/6 \\ \{0, 1\} & 0/6 < \alpha \leq 0/8 \\ \{0\} & 0/8 < \alpha \leq 1 \end{cases}$$

قضیه ۱.۱.۲. (اتحاد تجزیه ^۲ [۳۶]) فرض کنید A مجموعه‌ای فازی از \mathbb{R} با تابع عضویت $A(\cdot)$ باشد. در این صورت

$$A(x) = \sup_{\alpha \in [0,1]} \alpha I_{A_\alpha}(x)$$

که در آن $I_{A_\alpha}(\cdot)$ تابع نشانگر α -برش A است.

تعریف ۵.۱.۲. (اصل توسیع یا گسترش ^۳ [۵]) فرض کنید X_1, \dots, X_n مجموعه مرجع و همچنین A_1, \dots, A_n مجموعه فازی به ترتیب از X_1, \dots, X_n باشد. به علاوه $y = f(x_1, \dots, x_n)$ یک نگاشت از X به Y باشد. حاصل عمل f بر n مجموعه فازی A_1, \dots, A_n به صورت مجموعه فازی B از Y با تابع عضویت زیر تعریف می‌شود.

$$B(y) = f(A_1, \dots, A_n)(y) = \begin{cases} \sup_{y=f(x_1, \dots, x_n)} \min A_1(x_1), \dots, A_n(x_n) & f^{-1}(y) \neq \emptyset \\ 0 & f^{-1}(y) = \emptyset \end{cases}$$

^۲Resolution Identity

^۳Extension Principle

مثال ۴.۱.۲. فرض کنید X_1 و X_2 مجموعه اعداد حسابی، A_1 مجموعه فازی اعداد خیلی کوچک و A_2 مجموعه فازی اعداد تقریباً ۲ با توابع عضویت زیر باشند

$$A_1 = \left\{ \frac{1}{0}, \frac{0.8}{1}, \frac{0.6}{2}, \frac{0.4}{3}, \frac{0.2}{4} \right\}, \quad A_2 = \left\{ \frac{0.5}{1}, \frac{1}{2}, \frac{0.5}{3} \right\}$$

آنگاه براساس اصل توسیع، حاصل عمل $f(x_1, x_2) = x_1 + x_2$ بر A_1 و A_2 به صورت مجموعه فازی زیر به دست می آید.

$$\left\{ \frac{0.5}{1}, \frac{1}{2}, \frac{0.8}{3}, \frac{0.6}{4}, \frac{0.5}{5}, \frac{0.4}{6}, \frac{0.2}{7} \right\}$$

۱.۱.۲ معرفی تابع fapply

برای اعمال اصل گسترش در نرم افزارهای مختلف الگوریتم‌های مختلفی وجود دارد. این امر در نرم افزار R در بسته FuzzyNumber که در پیوست آ به طور مختصر معرفی شده است، توسط تابع fapply() میسر شده است. شایان ذکر است که عدد فازی مذکور باید از نوع قطعه ای باشد. با ذکر مثال زیر، روش کار این تابع را نشان می دهیم. [۱۴]

مثال ۵.۱.۲. عدد فازی مثلثی نامتقارن $A = (0.1, 0.2, 0.18)$ را در نظر بگیرید. برای اعمال تابع $y = \ln \frac{x}{1-x}$ بر این عدد مثلثی از دستور fapply() استفاده می کنیم (به پیوست آ مراجعه شود). عدد مثلثی متقارن حاصل از این تبدیل $A' = (-4/59512, -1/51635, -3/89182)$ می باشد.

۲.۲ اعداد فازی

اعداد فازی، که زیر مجموعه‌های خاصی از مجموعه اعداد حقیقی هستند، در بیشتر مسائل کاربردی می شوند. در ادامه، مفاهیم و نتایج اصلی درباره اعداد فازی براساس مرجع [۶] ارائه می شود.

تعریف ۱.۲.۲. مجموعه فازی N از \mathbb{R} (اعداد حقیقی) را یک عدد فازی (حقیقی) گوئیم، اگر

۱. N نرمال و تک نمایی باشد.

۲. α -برش‌های N ، به ازای هر $\alpha \in (0, 1]$ ، به صورت بازه‌های بسته باشند.

منظور از تک نمایی بودن این است که تنها یک x وجود دارد به طوری که $A(x) = 1$. مجموعه همه اعداد فازی را با $\mathcal{F}(\mathbb{R})$ نشان می دهیم.

مثال ۱.۲.۲. مجموعه فازی L با تابع عضویت زیر یک عدد فازی است. می توان L را یک مدل سازی برای اعداد فازی تقریباً صفر تعبیر کرد.

$$L(x) = 1 - |x| \quad -1 \leq x \leq 1$$

تعریف ۲.۲.۲. N یک عدد دقیق^۴ (غیر فازی) با مقدار m نامیده می‌شود هرگاه تابع عضویت آن به صورت زیر باشد

$$N(x) = \begin{cases} 1 & x = m \\ 0 & x \neq m \end{cases}$$

مجموعه همه اعداد فازی از \mathbb{R} را با $\mathcal{F}(\mathbb{R})$ نشان می‌دهیم.

تعریف ۳.۲.۲. عدد فازی N را مثبت (منفی) گوییم و آن را با نماد $N > 0$ ($N < 0$) نشان می‌دهیم، اگر برای هر $x \leq 0$ ($x \geq 0$)، $N(x) = 0$ ، عدد فازی N را نامنفی (نامثبت) گوییم اگر $x < 0$ ($x > 0$)، $N(x) = 0$.

نوع خاصی از اعداد فازی، اعداد فازی LR هستند که علاوه بر اینکه ساختار ویژه‌ای دارند، برخی اعمال حسابی بر آنها از قواعد خاصی پیروی می‌کند. این ویژگی‌ها باعث شده است که در کاربردها، عمدتاً از این نوع اعداد فازی استفاده می‌شود.

تعریف ۴.۲.۲. (عدد فازی LR [۱۷]) اگر ساختار تابع عضویت عدد فازی N به صورت زیر باشد

$$N(x) = \begin{cases} L\left(\frac{m-x}{\alpha}\right) & x \leq m \\ R\left(\frac{m-x}{\beta}\right) & x > m \end{cases}$$

که در آن توابع شکل L و R ، توابعی غیر صعودی از \mathbb{R}^+ به $[0, 1]$ هستند و $L(0) = R(0) = 1$ ، آنگاه N را یک عدد فازی LR نامیده و با نماد $N = (m, \alpha, \beta)_{LR}$ نشان می‌دهیم. عدد حقیقی m را مقدار نما (یا مرکز یا میانه) و اعداد مثبت α و β را به ترتیب پهنا چپ و پهنای راست N می‌نامیم.

در صورتی که $L = R$ و $\alpha = \beta$ ، آنگاه عدد فازی را متقارن نامیده و با $N = (m, \alpha)_{LL}$ نشان می‌دهیم.

برای $\alpha = 0$ قرار می‌دهیم $L\left(\frac{m-x}{\alpha}\right)$ و برای $\beta = 0$ قرار می‌دهیم $R\left(\frac{m-x}{\beta}\right)$ ، لذا منظور از $X = (m, 0, 0)_{LR}$ ، یک عدد دقیق با مقدار m است.

اعداد فازی LR متداول در تعریف زیر معرفی شده‌اند.

تعریف ۵.۲.۲. فرض کنید $N = (m, \alpha, \beta)_{LR}$ و $L = R$

۱. اگر $L(x) = \max\{0, 1 - |x|\}$ ، آنگاه N را یک عدد فازی مثلثی^۵ نامیده و با $N = (m, \alpha, \beta)_T$ نشان می‌دهیم.

۲. اگر $L(x) = e^{-x^2}$ ، آنگاه N را یک عدد فازی نرمال^۶ نامیده و با $N = (m, \alpha, \beta)_N$ نشان می‌دهیم.

^۴ crisp

^۵ Triangular

^۶ Normal

۳. اگر $L(x) = \max\{0, 1 - x\}$ ، آنگاه N را یک عدد فازی سهموی ^۷ نامیده و با $N = (m, \alpha, \beta)_P$ نشان می‌دهیم.

تعریف ۶.۲.۲. (بازه فازی LR) اگر در تعریف عدد فازی شرط تک‌نمایی بودن برداشته شود، آنگاه یک بازه فازی با نماد $N = (m_1, m_2, \alpha, \beta)_{LR}$ و با تابع عضویت زیر خواهیم داشت

$$N(x) = \begin{cases} L\left(\frac{m_1 - x}{\alpha}\right) & x \leq m_1 \\ 1 & m_1 \leq x < m_2 \\ R\left(\frac{x - m_2}{\beta}\right) & x \geq m_2 \end{cases}$$

که در آن اعداد حقیقی m_1 و m_2 به ترتیب نمای اول و نمای دوم و اعداد مثبت α و β به ترتیب پهنای چپ و پهنای راست N می‌باشند.

اگر $L(x) = R(x) = \max\{0, 1 - |x|\}$ آنگاه بازه فازی را عدد فازی دوزنقه‌ای ^۸ نیز می‌نامند و با $N = (m_1, m_2, \alpha, \beta)_{T\alpha}$ نشان می‌دهند.

بازه فازی می‌تواند به عنوان تعمیمی از عدد فازی تلقی شود. در واقع زمانی که $m_1 = m_2 = m$ آنگاه بازه فازی به عدد فازی تبدیل می‌شود.

۳.۲ حساب اعداد فازی

اکنون چند تعریف و قضیه را در مورد حساب اعداد فازی بیان می‌کنیم. نخست مفهوم T -نرم را یادآور می‌شویم.

تعریف ۱.۳.۲. [۵] تابع دو متغیره $I = [0, 1]$ را یک T -نرم گویند هرگاه ویژگی‌های زیر را دارا باشد.

۱. $\forall x \in [0, 1], T(x, 1) = x$

۲. یکنوایی: $x_1 \leq x_2, y_1 \leq y_2 \rightarrow T(x_1, y_1) \leq T(x_2, y_2)$

۳. جابه‌جایی: $T(x, y) = T(y, x)$

۴. شرکت‌پذیری: $T(x, T(y, z)) = T(T(x, y), z)$

نمونه‌هایی از T -نرم‌ها در زیر آورده شده است.

$$T_m(a, b) = \min(a, b)$$

$$T_p(a, b) = ab$$

$$T_L(a, b) = \max(0, a + b - 1)$$

^۷Parabolic

^۸Trapezoidal

گزاره ۱.۳.۲. [۵] برای هر T -نرم دلخواه داریم

$$T_w(x, y) \leq T(x, y) \leq T_m(x, y)$$

که $T_m(x, y) = \min(x, y)$ همچنین

$$T_w(x, y) = \begin{cases} x & y = 1 \\ y & x = 1 \\ \circ & o.w. \end{cases}$$

ضعیف‌ترین T -نرم است که به نام T -نرم دراستیک نیز نامیده می‌شود.

تعریف ۲.۳.۲. [۱۷] فرض کنید M و N دو عدد فازی و $\mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R} : *$ یک عملگر دوتایی بر اعداد حقیقی باشند. حاصل عمل تعمیم یافته \otimes بر M و N براساس اصل توسیع، به صورت یک مجموعه فازی از \mathbb{R} با تابع عضویت زیر تعریف می‌شود

$$(M \otimes N)(z) = \sup_{\substack{x, y \\ x * y = z}} T(M(x), N(y))$$

که در آن $T(\cdot, \cdot)$ یک T -نرم است.

در حالت خاص برای چهار عمل اصلی و با استفاده از T -نرم مینیمم در تعریف بالا داریم

$$(M \oplus N)(z) = \sup_{z=x+y} \min[M(x), N(y)]$$

$$(M \ominus N)(z) = \sup_{z=x-y} \min[M(x), N(y)]$$

$$(M \otimes N)(z) = \sup_{z=x*y} \min[M(x), N(y)]$$

$$(M \oslash N)(z) = \sup_{z=x/y} \min[M(x), N(y)]$$

به شرط اینکه در رابطه آخر $y \neq \circ$.

از این پس، مگر در مواردی که گفته شود، همواره از عملگر مینیمم برای حساب اعداد فازی استفاده می‌کنیم.

قضیه ۱.۳.۲. اگر $M = (m, \alpha, \beta)_{LR}$ و $\lambda \in \mathbb{R}$ آنگاه

$$\lambda M = M \lambda = \begin{cases} (\lambda m, \lambda \alpha, \lambda \beta)_{LR}, & \lambda > \circ \\ (\lambda m, -\lambda \beta, -\lambda \alpha)_{LR}, & \lambda < \circ \end{cases}$$

نتیجه ۱.۳.۲. اگر $M = (m, \alpha, \beta)_{LR}$ ، آنگاه قرینه M به صورت $(-M) = (-m, \beta, \alpha)_{LR}$ است.

قضیه ۲.۳.۲. اگر $M = (m, \alpha, \beta)_{LR}$ و $N = (n, \delta, \gamma)_{LR}$ ، آنگاه $M \oplus N$ یک عدد فازی LR به صورت زیر است.

$$M \oplus N = (m + n, \alpha + \delta, \beta + \gamma)_{LR}$$

نتیجه ۲.۳.۲. اگر $M = (m, \alpha, \beta)_{LR}$ و $N = (n, \gamma, \delta)_{RL}$ ، آنگاه $M \ominus N$ یک عدد فازی LR به صورت زیر است

$$M \ominus N = (m - n, \alpha + \delta, \beta + \gamma)_{LR}$$

درباره ضرب و تقسیم، رابطه دقیق مطابق بر آنچه که برای جمع و تفاضل گفته شده، وجود ندارد. اصولاً حاصل ضرب دو عدد فازی LR، یک عدد فازی LR نخواهد بود. همچنین در حالت‌هایی که حاصل تقسیم دو عدد فازی LR یک عدد فازی می‌شود، عدد حاصل از نوع LR نیست. برای ضرب و تقسیم، روابط تقریبی پیشنهاد شده که در ادامه به بیان آنها می‌پردازیم.

روابط تقریبی برای ضرب

الف- اگر $M > \circ$ و $N > \circ$ آنگاه

$$M \otimes N = (m, \alpha, \beta)_{LR} \otimes (n, \gamma, \delta)_{LR} \simeq (mn, m\gamma + n\alpha, m\delta + n\beta)_{LR}$$

ب- اگر $M < \circ$ و $N > \circ$ آنگاه

$$M \otimes N = (m, \alpha, \beta)_{LR} \otimes (n, \gamma, \delta)_{LR} \simeq (mn, n\alpha - m\delta, n\beta - m\gamma)_{LR}$$

ج- اگر $M > \circ$ و $N < \circ$ آنگاه

$$M \otimes N = (m, \alpha, \beta)_{LR} \otimes (n, \gamma, \delta)_{LR} \simeq (mn, m\gamma + n\beta, m\delta + n\alpha)_{LR}$$

د- اگر $M < \circ$ و $N < \circ$ آنگاه

$$M \otimes N = (m, \alpha, \beta)_{LR} \otimes (n, \gamma, \delta)_{LR} \simeq (mn, -n\beta - m\delta, -n\alpha - m\gamma)_{LR}$$

روابط تقریبی برای تقسیم

اگر $M = (m, \alpha, \beta)_{LR}$ یک عدد فازی مثبت باشد، آنگاه

$$M^{-1} \simeq (m^{-1}, \beta m^{-2}, \alpha m^{-2})_{RL}$$

اگر $M = (m, \alpha, \beta)_{LR}$ و $N = (n, \gamma, \delta)_{LR}$ دو عدد فازی مثبت باشند، آنگاه

$$(M \otimes N) \simeq \left(\frac{m}{n}, \frac{\delta m + \alpha n}{n^2}, \frac{\gamma m + \beta n}{n^2} \right)_{LR}$$

هرچه پهناهای اعداد فازی M و N نسبت به مقادیر نمای آنها کوچک باشند، روابط فوق دقیق‌ترند. همچنین در همسایگی مقادیر نما، این روابط دقت بیشتری دارند. چون در ادامه از اعداد فازی مثلثی برای ضرایب مدل رگرسیون، استفاده خواهیم کرد، لذا به یادآوری چند نکته درباره این اعداد می‌پردازیم.

می‌دانیم که هر عدد فازی مثلثی را می‌توان به صورت $\tilde{A} = (a, s^L, s^R)_T$ نشان داد که در آن a مقدار نما (یا میانه) و s^L و s^R به ترتیب پهناهای چپ و پهنای راست \tilde{A} هستند. اگر $s^R \neq s^L$ آن‌گاه عدد فازی مثلثی A را نامتقارن گوییم. در این حالت تابع \tilde{A} را با توجه به سه مشخصه a, s^L و s^R می‌توان به صورت زیر نوشت

$$\tilde{A}(x) = \begin{cases} 1 - \frac{a-x}{s^L} & a - s^L \leq x \leq a \\ 1 - \frac{x-a}{s^R} & a < x \leq a + s^R \end{cases} \quad (1.2)$$

به گونه‌ای دیگر نیز می‌توان این تابع عضویت را نمایش داد. یعنی پهناهای راست را بر حسب پهناهای چپ بیان کرد. به این صورت که در تابع عضویت بالا قرار دهیم $s^R = ks^L$ که در آن k ، که عددی حقیقی و مثبت است، ضریب کشیدگی نامیده می‌شود. بنابراین عدد فازی مثلثی نامتقارن \tilde{A} را می‌توان با سه تایی، $\tilde{A} = (a, s^L, s^R)_T$ نیز توصیف کرد. در این حالت تابع عضویت \tilde{A} به صورت زیر در می‌آید.

$$\tilde{A}(x) = \begin{cases} 1 - \frac{a-x}{s^L} & a - s^L \leq x \leq a \\ 1 - \frac{x-a}{s^R} & a < x \leq a + ks^R \end{cases} \quad (2.2)$$

اگر $s^R = s^L = s$ ، آن‌گاه \tilde{A} را عدد مثلثی متقارن نامیده و آن را با $\tilde{A} = (a, s)_T$ نمایش می‌دهیم. در این حالت تابع عضویت \tilde{A} با توجه به مشخصه a و s به صورت زیر خواهد بود.

$$\tilde{A}(x) = \begin{cases} 1 - \frac{a-x}{s} & a - s \leq x \leq a \\ 1 - \frac{x-a}{s} & a < x \leq a + s \end{cases} \quad (3.2)$$

در رابطه (۲.۲)، اگر $k = 1$ ، آن‌گاه تابع عضویت بالا (یعنی تابع عضویت عدد فازی متقارن) به دست می‌آید.

گزاره ۲.۳.۲. فرض کنید $\tilde{A}_1 = (a_1, s_1^L, s_1^R)_T$ و $\tilde{A}_2 = (a_2, s_2^L, s_2^R)_T$ دو عدد فازی مثلثی متقارن و c یک عدد حقیقی باشد،

۱. (ا) اگر $c > 0$ آنگاه $c\tilde{A}_1 = (ca_1, cs_1^L, cs_1^R)_T$

(ب) اگر $c < 0$ آنگاه $c\tilde{A}_2 = (ca_2, -cs_2^R, -cs_2^L)_T$

۲. $\tilde{A}_1 \oplus \tilde{A}_2 = (a_1 + a_2, s_1^L + s_2^L, s_1^R + s_2^R)_T$

۱.۳.۲ فاصله بین دو عدد فازی

در این بخش فرض می‌کنیم که \tilde{A} و \tilde{B} دو عدد فازی، و $\tilde{A}_\alpha = [a_1(\alpha), a_2(\alpha)]$ و $\tilde{B}_\alpha = [b_1(\alpha), b_2(\alpha)]$ به ترتیب α -برش‌های \tilde{A} و \tilde{B} باشند.

تعریف ۳.۳.۲. فاصله بین دو عدد فازی \tilde{A} و \tilde{B} بر پایه تابع وزنی $f(\alpha)$ به صورت زیر تعریف می‌شود

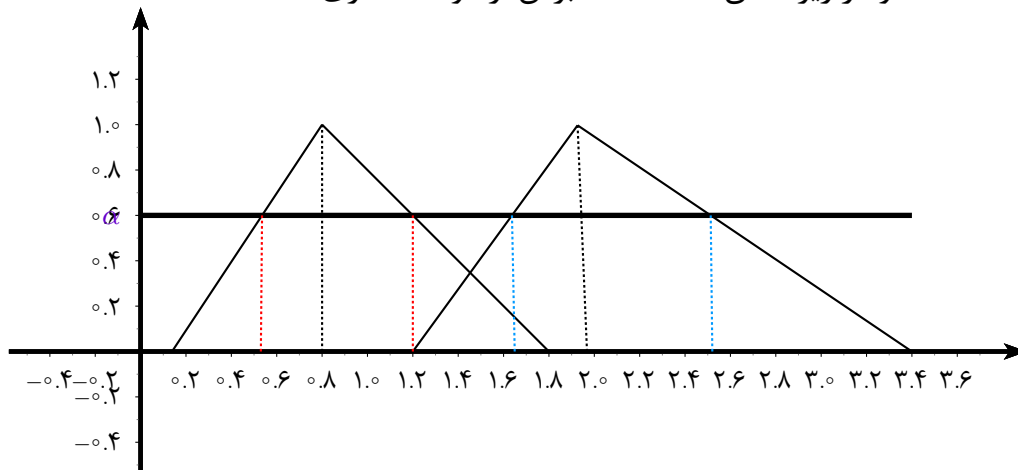
$$d(\tilde{A}, \tilde{B}) = \left[\int_0^1 f(\alpha) d^2(\tilde{A}_\alpha, \tilde{B}_\alpha) d\alpha \right]^{\frac{1}{2}} \quad (۴.۲)$$

که

$$d^2(\tilde{A}_\alpha, \tilde{B}_\alpha) = [a_1(\alpha) - b_1(\alpha)]^2 + [a_2(\alpha) - b_2(\alpha)]^2 \quad (۵.۲)$$

که در آن، $f(\alpha)$ یک تابع صعودی روی بازه $[0, 1]$ است که برای آن $f(0) = 0$ و $\int_0^1 f(\alpha) d\alpha = 1$ مقدار $d(\tilde{A}_\alpha, \tilde{B}_\alpha)$ فاصله بین α -برش‌های اعداد فازی \tilde{A} و \tilde{B} را اندازه می‌گیرد. تابع $f(\alpha)$ منجر می‌شود که درجات عضویت بالاتر، اهمیت بیشتری در تعیین فاصله بین دو عدد فازی \tilde{A} و \tilde{B} داشته باشند. شرط‌های $f(0) = 0$ و $\int_0^1 f(\alpha) d\alpha = 1$ نیز این اطمینان را می‌دهد که فاصله فوق، یک تعمیم معمولی از فاصله بین دو عدد حقیقی است.

مثال ۱.۳.۲. نمودار زیر نشان دهنده α -برش از دو عدد فازی است.



طبق نمودار فوق مقادیر برای دو عدد فازی \tilde{A} و \tilde{B} می‌توانیم روابط زیر را بیان کنیم.

$$\tilde{A} = (a, l_A, u_A) \quad \tilde{B} = (b, l_B, u_B) \quad (۶.۲)$$

حال معادله خط برای A را می‌نویسیم.

$$A_1 = \begin{cases} \frac{x - a - l_A}{l_A} & a - l_A < x < a \\ \frac{a + u_A - x}{u_A} & a < x < a + u_A \end{cases} \quad (۷.۲)$$

با تلاقی خط $y = \alpha$ با دو خط بالا داریم

$$\frac{x - a + l_A}{l_A} = \alpha \Rightarrow x - \alpha l_A - l_A + a = 0 \Rightarrow x = a - (1 - \alpha)l_A \quad (۸.۲)$$

$$\frac{a + u_A - x}{u_A} = \alpha \Rightarrow a + u_A - x - \alpha u_A \Rightarrow x = a - (1 - \alpha)u_A. \quad (۹.۲)$$

با جایگذاری (۸.۲) و (۹.۲) در فاصله تعریف شده توسط ژو [۳۵] ۴.۲ خواهیم داشت.

$$\begin{aligned} d(\tilde{A}, \tilde{B}) &= \left[\int_0^1 f(\alpha) d^{\vee}(\tilde{A}, \tilde{B}) d\alpha \right]^{\vee} \\ &= \int_0^1 \alpha (a_{\vee}(\alpha) - b_{\vee}(\alpha))^{\vee} d\alpha + \int_0^1 \alpha (a_{\wedge}(\alpha) - b_{\wedge}(\alpha))^{\vee} d\alpha \\ &= \int_0^1 \alpha [a - (1 - \alpha)l_A - b + (1 - \alpha)l_B]^{\vee} d\alpha \\ &\quad + \int_0^1 \alpha [a + (1 - \alpha)u_A - b - (1 - \alpha)u_B]^{\vee} d\alpha \\ &= \int_0^1 \alpha [(a - b) + (1 - \alpha)(l_A - l_B)] + \int_0^1 \alpha [(a - b) - (1 - \alpha)(u_A - u_B)] \\ &= \int_0^1 \alpha (a - b)^{\vee} d\alpha + \int_0^1 \alpha (1 - \alpha)^{\vee} (l_B - l_A)^{\vee} d\alpha \\ &\quad + \int_0^1 \alpha (1 - \alpha)^{\vee} (a - b)(l_B - l_A) d\alpha + \int_0^1 \alpha (a - b)^{\vee} d\alpha \\ &\quad + \int_0^1 \alpha (1 - \alpha)^{\vee} (u_B - u_A)^{\vee} d\alpha + \int_0^1 \alpha (1 - \alpha)^{\vee} (a - b)(u_B - u_A) d\alpha \\ &= (a - b)^{\vee} + \frac{1}{12} [(l_B - l_A)^{\vee} + (u_B - u_A)^{\vee}] \\ &\quad + \frac{1}{3} (a - b) [(l_B - l_A) - (u_B - u_A)] \quad (۱۰.۲) \end{aligned}$$

۴.۲ رگرسیون فازی

همان گونه که گفته شد، مدل‌های رگرسیونی برای ایجاد ارتباط بین یک متغیر وابسته و تعدادی متغیر مستقل به کار می‌روند. برای ساختن این مدل‌ها، نیاز به مشاهداتی از متغیرهای مورد مطالعه داریم. در رگرسیون کلاسیک فرض می‌شود که این متغیرها و مشاهدات مربوط به آنها، دقیق هستند. همچنین به علت عدم برازش کامل داده‌ها، در رگرسیون کلاسیک، یک جمله خطای تصادفی به مدل افزوده می‌شود. درباره این جمله خطای تصادفی و توزیع احتمالی آن مفروضاتی مانند نرمال بودن، ناهمبسته بودن، ثبات واریانس و ... در نظر گرفته می‌شود، به گونه‌ای که بتوان بر پایه این فرض‌ها، استنباط‌های آماری مانند برآورد پارامترها، پیش‌بینی

مقدار متغیر وابسته، آزمون فرض‌ها و مواردی دیگر را انجام داد. در این باره، به کارگیری تبدیلات مختلفی بر روی مشاهدات ممکن است کمک کنند تا فرض‌های اولیه مدل برقرار باشند.

اما در بسیاری اوقات ممکن است یک یا چند فرض از فرض‌های بالا برقرار نباشد یا اینکه نتوان از درستی بعضی فرض‌ها اطمینان حاصل کرد. مسلماً در این موارد مدل‌های رایج اعتبار و کارایی لازم را ندارند. یکی از راه‌حل‌ها در این زمینه، رگرسیون نیرومند است که به‌ویژه زمانی که حتی با انجام تبدیلات نتوان این فرضیات را برقرار کرد، کارایی دارند.

مشکل دیگر اینکه ممکن است در بررسی، مشاهدات مربوط به یک یا چند متغیر، نادقیق باشند و یا نادقیق گزارش شده باشند. همچنین ممکن است که متغیرهای مورد مطالعه، ذاتاً دارای ارتباطی نادقیق و مبهم باشند. یکی از شیوه‌های مهم جایگزین رگرسیون در چنین مواقعی، استفاده از رگرسیون فازی است. البته زمانی که فرضیات مربوط به خطای تصادفی برقرار نباشد یا مثلاً به دلیل حجم نمونه کم نتوان از درستی آنها اطمینان حاصل کرد، نیز می‌توان از رگرسیون فازی استفاده کرد.

بعد از معرفی نظریه مجموعه‌های فازی توسط زاده، از اوایل دهه هشتاد میلادی، گسترش آن در آمار شروع شد. یکی از مهم‌ترین و کاربردی‌ترین ارتباطات بین این نظریه و آمار کلاسیک، رگرسیون فازی است.

در یک تقسیم‌بندی کلی، بسته به اینکه هر کدام از متغیرهای وابسته، مستقل و ضرایب مدل، دقیق و یا فازی در نظر گرفته شوند، انواع رگرسیون فازی را می‌توان به حالات زیر تقسیم کرد.

[۵]

- رگرسیون فازی در حالتی که ارتباط بین متغیرهای دقیق، فازی فرض شود. به عبارت دیگر، ضرایب معادله رگرسیونی فازی در نظر گرفته شوند. این مدل‌ها به مدل‌های رگرسیون امکانی معروفند.

- رگرسیون فازی در حالتی که متغیرها و یا مشاهدات مربوط به آنها نادقیق و فازی باشند.

- مدل با ورودی و خروجی فازی و ضرایب غیر فازی. (به این مدل رگرسیون الزامی نیز گفته می‌شود).

- مدل با خروجی و ضرایب فازی. (اکثر مقالات ارائه شده در رگرسیون فازی در ارتباط با این گونه مدل‌ها می‌باشند).

- مدل با ورودی‌های دقیق و خروجی فازی و ضرایب دقیق به‌علاوه یک خطای فازی.

- رگرسیون در حالتی که هم متغیرها و هم ضرایب مدل فازی در نظر گرفته شوند.

علاوه بر تقسیم‌بندی ذکر شده، انواع رگرسیون فازی با در نظر گرفتن جنبه‌های دیگر،

به‌صورت زیر تقسیم بندی می‌شوند: [۸]

– از نظر نوع اعداد فازی
 عامل دیگری که به تنوع رگرسیون فازی می‌افزاید، نوع توابعی است که برای نشان دادن تابع عضویت کمیت‌های فازی انتخاب می‌شوند. همان‌طور که در قبل ذکر شد انواع اعداد فازی عبارت‌اند از: اعداد فازی مثلثی، اعداد فازی سهموی، اعداد فازی نرمال.
 – از نظر انواع روش
 مشابه با رگرسیون آماری، اصل یا اصولی که برآورد پارامترهای مدل براساس آن‌ها انجام می‌پذیرند نیز عاملی در ایجاد تنوع در رگرسیون فازی است. در رگرسیون فازی، دو روش مهم، رگرسیون امکانی و رگرسیون کمترین توان‌های دوم است.
 – از نظر تعامل ضرایب
 وجود یا عدم وجود تعامل بین ضرایب فازی مدل، رگرسیون فازی را به دو نوع رگرسیون فازی رگرسیون فازی با تاثیر متقابل و رگرسیون فازی بدون تاثیر متقابل تقسیم می‌کند.

۵.۲ رگرسیون امکانی

رگرسیون امکانی نخستین بار توسط تاناکا و همکاران [۳۳] پیشنهاد شد. آنها یک مدل برنامه‌ریزی خطی برای محاسبه ضرایب فازی مدل رگرسیونی با ورودی و خروجی دقیق ارائه کردند که تابع هدف آن که باید حداقل شود به صورت مجموع پهناهای ضرایب (که اعداد فازی مثلثی متقارن هستند) است و محدودیت‌های آن باعث می‌شود که به ازای مقدار مشخصی از h ، مقدار متغیر پاسخ دست کم به اندازه h در خروجی فازی برآورد شده عضو باشد. این روش سپس به مدلی با خروجی فازی نیز تعمیم یافت. در این رویکرد چون تابع عضویت مجموعه‌های فازی اغلب به عنوان توزیع‌های امکانی توصیف می‌شود، روش مرتبط را رگرسیون امکانی می‌نامند.

از جمله مشکلات رویکرد اولیه تاناکا این است که جواب مساله وابسته به مقیاس متغیرهای ورودی x_j است. ایراد دیگر اینکه در برخی مواقع تعداد زیادی از ضرایب فازی مدل به صورت دقیق برآورد می‌شوند. تاناکا و همکاران [۳۰] این نقیصه را بدین صورت اصلاح کردند که به جای مجموع پهناهای ضرایب، مجموع پهناهای پاسخ‌های برآورد شده حداقل شوند [۱۹].
 مشکل دیگر مساله تاناکا این است که اگر به جای x_j از $(x_j - \bar{x}_j)$ (که میانگین متغیر x_j است)، استفاده شود، تابع برآورد شده بسیار متفاوت خواهد بود. نقیصه دیگر حساسیت نسبت به داده‌های پرت است. برای حل این مشکل، تاناکا و همکاران [۳۳] از مدل عطفی^۹ استفاده کردند که به جای اینکه $-\alpha$ برش‌های پاسخ‌های برآورد شده از مدل، شامل $-\alpha$ برش‌های پاسخ‌های داده شده متناظرشان باشند، فقط نیاز دارد که هر $-\alpha$ برش پاسخ مشاهده شده، اشتراک داشته باشد.

^۹Conjunctive

محققان دیگر، تغییرات اندک یا قابل توجه در روش رگرسیون امکانی ایجاد کردند. از بین آنها، ساویک و پدريچ [۲۸] به دلیل اینکه مدل تاناکا نماهای داده‌ها را در محاسبه ضرایب رگرسیون به حساب نمی‌آورد، حال آنکه این نماها به دلیل داشتن درجه عضویت از اهمیت بیشتری برخوردارند، یک روش دو مرحله‌ای پیشنهاد کردند. ابتدا نماهای ضرایب فازی را با استفاده از رگرسیون کمترین مربعات معمولی به دست آوردند و سپس پهناهای آنها را با روش تاناکا و با استفاده مراکز برآورده شده، برآورد کردند. تاناکا و ایشیبوجی [۳۱] پیشنهاد کردند که نماهای ضرایب با استفاده از رگرسیون کمترین مربعات معمولی برآورد شوند اما برای به دست آوردن پهناهای ضرایب، از یک نوع برنامه‌ریزی درجه دویی مشابه مساله قبل استفاده می‌شود. چانگ و لی [۲۱] خاطر نشان کردند که زمانی که روند نماهای اعداد فازی در مغایرت با روند پهناها باشد، روش تاناکا اکثراً منجر به تفسیر اشتباه می‌شود. برای اجتناب از این مشکل، آنها قید مثبت بودن پهناهای ضرایب را از روش تاناکا حذف کردند.

۶.۲ رگرسیون کمترین توان‌های دوم فازی

رهیافت دیگری در رگرسیون فازی، روش کمترین توان‌های دوم فازی است که اولین بار توسط دیاموند [۱۲] و کلمینس [۱۰] ارائه شد. این روش تعمیمی از رگرسیون کمترین توان‌های دوم کلاسیک است که معمولاً بر پایه تعریف‌هایی برای فاصله بین اعداد فازی است. شیوه‌های متنوعی برپایه روش کمترین توان‌های دوم فازی پیشنهاد شده است که این تنوع عوامل مختلفی دارد مانند دقیق یا فازی بودن مشاهدات مربوط به متغیرهای تبیینی، دقیق یا فازی بودن مشاهدات متغیر پاسخ، نوع تعریف فاصله بین دو عدد فازی، معیارهای نیکویی برازش و برخی موارد دیگر. نمونه‌هایی از کارهای صورت گرفته در این زمینه به شرح زیر است.

دیاموند [۱۲] با معرفی فاصله‌ای بر اعداد فازی، مدل‌هایی را برای برازش کمترین توان‌های دوم فازی برای مدل با ورودی دقیق و خروجی فازی و مدل با ورودی و خروجی فازی پیشنهاد داد. دیاموند و کورنر [۱۶] رگرسیون فازی را براساس متغیرهای تصادف فازی ارائه کردند. کلمینس [۱۰] رگرسیون کمترین توان‌های دوم را برای توابع عضویت درجه دو مطرح کرد. ژو و لی [۳۵] یک روش کمترین توان‌های دوم را برای برآورد ضرایب فازی مدل با استفاده از متری بر مبنای α -برش‌ها ارائه کردند و با یک معیار نیکویی برازش، خوبی مدل رگرسیونی را اندازه گرفتند. محمدی و طاهری [۲۴] این روش را برای برازش مدل‌های رگرسیونی در خاک شناسی موسوم به پدوترانسفر به کار بردند.

کائو چیو [۲۰] یک روش دو مرحله‌ای برای برآورد مدل رگرسیون خطی فازی پیشنهاد نمودند. در مرحله اول با مقادیر غیر فازی شده مشاهدات فازی و به روش کمترین مربعات معمولی، ضرایب دقیق مدل را برآورد کردند و سپس در مرحله دوم به منظور افزایش قدرت توضیح دهنده‌گی داده‌ها، یک عبارت خطای فازی به مدل برآورد شده اضافه کردند. کاپی و همکاران [۱۱] یک مدل رگرسیون خطی کلی را برای مطالعه وابستگی یک متغیر پاسخ فازی از

نوع LR، روی یک مجموعه متغیرهای دقیق تبیینی از طریق یک روند برآورد کمترین توان‌های دوم تکرار شونده، معرفی کردند.

۱.۶.۲ مدل رگرسیون خطی با ضرایب فازی

در یک تقسیم‌بندی کلی، انواع رگرسیون فازی را می‌توان به سه حالت زیر تقسیم کرد:

الف: رگرسیون فازی در حالتی که ارتباط بین متغیرها فازی فرض می‌شود. به عبارت دیگر ضرائب معادله رگرسیونی، فازی در نظر گرفته می‌شوند.

ب: رگرسیون فازی در حالتی که متغیرها و یا مشاهدات مربوط به متغیرها نادقیق و فازی هستند.

ج: رگرسیون فازی در حالتی که هم متغیرها و هم ضرائب مدل، فازی در نظر گرفته می‌شوند

در اینجا مدل رگرسیون خطی با ضرایب فازی (حالت الف) را مورد مطالعه قرار داده‌ایم. به مدل‌های رگرسیون با ضرایب فازی، گاهی مدل‌های رگرسیون امکانی هم گفته می‌شود. صورت کلی مدلی که مورد بحث قرار می‌گیرد عبارت است از

$$\tilde{Y} = f(\mathbf{x}, A) = \tilde{A}_0 + \tilde{A}_1 x_1 + \tilde{A}_2 x_2 + \dots + \tilde{A}_{p-1} x_{p-1} \quad (11.2)$$

که در آن \tilde{Y} متغیر وابسته یا اصطلاحاً خروجی فازی است، $\mathbf{x} = (x_1, x_2, \dots, x_{p-1})$ بردار متغیرهای پاسخ یا اصطلاحاً بردار ورودی (با مقادیر حقیقی) و $A = \{\tilde{A}_0, \tilde{A}_1, \dots, \tilde{A}_{p-1}\}$ یک مجموعه از اعداد فازی است. مسئله‌ای که درصدد حل آن هستیم به این صورت مطرح می‌شود: مجموعه‌ای از داده‌های معمولی به صورت $(y_1, x_1), (y_2, x_2), \dots, (y_n, x_n)$ در اختیار داریم. می‌خواهیم پارامترهای فازی $\tilde{A}_0, \tilde{A}_1, \dots, \tilde{A}_{p-1}$ را به گونه‌ای تعیین کنیم که مدل (۱۱.۲) براساس برخی از معیارهای نیکویی برازش، بهترین برازش را به داده‌های مذکور داشته باشد.

تابع عضویت مدل با ضرایب نامتقارن

اگر A_j ها $j = 0, 1, \dots, p-1$ اعداد فازی نامتقارن و x_i ها نیز اعداد حقیقی و مثبت باشند، آن‌گاه بنا به رابطه ۱۱.۲ و گزاره (۲.۳.۲)، \tilde{Y} یعنی خروجی فازی نیز یک عدد فازی مثلثی نامتقارن به صورت $\tilde{Y} = (f^c(\mathbf{x}), f_s^L(\mathbf{x}), f_s^R(\mathbf{x}))$ که در آن $f^c(\mathbf{x})$ نما و $f_s^L(\mathbf{x})$ پهنای چپ و $f_s^R(\mathbf{x})$ پهنای راست \tilde{Y} می‌باشد و به صورت زیر به دست می‌آیند.

$$f^c(\mathbf{x}) = a_0 + a_1 x_1 + \dots + a_{p-1} x_{p-1}$$

$$f_s^L(\mathbf{x}) = s_0^L + s_1^L x_1 + \dots + s_{p-1}^L x_{p-1}$$

$$f_s^R(\mathbf{x}) = s_0^R + s_1^R x_1 + \dots + s_{p-1}^R x_{p-1}$$

به بیان دیگر تابع عضویت \tilde{Y} عبارت است از

$$\tilde{Y}(y) = \begin{cases} 1 - \frac{f^c(\mathbf{x}) - y}{f_s^L(\mathbf{x})} & f^c(\mathbf{x}) - f_s^L(\mathbf{x}) \leq y \leq f^c(\mathbf{x}) \\ 1 - \frac{y - f^c(\mathbf{x})}{f_s^R(\mathbf{x})} & f^c(\mathbf{x}) < y \leq f^c(\mathbf{x}) + f_s^R(\mathbf{x}) \end{cases} \quad (12.2)$$

اگر بخواهیم تابع عضویت بالا را برحسب ضرایب کشیدگی بیان کنیم، قرار می‌دهیم $s_i^R = k_i s_i^L$ ، بدین ترتیب $f_s^R(\mathbf{x})$ به صورت زیر تغییر می‌یابد

$$f_s^R(\mathbf{x}) = k_0 s_0^L + k_1 s_1^L x_1 + \dots + k_{p-1} s_{p-1}^L x_{p-1}$$

فصل ۳

رگرسیون لجستیک فازی با رویکرد امکانی

در این فصل می‌خواهیم رگرسیون لجستیک فازی در حالتی که مقادیر متغیر پاسخ مبهم و به صورت اعداد فازی هستند و نشان دهنده امکان موفقیت در موقعیت می‌باشند را معرفی کنیم. در همین راستا بخت امکانی و ترم زبانی را معرفی می‌کنیم و سپس مدل را صورت‌بندی و شیوه برآورد پارامتر به روش پوراحمد و همکاران [۱] را ارائه خواهیم کرد.

۱.۳ معرفی مدل رگرسیون لجستیک فازی

همان‌طور که در فصل اول بیان شد، مدل رگرسیون لجستیک مبتنی بر چند فرض اساسی هستند. برخی از این مفروضات را به‌طور خلاصه بیان می‌کنیم. نتایج مربوط به متغیر وابسته به صورت 0 و 1 است، رابطه بین متغیرها دقیق است، خطای مدل یک خطای تصادفی است. همچنین برای آنکه بتوانیم برخی تحلیل‌های پارامتری آماری را در مورد مدل انجام دهیم بعضاً حجم نمونه باید به اندازه کافی بزرگ باشد.

اما در بعضی موارد در عمل مفروضات بیان شده برقرار نیستند. مانند اینکه مشاهده متغیر پاسخ، تخصیص مربوط به موفقیت یا شکست، منجر به نتایج کاملاً دقیق نشود، یا اینکه ممکن است عدم اطمینان بیان شده با خطای مدل، از نوع تصادفی نباشد. علاوه بر این موارد، در

بعضی موارد به دلیل هزینه‌های بالا و ملاحظات اخلاقی در جمع‌آوری اطلاعات، اخذ تعداد زیاد نمونه غیر ممکن یا دشوار است. برای غلبه بر محدودیت‌هایی که گفته شدو برای تحلیل داده‌ها و مدل‌بندی داده‌های مبهم نظریه مجموعه‌های فازی و نظریه امکان می‌تواند کار ساز باشد.

برای مدل‌بندی رابطه بین متغیر پاسخ دودویی و متغیرهای تبیینی فرض می‌کنیم که متغیرهای تبیینی دقیق و پاسخ‌ها فازی هستند و بردار متغیرها به صورت $(x_{i0}, x_{i1}, \dots, x_{in}, \tilde{Y}_i)$ باشد. از آنجایی که متغیرهای پاسخ دودویی غیر دقیق (فازی) جمع‌آوری می‌شوند، بنابراین فرض احتمالات توزیع برنولی را نمی‌توان برای این داده‌ها در نظر گرفت. لذا احتمال $P(Y_i = \pi_i) = 1$ را نمی‌توان محاسبه کرد. به‌عنوان راه حل به‌جای آنکه اندازه پیشامد محاسبه شود، اندازه امکان آن را می‌توان در نظر گرفت. امکان وجه دیگری از ابهام است.

تعریف ۱.۱.۳. [۴] فرض کنید $\mu_i, i = 1, 2, \dots, m$ نشان دهنده امکان موفقیت باشد، آنگاه در نظر می‌گیریم:

الف یک مقدار واقعی دقیق μ_i به‌طوری که $0 \leq \mu_i \leq 1$ که $\mu_i \in \mathbb{R}$;

ب یک عبارت توصیفی به‌عنوان مثال $\mu_i \in \{\dots, low, medium, high, \dots\}$ (توابع عضویت) توابع فوق باید به‌نحوی بیان شوند که کل دامنه $(0, 1)$ را در برگیرند نسبت $\frac{\mu_i}{1 - \mu_i}$ به‌عنوان بخت امکانی^۱ برای فرد i ام در نظر گرفته می‌شود و نسبت امکان موفقیت به امکان شکست بیان می‌کند.

به دلیل اهمیت مدل‌بندی لجستیک در داده‌های پزشکی، مبهم و نادقیق بودن یافته‌ها و حجم نمونه کم عمدتاً در این مطالعات با آن روبرو می‌شویم. هدف این فصل ارائه مدلی است که امکان ابتلا به بیماری را در مواردی که افراد مورد مطالعه مشکوک به بیماری هستند مدل‌بندی کند. این مدل رگرسیون لجستیک فازی است که در آن مقادیر متغیر پاسخ مبهم است و جای اعداد صفر و یک، مجموعه‌های فازی از بازه واحد هستند. در این مدل متغیرهای تبیینی، دقیق فرض می‌شوند. مدل مذکور قادر است امکان بیمار بودن فرد مشکوک به بیماری را در زمان حال و یا ابتلای او در آینده را بیان کند.

فرض کنید می‌خواهیم ابتلا به یک بیماری خاص را براساس مجموعه‌ای از عوامل خطر مدل‌بندی کنیم. اگر ابتلا به بیماری را متغیر پاسخ دودویی (بیمار/ سالم) در نظر بگیریم، مدل‌بندی رابطه بین این متغیرها با استفاده از مدل رگرسیون لجستیک صورت می‌گیرد. برای برآورد پارامترهای مدل مذکور به نمونه‌ای به حجم n نیاز داریم که شامل تعدادی افراد بیمار و سالم با مقادیر معلوم عوامل خطر است. معمولاً پزشک در تشخیص وجود بیماری در فرد تردید دارد و امکان ابتلا به بیماری برای افراد مبهم است. در مدل‌های آماری این نمونه‌های مبهم از نمونه اصلی کنار گذاشته می‌شود و در مدل‌بندی شرکت داده نمی‌شوند، در واقع این

^۱Possibilistic odds

نمونه‌ها دور ریخته می‌شوند. بنابراین برای محاسبه احتمال $\pi = P(Y = 1)$ از نظریه فازی کمک می‌گیریم. و همان‌طور که گفته شد به‌جای احتمال بیمار بودن، امکان ابتلا به بیماری براساس عوامل خطر را مدل‌بندی می‌کنیم. در مدل، امکان ابتلا را یک متغیر زبانی در نظر می‌گیریم. قبل از ارائه مدل، مفهوم متغیر زبانی^۲ را تعریف می‌کنیم.

تعریف ۲.۱.۳. [۱] یک متغیر زبانی توسط یک پنج تایی مرتب $(X, T(X), U, G, M)$ تعریف می‌شود که در آن X نام متغیر است، U مجموعه مرجع است، $T(X)$ مجموعه ترم (واژه‌های مربوط به متغیر است و هر ترم یک مجموعه فازی است که توسط قاعده نحوی G تولید می‌شود و سرانجام M یک قاعده معنایی است که به هر ترم $T(X)$ معنای آن را مربوط می‌سازد، یعنی تابع عضویت آن ترم را مشخص می‌کند.

مثال ۱.۱.۳. متغیر سن را با مجموعه $U = [0, 200]$ در نظر بگیرید. مقادیر این متغیر اعدادی مانند ۳۵، ۷۰ و ۱۱۵ هستند. با این مقادیر، متغیر سن یک متغیر معمولی است. اما در زبان معمولی از مفاهیمی استفاده می‌کنیم که در واقع مفاهیم مبهمی از متغیر سن هستند، مانند خردسال، نوجوان، جوان، پیر و در این مثال عنصر X سن است، مجموعه مرجع $U = [0, 200]$ تعریف می‌شود، همچنین داریم:

M مجموعه قواعدی است که به هر مقدار زبانی سن، معنای آن را نسبت می‌دهد. برای مثال ترم «جوان» را به‌صورت یک مجموعه فازی با تابع عضویت زیر می‌توان نوشت.

$$A(x) = \begin{cases} 1 & 0 \leq x < 30 \\ 1 - 2 \left(\frac{x - 30}{30} \right)^2 & 30 \leq x < 45 \\ 2 \left(\frac{60 - x}{30} \right)^2 & 45 \leq x < 60 \\ 0 & 60 \leq x \end{cases}$$

حال به کمک تعریف‌هایی که بیان شد به صورت‌بندی و برازش مدل و در نهایت برآورد پارامترهای مجهول مدل می‌پردازیم.

۱.۱.۳. برازش مدل رگرسیون لجستیک با رویکرد امکانی

در فصل قبل تاریخچه مختصری رگرسیون امکانی معرفی شد. در این قسمت به معرفی این روش می‌پردازیم.

در ابتدا از خبره می‌خواهیم با توجه به تجربه خود ترم‌های زبانی را برای هر موقعیت بیان کند، آنگاه ترم‌های زبانی را به‌صورت یک مجموعه فازی با تابع عضویتی بر بازه $(0, 1)$ تعریف می‌کنیم. آنگاه با استفاده از اصل گسترش، تبدیل یافته هر یک از ترم‌های زبانی را به‌عنوان

^۲Linguistic Approach

ورودی‌های متغیر پاسخ در مدل پیشنهادی در نظر می‌گیریم. به عبارت دیگر $\tilde{\mu}_i$ یک متغیر زبانی است، که امکان ابتلا به بیماری را تعریف می‌کنیم. تبدیل یافته $\tilde{\mu}_i$ یعنی $\frac{\tilde{\mu}_i}{\sqrt{-\tilde{\mu}_i}}$ را بخت امکانی می‌نامیم و لجیت آن، $\tilde{y}_i = \ln \frac{\tilde{\mu}_i}{\sqrt{-\tilde{\mu}_i}}$ را به عنوان مشاهده متغیر پاسخ در نظر می‌گیریم. با استفاده از تبدیل لجیت \tilde{y}_i به صورت خطی با پارامترهای مدل در ارتباط است و در نتیجه برای برآورد پارامترهای فازی مدل از رگرسیون خطی فازی می‌توان استفاده کرد. با توجه به تعریف \tilde{y}_i و با استفاده از اصل گسترش، تابع عضویت \tilde{y}_i عبارت است از

$$\tilde{y}_i(y) = \sup_{\forall x \in \mathbb{R}, x \neq 1; \ln \frac{x}{\sqrt{1-x}} = y} \mu_i(x) \quad , i = 1, 2, \dots, n \quad (1.3)$$

طبق معادله (۱.۳) \tilde{y}_i یک عدد مثلثی متقارن و برآورد تبدیل لگاریتم بخت امکانی است.

بنابراین داریم

$$\tilde{y}_i = (f_i^c(x), f_{is}^L(x), f_{is}^R(x))_T \quad (2.3)$$

که در آن

$$\begin{aligned} f_i^c(x) &= a_0^c + a_1^c x_{i1} + \dots + a_{p-1}^c x_{ip-1} \\ f_{is}^L(x) &= s_0^L + s_1^L x_{i1} + \dots + s_{p-1}^L x_{ip-1} \\ f_{is}^R(x) &= s_0^R + s_1^R x_{i1} + \dots + s_{p-1}^R x_{ip-1} \end{aligned}$$

به ترتیب مرکز، پهناهای چپ و پهناهای راست عدد فازی هستند. تابع عضویت برای برآورد فازی خروجی به صورت زیر است: [۴]

$$\tilde{Y}_i(\tilde{y}_i) = \begin{cases} 1 - \frac{f_i^c(x) - y_i}{f_{is}^L(x)} & f_i^c(x) - f_{is}^L(x) \leq y_i \leq f_i^c(x) \\ 1 - \frac{y_i - f_i^c(x)}{f_{is}^R(x)} & f_i^c(x) \leq y_i \leq f_i^c(x) + f_{is}^R(x). \end{cases} \quad (3.3)$$

اگر $s_i^L = s_i^R = s_i$ باشد، عدد فازی مثلثی متقارن خواهیم داشت. در این صورت برای \hat{Y}_i شرط زیر را برای پهناهای چپ و راست عدد فازی خواهیم داشت.

$$f_i^L(x) = f_i^R(x) = f_i(x).$$

در این پایان نامه اعداد فازی مورد بررسی ما، اعداد فازی مثلثی نامتقارن هستند. همان طور که می‌دانیم، لگاریتم بخت امکانی \tilde{Y}_i ، شانس گرفتن یا داشتن ویژگی تعریف شده در عضو i ام داده‌ها می‌باشد. با توجه به اصل گسترش اگر \tilde{M} یک عدد فازی با تابع عضویت \tilde{Y}_i باشد و $f(x) = \exp(x)$ ، در این صورت $f(\tilde{M}) = \exp(\tilde{M})$ یک عدد فازی با تابع عضویت زیر است

$$\exp(\tilde{M}(x)) = \begin{cases} \tilde{M}(\ln x) & x > 0 \\ 0 & o.w. \end{cases} \quad (4.3)$$

بنابراین، پس از برآورد ضرایب مدل، می‌توان تابع عضویت بخت امکانی را تعیین کرد. (۵.۳)

$$\exp(\tilde{Y}(x)) = \tilde{Y}_i(\ln(x)) = \begin{cases} 1 - \frac{f_i^c(x) - \ln(x)}{f_{is}^L(x)} & f_i^c(x) - f_{is}^L(x) \leq \ln(x) \leq f_i^c(x) \\ 1 - \frac{y_i - f_i^c(x)}{f_{is}^R(x)} & f_i^c(x) \leq \ln(x) \leq f_i^c(x) + f_{is}^R(x). \end{cases}$$

از این رو، برای یک مورد جدید فازی، مدل می‌تواند بخت امکانی خود را به عنوان یک عدد فازی با استفاده از برداری از مشاهدات ورودی (مشاهدات تبیینی) دقیق پیش‌بینی کند.

۲.۱.۳ برآورد ضرایب فازی

در این بخش قصد داریم ضرایب $n, \dots, 1, j = \bar{b}_j$ را با استفاده از رویکرد امکانی در مدل رگرسیون خطی فازی با متغیر پاسخ فازی، ضرایب فازی و متغیر تبیینی دقیق (غیر فازی) برآورد کنیم. ایده اساسی در این روش این است که مدل فازی به دست آمده را به حداقل برسانیم که این کار با به حداقل رساندن ضرایب فازی مقدور می‌شود. بنابراین فرض می‌شود که:

۱. برای مجموعه مشاهدات، دارای درجه عضویت به بزرگی h در تابع خروجی برآورد فازی است. $\tilde{Y}_i(y_i) \geq h$ که ر آن

$$y_i = \ln\left(\frac{\mu_i}{1 - \mu_i}\right), h \in (0, 1). \quad (6.3)$$

۲. ضرایب فازی با به حداقل رساندن مدل فازی به دست می‌آیند. از آنجایی که فازی بودن اعداد فازی با پهنای آن افزایش می‌یابد، با به حداقل رساندن مجموع پهنای خروجی فازی، مقدار فازی مدل حداقل می‌شود.

تعیین ضریب فازی منجر به یک مسئله برنامه‌ریزی خطی می‌شود که در آن تابع هدف مجموع پهنای خروجی فازی است.

$$Z = n(s_0^L + s_0^R) + \sum_{j=1}^n [(s_j^L + s_j^R) \sum_{i=1}^n x_{ij}]. \quad (7.3)$$

که در آن x_{ij} ، i امین مقدار مشاهده شده برای j امین متغیر است. از سوی دیگر، با توجه به معادله (۵.۳) هر ضابطه را با توجه به $\tilde{Y}_i(y_i) \geq h$ را می‌توان به صورت زیر نوشت.

$$1 - \frac{f_i^c(x) - y_i}{f_{is}^L(x)} \geq h \Rightarrow (1 - h)s_0^L + (1 - h) \sum_{j=1}^m s_j^L x_{ij} - a_0^c - \sum_{j=1}^m a_j^c x_{ij} \geq -y_i$$

$$1 - \frac{f_i^c(x) - y_i}{f_{is}^R(x)} \geq h \Rightarrow (1-h)s_o^R + (1-h) \sum_{j=1}^m s_i^R x_{ij} - a_o^c - \sum_{j=1}^m a_j^c x_{ij} \geq -y_i$$

با توجه به نامعادلات بالا، $2m$ قید داریم. از این رو تابع هدف (۷.۳) را با استفاده از الگوریتم‌های برنامه‌نویسی خطی مانند سیمپلکس، به منظور برآورد مقدار m و پهنای چپ و راست هر ضریب، به حداقل رساند.

۳.۱.۳ معیار نیکویی برازش

همانند سایر مدل‌سازی‌های آماری، مدل‌های مبتنی بر قوانین فازی، باید توسط برخی از روش‌ها (که بهترین مدل را برای داده‌ها برازش می‌کند) ارزیابی شوند. چندین روش خوب برای مدل‌های فازی ارائه شده است [۱۵] در اینجا دو روش برای ارزیابی مدل‌های رگرسیون لجستیک فازی پیشنهاد می‌کنیم.

میانگین درجه عضویت MDM

تعریف ۳.۱.۳. مدل رگرسیون لجستیک فازی که براساس مشاهدات دقیق است را در نظر بگیرید. میانگین درجه عضویت (MDM) برای مقادیر مشاهده شده در تابع عضویت برآورد شده به صورت زیر تعریف می‌شود.

$$MDM = \frac{1}{n} \sum_{i=1}^n \tilde{Y}(y_i) = \frac{1}{n} \sum_{i=1}^n \exp\left(\tilde{Y}_i\left(\frac{\mu_i}{1-\mu_i}\right)\right).$$

این یک معیار برای ارزیابی مدل است.

مقادیر بزرگ عضویت در مقادیر مشاهده شده تایید می‌کند که مدل ساخته شده از این داده‌ها، به خوبی از داده‌ها پشتیبانی می‌کند. بیشترین مقدار MDM برابر ۱ و کمترین مقدار آن ۰ است. بنابراین مقدار نزدیک به یک نشان می‌دهد که مدل مناسبی برازش داده شده است.

میانگین مجموع توان‌های دوم خطای پیش بینی

روشی دیگر برای نیکویی برازش، اندازه‌گیری فاصله بین دو مقدار پاسخ مشاهده شده و پاسخ پیش‌بینی شده است. هر چه مقدار خروجی پیش‌بینی شده به مدل مشاهده شده نزدیکتر باشد، توان مدل برای پیش‌بینی وضعیت واقعی نمونه‌ها بالاتر می‌رود. با این حال، از آنجا که پاسخ مشاهده شده یک مقدار دقیق است، در حالی که مقدار پیش‌بینی یک عدد فازی می‌باشد، به همین علت لازم است که عدد فازی ابتدا با استفاده از روش‌های تبدیل فازی به دقیق، مقادیر پاسخ را به اعداد دقیق تبدیل کنیم [۱۳].

تعریف ۴.۱.۳. مدل رگرسیون لجستیک فازی با مشاهدات ورودی دقیق، $\frac{\mu_i}{1-\mu_i}$ ، و مقادیر برآورد شده فازی، $exp(\tilde{Y}_i)$ ، را در نظر بگیرید. میانگین توان دوم خطای پیشگویی، MSE، به صورت زیر به دست می‌آید.

$$MSPE = \frac{1}{n} \sum_{i=1}^n [def(exp(\tilde{W}_i)) - \frac{\mu_i}{1-\mu_i}]$$

در معادله بالا، $def(exp(\tilde{W}_i))$ تبدیل شده فازی به دقیق^۳ از $exp(\tilde{Y}_i)$ است، که در تعریف زیر به آن پرداخته‌ایم.

تعریف ۵.۱.۳. فرض کنید که \tilde{Y} یک عدد فازی در مجموعه \mathbb{R} است، برای تبدیل عدد فازی \tilde{Y} به عدد دقیق، به روش زیر عمل می‌کنیم. [۳۴]

$$def_{COG}(\tilde{Y}) = \frac{\int_x x \tilde{Y} dx}{\int_x \tilde{Y} dx} \quad (۸.۳)$$

۲.۳ انتخاب متغیر در رگرسیون لجستیک فازی

معمولاً در انتخاب مدل، دو معیار متقابل دخالت دارند.

۱. برای آن که معادله برای هدف‌های پیشگویی مفید باشد، ممکن است مایل باشیم مدل شامل حداکثر متغیرهای پیشگو باشد.

۲. به‌علت هزینه دستیابی به اطلاعات درباره تعداد زیاد متغیر پیشگو و ارائه آنها، ممکن است تمایل داشته باشیم معادله دارای حداقل تعداد متغیر پیشگو باشد.

شیوه یکتایی برای انتخاب مدل وجود ندارد. حتی معیارهای مختلف انتخاب مدل در یک مساله خاص هم لزوماً به یک مدل رگرسیونی یکتا و منحصر به فرد منجر نمی‌شود. حال برای ارزیابی صحیح رگرسیون لجستیک فازی برخی از معیارها را معرفی می‌کنیم.

همان‌طور که در فصل قبل بیان شد، مجموع مربعات خطا براساس فاصله بین مقادیر مشاهده شده و مقادیر پیشگویی شده محاسبه می‌شود. فاصله ژو (۴.۲)، یک فاصله رایج در تجزیه و تحلیل رگرسیون فازی است. در این فصل نیز برای انتخاب مدل مناسب از این فاصله استفاده می‌کنیم. میانگین توان‌های دوم خطا براساس متر ژو [۳۵] به صورت زیر تعریف می‌شود.

$$MPE = \frac{\sum_{i=1}^n d^{\tilde{Y}}(\tilde{y}_i, \tilde{Y}_i)}{n} \quad (۹.۳)$$

^۳ defuzzification

^۴ Center of Gravity Defuzzification Method

معیار اطلاعات آکائیک (AIC) در محیط فازی توسط آکائیک^۵ [۹]، تعریف شده است که معمولاً به عنوان یک معیار مناسب برای مقایسه مدل‌ها محسوب می‌شود. در تعریف زیر این معیار در محیط فازی و برای رگرسیون لجستیک فازی تعمیم داده شده است.

تعریف ۱.۲.۳. در محیط فازی، معیار ACI به صورت زیر تعریف می‌شود.

$$AIC = n \ln(SSE) - n \ln(n) - 2(p + 1) \quad (10.3)$$

که n تعداد نمونه‌ها و p تعداد متغیرهای مدل و $SSE = \sum_{i=1}^n d^2(\tilde{y}_i, \tilde{Y}_i)$ است.

تعریف ۲.۲.۳. در دامنه اعداد فازی، C_p به صورت زیر تعریف می‌شود:

$$C_p = \frac{SSE_{model}}{SSE_{saturated.model}} + 2(p + 3) - n \quad (11.3)$$

که در آن n حجم نمونه، p تعداد ضرائب، SSE_{model} مجموع مربعات خطای مدل مورد بررسی و $SSE_{saturated.model}$ مجموع مربعات مدل که شامل همه متغیرهای پیشگو می‌باشد.

۱.۲.۳ روش پیشرو برای انتخاب متغیر به شیوه سلمانی و همکاران

سلمانی و همکاران روش پیشرو برای انتخاب بهترین متغیر در رگرسیون لجستیک فازی را معرفی کردند. در روش پیشرو در حالت غیر فازی، متغیرهای پیشگو براساس کمترین میانگین مربعات خطا به مدل اضافه می‌شوند و تا زمانی که بهبود در مدل ایجاد شود روند اضافه شدن متغیرهای پیشگو به مدل ادامه دارد. مزیت این روش در این است که تعداد تکرار کمی دارد. آنها از معیار سطح اثر بخشی^۶ (LE) به عنوان یک معیار برای ارزیابی هر مدل در روش پیشرو استفاده کردند.

۲.۲.۳ معیار LE

اهمیت انتخاب سطح ورود و خروج به منظور قضاوت در مورد ارزش متغیرهای کمکی در روش انتخاب پیشرو مهم است. معیار سطح اثر بخشی برای تصمیم‌گیری در مورد وجود یا عدم وجود متغیرها در مدل معرفی شده است. با توجه به LE، متغیری که وارد مدل می‌شود، می‌تواند حداقل درصد LE از مدل قبلی را بهبود بخشد. اگر سطح کارایی متغیر وارد شده به مدل کم باشد منجر به افزایش ورود متغیرها به مدل می‌شود. این پیچیدگی مدل را افزایش می‌دهد و آنالیز را دشوار می‌سازد، لذا پیش‌بینی با عدم اطمینان بیشتری انجام می‌شود. از سوی دیگر، افزایش تعداد متغیرها در مطالعات بالینی شامل صرف هزینه زیادی برای جمع‌آوری اطلاعات در مورد متغیرها می‌شود. بنابراین انتخاب این معیار بستگی زیادی به شرایط مطالعه دارد. سلمانی و همکاران [۲۷]، یک الگوریتم پیشنهادی را برای انتخاب مدل مناسب معرفی کردند.

^۵Akaike

^۶Level of efficacy

۳.۲.۳ الگوریتم انتخاب متغیر به روش پیشرو

الگوریتم سلمانی و همکاران [۲۷] با این فرض آغاز می‌شود که هیچ متغیر تبیینی‌ای در مدل حضور نداشته باشد و فقط عرض از مبدا وجود دارد. اولین متغیر پیشگوی فازی که وارد مدل می‌شود، متغیری است که بیشترین تغییرات مدل را توصیف کرده باشد. در گام بعدی، شرایط هر متغیر برای وارد شدن به مدل، سطح کارایی حداقل LE است. گام‌های اساسی این شیوه عبارتند از: (شکل ۲.۳)

۱. مدل‌های تک متغیره برازش داده می‌شوند و از بین آنها بهترین مدل انتخاب می‌شود.
۲. متغیر بعدی وارد مدل می‌شود و MSE محاسبه می‌شوند. اگر درصد بهبود مدل بیشتر از LE درصد باشد، متغیر در مدل می‌ماند.
۳. این روند تا جایی که افزودن متغیر به مدل باعث تغییری در شاخص نشود، ادامه پیدا می‌کند.

۴.۲.۳ کاربرد رگرسیون لجستیک فازی در بیماری لوپوس

بیماری لوپوس ارتیماتوز سیستمیک^۷ که در علم پزشکی به اختصار SLE نامیده می‌شود، گونه‌ای بیماری خود ایمنی است که در آن سیستم دفاعی بدن، علیه ارگان‌ها و بافت‌های پیوندی خودی عمل کرده و به آنها آسیب می‌رساند. زنان در مقابل این بیماری ده برابر بیشتر از مردان آسیب پذیرند. این بیماری اغلب پوست و چندین اندام داخلی را درگیر می‌کند و با وجود اتوانتی بادی در خون همراه است. سیر بالینی لوپوس شامل دوره‌های فعالیت و بهبود است. علت بیماری، نقص ایمنی پوست و ارگان‌های داخلی است که منشأ آن تاکنون شناخته شده نیست. نوع حاد آن به ندرت به مرگ منجر می‌شود و بیشتر در زنان و درده سه سوم زندگی بروز می‌کند. [۱۸] تشخیص بیماری براساس علائم آن صورت می‌گیرد. تشخیص در مراحل اولیه باعث تسریع در درمان و مانع از پیشرفت بیماری می‌شود.

داده‌ها متشکل از ۱۵ نفر زن مشکوک به بیماری لوپوس در بازه سنی ۴۰ - ۱۸ سال از درمانگاه‌های شیراز جمع‌آوری شده است. در این تحقیق می‌خواهیم وضعیت افراد مشکوک (سالم/بیمار) به بیماری لوپوس را براساس تعدادی عوامل خطر، مدل‌سازی کنیم، به گونه‌ای که مدل برازش داده شده قادر باشد بخت امکانی ابتلا به بیماری در هر فرد را برآورد کند. عوامل خطری که در این مطالعه استفاده شده است شامل تست‌های آزمایشگاهی Anti، ESR، DNA، ANA، مواجه فرد با نور خورشید و سابقه خانوادگی فرد است. آزمایش‌هایی که نام برده شد، آزمایش‌های ویژه خونی هستند.

^۷Systemic Lupus Erythematosus

با توجه به ترم‌های زبانی ((خیلی کم، کم، متوسط، زیاد، خیلی زیاد)) وضعیت عمومی هر یک از ۱۵ فرد توسط پزشک بیان شده است. در همین راستا توابع عضویت هر ترم زبانی به صورت زیر تعریف شدند.

$$VeryLow = \begin{cases} 1 - \frac{0.02 - x}{0.01} & 0.01 \leq x \leq 0.02 \\ 1 - \frac{x - 0.02}{0.16} & 0.02 \leq x \leq 0.18 \end{cases} \quad Low = \begin{cases} 1 - \frac{0.25 - x}{0.15} & 0.1 \leq x \leq 0.25 \\ 1 - \frac{x - 0.02}{0.15} & 0.25 \leq x \leq 0.4 \end{cases}$$

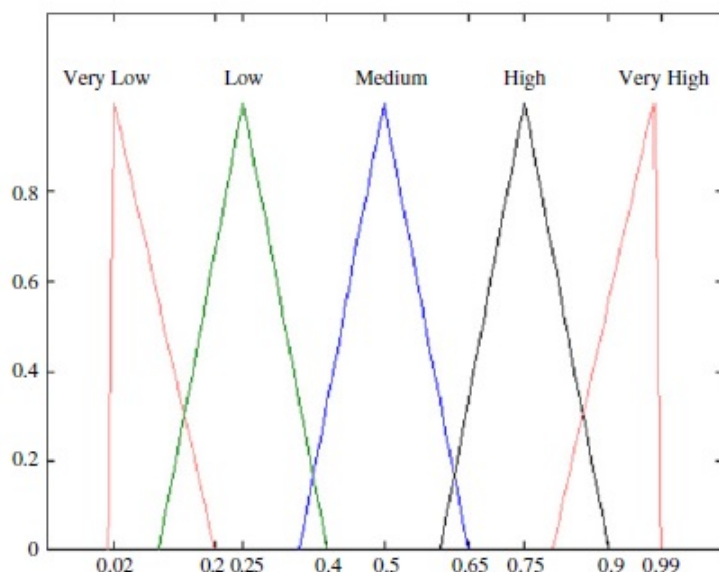
$$Medium = \begin{cases} 1 - \frac{0.05 - x}{0.15} & 0.35 \leq x \leq 0.5 \\ 1 - \frac{x - 0.5}{0.15} & 0.5 \leq x \leq 0.65 \end{cases} \quad High = \begin{cases} 1 - \frac{0.75 - x}{0.15} & 0.6 \leq x \leq 0.75 \\ 1 - \frac{x - 0.75}{0.15} & 0.75 \leq x \leq 0.9 \end{cases}$$

$$VeryHigh = \begin{cases} 1 - \frac{0.98 - x}{0.18} & 0.8 \leq x \leq 0.98 \\ 1 - \frac{x - 0.98}{0.01} & 0.98 \leq x \leq 0.99 \end{cases} \quad (12.3)$$

داده‌های این تحقیق در جدول ۱.۳ آمده است. مدل کلی با شرکت همه متغیرهای پیش‌بینی و پاسخ به صورت زیر است.

$$\hat{Y}_i = \tilde{b}_0 + \tilde{b}_1 ESRtest_i + \tilde{b}_2 AntiDNAtest_i + \tilde{b}_3 ANAtest_i + \tilde{b}_4 SunExposure_i + \tilde{b}_5 FamilyHistory_i \quad (13.3)$$

همچنین مقادیر مشاهده شده \tilde{y}_i را با استفاده از توابع عضویت و تبدیل لجیت $\tilde{y}_i = \ln \frac{\mu_i}{1-\mu_i}$ به دست آوردیم. برای این کار از نرم افزار R و از بسته FuzzyNumber کمک گرفتیم. کدها مربوط به R در پیوست ۱.ب و مقادیر به دست آمده در جدول ۲.۳ آورده شد. با روش پیشرو و به کمک الگوریتمی که توضیح دادیم با معیار LE بهترین مدل را انتخاب می‌کنیم.



شکل ۱.۳: نمودار پیشنهادی μ_i

جدول ۱.۳: موارد مشکوک به لوپوس و عوامل خطر مربوط به آن

امکان μ_i	سابقه خانوادگی	مواجه با نور خورشید	ANA test	AntiDNA test	ESR test	شماره فرد
high	۱	۱	۱۱۲	۱۰۵	۱	۱
medium	۰	۱	۸۰	۲۳	۰	۲
high	۰	۱	۱۱۵	۱۵	۰	۳
high	۰	۱	۱۰۵	۱۰۷	۱	۴
medium	۰	۰	۸۹	۱۵۰	۱	۵
veryhigh	۱	۱	۱۶۰	۱۰	۱	۶
medium	۰	۱	۱۰۰	۲۳	۰	۷
high	۰	۰	۱۰۰	۸۵	۱	۸
low	۰	۱	۴۸	۸۳	۰	۹
low very	۱	۰	۱۵	۱۹	۱	۱۰
low	۰	۰	۵۰	۹۱	۰	۱۱
medium	۰	۱	۵۹	۲۰۰	۱	۱۲
low	۰	۱	۸۳	۲۰	۱	۱۳
low	۰	۰	۱۵	۲۰۰	۰	۱۴
medium	۱	۰	۸۵	۱۵	۱	۱۵

جدول ۲.۳: مقادیر مشاهده شده \tilde{y}_i

ردیف	مرکز	پهنای چپ	پهنای راست
۱	۱.۰۶۳۷۲۱	۰.۶۵۱۷۵۲۱	۱.۰۳۵۹۳۴۵
۲	۰.۰۰۰۰۰۰	۰.۶۰۹۳۴۱۵	۰.۶۰۹۳۴۱۵
۳	۱.۰۶۳۷۲۱	۰.۶۵۱۷۵۲۱	۱.۰۳۵۹۳۴۵
۴	۱.۰۶۳۷۲۱	۰.۶۵۱۷۵۲۱	۱.۰۳۵۹۳۴۵
۵	۰.۰۰۰۰۰۰	۰.۶۰۹۳۴۱۵	۰.۶۰۹۳۴۱۵
۶	۳.۳۱۷۴۰۳	۲.۱۶۳۹۷۰۲	۱.۲۳۲۵۵۵۲
۷	۰.۰۰۰۰۰۰	۰.۶۰۹۳۴۱۵	۰.۶۰۹۳۴۱۵
۸	۱.۰۶۳۷۲۱	۰.۶۵۱۷۵۲۱	۱.۰۳۵۹۳۴۵
۹	-۱/۰۶۳۷۲۱	۱.۰۳۵۹۳۴۵	۰.۶۵۱۷۵۲۱
۱۰	-۳/۸۵۷۵۶۶	۰.۶۹۲۳۹۱۷	۲.۵۶۱۲۲۹۳
۱۱	-۱/۰۶۳۷۲۱	۱.۰۳۵۹۳۴۵	۰.۶۵۱۷۵۲۱
۱۲	۰.۰۰۰۰۰۰	۰.۶۰۹۳۴۱۵	۰.۶۰۹۳۴۱۵
۱۳	-۱/۰۶۳۷۲۱	۱.۰۳۵۹۳۴۵	۰.۶۵۱۷۵۲۱
۱۴	-۱/۰۶۳۷۲۱	۱.۰۳۵۹۳۴۵	۰.۶۵۱۷۵۲۱
۱۵	۰.۰۰۰۰۰۰	۰.۶۰۹۳۴۱۵	۰.۶۰۹۳۴۱۵

جدول ۳.۳: مرحله اول انتخاب بهترین مدل

شماره	مدل	MSE
۱	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ESRtest_i$	۹۵.۴۹۰۳
۲	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 AntiDNAtest_i$	۹۸.۷۷۴۹
۳	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i$	۱۸.۸۱۵۶
۴	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 SunExposure_i$	۸۲.۲۱۹
۵	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 FamilyHistory_i$	۹۸.۲۱۹۳

انتخاب بهترین مدل

مرحله ۱. طبق الگوریتمی که گفته شد، برای تمامی مدل‌ها با یک متغیر، MSE آنها را حساب می‌کنیم. همه مقادیر محاسبه شده را در جدول ۳.۳ آورده‌ایم. بهترین مدل در این مرحله

جدول ۴.۳: مرحله دوم انتخاب بهترین مدل

شماره	مدل	MSE
۱	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i + \tilde{b}_2 ESRtest_i$	۱۸.۳۰۵۵
۲	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i + \tilde{b}_2 AntiDNAtest_i$	۹.۵۹۵۹
۳	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i + \tilde{b}_2 SunExposure_i$	۱۸.۷۶۷۱
۴	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i + \tilde{b}_2 FamilyHistory_i$	۱۷.۵۷۰۰

$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i$ است که کمترین MSE با مقدار ۱۸/۸۱۵۶ را دارد.

مرحله ۲. در مرحله بعد با حفظ متغیر $ANAtest$ در مدل، متغیرهای دیگر را یک به یک وارد مدل کردیم. با توجه به مقدار MSE، بهترین مدل $\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i + \tilde{b}_2 AntiDNAtest_i$ است. اطلاعات مربوط به این مرحله در جدول ۴.۳ آمده است.

مرحله ۳. این کار را تا جایی ادامه دادیم که با اضافه شدن متغیر جدید در مدل، هیچ بهبود در آن ایجاد نشود. اطلاعات مربوط به این مراحل در جداول ۵.۳ و ۶.۳ آمده است. در نهایت بهترین مدلی که با معیار $LE = 0.7$ انتخاب می‌کنیم مدل زیر است.

$$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i + \tilde{b}_2 AntiDNAtest_i + \tilde{b}_3 ESRtest_i + \tilde{b}_4 FamilyHistory_i \quad (14.3)$$

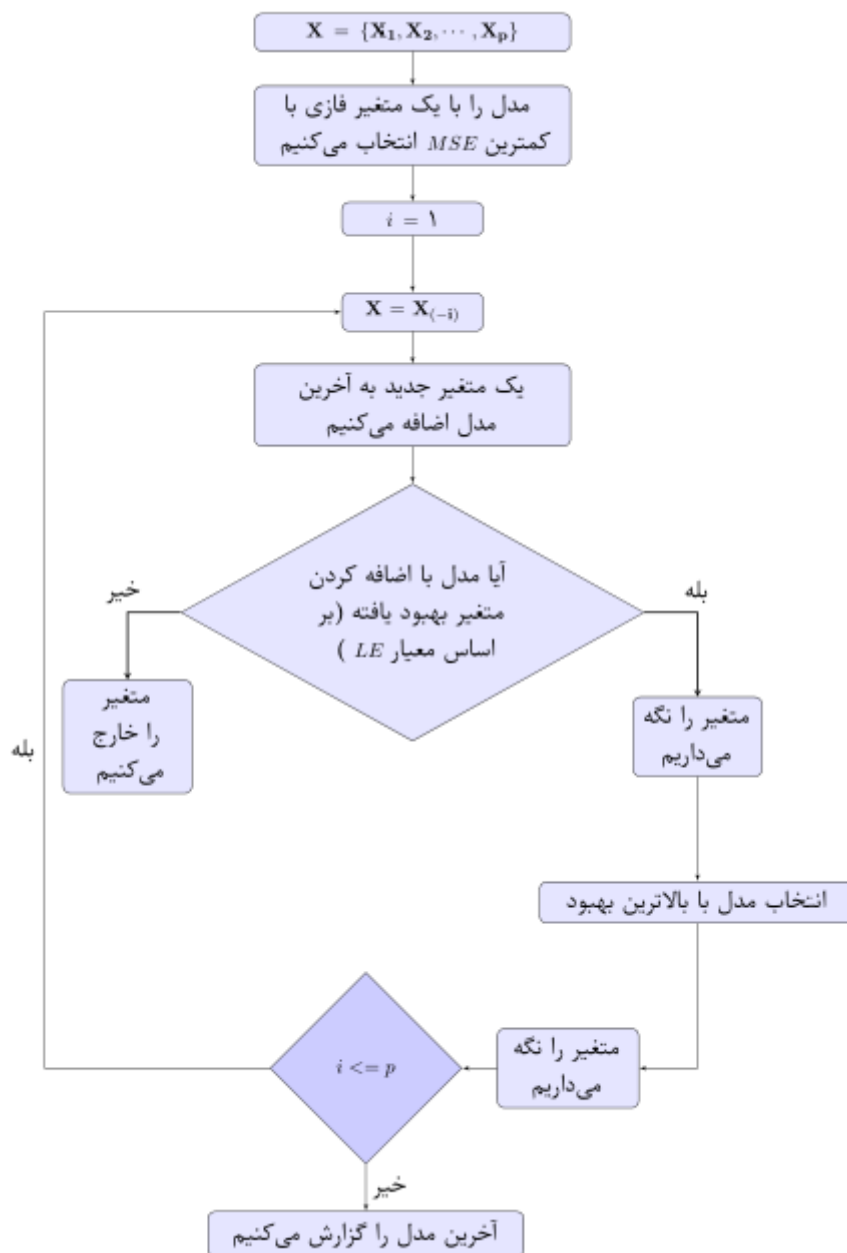
مدل (۱۴.۳) نشان می‌دهد که تست‌های آزمایشگاهی ANA ، $AntiDNA$ و ESR و سابقه خانوادگی عوامل تاثیر گذار در امکان ابتلای افراد به بیماری لوپوس را دارند. بنابراین \tilde{Y}_i لگاریتم بخت امکانی بیماری برای فرد i ام است. با استفاده از اصل گسترش و با استفاده از تابع آنتی لگاریتم $f(\tilde{Y}_i) = e^{(\tilde{Y}_i)}$ امکان ابتلای هر فرد به بیماری را می‌توان محاسبه کرد.

جدول ۵.۳: مرحله سوم انتخاب بهترین مدل

شماره	مدل	MSE
۱	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i + \tilde{b}_2 AntiDNAtest_i + \tilde{b}_3 ESRtest_i$	۷.۹۹۶۳
۲	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i + \tilde{b}_2 AntiDNAtest_i + \tilde{b}_3 SunExposure_i$	۹.۵۸۱۹
۳	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i + \tilde{b}_2 AntiDNAtest_i + \tilde{b}_3 FamilyHistory_i$	۹.۳۰۴۹

جدول ۶.۳: مرحله چهارم انتخاب بهترین مدل

شماره	مدل	MSE
۱	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i + \tilde{b}_2 AntiDNAtest_i + \tilde{b}_3 ESRtest_i + \tilde{b}_4 SunExposure_i$	۹.۵۸۱۹
۲	$\tilde{Y}_i = \tilde{b}_0 + \tilde{b}_1 ANAtest_i + \tilde{b}_2 AntiDNAtest_i + \tilde{b}_3 ESRtest_i + \tilde{b}_4 FamilyHistory_i$	۹.۳۰۴۹



شکل ۲.۳: الگوریتم انتخاب متغیر سلمانی در مدل رگرسیون لجستیک فازی

فصل ۴

رگرسیون لجستیک فازی با رویکرد کمترین توان دوم

در این فصل برای برآورد پارامترهای رگرسیون لجستیک با رویکرد کمترین توان دوم، از دو متر، یکی متر ژو [۳۵] با استفاده از مقاله پورا احمد و همکاران [۳۲] و متر دورسو [۱۱] استفاده می‌کنیم. سپس به کمک مثال عددی و با استفاده از نرم افزار R و Mathematica مطلب را کامل تر توضیح می‌دهیم.

۱.۴ روش مجموع توان دوم برای داده‌های فازی

رویکرد کمترین توان دوم فازی که اولین بار توسط دیاموند [۱۲] و کلمینس [۱۰] به طور هم‌زمان پیشنهاد شد، این روش در واقع یک تعمیم فازی از روش کمترین توان دوم با داده‌های دقیق است و مستلزم یک مفهوم فاصله بین مقادیر فازی مشاهده شده و مقادیر فازی برآورد شده است. همان‌طور که بیان شد در روش مجموع توان دوم، فاصله بین مقادیر مشاهده شده و مقادیر برآورد شده فازی را به حداقل می‌رسانیم که مزیت این روش است، زیرا باقیمانده‌ها برای سنجش دقت مدل مناسب هستند.

روش کمترین توان دوم خطی برای m متغیر دقیق (x_i, y_j) ، $1 \leq m \leq i$ برای پیدا کردن

مقادیر مجهول $a, b \in \mathbb{R}$ با حداقل کردن مجموع زیر است.

$$r(a, b) = \sum_i (a + bx_i + y_i)^2. \quad (1.4)$$

در مواقعی که داده‌ها فازی هستند ما برای اعداد فازی مانند \tilde{A}_1 و \tilde{A}_p به دنبال آن هستیم تا فاصله بین این دو عدد را حداقل کنیم.

۱.۱.۴ برآورد پارامترها با متر ژو

برای تعریف بهترین مدل بهینه و برای برآورد پارامترهای مدل از روش کمترین توان دوم خطا استفاده می‌کنیم. در این روش، مجموع توان دوم خطای بین \tilde{Y}_i و \tilde{y}_i را باید حداقل کرد. یعنی مجموع توان دوم خطا بین مقدار مشاهده شده و مقادیر برآورد شده را حداقل می‌کنیم. برای این منظور SSE به صورت زیر تعریف می‌شود

$$SSE = \sum_{i=1}^n d^2(\tilde{Y}_i, \tilde{y}_i) \quad (2.4)$$

در این تعریف $d^2(\tilde{Y}_i, \tilde{y}_i)$ فاصله بین دو مجموعه فازی است. یک تابع فاصله مناسب، تابع وزنی است که توسط ژو [۳۵] در رابطه (۴.۲) معرفی شده است. همان طور که گفته شد مدل پیشنهادی رگرسیون لجستیک فازی

$$\tilde{Y}_i = \tilde{A}_0 + \tilde{A}_1 x_{i1} + \dots + \tilde{A}_{p-1} x_{ip-1} \quad (3.4)$$

می‌باشد، که در آن متغیر پاسخ فازی و متغیر تبیینی دقیق است و همچنین ضرایب مدل فازی هستند. گفته شد که برای برآورد ضرائب باید مجموع توان دوم خطاها (۲.۴) حداقل شود. حال فرض کنید $\tilde{A}_j = (a_j, s_j^L, s_j^R)_T, j = 1, \dots, p-1$ ضرائب مدل (۳.۴) باشد، در نتیجه پاسخ برآورد شده، مقدار عدد فازی مثلثی نامتقارن $\tilde{Y}_i = (f_i(a), f_i^L(s), f_i^R(s))_T$ است که در آن مرکز، پهنای چپ و پهنای راست به ترتیب

$$\begin{aligned} f_i^a(x) &= a_0^c + a_1^c x_{i1} + \dots + a_{p-1}^c x_{ip-1} \\ f_i^L(s) &= s_0^L + s_1^L x_{i1} + \dots + s_{p-1}^L x_{ip-1} \\ f_i^R(s) &= s_0^R + s_1^R x_{i1} + \dots + s_{p-1}^R x_{ip-1} \end{aligned} \quad (4.4)$$

می‌باشند.

هدف ما به دست آوردن فاصله بین $(\tilde{Y}_i)_\alpha$ و $(\tilde{y}_i)_\alpha$ ، یعنی فاصله بین مقادیر برآورد شده و مقادیر مشاهده شده و حداقل سازی بین این دو فاصله است. یعنی در نهایت مقدار زیر را به دست آورده تا بتوان پارامترهای مجهول مدل را برآورد کرد. در رابطه (۴.۲):

$$d(\tilde{y}_i, \tilde{Y}_i) = \left[\int_0^1 f(\alpha) d^2((\tilde{Y}_i)_\alpha, (\tilde{y}_i)_\alpha) d\alpha \right]^{\frac{1}{\alpha}} \quad (5.4)$$

برای حالتی که ضرایب اعداد فازی مثلثی نامتقارن هستند، مقدار α -برش $(\tilde{Y}_i)_\alpha$ به صورت

$$(\tilde{Y}_i)_\alpha = [f_i(\mathbf{a}) - (1 - \alpha)f_i^L(s), f_i(\mathbf{a}) + (1 - \alpha)f_i^R(s)] \quad (6.4)$$

به دست می‌آید. همچنین α -برش $(\tilde{y}_i)_\alpha$ براساس $[b_1, b_2]$ به صورت زیر:

$$(\tilde{y}_i)_\alpha = \left[\ln \frac{b_1}{1 - b_1}, \ln \frac{b_2}{1 - b_2} \right]$$

می‌باشد. مقدار b_1 و b_2 از α -برش توابع عضویت براساس توابع عضویت بیان شده در (۱۳.۳) به صورت زیر به دست می‌آید:

$$(\tilde{\mu}_i)_\alpha = [f_i^L(s) + (\alpha - 1)f_i(\mathbf{a}), f_i^R(s) + (1 - \alpha)f_i(\mathbf{a})]$$

بنابراین فاصله بین $(\tilde{Y}_i)_\alpha$ و $(\tilde{y}_i)_\alpha$ بر اساس رابطه (۵.۲) به صورت

$$(7.4)$$

$$d^{\chi}((\tilde{y}_i)_\alpha, (\tilde{Y}_i)_\alpha) = \left[\ln \frac{b_1}{1 - b_1} - f_i(\mathbf{a}) + (1 - \alpha)f_i^L(s) \right]^2 + \left[\ln \frac{b_2}{1 - b_2} - f_i(\mathbf{a}) - (1 - \alpha)f_i^R(s) \right]^2$$

خواهد بود. بنابراین داریم:

$$(8.4)$$

$$SSE = \sum_{i=1}^n \int_0^1 f(\alpha) \left[\left[\ln \frac{b_1}{1 - b_1} - f_i(\mathbf{a}) + (1 - \alpha)f_i^L(s) \right]^2 + \left[\ln \frac{b_2}{1 - b_2} - f_i(\mathbf{a}) - (1 - \alpha)f_i^R(s) \right]^2 \right] d\alpha$$

این تابع از طریق $f_i(\mathbf{a})$ ، $f_i^L(s)$ و $f_i^R(s)$ به ضرایب مدل وابسته است. با حداقل کردن (۸.۴) مقادیر مجهول a_j ، s_j^L و s_j^R با شرط مثبت بودن پهناهای چپ و پهناهای راست، یعنی $s_j^L > 0$ و $s_j^R > 0$ به دست می‌آید.

برای حالتی که $s_j = s_j^L = s_j^R$ باشد، پوراحمد و همکاران [۳۲]

برای به دست آوردن پارامترها از SSE نسبت به a_j و s_j مشتق گرفته شده است.

$$(9.4)$$

$$SSE = \sum_{i=1}^n \int_0^1 f(\alpha) \left[\left[\ln \frac{b_1}{1 - b_1} - f_i(\mathbf{a}) + (1 - \alpha)f_i(s) \right]^2 + \left[\ln \frac{b_2}{1 - b_2} - f_i(\mathbf{a}) - (1 - \alpha)f_i(s) \right]^2 \right] d\alpha.$$

با در نظر گرفتن $f(\alpha) = \alpha$ و

$$\frac{\partial SSE}{\partial a_j} = 0 \quad \frac{\partial SSE}{\partial s_j} = 0$$

داریم:

$$\sum_{i=1}^n \left(\int_0^1 \alpha x_{ij} \left[\alpha f_i(\mathbf{a}) - \ln \frac{b_1}{1 - b_1} - \ln \frac{b_2}{1 - b_2} \right] d\alpha \right) = 0 \quad j = 0, 1, \dots, p-1$$

$$\sum_{i=1}^n \left(\int_0^1 \alpha(1 - \alpha) x_{ij} \left[\alpha(1 - \alpha) f_i(\mathbf{a}) - \ln \frac{b_1}{1 - b_1} - \ln \frac{b_2}{1 - b_2} \right] d\alpha \right) = 0 \quad j = 0, 1, \dots, p-1$$

مقدار معادلات فوق به تعریف امکان موفقیت (μ_i) برای هر فرد توسط خیره بستگی دارد. با جایگذاری این مقادیر و محاسبه انتگرال مربوط، معادلات زیر حاصل می‌شوند.

$$a_0 \sum_{i=1}^n x_{i0} x_{ij} + a_1 \sum_{i=1}^n x_{i1} x_{ij} + \dots + a_{p-1} \sum_{i=1}^n x_{i,p-1} x_{ij} = \sum_{i=1}^n z_{ij} x_{ij}$$

$$s_0 \sum_{i=1}^n x_{i0} x_{ij} + s_1 \sum_{i=1}^n x_{i1} x_{ij} + \dots + s_{p-1} \sum_{i=1}^n x_{i,p-1} x_{ij} = \sum_{i=1}^n k_{ij} x_{ij}$$

اگر $x_{i0} = 1$ و z_{ij} و k_{ij} نتایج محاسبه انتگرال برای هر مورد باشد، بنابراین می‌توانیم معادلات بالا را به صورت ساده‌تر نوشت.

$$Aa = Z \quad , \quad As = K$$

که در آن

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1,p-1} \\ 1 & x_{21} & \dots & x_{2,p-1} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{n,p-1} \end{bmatrix}, \quad A = X^T X$$

$$\mathbf{a} = (a_0, a_1, \dots, a_{p-1})^T \quad , \quad \mathbf{s} = (s_0, s_1, \dots, s_{p-1})^T$$

$$Z = \left(\sum_{i=1}^n z_i x_{i0}, \sum_{i=1}^n z_i x_{i1}, \dots, \sum_{i=1}^n z_i x_{i,p-1} \right)^T$$

$$K = \left(\sum_{i=1}^n k_i x_{i0}, \sum_{i=1}^n k_i x_{i1}, \dots, \sum_{i=1}^n k_i x_{i,p-1} \right)^T$$

هستند. در نهایت مقادیر مجهول به صورت $\mathbf{a} = A^{-1}Z$ و $\mathbf{s} = A^{-1}K$ به دست می‌آید. همان طور که بیان شده، در این پایان‌نامه، داده‌ها به صورت اعداد فازی مثلثی نامتقارن مفروضند. برای برآورد پارامترهای مجهول بایستی (۸.۴) را حداقل کنیم، یعنی نسبت به a_j ، مشتق گرفته و مساوی صفر قرار دهیم. در این صورت با عملیات ریاضی دشواری مواجه می‌شویم. برای حداقل کردن مجموع مربعات از (۱۰.۲) استفاده می‌نماییم. در ادامه با ذکر مثال، به شرح آن می‌پردازیم.

۲.۱.۴ مثال کاربردی

داده‌های بیماری لوپوس در بخش ۴.۲.۳ را در نظر بگیرید. برای برآورد پارامترهای مجهول مدل (۳.۴) در حالت نامتقارن بودن اعداد فازی مثلثی، باید مجموع مربعات (۸.۴) را حداقل

کنیم. برای این منظور از نرم افزار Mathematica استفاده کردیم، که کدهای آن در پیوست ب.۲ آورده شده است. مدل برازش شده به صورت زیر به دست می‌آید.

$$\begin{aligned} \hat{Y}_i = \ln \left(\frac{\tilde{\mu}}{1 - \tilde{\mu}} \right) &= (0.474, 0.996, 1.113) + (0.348, 1.099, 1.207)ESRtest_i \\ &+ (0.559, 1.324, 0.445)AntiDNAtest_i + (0.517, 1.488, 1.544)ANAtest_i \\ &+ (-0.545, 0.807, 0.666)SunExposure_i + (0.324, 1.178, 0.466)FamilyHistory_i \end{aligned} \quad (10.4)$$

با استفاده از مدل (۱۰.۴) می‌توان امکان وجود بیماری در یک فرد را پیش بینی کرد. برای مثال با توجه به داده‌های جدول ۱.۳ برای فرد هشتم، از رابطه (۱۰.۴) داریم:

$$\begin{aligned} \hat{Y}_8 = \ln \left(\frac{\tilde{\mu}_8}{1 - \tilde{\mu}_8} \right) &= (0.474, 0.996, 1.113) + (0.348, 1.099, 1.207) \times 1 \\ &+ (0.559, 1.324, 0.445) \times 85 + (0.517, 1.488, 1.544) \times 100 \\ &+ (-0.545, 0.807, 0.666) \times 0 + (0.324, 1.178, 0.466) \times 0 \\ &= (100.037, 263.435, 194.545) \end{aligned}$$

بنابراین، لگاریتم بخت امکانی بیماری برای فرد هشتم، 100.037 برآورد می‌شود و با استفاده از اصل گسترش و با تبدیل $f(x) = e^x$ ، امکان بیماری برای فرد به صورت زیر محاسبه می‌شود.

$$\begin{aligned} \left(\frac{\tilde{\mu}_8}{1 - \tilde{\mu}_8} \right) (x) &= \exp(\hat{Y}_8(x)) = \hat{Y}_8(\ln x), \quad x > 0 \\ &= \begin{cases} 1 - \frac{100.037 - \ln x}{263.435} & 163.398 \leq \ln x \leq 100.037 \\ 1 - \frac{\ln x - 100.037}{194.545} & 100.037 \leq \ln x \leq 294.582 \end{cases} \end{aligned}$$

از مدل (۱۰.۴) برای پیش بینی امکان وجود بیماری در افراد جدید نیز می‌توان استفاده کرد.

مدل برازش شده در پور احمد و همکاران [۳۲] با فرض متقارن بودن اعداد فازی مثلثی ضرایب و پاسخ برآورد شده به صورت

$$\begin{aligned} \hat{Y}_i &= (-3.8591, 1.1617) + (-0.6083, 0.1451)ESRtest_i + 0.0091AntiDNAtest_i \\ &+ 0.0431ANAtest_i + (-0.1309, 0.366)SunExposure_i \\ &+ (0.4248, 0.3832)FamilyHistory_i \end{aligned} \quad (11.4)$$

است. به طور مثال برای فرد سوم، از رابطه (۱۰.۴) داریم:

$$\begin{aligned}\hat{Y}_3 &= \ln\left(\frac{\tilde{\mu}_3}{1-\tilde{\mu}_3}\right) = (-3/8591, 1/1617) + (-0/6083, 0/1451) \times 0 \\ &+ 0/0091 \times 15 + 0/0431 \times 115 \\ &+ (-0/1309, 0/0366) \times 1 + (0/4248, 0/3832) \times 0 \\ &= (1/1030, 1/1913) = (1/10, 1/20)\end{aligned}$$

بنابراین، لگاریتم بخت امکانی بیماری برای فرد سوم، با استفاده از اصل گسترش، امکان بیماری برای فرد به صورت زیر محاسبه می‌شود.

$$\begin{aligned}\left(\frac{\tilde{\mu}_3}{1-\tilde{\mu}_3}\right)(x) &= \exp(\tilde{Y}_3(x)) = \tilde{Y}_3(\ln x), \quad x > 0 \\ &= \begin{cases} 1 - \frac{1/10 - \ln x}{1/20} & -0/24 \leq \ln x \leq 1/10 \\ & (0/79 \leq x \leq 3/00) \\ 1 - \frac{\ln x - 1/10}{1/20} & 1/10 \leq \ln x \leq 2/45 \\ & (3/00 \leq x \leq 11/59) \end{cases}\end{aligned}$$

حال از مدل (۱۱.۴) برای پیش بینی امکان وجود بیماری در فرد جدید استفاده می‌کنیم. برای مثال فرض می‌کنیم برای یک فرد جدید مقادیر متغیرهای تبیینی، به صورت زیر مشاهده شده باشند.

$$ESR_{test_{new}} = 0, AntiDNA_{test_{new}} = 20, ANA_{test_{new}} = 100$$

$$SunExposure_{new} = 1, FamilyHistory_{new} = 1$$

در این صورت با استفاده از مدل (۱۱.۴) داریم:

$$\begin{aligned}\hat{Y}_{new} &= \ln\left(\frac{\tilde{\mu}_{new}}{1-\tilde{\mu}_{new}}\right) = (-3/8591, 1/1617) + (-0/6083, 0/1451) \times 0 \\ &+ 0/0091 \times 20 + 0/0431 \times 100 \\ &+ (-0/1309, 0/0366) \times 1 + (0/4248, 0/3832) \times 1 \\ &= (0/9268, 1/5815)\end{aligned}$$

بنابراین، لگاریتم بخت امکانی بیماری برای این فرد، حدوداً ۰/۹۳ برآورد می‌شود. تابع عضویت

برای فرد جدید با توجه به اصل گسترش به صورت زیر است:

$$\left(\frac{\tilde{\mu}_{new}}{1 - \tilde{\mu}_{new}}\right)(x) = \exp(\tilde{Y}_{new}(x)) = \tilde{Y}_{new}(\ln x), \quad x > 0$$

$$= \begin{cases} 1 - \frac{0.9268 - \ln x}{1.5815} & -0.6547 \leq \ln x \leq 0.9268 \\ 0.5196 \leq x \leq 2.5264 \\ 1 - \frac{\ln x - 0.9268}{1.5815} & 0.9268 \leq \ln x \leq 2.5083 \\ 2.5264 \leq x \leq 12.2840 \end{cases}$$

در ادامه این فصل به برآورد پارامترهای رگرسیون لجستیک در محیط فازی به روشی دیگر می‌پردازیم. روش ارائه شده در فصل قسمت برای برآورد پارامترهای رگرسیون لجستیک فازی با استفاده از فاصله ژو [۳۵] بود، در این فصل می‌خواهیم پارامترهای مدل لجستیک فازی را با استفاده از فاصله تعریف شده به شیوه کپی و دورسو [۱۱] برآورد کنیم.

۲.۴ روش پیشنهادی برای برآورد پارامترها

برای برآورد پارامترهای مدل از روش کمترین توان‌های دوم به شیوه کپی و همکاران [۱۱] استفاده می‌کنیم. [۲] در این روش ملاک برآورد پارامترها، مینیمم سازی مجموع توان دوم فاصله بین توابع عضویت ترم‌های زبانی مشاهده شده و توابع عضویت ترم‌های زبانی پیش بینی شده توسط مدل است، یعنی

$$SSE = \sum_{i=1}^n d^2(\tilde{Y}_i, \tilde{y}_i)$$

که در آن مقدار تبدیل یافته (رجوع کنید به بخش ۱.۱.۲) مشاهده ترم‌های زبانی مربوط به هر ویژگی و \tilde{Y}_i مقدار برآورد شده متناظر با \tilde{y}_i است که در (۳.۴) تعریف شده است [۱]. همچنین $d^2(\tilde{y}_i, \tilde{Y}_i)$ فاصله تعریف شده به شیوه کپی و دورسو [۱۱] است. ایده اساسی در اینجا مدل سازی مراکز متغیر پاسخ فازی با استفاده از یک مدل رگرسیون خطی است. این ایده را می‌توان تحت شرایط زیر بیان کرد. ابتدا مدل مراکز مشاهده شده و مرزهای پایین و بالای متغیرهای پاسخ را با استفاده از مجموعه‌های مربوط به باقی‌مانده‌ها و مقادیر دیگر تعریف می‌کنیم [۱۱].

$$\mathbf{m} = \mu + \epsilon$$

$$\mathbf{m} - \mathbf{l} = (\mu - \delta_L) + \epsilon_L$$

$$\mathbf{m} - \mathbf{u} = (\mu \delta_U) + \epsilon_U.$$

در فرمول‌های بالا ϵ ، ϵ_L و ϵ_U بردار خطاها و μ ، δ_L و δ_U بردار مقادیر مراکز و پهناهای متغیر پاسخ است. بنابراین بردارهای معرفی شده براساس مدل رگرسیون به صورت زیر بررسی

می‌شوند.

$$\begin{aligned}\mu &= X\beta \\ L &= \eta_L \mu + \xi_L \mathbf{1} \\ \delta_U &= \eta_U \mu + \xi_U \mathbf{1}.\end{aligned}\tag{۱۲.۴}$$

ماتریس X ، ماتریس داده‌های مشاهده شده x_{ij} ها است، پارامترهای β و $(\eta_L, \eta_U, \xi_L, \xi_U)$ به ترتیب ضرایب خطی بین مراکز و توابع مربوط به X_j و رابطه خطی بین پهنای مشاهده شده و مراکز برآورد شده هستند و 1 نشان دهنده برداری $(n \times 1)$ با مقدار تکراری یک است.

روشی که برای برآورد پارامترهای فوق استفاده می‌شود مبتنی بر الگوریتم استاندارد است که شامل مینیم کردن فاصله بین مقادیر مشاهده شده متغیر \tilde{y}_i و مقادیر برآورده شده \tilde{Y}_i که به صورت $\tilde{Y}_i = (\mu_i, \delta_{L_i}, \delta_{U_i}), i = 1, \dots, n$ تعریف می‌شود، مقادیر μ_i ، δ_{L_i} و δ_{U_i} در مدل (۱۲.۴) تعریف شده است.

حال با توجه به معیار کمترین توان دوم، پارامترهای مربوط به (۱۲.۴) باید با به حداقل رساندن فاصله توان دوم بین مقادیر مشاهده شده متغیر پاسخ \tilde{y}_i و مقادیر \tilde{Y}_i برآورد شوند. برای این منظور از فاصله کیپی و همکاران [۱۱] استفاده می‌کنیم.

$$\begin{aligned}d^\lambda(\tilde{y}_i, \tilde{Y}_i) &= d^\lambda((\mathbf{m}, \mathbf{1}, \mathbf{u})_{LR}, (\boldsymbol{\mu}, \boldsymbol{\delta}_L, \boldsymbol{\delta}_U)_{LR}) \\ &= \Delta_{LR}^\lambda = \|\mathbf{m} - \boldsymbol{\mu}\|^\lambda + \|(\mathbf{m} - \lambda \mathbf{1}) - (\boldsymbol{\mu} - \lambda \boldsymbol{\delta}_L)\|^\lambda + \|(\mathbf{m} + \rho \mathbf{u}) - (\boldsymbol{\mu} + \rho \boldsymbol{\delta}_U)\|^\lambda \\ &= \lambda (\mathbf{m} - \boldsymbol{\mu})^\top (\mathbf{m} - \boldsymbol{\mu}) - 2\lambda (\mathbf{m} - \boldsymbol{\mu})^\top (\mathbf{1} - \boldsymbol{\delta}_L) + \lambda (\mathbf{1} - \boldsymbol{\delta}_L)^\top (\mathbf{1} - \boldsymbol{\delta}_L) \\ &\quad + 2\rho (\mathbf{m} - \boldsymbol{\mu})^\top (\mathbf{u} - \boldsymbol{\delta}_U) + \rho (\mathbf{u} - \boldsymbol{\delta}_U)^\top (\mathbf{u} - \boldsymbol{\delta}_U)\end{aligned}\tag{۱۳.۴}$$

در این فاصله توابع وزن λ و ρ با توجه به تابع عضویت LR به صورت $\lambda = \int_0^1 L^{-1}(\omega) d\omega$ و $\rho = \int_0^1 U^{-1}(\omega) d\omega$ تعریف می‌شوند، مقادیر مربوط به این توابع وزن در تنظیم مناسب پهنای چپ و راست در هنگام محاسبه مرزهای پایین و بالا اعداد فازی \tilde{y}_i و \tilde{Y}_i کمک می‌کند [۱۱].

در فاصله (۱۳.۴) می‌توانیم تابع کمترین توان دوم را براساس پارامترهای مدل (۱۲.۴)

جایگذاری کنیم: ξ_U و $\xi_L, \eta_U, \eta_L, \beta$

$$\begin{aligned}
 & \min_{\beta, \eta_L, \eta_U, \xi_L, \xi_U} \Delta_{LR}^{\checkmark}(\beta, \eta_L, \eta_U, \xi_L, \xi_U) \\
 &= \checkmark(\mathbf{m} - \mathbf{X}\beta)^\top (\mathbf{m} - \mathbf{X}\beta) - \checkmark\lambda(\mathbf{m} - \mathbf{X}\beta)^\top (\mathbf{1} - \mathbf{X}\beta\eta_L - \mathbf{1}\xi_L) \\
 & \quad + \lambda^\checkmark(\mathbf{m} - \mathbf{X}\beta\eta_L - \mathbf{1}\xi_L)^\top (\mathbf{1} - \mathbf{X}\beta\eta_L - \mathbf{1}\xi_L) \\
 & \quad + \checkmark\rho(\mathbf{m} - \mathbf{X}\beta)^\top (\mathbf{u} - \mathbf{X}\beta\eta_U - \mathbf{1}\xi_U) + \rho^\checkmark(\mathbf{u} - \mathbf{X}\beta\eta_U - \mathbf{1}\xi_U)^\top (\mathbf{u} - \mathbf{X}\beta\eta_U - \mathbf{1}\xi_U) \\
 &= \checkmark(\mathbf{m}^\top \mathbf{m} - \checkmark\mathbf{m}^\top \mathbf{X}\beta + \beta^\top \mathbf{X}^\top \mathbf{X}\beta) - \checkmark\lambda(\mathbf{m}^\top \mathbf{1} - \mathbf{m}^\top \mathbf{X}\beta\eta_L - \mathbf{m}^\top \mathbf{1}\xi_L) \\
 & \quad + \beta^\top \mathbf{X}^\top \mathbf{1} + \beta^\top \mathbf{X}^\top \mathbf{X}\beta\eta_L + \beta^\top \mathbf{X}^\top \mathbf{1}\xi_L) \\
 & \quad + \lambda^\checkmark(\mathbf{1}^\top \mathbf{1} - \checkmark\mathbf{1}^\top \mathbf{X}\beta\eta_L - \checkmark\mathbf{1}^\top \mathbf{1}\xi_L + \beta^\top \mathbf{X}^\top \mathbf{X}\beta\eta_L^\checkmark + \checkmark\beta^\top \mathbf{X}^\top \mathbf{1}\eta_L\xi_L + n\xi_L^\checkmark) \\
 & \quad + \checkmark\rho(\mathbf{m}^\top \mathbf{u} - \mathbf{m}^\top \mathbf{X}\beta\eta_U - \mathbf{m}^\top \mathbf{1}\xi_U + \beta^\top \mathbf{X}^\top \mathbf{u} + \beta^\top \mathbf{X}^\top \mathbf{X}\beta\eta_U + \beta^\top \mathbf{X}^\top \mathbf{1}\xi_U) \\
 & \quad + \rho^\checkmark(\mathbf{u}^\top \mathbf{u} - \checkmark\mathbf{u}^\top \mathbf{X}\beta\eta_U - \checkmark\mathbf{u}^\top \mathbf{1}\xi_U + \beta^\top \mathbf{X}^\top \mathbf{X}\beta\eta_U^\checkmark + \checkmark\beta^\top \mathbf{X}^\top \mathbf{1}\eta_U\xi_U + n\xi_U^\checkmark). \quad (14.4)
 \end{aligned}$$

از تابع فوق نسبت به پارامترهای مجهول مشتق گرفته و مساوی صفر قرار می‌دهیم برآوردهای پارامترها به صورت زیر حاصل می‌شود.

$$\begin{aligned}
 \eta_L &= \lambda^{-1}(\beta^\top \mathbf{X}^\top \mathbf{X}\beta)^{-1} \left[\lambda(\beta^\top \mathbf{X}^\top \mathbf{1} - \beta^\top \mathbf{X}^\top \mathbf{1}\xi_L) - (\beta^\top \mathbf{X}^\top \mathbf{m} - \beta^\top \mathbf{X}^\top \mathbf{X}\beta) \right], \\
 \eta_U &= \rho^{-1}(\beta^\top \mathbf{X}^\top \mathbf{X}\beta)^{-1} \left[\rho(\beta^\top \mathbf{X}^\top \mathbf{u} - \beta^\top \mathbf{X}^\top \mathbf{1}\xi_U) - (\beta^\top \mathbf{X}^\top \mathbf{m} - \beta^\top \mathbf{X}^\top \mathbf{X}\beta) \right], \\
 \xi_L &= (n\lambda)^{-1} \left[\lambda\mathbf{1}^\top (\mathbf{1} - \mathbf{X}\beta\eta_L) - \mathbf{1}^\top (\mathbf{m} - \mathbf{X}\beta) \right], \\
 \xi_U &= (n\rho)^{-1} \left[\rho\mathbf{1}^\top (\mathbf{u} - \mathbf{X}\beta\eta_U) - \mathbf{1}^\top (\mathbf{m} - \mathbf{X}\beta) \right], \\
 \beta &= [\checkmark - \lambda\eta_L(\checkmark - \lambda\eta_L) + \rho\eta_U(\checkmark + \rho\eta_U)]^{-1} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \\
 & \quad \times \left[\checkmark\mathbf{m} - \lambda(\mathbf{m}\eta_L + \mathbf{1} + \mathbf{1}\xi_L) + \lambda^\checkmark(\mathbf{1}\eta_L - \mathbf{1}\eta_L\xi_L) \right. \\
 & \quad \left. + \rho(\mathbf{m}\eta_U + \mathbf{u} + \mathbf{1}\xi_U) + \rho^\checkmark(\mathbf{u}\eta_U - \mathbf{1}\eta_U\xi_U) \right].
 \end{aligned}$$

باید به این نکته توجه داشت که برای به دست آوردن برآوردهای فوق، پیشنهاد می‌کنیم از یک الگوریتم تکراری با استفاده از چندین نقطه شروع مختلف برای ارزیابی پایداری جواب‌ها، استفاده کنید.

۱.۲.۴ مثال کاربردی

در داده‌های بیماری لوپوس (رجوع شود به بخش ۴.۲.۳) برای بدست آوردن ضرایب مدل با استفاده از روش پیشنهادی از نرم افزار Mathematica کمک گرفتیم. کدهای مربوط به آن را در

پیوست ب.۳ آورده‌ایم. مدل به دست آمده با استفاده از آنها به صورت زیر است.

$$\hat{Y}_i = (-۳/۹۹۶, ۰/۶۹۳, ۰/۷۳) + (-۰/۶۴۲, ۰, ۰/۰۲۸)ESRtest + ۰/۰۰۹AntiDNAtest \\ + ۰/۰۴۴ANAtest + (۰/۱۲۴۰/۷۴۰, ۰)SunExposure \\ + (۰/۳۴۳, ۰/۲۵۸, ۰/۵۹۳)FamilyHistory$$

۲.۲.۴ معیار نیکویی برازش

در این بخش ضریب تعیین که یکی از معیارهای نیکویی برازش برای انتخاب متغیر است را تعریف می‌کنیم. برای این منظور ابتدا مفاهیم زیر را تعریف می‌کنیم.

تعریف ۱.۲.۴. برای متغیرهای خروجی LR می‌توان مفاهیم زیر را تعریف کرد:

◀ برای متغیر خروجی عدد فازی LR، مجموع مربعات کل از فاصله کپی و همکاران [۱۱] بین مقادیر مشاهده شده $\tilde{\mathbf{y}} = (\mathbf{m}, \mathbf{l}, \mathbf{u})_{LR}$ و میانگین مشاهدات $\bar{\mathbf{y}} = (\bar{m}, \bar{l}, \bar{u})_{LR}$ به دست آید، به عبارت دیگر

$$d^{\chi}(\tilde{\mathbf{y}}, \bar{\mathbf{y}}) = SST = \|\mathbf{m} - \bar{m}\|^2 + \|(\mathbf{m} - \lambda\mathbf{l}) - (\bar{m} - \lambda\bar{l})\|^2 + \|(\mathbf{m} + \rho\mathbf{u}) - (\bar{m} + \rho\bar{u})\|^2 \\ = \mathfrak{V}(\mathbf{m} - \bar{m})^{\top}(\mathbf{m} - \bar{m}) - 2\lambda(\mathbf{m} - \bar{m})^{\top}(\mathbf{l} - \bar{l}) + \lambda^2(\mathbf{l} - \bar{l})^{\top}(\mathbf{l} - \bar{l}) \\ + 2\rho(\mathbf{m} - \bar{m})^{\top}(\mathbf{u} - \bar{u}) + \rho^2(\mathbf{u} - \bar{u})^{\top}(\mathbf{u} - \bar{u})$$

که در آن \bar{m} ، \bar{l} و \bar{u} میانگین مقادیر بردارهای \mathbf{m} ، \mathbf{l} و \mathbf{u} هستند.

◀ مجموع مربعات رگرسیونی (SSR)، یعنی تغییراتی که توسط مدل مطرح شده اند، به صورت زیر تعریف می‌شوند.

$$d^{\chi}(\hat{\mathbf{y}}, \bar{\mathbf{y}}) = SSR = \|\hat{\mu} - \bar{m}\|^2 + \|(\hat{\mu} - \lambda\hat{\delta}_L) - (\bar{m} - \lambda\bar{l})\|^2 \\ + \|(\hat{\mu} + \rho\mathbf{u}) - (\bar{m} + \rho\hat{\delta}_U) - (\bar{m} + \rho\bar{u})\|^2.$$

◀ مجموع مربعات خطا (باقیمانده) (SSE)، یعنی تغییراتی که توسط مدل مطرح نشده اند، به صورت زیر تعریف می‌شوند.

$$d^{\chi}(\tilde{\mathbf{y}}, \hat{\mathbf{y}}) = SSE = \|\mathbf{m} - \hat{\mu}\|^2 + \|(\mathbf{m} - \lambda\mathbf{l}) - (\hat{\mu} - \lambda\hat{\delta}_L)\|^2 + \|(\mathbf{m} + \rho\mathbf{u}) - (\hat{\mu} + \rho\hat{\delta}_R)\|^2 \\ = \mathfrak{V}(\mathbf{m} - \hat{\mu})^{\top}(\mathbf{m} - \hat{\mu}) - 2\lambda(\mathbf{m} - \hat{\mu})^{\top}(\mathbf{l} - \hat{\delta}_L) + \lambda^2(\mathbf{l} - \hat{\delta}_L)^{\top}(\mathbf{l} - \hat{\delta}_L) \\ + 2\rho(\mathbf{m} - \hat{\mu})^{\top}(\mathbf{u} - \hat{\delta}_U) + \rho^2(\mathbf{u} - \hat{\delta}_U)^{\top}(\mathbf{u} - \hat{\delta}_U)$$

برای اندازه‌گیری نیکویی برازش مدل رگرسیون لجستیک با پاسخ فازی، از ضریب تعیین R^2 و ضریب تعیین تعدیل یافته R^2_{adj} استفاده می‌کنیم. برای این منظور به تعریف آنها می‌پردازیم.

تعریف ۲.۲.۴. ضریب تعیین چندگانه^۱ به صورت زیر تعریف می‌شود.

$$R^2 = 1 - \frac{SSE}{SST} \quad (15.4)$$

تعریف ۳.۲.۴. ضریب تعیین چندگانه تعمیم یافته^۲ R_{adj}^2 به صورت زیر تعریف می‌شود.

$$R_{adj}^2 = 1 - (1 - R^2) \frac{n-1}{n-p-4}. \quad (16.4)$$

این شاخص شامل یک عامل تعدیل یافته بر مبنای تعداد ضرایب رگرسیون در مدل مرکزی، دو ضریب در مدل پهنای چپ و دو ضریب در مدل پهنای راست است.

۳.۲.۴ روش پیشرو برای انتخاب متغیر به شیوه کپی و همکاران

با توجه به تعریف‌هایی که قبل انجام دادیم، می‌توانیم تخمین یک مدل رگرسیونی مناسب را انجام دهیم. برای انتخاب مدل مناسب از الگوریتم زیر استفاده می‌کنیم: (شکل ۱.۴)

۱. ابتدا مدل‌های تک متغیره می‌سازیم، از بین آنها مدلی که کمترین MPE را دارد انتخاب کرده و متغیر وارد شده در مدل را حفظ می‌کنیم.

۲. با حفظ متغیر در مدل، از ما بقی متغیرهای باقیمانده، متغیر جدیدی به مدل اضافه کرده و MPE هر یک از آنها را محاسبه می‌کنیم. مجدداً هر یک از مدل‌های که MPE کمتری داشت انتخاب می‌شود. به این ترتیب دو متغیر در مدل انتخاب می‌شود. این کار را تا جایی ادامه می‌دهیم که همه متغیرها وارد مدل شده باشند.

۳. در هر مرحله بعد از انتخاب مدل با کمترین MPE، مقدار R^2 و R_{adj}^2 را انتخاب می‌کنیم. زمانی که همه متغیرها وارد مدل شدند، از میان مدل‌های با کمترین MPE، مدلی را به عنوان مدل مناسب انتخاب می‌کنیم که بیشترین مقدار R_{adj}^2 را داشته باشد.

در داده‌های مربوط بیماری لوپوس، ۵ متغیر تبیینی برای پیش‌بینی متغیر پاسخ دودویی داریم. مدل با همه متغیرها (۱۳.۳) است. در اینجا می‌خواهیم با توجه به الگوریتمی که توضیح دادیم، انتخاب مناسب‌ترین مدل به کمک متر دورسو را انجام دهیم.

بیشترین مقدار R_{adj}^2 در مرحله ۴ است، این مقدار بیان می‌کند که بهترین مدل، مدلی است که در آن متغیرهای $(ESRtest)x_1$ ، $(AntiDNAtest)x_2$ ، $(ANAtest)x_3$ و $(FamilyHistory)x_5$ حضور داشته باشند. مدل مناسب با استفاده از این روش به صورت زیر است.

$$\hat{Y} = \tilde{b}_0 + \tilde{b}_1 ESRtest + \tilde{b}_2 AntiDNAtest + \tilde{b}_3 ANAtest + \tilde{b}_4 FamilyHistory \quad (17.4)$$

^۱ Coefficient of multiple determination

^۲ Adjusted coefficient of multiple determination

جدول ۱.۴: مراحل انتخاب مدل

شماره	MPE	R^2	R^2_{adj}	متغیرهای تبیینی
۱	۱.۲۴۵	۰.۸۱۰	۰.۷۰۵	x_3
۲	۰.۶۳۹	۰.۹۰۳	۰.۸۳۰	x_3, x_2
۳	۰.۵۳۳	۰.۹۱۹	۰.۸۳۸	x_3, x_2, x_1
۴	۰.۴۵۶	۰.۹۳۰	۰.۸۳۹	x_3, x_2, x_1, x_5
۵	۰.۴۴۷	۰.۹۳۲	۰.۸۱۰	x_3, x_2, x_1, x_5, x_6

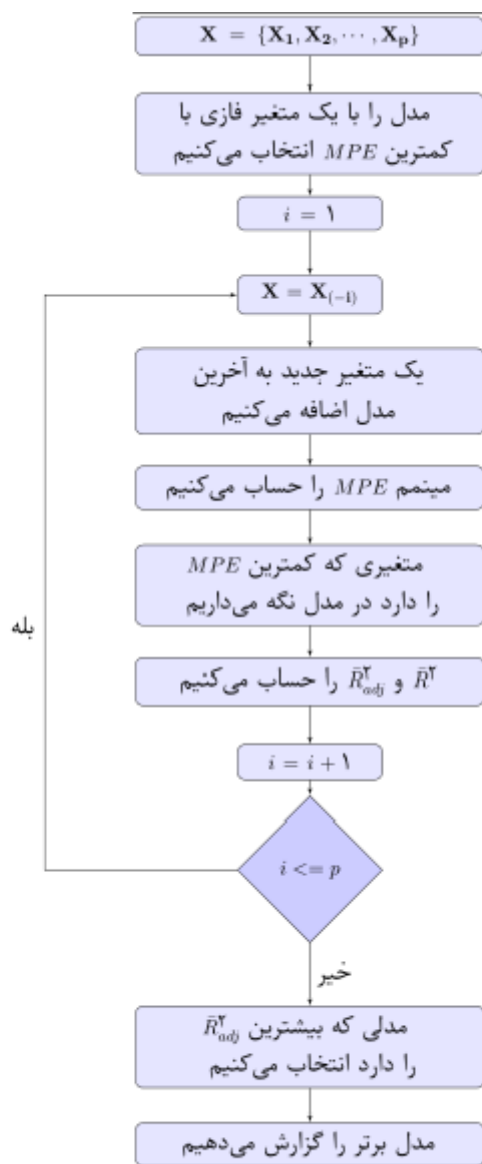
دلیل اینکه از ضریب تعیین تعدیل یافته برای انتخاب مدل مناسب انتخاب می‌کنیم این است که درصد نشان داده شده در ضریب تعیین تأثیر همه متغیرهای تبیینی بر متغیر وابسته است، به همین علت با اضافه کردن متغیر جدید به مدل مقدار ضریب تعیین افزایش می‌یابد. در حالی که درصد داده شده در ضریب تعیین تعدیل یافته فقط حاصل از تأثیر واقعی متغیرهای تبیینی مدل بر متغیر پاسخ است.

مدل (۱۷.۴) که با استفاده از روش انتخاب مدل متغیر کپی و همکاران [۱۱] انتخاب شده است، همانند مدل (۱۴.۳) که با شیوه سلمانی و همکاران [۲۷] با قرار دادن $LE = 0.7$ ، می‌باشد. به عبارت دیگر می‌توان گفت در امکان ابتلا به بیماری لوپوس عوامل $ANAtest$ ، $ESRtest$ و $AntiDNAtest$ و سابقه خانوادگی تأثیر گذار تر از سایر عوامل هستند.

مدل (۱۷.۴) را با استفاده از مدل پیشنهادی کپی و همکاران [۱۱] برازش می‌دهیم، در این صورت خواهیم داشت:

$$\begin{aligned} \hat{Y} = & (-4/0.23, 0.7, 0.733)_T + (-0.614, 0, 0.019, 0.019)_T ESRtest \\ & + 0.009 AntiDNAtest + (0.435, 0.0009, 0)_T ANAtest \\ & + (0.361, 0.237, 0.606)_T FamilyHistory \end{aligned} \quad (18.4)$$

مدل (۱۸.۴) نشان دهنده لگاریتم بخت امکانی فرد i ام به بیماری لوپوس است.



شکل ۱.۴: الگوریتم انتخاب متغیر کپی و همکاران در مدل رگرسیون لجستیک فازی

بحث و نتیجه گیری و پیشنهادات برای آینده تحقیق

بحث و نتیجه گیری

در این پایان نامه، مدلی که برای برآورد ضرایب با پاسخ دودویی ارائه کردیم، مدل رگرسیون لجستیک فازی است. رگرسیون لجستیک فازی همان رگرسیون لجستیک است که در فصل یک آن را معرفی کردیم. تفاوت آنها در این است که در رگرسیون لجستیک فازی متغیرهای تبیینی دقیق هستند اما متغیرهای پاسخ به جای اعداد صفر و یک به صورت مجموعه‌های فازی بین صفر و یک [۰, ۱] هستند. رگرسیون لجستیک فازی نیز با استفاده از تبدیل لجیت خطی می‌شود. در این مدل متغیر پاسخ و ضرایب فازی و متغیر تبیینی (پیشگو) دقیق هستند. دو روش استفاده شده در این پایان نامه برای برآورد ضرایب رگرسیون لجستیک فازی، روش امکانی و روش کمترین توان دوم خطا می‌باشد.

رگرسیون با رویکرد امکانی به دلیل اینکه توابع عضویت مجموعه‌های فازی اغلب به عنوان توزیع‌های امکانی توصیف می‌شوند، با این اسم نام‌گذاری شده‌اند. برای برآورد ضرایب فازی باید مدل فازی را به حداقل برسانیم، که یک مسئله برنامه‌ریزی خطی است و در آن تابع هدف مجموع پهنای فازی است. برای این کار از روش‌های عددی استفاده می‌کنیم.

روش دیگر که برای برآورد ضرایب استفاده می‌شود، روش کمترین توان دوم خطا است. در این روش فاصله بین مقادیر مشاهده شده و مقادیر برآورد شده حداقل می‌شود. از دو فاصله متفاوت برای این منظور استفاده شد. فاصله که توسط ژو و همکاران [۳۵] معرفی شده است، یکی از این فواصل است. این فاصله در تحقیقات زیادی استفاده شده است. علاوه بر این چون این فاصله براساس α - برش‌ها به دست می‌آید حایز اهمیت است. همچنین تابع وزنی که در تعریف این فاصله وجود دارد، باعث انعطاف‌پذیری بیشتر این فاصله شده است.

از فاصله کپی و دورسو [۱۱] به عنوان یک روش پیشنهادی برای برآورد ضرایب مدل با رویکرد کمترین توان دوم استفاده شده است. برای حداقل کردن این فاصله، مراکز متغیر پاسخ فازی با استفاده از یک مدل رگرسیون خطی مدل‌سازی شدند. در این فاصله نیز مانند فاصله ژو [۳۵]، از توابع وزن استفاده می‌شود که به تنظیم مناسب پهنای چپ و راست در

محاسبه فاصله بین \bar{y} و \bar{Y} کمک می کند.

در هر دو شیوه امکانی و کمترین توان دوم، مناسبترین مدل را انتخاب کردیم. در این روش از میانگین توان دوم خطا و ضریب تعیین تعدیل یافته و معیار LE استفاده می شود. با انتخاب $LE = 0.7$ مدل انتخابی به روش امکانی و کمترین توان دوم از نظر انتخاب متغیر، یکسان شدند.

در نهایت باید گفت دلیل اصلی انجام این پایان نامه با موضوع مذکور، ابهام موجود در یافته های پزشکی است زیرا برای مدل سازی روابط بین داده های پزشکی، فرضیات آماری کافی نیستند. مجموعه های فازی می توانند جواب مفیدی برای ابهامات پزشکی راجع به امکان ابتلای یک فرد مشکوک به بیماری باشد.

پیشنهادات

۱. بررسی رگرسیون لجستیک، با متغیرهای تبیینی فازی و متغیر پاسخ و ضرایب دقیق
۲. بررسی رگرسیون لجستیک، با متغیر تبیینی، متغیر پاسخ و ضرایب فازی
۳. تشخیص داده های پرت در رگرسیون لجستیک فازی
۴. انجام آزمون فرض های معنی داری پارامترها برای رگرسیون لجستیک فازی

مراجع

- [۱] پوراحمد، س. آیت‌اللهی، س.م.ت. طاهری، س.م. حبیب‌آگهی، ز (۱۳۹۰). مدل‌سازی موقعیت‌های مبهم تشخیص در پزشکی با استفاده از رگرسیون لجستیک به روش کمترین مربعات فازی، **سری سیستم‌های فازس و محاسبات نرم**، جلد ۲، ص ۱۹
- [۲] خراباتی، میترا. ربیعی، محمدرضا، (۱۳۹۸). شیوه‌ای جدید در تحلیل رگرسیون لجستیک فازی، **نهمین سمینار آمار و احتمال فازی، بابلسر**
- [۳] ربیعی، محمدرضا. خراباتی، میترا (۱۳۹۸). معرفی بسته fuzzyreg برای تحلیل برخی از روشهای رگرسیون فازی در R، **نهمین سمینار آمار و احتمال فازی، بابلسر**
- [۴] سلمانی ف، (۱۳۹۵)، پایان‌نامه دکتری: مدل مارکف فازی وابسته به پیشگو و کاربرد آن در تحلیل پاسخ‌های مبهم طولی، دانشگاه علوم پزشکی و خدمات بهداشتی درمانی شهید بهشتی
- [۵] طاهری م، (۱۳۷۵)، **آشنایی با نظریه مجموعه‌های فازی** انتشارات جهاد دانشگاهی مشهد
- [۶] طاهری م، ماشین چی م، (۱۳۹۲)، **مقدمه‌ای بر احتمال و آمار فازی** چاپ دوم، انتشارات دانشگاه شهید باهنر کرمان
- [۷] مشکانی علی، (۱۳۸۳)، **تحلیل آماری داده‌های رسته‌ای** چاپ سوم، انتشارات دانشگاه فردوسی مشهد
- [۸] میرزایی یگانه ش، ارقامی ن، (۱۳۸۶)، **رگرسیون فازی: مروری بر چند رویکرد، اندیشه آماری ۲۳:۳۵-۴۷**

[9] Akaike H.,(1974), A new look at the statistical model identification, **IEEE Trans Auto Control** 19, (6), pp 716

[10] Clemins A.,(1987), Least squares mthod fitting to fuzzy vectr data, **Fuzzy Sets and Systems** 22, pp 260-269

- [11] Coppi A., Dourso P., Giordani P. and Santaro A.,(2006), Least squares estimatin of a linear regression model with LR fuzzy response, **Comp. Stat. Data. Anal** 51, pp 267-186
- [12] Diamond P., (1987), Least squares method fitting of several fuzzy variable, Proc of second IFSA congres, Tokyo, pp 20
- [13] Dubois D., Kerre E., Mesiar R., and Prade H., (2000), **Fuzzy interval analysis. In Fundamentals of fuzzy sets**, pp 483-581
- [14] Gagolewski M., and Caha J., (2019), A Guide to the FuzzyNumbers Package for R.
- [15] Gildeh B., Sand Gien D. E. N. I. S., (2002), A goodness of fit index to reliability analysis in fuzzy model, **3rd WSEAS international conference on fuzzy sets and fuzzy systems. Iterlaken, Switzerland**, pp 11-14
- [16] Diamond P. and Korner R., (1997), Extaded fuzzy linear models and least squares estimated, **Comput. Math. Appl** 33, (9), pp 15-32
- [17] Dubias D., (1980), **Fuzzy set system: Theory and Application** , Academic Press, New York
- [18] Fauci A., Braunwald E., Kasper D., Hauser S., Longo D., Jameson J. and Loscalzo J., (2008), **Harrison's Principles of Internal Medicine** val:II, J.Weily, New York
- [19] Hojati M., Bector C. R. and Smimou K., (2005), A simple method for computation of fuzzy linear regression, **Eyropean Journal of Operational Research** 166 , pp 172-182
- [20] Kao C. and Chyu C. L., (2002), A fuzzy linear regression model with explantory power, **Fuzzy Sets and Systems** 126, pp 401-409
- [21] Lee E. S. and Chang P. T., (1994), Fuzzy linear regression analysis with spread unconstraines in sign, **Comp. Math. Appl** 28, (4), pp 61-79
- [22] Meye J. P., Stanley D. J., Herscovitch, L. and Topolnytsky, L., (2002), Affective, Continuance, and Normative Commitment to the Organization: A Metaanalysis of Antecedents, Correlates, and Consequences, **Journal of Vocational Behavior**, 61, pp 20–52
- [23] McCullagh, P. and Nelder, J. A., (1989), **Generalized Linear Models**, 2nd ed, Chapman and Hall, London
- [24] Mohammadi j. Taheri S. M., (2004), Pedomodels fitting with fuzzy least squares regression, **Iranian Journal of Fuzzy System**, (2), pp 45-61

- [25] Nelder, J. A., Wedderburn, R. W. M., (1972), Generalized Linear Models, **Journal of the Royal Statistical Society, A** , 135 pp 370–384.
- [26] Neter, J., Kutner M. H., Nachtsheim C. J. and Wasserman W., (1996), **Applied Linear Statistical Models**, Chicago: Irwin
- [27] Salmani F., Taheri S. M. and Abad, A., (2019), A forward variable selection method for fuzzy logistic regression, **International Journal of Fuzzy Systems**, **21(4)**, pp 1259-1269.
- [28] Savic D.A. and Pedrycz W., (1991), Evaluation of fuzzy linear regression models, **Fuzzy Sets and Systems** 47, (2), pp 173-181
- [29] Tabaei B. P. and Herman W. H., (2002), A multivariate logistic regression equation to screen for diabetes: development and validation, **Diabetes Care** 25(11), pp 1999-2003
- [30] Tanaka H., Hayashi I. and Watada J., (1989), Possibilitic linear regression for fuzzy data, **Journal of Operational Research** 40, pp 389-396
- [31] Tanaka H. and Ishibuchi H., (1991), Identification of possibilitic linear systems by quadratic membership function of fuzzy parameters, **Fuzzy Sets and Systems** 41, pp 145-160
- [32] Pourahmad S., Ayatollahi S. M. T., Taheri S. M. and Agahi Z. H., (2011), Fuzzy logistic regression based on the least squares approach with application in clinical studies. **Computers and Mathematics with Applications**, 62(9), pp 3353-3365.
- [33] Tanaka H., vejima D. and Asia K., (1982), Linear regression analysis with fuzzy model, **IEEE Trans, system Man cyberent** 12, pp 903-907
- [34] Van Broekhoven E. and De Baets B., (2006), Fast and accurate center of gravity defuzzification of fuzzy system outputs defined on trapezoidal fuzzy partitions. **Fuzzy sets and systems**, 157(7), pp 904-918
- [35] Xu R. Li C. (2001), Multicimentional leas-squares fitting with a fuzzy model **Fuzzy Sets and Systems** 119, pp 215-223
- [36] Zadeh L.A., (1975) , The concept of linguistic variable and its application to approximate reasonig I , II and III, **Information Sciences**, (8) pp 199-249
- [37] Zimmermann H.J., (2001), **Fuzzy Set Theory and Its Appliction**, Kluwer, Boston 4rd ed

پیوست آ

معرفی بسته FuzzyNumbers

نام کامل این بسته Tools to Deal with Fuzzy Numbers است. [۳] این بسته از قابلیت معرفی اعداد فازی به روش‌های گوناگون برخوردار است. برای نصب نرم افزار از دستور

```
install.packages('FuzzyNumbers')
```

استفاده می‌کنیم. با دستور

```
library(FuzzyNumbers)
```

بسته را فراخوانی می‌کنیم.

عدد فازی LR

عدد فازی A یک زیر مجموعه از \mathbb{R} با تابع عضویت زیر است:

$$\mu_A(x) = \begin{cases} 0 & x < a_1 \\ \text{left}\left(\frac{x-a_1}{a_2-a_1}\right) & a_1 \leq x < a_2 \\ 1 & a_2 \leq x \leq a_3 \\ \text{right}\left(\frac{x-a_3}{a_4-a_3}\right) & a_3 < x < a_4 \\ 0 & a_4 < x \end{cases} \quad (1.آ)$$

که $a_1 \leq a_2 \leq a_3 \leq a_4$ و $a_1, a_2, a_3, a_4 \in \mathbb{R}$ است. تابع $left : [0, 1] \rightarrow [0, 1]$ و $left(0) \geq 1$ و $right(1) \leq 0$ و تابعی غیر نزولی و $right : [0, 1] \rightarrow [0, 1]$ و $right(0) \geq 1$ و $right(1) \leq 0$ تابعی غیر صعودی است.

مثال: A_1 یک عدد فازی دوزنقه‌ای است، در نرم افزار به صورت زیر تعریف می‌شود:

```
A1 <- FuzzyNumber(1, 2, 4, 7,
left=function(x) x,
right=function(x) 1-x
)
```

دستورات زیر جزییات مربوط به عدد فازی را نمایش می‌دهند.

```
print(A1)
Fuzzy number with:
  support=[1,7],
  core=[2,4].

class(A1)
## [1] "FuzzyNumber"
## attr(,"package")
## [1] "FuzzyNumbers"

A1
## Fuzzy number with:
## support=[1,7],
## core=[2,4].

summary(A1)
      Length      Class      Mode
      1 FuzzyNumber      S4
```

عدد فازی مثلثی و دوزنقه‌ای

اعداد فازی مثلثی با دستور

```
TrapezoidalFuzzyNumber(a_1, a_2, a_3, a_4)
```

```
TriangularFuzzyNumbe(a_1, a_2, a_3)
```

نمایش داده می‌شوند تحت شرط $a_1 < a_2 < a_3 < a_4$.
نمایش‌های دیگری برای اعداد فازی ذوزنقه‌ای و مثلثی وجود دارند که در زیر به آنها اشاره شده است.

```
T1 <- TrapezoidalFuzzyNumber(a_1, a_2, a_3, a_4)
T2 <- TrapezoidalFuzzyNumber(a_1, a_2, a_2, a_3)
# or T2 <- TriangularFuzzyNumber(a_1, a_2, a_3)
TrapezoidalFuzzyNumber(a_1, a_1, a_2, a_2) # crisp interval
as.TrapezoidalFuzzyNumber(c(a_1, a_2)) # is the same
TrapezoidalFuzzyNumber(a_1, a_1, a_1, a_1) # crisp real
as.TrapezoidalFuzzyNumber(a_1) # is the same

T1 <- TrapezoidalFuzzyNumber(1, 1.5, 4, 7)
```

عدد فازی قطعه‌ای خطی

عدد فازی قطعه‌ای خطی را نیز می‌توان در بسته FuzzyNumbers تولید کرد، برای این منظور از دستور زیر استفاده می‌کنیم.

```
PiecewiseLinearFuzzyNumber(a1, a2, a3, a4, knot.n = 0,
  knot.alpha = numeric(0), knot.left = numeric(0),
  knot.right = numeric(0))
```

knot.n تعداد برش (گره)‌های قطعه‌کننده عدد فازی را مشخص می‌کند. knot.left ابتدای برش‌ها، knot.right انتهای برش‌ها و knot.alpha ارتفاع برش‌ها را مشخص می‌کند.

```
P1 <- PiecewiseLinearFuzzyNumber(1, 2, 3, 4,
  knot.n=1, knot.alpha=0.25, knot.left=1.5, knot.right=3.25)
P1
Piecewise linear fuzzy number with 1 knot(s),
  support=[1,4],
  core=[2,3].
```

```
P2 <- PiecewiseLinearFuzzyNumber(1, 2, 3, 4,
knot.n=2, knot.alpha=c(0.25,0.6),
knot.left=c(1.5,1.8), knot.right=c(3.25, 3.5))
P2
Piecewise linear fuzzy number with 2 knot(s),
  support=[1,4],
  core=[2,3].
```

اعمال تابع یکنوا بر عدد فازی

به منظور اعمال یک تابع یکنوا بر اعداد فازی از دستور `fapply` استفاده می‌شود. نکته قابل توجه این است که عدد فازی در اینجا باید از نوع قطعه‌ای باشد. در مثال زیر تابع $\sqrt{\log A}$ اعمال شده است.

```
A <- as.PiecewiseLinearFuzzyNumber(
TrapezoidalFuzzyNumber(0,1,2,3), knot.n=100)
plot( fapply(A, function(x) sqrt(log(x+1))))
```

محاسبه هسته، دامنه و برش‌های عدد فازی

هسته، دامنه و برش‌های یک عدد فازی به ترتیب با دستورهای `supp` و `alphacut` قابل محاسبه‌اند.

```
A <- FuzzyNumber(-5, 3, 6, 20,
left=function(x) pbeta(x,0.4,3),
right=function(x) 1-x^(1/4),
lower=function(alpha) qbeta(alpha,0.4,3),
upper=function(alpha) (1-alpha)^4
)
```

```
supp(A)
[1] -5 20
```

```
core(A)
[1] 3 6
```

alphacut(A, 0)

L U

0 -5 20

عملگرهای حسابی

عملگرهای دوبعدی^۱ ریاضی برای اعداد فازی متقارن برای هر $\alpha \in [0, 1]$ به صورت زیر تعریف می شوند:

$$(A \otimes B)_\alpha = A_\alpha \otimes B_\alpha$$

که $\otimes = +, -, \times, /$ و A و B اعداد فازی متقارن هستند. به طور مثال جمع $A + B$ به صورت زیر تعریف می شود:

$$(A + B)_\alpha = A_\alpha + B_\alpha = [A_L(\alpha) + B_L(\alpha), A_U(\alpha) + B_U(\alpha)]$$

همچنین برای $\lambda \in \mathbb{R}$ ضرب اسکالر به صورت زیر تعریف می شود:

$$(\lambda \cdot A)_\alpha = \lambda A_\alpha = \begin{cases} [\lambda A_L(\alpha), \lambda A_U(\alpha)], & \lambda \geq 0 \\ [\lambda A_U(\alpha), \lambda A_L(\alpha)], & \lambda < 0 \end{cases}$$

```
A <- TrapezoidalFuzzyNumber(0, 1, 1, 2)
```

```
B <- TrapezoidalFuzzyNumber(1, 2, 2, 3)
```

```
A+B
```

```
Trapezoidal fuzzy number with:
```

```
  support=[1,5],
```

```
  core=[3,3].
```

```
A-B
```

```
Trapezoidal fuzzy number with:
```

```
  support=[-3,1],
```

```
  core=[-1,-1].
```

```
3*A
```

```
Trapezoidal fuzzy number with:
```

```
  support=[0,6],
```

```
  core=[3,3].
```

^۱binary

تخمین قطعه‌ای یک عدد فازی

در بسته FuzzyNumbers، ۳ روش مختلف برای تخمین قطعه‌ای یک عدد فازی وجود دارد:

۱. کمترین فاصله اقلیدسی (NearestEuclidean)

۲. محافظ هسته و دامنه (SupportCorePreserving)

۳. ساده (Naive)

برای تخمین قطعه‌ای اعداد فازی از دستور زیر در نرم افزار استفاده می‌کنیم:

```
piecewiseLinearApproximation(object,
  method=c("NearestEuclidean", "SupportCorePreserving",
    "Naive"),
  knot.n=1, knot.alpha=seq(0, 1, length.out=knot.n+2)[-c(1,knot.n+2)],
  ..., verbose=FALSE)
```

حال به اختصار به تعریف روش‌ها می‌پردازیم.

۱. کمترین فاصله اقلیدسی در این روش، تخمین یک عدد فزی به وسیله یک عدد فازی قطعه‌ای به گونه‌ای انجام می‌گیرد که کمترین فاصله اقلیدسی را از عدد فازی اصلی داشته باشد.

۲. محافظ هسته و دامنه هسته و دامنه تخمین قطعه‌ای یک عدد فازی به این روش، معادل هسته و دامنه همان عدد فازی اصلی است.

۳. ساده در این روش، تخمین قطعه‌ای یک عدد فازی به گونه‌ای صورت می‌گیرد که:

(آ) هسته عدد فازی اصلی برابر با هسته عدد فازی تخمینی

(ب) دامنه عدد فازی اصلی برابر با دامنه عدد فازی تخمینی

(ج) برش‌های معرفی شده در دستور برای عدد فازی اصلی، برابر با برش‌های متناظر عددهای فازی تخمینی

```
A <- FuzzyNumber(-5, 3, 6, 20,
  left=function(x) pbeta(x,0.4,3),
  right=function(x) 1-x^(1/4),
  lower=function(alpha) qbeta(alpha,0.4,3),
  upper=function(alpha) (1-alpha)^4)
```

)

```
(T1 <- trapezoidalApproximation(A, method='Naive'))
```

Trapezoidal fuzzy number with:

```
support=[-5,20],
```

```
core=[3,6].
```

```
(T2 <- trapezoidalApproximation(A, method='NearestEuclidean'))
```

Trapezoidal fuzzy number with:

```
support=[-5.85235,14.4],
```

```
core=[-2.26529,3.2].
```

```
(T3 <- trapezoidalApproximation(A, method='ExpectedIntervalPreserving'))
```

Trapezoidal fuzzy number with:

```
support=[-5.85235,14.4],
```

```
core=[-2.26529,3.2].
```

مینیمم و ماکزیمم عدد فازی

مینیمم عدد فازی برای برش‌ها به صورت

$$A_\alpha \wedge B_\alpha = [A_L(\alpha) \wedge B_L(\alpha), A_U(\alpha) \wedge B_U(\alpha)]$$

و ماکزیمم عدد فازی به صورت

$$A_\alpha \vee B_\alpha = [A_L(\alpha) \vee B_L(\alpha), A_U(\alpha) \vee B_U(\alpha)]$$

تعریف می‌شود.

مقدار مینیمم و ماکزیمم عدد فازی به صورت زیر قابل استفاده است. تعریف زیر فقط برای اعداد فازی خطی قطعه‌ای تعریف شده است.

```
x = as.PiecewiseLinearFuzzyNumber(TriangularFuzzyNumber
```

```
(-4.8, -3, -1.5), knot.n = 9)
```

```
y = as.PiecewiseLinearFuzzyNumber(TriangularFuzzyNumber
```

```
(-5.5, -2.5, -1.1), knot.n = 9)
```

```
min = min(x@a1,y@a1)
```

```
max = max(x@a4,y@a4)
```


پیوست ب

کدهای استفاده شده در پایان نامه

ب.۱ کدهای انتخاب مدل با معیارر LE

کدهای R

```
library(FuzzyNumbers)
verylow<- as.PiecewiseLinearFuzzyNumber(TrapezoidalFuzzyNumber
(0.01,0.02,0.02,0.18), knot.n=100)
y1<-fapply(verylow, function(x) log(x/(1-x)))
T1 <- trapezoidalApproximation(y1, method='ExpectedIntervalPreserving')
Y1<- as.PiecewiseLinearFuzzyNumber(TrapezoidalFuzzyNumber
(supp(T1)[1],core(T1)[1],core(T1)[1],supp(T1)[2]), knot.n=100)
Y11=c(core(Y1)[1],core(Y1)[1]-supp(Y1)[1],supp(Y1)[2]-core(Y1)[2])
low<- as.PiecewiseLinearFuzzyNumber(TrapezoidalFuzzyNumber
(0.1, 0.25, 0.25, 0.4), knot.n=100)
y2<-fapply(low, function(x) log(x/(1-x)))
T2 <- trapezoidalApproximation(y2, method='ExpectedIntervalPreserving')
Y2<- as.PiecewiseLinearFuzzyNumber(TrapezoidalFuzzyNumber
```

```
(supp(T2) [1], core(T2) [1], core(T2) [1], supp(T2) [2]), knot.n=100)
Y22=c(core(Y2) [1], core(Y2) [1]-supp(Y2) [1], supp(Y2) [2]-core(Y2) [2])
medium<- as.PiecewiseLinearFuzzyNumber(TrapezoidalFuzzyNumber
(0.35, 0.5, 0.5, 0.65), knot.n=100)
y3<-fapply(medium, function(x) log(x/(1-x)))
T3 <- trapezoidalApproximation(y3, method='ExpectedIntervalPreserving')
Y3<- as.PiecewiseLinearFuzzyNumber(TrapezoidalFuzzyNumber
(supp(T3) [1], core(T3) [1], core(T3) [1], supp(T3) [2]), knot.n=100)
Y33=c(core(Y3) [1], core(Y3) [1]-supp(Y3) [1], supp(Y3) [2]-core(Y3) [2])
high<- as.PiecewiseLinearFuzzyNumber(TrapezoidalFuzzyNumber
(0.6 , 0.75, 0.75, 0.9), knot.n=100)
y4<-fapply(high, function(x) log(x/(1-x)))
T4 <- trapezoidalApproximation(y4, method='ExpectedIntervalPreserving')
Y4<- as.PiecewiseLinearFuzzyNumber(TrapezoidalFuzzyNumber
(supp(T4) [1], core(T4) [1], core(T4) [1], supp(T4) [2]), knot.n=100)
Y44=c(core(Y4) [1], core(Y4) [1]-supp(Y4) [1], supp(Y4) [2]-core(Y4) [2])
veryhigh<- as.PiecewiseLinearFuzzyNumber(TrapezoidalFuzzyNumber
(0.8, 0.98, 0.98, 0.99), knot.n=100)
y5<-fapply(veryhigh, function(x) log(x/(1-x)))
T5 <- trapezoidalApproximation(y5, method='ExpectedIntervalPreserving')
Y5<- as.PiecewiseLinearFuzzyNumber(TrapezoidalFuzzyNumber
(supp(T5) [1], core(T5) [1], core(T5) [1], supp(T5) [2]), knot.n=100)
Y55=c(core(Y5) [1], core(Y5) [1]-supp(Y5) [1], supp(Y5) [2]-core(Y5) [2])
Y=rbind(Y44,Y33,Y44,Y44,Y33,Y55,Y33,Y44,Y22,Y11,Y22,Y33,Y22,Y22,Y33)
```

کدهای Mathematica

```
Dade = Import["F:\\\lopos.xlsx", {"Data", 1}];
n = Length[Transpose[Dade][[1]]];
p = Length[Dade[[1]]] - 3;
Y = Dade[[1 ;;, {p + 1, p + 2, p + 3}]];
Ybar = Mean[Y];
X = Dade[[1 ;;, Table[i, {i, p}]]];
t = Table[i, {i, 2, p}];
```

```

LE = 0.07;

MSE = Array[mse, p - 1];
For[i = 1, i <= p - 1, i++,
  X = Dade[[1 ;;, {1, i + 1}]];
  j = Length[Transpose[X]];
  A = Array[\[Beta], j, 0];
  A1 = Array[\[Sigma]1, j, 0];
  Au = Array[\[Sigma]u, j, 0];
  Yhat = Array[z, {n, 3}];
  For[h = 1, h <= n, h++,
    Yhat[[h, 1]] = X[[h]].A;
    Yhat[[h, 2]] = X[[h]].A1;
    Yhat[[h, 3]] = X[[h]].Au;
  ];
  d = {};
  For[k = 1, k <= j, k++,
    d = Join[d, {A1[[k]] >= 0, Au[[k]] >= 0}];
  ];
  B = Join[A, A1, Au];
  \[Lambda] = 0.5;
  \[Rho] = 0.5;
  Q1 = NMinimize[{3*(Y[[A11, 1]] - Yhat[[A11, 1]]).(Y[[A11, 1]] -
    Yhat[[A11, 1]]) -
    2 \[Lambda]*(Y[[A11, 1]] - Yhat[[A11, 1]]).(Y[[A11, 2]] -
    Yhat[[A11, 2]]) + \[Lambda]^2*(Y[[A11, 2]] -
    Yhat[[A11, 2]]).(Y[[A11, 2]] - Yhat[[A11, 2]]) +
    2 \[Rho]*(Y[[A11, 1]] - Yhat[[A11, 1]]).(Y[[A11, 3]] -
    Yhat[[A11, 3]]) + \[Lambda]^2*(Y[[A11, 3]] -
    Yhat[[A11, 3]]).(Y[[A11, 3]] - Yhat[[A11, 3]]), d}, B];
  MSE[[i]] = Q1[[1]];
MSE
a1 = Min[MSE]
var1 = Which[a1 == MSE[[1]], t[[1]], a1 == MSE[[2]], t[[2]],

```

```

a1 == MSE[[3]], t[[3]], a1 == MSE[[4]], t[[4]], a1 == MSE[[5]], t[[5]]

MSE = Array[mse, p - 2];
bb = Delete[t, var1 - 1]
For[i = 1, i <= p - 2, i++,
  X = Dade[[1 ;;, {1, var1, bb[[i]]}]];
  j = Length[Transpose[X]];
  A = Array[\[Beta], j, 0];
  A1 = Array[\[Sigma]1, j, 0];
  Au = Array[\[Sigma]u, j, 0];
  Yhat = Array[z, {n, 3}];
  For[h = 1, h <= n, h++,
    Yhat[[h, 1]] = X[[h]].A;
    Yhat[[h, 2]] = X[[h]].A1;
    Yhat[[h, 3]] = X[[h]].Au;
  ];
  d = {};
  For[k = 1, k <= j, k++,
    d = Join[d, {A1[[k]] >= 0, Au[[k]] >= 0}]
  ];
  B = Join[A, A1, Au];
  \[Lambda] = 0.5;
  \[Rho] = 0.5;
  Q1 = NMinimize[{3*(Y[[A11, 1]] - Yhat[[A11, 1]]).(Y[[A11, 1]] - Yhat[[A11, 1]]) -
    2 \[Lambda]*(Y[[A11, 1]] - Yhat[[A11, 1]]).(Y[[A11, 2]] -
    Yhat[[A11, 2]]) + \[Lambda]^2*(Y[[A11, 2]] -
    Yhat[[A11, 2]]).(Y[[A11, 2]] - Yhat[[A11, 2]]) +
    2 \[Rho]*(Y[[A11, 1]] - Yhat[[A11, 1]]).(Y[[A11, 3]] -
    Yhat[[A11, 3]]) + \[Lambda]^2*(Y[[A11, 3]] -
    Yhat[[A11, 3]]).(Y[[A11, 3]] - Yhat[[A11, 3]]), d}, B];
  MSE[[i]] = Q1[[1]];]
MSE
a2 = Min[MSE]
var2 = Which[a2 == MSE[[1]], bb[[1]], a2 == MSE[[2]], bb[[2]], a2 == MSE[[3]], bb[[3]],

```

۸۱ کدهای انتخاب مدل با معیارر LE

```
a2 == MSE[[4]],bb[[4]]
a2 - a1 > LE

MSE = Array[mse, p - 3];
cc = Delete[t, {{var1 - 1}, {var2 - 1}}]
For[i = 1, i <= p - 3, i++,
  X = Dade[[1 ;;, {1, var1, var2, cc[[i]]}]];
  j = Length[Transpose[X]];
  A = Array[\[Beta], j, 0];
  A1 = Array[\[Sigma]1, j, 0];
  Au = Array[\[Sigma]u, j, 0];
  Yhat = Array[z, {n, 3}];
  For[h = 1, h <= n, h++,
    Yhat[[h, 1]] = X[[h]].A;
    Yhat[[h, 2]] = X[[h]].A1;
    Yhat[[h, 3]] = X[[h]].Au;
  ];
  d = {};
  For[k = 1, k <= j, k++,
    d = Join[d, {A1[[k]] >= 0, Au[[k]] >= 0}]
  ];
  B = Join[A, A1, Au];
  \[Lambda] = 0.5;
  \[Rho] = 0.5;
  Q1 = NMinimize[{3*(Y[[A11, 1]] - Yhat[[A11, 1]]).(Y[[A11, 1]] -
    Yhat[[A11, 1]]) -
    2 \[Lambda]*(Y[[A11, 1]] - Yhat[[A11, 1]]).(Y[[A11, 2]] -
    Yhat[[A11, 2]]) + \[Lambda]^2*(Y[[A11, 2]] -
    Yhat[[A11, 2]]).(Y[[A11, 2]] - Yhat[[A11, 2]]) +
    2 \[Rho]*(Y[[A11, 1]] - Yhat[[A11, 1]]).(Y[[A11, 3]] -
    Yhat[[A11, 3]]) + \[Lambda]^2*(Y[[A11, 3]] -
    Yhat[[A11, 3]]).(Y[[A11, 3]] - Yhat[[A11, 3]]), d}, B];
  MSE[[i]] = Q1[[1]];]
MSE
```

```

a3 = Min[MSE]
var3 = Which[a3 == MSE[[1]], cc[[1]], a3 == MSE[[2]], cc[[2]],
  a3 == MSE[[3]], cc[[3]]]
a3-a2>LE

MSE = Array[mse, p - 4];
dd = Delete[t, {{var1 - 1}, {var2 - 1}, {var3 - 1}}]
For[i = 1, i <= p - 4, i++,
  X = Dade[[1 ;;, {1, var1, var2, var3, dd[[i]]}]];
  j = Length[Transpose[X]];
  A = Array[\[Beta], j, 0];
  A1 = Array[\[Sigma]1, j, 0];
  Au = Array[\[Sigma]u, j, 0];
  Yhat = Array[z, {n, 3}];
  For[h = 1, h <= n, h++,
    Yhat[[h, 1]] = X[[h]].A;
    Yhat[[h, 2]] = X[[h]].A1;
    Yhat[[h, 3]] = X[[h]].Au;
  ];
  d = {};
  For[k = 1, k <= j, k++,
    d = Join[d, {A1[[k]] >= 0, Au[[k]] >= 0}]
  ];
  B = Join[A, A1, Au];
  \[Lambda] = 0.5;
  \[Rho] = 0.5;
  Q1 = NMinimize[{3*(Y[[A11, 1]] - Yhat[[A11, 1]])*(Y[[A11, 1]] -
    Yhat[[A11, 1]]) -
    2 \[Lambda]*(Y[[A11, 1]] - Yhat[[A11, 1]])*(Y[[A11, 2]] -
    Yhat[[A11, 2]]) + \[Lambda]^2*(Y[[A11, 2]] -
    Yhat[[A11, 2]])*(Y[[A11, 2]] - Yhat[[A11, 2]]) +
    2 \[Rho]*(Y[[A11, 1]] - Yhat[[A11, 1]])*(Y[[A11, 3]] -
    Yhat[[A11, 3]]) + \[Lambda]^2*(Y[[A11, 3]] -
    Yhat[[A11, 3]])*(Y[[A11, 3]] - Yhat[[A11, 3]]), d}, B];

```

```

MSE[[i]] = Q1[[1]];]
MSE
a4 = Min[MSE]
var4 = Which[a4 == MSE[[1]], dd[[1]], a4 == MSE[[2]], dd[[2]]]
a4-a3>LE

```

ب.۲ کدهای مترژو

```

Dade = Import["F:\\\lopos.xlsx", {"Data", 1}];
n = Length[Transpose[Dade][[1]]];
p = Length[Dade[[1]]] - 3;
X = Dade[[1 ;;, Table[i, {i, p}]]];
AA = Dade[[1 ;;, {p + 1, p + 2, p + 3}]];
A = Array[\[Beta], p, 0];
S1 = Array[\[Sigma]1, p, 0];
Sr = Array[\[Sigma]r, p, 0];
What = Array[z, {n, 3}];
For[h = 1, h <= n, h++,
  What[[h, 1]] = X[[h]].A;
  What[[h, 2]] = X[[h]].S1;
  What[[h, 3]] = X[[h]].Sr;
];
dis = Array[q, {n, 2}];

For[i = 1, i <= n, i++,
  dis[[i, 1]] =
  Log[(AA[[i, 1]] - (AA[[i, 2]]*(1 - \[Alpha])))/(
  1 - AA[[i, 1]] + (AA[[i, 2]]*(1 - \[Alpha])));
  dis[[i, 2]] =
  Log[(AA[[i, 1]] + (AA[[i, 3]]*(1 - \[Alpha])))/(
  1 - AA[[i, 1]] - (AA[[i, 3]]*(1 - \[Alpha])));
];
d = {};
For[k = 1, k <= p, k++,

```

```
d = Join[d, {Sl[[k]] >= 0, Sr[[k]] >= 0}]
];
B = Join[A, Sl, Sr];
Z = \!\(
\*UnderoverscriptBox[\(\[Sum]\), \((t =
1\), \((15\))\)\(\([Alpha] . \((\((dis[\([t,
1\)]) - \((\((\([Alpha] - 1)\)*What[\([t, 2\)])\))\)) -
What[\([t, 1\)])\)^2 + \((dis[\([t,
2\)]) - \((\((1 - \[Alpha])\)*What[\([t, 3\)])\))\)) -
What[\([t, 1\)])\)^2\))\)\)\)
intt = \!\(
\*SubsuperscriptBox[\(\[Integral]\), \((0\), \((1\))\)\(Z \
\[DifferentialD]\[Alpha]\)\)
NMinimize[{intt, d}, B, Method -> "DifferentialEvolution"];
```

ب.۳ کدهای مترکیبی و همکاران

برآورد ضرایب

```
Dade = Import["F:\\\lopos.xlsx", {"Data", 1}];
n = Length[Transpose[Dade][[1]]];
p = Length[Dade[[1]]] - 3;
Y1 = {"MPE", "Beta"};
Y = Dade[[1 ;;, {p + 1, p + 2, p + 3}]];

A = Array[\[Beta], p, 0];
A1 = Array[\[Sigma]1, p, 0];
Au = Array[\[Sigma]u, p, 0];
Yhat = Array[z, {n, 3}];
X = Dade[[1 ;;, Table[i, {i, p}]]];
For[i = 1, i <= n, i++,
Yhat[[i, 1]] = X[[i]].A;
```



```

Yhat[[i, 2]] = X[[i]].A1;
Yhat[[i, 3]] = X[[i]].Au;
];
d = {};
For[k = 1, k <= p, k++,
  d = Join[d, {A1[[k]] >= 0, Au[[k]] >= 0}]
];
B = Join[A, A1, Au];
\[Lambda] = 0.5;
\[Rho] = 0.5;
Q1 = NMinimize[{3*(Y[[A11, 1]] - Yhat[[A11, 1]]).(Y[[A11, 1]] -
  Yhat[[A11, 1]]) -
  2 \[Lambda]*(Y[[A11, 1]] - Yhat[[A11, 1]]).(Y[[A11, 2]] -
  Yhat[[A11, 2]]) + \[Lambda]^2*(Y[[A11, 2]] -
  Yhat[[A11, 2]]).(Y[[A11, 2]] - Yhat[[A11, 2]]) +
  2 \[Rho]*(Y[[A11, 1]] - Yhat[[A11, 1]]).(Y[[A11, 3]] -
  Yhat[[A11, 3]]) + \[Lambda]^2*(Y[[A11, 3]] -
  Yhat[[A11, 3]]).(Y[[A11, 3]] - Yhat[[A11, 3]]), d} , B]

```

Abstract

Logistic regression helps in the modeling of discrete statistical data. These models are very useful in medical sciences (such as death / life, the presence of disease / non-disease, etc.). In diagnosis, doctors use various information sources, such as clinical history, medical examinations, laboratory tests, and so on. But placing people in both the sick and healthy groups is basically ambiguous, and obscure observations in clinical diagnosis are abundant. In such cases, logistic regression is not appropriate due to the lack of accuracy of the response variable. For fuzzy binary fuzzy response modeling, the fuzzy logistic regression model is presented in this thesis. To calculate the success of a fuzzy logistic model, the term linguistic terminology such as ..., low, medium, high, ... is defined. Then, using the principle of expansion, the transformation of the logic "odds ratio" is modeled on a set of observations of precise explanatory variables. also to estimate the coefficients of the proposed model, we use the probabilistic method and the least squares fuzzy method. We also offer two model selection methods to evaluate the model.

Keywords : Fuzzy Regression, Logistic Regression, Fuzzy Logistic Regression, Least Squares Error, Model Selection



Shahrood University of Technology

Faculty Of Mathematical Sciences

MSc Thesis in: Mathematical Statistics

**Logistic Regression in Fuzzy Environment
and its application in medical sciences**

By: Mitra Kharabati

Supervisor

Mohammad Reza Rabiee

Advisor

Fatemeh Salmani

September 2019