

مقایسه مدل‌های رگرسیون لجستیک با تحلیل جداسازی در پیش‌بینی دیابت نوع ۲

محمد آرام احمدی^۱، عباس بهرامپور^۲

^۱ فوق‌لیسانس آمار زیستی، مرکز مدل‌سازی در سلامت و گروه آمار زیستی و اپیدمیولوژی دانشکده بهداشت، دانشگاه علوم پزشکی کرمان

^۲ دکترای آمار زیستی، استاد مرکز مدل‌سازی سلامت پژوهشکده آینده پژوهی و گروه آمار زیستی و اپیدمیولوژی دانشکده بهداشت، دانشگاه علوم پزشکی کرمان

نویسنده رابط: عباس بهرامپور، نشانی: دانشکده بهداشت دانشگاه علوم پزشکی کرمان، تلفن: ۰۹۱۳۱۴۰۴۵۱۲، پست الکترونیک: abahrampour@yahoo.com

تاریخ دریافت: ۹۳/۱۰/۰۴؛ پذیرش: ۹۴/۰۳/۰۲

مقدمه و اهداف: بیماری دیابت از جمله بیماری‌های مزمن بوده، که درمان قطعی ندارد و شایع‌ترین علت قطع اندام، نابینایی و نارسایی کلیوی و از عوامل خطر در ایجاد بیماری‌های قلبی است. رگرسیون لجستیک و تحلیل جداسازی از مدل‌های تحلیل آماری برای امر پیش‌بینی و جداسازی چند متغیره می‌باشند. هدف تعیین متغیرهای تأثیرگذار بر دیابت نوع ۲ و مقایسه مدل‌های رگرسیون لجستیک و تحلیل جداسازی می‌باشد.

روش کار: داده‌ها شامل اطلاعات ۵۳۵۷ نفر از مرکز فیزیولوژی دانشگاه علوم پزشکی کرمان می‌باشد. متغیر پاسخ دیابت و متغیرهای وزن، قد، (BMI) (Body Mass Index)، دور کمر، دور باسن، نسبت کمر به باسن (WHR)، (Waist hip Ratio) کلاسترول و ... در مدل در نظر گرفته شدند. برای مقایسه از حساسیت، ویژگی، دقت، منحنی راک و شبیه‌سازی استفاده شد.

نتایج: حساسیت، ویژگی و دقت پیش‌بینی در مدل لجستیک و تحلیل جداسازی به ترتیب ۰.۷۴ و ۷۱/۱ و ۷۱/۵ و ۹۵/۲۲، ۲۲/۴ و ۸۵/۳ و ۹۹/۱۹، ۹۹/۲۶ و ۹۸/۲۶ و ۹۹/۵۶ شدند.

نتیجه‌گیری: نتایج نشان داد که دور کمر، سن، جنس، مصرف داروی کاهنده فشار خون، اندازه فشار خون سیستولیک و سطح LDL مهم می‌باشند. حساسیت در مدل لجستیک بیشتر، اما ویژگی و دقت پیش‌بینی در تحلیل جداسازی بالاتر بود. منحنی‌های راک نشان داد که مقادیر پیش‌بینی به طور مجانبی یکسان بودند.

واژگان کلیدی: حساسیت، ویژگی، رگرسیون لجستیک، تحلیل جداسازی، منحنی راک

مقدمه

دیابت به مجموعه اختلال متابولیسم کربوهیدرات، چربی و پروتئین که به وسیله یا فقدان ترشح انسولین یا کاهش حساسیت بافت‌ها به انسولین صورت می‌گیرد؛ اطلاق می‌شود. از میان انواع دیابت، دیابت نوع ۲ شایع‌ترین نوع است (۱). میزان وقوع جهانی دیابت به دلیل افزایش شیوع چاقی و کاهش میزان فعالیت بدنی در حال افزایش است (۲). شیوع چاقی در کشورهای صنعتی و توسعه یافته افزایش روزافزون دارد، به طوری که سازمان جهانی بهداشت، چاقی را اپیدمی جهانی گزارش کرده است (۳).

تشخیص الگوها و طبقه‌بندی، یکی از مهم‌ترین کاربردهای روش‌های آماری در علوم مختلف است. از جمله اهداف عمده طبقه‌بندی و مدل‌سازی در علوم آمار، پیش‌بینی بر اساس شواهد و متغیرها و اطلاعات موجود از یک موضوع خاص است. این امر در علوم آماری توسط روش‌هایی مانند رگرسیون، تحلیل ممیزی

در بین این روش‌ها رگرسیون لجستیک فرضیه‌های قابل توجهی نمی‌خواهد، اما در تحلیل جداسازی برقراری فرض نرمال بودن چند متغیره برای متغیرهای پیش‌بینی کننده علاوه بر استقلال و تصادفی بودن آن‌ها و برقراری فرض تساوی ماتریس کوواریانس لازم و ضروری می‌باشد. روش‌های طبقه‌بندی آماری سنتی مانند تحلیل جداسازی خطی فیشر و رگرسیون لجستیک به طور

بهداشت، چاقی را اپیدمی جهانی گزارش کرده است (۳). تشخیص الگوها و طبقه‌بندی، یکی از مهم‌ترین کاربردهای روش‌های آماری در علوم مختلف است. از جمله اهداف عمده طبقه‌بندی و مدل‌سازی در علوم آمار، پیش‌بینی بر اساس شواهد و متغیرها و اطلاعات موجود از یک موضوع خاص است. این امر در علوم آماری توسط روش‌هایی مانند رگرسیون، تحلیل ممیزی

رگرسیون لجستیک و تحلیل جداسازی به این نتیجه رسیدند که میزان حساسیت برای مدل‌ها به ترتیب $0/483$ و $0/677$ بود. ویژگی به ترتیب $0/857$ و $0/66$ به دست آمد و مساحت زیر منحنی راک (ROC)^۳ برای دو مدل به ترتیب $0/749$ و $0/739$ بود (۴). در مطالعه‌ای که توسط جورج آنتونوجورجس و همکاران (۲۰۰۹) انجام گرفت، دو مدل رگرسیون لجستیک و تحلیل جداسازی خطی را برای ارزیابی فاکتورهای مرتبط با آسم مقایسه کردند و میزان دقت، حساسیت و ویژگی مدل رگرسیون لجستیک و تحلیل جداسازی خطی را به ازای نقطه برش‌های مختلف به دست آوردند، که در نهایت به این نتیجه رسیدند که به طور کلی تابع جداسازی خطی زمانی که فرضیه‌های نرمال بودن برقرار باشد، روشی بهتر از رگرسیون لجستیک می‌باشد و اختلاف‌های بین دو روش زمانی که حجم نمونه به اندازه‌ی کافی بزرگ باشد، بسیار ناچیز است (۱۰).

هدف از این مطالعه تعیین متغیرهای تأثیرگذار بر ابتلا به دیابت نوع ۲ و مقایسه توانایی مدل‌های رگرسیون لجستیک و تحلیل جداسازی برای پیش‌بینی دیابت در نمونه‌ای از افراد که از قبل داده‌های آن از مرکز تحقیقات فیزیولوژی دانشگاه علوم پزشکی کرمان جمع‌آوری شده‌اند، می‌باشد.

روش کار

رگرسیون لجستیک

فرمی از رگرسیون می‌باشد و زمانی مورد استفاده قرار می‌گیرد، که متغیر وابسته به صورت دو یا چندسطحی بوده و متغیرهای مستقل از هر نوع دیگر باشند. در علوم پزشکی متغیر پیشامد معمولاً حضور یا غیاب یک وضع بیان شده یا یک بیماری می‌باشد (۱۰). برای این مدل اگر دو گروه وجود داشته باشد، رگرسیون لجستیک باینری مورد استفاده قرار می‌گیرد و چنانچه سه و بیش‌تر از سه گروه موجود باشد بین رگرسیون لجستیک اسمی و ترتیبی باید انتخابی صورت گیرد (۱۱).

از آنجایی که احتمال پیش‌بینی شده باید بین اعداد ۰ و ۱ قرار گیرد، روش‌های رگرسیون خطی ساده برای دستیابی به آن کفایت نمی‌کند، به این دلیل که آن‌ها به متغیر وابسته اجازه می‌دهند که از این محدودیت‌ها گذشته و نتایج ناسازگار تولید کنند. با تعریف P به عنوان احتمال تعلق یک مشاهده به گروه

وسعی در مسائل طبقه‌بندی پزشکی وقتی متغیر اصلی باینری (دو سطحی) باشد، مورد استفاده قرار می‌گیرد (۶). رگرسیون لجستیک و تحلیل جداسازی از جمله روش‌های آماری چند متغیره‌ای هستند که می‌توانند برای ارزیابی ارتباط بین متغیرهای مستقل هرچند مخدوش کننده و یک متغیر وابسته (دو یا چند سطحی) مورد استفاده قرار گیرد.

در مطالعه‌هایی که قبلاً در این زمینه انجام شده، ویلسون و همکاران در سال ۲۰۰۷ میلادی در مطالعه‌ای برای پیش‌بینی رخداد دیابت در افراد بالای ۵۰ سال، عوامل خطر ساز شامل سن بالا، دور کمر بالا، سابقه‌ی فامیلی دیابت، اختلال تحمل قند خون ناشتا، تری‌گلیسرید بالا و کلسترول HDL پایین را به‌عنوان متغیرهای پیش‌بینی کننده معرفی می‌کنند. در مطالعه‌ی دیگر یک‌هپارک و همکاران در سال ۱۹۹۸ میلادی انجام دادند، نژاد، چاقی، بیماری‌های قلبی، تری‌گلیسرید بالا و اختلال تحمل قند خون ناشتا و دو ساعته را به‌عنوان عوامل خطر ساز مؤثر در رخداد دیابت معرفی کرده‌اند (۷). در مطالعه‌ای که توسط جیمز پرس و ساندرا ویلسون (۱۹۷۸) انجام گرفته دو روش رگرسیون لجستیک و تحلیل جداسازی خطی را بدون در نظر گرفتن پیش‌فرض توزیع نرمال چند متغیره با هم مقایسه نمودند و دریافتند که روش رگرسیون لجستیک برآورد و نتایج بهتری - حتی با نمونه‌های نه چندان بزرگ - نسبت به روش تحلیل جداسازی می‌دهد (۸).

در مطالعه‌ای که توسط ماجا پوهار و همکاران (۲۰۰۴) انجام شده در مقایسه دو روش رگرسیون لجستیک (LR)^۱ و تحلیل جداسازی خطی (LDA)^۲ به این نتایج رسیده‌اند که با وجود شرایط نرمال بودن برای LDA این روش نتایج بهتری از LR ارائه می‌دهد و به هر حال برای حجم نمونه‌های بزرگ نتایج دو روش خیلی به هم نزدیک می‌شوند و هم‌چنین دریافتند که در روش LDA هرچه تعداد سطوح متغیر وابسته بیشتر باشد (از ۵ به بالا) توان پیش‌بینی بالاتر خواهد بود و بهتر از LR می‌باشد، اما این قضیه هنگامی که تعداد طبقات کم باشد متفاوت بوده و در واقع LR تنها زمان انتخاب مناسبی است که توزیع به صورت باینری یا دو طبقه‌ای باشد (۹).

مرتضی سدهی و همکاران با مطالعه‌ای که برای پیش‌بینی سندرم متابولیک انجام داده‌اند (۱۳۸۸). با مقایسه دو مدل

^۱ Logistic Regression; LR

^۲ Linear Discriminant Analysis; LDA

^۳ Receiver Operating Characteristic; ROC

از متغیرهای مستقل متمرکز است و می‌تواند یک مدل طبقه‌بندی را برای پیش‌بینی مشاهده‌های جدید به گروه‌های از پیش تعیین شده از آن به‌دست آورد. در ساده‌ترین نوع تحلیل جداسازی خطی یعنی تحلیل جداسازی خطی دو گروهی، یک تابع جداسازی خطی^۳ که از میان نقاط ثقل (مراکز هندسی) دو گروه می‌گذرد؛ می‌تواند برای جداسازی بین دو گروه مورد استفاده قرار گیرد (۱۰).

تابع جداسازی خطی به صورت زیر نمایش داده می‌شود:

$$LDF = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

که β_0 عرض از مبدأ و β_1 تا β_k ضرایب رگرسیونی k متغیر موجود هستند.

برای محاسبه‌ی احتمال این‌که یک مشاهده داده شده به یک گروه تعلق داشته باشد، از فاصله‌ی مایلانوبیس استفاده می‌شود. فاصله‌ی مایلانوبیس، فاصله‌ی مشاهده از مرکز ثقل گروه در فضای چند بعدی تعریف شده توسط متغیرهای پیش‌بینی کننده است. فاصله مایلانوبیس وقتی که متغیرهای فضای چند بعدی همبسته هستند؛ معیاری مناسب از فاصله می‌باشد. اگر متغیرهای پیش‌بینی کننده ناهمبسته باشند؛ از فاصله‌ی اقلیدسی استفاده می‌شود (۱۱).

وقتی که گروه‌ها (سطوح متغیر وابسته) بیش از دو گروه باشند، به «تعداد گروه‌ها منهای یک» تابع برای طبقه‌بندی یک مشاهده به میان آن‌ها مورد نیاز است. برای هر یک از گروه‌ها، تحلیل جداسازی خطی متغیرهای توضیحی را با توزیع نرمال و ماتریس کوواریانس‌های مساوی فرض می‌گیرد. برای هر مشاهده، ضریب برآورد شده برای یک متغیر مستقل در مقدار مشاهده‌ی آن متغیر ضرب می‌شود. این حاصل ضرب‌ها با یک مقدار ثابت جمع می‌شوند و نتیجه یک نمره ترکیبی خواهد بود که امتیاز جداسازی برای آن مشاهده می‌باشد (۱۰).

تکنیک‌های روش‌های تحلیل جداسازی و رگرسیون لجستیک دو مقوله متمایز از هم هستند و همان‌طور که انتظار می‌رود اگرچه آن‌ها باهم مرتبط‌اند، اما به طور کلی راه‌حل‌های پیشنهاد شده برای یکی متفاوت از دیگری است (۸).

هدف رگرسیون لجستیک پیدا کردن بهترین برازش و به صرفه‌ترین مدل برای توصیف ارتباط بین متغیر پیشامد (پاسخ یا وابسته) و مجموعه‌ای از متغیرهای مستقل (پیش‌بینی کننده یا توضیحی) است. روشی نسبتاً قوی انعطاف‌پذیر و با کاربردی آسان

بیمار و $1-P$ احتمال تعلق یک مشاهده به گروه غیر بیمار، مدل رگرسیون لجستیک به صورت زیر است:

$$\text{Log}\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k = \alpha + \sum_{i=1}^k \beta_i \cdot x_i$$

که p احتمال تخصیص مشاهده‌ای به سطحی از متغیر وابسته (گروه بیمار)، x_i آمین متغیر مستقل و β_i ضریب برآورد شده مدل برای آمین متغیر مستقل و β_0 یک عرض از مبدأ می‌باشد. پارامترهای β_0 تا β_k از مدل رگرسیون لجستیک توسط روش حداکثر درست‌نمایی برآورد شده‌اند (۱۰، ۴).

در تساوی بالا تبدیل logit برای مرتبط ساختن احتمالات عضویت گروه به یک تابع خطی از متغیرهای پیش‌بینی کننده مورد استفاده قرار گرفته است (۱۱).

روش حداکثر درست‌نمایی (ML) برای حداکثر کردن تابع درست‌نمایی داده‌های گرفته شده از برآورد پارامترها طراحی شده است. این الگوریتم‌ها با کم‌ترین مربعات رگرسیونی تفاوت دارد؛ زیرا داده‌های ورودی متغیر وابسته از نوع ترتیبی، طبقه‌ای و باینری هستند (۱۲، ۱۳).

از مزایای استفاده از مدل رگرسیون لجستیک علاوه بر مدل‌سازی مشاهده‌ها، امکان پیش‌بینی احتمال تعلق هر فرد به هر یک از سطوح متغیر وابسته و همچنین امکان محاسبه‌ی مستقیم نسبت شانس^۴ با استفاده از ضرایب مدل است (۴).

تحلیل جداسازی

از جمله روش‌های چند متغیره است که با تفکیک کردن مجموعه‌های متمایز مشاهده‌ها و نیز تخصیص مشاهده‌های جدید به گروه‌های تعریف شده که همان سطوح یا طبقات متغیر وابسته می‌باشند؛ سروکار دارد. معیاری که در این‌جا مطرح است، به وجود آوردن یک قانون یا تابع تشخیص بر مبنای اندازه‌های حاصل از مشاهده‌ها می‌باشد. مشاهده‌های جدیدی را که معلوم نیست از کدام گروه هستند به وسیله این تابع به‌دست آمده، می‌توان به یکی از سطوح یا گروه‌های متغیر وابسته نسبت داد. از معروف‌ترین توابع مورد استفاده در تحلیل جداسازی خطی می‌توان به تابع جداسازی فیشر اشاره کرد (۴).

تحلیل جداسازی خطی روی ارتباط بین متغیرهای مستقل چندگانه و یک متغیر وابسته چندسطحی توسط تشکیل ترکیبی

^۱ Maximum Likelihood; ML

^۲ Odds Ratio; OR

^۳ Linear Discriminant Function; LDF

ترکیب منحنی‌های راک دو مدل از نرم‌افزارهای MINITAB نسخه ۱۶، EASYFIT نسخه ۵/۵، SPSS نسخه ۲۰، و STATA نسخه ۱۱ استفاده شدند.

یافته‌ها

متغیر وابسته به دو سطح قند خون ناشتای کم‌تر و مساوی ۱۲۶ (غیر دیابتی) و قند خون ناشتای بیش‌تر از ۱۲۶ هم‌چنین افرادی که داروی دیابت مصرف می‌کنند (دیابتی) تقسیم شدند. فراوانی‌ها نشان داد که برای همه ۵۳۵۷ نفر، ۴۶۱۷ نفر غیر دیابتی (۸۶/۲ درصد) و ۷۴۰ نفر دیابتی (۱۳/۸ درصد) بودند که عدد ۰/۱۳۸ به عنوان نقطه‌ی برش احتمال برای محاسبه‌ی مدل‌های تحلیل جداسازی و رگرسیون لجستیک به عنوان احتمال پیشین در نرم‌افزار SPSS تعریف شد. در قدم نخست، مدل‌های رگرسیون لجستیک و تحلیل جداسازی به داده‌ها برازش داده شدند و متغیرهای تحصیلات، وضع شغل و وضع مصرف مواد با در نظر گرفتن سطح آخر به عنوان سطح مبدأ طبقه‌بندی شدند و مدل کامل ساخته شد. در مدل کامل با کنترل اثر متقابل متغیرهای HDL و تری‌گلیسرید، متغیرهای معنی‌دار با مقدار p کمتر یا مساوی از ۰/۰۵ استخراج شدند که شامل قد، محیط دور کمر، سن، جنس، اندازه فشار خون سیستولیک، مصرف داروی کاهنده فشار خون و LDL بودند.

پس از استخراج این متغیرهای معنی‌دار و ساختن مدل کاهش یافته، اثرات اصلی متغیرها و عوامل خطر معنی‌دار دیابت در جدول شماره ۱ نشان داده شده است.

برای مدل تحلیل جداسازی کاهش یافته با اعمال این مدل به متغیرهای معنی‌دار استخراج شده از مدل کامل، یک متغیر به نام «اسکور» به دست می‌آید. متغیر امتیاز از طریق ایجاد یک تابع تحلیل جداسازی ساخته می‌شود و می‌تواند ویژگی همه‌ی مشاهده‌ها باشد. بنابراین با به کار بستن مدل رگرسیون لجستیک به عنوان مدل تک‌متغیره بین امتیاز و متغیر وابسته، اثر امتیاز بر متغیر وابسته به صورت ویژه‌ای ساخته می‌شود.

ضریبی که از مدل تحلیل جداسازی به دست می‌آید با مدل رگرسیون لجستیک متفاوت خواهد بود. نتایج در جدول شماره ۲ گزارش شده‌اند.

جدول طبقه‌بندی به دست آمده در جدول شماره ۳، برای رگرسیون لجستیک حساسیت ۰/۷۴، ویژگی ۰/۷۱/۱، و دقت ۰/۷۱/۵ درصد و برای تحلیل جداسازی حساسیت ۰/۲۲/۴، ویژگی ۰/۹۵/۴، و

است. در رگرسیون لجستیک برخلاف تحلیل جداسازی خطی هیچ فرضیه‌هایی در توزیع متغیرها لازم نیست.

تحلیل جداسازی خطی در هر داده منحصر به فردی واریانس نسبت بین و داخل دو گروه را حداکثر می‌کند. بدین‌سان تفکیک‌پذیری حداکثری تضمین می‌شود. استفاده از تحلیل جداسازی خطی برای طبقه‌بندی داده‌ها برای مسأله طبقه‌بندی در بحث تشخیص به کار برده شده است (۱۴).

هدف تحلیل جداسازی خطی برای تعیین این‌که متغیر بین دو یا چند کلاس جداسازی شود، مورد استفاده قرار می‌گیرد و نیز برای به دست آوردن مدل طبقه‌بندی برای پیش‌بینی عضویت گروه‌ها برای مشاهده‌های جدید به کار می‌رود. به این معنی که مشاهده‌های جدید به کدام گروه (سطح متغیر وابسته) تخصیص یابند. برای هر کدام از گروه‌ها، تحلیل جداسازی خطی متغیرهای توضیحی را با توزیع نرمال و ماتریس کوواریانس مساوی در نظر می‌گیرد (۹).

از تحلیل راک به عنوان یک مقیاس اندازه‌گیری توانایی جداسازی یک مدل استفاده می‌شود که بیش‌ترین ناحیه‌ی زیر منحنی نشان دهنده‌ی توانایی پیش‌بینی بهتر برای مقایسه مدل‌هاست. در نهایت هرچه پیش‌بینی‌ها به واقعیت نزدیک‌تر باشد، مبنای تصمیم‌های صحیح‌تری قرار خواهند گرفت (۱۶، ۱۵). داده‌های پژوهش حاضر از یک نوع مطالعه آینده‌نگر در کرمان بین سال‌های ۹۰-۱۳۸۸ به دست آمد. متغیر دیابت (قند خون ناشتا) به عنوان متغیر وابسته در نظر گرفته شد. متغیرهای مستقل پس از حذف هم‌خطی و همبستگی بین برخی متغیرها شامل قد، محیط دور کمر، سن، جنس، شغل، تحصیلات، مصرف داروهای کاهنده فشار خون، اندازه فشار خون سیستولیک، HDL، LDL، وضع مصرف مواد مخدر، فعالیت‌هایی که باعث بالا رفتن ضربان قلب شود و تری‌گلیسرید بودند. اثر متقابل بین HDL و تری‌گلیسرید در فرایند تحلیل کنترل شدند. دقت، حساسیت، ویژگی و منحنی راک برای مقایسه‌ی قدرت پیش‌بینی مدل‌های تحلیل جداسازی و رگرسیون لجستیک محاسبه شدند. فرایند شبیه‌سازی به این صورت انجام گرفت که، توزیع هر متغیر را با توجه به پارامترهای آن که به‌طور عمده با پارامتر توزیع هایگاما و لگ نرمال نزدیک بود، توسط نرم‌افزار EASYFIT نسخه ۵/۵ برازش یافته، سپس از طریق نرم‌افزار MINITAB نسخه ۱۶ داده‌های مورد نظر با توزیع تعیین شده، تولید و شبیه‌سازی می‌شدند.

در کل برای محاسبه‌های یاد شده، ساختن فرایند شبیه‌سازی و

EASYFIT، داده‌ها که شامل ۱۰ سری با حجم ۲۰۰۰۰ نفر بودند، توسط نرم‌افزار MINITAB شبیه‌سازی شدند. برای هر مجموعه از مقادیر حساسیت، ویژگی، دقت و منحنی راک برای دو مدل محاسبه شدند. با میانگین‌گیری از نتایج ۱۰ مجموعه داده، نتایج در جدول شماره ۴ نشان داده شده است.

دقت ۸۵/۳ درصد به دست آمد. مقادیر منحنی راک در مدل رگرسیون لجستیک و تحلیل جداسازی به ترتیب ۸۰/۳ و ۸۰/۱ درصد بود.

شبیه‌سازی

با اطلاعات اصلی در ابتدای مطالعه و تمامی متغیرهای آن، پس از تعیین توزیع مورد نظر برای تولید داده‌ها توسط نرم‌افزار

جدول شماره ۱- مقادیر ضرایب، نسبت شانس، معنی‌داری و فواصل اطمینان (۹۵ درصد) برای متغیرهای استخراج شده مدل کاهش یافته رگرسیون لجستیک

رگرسیون لجستیک (مدل کاهش یافته)				مدل
CI (OR) 95%	p- value	OR	B	متغیرها
(۱/۰۰۰ و ۱/۰۲۸)	۰/۰۵۳	۱/۰۱۴	۰/۰۱۴	قد
(۱/۰۲۸ و ۱/۰۴۳)	۰/۰۰۰۱	۱/۰۳۶	۰/۰۳۵	محیط دور کمر
(۱/۰۳۹ و ۱/۰۵۶)	۰/۰۰۰۱	۱/۰۴۷	۰/۰۴۶	سن
(۱/۲۷۷ و ۲/۲۱۶)	۰/۰۰۰۱	۱/۶۴۸	۰/۴۹۹	جنسیت
(۱/۶۳ و ۲/۴۴)	۰/۰۰۰۱	۱/۹۹۴	۰/۶۹۰	مصرف داروی کاهنده فشار خون
(۱/۰۰۲ و ۱/۰۱۲)	۰/۰۰۲	۱/۰۰۷	۰/۰۰۷	سیستولیک
(۰/۹۹۵ و ۰/۹۹۹)	۰/۰۱۲	۰/۹۹۷	-۰/۰۰۳	LDL

جدول شماره ۲ - مقادیر ضریب، نسبت شانس، معنی‌داری و فواصل اطمینان (۹۵ درصد) برای متغیر امتیاز به دست آمده در مدل کاهش یافته تحلیل جداسازی

تحلیل جداسازی (مدل کاهش یافته)				مدل
CI (OR) 95%	p- value	OR	β	متغیر
(۲/۴۵۱ و ۲/۸۵۸)	۰/۰۰۰۱	۲/۶۴۶	۰/۹۷۳	امتیاز

جدول شماره ۳- جدول طبقه‌بندی مدل کاهش یافته برای رگرسیون لجستیک و تحلیل جداسازی (داخل پرانتز)

پیش‌بینی شده				مشاهده شده	
جمع	دیابت		دیابت		جمع
	خیر	بلی	خیر	بلی	
۴۶۱۷	۳۲۸۴ (۴۴۰۷)	۱۳۳۳ (۲۱۰)	۳۲۸۴ (۴۴۰۷)	۱۳۳۳ (۲۱۰)	دیابت
۷۴۰	۱۹۲ (۵۷۴)	۵۴۸ (۱۶۶)	۱۹۲ (۵۷۴)	۵۴۸ (۱۶۶)	بلی
۵۳۵۷					جمع

جدول شماره ۴- نتایج نهایی فرایند شبیه‌سازی برای دو مدل تحلیل جداسازی و رگرسیون لجستیک

رگرسیون لجستیک				تحلیل جداسازی				میانگین‌های ۱۰ سری داده
حساسیت	ویژگی	دقت	ROC	حساسیت	ویژگی	دقت	ROC	
۹۹/۱۸	۹۸/۴۹	۹۸/۵۹	۹۹/۹	۹۲/۶۲	۹۹/۱۹	۹۸/۲۶	۹۹/۵۶	

بحث

و ۰/۶۷۷ بود. ویژگی به ترتیب ۰/۸۵۷ و ۰/۶۶ به دست آمد و مساحت زیر منحنی راک (ROC) برای دو مدل به ترتیب ۰/۷۴۹ و ۰/۷۳۹ بود (۴) که با وجود مقداری تفاوت در دیگر شاخص‌ها، اما در پیش‌بینی منحنی راک، مطالعه حاضر دارای دقت بالاتری می‌باشد.

هم‌چنین در مطالعه‌ای که توسط آنتونوجرجوس و همکاران (۲۰۰۹) انجام شد در نهایت به این نتیجه رسیدند که به طور کلی تابع جداسازی خطی زمانی که فرضیه‌های نرمال بودن برقرار باشد، روشی بهتر از رگرسیون لجستیک می‌باشد و اختلاف‌های بین دو روش زمانی که حجم نمونه به اندازه کافی بزرگ باشد، بسیار ناچیز است (۱۰)، که این مطلب نیز با ادعای مطالعه حاضر هم‌خوانی دارد و همان‌طور که گزارش شد با افزایش حجم نمونه در فرایند شبیه‌سازی دقت پیش‌بینی مدل تحلیل جداسازی بالاتر رفته و به مدل رگرسیون لجستیک نزدیک شد که اگر فرضیه‌های نرمال بودن نیز برقرار باشد چه بسا از مدل رگرسیون لجستیک نیز بهتر نتیجه دهد.

در وضع فعلی جامعه و زندگی نوین اجتماعی و هم‌چنین ویژگی‌های رژیم‌های غذایی می‌توان با تغذیه‌ای مناسب به دور از فست‌فودها و غذاهای چاق کننده و هم‌چنین برنامه‌ریزی برای فعالیت‌های فیزیکی مستمر از جمله نرمش‌های بدنی و هوازی از افزایش وزن و بالارفتن کلسترول و اندازه فشار خون جلوگیری کرد، که این مسأله می‌تواند با توجه به نتایج مطالعه حاضر و مطالعه‌های گذشته به شکلی منجر به کاهش رشد ابتلا به دیابت در جامعه شود.

نتیجه‌گیری

بر اساس نتایج به دست آمده از داده‌های شبیه‌سازی مدل رگرسیون لجستیک دقت، حساسیت و منحنی راک بیش‌تر از تحلیل جداسازی بود، اما در تحلیل جداسازی ویژگی بهتر از رگرسیون لجستیک به دست آمد.

در نهایت، نتایج در مدل‌های کاهش یافته نشان داد که مقدار حساسیت در رگرسیون لجستیک بهتر از تحلیل جداسازی بود، اما مقادیر ویژگی و دقت در تحلیل جداسازی بهتر از رگرسیون لجستیک بود و منحنی راک در دو روش تقریباً یکسان بود.

در این پژوهش مشخص شد که مقدار حساسیت در مدل تحلیل جداسازی با افزایش حجم نمونه بهبود یافت و به طور کلی

متغیرهای پیش‌گوکننده‌ی معنی‌دار در این مطالعه شامل قد، محیط دور کمر، سن، جنس، اندازه فشار خون سیستولیک، مصرف داروی کاهنده فشار خون و LDL بودند که بر اساس ویژگی‌های هر متغیر بر دیگر متغیرها فعالیت‌های شدیدی داشته و روی دیابت مؤثرتر بوده‌اند.

با مقایسه متغیرهای معنی‌دار این مطالعه و مطالعه‌های قبلی، مشخص می‌شود که عوامل خطر و جدی در دیابت می‌تواند به طور مشترک شامل چاقی، تری‌گلیسرید، سطح LDL بالا و سن باشد که اثرهای اصلی و دیگر عوامل خطرآمیز شامل قد، محیط دور کمر بالا، اندازه فشار خون سیستولیک در این مطالعه سهم بیش‌تری به خود اختصاص دادند.

در مطالعه‌ای که توسط جیمز پرس و ساندر و ویلسون (۱۹۷۸) انجام شد که دو تحلیل رگرسیون لجستیک و تحلیل جداسازی، بدون در نظر گرفتن توزیع نرمال چندمتغیره پیش‌فرض مقایسه شدند و دریافتند که استفاده از تحلیل رگرسیون لجستیک توسط روش برآورد حداکثر درست‌نمایی برای برآورد تحلیل جداساز حتی بدون نمونه‌های خیلی بزرگ نتایج بهتری می‌دهد. در مطالعه‌ای دیگر که توسط ماجاپوهار و همکاران (۲۰۰۴) انجام شد، مقایسه‌ای بین رگرسیون لجستیک و تحلیل جداسازی خطی نشان داد که با وجود این‌که حالت‌های نرمال برای تحلیل جداسازی خطی در نظر گرفته شد، این روش نتایج بهتری از رگرسیون لجستیک ارائه داد. نتایج دو روش وقتی که حجم نمونه بزرگ بودند، به هم نزدیک بوده و تفاوت‌های اصلی می‌تواند در نمونه‌های کوچک مشاهده شود (۹ و ۸). با در نظر گرفتن نتایج به دست آمده از این پژوهش، با استفاده از داده‌های واقعی، دقت تحلیل جداسازی نسبتاً بهتر از رگرسیون لجستیک بود هرچند در فرایند شبیه‌سازی با افزایش حجم نمونه، نتایج دو روش به هم نزدیک شدند.

در محاسبه‌هایی که برای دو مدل در این مطالعه انجام شد، طبق جدول شماره ۳ دقت، حساسیت، ویژگی و منحنی راک برای رگرسیون لجستیک به ترتیب ۰/۷۴، ۰/۷۱/۱، ۰/۷۱/۳ و ۰/۸۰/۳ درصد و برای تحلیل جداسازی به ترتیب ۰/۸۵/۳۶، ۰/۲۲/۴۳، ۰/۸۵/۳۶ و ۰/۸۰/۱ درصد به دست آمدند که در مقایسه با مطالعه‌ای که توسط مرتضی سدهی و همکاران (۱۳۸۸) انجام گرفته بود، میزان حساسیت برای مدل‌های رگرسیون لجستیک و تحلیل جداسازی به ترتیب ۰/۴۸۳

دانشگاه علوم پزشکی کرمان برای در اختیار گذاردن داده‌ها تشکر به عمل آورند. این مقاله بخشی از پایان‌نامه کارشناسی ارشد آمار زیستی بوده و بخشی از آن با حمایت مرکز تحقیقات مدل‌سازی در سلامت دانشگاه علوم پزشکی کرمان انجام شده است.

می‌توان نشان داده شود که در فرایند شبیه‌سازی هرچه حجم نمونه بزرگ باشد، نتایج مدل‌های رگرسیون لجستیک و تحلیل جداسازی به هم نزدیک‌تر خواهد بود.

تشکر و قدردانی

نگارندگان بر خود لازم می‌دانند که از مرکز تحقیقات فیزیولوژی

منابع

1. Meraci M, Feizi A, BagherNejad M. Investigating the Prevalence of High Blood Pressure, Type 2 Diabetes Mellitus and Related Risk Factors According to a Large General Study in Isfahan- Using Multivariate Logistic Regression model. *Journal of Health Systems Research*. 1391; 8: 193-203.
2. Esmaeilnasab N, Afkhamzadeh A, Ebrahimi A. Sectional study of risk factors in type 2 diabetes in the Diabetes Center Sanandaj. *Iran Research Journal of Epidemiology*. 1389; 6: 39-45.
3. Zabetian A, Hadaeq F, Harati H and Azizi F. Anthropometric indices forecasts of incident type 2 diabetes in adults: Tehran Lipid and Glucose Study. *Journal of Diabetes and Lipid Disorders*. 1384; 5: 143-51.
4. Sedehi M, Mehrabi Y, Kazemnejad A, Hadaegh F. Comparison of artificial neural network, logistic regression and discriminant analysis methods in prediction of metabolic syndrome. *Iranian journal of endocrinology and metabolism (IJEM)*. 2010; 11: 638-46.
5. Osborne J, Waters E. Four assumptions of multiple regression that researchers should always test. *Practical Assessment, Research & Evaluation*. 2002; 8: 1-9.
6. Maroco J, Silva D, Rodrigues A, Guerreiro M, Santana I and Mendonça A. Data mining methods in the prediction of Dementia: A real-data comparison of the accuracy, sensitivity and specificity of Linear Discriminant Analysis, Logistic Regression, Neural Networks, Support Vector Machines, Classification Trees and Random Forests. *BMC Research Notes* 2011, 4:299.
7. Mehrabi Y, Khadem-Maboudi A, Sarbakhsh P and Hadaegh F. Prediction of Diabetes Using Logic Regression. *Iranian Journal of Endocrinology and Metabolism*. 2010; 12: 16-24.
8. Press SJ, Wilson S. Choosing between logistic regression and discriminant analysis. *Journal of the American Statistical Association*. 1978; 699-705.
9. Pohar M, Blas M and Turk S. Comparison of logistic regression and linear discriminant analysis: a simulation study. *Metodolski Zvezki*. 2004; 1: 143-61.
10. Antonogeorgos G, Panagiotakos DB, Priftis KN and Tzonou A. Logistic regression and linear discriminant analyses in evaluating factors associated with asthma prevalence among 10-to 12-years-old children: divergence and similarity of the two statistical methods. *International Journal of Pediatrics*. Volume 2009, Article ID 952042, 6 pages.
11. Worth AP, Cronin MTD. The use of discriminant analysis, logistic regression and classification tree analysis in the development of classification models for human health effects. *Journal of Molecular Structure: THEOCHEM*. 2003; 622: 97-111.
12. Peng CYJ, Lee KL and Ingersoll GM. An introduction to logistic regression analysis and reporting. *The Journal of Educational Research*. 2002; 96: 3-14.
13. Harris DV, Pan G. Mineral favorability mapping: a comparison of artificial neural networks, logistic regression, and discriminant analysis. *Natural Resources Research* 1999; 8: 93-109.
14. Balakrishnama S, Ganapathiraju A. Linear discriminant analysis-a brief tutorial. *Institute for Signal and information Processing*. 1998
15. Abdolmaleki P, Yarmohammadi M and Gity M. Comparison of logistic regression and neural network models in predicting the outcome of biopsy in breast cancer from MRI findings. *Iranian Journal of Radiation Research*. 2004; 1: 217-28.
16. Makian N, Almodaresi M, Karimi T. Comparison of artificial neural network models with logistic regression and discriminant analysis methods in predict of Companies Bankruptcy. *Journal of Economic Research*. 1389; 10: 141-61.

Comparison of Logistic Regression and Discriminant Analysis in Predicting Type 2 Diabetes

Aram Ahmadi M¹, Bahrampour A²

1- Department of Biostatistics and Epidemiology, Faculty of Health, Kerman University of Medical Sciences, Kerman, Iran

2- Research Center for Modeling in Health, Institute for Future Studies in Health, Department of Biostatistics and Epidemiology, Faculty of Health, Kerman University of Medical Sciences, Kerman, Iran

Corresponding author: Bahrampour A, abahrampour@yahoo.com

Background and Objectives: Diabetes is a chronic and common metabolic disease which has no curative treatment. Logistic regression (LR) is a statistical model for the analysis and prediction in multivariate statistical techniques. Discriminant analysis is a method for separating observations in terms of dependent variable levels which can allocate any new observation after making discriminating functions. The aim of this study was to compare and determine the effective variables in type 2 diabetes.

Methods: The data included 5357 persons obtained through a cohort study in Kerman, southeastern Iran, in 2009-11. Diabetes was considered the response variable. The independent variables after deleting collinearity and correlated variables included height, waist circumference, age, gender, occupation, education, drugs, systolic blood pressure, HDL, LDL, drug abuse, activities, and triglyceride. Sensitivity, specificity, accuracy, and ROC curve were applied for determining and comparing the prediction power of the models.

Results: The results in the reduced model with extracted significant variables from the full model, the sensitivity of the LR model and DA was 74% and 22.4%, the specificity of the LR model and DA was 71.1 % and 95.4 %, the prediction accuracy of the LR model and DA was 71.5% and 85.3%, and the ROC curve of the LR model and DA was 80.3% and 80.1%, respectively. Simulation showed the sensitivity, specificity, accuracy, and ROC curve was 99.18%, 98.49%, 98.59%, and 99.9% for the LR model and 92.62%, 99.19%, 98.26%, and 99.56% for DA, respectively.

Conclusion: The results showed that the risk factors of diabetes in the logistic regression reduced model were waist circumference, age, gender, LDL level, systolic pressure, and drugs. Also, the sensitivity of the LR model was more than DA while DA had a higher specificity and prediction accuracy. Comparison of the ROC curve showed that the prediction estimated values were rather similar in both models, but the two models were the same asymptotically.

Keywords: Sensitivity, Specificity, Logistic regression, Discriminant analysis, ROC curve