

مقایسه نتایج مدل رگرسیون لجستیک معمولی و لجستیک نیرومند در مدل‌بندی عوامل مرتبط با بیماری پره‌دیابت

چکیده

دریافت: ۱۳۹۷/۰۱/۲۰ ویرایش: ۱۳۹۷/۰۱/۲۷ پذیرش: ۱۳۹۷/۰۷/۱۵ آنلاین: ۱۳۹۷/۰۷/۳۰

زمینه و هدف: با توجه به اینکه خطر ابتلا به دیابت در افراد پره‌دیابتیک بسیار بالا است، تعیین عوامل موثر بر پره‌دیابت دارای اهمیت می‌باشد. این مطالعه با هدف مقایسه نتایج مدل رگرسیون لجستیک معمولی و لجستیک نیرومند در مدل‌بندی عوامل مرتبط با بیماری پره‌دیابت انجام شد.

روش بررسی: این مطالعه که از نوع مقطعی-تحلیلی است روی ۶۴۶۰ نفر از افراد بالای ۳۰ سال، شرکت‌کننده در طرح غربالگری دیابت دانشگاه علوم پزشکی مشهد، از مهر تا آذر ۱۳۸۹ انجام شد. با توجه به میزان قندخون ناشتای افراد، ۵۴۱۴ نفر سالم و ۱۰۴۶ نفر به‌عنوان پره‌دیابتیک شناسایی شدند. سن، جنس، نمایه توده بدن، فشارخون سیستولیک، فشارخون دیاستولیک و نسبت کمر به باسن در مورد هر فرد اندازه‌گیری شد. مدل رگرسیون لجستیک معمولی روی داده‌ها برازش شد. سپس داده‌های پرت مشخص و سه مدل نیرومند Mallow، WBV و BY برازش شد. آنگاه مدل‌ها با هم مقایسه گردیدند.

یافته‌ها: متغیرهای سن، نمایه توده بدن و فشارخون سیستولیک در همه مدل‌ها از لحاظ آماری معنادار شدند ($P < 0.01$) و متغیر نسبت کمر به باسن معنادار نشد ($P > 0.1$). تعداد ۵۵۲ داده‌ی پرت با خطای بدرده‌بندی در مدل معمولی وجود داشت. مقادیر کای دو پیروسون و سطح زیرمنحنی راک در مدل Mallow به‌طور تقریبی فرقی با مدل معمولی نداشت. اما در مدل‌های WBV و BY به‌نسبت بیشتر بود.

نتیجه‌گیری: با توجه به نتایج این پژوهش سن بالا، نمایه توده بدنی و فشارخون بالا در ابتلا به بیماری پره‌دیابت موثر می‌باشند. همچنین مدل‌های رگرسیون نیرومند WBV و BY برازش بهتر و توان پیشگویی بالاتری نسبت به رگرسیون لجستیک معمولی در مدل‌بندی عوامل گفته‌شده در ارتباط با پره‌دیابت دارند.

کلمات کلیدی: نمایه توده بدنی، دیابت شیرین، مدل‌های لجستیک، مرحله پره‌دیابت.

فاطمه عزیزی مایوان^۱

مهدی جباری نوقابی^۲، علی تقی‌پور^۳
محمدتقی شاکری^{۴*}، مهسا مکرّم^۵

۱- گروه بهداشت عمومی، دانشگاه علوم پزشکی نیشابور، نیشابور، ایران.

۲- گروه آموزشی آمار، دانشکده علوم ریاضی، دانشگاه فردوسی، مشهد، ایران.

۳- گروه اپیدمیولوژی، مرکز تحقیقات عوامل اجتماعی موثر بر سلامت، دانشگاه علوم پزشکی مشهد، مشهد، ایران.

۴- گروه آمار زیستی، مرکز تحقیقات عوامل اجتماعی موثر بر سلامت، دانشگاه علوم پزشکی مشهد، مشهد، ایران.

۵- گروه جمعیت‌شناسی، کمیته تحقیقات دانشجویی، دانشگاه آزاد اسلامی، واحد تهران مرکزی، تهران، ایران.

* نویسنده مسئول: مشهد، خیابان دانشگاه، دانشکده بهداشت عمومی.
تلفن: ۰۵۱-۳۸۵۱۱۱۶

E-mail: shakerim@mums.ac.ir

مقدمه

بزرگسال بالای ۱۸ سال، ۹٪ برآورد شد.^۲ در سال ۲۰۱۲ تعداد مرگ‌هایی که در دنیا به‌طور مستقیم بر اثر دیابت رخ داده‌اند، ۱/۵ میلیون مورد برآورد شد. بیشتر از ۸۰٪ موارد مرگ ناشی از دیابت در کشورهای با درآمد کم یا متوسط رخ داده است.^۳ سازمان جهانی بهداشت پیش‌بینی کرده است که دیابت هفتمین علت مرگ در سال ۲۰۳۰ خواهد بود.^۴ شیوع دیابت در ایران در سال ۲۰۱۱، ۱۱/۱٪ برآورد شد که این میزان تا سال ۲۰۳۰ به ۱۲/۸٪ خواهد رسید.^۱ این بیماری خطر ابتلا به بیماری

دیابت یکی از شایع‌ترین بیماری‌های مزمن در تقریباً تمامی کشورها است و به‌خاطر تغییر الگوی زندگی مردم به سمت تحرک کمتر و چاقی بیشتر همچنان روند افزایشی دارد. دیابت نوع دو شایع‌ترین نوع دیابت در دنیا است و به‌طور تقریبی ۹۰٪ بیماران دیابتی را شامل می‌شود.^۱ شیوع جهانی دیابت در سال ۲۰۱۴ در میان افراد

مشهد از نیمه‌ی مهر تا اواخر آذر ۱۳۸۹ در دانشگاه علوم پزشکی مشهد انجام گرفت. تمامی افراد بالای ۳۰ سال و داوطلب ساکن در بافت مرکزی شهر مشهد وارد مطالعه شده و افراد مبتلا به دیابت از مطالعه خارج گردیدند. از مجموع داده‌های حاصل پس از مدیریت داده‌ها و حذف داده‌های ناقص حجم جمعیتی بالغ بر ۶۴۶۰ نفر در مطالعه کنونی وارد شدند. این طرح پژوهشی با کد اخلاق IR.MUMS.REC.1395.363 در دانشگاه علوم پزشکی مشهد به تصویب رسیده است.

از آنجایی که نتایج حاصل از مدل رگرسیون لجستیک در حجم نمونه‌های بالا دقیق‌تر می‌باشد، تمامی داده‌های مناسب و کامل این مطالعه مورد تحلیل آماری قرار گرفت. تمامی اندازه‌گیری‌ها شامل سن و قد، وزن، جنس، فشارخون سیستولیک، فشارخون دیاستولیک، قندخون ناشتا، دور کمر و دور باسن در مورد هر فرد اندازه‌گیری و ثبت شد. منظور از دور کمر، محیط دور کمر واقع در قسمت بالای ناف و دور باسن، محیط دور پهن‌ترین و برجسته‌ترین قسمت بدن در پایین خط کمر می‌باشد. به‌منظور استفاده از شاخص‌های چاقی مناسب به‌جای قد و وزن، از نمایه توده بدن (BMI) که ترکیبی از هر دو است، استفاده شد و به‌جای دور کمر و دور باسن از شاخص نسبت دور کمر به باسن استفاده گردید. میزان قندخون ناشتا به‌عنوان معیار تشخیص پره‌دیابت در نظر گرفته شد. به این صورت که افراد با قندخون ناشتای $100-125$ mg/dl به‌عنوان پره‌دیابتیک و کمتر از ۱۰۰ به‌عنوان سالم شناسایی شدند. برای انتخاب مدل، ارتباط آماری تک‌تک متغیرهای موردنظر با ابتلا به پره‌دیابت را با استفاده از Chi-square test و Independent samples t-test در سطح معناداری $P < 0.05$ بررسی کردیم. سپس مدل رگرسیون لجستیک معمولی را با استفاده از متغیرهای معنادار شده، روی داده‌ها برازش دادیم. آنگاه داده‌های پرت با توجه به شاخص باقیمانده‌ی استیودنتیته‌ی پیرسون مشخص گردیدند. سه مدل نیرومند Mallow، WBLY و BY که بر مبنای نیرومندتر کردن روش MLE است،^{۱۲} روی داده‌ها برازش شد. سطح معناداری $P < 0.05$ در نظر گرفته شد. سپس این مدل‌ها با استفاده از معیار Chi-square test پیرسون از لحاظ نیکویی برازش مورد مقایسه قرار گرفتند. همچنین سطح زیرمنحنی راک برای مقایسه‌ی توان پیش‌بینی مدل‌ها محاسبه گردید. نرم‌افزار مورد استفاده R Project for Statistical Computing، Version R 3.1.2 (www.r-project.org) بود.

قلبی و سکنه را افزایش می‌دهد. همچنین باعث صدمه‌ی کلیه، ایجاد زخم در پا و نیاز احتمالی به قطع عضو و آسیب در شبکیه چشم می‌شود.^۶ خطر مرگ در افراد دیابتی حداقل دو برابر افراد غیردیابتی است.^۷ پره‌دیابت یا پیش از دیابت مرحله‌ای است که در آن اگرچه قند خون بالاتر از حد نرمال است اما آنقدر بالا نیست که بتوان آن را در دسته‌ی دیابت نوع ۲ به‌حساب آورد. در این مرحله، قند خون ناشتا بین $100-125$ mg/dl و یا دو ساعت پس از آزمایش تحمل گلوکز با $140-199$ mg/dl می‌باشد. براساس گزارش انجمن دیابت آمریکا برخی از عوامل خطرزای دیابت نوع ۲ عبارتند از داشتن اضافه وزن، سن بیشتر از ۴۵ سال، عدم فعالیت بدنی، پیشینه‌ی خانوادگی دیابت، پرفشاری خون (فشارخون بالاتر یا مساوی mmHg $140/90$ یا دریافت درمان پرفشاری خون)، کلسترول HDL سرم برابر یا کمتر از 35 mg/dl و یا تری‌گلیسرید سرم برابر یا بالاتر از 175 mg/dl.^۸ احتمال اینکه افراد مبتلا به پره‌دیابت بدون رعایت پیشگیری طی ده سال یا کمتر مبتلا به دیابت شوند، بسیار بالا است مگر اینکه شیوه‌ی زندگی خود را با تلاش برای کاهش وزن، رعایت رژیم غذایی مناسب و تحرک بیشتر، تغییر دهند.^۹ بنابراین شناسایی افراد مبتلا به پره‌دیابت و تعیین الگوی ابتلا به این بیماری جهت درمان و پیشگیری از ابتلا به آن بسیار با اهمیت می‌باشد. انجام این کار نیازمند به روش‌های مناسب آماری است. رگرسیون لجستیک یک روش استاندارد برای مدل‌سازی و بیان ارتباط بین یک متغیر پاسخ کیفی با یک یا چند متغیر توضیحی کیفی یا کمی است. در واقع تفاوت اصلی رگرسیون لجستیک با رگرسیون خطی شرط کیفی بودن متغیر پاسخ است.^{۱۰}

روش معمول برآورد پارامترها در مدل رگرسیون لجستیک، روش ماکزیمم درست‌نمایی (MLE) است. اما برآورد پارامترها در این روش بسیار تحت تاثیر وجود داده‌های پرت قرار می‌گیرد و ممکن است در حضور این داده‌ها بسیار غیر واقعی ارایه شود.^{۱۱} این مطالعه با هدف مدل‌بندی عوامل مرتبط با بیماری پره‌دیابت با استفاده از مدل‌های رگرسیون لجستیک معمولی و نیرومند و مقایسه نتایج آن‌ها انجام گرفت.

روش بررسی

این مطالعه به‌صورت مقطعی - تحلیلی بر روی داده‌های مطالعه‌ای انجام گرفت که تحت عنوان غربالگری دیابت در سطح کلان شهر

یافته‌ها

سیستولیک و نسبت دور کمر به دور باسن، مقادیر شاخص تورم واریانس محاسبه شد. این مقادیر همگی نزدیک به یک بودند که بیانگر عدم وجود هم‌خطی قابل توجه بین متغیرها است. برای بررسی خطی بودن رابطه‌ی متغیرهای توضیحی معنادار شده در مدل که همگی به صورت کمی هستند، Lowess smoothed logit در برابر هر یک از این متغیرها رسم و خطی بودن رابطه‌ی آن‌ها تایید شد. آزمون نیکویی برازش Chi-square نیز برای این مدل معنادار شد (P=۰/۰۰۱).

در مدل نهایی رگرسیون لجستیک معمولی متغیر عرض از مبدا و متغیرهای توضیحی سن، نمایه توده بدن و فشارخون سیستولیک معنادار شدند (P=۰/۰۰۱)، اما متغیر کمر به باسن معنادار نشد (P=۰/۱۵). اندازه‌ی ضرایب رگرسیونی برآوردشده از این مدل در جدول ۲ آمده است.

در میان ۶۴۶۰ نفر شرکت‌کننده در این مطالعه ۴۵۰۷ نفر (۶۹/۸٪) مرد و ۱۹۵۳ نفر (۳۰/۲٪) زن بودند که ۱۰۴۶ نفر (۱۶/۲٪) از مجموع آنان پرهدیابتیک و ۵۴۱۴ نفر (۸۳/۸٪) سالم بودند. میانگین سن و ویژگی‌های بالینی افراد سالم و پرهدیابتیک در جدول ۱ آمده است. Chi-square test حاکی از همگنی دو گروه دیابتیک و پرهدیابتیک به لحاظ جنس بود (P=۰/۴۳۱). اختلاف میانگین‌های هر یک از متغیرهای توضیحی کمی در افراد سالم و پرهدیابتیک با توجه به Independent samples t-test معنادار شد (P=۰/۰۰۱). در مدل رگرسیون لجستیک چندگانه متغیر فشارخون دیاستولیک را به‌خاطر معنادار نشدن (P=۰/۲۸۱) از مدل حذف شد. برای بررسی هم‌خطی میان متغیرهای توضیحی سن، نمایه توده بدن، فشارخون

جدول ۱: مقایسه توزیع میانگین متغیرهای فردی و بالینی بیماران پرهدیابتیک با افراد سالم

متغیرها	سالم (انحراف معیار ± میانگین)	پرهدیابتیک (انحراف معیار ± میانگین)	P*
سن (سال)	۴۳/۹۰ ± ۱۰/۳۵	۴۶/۴۰ ± ۱۰/۴۶	۰/۰۰۱
نمایه توده بدن (kg/m ²)	۲۶/۵۲ ± ۴/۱۳	۲۹/۱۷ ± ۴/۳۱	۰/۰۰۱
فشارخون سیستولیک (mmHg)	۱۱۵/۸۹ ± ۱۴/۱۸	۱۲۰/۰۹ ± ۱۵/۹۵	۰/۰۱۳
فشارخون دیاستولیک (mmHg)	۷۴/۲۶ ± ۹/۹۰	۷۶/۵۳ ± ۱۰/۴۳	۰/۸۰۵
نسبت کمر به باسن	۰/۹۲ ± ۰/۰۹	۰/۹۳ ± ۰/۰۸	۰/۱۵۶

* گزارش شده از مدل رگرسیون لجستیک چندگانه، P < ۰/۰۵ معنادار در نظر گرفته شد.

جدول ۲: برآورد ضرایب رگرسیون در مدل‌های مختلف مورد مطالعه

متغیرها	مدل رگرسیون لجستیک معمولی (انحراف معیار ± ضریب رگرسیونی)	مدل Mallow (انحراف معیار ± ضریب رگرسیونی)	مدل WBY (انحراف معیار ± ضریب رگرسیونی)	مدل BY (انحراف معیار ± ضریب رگرسیونی)
سن (سال)	۰/۰۱۹۱ ± ۰/۰۰۳۴	۰/۰۱۹۱ ± ۰/۰۰۳۴	۰/۰۲۰۴ ± ۰/۰۰۳۳	۰/۰۱۹۱ ± ۰/۰۰۳۳
نمایه توده بدن (kg/m ²)	۰/۱۳۹۲ ± ۰/۰۰۸۴	۰/۱۴۰۱ ± ۰/۰۰۸۴	۰/۱۴۶۳ ± ۰/۰۰۸۶	۰/۱۳۸۷ ± ۰/۰۰۸۴
فشارخون سیستولیک (mmHg)	۰/۰۰۸۲ ± ۰/۰۰۲۵	۰/۰۰۸۲ ± ۰/۰۰۲۵	۰/۰۰۶۸ ± ۰/۰۰۲۶	۰/۰۰۸۱ ± ۰/۰۰۲۵
نسبت کمر به باسن	-۰/۶۱۱۳ ± ۰/۴۲۴۸	-۰/۵۷۶۹ ± ۰/۴۲۵۷	-۰/۴۷۷۵ ± ۰/۴۲۲۸	-۰/۶۱۱۶ ± ۰/۴۲۳۴

جدول ۳. مقادیر Chi-square پیرسون و سطح زیرمنحنی راک برای مدل رگرسیون لجستیک معمولی و سه مدل نیرومند Mallow, WBY و BY

مدل رگرسیون لجستیک	معمولی	Mallow	WBY	BY
Chi-square پیرسون	۶۴۸۳	۶۴۷۶	۶۰۸۶	۶۰۴۵
سطح زیرمنحنی راک (درصد)	۶۷/۸۱	۶۸/۸۲	۷۲/۲۶	۷۲/۳۳

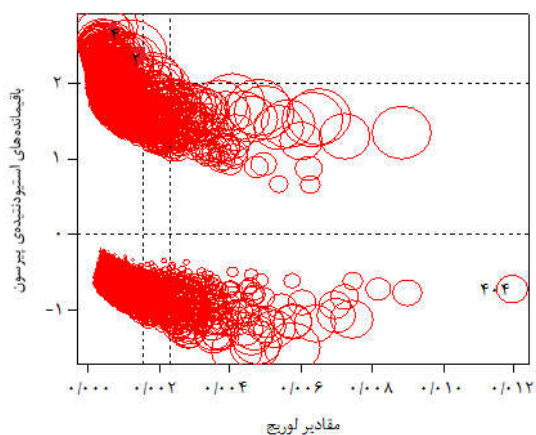
در جدول ۳ معیار Chi-square پیرسون برای مقایسه‌ی نیکویی برازش و سطح زیرمنحنی راک برای مقایسه‌ی توان پیشگویی مدل‌ها به کار رفته است. با توجه به جدول، مقدار Chi-square پیرسون برای هر سه مدل نیرومند کمتر از مدل معمولی است. در مدل Mallow این مقدار هفت واحد کمتر از مدل معمولی است. در مدل WBY، ۳۹۷ واحد و در مدل BY، ۴۳۵ واحد کمتر از مدل معمولی است. مقدار سطح زیرمنحنی راک برای مدل Mallow کمابیش برابر با مدل معمولی است. اما برای مدل‌های WBY و BY ۴-۳٪ بیشتر است.

بحث

در همه مدل‌های برازش شده در این مطالعه متغیرهای سن، نمایه توده بدن و فشارخون سیستمولیک معنادار شدند که این نتیجه منطبق با گزارش سالانه‌ی انجمن دیابت آمریکا در سال ۲۰۱۵ در رابطه با عوامل خطرزای دیابت و پره‌دیابت می‌باشد.^۸

در مطالعه Díaz-Redondo و همکاران که در میان جامعه اسپانیایی انجام شده است، ارتباط فشارخون با پره‌دیابت تایید شد.^{۱۳} در پژوهش کنونی متغیر نسبت کمر به باسن معنادار نشد درحالی‌که در مطالعات دیگر اثر این متغیر در ابتلا به دیابت تایید شده است.^{۱۴، ۱۵} در بسیاری از مطالعاتی که تاکنون در زمینه‌ی دیابت و پره‌دیابت انجام شده است، از رگرسیون لجستیک استفاده شده است.^{۱۸-۱۶} اما مطالعه‌ای که در این زمینه از رگرسیون لجستیک نیرومند استفاده کرده باشد، یافت نشد.

تعداد ۵۵۲ داده‌ی پرت از نوع بدرده‌بندی در مدل رگرسیون لجستیک معمولی وجود داشت. داده‌های بدرده‌بندی شده دارای لوریج کوچک و باقیمانده استیون‌تیده‌ی بزرگ می‌باشند.^{۱۲} تمام داده‌های پرت مربوط به افراد پره‌دیابتیک بود. علت آن می‌تواند وجود برخی عوامل دیگر مرتبط با پره‌دیابت باشد که در مطالعه کنونی وارد نشده



نمودار ۱: نمودار تاثیر

نمودارهای هموارشده‌ی Lowess مربوط به مقادیر احتمال برآوردشده از مدل در برابر باقیمانده‌های استیون‌تیده‌ی پیرسون و استیون‌تیده‌ی انحراف رسم شد، که در آن‌ها شیب خط به‌طور معناداری غیرصفر شد. این نکته بیانگر عدم کفایت برازش مدل و اثر معنادار داده‌های پرت در مدل بود.

در نمودار تاثیر که در نمودار ۱ رسم شده است، محور افقی آن مربوط به مقادیر لوریج (Leverage) و محور عمودی آن مربوط به مقادیر باقیمانده‌های استیون‌تیده‌ی پیرسون و اندازه‌ی دایره‌ها متناسب با مقادیر فاصله‌ی کوک (Cook's distance) مشاهده‌ی مربوطه می‌باشد.

مراحل برازش مدل‌های رگرسیون لجستیک نیرومند مانند مدل رگرسیون لجستیک معمولی انجام شد. خطی بودن رابطه‌ی متغیرهای توضیحی معنادار شده با لوجیت‌ها در هر سه مدل بررسی و تایید گردید. نتایج حاصل از برازش سه مدل نیرومند Mallow, WBY و BY به همراه نتایج مدل رگرسیون لجستیک معمولی در جدول ۲ آورده شده است.

در آن مطالعه، مدل WBY و سپس مدل Mallow عملکرد بهتری نسبت به بقیه مدل‌ها داشتند.^{۱۲}

در رگرسیون لجستیک، نقاط پرت خطرناک از نوع نقاط نافذ بد می‌باشد.^{۱۳} این داده‌ها را می‌توان به این صورت که مقدار لوریج و باقیمانده‌ی آن‌ها هر دو بزرگ باشد، تشخیص داد. با توجه به نمودار تاثیر رسم‌شده در شکل ۲، چنین داده‌ای در مطالعه کنونی وجود نداشت و همه داده‌های پرت از نوع تنها بدرده‌بندی بودند. بنابراین این انتظار وجود داشت که مدل نیرومند تفاوت خیلی زیادی با مدل معمولی نداشته باشد. همچنین با توجه به مطالعه Sarkar و همکاران با افزایش حجم نمونه، کنترل داده‌های پرت اثر کمتری در برآورد پارامترها دارد.^{۱۴} بنابراین حجم نمونه‌ی زیاد عامل دیگری برای چندان بهتر نشدن مدل نیرومند نسبت به مدل معمولی بوده است. با توجه به نتایج این پژوهش، سن بالا، نمایه توده بدنی و فشارخون بالا در ابتلا به بیماری پره‌دیابت موثر می‌باشند. همچنین مدل‌های رگرسیون نیرومند WBY و BY برازش بهتر و توان پیشگویی بالاتری نسبت به مدل رگرسیون لجستیک معمولی در مدل‌بندی عوامل گفته‌شده در ارتباط با پره‌دیابت دارند.

سپاسگزاری: این مقاله حاصل پایان‌نامه تحت عنوان "مقایسه نتایج مدل رگرسیون لجستیک معمولی و مدل رگرسیون لجستیک نیرومند در بیماران پره‌دیابتیک" در مقطع کارشناسی ارشد در سال ۱۳۹۵ و کد IR.MUMS.REC.1395.363 می‌باشد که با حمایت دانشگاه علوم پزشکی و خدمات بهداشتی درمانی مشهد اجرا شده است.

است. Hosmer و همکاران نیز در مطالعه خود یکی از عوامل وجود خطای بدرده‌بندی در رگرسیون لجستیک را حذف حداقل یک متغیر توضیحی مهم در مدل بیان کرده‌اند.^{۱۹}

در مطالعات متعددی روش‌های مختلف رگرسیون لجستیک نیرومند معرفی و ارزیابی شده است.^{۲۰-۲۲} که در پژوهش کنونی سه تا از روش‌های ارزیابی شده برای نیرومندتر کردن روش MLE، یعنی روش‌های Mallow، WBY و BY برای برازش روی داده‌ها و مقایسه با روش رگرسیون لجستیک معمولی، انتخاب شد. اندازه‌ی ضرایب در این سه مدل تفاوت چندانی با مدل رگرسیون لجستیک معمولی نداشت. با توجه به مقادیر Chi-square پیروسون، مدل Mallow از لحاظ نیکویی برازش تفاوت چندانی با مدل معمولی نداشت، درحالی‌که دو مدل WBY و BY مدل‌های به‌نسبت بهتری بودند. همچنین قدرت پیشگویی مدل Mallow با توجه به سطح زیرمنحنی راک به‌طور تقریبی تفاوتی با مدل معمولی نداشت. اما دو مدل WBY و BY توان پیشگویی بهتری داشتند (۴-۳٪ بیشتر). با توجه به این نتایج، استفاده از مدل Mallow برای این داده‌ها توصیه نمی‌شود. اما دو مدل WBY و BY را می‌توان به مدل رگرسیون لجستیک معمولی ترجیح داد. در مطالعه Ahmad و همکاران مدل رگرسیون لجستیک معمولی و چهار مدل نیرومند Mallow، شوپه، WBY و BY، با استفاده از داده‌های شبیه‌سازی‌شده و داده‌های واقعی مورد مقایسه قرار گرفتند و عملکرد آن‌ها در حضور انواع داده‌های پرت ارزیابی شد که

References

- Whiting DR, Guariguata L, Weil C, Shaw J. IDF diabetes atlas: global estimates of the prevalence of diabetes for 2011 and 2030. *Diabetes Res Clin Pract* 2011;94(3):311-21.
- Mendis S, Davis S, Norrving B. Organizational update: the world health organization global status report on noncommunicable diseases 2014; one more landmark step in the combat against stroke and vascular disease. *Stroke* 2015;46(5):e121-2.
- World Health Organization (WHO). Global Health Estimates: Deaths by Cause, Age, Sex and Country, 2000-2012. Geneva: WHO; 2014.
- Mathers CD, Loncar D. Projections of global mortality and burden of disease from 2002 to 2030. *PLoS Med* 2006;3(11):e442.
- World Health Organization (WHO). Global status report on noncommunicable diseases 2010. Description of the global burden of NCDs, their risk factors and determinants. [Internet] Geneva: WHO; 2011 [cited 2018 Aug 15]. Available from: http://www.who.int/nmh/publications/ncd_report2010/en/
- Morrish NJ, Wang SL, Stevens LK, Fuller JH, Keen H. Mortality and causes of death in the WHO Multinational Study of Vascular Disease in Diabetes. *Diabetologia* 2001;44 Suppl 2:S14-21.
- Roglic G, Unwin N, Bennett PH, Mathers C, Tuomilehto J, Nag S, et al. The burden of mortality attributable to diabetes: realistic estimates for the year 2000. *Diabetes Care* 2005;28(9):2130-5.
- American Diabetes Association. Classification and diagnosis of diabetes. *Diabetes Care* 2015;38(Suppl 1):S8-16.
- Alberti KG, Zimmet PZ. Definition, diagnosis and classification of diabetes mellitus and its complications. Part 1: diagnosis and classification of diabetes mellitus provisional report of a WHO consultation. *Diabet Med* 1998;15(7):539-53.
- Cabrera AF. Logistic regression analysis in higher education: An applied perspective. In: Smart JC, editor. Higher Education: Handbook of Theory and Research. Vol. 10. New York, NY: Agathon Press; 1994. P. 225-56.
- Kalina J. Highly robust methods in data mining. *Serbian J Manag* 2013;8(1):9-24.

12. Ahmad S, Ramli NM, Midi H. Robust estimators in logistic regression: a comparative simulation study. *J Modern Appl Stat Methods* 2010;9(2):18.
13. Diaz-Redondo A, Giráldez-García C, Carrillo L, Serrano R, García-Soidán FJ, Artola S, et al. Modifiable risk factors associated with prediabetes in men and women: a cross-sectional analysis of the cohort study in primary health care on the evolution of patients with prediabetes (PREDAPS-Study). *BMC Fam Pract* 2015;16:5.
14. Hadaegh F, Bozorgmanesh M, Safarkhani M, Khalili D, Azizi F. Predictability of body mass index for diabetes: Affected by the presence of metabolic syndrome? *BMC Public Health* 2011;11(1):1.
15. Wei W, Xin X, Shao B, Zeng FF, Love EJ, Wang BY. The relationship between anthropometric indices and type 2 diabetes mellitus among adults in north-east China. *Public Health Nutr* 2015;18(9):1675-83.
16. Aram-Ahmaddi M, Bahrampour A, Najafipour H. Predicting type two diabetes and determination of effectiveness of risk factors applying logistic regression model. *J Kerman Univ Med Sci* 2014;21(2):103-13.
17. Heikes KE, Eddy DM, Arondekar B, Schlessinger L. Diabetes risk calculator. *Diabetes Care* 2008;31(5):1040-5.
18. Tabaei BP, Herman WH. A multivariate logistic regression equation to screen for diabetes: development and validation. *Diabetes Care* 2002;25(11):1999-2003.
19. Hosmer Jr DW, Lemeshow S. *Applied Logistic Regression*. New York, NY: John Wiley & Sons, Inc.; 2004.
20. Feng J, Xu H, Mannor S, Yan S, editors. *Robust Logistic Regression and Classification. Adv Neural Inf Process Syst*; 2014.
21. Croux C, Haesbroeck G. Implementing the Bianco and Yohai estimator for logistic regression. *Comput Stat Data Anal* 2003;44(1):273-95.
22. Carroll RJ, Pederson S. On robustness in the logistic regression model. *J Roy Stat Soc B (Methodological)* 1993;693-706.
23. Copas J. Binary regression models for contaminated data. *J Roy Stat Soc B (Methodological)* 1988;50:225-65.
24. Sarkar SK, Midi H, Rana M. Detection of outliers and influential observations in binary logistic regression: An empirical study. *J Appl Sci* 2011;11(1):26-35.

Comparison of ordinary logistic regression and robust logistic regression models in modeling of pre-diabetes risk factors

Fateme Azizi Mayvan M.Sc.¹
Mehdi Jabbari Nooghabi Ph.D.²
Ali Taghipour Ph.D.³
Mohammad Taghi Shakeri
Ph.D.^{4*}
Mahsa Mokarram Ph.D. Student⁵

1- Department of Public Health,
Neyshabur University of Medical
Sciences, Neyshabur, Iran.

2- Department of Statistics, School
of Mathematics, Ferdowsi University,
Mashhad, Iran.

3- Department of Epidemiology,
Social Determinants of Health
Research Center, Mashhad University
of Medical Sciences, Mashhad,
Iran.

4- Department of Biostatistics,
Social Determinants of Health
Research Center, Mashhad University
of Medical Sciences, Mashhad,
Iran.

5- Department of Demographics,
Student Research Committee, Islamic
Azad University, Central Tehran
Branch, Tehran, Iran.

* Corresponding author: School of Public
Health, Daneshgah St., Mashhad, Iran.
Tel: +98 51 38515116
E-mail: shakerim@mums.ac.ir

Abstract

Received: 09 Apr. 2018 Revised: 16 Apr. 2018 Accepted: 07 Oct. 2018 Available online: 22 Oct. 2018

Background: Regarding the increased risk of developing type 2 diabetes in pre-diabetic people, identifying pre-diabetes and determining of its risk factors seems so necessary. In this study, it is aimed to compare ordinary logistic regression and robust logistic regression models in modeling pre-diabetes risk factors.

Methods: This is a cross-sectional study and conducted on 6460 people, over 30 years old, who have participated in the screening of diabetes plan in Mashhad city that it was done by Mashhad University of Medical Sciences from October to December 2010. According to the fasting blood sugar criteria, 5414 individuals were identified as healthy and 1046 individuals were identified as pre-diabetic. Age, gender, body mass index, systolic blood pressure, diastolic blood pressure and waist-to-hip ratio were measured for every participant. The data was entered into the Microsoft Excel 2013 (Microsoft Corp., Redmond, WA, USA) and then analysis of the data was done in R Project for Statistical Computing, Version R 3.1.2 (www.r-project.org). Ordinary logistic regression model was fitted on the data. The outliers were identified. Then Mallows, WBY and BY robust logistic regression models were fitted on the data. And then, the robust models were compared with each other and with ordinary logistic regression model according to goodness of fit and prediction ability using Pearson's chi-square and area under the receiver operating characteristic (ROC) curve respectively.

Results: Among the variables that were included in the ordinary logistic regression model and three robust logistic models, age, body mass index and systolic blood pressure were statistically significant ($P < 0.01$) but waist-to-hip ratio was not statistically significant ($P > 0.1$). There were 552 outliers with misclassification error in the ordinary logistic regression model. Pearson's chi-square value and area under the ROC curve value in the Mallows model were almost the same as for ordinary logistic regression model. But it was relatively higher in BY and WBY models.

Conclusion: Based on results of this study age, overweight and hypertension are risk factors of prediabetes. Also, WBY and BY models were better than ordinary logistic regression model, according to goodness of fit criteria and prediction ability.

Keywords: body mass index, diabetes mellitus, logistic models, prediabetic state.