

به کارگیری متغیرهای پنهان در مدل رگرسیون بجستیک برای حذف اثر همخطی چندگانه در تحلیل برخی عوامل مرتبط با سرطان پستان

محمد امین پورحسینقلی^۱، یدا... محابی^۲، حمید علوی جد^۳، پروین یاوری^۴

^۱کارشناس ارشد آمار زیستی، دانشگاه علوم پزشکی شهید بهشتی، تهران.

^۲دانشیار آمار زیستی، دانشگاه علوم پزشکی شهید بهشتی، تهران.

^۳استادیار آمار زیستی، دانشگاه علوم پزشکی شهید بهشتی، تهران.

^۴استاد اپیدمیولوژی، دانشگاه علوم پزشکی شهید بهشتی، تهران.

نویسنده‌ی رایط: دکتر یدا... محابی، گروه پزشکی اجتماعی و بهداشت دانشکده پزشکی، دانشگاه علوم

پزشکی شهید بهشتی، تهران، اوین، کد پستی ۱۹۳۹۰، تلفن: ۰۲۱-۲۳۸۷۷۵۶۷-۸

نمبر: ۰۲۱-۲۲۴۱۴۱۰۸، پست الکترونیک: mehrabi@sbmu.ac.ir ; ymehrabi@gmail.com

تاریخ دریافت: ۸۵/۲/۲، پذیرش: ۸۴/۱۰/۲۶

مقدمه و اهداف: رگرسیون بجستیک یکی از کاربردی‌ترین مدل‌های خطی تعمیم‌یافته برای تحلیل رابطه‌ی یک یا چند متغیر توضیحی بر متغیر پاسخ رسته‌ای است. زمانی که بین متغیرهای توضیحی همبستگی‌های نسبتاً قوی وجود داشته باشد همخطی چندگانه اجاد شده، ممکن است به کاهش کارآیی مدل منجر شود. هدف این تحقیق استفاده از متغیرهای پنهان برای کاهش اثر همخطی چندگانه در تحلیل یک مطالعه مورد - شاهدی است.

روش کار: داده‌های مورد استفاده در این تحقیق متعلق به یک مطالعه مورد - شاهدی است که در آن ۳۰۰ نفر زن مبتلا به سرطان پستان با ۳۰۰ زن شاهد از نظر عوامل خطر مورد مقایسه قرار گرفتند. برای بررسی اثر همخطی، پنج متغیر کمی که بین آن‌ها همبستگی بالای وجود داشت، در نظر گرفته شدند. ابتدا مدل بجستیک به متغیرهای فوق برآرازش داده شد. سپس به منظور حذف اثر همخطی، دو متغیر پنهان با استفاده از هرکدام از دو روش تحلیل عاملی و تحلیل مؤلفه‌های اصلی به دست آورده، بر مبنای آن‌ها پارامترهای مدل‌های بجستیک جدد آخاسبه شدند. کارآیی مدل‌ها، با استفاده از خطای استاندارد پارامترها مقایسه گردید.

نتایج: مدل رگرسیون بجستیک برآساس متغیرهای اولیه حاکی از مقادیر غیرعادی نسبت شانس برای سن در اولین زایمان زنده ($OR=67960$) و سن در اولین حاملگی ($OR=0/00029$) بود. درحالی که پارامترهای مدل‌های بجستیک حاصل از متغیرهای پنهان به دست آمده از هر دو روش تحلیل عامل و تحلیل مؤلفه‌های اصلی، از نظر آماری معنی دار ($p<0.003$) و خطای استاندارد همه‌ی آن‌ها کوچکتر از خطای استاندارد مربوط به رگرسیون بجستیک معمولی بود. فاکتورها و مولفه‌های اصلی تولید شده توسط دو روش حداقل ۸۵ درصد کل واریانس را تبیین کردند.

نتیجه‌گیری: تحقیق نشان داد اخراج اسitanدارد پارامترهای برآورد شده در رگرسیون بجستیک برآساس متغیرهای پنهان از رگرسیون بجستیک برآساس مشاهدات اولیه کوچکتر بوده و در نتیجه این‌گونه مدل‌بندی در تحلیل برخی عوامل خطر سرطان پستان که همخطی چندگانه، متغیر پنهان، تحلیل عاملی، تحلیل مؤلفه‌های

واژگان کلیدی: همخطی چندگانه، متغیر پنهان، تحلیل عاملی، تحلیل مؤلفه‌های اصلی، رگرسیون بجستیک، سرطان پستان.

مقدمه

متغیرهای توضیحی افزایش می‌یابد، مدل‌سازی مشکل شده و کارآیی آن نیز کاهش می‌یابد؛ به خصوص اگر برخی از متغیرها علی‌رغم فرض استقلال در مدل‌سازی، با یکدیگر همبستگی قوی

رگرسیون بجستیک یکی از کاربردی‌ترین مدل‌های خطی تعمیم‌یافته است که برای تحلیل رابطه‌ی یک یا چند متغیر توضیحی بر متغیر پاسخ رسته‌ای به کار می‌رود (۱). زمانی که تعداد

استفاده از ماتریس مقادیر ویژه (Eigen Values) (۱۰)، مؤلفه‌های اصلی به صورت ترکیب خطی از متغیرهای اولیه و مستقل از یکدیگر ساخته می‌شوند و در آنالیز داده‌ها، به جای متغیرهای اولیه مورد استفاده قرار می‌گیرند (۱۱).

تحلیل عاملی (Factor Analysis) از دیگر روش‌های کاهش ابعاد داده‌ها است، که خستین بار توسط اسپیرمن (Spearman) معرفی شد. در این روش با فرض وجود یک مدل مبنایی مشخص برای کل داده‌ها، و براساس ماتریس واریانس-کوواریانس یا ماتریس ضرایب همبستگی، عامل‌های مستقل از یکدیگر، از روی متغیرهای اولیه به دست می‌آیند (۱۲). در این تحقیق، برای کاهش ابعاد مدل رگرسیون جستیک با متغیرهای توضیحی هم خط، در تحلیل داده‌های یک مطالعه مورد - شاهدی پیرامون عوامل خطرسرطان پستان از دو روش تحلیل عاملی و تحلیل مؤلفه‌های اصلی استفاده شده است.

روش‌ها

برای بررسی خواهی کاربرد دو روش تحلیل مؤلفه‌های اصلی و تحلیل عاملی در کاهش ابعاد مدل و ایجاد متغیرهای پنهان، از داده‌های مطالعه‌ی مورد - شاهدی مربوط به عوامل خطر سرطان پستان استفاده شد (۱۳). در مطالعه‌ی مذکور که در فاصله‌ی زمانی بهمن ۸۲ تا آذر ۸۳ در مرکز پزشکی - آموزشی، درمانی شهدای تجریش انجام شد گروه مورد، بیمارانی بودند که بیماری سرطان پستان آن‌ها با استفاده از آزمایش‌های پاتولوژیک، تشخیص قطعی داده شده و یا برای درمان یا پیگیری به درمانگاه بیمارستان شهدای تجریش مراجعه کرده بودند. گروه شاهد زنانی بودند که به دلایل دیگری غیر از سرطان پستان و به طور هم‌زمان در چشم‌های دیگر بیمارستان شهداء، مثل جراحی، پوست، داخلی و غیره بسته و یا برای پیگیری یا درمان به درمانگاه بیمارستان مراجعه کرده بودند و از نظر سی با گروه مورد با حداقل ۲ سال اختلاف مشابه‌سازی شدند. با اطمینان ۹۵ درصد و توان

داشته باشد و به عبارت دیگر هم خطی چندگانه (Multicollinearity) ایجاد شده باشد (۲). هم خطی چندگانه یکی از دلایل افزایش خطای استاندارد برآورد ضرایب رگرسیونی و درنتیجه کاهش کارآیی مدل بوده و ممکن است منجر به پیش‌بینی‌هایی خارج از دامنه‌ی مورد انتظار شود (۳).

مسئله‌ی هم خطی در مدل‌های رگرسیون خطی مورد توجه بسیاری از محققان قرار گرفته و روش‌های گوناگونی برای مقابله با اثرات نامطلوب آن ابداع شده است (۲). از جمله این روش‌ها، کاهش ابعاد مدل با استفاده از متغیرهای پنهان (Latent Variables) است. این نوع متغیرها مستقیماً مشاهده نمی‌شوند؛ بلکه از ترکیب سایر متغیرهای مشاهده شده قابل دستیابی بوده، به عنوان نماینده‌ی برخی از متغیرهای همبسته در مدل به کار می‌روند (۴).

اگر چه استفاده از متغیرهای پنهان برای کاهش ابعاد مدل، در عمل بیشترین کاربرد را در مطالعات مربوط به علوم اجتماعی و روانشناسی داشته است، ولی به دلیل نوع مطالعات انجام شده در علوم پزشکی و بهداشت که مستلزم جمع‌آوری تعداد قابل توجهی متغیرهای مرتبط با یکدیگر است، مشکل هم خطی در بسیاری از مدل‌های آماری این مطالعات قابل انتظار است (۵) و علیرغم این‌که هم خطی چندگانه در مدل رگرسیون جستیک نیز ایجاد مشکل می‌کند (۶، ۷)، تاکنون توجه محققان بیشتر بر رگرسیون خطی با متغیر پاسخ دارای توزیع نرمال متمرکز بوده است.

تحلیل مؤلفه‌های اصلی (Principal Component Analysis) یکی از کاربردی‌ترین روش‌های کاهش ابعاد در روش‌های چند متغیری است. تاریخچه‌ی ابداع این روش به ابداعات پیرسون (Pearson) در برآشندگان مربعات متعامد برمی‌گردد؛ ولی بسط عمدی تئوری به وسیله‌ی هتلینگ (Hotteling) انجام شده است (۸). مؤلفه‌های اصلی با توجه به خصوصیاتی که دارند برای مقابله با مشکل هم خطی و کاهش ابعاد مدل در رگرسیون‌های خطی مورد استفاده قرار می‌گیرند (۲، ۹). در این روش با

در اولین حاملگی و سن در اولین زایمان زنده همبستگی خطی بالای مشاهده می‌شود. همچنین مقادیر عامل تورم و اریانس (VIF: Variance Inflation Factor) نیز که میزان بروز همخطی را نشان می‌دهد، محاسبه شد که همه مقادیر بیشتر از یک و نشان دهنده وجود همخطی در متغیرهاست. به ویژه اینکه دو متغیر سن در اولین حاملگی و سن در اولین زایمان زنده دارای VIF بیشتر از ۱۵ بوده و بنابراین برآساس این معیار همخطی شدیدی دارند. آزمون بارتلت نیز نشان داد ماتریس ضرایب همبستگی متغیرهای توضیحی با صفر اختلاف معنی‌داری دارد ($P<0.01$). بنابراین بین متغیرهای توضیحی مورد بررسی همخطی چندگانه وجود دارد.

نتایج حاصل از رگرسیون جستیک بدون در نظر گرفتن وجود این همخطی (جدول ۲)، نشان میدهد که فقط دو متغیر سن در اولین حاملگی و سن در اولین زایمان زنده معنیدار شده اند ($P<0.001$). همچنین نسبت شانس به دست آمده برای سن در اولین زایمان زنده بسیار بزرگ ($OR=67960$) و بر عکس برای سن در اولین حاملگی بسیار کوچک ($OR=0.00029$) بود که هردو غیرعادی هستند. در مرحله‌ی بعد با روش تحلیل عاملی دو عامل به صورت ترکیب خطی زیر از متغیرهای اولیه بدست آمد:

$$\text{Factor1} = 0.85\text{NP} + 0.97\text{NLB} + 0.68\text{TLBF} - 0.26$$

$$\text{AFP} - 0.26\text{AFLB}$$

$$\text{Factor2} = 0.21\text{NP} - 0.23\text{NLB} - 0.25$$

$$\text{TLBF} + 0.94\text{AFLB} + 0.95\text{AFLB}$$

در جمیع $\%84/79$ واریانس توسط عامل اول ($\%45/67$) و عامل دوم ($\%39/12$) تبیین شده است. وارد کردن عوامل فوق به عنوان متغیرهای توضیحی در مدل رگرسیون جستیک، نشان داد که عامل اول ($P<0.002$) و عامل دوم ($P<0.001$) هر دو تاثیر معنیداری بر متغیر وابسته دارند (جدول ۳). در مرحله‌ی سوم، با روش تحلیل مؤلفه‌های اصلی دو مؤلفه به صورت ترکیب خطی زیر از متغیرهای اولیه بدست آمد:

$$\text{Component1} = 0.90\text{NP} + 0.92\text{NLB} + 0.82\text{TLBF}$$

$$- 0.25\text{AFP}$$

$$- 0.24\text{AFLB}$$

آزمون ۸۰ درصد، تعداد نمونه برای هر گروه ۳۰۰ نفر انتخاب شد (۱۴). متغیرهای مختلفی به عنوان عوامل خطر یا متغیر کنترل جمع‌آوری شدند که در این مقاله، پنج متغیر که بین آن‌ها همبستگی بالای وجود داشت در نظر گرفته شدند. این متغیرها عبارتند از: تعداد حاملگی (NP: Number of Pregnancy)، تعداد فرزندان زنده به دنیا آورده (NLB: Number of Live Birth)، کل طول مدت شیردهی به فرزندان (TLBF: Total Length of Breast Feeding) حاملگی (Afp: Age at First Pregnancy) و سن AFLB: Age at First Live Birth). برای بررسی میزان بروز همخطی در این مشاهدات از ماتریس ضرایب همبستگی استفاده است (۲).

ابتدا بدون در نظر گرفتن وجود همخطی، مدل رگرسیون جستیک به داده‌ها برآذش داده شد. سپس با ترکیب پنج متغیر مورد بررسی، یک بار با روش تحلیل مؤلفه‌های اصلی و بار جدول ۱- ماتریس ضرایب همبستگی متغیرهای توضیحی

	NP	NLB	TLBF	AFP	AFLB	
۱	۰/۹۰	۰/۶۷	-۰/۴۴	-۰/۴۲	NP	
		۰/۷۵	۰/۴۷	-۰/۰۴۷	NLB	
			-۰/۴۲	-۰/۴۱	TLBF	
				۰/۹۷	AFP	
					AFLB	

NP: تعداد حاملگی؛ NLB: تعداد فرزندان زنده به دنیا آورده؛ TLBF: کل طول مدت شیردهی به فرزندان؛ AFP: سن در اولین حاملگی؛ AFLB: سن در اولین زایمان زنده

دیگر به طریق تحلیل عاملی، دو متغیر پنهان به دست آمد و بر اساس آن‌ها، پارامترهای مدل رگرسیون جستیک برآورد شد. مدل‌های حاصل، برآسانه واریانس‌های تبیین شده به وسیله‌ی دو روش و خطای استاندارد پارامترهای برآورد شده، مورد مقایسه قرار گرفتند.

یافته‌ها

جدول یک نشان میدهد بین سه متغیر تعداد حاملگی، تعداد فرزندان زنده به دنیا آورده و کل طول مدت شیردهی به فرزندان و نیز بین دو متغیر سن

جدول ۳- برآورد پارامترهای رگرسیون جستیک مدل در هر دو روش تحلیل مؤلفه‌های بر اساس متغیرهای پنهان اجتاد شده به وسیله روش تحلیل عاملی و مؤلفه‌های اصلی

متغیرهای اصلی رگرسیونی استاندارد	ضرایب خطای استاندارد	P-value	(فاصله اطمینان) شانس	نسبت
روش تحلیل عاملی عرض از مبداء	-۰/۰۴	۰/۰۸	۰/۶۴	۰/۹۶
عامل ۱	-۰/۲۷	۰/۰۹	۰/۰۰۲	(۰/۶۴)
عامل ۲	۰/۲۹	۰/۰۹	۰/۰۰۱	(۱/۱۲)
روش تحلیل مؤلفه‌های اصلی عرض از مبداء	-۰/۰۴	۰/۰۸	۰/۶۴	۰/۹۶

Component2= -۰/۲۱NP-۰/۲۵NLP-۰/۲۵
TLBF+۰/۹۶AFP+۰/۹۶AFLB

مؤلفه‌ی اول %۴۹/۲۶ و مؤلفه‌ی دوم %۸۹/۲۶ کل واریانس را تبیین می‌کنند. در برآورد پارامترهای رگرسیون جستیک براساس این مؤلفه‌های اصلی، نتایج مشابهی حاصل شد. به عبارت دیگر مؤلفه اول با ($P<0/02$) و مؤلفه دوم با ($P<0/003$) معنیدار شدند و نسبت شانس‌ها نیز تقریباً مشابه روش تحلیل عاملی به دست آمدند (جدول ۳).

بحث

هدف پژوهش حاضر استفاده از متغیرهای پنهان در مدل رگرسیون جستیک به منظور کاهش اثر در حالت بروز همخطی چندگانه بود. نتایج حاصل نشان داد برآورد پارامترهای

جدول ۲- برآورد پارامترهای مدل رگرسیون جستیک بدون درنظر گرفتن همخطی بین متغیرهای توضیحی

متغیرهای اصلی	ضرایب رگرسیونی استاندارد	خطای استاندارد	P-value	(فاصله اطمینان) شانس	نسبت	عامل تورم واریانس * (VIF)
عرض از مبداء	۰/۱۰	۰/۲۰	۰/۶۲	-۴۵۳۵۰۳) (۱۰۱۸۴	-	-
تعداد حاملگی (NP)	۰/۰۶	۰/۰۹	۰/۵۰	(۰/۸۹-۱/۲۲) ۱/۰۶	۰/۰۷	۴/۵۷
تعداد فرزندان زنده به دنیا آورده (NLB)	-۰/۱۷	۰/۲۶	۰/۵۱	(۰/۵۰-۱/۴۰) ۰/۸۴	۰/۸۳	۵/۸۳
کل طول مدت شیردهی به فرزندان (TLBF)	-۰/۲۲	۰/۱۵	۰/۱۴	(۰/۵۹-۱/۰۸) ۰/۸۰	۰/۱۳	۲/۱۳
سن در اولین حاملگی (AFP)	-۱۰/۴۶	۰/۹۲	<۰/۰۰۱	۰/۰۰۰۲۹	۰/۰۰۰۲۹	۱۰/۴۱
سن در اولین زایمان زنده (AFLB)	۱۱/۱۲	۰/۹۷	<۰/۰۰۱	-۴۵۳۵۰۳) (۱۰۱۸۴	۶۷۹۶۰	۱۰/۳۴

نتایج مشابهی داشته، نسبت به مدل جستیک با متغیرهای همخط اولیه از کارآیی بالاتری برخوردار هستند.

تشکر و قدردانی

در این مقاله از داده‌های طرح تحقیقاتی عوامل خطر سرطان پستان، مصوب دانشگاه علوم پزشکی شهید بهشتی استفاده شده است که به این وسیله از کلیه همکاران طرح مذکور و نیز از معنویت پژوهشی دانشکده پزشکی سپاسگزاری به عمل می‌آید.

منابع

1. Myers R.H., Montgomery D.C. and Vining G.G., Generalized linear models with application in engineering and sciences, 2002, John Wiley & Sons.
2. Chatterjee, S., Hadi, A.S. and Price, B. (2000). Regression analysis by example, 2002, John Wiley & Sons, USA. PP: 225-258.
3. Myers, R.H. (1990). Classical and modern regression with applications., 1990, Pws-Kent publishing company. PP: 123-129.
4. Van Eye, A., Clogg, C.C., Latent variables analysis; application for developing research. 1994, SAGE publication. PP: 3-35.
5. Hazard munro, B. Statistical methods for health care research., 2001, Philadelphia: Lippincott. PP: 287-288.
6. Kleinbaum, D. Logistic Regression., 1994, Springer, New York. PP: 168.
7. Hosmer, D.W., Lemeshow, S. Applied logistic regression., 1989, John Wiley & Sons.
8. Morrison, D. F. Multivariate statistical methods. 2002, John Wiley & Sons. PP: 312-398.
9. Rawlings, J. O. Applied regression analysis: A research tools., 1988, Belmont: Wadsworth. PP: 327-356.
10. Schott, J. R. Matrix analysis for statistics., 1997, John Wiley & Sons. PP: 84-131.
11. Jolliffe, I.T. PrinCI95%pal component analysis., 1986, Springer. PP: 129-141.
12. Srivastava, M. S. Methods of multivariate statistics, 2002, John Wiley & Sons. New York. PP: 397-450.
13. Yavari, P., Mousavizadeh, M., Sadrol-Hafezi, B. and Mehrabi, Y., Reproductive characteristics and the risk of breast cancer, A case-control study. Asian PaCI95%fic J Cancer Prev, 2005, 6, 370-375.
14. Lemeshow, S., Hosmer, D. W. and Klar, J. Adequacy of sample size in Health studies. World Health Organization, 1998, John Wiley & Sons. PP: 19.
15. Aguilera, A.M. and Escabias, M., PrinCI95%pal component logistic regression. Proceedings in computational statistics, 2000, 175-180. Physica-Verlag.
16. Escabias, M., Aguilera, A. M. and Valderrama, M. J., Modeling climatological data by functional logistic regression. The ISI International Conference on Environmental Statistics and Health, 2003.
17. Wall, M. M. and Li, R., Comparison of multiple regression to two latent variable techniques for estimation and prediction. Statistics in Medicine; 2003, 22:3671-3685.
18. Sobel, M. E. Causal inference in latent variable models. In Latent variables analysis; application for developing research. By Van Eye, A., Clogg, 1994, C.C SAGE publication. PP: 3-35.

حاصل از دو روش تحلیل عامل و معادله‌ی مدل‌سازی ساختاری (Structural Equation Modeling) را با رگرسیون کلاسیک برآسانس متغیرهای اولیه همخط مقایسه کرده نشان دادند که متغیرهای پنهان، پارامترهایی با خطاهای استاندارد کوچکتر تولید می‌کنند (۱۷). نتایج تحقیق حاضر از این حاظ با مطالعه آنان هموارانی دارد.

ایده‌ی استفاده از متغیرهای پنهان به جای متغیرهای اصلی، با هدف کاهش ابعاد داده‌ها از این حقیقت ناشی می‌شود که این متغیرها می‌توانند بازتابدهنده‌ی ارتباط بین مشاهدات باشند (۱۸). با این حال استفاده از مدل‌های دربرگیرنده‌ی متغیرهای پنهان تبعاً مزايا و محدودیت‌هایی دارد. یکی از اهداف اصلی در ساختن مدل‌های آماری تفسیر مدل با توجه به پارامترهای برآوردشده می‌باشد؛ ولی تفسیر مدل‌هایی که برآسانس عامل‌ها یا تحلیل مؤلفه‌های اصلی به دست می‌آیند قدری پیچیده است (۱۱، ۱۹). برای این کار استفاده از روش تحلیل مؤلفه‌های اصلی بهتر از تحلیل عاملی است؛ زیرا مؤلفه‌های اصلی صرفاً ترکیبی خطی از متغیرهای اولیه هستند و بر خلاف روش تحلیل عامل، مدلی برای داده‌ها فرض نمی‌کند (۱۱، ۲۰، ۱۹). در نتیجه از طریق معکوس ماتریس دوران می‌توان برآوردهای اولیه را به دست آورد (۱۱). هم‌چنین روش‌هایی نیز برای تفسیر این مؤلفه‌ها در مدل کاهشیافته پیشنهاد شده است (۲۱). به هر حال سودمندی‌های حاصل از کاهش ابعاد مدل و کاستن تعداد متغیرها آنچنان قابل ملاحظه است که علی‌رغم مشکلات حاصل در تفسیر پارامتر، برخی تکنیک‌های جدید علاوه بر تولید متغیرهای پنهان برای متغیرهای توضیحی، اکنون بر تولید این متغیرها برای متغیرهای پاسخ توجه دارند (۲۲).

نتیجه‌گیری

برآسانس یافته‌های این تحقیق می‌توان نتیجه‌گیری کرد که در بررسی برخی عوامل خطر سرطان پستان، دو روش تحلیل عاملی و تحلیل مؤلفه‌های اصلی

- 2002, <http://ace.acadiau.ca/math/chipmanh/publications.html>.
22. Guo, J., Wall, M. M. and Amemiya Y. Latent class regression on latent factors to appear in Biostatistics .
19. Rencher, A. C. *Methods of multivariate analysis*, 2002, John Wiley & Sons.
20. Armitage, P. and Colton, T., *Encyclopedia of Biostatistics*. Volume 2. Chichester: 1998, John Wiley & Sons. PP: 1480-1481.
21. Chipman HA and Gu H. Interpretable dimension reduction.

εγ /