

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

محمد امین پورحسینقلی^۱، یدا... محرابی^۲، حمید علوی مجد^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... محرابی، گروه پزشکی اجتماعی و بهداشت دانشکده پزشکی، دانشگاه علوم پزشکی شهید بهشتی، تهران، اوین، کد پستی ۱۹۳۹۵، تلفن: ۰۲۱-۲۳۸۷۲۵۶۷-۸

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

تاریخ دریافت: ۱۰/۰۴/۸۴، پذیرش: ۲۳/۰۲/۸۵

مقدمه و اهداف: رگرسیون لجستیک یکی از کاربردی‌ترین مدل‌های خطی تعمیم‌یافته برای تحلیل رابطه‌ی یک یا چند متغیر توضیحی بر متغیر پاسخ رسته‌ای است. زمانی که بین متغیرهای توضیحی همبستگی های نسبتاً قوی وجود داشته باشد هم خطی چندگانه ایجاد شده، ممکن است به کاهش کارآیی مدل منجر شود. هدف این تحقیق استفاده از متغیرهای پنهان برای کاهش اثر هم خطی 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 هم خطی، پنج متغیر کمی که بین آن‌ها همبستگی بالایی وجود داشت، در نظر گرفته شدند. ابتدا مدل لجستیک به متغیرهای فوق برآش داده شد. سپس به منظور حذف اثر هم خطی، دو متغیر پنهان با استفاده از هر کدام از دو روش تحلیل عاملی و تحلیل مؤلفه‌های اصلی به دست آورده، بر مبنای آن‌ها پارامترهای مدل‌های لجستیک مجدداً محاسبه شدند. کارآیی مدل‌ها، با استفاده از خطای استاندارد پارامترها مقایسه گردید.

نتایج: مدل رگرسیون لجستیک براساس متغیرهای اولیه حاکی از مقادیر غیرعادی نسبت شانس برای سن در اولین زایمان زنده (۴۵۳۵۰۳٪ = ۱۰۱۸۴ و OR = ۶۷۹۶۰ CI ۹۵٪ = ۰/۰۰۰۰۲۹) بود. در حالی که پارامترهای مدل‌های لجستیک حاصل از متغیرهای پنهان به دست آمده از هر دو روش تحلیل عاملی و تحلیل مؤلفه‌های اصلی، از نظر آماری معنی دار ($p < 0/003$) و خطای استاندارد همه‌ی آن‌ها کوچکتر از خطای استاندارد مربوط به رگرسیون لجستیک معمولی بود. فاکتورها و مؤلفه‌های اصلی تولید شده توسط دو روش حداقل ۸۵ درصد کل واریانس را تبیین کردند.

نتیجه‌گیری: تحقیق نشان داد انحراف استاندارد پارامترهای برآورده شده در رگرسیون لجستیک براساس متغیرهای پنهان از رگرسیون لجستیک براساس مشاهدات اولیه کوچکتر بوده و در نتیجه این‌گونه مدل‌بندی در تحلیل برخی عوامل خطر سرطان پستان که هم خطی دارند کارآتر است.

واژگان کلیدی: هم خطی چندگانه، متغیر پنهان، تحلیل عاملی، تحلیل مؤلفه‌های اصلی، رگرسیون لجستیک، سرطان پستان.

مقدمه

تعداد 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 لجستیک 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

ماتریس ضرایب همبستگی، عامل‌های مستقل از یکدیگر، از روی متغیرهای اولیه به دست می‌آیند (۱۲).

در این تحقیق، برای کاهش ابعاد مدل رگرسیون لجستیک 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 نحوه کاربرد دو روش تحلیل مؤلفه‌های اصلی و تحلیل عاملی در کاهش ابعاد مدل و ایجاد متغیرهای پنهان، از داده‌های مطالعه مورد شاهدی مربوط به عوامل خطر سرطان پستان استفاده شد (۱۳). در مطالعه مذکور که در فاصله زمانی بهمن ۸۲ تا آذر ۸۳ در مرکز پژوهشی - 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 بخش‌های دیگر بیمارستان شهداء، مثل جراحی، پوست، داخلی و غیره بستری و یا برای پی‌گیری یا درمان به درمانگاه بیمارستان مراجعه کرده بودند و از نظر سنی با گروه مورد با حداقل ۲ سال اختلاف مشابه‌سازی شدند. با اطمینان ۹۵ درصد و توان 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 بالایی 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)، سن در اولین حاملگی (Feeding AFP: Age at First)، سن در اولین زایمان زنده (Pregnancy AFLB: Age at First) و سن در اولین زایمان زنده (Live Birth AFLB: Age at First). برای بررسی میزان بروز هم خطی در این مشاهدات از ماتریس ضرایب همبستگی استفاده است (۲).

ابتدا بدون در نظر گرفتن وجود هم خطی، مدل رگرسیون لجستیک به داده‌ها برازش داده شد. سپس با ترکیب پنج متغیر مورد بررسی، یک بار با روش تحلیل مؤلفه‌های اصلی و بار دیگر به طریق تحلیل عاملی، دو متغیر پنهان به دست آمد و بر اساس آن‌ها، پارامترهای مدل رگرسیون لجستیک برآورد شد. مدل‌های حاصل،

داشته باشند 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 دامنه مورد انتظار شود (۳).

مسئله 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 می‌رond (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 این که هم خطی چندگانه در مدل رگرسیون لجستیک نیز ایجاد مشکل می‌کند (۶،۷)، تاکنون توجه محققان بیشتر بر رگرسیون خطی با متغیر پاسخ دارای توزیع نرمال متتمرکز 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 پیرسن (Pearson) در برآش حداقل مربعات متعامد برمی‌گردد؛ ولی بسط عمدی تئوری بهوسیله‌ی هتلینگ (Hotteling) (۸) انجام شده است. 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 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 نخستین 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

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

$$\text{Factor 1} = +0.85\text{NP} + 0.97\text{NLB} + 0.68\text{TLBF} - 0.26\text{AFP} - 0.26\text{AFLB}$$

$$\text{Factor 2} = +0.21\text{NP} - 0.23\text{NLB} - 0.25\text{TLBF} + 0.94\text{AFP} + 0.95\text{AFLB}$$

در مجموع $84/79\%$ واریانس توسط عامل اول ($45/67\%$) و عامل دوم ($39/12\%$) تبیین شده است. وارد کردن عوامل فوق به عنوان عامل اول ($P<0.002$) و عامل دوم ($P<0.001$) هر دو تاثیر معنی داری بر متغیر وابسته دارند (جدول ۳). در مرحله ای سوم، با روش تحلیل مؤلفه های اصلی دو مؤلفه به صورت ترکیب خطی زیر از متغیرهای اولیه بدست آمد:

$$\text{Component 1} = +0.90\text{NP} + 0.92\text{NLB} + 0.82\text{TLBF} - 0.25\text{AFLB}$$

$$- 0.24\text{AFLB}$$

$$\text{Component 2} = -0.21\text{NP} - 0.25\text{NLB} - 0.25\text{TLBF} + 0.96\text{AFP} + 0.96\text{AFLB}$$

مؤلفه ای اول $49/26\%$ و مؤلفه ای دوم 40% و در مجموع $89/26\%$ کل واریانس را تبیین می کنند. در برآورد پارامترهای رگرسیون لجستیک براساس این مؤلفه های اصلی، نتایج مشابهی حاصل شد. به عبارت دیگر مؤلفه اول با $(P<0.002)$ و مؤلفه دوم با $(P<0.003)$ معنی دار شدند و نسبت شانس ها نیز تقریباً مشابه

نحوی اند.

براساس واریانس های تبیین شده به وسیله‌ی دو روش و خطای استاندارد پارامترهای برآورده شده، مورد مقایسه قرار گرفتند.

یافته ها

جدول یک نشان می دهد بین سه متغیر تعداد حاملگی، تعداد فرزندان زنده به دنیا آورده و کل طول مدت شیردهی به فرزندان و نیز بین دو متغیر سن در اولین حاملگی و سن در اولین زایمان زنده همبستگی خطی بالای مشاهده می شود. همچنین مقادیر عامل تورم واریانس (VIF: Variance Inflation Factor) نیز که میزان بروز هم خطی را نشان می دهد، محاسبه شد که همه‌ی مقادیر بیشتر از یک و نشان دهنده وجود هم خطی در متغیرهاست. به ویژه این که دو متغیر سن در اولین حاملگی و سن در اولین زایمان زنده دارای VIF بیشتر از ۱۵ بوده و بنابراین براساس این معیار هم خطی شدیدی دارند. آزمون بارتلت نیز نشان داد ماتریس ضرایب همبستگی متغیرهای توضیحی با صفر اختلاف معنی داری دارد ($P<0.001$). بنابراین بین متغیرهای توضیحی مورد بررسی هم خطی چندگانه وجود دارد.

جدول ۱- ماتریس ضرایب همبستگی متغیرهای توضیحی

NP	NLB	TLBF	AFP	AFLB	
1	+0.90	-0.67	-0.44	-0.42	NP
1	+0.75	+0.47	-0.47	-0.47	NLB
		1	-0.42	-0.41	TLBF
			1	0.97	AFP
				1	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

جدول ۲- برآورده شده پارامترهای مدل رگرسیون لجستیک بدون درنظر گرفتن همخطی بین متغیرهای توضیحی

متغیرهای اصلی	ضرایب رگرسیونی	استاندارد	خطای	P-value	(فاصله اطمینان ۹۵%)	نسبت شانس	عامل تورم واریانس * (VIF)
عرض از مبداء	0.10	0.20	0.62	(-0.184-45350.3)			
تعداد حاملگی (NP)	0.06	0.09	0.50	(1.06/0.89-1.27)			4/57
تعداد فرزندان زنده به دنیا آورده (NLB)	-0.17	0.26	0.51	(0.84/0.50-1.40)			5/83
کل طول مدت شیردهی به فرزندان (TLBF)	-0.22	0.15	0.14	(0.80/0.59-1.08)			2/13
سن در اولین حاملگی (AFP)	-1.046	0.92	<0.001	(0.000029)			15/41
سن در اولین زایمان زنده (AFLB)	11.13	0.97	<0.001	67960			15/34

Variance Inflation Factor : VIF*

توصیفی میزان مرگ و میر بیماری‌های تنفسی بود، نتایج حاصل از دو روش تحلیل عامل و معادله‌ی مدل‌سازی ساختاری (Structural Equation Modeling) را با رگرسیون کلاسیک براساس متغیرهای اولیه هم خط مقایسه کرده نشان دادند که متغیرهای پنهان، پارامترهایی با خطاها استاندارد کوچکتر تولید می‌کنند (۱۷). نتایج تحقیق حاضر از این لحاظ با مطالعه‌ی آنان هم‌خوانی 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 برآورده شده می‌باشد؛ ولی تفسیر مدل‌هایی که براساس عامل‌ها یا تحلیل مؤلفه‌های اصلی به دست می‌آیند قدری پیچیده است (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 (۱۹، ۲۰، ۱۱). در نتیجه از طریق معکوس ماتریس دوران می‌توان برآوردهای تصحیح شده پارامترهای متغیرهای 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 تکیک‌های جدید علاوه بر تولید متغیرهای پنهان برای متغیرهای توضیحی، اکنون بر تولید این متغیرها برای متغیرهای پاسخ توجه دارند (۲۲).

نتیجه‌گیری

براساس یافته‌های این تحقیق می‌توان نتیجه‌گیری کرد که در بررسی برخی عوامل خطر سرطان پستان، دو روش تحلیل عاملی و تحلیل مؤلفه‌های اصلی نتایج مشابهی داشته، نسبت به مدل لجستیک با متغیرهای هم خط اولیه از کارآیی بالاتری برخوردار هستند.

تشکر و قدردانی

در این مقاله از داده‌های طرح تحقیقاتی عوامل خطر سرطان

جدول ۳- برآورد پارامترهای رگرسیون لجستیک بر اساس متغیرهای پنهان ایجاد شده به‌وسیله روش تحلیل عاملی و مؤلفه‌های اصلی

متغیرهای پنهان	نسبت شانس (فاصله اطمینان ۹۵%)	p-value	ضرایب خطای استاندارد	روش تحلیل عاملی
عرض از مبداء	۰/۹۶	۰/۶۴	۰/۰۸	روش تحلیل عاملی
عامل ۱	۰/۷۶ (۰/۶۴-۰/۹۱)	۰/۰۰۲	۰/۰۹	-۰/۲۷
عامل ۲	۱/۳۳ (۱/۱۲-۱/۵۹)	۰/۰۰۱	۰/۰۹	۰/۲۹
روش تحلیل مؤلفه‌های اصلی				
عرض از مبداء	۰/۹۶	۰/۶۴	۰/۰۹	روش تحلیل مؤلفه‌های اصلی
مؤلفه ۱	۰/۷۷ (۰/۶۵-۰/۹۱)	۰/۰۰۲	۰/۰۹	-۰/۲۶
مؤلفه ۲	۱/۳۱ (۱/۱۰-۱/۵۶)	۰/۰۰۳	۰/۰۹	۰/۲۷

روش تحلیل عاملی به دست آمدند (جدول ۳).

بحث

هدف پژوهش حاضر استفاده از متغیرهای پنهان در مدل رگرسیون لجستیک به منظور کاهش اثر در حالت بروز هم‌خطی چندگانه بود. نتایج حاصل نشان داد برآورد پارامترهای مدل در هر دو روش تحلیل مؤلفه‌های اصلی و تحلیل عامل، مشابه هستند و خطای استاندارد آن‌ها نسبت به مدل اصلی بسیار کوچک‌تر است. بنابراین دو روش به کار رفته برای تولید متغیرهای پنهان نسبت به مدل رگرسیون لجستیک براساس مشاهدات اولیه، در تحلیل برخی عوامل مرتبط با سرطان پستان از کارایی بالاتری برخوردار هستند. در مدل اولیه به دلیل وجود هم‌خطی و بالا بودن خطای استاندارد برآوردها، تعدادی از متغیرها معنی‌دار نشده‌اند و از بین دو متغیر معنی‌دار، متغیر سن در اولین زایمان زنده، نسبت شانس بسیار بزرگ و غیر معمول و متغیر سن در اولین حاملگی، نسبت شانسی نزدیک به صفر را نشان داد.

یافته‌های این تحقیق با نظر آگویلرا و اسکابایاس (Aguilera & Escabias) که در مقاله‌ی خود نشان داده‌اند استفاده از تحلیل مؤلفه‌های اصلی در رگرسیون لجستیک با داده‌های هم‌خط می‌تواند برآورد پارامترها را بهبود بخشد، مطابق است (۱۵). این روش در مقاله‌ی کاربردی اسکابایاس، آگویلرا و والدراما (Valderrama) که در مورد مدل‌سازی داده‌های هواشناسی بود نیز مورد استفاده قرار گرفته است. آنان در این مطالعه بر نحوه‌ی انتخاب مؤلفه‌های اصلی برای بهبود برآورد پارامترها تأکید کرده‌اند (۱۶).

وال و لی (Li & Wall) در مطالعه‌ای که براساس متغیرهای

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.
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. 2002, <http://ace.acadiau.ca/math/chipmahn/publications.html>.
22. Guo, J., Wall, M. M. and Amemiya Y. Latent class regression on latent factors to appear in *Biostatistics*.

پستان، مصوب دانشگاه علوم پزشکی شهید بهشتی استفاده شده است که به این وسیله از کلیه همکاران طرح مذکور و نیز از معاونت 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.