

## تحلیل درستی‌نمایی ماکزیم مدل رگرسیون لجستیک در حالتی که داده‌های متغیرهای پیشگو کامل نیستند ولی متغیرهای کمکی وجود دارند

محمد امین پورحسینی: دانشجوی کارشناسی ارشد آمار زیستی، گروه آمار زیستی دانشکده پیراپزشکی، دانشگاه علوم پزشکی شهیدبهشتی : نویسنده رابط amin\_phg@yahoo.com

دکتر حمید علوی مجد: استادیار، گروه آمار زیستی، دانشکده پیراپزشکی، دانشگاه علوم پزشکی شهیدبهشتی  
دکتر علیرضا ابدی: استادیار گروه پزشکی اجتماعی و بهداشت، دانشکده پزشکی، دانشگاه علوم پزشکی شهیدبهشتی  
سیمین پروانه وار: دانش آموخته کارشناسی ارشد مامایی، دانشگاه علوم پزشکی شهیدبهشتی

دریافت: ۸۴/۴/۵ پذیرش: ۸۴/۵/۲۱

### چکیده:

**مقدمه و هدف:** داده‌های گمشده در بسیاری از مطالعات آماری از جمله مدل‌های رگرسیونی وجود دارند و باعث کاهش دقت برآورد می‌شوند. تا کنون روش‌های گوناگونی برای مقابله با مشکل داده‌های گمشده ابداع شده که عموماً بر داده‌های گمشده متغیر پاسخ متمرکز بوده است حال آنکه متغیرهای پیشگو نیز می‌توانند دستخوش تغییر و از دست رفتن اطلاعات شوند. مواد و روشها: در این تحقیق ضمن بررسی روش‌های گمشده با استفاده از الگوریتم EM و متغیر کمکی، نتایج حاصل از این روش را با روش تحلیل مورد کامل در یک مدل رگرسیون لجستیک پیرامون عوامل مؤثر بر انتخاب نوع زایمان مقایسه می‌کنیم.

**یافته‌ها:** داده‌های مورد استفاده در این مقاله از یک مطالعه توصیفی پیرامون عوامل مرتبط با انتخاب نوع زایمان در زنان مراجعه کننده به مراکز بهداشتی و درمانی شهر تهران بدست آمده است. حجم نمونه در این تحقیق ۳۸۵ نفر بوده و از روش نمونه‌گیری چند مرحله‌ای انتخاب شدند و مشخصات فردی، سوابق مامایی، نوع نگرش و عوامل اجتماعی نمونه‌ها از طریق پرسشنامه ثبت شدند. برای مقایسه میزان کارایی دو روش، برآورد انحراف معیار پارامترها مورد استناد قرار گرفت.

**بحث و نتیجه‌گیری:** نتایج حاصل نشان می‌دهد روش تحلیل درستی‌نمایی با الگوریتم EM در مقایسه با روش مورد کامل کارایی بهتری دارد. مشکل داده‌های گمشده در بسیاری از مطالعات آماری وجود دارد و موجب ارزیابی و کاهش کارایی می‌شوند. در این بررسی نشان داده‌ایم استفاده از الگوریتم EM برای جانهی گمشده‌ها در یک مدل رگرسیون لجستیک با متغیرهای توضیحی گسسته و سپس تحلیل مدل، از روش مورد کامل که مستلزم حذف گمشده‌ها به همراه قسمتهایی از اطلاعات است کاراتر است. از سوی دیگر اگر متغیر توضیحی ناکامل پیوسته باشد بدست آوردن مدل، روشی متفاوت می‌طلبد و یا می‌توان با تبدیل آن به متغیری گسسته از روش قبل استفاده کرد.

**کلید واژه‌ها:** مدل رگرسیون لجستیک، الگوریتم EM، تحلیل مورد کامل، داده گمشده، متغیر کمکی، سزارین

### مقدمه

دچار مشکل می‌کنند و بخصوص در برآورد به روش درست‌نمایی، ایجاد ارزیابی کرده و کارایی را کاهش می‌دهند (۱، ۲، ۳). روش‌های گوناگونی برای مقابله با مشکل داده‌های گمشده ابداع شده که بر اساس آنها گمشده‌ها جانهی می‌شوند. یکی از این روشها استفاده از الگوریتم EM است (۴). الگوریتم EM یک روش تکرار شونده است که در هر تکرار دو گام را شامل می‌شود. گام E (گام امید ریاضی) که داده گمشده به

در بسیاری از تحقیقات پزشکی با متغیرهایی مواجه می‌شویم که قسمتی از اطلاعاتشان به دلایل مختلفی چون عدم پاسخ، ناکامل بودن چارچوب بررسی، از دست رفتن اطلاعات موجود در پرونده‌ها و ... از دست رفته‌اند. به این گونه داده‌ها در آمار داده گمشده گویند و در صورت نادیده گرفتن آنها استنباط آماری را

نمونه به صورتی قابل ملاحظه کوچک شده و منجر به کاهش دقت می شود. همچنین اگر واحدهایی که از تحلیل حذف می شوند با آنهایی که باقی می ماند تفاوت زیادی داشته باشند، ممکن است برآوردهای حاصل به شدت اریب شوند (۸،۲).

۲- تحلیل درستنمایی با استفاده از الگوریتم EM و متغیر کمکی

۱-۲- الگوریتم EM: یکی از روش هایی که برای برآورد داده های گمشده ابداع شده الگوریتم EM است.

الگوریتم EM یک روش محاسباتی عمومی برای برآوردهای حداکثر درستنمایی تحت داده های ناکامل است. تاریخچه EM به سال ۱۹۲۶ و مقاله ایی از مک کندیگ باز می گردد ولی اولین بار بوسیله دمپستر و همکاران ابداع شد (۴). نام این الگوریتم برگرفته از دو گام این الگوریتم است: گام (E) که محاسبه مقادیر مورد انتظار برای گمشده هاست و گام (M) که محاسبه برآوردهای حداکثر درستنمایی پارامترها با فرض کامل بودن داده هاست که براساس دو مرحله پایه گذاری شد.

الف) اگر ما مقادیر گمشده را بدانیم می توانیم پارامترها را برآورد کنیم.

ب) اگر پارامترها را بدانیم می توانیم مقادیر گمشده را با مقادیر مورد انتظار جایگذاری کنیم.

۲-۲- مدل رگرسیون لجستیک با داده های کمکی

برای مشاهدات دو حالتی  $y_1, y_2, \dots, y_n$  و ماتریس  $x$  که مجموعه ایی از متغیرهای گسسته مستقل است مدل لجستیک زیر را داریم.

$$\logit(E[Y_i|x_i]) = x_i^T \beta \quad (1)$$

که بعد  $\beta$ ،  $P \times 1$  است. حداکثر درستنمایی براساس احتمال شرطی  $f(y|x, \beta)$  بدست می آید ولی اگر  $x$  دچار داده های گمشده باشد روش ML مؤثر نیست. ابراهیم روش درستنمایی ماکزیمم را پیشنهاد می کند که با استفاده از الگوریتم EM انجام شده و بجای مدل سازی احتمال  $Y$  به شرط  $X$  از احتمال توأم آنها نیز استفاده می کند (۵).

شرط داده های مشاهده شده محاسبه می شوند این امیدهای ریاضی را به جای داده های گمشده قرار می دهند و پارامترهای مورد نظر برآورد می شوند. در گام بعدی یعنی گام M (گام ماکزیمم کردن) بعد از جایگذاری اعداد اولیه بجای داده های گمشده به شرط داده های مشاهده شده لگاریتم تابع درستنمایی را حداکثر می کنیم. این مکانیسم آنقدر تکرار می شود تا به همگرایی میان پارامترهای برآورد شده در تکرارها برسیم.

با این حال بیشترین روشهای جانهی متوجه داده های گمشده در متغیر پاسخ بوده است حال آنکه متغیرهای پیشگو نیز می توانند دستخوش تغییر و گمشدگی شوند. از اینرو داده های گمشده در متغیرهای پیشگو نیز مورد توجه برخی آمار دانان قرار گرفته است.

ابراهیم روش درستنمایی ماکزیمم را برای مدل های عمومی رگرسیونی پیشنهاد می کند که با بهره گیری از الگوریتم EM گمشده های متغیرهای مستقل گسسته برآورد می شوند (۵).

در اکثر مطالعات آماری معمولاً متغیرهایی بیش از آنچه برای طراحی مدل مورد نیاز است در مورد آزمودنی ها جمع آوری می شوند. به این متغیرها که مقادیر آنها جمع آوری و ثبت شده ولی به عنوان یک متغیر پیشگو در مدل استفاده نشده اند متغیرهای کمکی گویند. در حالی که متغیر پیشگو کامل نیست این متغیرها ممکن است کاملاً مشاهده شده باشند. استفاده از متغیرهای کمکی که داده هایشان کامل است در حالتی که متغیرهای پیشگو کامل نیستند می تواند موجب افزایش کارایی آنالیز مدل شود.

در این تحقیق، برآورد درستنمایی را در مدل رگرسیون لجستیک بررسی می کنیم که متغیرهای مستقل گسسته آن ناکامل اند ولی متغیر کمکی گسسته آن کامل است (۶) و با استفاده از داده های تحقیقی که درباره علل تمایل زنان باردار به نوع زایمان است (۷) کارایی این روش را براساس انحراف معیار برآوردها با روش تحلیل مورد کامل مقایسه می کنیم.

۱- تحلیل مورد کامل: در این روش همه واحدهایی که دارای مقادیر گمشده اند کنار گذاشته می شوند. شاید این متداول ترین روش برای حل مشکل داده های گمشده باشد ولی در کل شیوه خوبی تلقی نمی شود زیرا با حذف همه واحدهای دارای داده های گمشده، اندازه

$$f(Y, X | \Omega) = f(Y|X, \beta) f(X|\gamma) \quad (2)$$

که  $\Omega = (\beta, \gamma)$

اگر  $X$  به طور کامل مشاهده شده باشد  $f(X; Y)$  در دستنمایی برای  $\beta$  شرکت نمی کند و اگر داده گمشده در  $X$  باشد برآورد  $\beta$  ممکن است دچار اربسی و کاهش کارایی شود. از سوی دیگر معمولاً محققین مجموعه ایی از متغیرهای مستقل را در شروع تحقیق بررسی و اطلاعات آنها را ثبت می کنند ولی فقط

$$f(Y, X, A | \Omega^*) = f(Y|X, A, \beta^*) f(X, A | \gamma^*) \quad (3)$$

که  $\Omega^* = (\beta^*, \gamma^*)$

زمانی بکار می رود که استقلال شرطی برقرار باشد (۹). اما استقلال شرطی همیشه برقرار نیست. اگر استقلال شرطی برقرار نباشد فاکتور بندی رابطه (۳) طبیعی نیست زیرا ضرایب رگرسیونی برای  $X$  در  $E(Y|X, A)$  عموماً برابر  $\beta$  در مدل (۱) نیستند. روش پیشنهاد شده دیگری که به وسیله واچ برای داده های کامل بکار رفته بصورت زیر است (۱۰):

$$f(Y, X, A | \theta) = f(A|Y, X, \alpha) f(Y|X, \beta) f(X|\gamma) \quad (4)$$

مجموعه این از پارامترهای دستنمایی لگاریتمی تابع  $\theta = (\alpha, \beta, \gamma)$

$$\sum_i l_{a,y,x}(\theta | a_i, y_i, x_i) = \sum_i \{ l_{a|y,x}(\alpha | a_i, y_i, x_i) + l_{y|x}(\beta | y_i, x_i) + l_x(\gamma | x_i) \} \quad (5)$$

برآورد می شود. اما اگر  $X_i$  های گمشده باشند هر سه قسمت دستنمایی باید برآورد شوند. با فرض اینکه گمشده ها به صورت تصادفی حادث شده اند دستنمایی لگاریتمی زیر را داریم (۶):

$$\sum_i \log \sum_{x_{miss}} \{ L_{a|y,x}(\alpha | a_i, y_i, x_i) \times L_{y|x}(\beta | y_i, x_i) L_x(\gamma | x_i) \} L_{a,y,x}^0(\theta)$$

بوسیله این تابع برآوردهای برای  $\beta$  بدست می آوریم که تحت متغیر کمکی و بدون فرض استقلال شرطی ماکزیم شده است.

بر اساس روش ابراهیم (۵) متغیرهای پیشگو  $X = (x_1, x_2, \dots, x_p)$  متغیرهای گسسته تصادفی با توزیع چندگانه تحت پارامتر

زیر مجموعه ایی از این متغیرها را وارد مدل می کنند و عموماً تعدادی از متغیرها علی رغم گردآوری اطلاعاتشان، در مدل استفاده نمی شوند. به این گونه متغیرها، متغیرهای کمکی گویند (۶).

اگر متغیر کمکی  $A$  موجود باشد می توانیم تابع دستنمایی داده های کامل را به صورت زیر بنویسیم:

\* را برای نشان دادن وابستگی به توزیع  $A$  بکار می بریم.

اگر فرض کنیم  $f(Y|X, A, \beta^*) = f(Y|X, \beta)$  (یعنی استقلال شرطی  $Y$  و  $A$  تحت  $X$  روش ابراهیم را برای مدل می توان بکار برد.

واچ این برآوردگر را وقتی متغیر کمکی برای بهبود بخشیدن به مدل استفاده می شود پیشنهاد کرده است و

۲-۳- برآورد دستنمایی ماکزیم تحت متغیرهای

مستقل ناکامل و اطلاعات کمکی اگر داده ها کامل باشند  $\prod f(Y_i|X_i, \beta)$  برای برآورد  $\beta$  براحتی ماکزیم می شود. در این حالت  $\beta$  جدای از  $\alpha$  و  $\gamma$  با رگرسیون لجستیک

که  $X_{miss,i}$  نشان دهنده قسمتهایی از  $X_i$  است که گمشده اند و  $\sum_{x_{miss,i}}$  مجموع مکان نمونه ایی  $X_{miss,i}$  است.

$X$  کاملاً مشاهده شده است برای برآورد  $\beta$  بکار رود.

اما در حالتی که  $X$  ناکامل است از الگوریتم EM (۴) برای برآورد پارامترها در درستنمایی داده های مشاهده شده استفاده می کنیم. این الگوریتم تکرار شونده شامل دو مرحله است. در مرحله E امید داده های مشاهده شده محاسبه می شوند.

در این وضعیت، داده های مشاهده شده شامل  $(y_i, x_{obs,i}, a_i)$  اند که  $X_{obs,i}$  نشان دهنده مجموعه ایی با بعد کوچکتر یا مساوی  $P$  متغیر مستقل در  $X_i$  است و اگر  $i$  امین عضو مشاهده شده در متغیرهای مستقل باشد معادل  $X_i$  خواهد بود.

این امید با  $Q(\theta|\theta^{(t)})$  نشان داده می شود و داریم (۶):

$$Q(\theta|\theta^{(t)}) = \sum_{i=1}^n \sum_{j=1}^{r+1} w_{ij}^{(t)} L_{a,y,x}(\theta|a_i, y_i, x^j) \\ = \sum_{i=1}^n \sum_{j=1}^{r+1} w_{ij}^{(t)} \{l_{a|y,x}(\alpha|a_i, y_i, x^j) + l_{y|x}(\beta|y_i, x^j) + L_x(\gamma|x^j)\} \quad (6)$$

و  $X_{obs,i}$  زمانی متناسب اند که اعضای متناظر با ستون  $x^j$  برابر داده های مشاهده شده  $X_{obs,i}$  باشد.

اگر همه متغیرهای مستقل برای امین عضو، مشاهده شده باشند فقط یک قسمت مخالف صفر وجود خواهد داشت و وزنهای برای مقدار مشاهده شده  $x^j = x_{obs,i}$  برابر یک خواهد بود.

$\gamma = (\gamma_1, \gamma_2, \dots, \gamma_r)$  مثلاً اگر سه متغیر دو حالتی وجود داشته باشد بعد  $\gamma$  برابر است با  $r = 2^3 - 1 = 7$ . در حالت کلی فرض کنید  $c_1, \dots, c_p$  نشان دهنده تعداد رده ها برای هر متغیر باشد. در نتیجه  $\gamma$  دارای بعد  $r = c_1 \times \dots \times c_p - 1$  است.

فرض کنید  $f(Y|X, \beta)$  چگالی متغیر پاسخ و  $f(X|Y)$  و  $f(A|X, Y, \alpha)$  توزیع های چند جمله ایی با پارامترهای  $\gamma, \alpha$  (مجزا از  $\beta$ ) باشند. اگر هیچ محدودیتی به توزیع  $X$  تحمیل نشود تحت داده های کامل،  $\gamma$  بوسیله  $r+1$  خانه شمارشی مشاهده شده از رده های  $x$  در یک جدول توافقی  $p$  بعدی بدست می آید. به طور مشابه مدل چند جمله ایی برای برآورد  $f(A|X, Y, \alpha)$  می تواند مورد استفاده قرار گیرد و رگرسیون لجستیک را وقتی

که  $x^j$ ، امین نمونه ممکن از متغیرهای مستقل،  $L_{a,y,x}(\theta|a_i, y_i, x_i)$  درستنمایی لگاریتمی داده های کامل برای  $\theta$  با  $i$  امین مشاهده تحت  $X_i$  ارزیابی شده با  $x^j$  است و  $w_{ij}^{(t)} = p(x^j|a_i, y_i, x_{obs,i}, \theta^{(t)})$  میتواند وزنی از  $j$  امین نمونه برای  $i$  امین مشاهده در  $t$  امین تکرار باشد. در نتیجه  $W_{it}^{(t)} = 1$ . البته بیشتر این  $W_{ij}$  ها صفر خواهند شد.

پس هر مقدار  $X_{obs,i}$  تحت  $x^j$  مناسب تشخیص داده نشود مورد استفاده قرار نمی گیرد. (ستون  $x^j$

وزنها را می توان از قانون بیز به صورت زیر محاسبه کرد (۶):

$$W_{ij}^{(t)} = P(x^j|a_i, y_i, x_i, \theta^{(t)})$$

$$= \begin{cases} 0 & \text{if } x^j \text{ is not compatible with } x_i \\ \frac{p(y_i|x_i^j)P(a_i|x_i^j, y_i)p(x_i^j)}{\sum_{k \in \text{obs}_i} p(y_i|x_i^k)p(a_i|x_i^k, y_i)p(x_i^k)} & \text{if } x^j \text{ is compatible with } x_i \end{cases} \quad (7)$$

مرحله ای بوده و مشخصات فردی، سوابق مامایی، نوع نگرش و عوامل اجتماعی نمونه‌ها از طریق پرسشنامه ثبت شدند. متغیر پاسخ در این تحقیق نوع زایمان انتخاب شده بوسیله زنان باردار است که به دو سطح زایمان طبیعی و سزارین تقسیم شده است. متغیرهای پیشگویی که برای ساختن مدل انتخاب شده اند شش متغیر دوحالتی هستند شامل شغل آزمودنی، شرکت در کلاسهای آموزشی دوران بارداری، نوع حاملگی، نگرش زنان باردار به زایمان طبیعی، داشتن سابقه تولد نوزاد باوزن بیشتر از ۴۰۰۰ گرم و داشتن سابقه تولد نوزاد باوزن کمتر از ۲۵۰۰ گرم (جدول یک).

در متغیر داشتن سابقه تولد نوزاد باوزن کمتر از ۲۵۰۰، ۲۵ درصد از داده‌ها به عنوان داده گمشده به صورت تصادفی حذف شدند. همچنین داشتن سابقه زایمان پس از موعد نیز به همراه سایر متغیرها در پرسشنامه ثبت شده و اطلاعات آن کامل بود که در مدل وارد نشده و به عنوان یک متغیر کمکی مورد استفاده قرار گرفت (این متغیر نیز دو حالتی است).

ابتدا با روش تحلیل مورد کامل، قسمتهای گمشده از مطالعه حذف و پارامترها برآورد شدند. درمرحله بعد، از برآورد درستی‌نمایی ماکزیمم تحت متغیر کمکی برای برآورد پارامترهای مدل استفاده کردیم. یک برنامه کامپیوتری تحت نرم افزار s-plus برای این منظور مورد استفاده قرار گرفته و برای مقایسه میزان کارایی دو روش برآورد انحراف معیار پارامترها مورد استناد قرار گرفت که انحراف معیارها در روش تحلیل درستی‌نمایی با الگوریتم EM از روش لوئیس (۱۱) برآورد شدند.

### یافته‌ها

نتایج حاصل از برآورد انحراف معیارهای دو روش تحلیل مورد کامل و جانهای با الگوریتم EM نشان می‌دهد که انحراف معیارهای برآورد شده برای

که  $x^j$  نشان دهنده یک بردار  $p$  تایی است که اجزاء مشاهده شده اش  $X_{obs,i}$  و اجزاء باقی مانده آن بر اساس  $j$  امین نمونه برای متغیرهای گمشده است.

برای مرحله M معادله (۶) را به عنوان یک تابع از  $\theta$  و بوسیله یافتن راه حل معادله درستی‌نمایی لگاریتمی داده‌های کامل با استفاده از این وزن‌ها ماکزیمم می‌کنیم. اکنون می‌توانیم از رگرسیون لجستیک وزن داده شده برای برآورد  $\beta$  استفاده کنیم و برای برآورد  $\alpha, \gamma$  نیز از خانه‌های شمارش شده مورد انتظار بهره می‌گیریم.

برای اولین تکرار، وزن‌ها را می‌توان از رابطه (۷) و نخستین برآوردها برای  $\theta$  را با فرض کامل بودن داده‌ها بدست آورد. این مراحل آنقدر تکرار می‌شوند که به برآوردهای درستی‌نمایی ماکزیمم همگرا شوند و اختلاف میان پارامترها در تکرارهای بعد معنی دار نباشد.

### روش کار

برای مقایسه کارایی دو روش تحلیل مورد کامل و تحلیل درستی‌نمایی با الگوریتم EM از اطلاعات تحقیقی استفاده کرده ایم که در مورد عوامل مرتبط با انتخاب نوع زایمان در زنان مراجعه کننده به مراکز بهداشتی و درمانی شهر تهران (۷) است. این تحقیق یک مطالعه توصیفی بوده است و در آن نمونه ای به حجم ۳۸۵ نفر از مادران باردار با سن حاملگی ۲۸

هفته و بالاتر از مراکز بهداشتی و درمانی شهر تهران انتخاب شدند. روش نمونه گیری از نوع چند

وجود متغیری کمکی که اطلاعات آن به طور کامل موجود است ولی مستقیماً به عنوان یک متغیر پیشگو در مدل مورد استفاده قرار نگرفته می‌تواند برای افزایش کارایی مدل و برآورد دقیقتر پارامترها مفید باشد. انتخاب متغیر کمکی بستگی به اطلاعات موجود در تحقیق و نظر محقق دارد ولی هرچقدر ارتباط میان متغیر کمکی و متغیری که دچار گمشدگی شده بیشتر باشد این روش مؤثرتر خواهد بود (۶). از طرف دیگر برآورد انحراف معیار پارامترهایی که با EM بدست آمده اند مشکل است. لوئیس روشی برای برآورد انحراف معیار پیشنهاد کرده که در این تحقیق مورد استفاده قرار گرفته است (۱۱). همچنین از روش بوت استراپ نیز می‌توان برای برآورد انحراف معیارها استفاده کرد.

توزیعهای  $f(A|X, Y, \alpha)$ ,  $f(X; \delta)$  از طریق برآورد پارامترهای توزیع چند جمله ایی برآورد می‌شوند هرچند گمشده‌ها می‌توانند تاثیر نامطلوبی بر برآورد داشته باشند.

از سوی دیگر اگر  $Y$  متغیری پیوسته باشد بدست آوردن مدلی برای  $f(A|Y, X)$  روشی متفاوت می‌طلبد یا می‌توان با تبدیل کردن  $Y$  به متغیری گسسته از روش قبل استفاده کرد.

توسعه این روش به سایر مدل‌های رگرسیونی قابل بررسی است و برخی از این مدل‌ها - بدون استفاده از متغیر کمکی - بوسیله ابراهیم و همکاران مورد توجه قرار گرفته است. همچنین توجه به نوع مکانیسم گمشدگی (۱۲) در استفاده از روش جانهای با الگوریتم EM مهم است زیرا این روش مانند اکثر روشهای جانهای (۱۳، ۲) صرفاً در مکانیسمهای گمشدگی تصادفی مناسب عمل می‌کند.

تمام پارامترها در روش جانهای با الگوریتم EM کمتر از انحراف معیارهایی هستند که در روش تحلیل مورد کامل بدست آمده است که نشان می‌دهد کارایی برآورد پارامترها در این روش بهتر از روش تحلیل مورد کامل است و علی‌رغم اینکه هیچ قسمت از اطلاعات از مدل حذف نمی‌شوند پارامترهایی با انحراف معیارهای پایینتر برآورد میشوند و در نتیجه فاصله اطمینانهای کوتاهتری برای پارامترها بدست می‌آید (جدول دو). اگر چه اختلافات قابل توجهی نیز در برآورد پارامترهای لجستیک در دو روش مشاهده می‌شود که احتمالاً ناشی از اریبی حاصل از حذف قسمتی از اطلاعات در روش تحلیل مورد کامل است (۲).

## بحث

مشکل داده‌های گمشده در بسیاری از مطالعات آماری وجود دارد و اگر بدون توجه به وجود گمشده‌ها، مدل رابرازش دهیم و ضرایب رگرسیونی را با روش درستنمایی ماکزیمم برآورد کنیم گمشده‌ها بر برآورد ها اثر منفی گذاشته و موجب اریبی و کاهش کارایی می‌شوند (۱، ۲، ۳). در این تحقیق روشی را بررسی کردیم که با بهره‌گیری از متغیری کمکی آنالیز درستنمایی ماکزیمم را در مدل لجستیک انجام داده و با استفاده از الگوریتم EM برآورد پارامترها را در حالتیکه متغیر پاسخ کامل نیست محاسبه می‌کند. نتایج حاصل از بررسی این روش در مدل مربوط به عوامل موثر بر انتخاب نوع زایمان نشان داد که برآورد انحراف معیار پارامترها در این روش از تحلیل مورد کامل کمتر است که این مساله با یافته‌های مدل کاربردی هورتون و لایرد تطابق دارد (۶).

## جدول ۱: متغیرهای توضیحی مدل

متغیر	رده‌ها
-------	--------

شغل آزمودنی	شاغل=۱ ، خانه دار=۰
شرکت در کلاسهای آموزشی دوران بارداری	شرکت = ۱ ، عدم شرکت = ۰
نوع حاملگی	خواسته=۱ ، ناخواسته=۰
نگرش زنان باردار به زایمان طبیعی	نگرش مثبت = ۰ نگرش منفی = ۱
سابقه تولد نوزاد باوزن بیشتر از ۴۰۰۰ گرم	سابقه دارد=۱ ، سابقه ندارد=۰
سابقه تولد نوزاد باوزن کمتر از ۲۵۰۰ گرم	سابقه دارد=۱ ، سابقه ندارد=۰

جدول ۲- برآورد پارامترها و انحراف معیار آنها از دو روش در دستنمایی الگوریتم EM و تحلیل مورد کامل

پارامترهای مدل	انحراف معیار	برآورد ضریب رگرسیونی	EM و متغیر کمکی	انحراف معیار	برآورد ضریب رگرسیونی	روش تحلیل مورد کامل
عرض از مبدا	۰/۲۷۲۸	-۱/۳۱۷۲	EM و متغیر کمکی	۰/۴۶۲	-۲/۲۴۲	روش تحلیل مورد کامل
شغل آزمودنی	۰/۴۰۰۶	-۱/۲۴۹۶	برآورد ضریب رگرسیونی	۰/۹۳۵	۰/۹۸۵	برآورد ضریب رگرسیونی
شرکت در کلاسهای آموزشی دوران بارداری	۰/۲۶۸۳	-۰/۵۳۳۹	انحراف معیار	۰/۴۵۵	۰/۱۶۶	برآورد ضریب رگرسیونی
نوع حاملگی	۰/۲۲۷	-۰/۳۲۲۱	برآورد ضریب رگرسیونی	۰/۳۲۱	-۰/۰۸۳	برآورد ضریب رگرسیونی
نگرش زنان باردار به زایمان طبیعی	۰/۲۸۶۴	-۰/۵۷۹۴	انحراف معیار	۰/۳۵۴	۰/۴۶۰	برآورد ضریب رگرسیونی
سابقه تولد نوزاد باوزن بیشتر از ۴۰۰۰ گرم	۰/۲۸۴۳	۲/۴۸۳۸	برآورد ضریب رگرسیونی	۰/۳۸۳	۲/۵۳۱	برآورد ضریب رگرسیونی
سابقه تولد نوزاد باوزن کمتر از ۲۵۰۰ گرم	۰/۳۸۶۷	-۰/۳۴۸۸	انحراف معیار	۰/۸۱۳	۰/۲۱۰	برآورد ضریب رگرسیونی

## References

- 7- Horton. N. J. and Laird. N. M. (2001) Maximum Likelihood Analysis of Logistic Regression Models with Incomplete Covariate Data and Auxiliary Information. *Biometrics* 2001, 57, 34-42.
- 8- Glynn, R. J. and Laird N. M. Regression Estimates and Missing Data: Complete Case Analysis. Unpublished Manuscript, Department of Biostatistics, Harvard University 1983.
- 9- Vach, W. Some Issues in Estimating the Effect of Prognostic Factors from Incomplete Covariate Data. *Statistics in Medicine* 1997 16, 57-72.
- 10- Vach, W. Logistic Regression with Missing Values in the Covariates. Berlin: Springer-Verlag 1994.
- 11- Louis, T. A. Finding the Observed Information Matrix When Using the EM Algorithm. *Journal of the Royal Statistical Society*, 1982 Series B 44, 226-233.
- 12- Rubin, D. B. Inference and Missing Data. *Biometrika* 1976 63, 581-592.
- 13- Saleh A. M. Some Methods for Dealing with Missing Data in Sample Surveys. Invited Papers Proceedings of the 7<sup>th</sup> Iranian Statistical Conference, 2004, 313-324
- ۱- پروانه وار سیمین. بررسی عوامل مرتبط با انتخاب نوع زایمان در زنان باردار مراجعه کننده به مراکز بهداشتی و درمانی شهر تهران سال ۱۳۸۲. پایان نامه کارشناسی ارشد رشته مامایی، دانشگاه علوم پزشکی شهید بهشتی ۱۳۸۲
- 2- Little, R. J. Biostatistical Analysis with Missing Data. In *the Encyclopedia of Biostatistics*, Armitage, P. A. and Colton, T. , Eds., Wiley, Chichester U. K. , 1998
- 3- Levy, P. S. and Lemeshow, S. Sampling of Populations Methods and Applications. Third Edition John Wiley and Sons 1999;393-416.
- 4- Little, R. J. and Rubin, D. B. Statistical Analysis with Missing Data. Newyork: John Wiley and Sons.1987 .
- 5- Dempster, A. P., Laird, N. M., and Rubin, D. B. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*,1977; Series B39: 1-22.
- 6- Ibrahim, J. G. Incomplete data in generalized linear models. *Journal of the American Statistical Association*, 1990, 85, 765-769.