

بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِیْمِ



دانشگاه علوم پزشکی و توانبخشی

گروه آمار و کامپیوتر

پایان نامه کارشناسی ارشد

رشته آمار زیستی

عنوان

کاربرد روش بیزی در برآورد پارامترهای مدل رگرسیون لوجستیک با مقادیر گمشده تصادفی
(MAR) در متغیر کمکی

نگارنده

الهه کاظمی

استاد راهنما

دکتر مسعود کریملو

استاد مشاور

دکتر مهدی رهگذر

تابستان ۱۳۹۰

شماره ثبت ۱۰۳-۶۰۰۰



University of Social Welfare and Rehabilitation Sciences

Department of Biostatistics and Computer

M. Sc. Thesis

**Application of Bayesian Method in Parameters Estimation of Logistic
Regression Model with Missing at Random (MAR) Covariate**

By:

Elaheh Kazemi

Supervisor:

Dr. Masoud Karimlou

Dr. Mehdi Rahgozar

Summer 2011

تقدیم به :

پدر و مادر مهربانم

که در سختی ها و مشکلات زندگی همواره در کنارم بوده اند و مرا در تمامی مراحل زندگی حمایت

کرده اند.

تقدیر و تشکر :

بدینوسیله بر خود واجب می دانم که از اساتید محترم گروه نهایت قدردانی و سپاس را داشته باشم.

استاد ارجمندم جناب آقای دکتر مسعود کریملو استاد راهنمای اینجانب که در کلیه مراحل انجام این پایان نامه از راهنمایی های فکری و علمی ایشان بهره مند بوده ام.

استاد ارجمندم جناب آقای دکتر مهدی رهگذر استاد مشاور اینجانب که افتخار شاگردی ایشان موهبت بزرگی برای من بوده است.

استاد ارجمندم جناب آقای دکتر عنایت الله بخشی که از نظرات ارزشمند ایشان در انجام این پایان نامه بهره برده ام.

سرکار خانم دکتر ایمانه عسگری و سرکار خانم مینو دیانت خواه که داده های مورد استفاده در این پژوهش را در اختیار اینجانب قرار داده اند.

دوستان هم دوره ای سرکار خانم زهره رزاقی و سرکار خانم لیلا چراغی که دوران سخت تحصیل در کنار ایشان بسیار خاطره انگیز بود.

چکیده :

رگرسیون لوجستیک مدلی عمومی برای تحلیل داده های پزشکی و اپیدمیولوژیکی می باشد و اخیراً محققین معدودی تحقیقات خود را به تحلیل مدل های رگرسیون لوجستیک با وجود مقادیر گمشده در متغیرهای کمکی معطوف داشته اند. در بسیاری از پژوهش ها محققین با مجموعه داده هایی مواجه هستند که دارای مقادیر گمشده است. گمشدگی تهدید عمده ای برای درستی نتایج حاصل از مجموعه داده ها محسوب می شوند و اجتناب از آن بسیار مشکل است.

تمامی روش های برآورد پارامترها بر پایه فرض کامل بودن مجموعه داده ها استوار است و تحت برقراری این شرایط منجر به برآوردهایی نا اریب می شوند. در مطالعات انجام شده ثابت شده است که در صورت وجود مقادیر گمشده تصادفی در مجموعه داده ها برآوردهای حاصل، دیگر نارایب نخواهند بود و با افزایش نسبت گمشدگی، مقدار اریبی نیز افزایش خواهد یافت.

ساتن و کارول تابع درستنمایی ویژه ای را برای برآورد پارامترهای مدل رگرسیون لوجستیک وقتی که مقادیر کمکی Z به طور کامل مشاهده شده و مقادیر متغیر کمکی X دارای مقادیر گمشده از نوع مکانیسم گمشدگی تصادفی (MAR) باشند، معرفی کرده اند. در این پژوهش از این تابع درستنمایی در تحلیل بیزی برای برآورد پارامترهای مدل رگرسیون لوجستیک استفاده شده است و نتایج به دست آمده با روش های جانهای چندگانه و واحد کامل مقایسه شده است.

روش های مذکور را بر روی داده های شبیه سازی شده و داده های دندانپزشکی اجرا کرده و پس از مقایسه نتایج حاصل از سه روش مذکور نتیجه گرفته شد که اگر مکانیسم گمشدگی تصادفی باشد، به کارگیری تحلیل بیزی با تکنیک MCMC منجر به برآوردهای دقیق تری نسبت به روش جانهای چند گانه و روش واحد کامل می شود.

واژه های کلیدی : رگرسیون لوجستیک، مکانیسم گمشدگی تصادفی (MAR)، روش بیزی، زنجیرهای مارکوف مونت کارلویی (MCMC)، جانهای چندگانه، روش واحد کامل، DMFT.

فهرست مطالب

فصل اول: کلیات تحقیق

- ۱-۱ بیان مسئله..... ۱
- ۱-۱-۱ مقدمه..... ۱
- ۲-۱ اهمیت و ضرورت..... ۴
- ۳-۱ اهداف پژوهش..... ۵

فصل دوم: چارچوب نظری، پنداشتی، مفهومی و تاریخی موضوع پژوهش

- ۱-۲ مقدمه..... ۷
- ۲-۲ روش های مورد استفاده در بررسی داده های گمشده..... ۸
- ۱-۲-۲ روش های مبتنی بر واحد کامل..... ۸
- ۲-۲-۲ روش های مبتنی بر جانمایی..... ۹
- ۳-۲-۲ روش های مبتنی بر مدل..... ۱۰
- ۳-۲ الگوهای گمشدگی داده ها..... ۱۱
- ۴-۲ مکانیسم گمشدن داده ها..... ۱۶
- ۱-۴-۲ انواع مکانیسم گمشدن داده ها..... ۱۷
- ۲-۴-۲ توزیع داده های گمشده..... ۲۳
- ۳-۴-۲ آزمون کردن مکانیسم گمشدگی کاملاً تصادفی..... ۲۸
- ۱-۳-۴-۲ مقایسه های آزمون t تک متغیره..... ۲۸
- ۲-۳-۴-۲ آزمون گمشدگی کاملاً تصادفی لیتل..... ۳۰
- ۴-۴-۲ آیا فرض مکانیسم گمشدگی تصادفی پذیرفتنی است..... ۳۱
- ۱-۴-۴-۲ روشی برای تشخیص مکانیسم گمشدگی تصادفی..... ۳۲

۳۶..... ۵-۲ گمشدگی تصادفی داده ها در جدول ۲×۲.....

۳۷..... ۱-۵-۲ برآورد نقطه ای.....

۴۸..... ۲-۵-۲ مقایسه دو روش.....

۴۹..... ۶-۲ گمشدگی تصادفی داده در چندین جدول ۲×۲.....

۴۹..... ۷-۲ بررسی متون.....

فصل سوم: روش شناسی تحقیق

۵۵..... ۱-۳ روش گردآوری داده ها.....

۵۵..... ۲-۳ متغیرهای تحقیق.....

۵۷..... ۳-۳ روش تجزیه و تحلیل داده ها.....

۵۷..... ۱-۳-۳ تحلیل بیزی.....

۵۹..... ۱-۱-۳-۳ توزیع پیشین پارامترها.....

۶۰..... ۲-۱-۳-۳ توزیع پسین پارامترها.....

۶۱..... ۳-۱-۳-۳ نمونه گیری از پارامترها.....

۶۲..... ۴-۱-۳-۳ نمونه گیری گیبس.....

۶۳..... ۵-۱-۳-۳ روش تشخیص همگرایی.....

۶۵..... ۶-۱-۳-۳ تحلیل بیزی با استفاده از نرم افزار WinBUGS.....

۶۸..... ۷-۱-۳-۳ تعیین مدل.....

۶۸..... ۸-۱-۳-۳ مدل های گرافیکی.....

۷۰..... ۲-۳-۳ مدل های خطی تعمیم یافته.....

۷۲..... ۱-۲-۳-۳ مدل لوجیت.....

۷۳..... ۲-۲-۳-۳ رگرسیون لوجستیک.....

۷۳ رگرسیون لوجستیک ساده.....
۷۴ تابع درستنمایی رگرسیون لوجستیک ساده.....
۷۶ رگرسیون لوجستیک چندگانه.....
۷۶ مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر کمکی.....
۸۱ تابع درستنمایی رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر کمکی.....
۸۵ تحلیل بیزی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر کمکی.....
۸۵ توزیع های پیشین.....
۸۶ توزیع های شرطی کامل پارامترها.....
فصل چهارم : نتایج و یافته ها	
۹۱ مقدمه.....
۹۱ تحلیل داده های شبیه سازی شده.....
۹۱ داده های شبیه سازی شده.....
۹۲ اجرای مدل روی داده های شبیه سازی شده.....
۱۱۳ بررسی داده ها.....
۱۱۷ بررسی مکانیسم گمشدگی داده ها.....
۱۱۷ بررسی مکانیسم گمشدگی داده ها با استفاده از روش های غیر تحلیلی.....
۱۲۰ بررسی مکانیسم گمشدگی داده ها با استفاده از روش های تحلیلی.....
۱۲۰ بررسی مکانیسم گمشدگی کاملاً تصادفی.....
۱۲۵ بررسی مکانیسم گمشدگی تصادفی.....
۱۲۷ تحلیل داده ها.....

فصل پنجم : بحث، نتیجه گیری و پیشنهادات

۱-۵	مقدمه.....	۱۴۶
۲-۵	خلاصه نتایج و مقایسه آن با پژوهش های گذشته.....	۱۴۶
۳-۵	محدودیت پژوهش.....	۱۴۹
۴-۵	پیشنهادات.....	۱۴۹

فهرست منابع

ضمائم و پیوست ها

فهرست جداول

- ۱-۲ جدول مجموعه داده انتخاب کارمند..... ۱۵
- ۲-۲ جدول توزیع داده ها..... ۱۸
- ۳-۲ جدول داده های کامل و داده های گمشده..... ۱۸
- ۴-۲ جدول نمره کارایی شغلی..... ۲۱
- ۵-۲ جدول توزیع داده ها و نشانگر مشاهدات..... ۲۵
- ۶-۲ جدول متغیر به طور کامل مشاهده شده S (سیگار کشیدن) و متغیر نشانگر Δ ۲۷
- ۷-۲ جدول داده ها..... ۳۴
- ۸-۲ جدول داده های تکثیر شده..... ۳۵
- ۹-۲ نتایج به دست آمده از خروجی SAS..... ۳۶
- ۱-۳ جدول متغیرها..... ۵۶
- ۱-۴ جدول داده های شبیه سازی شده با درصد گمشدگی متفاوت به تفکیک طبقات داده های به طور کامل مشاهده شده..... ۹۱
- ۲-۴ (الف) نتایج به دست آمده برای برآورد پارامتر β_0 برای داده های کامل و داده های دارای مقادیر گمشده با استفاده از روش های واحد کامل، جانهی چند گانه و روش بیز درستنمایی ساتن و کارول..... ۹۶
- ۲-۴ (ب) جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_0 ۹۶
- ۳-۴ (الف) نتایج به دست آمده برای برآورد پارامتر β_1 برای داده های کامل و داده های دارای مقادیر گمشده با استفاده از روش های واحد کامل، جانهی چند گانه و روش بیز درستنمایی ساتن و کارول..... ۹۷
- ۳-۴ (ب) جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_1 ۹۸
- ۴-۴ (الف) نتایج به دست آمده برای برآورد پارامتر β_2 برای داده های کامل و داده های دارای مقادیر گمشده با استفاده از روش های واحد کامل، جانهی چند گانه و روش بیز درستنمایی ساتن و کارول..... ۹۹
- ۴-۴ (ب) جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_2 ۹۹

- ۴-۵ (الف) نتایج به دست آمده برای برآورد پارامتر β_{12} برای داده های کامل و داده های دارای مقادیر گمشده با استفاده از روش های واحد کامل، جانپهی چند گانه و روش بیز درستنمایی ساتن و کارول..... ۱۰۰
- ۴-۵ (ب) جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_{12} ۱۰۰
- ۴-۶ جدول داده های شبیه سازی شده با درصد گمشدگی به تفکیک طبقات داده های به طور کامل مشاهده شده..... ۱۰۶
- ۴-۷ (الف) نتایج به دست آمده برای برآورد پارامتر β_0 برای داده های کامل و داده های دارای مقادیر گمشده با استفاده از روش های واحد کامل، جانپهی چند گانه و روش بیز درستنمایی ساتن و کارول..... ۱۰۷
- ۴-۷ (ب) جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_0 ۱۰۷
- ۴-۸ (الف) نتایج به دست آمده برای برآورد پارامتر β_1 برای داده های کامل و داده های دارای مقادیر گمشده با استفاده از روش های واحد کامل، جانپهی چند گانه و روش بیز درستنمایی ساتن و کارول..... ۱۰۸
- ۴-۸ (ب) جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_1 ۱۰۸
- ۴-۹ (الف) نتایج به دست آمده برای برآورد پارامتر β_2 برای داده های کامل و داده های دارای مقادیر گمشده با استفاده از روش های واحد کامل، جانپهی چند گانه و روش بیز درستنمایی ساتن و کارول..... ۱۰۹
- ۴-۹ (ب) جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_2 ۱۰۹
- ۴-۱۰ (الف) نتایج به دست آمده برای برآورد پارامتر β_3 برای داده های کامل و داده های دارای مقادیر گمشده با استفاده از روش های واحد کامل، جانپهی چند گانه و روش بیز درستنمایی ساتن و کارول..... ۱۱۰
- ۴-۱۰ (ب) جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_3 ۱۱۰
- ۴-۱۱ (الف) نتایج به دست آمده برای برآورد پارامتر β_{12} برای داده های کامل و داده های دارای مقادیر گمشده با استفاده از روش های واحد کامل، جانپهی چند گانه و روش بیز درستنمایی ساتن و کارول..... ۱۱۱
- ۴-۱۱ (ب) جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_{12} ۱۱۱
- ۴-۱۲ (الف) نتایج به دست آمده برای برآورد پارامتر β_{13} برای داده های کامل و داده های دارای مقادیر گمشده با استفاده از روش های واحد کامل، جانپهی چند گانه و روش بیز درستنمایی ساتن و کارول..... ۱۱۲
- ۴-۱۲ (ب) جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_{13} ۱۱۲
- ۴-۱۳ درصد گمشدگی به تفکیک طبقات سایر متغیرها در متغیر متراژ..... ۱۱۷

- ۱۴-۴ جدول درصد گمشدگی به تفکیک طبقات سایر متغیرها در متغیر متراژ.....۱۱۸
- ۱۵-۴ جدول درصد گمشدگی به تفکیک طبقات سایر متغیرها در متغیر وضعیت اقتصادی-اجتماعی.....۱۱۹
- ۱۶-۴ جدول درصد گمشدگی به تفکیک طبقات سایر متغیرها در متغیر وضعیت اقتصادی-اجتماعی.....۱۲۰
- ۱۷-۴ جدول میانگین های دو گروه مشاهده شده و گمشده برای متغیر متراژ منزل.....۱۲۱
- ۱۸-۴ جدول نتایج آزمون t برای بررسی مکانیسم گمشدگی تصادفی.....۱۲۱
- ۱۹-۴ جدول میانگین های دو گروه مشاهده شده و گمشده برای متغیر وضعیت اقتصادی-اجتماعی.....۱۲۲
- ۲۰-۴ جدول نتایج آزمون t برای بررسی مکانیسم گمشدگی تصادفی.....۱۲۳
- ۲۱-۴ جدول نتایج به دست آمده از خروجی SAS برای متغیر متراژ منزل.....۱۲۶
- ۲۲-۴ جدول نتایج به دست آمده از خروجی SAS برای متغیر متراژ منزل.....۱۲۷
- ۲۳-۴ جدول نتایج به دست آمده برای برآورد پارامترها برای داده های گمشده برای گروه سنی، متراژ و DMF.....۱۲۹
- ۲۴-۴ جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_0۱۳۰
- ۲۵-۴ جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_1۱۳۰
- ۲۶-۴ جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_2۱۳۰
- ۲۷-۴ جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_{12}۱۳۰
- ۲۸-۴ جدول نتایج به دست آمده برای برآورد پارامترها برای داده های گمشده با استفاده از روش SCMCMC برای گروه سنی، متراژ و DMF.....۱۳۵
- ۲۹-۴ جداول نتایج به دست آمده برای برآورد پارامترها برای داده های گمشده به روش های مختلف برای جنسیت، گروه سنی، متراژ و DMF.....۱۳۶
- ۳۰-۴ جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_0۱۳۷
- ۳۱-۴ جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_1۱۳۷
- ۳۲-۴ جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_2۱۳۸

- ۳۳-۴ جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_3 ۱۳۸
- ۳۴-۴ جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_{12} ۱۳۸
- ۳۵-۴ جدول نسبت بخت ها و فواصل اطمینان برای برآورد پارامتر β_{13} ۱۳۸
- ۳۶-۴ نتایج به دست آمده برای برآورد پارامترها برای داده های گمشده با استفاده از روش SCMCMC
برای جنسیت، گروه سنی، متراژ و DMF ۱۴۵

فهرست نمودارها و تصاویر

- ۱-۲ شکل شش طرح الگوهای گمشدگی داده در مجموعه داده ها با چهار متغیر.....۱۳
- ۲-۲ نمایش مکانیسم های گمشدگی داده۲۶
- ۱-۳ شکل گرافیکی مثال محاسبه خطر نسبی.....۷۰
- ۲-۳ داده های مشاهده شده که در آن متغیر کمکی x ، $n-m$ داده گمشده دارد.....۸۲
- ۳-۳ نمودار گرافیکی مدل رگرسیون لوجستیک با مقادیر گمشده در متغیر کمکی x۸۷
- ۱-۴ نمودار اثر پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر x۱۰۲
- ۲-۴ نمودار تاریخچه پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر x۱۰۳
- ۳-۴ نمودار چگالی پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر x۱۰۴
- ۴-۴ نمودار چارک پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر x۱۰۵
- ۵-۴ نمودار خود همبستگی پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر x۱۰۶
- ۶-۴ نمودار دایره ای مقادیر گمشده در متغیر متراژ.....۱۱۴
- ۶-۴ نمودار دایره ای مقادیر گمشده در متغیر وضعیت اجتماعی - اقتصادی.....۱۱۵
- ۸-۴ شکل الگوی گمشدگی.....۱۱۶
- ۹-۴ نمودار اثر پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر متراژ.....۱۳۱
- ۱۰-۴ نمودار تاریخچه پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر متراژ.....۱۳۲
- ۱۱-۴ نمودار چگالی پسین پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر متراژ.....۱۳۳
- ۱۲-۴ نمودار چارک پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر متراژ.....۱۳۴
- ۱۳-۴ نمودار خود همبستگی پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر متراژ.....۱۳۵
- ۱۴-۴ نمودار خود اثر پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر متراژ.....۱۳۹

- ۱۵-۴ نمودار تاریخچه پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر متراژ.....۱۴۱
- ۱۶-۴ نمودار چگالی پسین پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر متراژ.....۱۴۲
- ۱۷-۴ نمودار چارک پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر متراژ.....۱۴۳
- ۱۸-۴ نمودار خود همبستگی پارامترها ی مدل رگرسیون لوجستیک با وجود مقادیر گمشده در متغیر متراژ.....۱۴۴

فصل اول :

کلیات تحقیق

۱-۱ : بیان مسئله

۱-۱-۱ : مقدمه

مدل رگرسیون لجستیک روشی تحلیلی است که به طور وسیعی در تحقیقات پزشکی و اپیدمیولوژیکی کاربرد دارد. با گسترش و تنوع این مدل ها تجزیه و تحلیل داده های حاصل از تحقیقات در علوم مختلف با به کار گیری این مدل ها روز به روز افزایش یافته و به لحاظ نظری لزوم تحقیق در زمینه های گوناگون این مدل ها را بیش از پیش فراهم نموده است.

هدف از تحلیل رگرسیون لجستیک همانند مدل های رگرسیون معمولی دستیابی به مدلی مناسب و در عین حال ساده جهت بررسی ارتباط بین متغیر پاسخ (وابسته)^۱ با یک یا مجموعه ای از متغیر های مستقل (کمکی)^۲ است. با این ویژگی که در این گونه مدل ها متغیر پاسخ بر خلاف رگرسیون معمولی عموماً از نوع رسته ای دو یا چند حالتی می باشد.

از آن جا که مدل رگرسیون لجستیک مدلی غیر خطی است، در آمار کلاسیک، تحلیل این گونه مدل ها مبتنی بر برآورد پارامترها از طریق ماکسیمم درستنمایی^۳ است، این روش ایده آل ترین روش است، چرا که در آن هیچ گونه محدودیتی برای متغیرهای مستقل در نظر گرفته نمی شود. به طور معمول دو روش متفاوت ماکسیمم درستنمایی برای برآورد پارامترهای مدل رگرسیون لجستیک وجود دارد که ماکسیمم درستنمایی غیر شرطی^۴ و شرطی^۵ نامیده می شود.

1. Response variable ، dependent variable
2. Independent (Covariate) variable
3. Maximum likelihood
4. Unconditional
5. Conditional

یک روش استنباط آماری، روش بیزی^۱ است که در آن اطلاعات ناشی از مطالعات قبلی و یا تجارب شخصی را می توان به کمک توزیع پیشین پارامترها در مدل لحاظ نمود. در روش بیزی لازم است با تلفیق داده های مشاهده شده و توزیع های پیشین^۲، توزیع های پسین^۳ پارامترهای مدل را به دست آورد و بر اساس آن ها استنباط درباره پارامترهای مدل را انجام داد. هدف این پژوهش بررسی و تحلیل بیزی با استفاده از تابع درستنمایی ساتن و کارول و مقایسه این روش با روش های واحد کامل^۴ و جانهای چندگانه^۵ می باشد. انتظار بر این است که با استفاده از روش روش بیزی برآوردهای دقیق تر و نتایج بهتری نسبت به روش های معمول به دست بیاوریم.

در بسیاری از مطالعات با مجموعه داده هایی مواجه می شویم که بخشی از آن ها گزارش نشده اند از قبیل خود داری از پاسخ، عدم تکمیل کامل پرسشنامه ها یا پرونده ها، ناقص بودن چارچوب مطالعه و غیره. در این صورت با داده های گمشده سروکار داریم که می تواند در متغیر پاسخ یا در متغیرهای کمکی بوجود آید. در این پژوهش گمشدگی در متغیرهای کمکی مورد نظر می باشد. مواجهه با داده گمشده مشکل بومی تحقیقات اجتماعی، پزشکی و اپیدمیولوژیکی است و در هنگام تحلیل، وجود این گونه موارد مشکلات عدیده ای را فراهم می سازد و عملاً تجزیه و تحلیل آماری را به سوی نتایج اریب سوق داده و نهایتاً دستیابی به یک نتیجه گیری مفید از داده های جمع آوری شده را با مشکل مواجه می سازد. ساده ترین روش برای تجزیه و تحلیل چنین داده هایی صرفه نظر کردن از موردهای دارای مقادیر گمشده و انجام آنالیز با داده های کامل می باشد (روش واحد کامل)، که این روش در عمل کارا نیست [1].

به طور کلی سه روش در نحوه بررسی داده های گمشده مورد استفاده قرار می گیرد [2] :

۱: روش های مبتنی بر واحد های کامل

۲: روش های مبتنی بر جانهای^۶

1 .Bayesian method
2 .Prior distributions
3 .Posterior distributions
4 .Complete Case
5 .Multiple Imputation
6 .Imputation

۳: روش های مبتنی بر مدل

تعریف :

داده ی گمشده : در تحلیل های آماری معمولاً داده ها به صورت مجموعه ای ماتریسی اند. سطرهای ماتریس، مشاهدات، آزمودنی ها^۱ یا واحدهای مورد مطالعه و ستون های آن متغیرها هستند. متغیرها به دو گروه متغیرهای مستقل یا کمکی و متغیرهای پاسخ تقسیم بندی می شوند. هر گاه برخی از مشاهدات این ماتریس اندازه گیری نشده و یا مقدار آن در دسترس نباشد با موضوع مقادیر گمشده سروکار داریم. گمشدگی مقادیر ممکن است در متغیر پاسخ یا متغیرهای کمکی رخ دهد. گم شدن در متغیرهای کمکی به دلایل مختلفی می تواند رخ دهد. به عنوان مثال در یک مطالعه گذشته نگر به علت نقص مدارک و سوابق ممکن است برخی از اطلاعات در دسترس نباشد، مانند مطالعات اپیدمیولوژیکی که از موارد ثبت شده در پرونده پزشکی بیماران برای تعیین مقادیر متغیرها استفاده می شود ممکن است پرونده ها برای همه بیماران شرکت کننده در پژوهش به طور کامل پر نشده باشند. در برخی مطالعات گمشدگی در مقادیر برخی از متغیرها ممکن است به علت نوع طرح پژوهشی رخ داده باشد، به طور مثال در مطالعه ای برای جمع آوری داده های مورد نظر، یک نمونه گیری دو مرحله ای انجام می شود، در مرحله اول مقادیر متغیرهایی که اندازه گیری آن ها آسان و ارزان است برای تمامی افراد شرکت کننده در مطالعه جمع آوری می شود و سپس در مرحله دوم مقادیر متغیرهایی که اندازه گیری آن ها پر هزینه و پیچیده می باشد برای زیر مجموعه ای از افراد شرکت کننده در مطالعه جمع آوری می شود [3]. همچنین این احتمال هست که به علت نقص یا ضعف دستگاه و تجهیزات، امکان مشاهده و اندازه گیری وجود نداشته باشد و یا در بعضی از مطالعات نظر سنجی، افرادی قادر به اظهار نظر دقیق نباشند. در این حالت ها با داده های گمشده روبرو هستیم [1,4].