



پژوهشکده آمار

جمهوری اسلامی ایران
مرکز آمار ایران
پژوهشکده‌ی آمار

معرفی متغیرهای نیم‌پویسته در آمارگیری‌های بزرگ

جانسی مقادیر کم‌شده‌ی متغیرهای پویسته، نیم‌پویسته و رسته‌آ

باروش مونت کارلومی بنحیر مارکوفی و تحلیل آماري این گونه داده‌ها

بسمه تعالی



پژوهشکده آمار

گروه پژوهشی پردازش داده‌ها و اطلاع‌رسانی

معرفی متغیرهای نیم‌پیوسته در آمارگیری‌های بزرگ، جانهی مقادیر گم‌شده‌ی
متغیرهای پیوسته، نیم‌پیوسته و رسته‌ای با روش مونت کارلوی زنجیر مارکوفی و
تحلیل آماری این‌گونه داده‌ها

مجتبی گنجعلی (مجری)

محمد رضا مشکانی

امید حمیدی

حسن رنجی

مرجان نورینی

فهرست مندرجات

۱	معرفی انواع متغیرهای رسته‌ای، پیوسته و نیم پیوسته با گم شدگی	۱
۱-۱	متغیرهای رسته‌ای، پیوسته و آمیخته‌ای از آنها (نیم پیوسته)	۱
۱-۱-۱	متغیرهای رسته‌ای و پیوسته	۲
۱-۱-۲	توابع توزیع آمیخته	۳
۲-۱	مسئله‌ی گم شدگی و جانهی در داده‌هایی با متغیرهای رسته‌ای، پیوسته و نیم پیوسته	۱۰
۱-۲-۱	مکانیسم‌های گم شدگی	۱۲
۲-۲-۱	تابع درست‌نمایی و توزیع پسین داده‌های مشاهده شده	۱۴
۳-۲-۱	روش‌های جانهی موجود برای داده‌های نیم پیوسته	۲۳
۳-۱	عملگر روفتن	۲۸
۴-۱	توزیع ویشارت وارونه	۳۰
۵-۱	توزیع دیریکله	۳۲

۳۳	۱-۵-۱ رابطه‌ی توزیع دیریکله با توزیع گاما
۳۵	۲ الگوریتم‌های EM و داده‌افزایی و مروری بر جانهی چندگانه بر اساس روش‌های مونت کارلوی زنجیر مارکوفی
۳۶	۱-۲ الگوریتم EM
۳۸	۱-۱-۲ گام E و گام M از الگوریتم EM
۳۹	۲-۱-۲ الگوریتم EM در خانواده‌ی نمایی
۴۷	۲-۲ روش‌های مونت کارلوی زنجیر مارکوفی (MCMC)
۴۹	۱-۲-۲ نمونه‌گیری گیبس
۵۰	۲-۲-۲ داده‌افزایی
۵۲	۳-۲-۲ گام I
۵۷	۳-۲ مروری بر جانهی چندگانه
۶۳	۴-۲ استفاده از نرم‌افزارهای SAS 9.1، SPSS 13 و Norm 2.03 در الگوریتم‌های EM ، داده‌افزایی و جانهی چندگانه
۶۳	۱-۴-۲ استفاده از نرم‌افزار SAS 9.1 در جانهی چندگانه
۷۴	۲-۴-۲ استفاده از نرم‌افزار Norm 2.03 در شیوه‌های EM و داده‌افزایی
۷۷	۳-۴-۲ استفاده از نسخه‌ی ۱۳ نرم‌افزار SPSS برای الگوریتم EM
۸۵	۳ جانهی متغیرهای رسته‌ای و پیوسته با استفاده از روش‌های مونت کارلوی زنجیر مارکوفی
۸۶	۱-۳ معرفی مدل مکانی عام
۸۸	۱-۱-۳ درست‌نمایی داده‌های کامل

- ۹۰ ۲-۱-۳ استنباط بیزی از داده‌های کامل
- ۹۳ ۲-۳ الگوریتم‌های EM و داده افزایشی برای مدل مکانی عام
- ۹۳ ۱-۲-۳ توزیع‌های پیشگو
- ۹۷ ۲-۲-۳ الگوریتم EM برای مدل مکانی عام با داده‌های ناکامل
- ۱۰۱ ۳-۲-۳ الگوریتم داده افزایشی برای مدل مکانی عام با داده‌های ناکامل
- ۱۰۶ ۳-۳ نیکویی برازش مدل مکانی عام
- ۱۰۶ ۱-۳-۳ نیکویی برازش مدل مکانی عام با داده‌های کامل
- ۱۱۲ ۲-۳-۳ مدل مکانی عام با داده‌های گم‌شده

۴ جانمایی متغیرهای رسته‌ای، پیوسته و نیم‌پیوسته با استفاده از روش‌های

۱۱۷ مونت کارلوی زنجیر مارکوفی

- ۱۱۸ ۱-۴ معرفی مدل مکانی عام بلوک شده
- ۱۱۹ ۱-۱-۴ تعمیم مدل مکانی عام
- ۱۲۰ ۲-۱-۴ تعریف بلوک‌ها
- ۱۲۱ ۳-۱-۴ فرمول بندی مدل مکانی عام بلوک شده
- ۱۲۵ ۴-۱-۴ تولید جانمایی تحت مدل مکانی عام بلوک شده
- ۱۲۶ ۲-۴ الگوریتم‌های EM و داده افزایشی برای مدل مکانی عام بلوک شده
- ۱۲۶ ۱-۲-۴ الگوریتم EM برای مدل مکانی عام بلوک شده
- ۱۳۰ ۲-۲-۴ الگوریتم داده افزایشی برای مدل مکانی عام بلوک شده
- ۱۳۱ ۳-۴ آزمایش شبیه سازی

- ۵ بررسی وضعیت موجود مرکز آمار ایران در خصوص جانمایی داده‌های گم‌شده و تحلیل داده‌های طرح منتخب ۱۳۵
- ۱-۵ بررسی وضعیت موجود مرکز آمار ایران در زمینه‌ی جانمایی متغیرهای نیم‌پیوسته . ۱۳۵
- ۱-۱-۵ جانمایی در طرح‌های آمارگیری از کارگاه‌های حمل و نقل زمینی . . ۱۳۶
- ۲-۱-۵ جانمایی در طرح‌های صنعت و معدن ۱۳۸
- ۳-۱-۵ طرح آمارگیری از کارگاه‌های بازرگانی (خرده‌فروشی و تعمیر کالاهای شخصی و خانگی، فروش، نگهداری و تعمیر وسایل نقلیه‌ی موتوری و . . . و عمده‌فروشی و حق‌العمل‌کاری) ۱۳۹
- ۲-۵ تحلیل داده‌های آمارگیری از فعالیت‌های تحقیق و توسعه‌ی کشور در سال ۱۳۸۴ ۱۴۱

فصل ۱

معرفی انواع متغیرهای رسته‌ای، پیوسته و نیم پیوسته با گم شدگی

در این فصل در ابتدا انواع متغیرهای رسته‌ای، پیوسته و نیم پیوسته تعریف شده‌اند. در ادامه مسئله گم شدگی و جانهی در داده‌هایی با متغیرهای آمیخته (شامل رسته‌ای، پیوسته و نیم پیوسته) بررسی شده و مکانیسم‌های گم شدگی و روش‌های موجود جانهی برای این گونه داده‌ها آورده شده‌اند. برخی تعاریف لازم مورد استفاده در فصل‌های بعدی نیز در انتها ذکر شده‌اند.

۱-۱ متغیرهای رسته‌ای، پیوسته و آمیخته‌ای از آنها (نیم پیوسته)

در این بخش ابتدا تعریف متغیرهای رسته‌ای، ترتیبی و پیوسته آمده است. سپس توابع توزیع آمیخته و انواع آن تعریف شده‌اند.

۱-۱-۱ متغیرهای رسته‌ای و پیوسته

متغیرها را با چهار نوع مقیاس، اندازه‌گیری می‌کنند. متغیرهای کیفی که بدون هیچ ترتیب خاصی دسته‌بندی شده‌اند، اسمی نامیده می‌شوند. از این نوع متغیرها می‌توان مذهب (مسلمان و غیر مسلمان)، نوع وسیله‌ی نقلیه مورد استفاده (اتومبیل، دوچرخه، موتورسیکلت، مترو، اتوبوس و...) و نوع محل اقامت (آپارتمان، خانه و...) را نام برد. برای متغیرهای اسمی ترتیب فهرست کردن رسته‌ها مهم نیست و تحلیل آماری وابسته به آن ترتیب نیست. برخی متغیرها در مقیاس خود دارای نوعی ترتیب هستند چنین متغیرهایی ترتیبی نامیده می‌شوند. برای مثال وضعیت اجتماعی (بالا، متوسط و پایین)، شرایط بیمار (خوب، متوسط و خیلی بد) متغیرهای ترتیبی هستند. در این مقیاس فاصله‌ی بین مقادیر معلوم نیست، هرچند شخصی که برای مثال وضعیت اجتماعی متوسط دارد، وضعیت اجتماعی وی بهتر از شخصی است که وضعیت اجتماعی پایین دارد، اما هیچ مقدار عددی نمی‌تواند فاصله‌ی بین متوسط و بد را اندازه‌گیری کند. یک متغیر بازه‌ای متغیری است که بین هر دو مقدار آن بازه‌ای از اعداد وجود دارد ولی این متغیرها لزوماً یک مبدأ خاص تعریف شده ندارند. برای مثال سطح فشار خون و درجه حرارت از این نوع متغیرها هستند. متغیرهای نسبتی متغیرهایی هستند که در مقیاس بازه‌ای تعریف شده‌اند اما مقیاس آنها دارای یک مبدأ تعریف شده است برای مثال درجه حرارت در مقیاس بازه‌ای اندازه‌گیری می‌شود ولی دارای صفر تعریف شده نیست. ولی متغیرهای طول عمر، درآمد و قد در مقیاس نسبتی اندازه‌گیری می‌شوند و مبدأ در آنها تعریف شده است (به طور عام صفر).

طریقه‌ای که یک متغیر را اندازه‌گیری می‌گیرد نوع آن را مشخص می‌کند. برای مثال آموزش متغیری اسمی است اگر به صورت نوع آموزشگاه (انتفاعی و غیرانتفاعی) رسته‌بندی شود، اگر متغیر آموزش به صورت بی‌سواد، سیکل، دیپلم و... ثبت شود متغیری ترتیبی است و چنانچه این متغیر طول آموزش فرد را نشان دهد متغیری نسبتی خواهد بود. مقیاس اندازه‌گیری متغیر است که تعیین می‌کند چه روش آماری برای تحلیل در مورد آن متغیر مناسب است. روش‌های آماری مورد استفاده برای یک متغیر با مقیاس اندازه‌گیری خاص، نمی‌تواند برای یک متغیر با مقیاس اندازه‌گیری دیگر استفاده شود. متغیرهای اسمی و ترتیبی از نوع رسته‌ای و متغیرهای پیوسته از نوع بازه‌ای یا نسبتی هستند. در این طرح متغیرهایی را مورد بررسی قرار می‌دهیم که از نوع اسمی یا بازه‌ای یا نسبتی باشند، متغیرهایی را

نیز در نظر خواهیم گرفت که از نوع نیم پیوسته هستند (به تعریف آن در بخش بعد توجه کنید).

۱-۱-۲ توابع توزیع آمیخته

فرض کنید X متغیری تصادفی است که نتیجه‌ی یک آزمایش را نشان می‌دهد و مقادیری که می‌تواند اختیار کند از مجموعه‌ای مانند S زیر مجموعه‌ای از \mathbf{R} است. متغیر X دارای تابع توزیع آمیخته است اگر بتوان S را به دو زیر مجموعه مانند C و D افراز کرد به گونه‌ای که:

$$(۱) \quad 0 < P(X \in D) < 1$$

$$(۲) \quad P(X = x) = 0 \text{ برای هر } x \text{ متعلق به } C$$

$$(۳) \quad P(X \in D) + P(X \in C) = 1$$

بنابراین بخشی از جرم احتمال X روی نقاط متعلق به مجموعه‌ی گسسته D متمرکز می‌شود و بقیه‌ی آن به طور پیوسته روی C پخش می‌شود.

فرض کنید $p = P(X \in D)$ و $0 < p < 1$ ، در این صورت می‌توانیم یک تابع روی D تعریف کنیم. این تابع را تابع احتمال گسسته‌ی جزئی^۱ نامند. فرض کنید به ازای هر $x \in D$ ، $g(x) = P(X = x)$ باشد. در این صورت:

$$g(x) > 0, \quad \forall x \in D$$

$$\sum_{x \in D} g(x) = p$$

$$P(X \in A) = \sum_{x \in A} g(x), \quad A \subseteq D$$

قسمت پیوسته‌ی توزیع به وسیله‌ی یک تابع چگالی جزئی تعریف می‌شود. برای این فرض کنید یک تابع نامنفی مثل h روی C وجود داشته باشد به گونه‌ای که:

$$P(X \in A) = \int_A h(x) dx, \quad A \subseteq C$$

در این صورت:

^۱ Partial Discrete Probability Function

$$\int_C h(x)dx = 1 - p$$

دقت کنید که به ازای همه‌ی مقادیر X متعلق به D داریم $h(x) = 0$ و به ازای همه‌ی مقادیر x متعلق به C ، داریم $g(x) = 0$. فرض کنید که $A \subseteq S$ در این صورت:

$$P(X \in A) = \sum_{x \in A} g(x) + \int_A h(x)dx$$

به آسانی نتیجه می‌شود:

$$f_1(x | X \in D) = g(x)/p, \quad \forall x \in D$$

$$f_2(x | X \in C) = h(x)/(1 - p), \quad \forall x \in C$$

که در آن f_1 و f_2 چگالی‌های شرطی x هستند. بنابراین توزیع x آمیخته‌ای از یک توزیع گسسته و یک توزیع پیوسته است.

متغیرهای بریده^۱

متغیرهای آمیخته معمولاً زمانی رخ می‌دهند که یک متغیر تصادفی با توزیع پیوسته از نقطه‌ای به بعد بریده شود. برای مثال فرض کنید T طول عمر یک دستگاه و به ازای $t > 0$ تابع چگالی آن باشد. هنگام آزمایش این دستگاه نمی‌توان مدت زمان بی‌نهایتی منتظر ماند، بنابراین ممکن است یک عدد ثابت مثبت مانند a را انتخاب کنیم و مقدار متغیر زیر را ثبت کنیم. اگر مقادیر متغیر تصادفی U را

^۱ Truncated

با استفاده از متغیر تصادفی T به صورت

$$U = \begin{cases} T, & T < a \\ a, & T \geq a \end{cases}$$

تعریف کنیم در این صورت U یک توزیع آمیخته دارد. لذا داریم:

$$D = \{a\}, \quad g(a) = \int_{\{t:t>a\}} f_T(t) dt$$

$$C = (0, a), \quad h(t) = f_T(t), \quad 0 < t < a$$

فرض کنید X یک توزیع پیوسته روی \mathbf{R} (مجموعه‌ی اعداد حقیقی) با تابع چگالی f_X باشد. اگر متغیر تصادفی X در a و b ، $(b > a)$ ، بریده شود متغیر جدید (Y) با مقادیر زیر به وجود می‌آید:

$$Y = \begin{cases} a, & X \leq a \\ X, & a < X < b \\ b, & X \geq b \end{cases}$$

آنگاه Y دارای توزیع آمیخته است و نتایج زیر را داریم:

$$D = \{a, b\}, \quad g(a) = \int_{\{x:x<a\}} f_X(x) dx, \quad g(b) = \int_{\{x:x>b\}} f_X(x) dx$$

$$C = (a, b), \quad h(x) = f_X(x), \quad a < x < b.$$

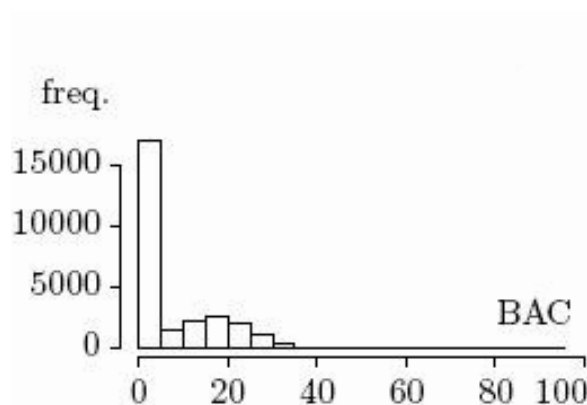
متغیرهای نیم پیوسته

متغیر نیم پیوسته متغیری است که مقداری تک (معمولا صفر) را با احتمالی مثبت و مقادیری پیوسته در یک بازه را طبق یک چگالی پیوسته می‌تواند اختیار کند. این متغیرها متفاوت از متغیرهای از چپ سانسور شده^۱ یا بریده هستند، چرا که در آنها احتمالی که به صفر نسبت داده می‌شود به جهت پاسخ‌های منفی یا داده‌های گم شده نیست. چنین داده‌هایی در زمینه‌های مختلف به وفور یافت

^۱ Left Censored

می‌شوند. متغیرها در بررسی‌های اقتصادی و مربوط به درآمد یا هزینه‌ها اکثراً نیم پیوسته هستند، چون بسیاری از افراد یک نوع خاص از هزینه را ندارند. در مطالعات مصرف مواد مخدر و یا الکل، بسیاری از افراد در طول دوره‌ی مطالعه، از این مواد استفاده نمی‌کنند.

بافت نگار نشان داده شده در شکل ۱-۱ را در نظر بگیرید. این بافت نگار نتایج آزمایش مقدار الکل خون (BAC)^۲ از ۲۷۶۳۳ راننده وسایل نقلیه و عابران پیاده در تصادف‌های شدید در بزرگراه‌های آمریکا در سال ۱۹۹۳ را نشان می‌دهد، این بررسی که توسط شی‌فر و اولسن^۳ (۱۹۹۹) انجام شده است، نشان می‌دهد که برای ۵۷/۳ درصد از موارد نشان داده شده در شکل ۱-۱، مقدار BAC صفر است که نمایانگر عدم مصرف الکل است، سایر مقادیر که از ۱ تا ۹۴ تغییر می‌کنند، نشان دهنده‌ی سطوح الکل خون از ۰/۰۱ تا ۰/۹۴ گرم در هر دسی لیتر (g/dl) است.



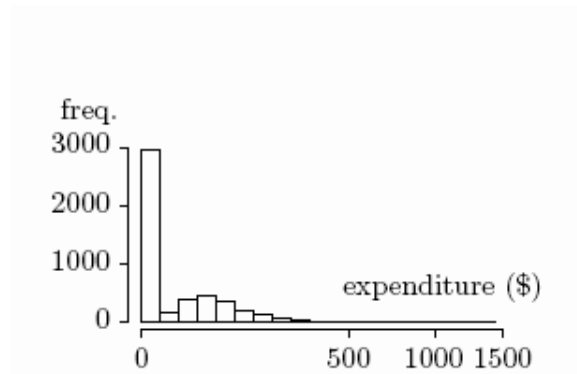
نمودار ۱-۱: بافت نگار مقدار الکل خون در بین رانندگان

بافت نگار شکل ۱-۲ نشان دهنده‌ی هزینه‌ی خرید کفش در فصل اول سال ۱۹۹۶ به وسیله‌ی ۴۸۷۶ واحد مصرف کننده است که توسط اداره‌ی آمار آمریکا در بررسی هزینه‌ی مصرف کننده جمع آوری شده است (شی‌فر و اولسن، ۱۹۹۹). یک واحد مصرف کننده یک خانوار یا گروهی از افراد است که در آمد مشترک دارند. برای کاهش چولگی ظاهر شده، جذر مقادیر هزینه رسم شده اند.

^۲ Blood alcohol content

^۳ Schafer and Olsen

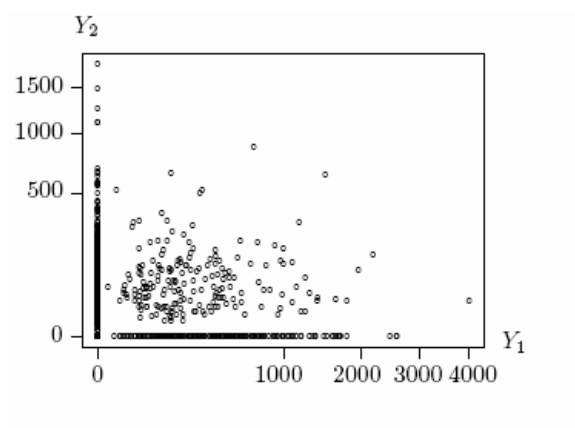
نزدیک ۶۱ درصد از مقادیر مشاهده شده صفر است و بقیه بین ۱ تا ۱۳۱۵ دلار توزیع شده‌اند.



نمودار ۱-۲: بافت نگار هزینه‌ی کفش

در هر یک از این مثال‌ها کمیت مورد نظر نیم پیوسته است، یعنی آمیخته‌ای از صفر و مقادیر مثبتی است که به طور پیوسته توزیع شده‌اند. متغیرهای نیم پیوسته برای محقق‌ی که پیوند بین این متغیرها را با متغیرهای دیگر ارزیابی می‌کند مشکلاتی به وجود می‌آورد، زیرا این روابط نسبتاً پیچیده هستند. برای نشان دادن انواع روابطی که یافت می‌شوند، در شکل ۱-۳ دو متغیر دیگر در همان دوره‌ی سه ماهه در بررسی هزینه‌ی مصرف کننده رسم شده‌اند که y_1 جذر هزینه‌های خرید لوازم بزرگ منزل (مثل یخچال) و y_2 جذر هزینه‌های خرید لوازم کوچک منزل است. در بین ۴۸۷۶ واحد مصرف کننده که مورد بررسی قرار گرفته‌اند $70/3$ درصد هیچ هزینه‌ای نداشتند ($y_1 = y_2 = 0$)، $5/3$ درصد لوازم بزرگ خریداری کرده‌اند اما لوازم کوچک خریداری نکرده‌اند ($y_1 > 0, y_2 = 0$)، $20/5$ درصد لوازم کوچک خریداری کرده‌اند اما لوازم بزرگ خریداری نکرده‌اند ($y_1 = 0, y_2 > 0$) و $3/9$ درصد هر دو را خریداری کرده‌اند ($y_1 > 0, y_2 > 0$). بین آنهایی که لوازم کوچک خریداری کرده‌اند ($y_2 > 0$) همبستگی بین $y_1^{\frac{1}{2}}$ و نشانگر دودویی برای $y_1 > 0$ برابر $r = +0/11$ است. میزان مخارج لوازم کوچک دارای یک رابطه‌ی معنی دار مثبت با خرید لوازم بزرگ است. از بین آنهایی که لوازم بزرگ خریداری کرده‌اند ($y_1 > 0$)، میزان مخارج لوازم بزرگ $y_1^{\frac{1}{2}}$ دارای رابطه‌ی معنی داری منفی با نشانگر دودویی برای خرید لوازم کوچک است ($r = -0/12$). بالاخره بین آنهایی که هر دو را خریداری

کرده‌اند ($y_1 > 0, y_2 > 0$) همبستگی بین $y_1^{\frac{1}{3}}, y_2^{\frac{1}{3}}$ تقریباً صفر است ($r = +0/0/3$).



نمودار ۱-۳: نمودار پراکنش هزینه‌های لوازم بزرگ و لوازم کوچک منزل

تابع توزیع تجمعی متغیر نیم پیوسته

تابع توزیع تجمعی $F^{sc}(x)$ برای متغیر نیم پیوسته را می‌توان به صورت

$$F^{sc}(x) = pF^{dg}(x) + (1-p)F^{ac}(x) \quad 0 < p < 1, \quad (1.1)$$

بیان کرد که $F^{dg}(\cdot)$ و $F^{ac}(\cdot)$ به ترتیب تابع توزیع تجمعی بخش گسسته و پیوسته است. $F^{dg}(\cdot)$ تباهیده در نقطه‌ی مفروض a و $F^{ac}(\cdot)$ مطلقاً پیوسته است. تابع توزیع $F_X^{dg}(\cdot)$ در نقطه‌ی x به صورت

$$F_X^{dg}(x) = \begin{cases} 0 & \text{اگر } x < a \\ 1 & \text{اگر } x \geq a \end{cases}$$

است. $F^{ac}(x)$ برای مقادیر کوچکتر یا مساوی a برابر صفر است.

مثال ۱.۱ متغیر نیم پیوسته‌ی نمایی. مدت توقف راننده‌ی یک وسیله‌ی نقلیه را در مقابل تابلوی ایست یک خیابان یک طرفه در نظر بگیرید. اگر X متغیری تصادفی باشد که مدت توقف او را تا پایان

ایست نشان می دهد احتمال معینی وجود دارد که هیچ وسیله‌ی نقلیه‌ای در حال گذر از تقاطع نباشد و این راننده بتواند بدون توقف به حرکت خود ادامه دهد. از سوی دیگر اگر مجبور باشد توقف کند وقفه‌ی او می تواند هر مقدار از زمان‌های ممکن باشد. این آزمایش را با فرض این که X دارای تابع توزیع معلوم به صورت $F_X^{sc}(x) = [1 - (1 - p)e^{-\lambda x}]$ که $x \geq 0$ می باشد می توان مدل بندی کرد. این $F_X^{sc}(x)$ در $x = 0$ دارای جهشی به اندازه‌ی p است اما برای $x > 0$ پیوسته با توزیع نمایی است.

روش مدل بندی داده‌های نیم پیوسته نشان دادن آنها به صورت یک آمیخته‌ی دو بخشی از یک توزیع پیوسته و یک جرم نقطه‌ای تباهیده می باشد. فرض کنید مشاهدات نیم پیوسته y_1, y_2, \dots, y_n به صورت دو متغیر (w_i, z_i) ، $i = 1, 2, \dots, n$ دوباره کدگذاری شده باشند که

$$w_i = \begin{cases} 1 & \text{اگر } y_i \neq 0 \\ 0 & \text{اگر } y_i = 0 \end{cases} \quad (2.1)$$

و

$$z_i = \begin{cases} g(y_i) & \text{اگر } w_i = 1 \\ \text{نامرتبط} & \text{اگر } w_i = 0 \end{cases} \quad (3.1)$$

و g تابع یکنوای افزایشی (برای مثال لگاریتم) باشد که برای نرمال کردن تقریبی مقادیر غیر صفر انتخاب شده است. فرض می شود نشانگرهای دودویی w_i دارای توزیع برنولی با $p(w_i = 1) = \pi$ و توزیع شرطی z_i به شرط $w_i = 1$ با پارامترهای μ و σ^2 باشد و وقتی $w_i = 0$ مقدار $y_i = 0$ بوده و در توزیع z_i ($z_i > 0$) نقشی نداشته باشد در این صورت گوییم z_i نامرتبط است. تحت این فرض‌ها، درست‌نمایی برای $\theta = (\pi, \mu, \sigma^2)$ به توابع مجزا برای π و (μ, σ^2) به صورت زیر تجزیه می شود:

$$L(\theta) \propto \prod_{i=1}^n \pi^{w_i} (1 - \pi)^{1-w_i} \prod_{i:y_i \neq 0} \sigma^{-1} \exp \left\{ -\frac{(z_i - \mu)^2}{2\sigma^2} \right\}, \quad (4.1)$$

و برآوردهای ML به فرم بسته برای پارامترها به صورت زیر به دست می آید:

$$\begin{aligned} \hat{\pi} &= n^{-1} \sum_{i=1}^n w_i \\ \hat{\mu} &= n_1^{-1} \sum_{i:y_i \neq 0} z_i \\ \hat{\sigma}^2 &= n_1^{-1} \sum_{i:y_i \neq 0} (z_i - \hat{\mu})^2 \end{aligned}$$

که در آن $n_1 = \sum_{i=1}^n w_i$.

۱-۲ مسئله‌ی گم شدگی و جانهی در داده‌هایی با متغیرهای رسته‌ای،

پیوسته و نیم پیوسته

تحلیل آماری داده‌های گم شده غالباً پیچیده است، به خصوص در بررسی‌هایی که موارد عدم پاسخ فراوانند. روش عمومی در برخورد با چنین داده‌هایی جانهی آنها قبل از انجام هرگونه تحلیل آماری است. این روش، تحلیل داده‌ها را به دو مرحله تقسیم می‌کند. نخست، به وسیله جانهی داده‌های گم شده به مجموعه‌ای از داده‌های کامل تبدیل می‌شود. دوم، در مرحله تحلیل، مجموعه‌ی داده‌های جانهی شده، همانند وقتی که داده‌ی گم شده وجود ندارد، تحلیل می‌شوند. معمولاً جانهی با کمک شخصی که داده‌ها را جمع‌آوری کرده و نیز کمک گرفتن از اطلاعات اضافی موجود و استفاده از ابزارهای آماری موجود (مثل، روش‌های مونت کارلو) انجام می‌شود. مزیت اصلی جانهی داده‌های گم شده این است که از پیچیدگی‌ها در تحلیل‌های آماری که به وسیله‌ی داده‌های گم شده به وجود می‌آیند، جلوگیری می‌کند. بنابراین از نقطه نظر محاسباتی، جانهی مرحله‌ی تحلیل را ساده می‌کند، به طوری که در تحلیل تنها روش‌های استاندارد (روش‌هایی مربوط به تحلیل داده‌های بدون مقادیر گم شده) مورد نیاز است. مهم‌ترین که، این روش از دانش و اطلاعات شخص جمع‌آوری کننده‌ی داده‌ها در جانهی استفاده می‌کند.

متأسفانه، قطعی بودن کامل جانهی نسبتاً غیر معمول است، یعنی حتی با اطلاعات جمع‌آوری کننده ماهر داده‌ها، در جانهی عدم اطمینان وجود دارد. استنباط معتبر بایستی این عدم اطمینان را در نظر بگیرد. جانهی چندگانه (MI)^۱ (رایین^۲، ۱۹۸۷) یک روش پذیرفته شده برای به دست آوردن چنین استنباط معتبر با وجود داده‌های گم شده می‌باشد (مروری بر جانهی چندگانه در فصل ۲، ارائه خواهد شد).

معمولاً در جانهی چندگانه یک روش، بر اساس مدل بندی، برای تولید m جانهی استفاده می‌شود

^۱ Multiple Imputation

^۲ Rubin

که برآمدهای تصادفی تکرار شده از یک مدل را برای عدم پاسخ نشان می‌دهد. در مرحله تحلیل، هر یک از m مجموعه داده‌ی جهانی شده همانند یک مجموعه داده‌ی کامل، تحلیل می‌شوند و m مجموعه‌ی برآوردها و خطاهای مربوط بر طبق قانون‌های ساده ترکیب می‌شوند. توان جهانی چندگانه در توانایی ارائه‌ی استنباط‌های معتبر می‌باشد که تنها نیاز به روش‌های آماری استاندارد در مرحله‌ی تحلیل دارد. یعنی این روش تغییرات اضافی در نتیجه‌ی مقادیر گم شده تحت مدل جهانی را منعکس می‌کند.

البته مهم است که مدل جهانی برای مجموعه داده‌های خاص مناسب باشد. برای انواع معین داده‌ها، مدل‌های موجود می‌توانند مدل‌های جهانی مناسبی باشند. به عنوان مثال مدل مکانی عام (GLOM)^۱ (اولکین و تیت^۲، ۱۹۶۱ و شی‌فر، ۱۹۹۷) یک مدل جهانی عام است که برای جهانی داده‌های ناقص با متغیرهای پیوسته و رسته‌ای استفاده می‌شود. این مدل فرض می‌کند که به شرط سطوح مختلف متغیرهای رسته‌ای، متغیرهای پیوسته یا تبدیلی از آنها دارای توزیع نرمال چند متغیره با میانگینی است که وابسته به سطح متغیرهای رسته‌ای می‌باشد ولی واریانس دارد که در سطوح مختلف متغیرهای رسته‌ای ثابت است. اگر این فروض برقرار باشند این مدل، مدل مناسبی برای استفاده در جهانی است. متغیرهای نیم پیوسته با مقادیر گم شده نیز جمع آوری کننده داده‌ها را که نیاز به جهانی مقادیر گم شده این متغیرها دارد دچار مشکل می‌کنند. روش‌های جهانی که برای داده‌های نیم پیوسته مناسب نیست، ممکن است به طور جدی توزیع کناری متغیرها یا روابط آن‌ها با سایر کمیت‌ها را تحریف کند. یک روش جهانی معقول بایستی شکل توزیعی همانند آنچه در شکل‌های ۱-۱ و ۱-۲ نشان داده شده است را حفظ کند. روش‌های جهانی جدید بر اساس مدل سازی شامل جهانی چندگانه برای داده‌های چند متغیره پیوسته و رسته‌ای توسط شی‌فر (۱۹۹۷) گسترش داده شده است.

توجه به دلایلی که موجب گم شدن داده‌ها می‌شوند اهمیت زیادی دارد زیرا با استفاده از آن‌ها می‌توان روش‌های تحلیل واقعی را مشخص کرد. دلایل گوناگونی برای داده‌های گم شده وجود دارند. برخی داده‌ها ممکن است به طور کاملاً تصادفی گم شده باشند در حالی که برخی ممکن است به طور تصادفی و یا غیر تصادفی گم شده باشند. اغلب راهی برای حصول اطمینان از دلایل گم شدن

^۱ The General Location Model

^۲ Olkin and Tate

داده‌ها وجود ندارد. در مطالعات طولی، ممکن است افراد از همکاری در آزمایش انصراف دهند و یا در دوره‌های جمع‌آوری داده‌ها در دسترس نباشند. وقتی داده‌ها با استفاده از پرسشنامه جمع‌آوری می‌شوند، ممکن است افراد برخی از سؤال‌ها را به علت نداشتن وقت و یا علاقه پاسخ ندهند و کامل نکنند. ممکن است افراد به جواب دادن به سؤال‌هایی که آنها را مضطرب یا شرمنده کند تمایل نداشته باشند. ممکن است برخی سؤال‌ها خوب بیان نشده باشند و پاسخ گو نداند که چگونه پاسخ دهد. حتی اگر پرسشنامه بسیار معتبر باشد و مشکلات فوق را نداشته باشد، ممکن است برخی سؤال‌ها در مورد بعضی از افراد مصداق نداشته باشد.

در تحلیل بعضی از داده‌ها عواملی که منجر به داده‌های گم شده می‌شوند، صراحتاً دخالت داده نمی‌شوند که در چنین مواردی فرضی تحت عنوان «قابل چشم‌پوشی» در نظر گرفته می‌شود. البته ممکن است با در نظر گرفتن یک متغیر نشانگر پاسخ (یک برای داده‌های ثبت شده و صفر برای داده‌های گم شده) عاملی را برای داده‌های گم شده در مدل وارد کنند.

عاملی که منجر به داده‌های گم شده می‌شود عموماً قابل چشم‌پوشی نیست. مثلاً در تحقیق درباره‌ی درآمد افراد ممکن است بی‌پاسخ ماندن، مربوط به مقدار خود درآمدها باشد یعنی مثلاً تمام درآمدهایی که از یک مقدار کمتر بوده‌اند مشاهده نشده‌اند و در این حالت دیگر نمی‌توان اظهار داشت که گم شدن داده‌ها تصادفی بوده بلکه عامل کم بودن درآمد باعث گم شدگی آن است و در نظر نگرفتن این عوامل در تحلیل داده‌ها می‌تواند ایجاد مشکل کند. این عوامل مکانیسم‌های گم شدگی را به وجود می‌آورند که در زیر آنها را مورد بحث قرار می‌دهیم.

۱-۲-۱ مکانیسم‌های گم شدگی

مجموعه‌ی داده‌ها را به صورت یک ماتریس $n \times p$ که سطرهای آن مشاهداتی از یک توزیع احتمالی چند متغیره‌اند و متغیرها از هم مستقل هستند، در نظر بگیرید. این ماتریس شامل n سطر (تعداد مشاهدات) و p ستون (تعداد متغیرهای ثبت شده برای این واحدها) است.

فرض کنید ماتریس Y نشان دهنده‌ی ماتریس $n \times p$ داده‌ها باشد، که به طور کامل مشاهده نشده‌اند، و y_i سطر i ام Y برای $i = 1, 2, \dots, n$ باشد. با فرض داشتن توزیع یکسان و مستقل، تابع

چگالی یا تابع احتمال داده‌های کامل را به صورت زیر نشان می‌دهیم

$$P(Y | \theta) = \prod_{i=1}^n f(y_i | \theta), \quad (5.1)$$

که f تابع چگالی یا تابع احتمال یک سطر است و θ بردار پارامترهای مجهول است. همچنین فرض کنید بخش مشاهده شده‌ی Y را به وسیله‌ی Y_{obs} و بخش گم شده‌ی Y را به وسیله‌ی Y_{mis} نشان دهیم، به طوری که $Y = (Y_{obs}, Y_{mis})$. طبق تعریف رابین (۱۹۷۶)، چنانچه گم شدگی بستگی به Y_{obs} داشته باشد، ولی به شرط Y_{obs} به Y_{mis} بستگی نداشته باشد، گم شدگی به صورت تصادفی است.

فرض کنید U و V دو متغیر باشند. همچنین فرض کنید واحدهایی از U که مشاهده شده‌اند و مقدار مشاهده‌ی آن u باشد، مد نظر هستند. گم شدگی تصادفی V به این معنی است که بین این واحدها، توزیع V در بین بخشی از V که مشاهده شده و بخشی که مشاهده نشده به شرط $U = u$ ، یکسان است. یعنی احتمال ثبت V فقط به مقادیر U بستگی داشته باشد ولی به مقادیر خود V بستگی نداشته باشد.

گم شدگی تصادفی به این معنی نیست که مقادیر داده‌های گم شده به صورت نمونه‌ی تصادفی ساده از کل داده‌ها هستند. حالت بعدی گم شدگی، گم شدگی کاملاً تصادفی نامیده می‌شود. گم شدگی کاملاً تصادفی تنها حالت خاصی از گم شدگی تصادفی است. گم شدگی تصادفی محدودیت کمتری نسبت به گم شدگی کاملاً تصادفی دارد، زیرا تنها نیاز به این دارد که مقادیر گم شده همانند یک نمونه‌ی تصادفی از همه‌ی مقادیر داخل زیررده‌های تعریف شده به وسیله‌ی مقادیر داده‌های مشاهده شده رفتار کنند. به عبارت دیگر، گم شدگی تصادفی اجازه می‌دهد که احتمال گم شدن یک داده به شرط تعداد کمیت‌هایی که مشاهده شده‌اند، بستگی به خود داده نداشته باشد.

رابین (۱۹۷۶) گم شدگی تصادفی را برحسب مدل احتمالاتی برای گم شدگی تعریف کرد. فرض

کنید $R = (R_{ij})$ یک ماتریس نشانگر $n \times p$ باشد که اعضای آن به صورت زیر تعریف می‌شوند:

$$R_{ij} = \begin{cases} 1 & \text{اگر } y_{ij} \text{ مشاهده شود} \\ 0 & \text{در غیر این صورت} \end{cases}$$

به طور کلی نباید انتظار داشت که توزیع R با Y نامرتب باشد، بنابراین می‌توان یک مدل احتمالاتی برای R به صورت $P(R | Y, \xi)$ در نظر گرفت، که بستگی به Y و پارامترهای مجهول ξ دارد. فرض گم شدگی تصادفی به این معنی است که این توزیع بستگی به Y_{mis} ندارد، یعنی

$$P(R | Y_{obs}, Y_{mis}, \xi) = P(R | Y_{obs}, \xi). \quad (6.1)$$

مکانیسم گم شدگی غیر تصادفی نه حالت گم شدگی تصادفی و نه حالت گم شدگی کاملاً تصادفی است در این حالت مکانیسم گم شدن به داده‌هایی که ممکن است گم شوند وابسته است. در دو حالت گم شدگی کاملاً تصادفی و گم شدگی تصادفی می‌توان از مکانیسم گم شدن داده‌ها چشم پوشی کرد ولی در گم شدن غیر تصادفی نمی‌توان از عامل گم شدن داده‌ها چشم پوشی کرد که این حالت به مدل غیر قابل چشم پوشی موسوم است.

مجزا بودن پارامترها

برای کنکاش بیشتر در این زمینه، نیاز به این فرض داریم که پارامترهای مدل داده‌ها θ و پارامترهای مکانیسم گم شدگی ξ مجزا باشند. از دیدگاه فراوانی‌گرا، این بدان معنی است که بایستی فضای توأم پارامتر (θ, ξ) حاصل ضرب دکارتی فضاهای پارامترهای θ و ξ باشد. به این معنی که θ به صورت تابعی به ξ وابسته نباشد. از دیدگاه بیزی، این بدان معنی است که بایستی هر توزیع توأم پیشین که برای (θ, ξ) بکار برده می‌شود، حاصل ضرب توزیع‌های حاشیه‌ای پیشین θ و ξ باشد. در بسیاری از شرایط، منطقی است که θ حاوی مقدار کمی اطلاعات در مورد ξ و بالعکس باشد. اگر فرض گم شدگی تصادفی و مجزا بودن همزمان اتفاق بیافتند، آنگاه گوییم مکانیسم گم شدن داده‌ها قابل چشم پوشی است (لیتل^۱ و رابین، ۲۰۰۲؛ رابین ۱۹۸۷).

۱-۲-۲ تابع درست‌نمایی و توزیع پسین داده‌های مشاهده شده

تابع درست‌نمایی داده‌های مشاهده شده

بنا بر رابین (۱۹۷۶) و لیتل و رابین (۲۰۰۲)، می‌توان نشان داد که تحت فرض گم شدگی قابل چشم پوشی، هنگامی که براساس تابع درست‌نمایی و یا روش بیزی در مورد θ استنباط می‌کنیم، نیازی به در نظر گرفتن مدل برای R_i و در نظر گرفتن پارامترهای ξ نداریم.

چون داده‌های مشاهده شده شامل Y_{obs} و R_i هستند، توزیع احتمالاتی داده‌های مشاهده شده به صورت

^۱ Little

زیر می باشد

$$\begin{aligned} P(R, Y_{obs} | \theta, \xi) &= \int P(R, Y | \theta, \xi) dY_{mis} \\ &= \int P(R | Y, \xi) P(Y | \theta) dY_{mis}, \end{aligned} \quad (7.1)$$

که برای توزیع های گسسته این انتگرال چندگانه به جمع چندگانه تبدیل می شود. تحت فرض گم شدگی تصادفی، رابطه ی فوق به رابطه ی زیر تبدیل می شود

$$\begin{aligned} P(R, Y_{obs} | \theta, \xi) &= P(R | Y_{obs}, \xi) \int P(Y | \theta) dY_{mis} \\ &= P(R | Y_{obs}, \xi) P(Y_{obs} | \theta). \end{aligned} \quad (8.1)$$

بنابراین می توان تحت فرض گم شدگی تصادفی، تابع درستنمایی داده های مشاهده شده را به دو جزء افراز کرد که یک جزء مربوط به پارامتر مورد علاقه θ و جزء دیگر مربوط به پارامتر مزاحم ξ می باشد. هنگامی که دو پارامتر مجزا هستند، استنباط بر اساس درستنمایی در مورد θ تحت تأثیر ξ یا $P(R | Y_{obs}, \xi)$ قرار نمی گیرد. بنابراین برآورد ماکسیمم درستنمایی θ ، آزمون های نسبت درستنمایی مربوط به θ و غیره را می توان بدون در نظر گرفتن مکانیسم گم شدن داده ها انجام داد؛ یعنی، از مکانیسم گم شدگی داده ها چشم پوشی می کنیم.

عامل مربوط به θ در رابطه ی (۸.۱) (یا به طور دقیق تر، هر تابع متناسب با این عامل) را همان طور که در لیتل و رابین (۲۰۰۲) ذکر شده، می توان به عنوان درستنمایی که از مکانیسم گم شدن داده ها چشم پوشی می کند در نظر گرفت، یعنی

$$L(\theta | Y_{obs}) \propto P(Y_{obs} | \theta) \quad (9.1)$$

برای اختصار، رابطه ی (۹.۱) را به عنوان درستنمایی داده های مشاهده شده در نظر می گیریم. به طور کلی چون فرض قابلیت چشم پوشی را در نظر می گیریم، نیاز به کار کردن با تابع کامل (۸.۱) نداریم. توجه کنید که در نگاه اول، به نظر می رسد تجزیه به عوامل (۸.۱) شامل فرض های ذکر شده در مورد مکانیسم گم شدگی نیست. توزیع توأم دو متغیر تصادفی مانند، Z_1 و Z_2 ، را همواره می توان به صورت حاصل ضرب توزیع حاشیه ای Z_1 و توزیع شرطی Z_2 به شرط Z_1 نوشت. یک تفاوت مهم بین این