

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشکده

رشته آمار ریاضی ، گرایش آمار ریاضی

پایان نامه کارشناسی ارشد

مدل‌های اثرات تصادفی فضایی و غیرفضایی بیزی برای برآورد نواحی کوچک

نگارنده: مطهره یوسفی

استادان راهنما

دکتر محمدرضا ربیعی
دکتر حسین باغیشنی

شهریور ۱۳۹۷

شماره:

تاریخ:

باسمه تعالی



مدیریت تحصیلات تکمیلی

فرم شماره (۳) صورتجلسه نهایی دفاع از پایان نامه دوره کارشناسی ارشد

با نام و یاد خداوند متعال، ارزیابی جلسه دفاع از پایان نامه کارشناسی ارشد خانم مطهره یوسفی با شماره دانشجویی ۹۴۱۹۳۰۴ رشته آمار گرایش آمارریاضی تحت عنوان مدل‌های اثرات تصادفی فضایی و غیرفضایی بیزی برای برآورد نواحی کوچک که در تاریخ ۹۷/۶/۱۳ با حضور هیأت محترم داوران در دانشگاه صنعتی شاهرود برگزار گردید به شرح ذیل اعلام می‌گردد:

قبول (با درجه: ...ک...ک...ک) مردود

نوع تحقیق: نظری عملی

عضو هیأت داوران	نام و نام خانوادگی	مرتبه علمی	امضاء
۱- استاد راهنمای اول	دکتر محمدرضا ربیعی	استادیار	
۲- استاد راهنمای دوم	دکتر حسین باغیشتی	استادیار	
۳- استاد مشاور			
۴- نماینده تحصیلات تکمیلی	دکتر احمد معتمدنژاد	دانشیار	
۵- استاد ممتحن اول	دکتر محمد آرشی	دانشیار	
۶- استاد ممتحن دوم	دکتر نگار اقبال	استادیار	

نام و نام خانوادگی رئیس دانشکده: دکتر ابراهیم هاشمی
تاریخ و امضاء و مهر دانشکده:

تبصره: در صورتی که کسی مردود شود حداکثر یکبار دیگر (در مدت مجاز تحصیل) می‌تواند از پایان نامه خود دفاع نماید (دفاع

مجدد نباید زودتر از ۴ ماه برگزار شود).

الهی

جان ما را صفای خودده و دل ما را هوای خودده، و چشم ما را ضیای خودده، و ما را از فضل

و کرم خود آن ده که آن به.

یار ب دل ما را توبه رحمت جان ده در دهمه راه صابری درمان ده

این بنده چه داند که چه می باید جست داننده تویی هر آنچه دانی آن ده

تقدیم به پدر و مادرم عزیزم

که از نگاهشان صلابت

از رفتارشان محبت

و از صبرشان ایستادگی را آموختم

سپاس‌گزاری

سپاس‌ خدای‌ راکه‌ سخ‌ن‌وران‌، در‌ ست‌ودن‌ او‌ ب‌ماند‌ و‌ شمار‌ندگان‌، ش‌م‌ردن‌ نعمت‌ های‌ او‌ ندانند‌ و‌ کوشندگان‌، حق‌ او‌ را‌ کزاردن‌ نتوانند‌. سلام‌ و‌ د‌ود‌ر‌ ب‌ر‌م‌حمد‌ و‌ خاندان‌ پاک‌ او‌، طاهران‌ معصوم‌، هم‌ آنان‌ که‌ وجودمان‌ و‌ امدار‌ وجودشان‌ است‌.

بدون‌ شک‌ جایگاه‌ و‌ منزلت‌ معلم‌، اجل‌ از‌ آن‌ است‌ که‌ در‌ مقام‌ قدردانی‌ از‌ زحمات‌ بی‌شائبه‌ی‌ او‌، بازبان‌ قاصر‌ و‌ دست‌ ناتوان‌، چیزی‌ ب‌ن‌گاریم‌. اما‌ از‌ آنجایی‌ که‌ تجلیل‌ از‌ معلم‌، سپاس‌ از‌ انسانی‌ است‌ که‌ هدف‌ و‌ غایت‌ آفرینش‌ را‌ تا‌ین‌ می‌کند‌ و‌ سلامت‌ امانت‌ بانی‌ راکه‌ به‌ دستش‌ سپرده‌اند‌، تضمین‌؛ بر‌ حسب‌ و‌ طیفه‌ و‌ از‌ باب‌ «من‌ لم‌ یشکر‌ المنعم‌ من‌ المخلوقین‌ لم‌ یشکر‌ الله‌ عز‌وجل‌»،؛ از‌ پ‌ر‌ و‌ ما‌د‌ عزیزم‌ این‌ دو‌ معلم‌ بزرگوارم‌ که‌ همواره‌ بر‌ کوتاهی‌ و‌ درشتی‌ من‌، قلم‌ عفو‌کننده‌ و‌ کرم‌اند‌ از‌ کنار‌ غفلت‌ بایم‌ گذشته‌اند‌ و‌ در‌ تمام‌ عرصه‌ های‌ زندگی‌ یار‌ و‌ یاور‌ بی‌ چشم‌ داشت‌ برای‌ من‌ بوده‌اند‌؛ از‌ اساتید‌ با‌ کمالات‌ و‌ شایسته‌؛ جناب‌ آقای‌ د‌ک‌تر‌ ربیع‌ی‌ و‌ د‌ک‌تر‌ باغ‌شینی‌ که‌ در‌ کمال‌ سع‌د‌ر‌، با‌ حسن‌ خلق‌ و‌ فروتنی‌، از‌ بیچ‌ لگی‌ در‌ این‌ عرصه‌ بر‌ من‌ دین‌ نمودند‌ و‌ زحمت‌ راه‌نمایی‌ این‌ رساله‌ را‌ بر‌ عهده‌ گرفتند‌؛ ش‌م‌کر‌ می‌کنم‌. باشد‌ که‌ این‌ خردترین‌، بخشی‌ از‌ زحمات‌ آنان‌ را‌ سپاس‌ گوید‌. از‌ اساتدان‌ دلسوز‌ گروه‌ آمار‌ و‌ ان‌ش‌گاه‌ صنعتی‌ شاهرود‌؛ سرکار‌ خانم‌ د‌ک‌تر‌ اقبال‌ و‌ جناب‌ آقای‌ د‌ک‌تر‌ آرشی‌ که‌ زحمت‌ د‌اور‌ی‌ این‌ پایان‌نامه‌ را‌ متقبل‌ شدند‌؛ کمال‌ ش‌م‌کر‌ و‌ قدردانی‌ را‌ دارم‌.

تعهد نامه

این جانب **مطهره یوسفی** دانشجوی کارشناسی ارشد رشته **آمار ریاضی** دانشگاه صنعتی شاهرود، نویسنده پایان نامه با عنوان **مدل های اثرات تصادفی فضایی و غیرفضایی بیزی برای برآورد نواحی کوچک**، تحت راهنمایی دکتر **محمد رضا ربیعی** و دکتر **حسین باغیثنی** متعهد می شوم:

- تحقیقات در این پایان نامه توسط این جانب انجام شده است و از صحت و اصالت برخوردار است.
- در استفاده از نتایج پژوهش های دیگر پژوهش گران، به مرجع مورد استفاده استناد شده است.
- مطالب این پایان نامه، تا کنون توسط خود، یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارایه نشده است.
- حقوق معنوی این اثر، به دانشگاه صنعتی شاهرود تعلق دارد، و مقالات مستخرج با نام “دانشگاه صنعتی شاهرود” یا “Shahrood University of Technology” به چاپ خواهند رسید.
- حقوق معنوی تمام افرادی که در به دست آوردن نتایج اصلی پایان نامه تاثیرگذار بوده اند، در مقالات مستخرج از پایان نامه رعایت می شود.
- در تمام مراحل انجام این پایان نامه، در مواردی که از موجود زنده (یا بافت های آنها) استفاده شده است، ضوابط و اصول اخلاقی رعایت شده است.
- در تمام مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته (یا استفاده شده است)، اصل رازداری و اصول اخلاق انسانی رعایت شده اند.

مطهره یوسفی

شهریور ۱۳۹۷

مالکیت نتایج و حق نشر

- تمام حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه های رایانه ای، نرم افزارها و تجهیزات ساخته شده) متعلق به دانشگاه صنعتی شاهرود می باشد. این مطلب باید به نحو مقتضی، در تولیدات علمی مربوطه ذکر شود.
- استفاده از اطلاعات و نتایج موجود در این پایان نامه بدون ذکر منبع مجاز نیست.

چکیده

در سال‌های اخیر مساله برآورد کوچک ناحیه‌ای به دلیل نیاز به آمارهای قابل اعتماد، بسیار مورد توجه قرار گرفته است. مشکل عمده در این زمینه، ناممکن بودن اندازه‌گیری متغیر هدف برای تک تک افراد حاضر در نواحی مورد نظر است. حتی، نمونه‌گیری از همه نواحی تحت مطالعه می‌تواند هزینه‌های مالی و زمانی سنگینی در پی داشته باشد. مدل‌های آماری متعددی برای نیل به این هدف پیشنهاد شده‌اند. هدف اصلی این مدل‌ها، به‌کارگیری اطلاعات کمکی برای ارتقا برآوردهای مستقیم است. به‌همین منظور، اهمیت استفاده از اطلاعات مکانی نواحی در چارچوب مدل‌های فضایی، نقش موثری در تحلیل نواحی کوچک دارند. مشابهت فضایی بین نواحی همجوار از جمله اطلاعات مفید پرکاربرد است که مدل‌های کوچک ناحیه‌ای فضایی برای استفاده از این اطلاعات، پیشنهاد و توسعه یافته‌اند. با توجه به این که در بسیاری از تحلیل‌های نواحی کوچک با پاسخ‌های غیر نرمال مانند پاسخ‌های شمارشی مواجه هستیم، در این پایان‌نامه رده مدل‌های آمیخته خطی تعمیم‌یافته فضایی را در چارچوب استنباط بیزی، برای تحلیل داده‌های نواحی کوچک مورد بررسی قرار می‌دهیم. با تحلیل مجموعه داده‌های بیمه استان گیلان، عملکرد این مدل‌ها را ارزیابی می‌کنیم.

کلمات کلیدی: برآورد کوچک ناحیه‌ای، مدل‌های فضایی، مدل‌های آمیخته خطی تعمیم‌یافته فضایی، استنباط بیزی.

پیش‌گفتار

روش‌های برآورد کوچک ناحیه‌ای به منظور جبران نقص مربوط به خطاهایی که از برآوردهای حاصل از داده‌های نمونه‌ای با تعدادی محدود به وجود آمده است، به کار گرفته می‌شود (رائو، ۲۰۰۳). در این روش سعی بر این است، تا با به‌کارگیری از مدل‌های آماری مناسب برای استفاده از سایر اطلاعات موجود در هر ناحیه که با متغیر تحت بررسی مرتبط هستند، برای بهبود این نوع از برآوردها استفاده شود. استفاده روزافزون از این برآوردها باعث شده است که قابل اطمینان بودن آن‌ها حائز اهمیت قرار گیرد. از آن‌جا که اندازه نمونه در این نواحی کوچک یا بسیار کوچک است، برآوردهای مستقیم با خطای استاندارد و ضریب تغییرات بزرگ همراه هستند. بنابراین لزوم استفاده از برآوردهای غیرمستقیم در ناحیه کوچک برای دسترسی به آماره‌های قابل اطمینان و دقیق اهمیت پیدا می‌کند. در عمل نواحی جغرافیایی، معمولاً، از یکدیگر مستقل نیستند و ویژگی‌های مورد مطالعه در هر ناحیه مشابه نواحی همسایه هستند. به‌کارگیری این وابستگی بین نواحی، که به آن وابستگی فضایی گویند، می‌تواند به‌عنوان اطلاعات کمکی در برآورد نواحی کوچک موثر باشد و به برآوردهای دقیق‌تری منجر شود.

هدف اصلی در این پایان‌نامه به‌کارگیری دیدگاه بیزی در برآورد کوچک ناحیه‌ای است. آرورا و لاهیری (۱۹۹۷) طبق قضیه‌ای نشان دادند که برآوردهای بیزی برابر بهترین پیشگوی نارایب خطی پارامترها هستند و آن‌ها نشان دادند که مزیت روش‌های بیزی، در کوچک بودن میانگین توان دوم خطاها است. برای تعیین برآوردهای بیزی، با در نظر گرفتن پیشین مناسب و استفاده از نمونه‌گیری گیبز، توزیع پسین پارامترها به‌دست می‌آید و پارامتر مورد علاقه به روش بیزی برآورد می‌شود. به‌عبارت دیگر با به‌دست آوردن توزیع پسین مورد نظر، برآوردها بهینه به‌دست می‌آید. همچنین برای بیان اثرات تصادفی فضایی حاصل از مجاورت نواحی، به علت بهینگی، معمولاً، از مدل اتورگرسیو شرطی استفاده می‌شود.

با این مقدمه، تمرکز این پایان‌نامه به استنباط بیزی در مدل‌های کوچک ناحیه‌ای (فضایی) و کاربرد آن برای پیش‌گویی برآورد کوچک ناحیه‌ای است. در فصل اول، برخی از مفاهیم پایه‌ای مرتبط با برآوردهای نواحی کوچک معرفی می‌شود. در فصل دوم، مدل‌های کوچک ناحیه‌ای را معرفی می‌کنیم. در فصل سوم، مدل‌بندی بیزی را بر روی نواحی کوچک بیان خواهیم کرد و در فصل چهارم نیز کاربرد مدل فضایی بیزی را برای برآورد کوچک ناحیه‌ای داده‌های بیمه استان گیلان نشان خواهیم داد. همچنین با ارائه پیشنهاداتی و یک نتیجه‌گیری کلی، ادامه روند پژوهش را در آینده تحقیق، مطرح می‌نماییم.

لیست مقالات مستخرج از پایان نامه

۱. یوسفی، م.، ربیعی، م.ر. و باغیثنی، ح. (۱۳۹۷). "استنباط بیزی فضایی نواحی کوچک برای بیمه‌شدگان اجباری استان گیلان"، مجموعه مقالات چهاردهمین کنفرانس آمار ایران، شاهرود (۶۸۴).

فهرست مطالب

ش	فهرست تصاویر
ث	فهرست جداول
۱	۱ مفاهیم پایه برآورد کوچک ناحیه‌ای
۱	۱.۱ مقدمه
۴	۱.۲ علل تقاضا برای برآورد کوچک ناحیه‌ای
۴	۱.۳ آشنایی با مفاهیم لازم در برآورد کوچک ناحیه‌ای
۷	۱.۳.۱ برآورد کوچک ناحیه‌ای
۸	۱.۴ برآوردهای غیرمستقیم سنتی
۸	۱.۵ روش برآورد حوزه‌ای مستقیم
۹	۱.۵.۱ رویکرد طرح پایه
۱۲	۱.۶ روش‌های برآورد حوزه‌ای غیرمستقیم
۱۳	۱.۷ دیدگاه‌های استنباطی در برآورد کوچک ناحیه‌ای
۱۳	۱.۸ روش‌های مبتنی بر نمونه‌گیری
۱۳	۱.۸.۱ مونت کارلوی زنجیره مارکوفی MCMC
۱۸	۱.۸.۲ نمونه‌گیر گیبز
۱۹	۱.۸.۳ الگوریتم متروپلیس-هستینگز درون‌گیری
۲۱	۱.۹ مقدمه‌ای بر آمار فضایی
۲۳	۱.۹.۱ مدل بندی ساختار وابستگی فضایی
۲۵	۱.۹.۲ وارد کردن ساختار همبستگی فضایی در مدل
۲۵	۱.۹.۳ ماتریس همسایگی
۲۹	۲ معرفی مدل‌های کوچک ناحیه‌ای
۲۹	۲.۱ برآوردهای کوچک ناحیه‌ای
۳۱	۲.۲ معرفی مدل‌های کوچک ناحیه‌ای

۳۳	مدل پایه‌ای در سطح ناحیه	۲.۳
۳۵	گسترش‌های مدل در سطح ناحیه (A)	۲.۳.۱
۴۲	مدل پایه‌ای در سطح واحد آماری	۲.۴
۴۵	گسترش‌های مدل در سطح واحد آماری (B)	۲.۴.۱
۴۵	مدل رانو-یو	۲.۵
۴۶	مدل‌های آمیخته‌ی خطی کلی	۲.۶
۴۶	مقدمه	۲.۶.۱
۴۷	مدل آمیخته‌ی خطی LMM	۲.۶.۲
۴۸	مدل آمیخته‌ی خطی تعمیم‌یافته GLMM	۲.۶.۳
۵۰	مدل‌های آمیخته‌ی خطی تعمیم‌یافته فضایی SGLMM	۲.۶.۴
۵۲	برآوردگرهای غیرمستقیم و روش‌های مبتنی بر مدل	۲.۷
۵۲	برآوردگر ترکیبی	۲.۷.۱
۵۳	برآوردگرهای مبتنی بر مدل صریح	۲.۸
۵۴	برآوردگرهای پیشگوی ناریب خطی تجربی BLUP	۲.۸.۱
۵۵	برآوردگرهای بهترین پیشگوی ناریب خطی تجربی EBLUP	۲.۸.۲
۵۷	برآوردگر بیز تجربی EB	۲.۸.۳
۵۸	برآوردگر بیز سلسله‌مراتبی HB	۲.۸.۴
۶۱	مدل بندی بیزی برآورد کوچک ناحیه‌ای	۳
۶۱	روش‌های بیزی	۳.۱
۶۳	ساختار سلسله‌مراتبی مدل بندی بیزی برای نواحی کوچک	۳.۲
۶۳	بررسی رویکرد بیزی	۳.۳
۶۵	مدل سلسله‌مراتبی بیزی فضایی	۳.۴
۶۶	داده‌های شبکه‌ای	۳.۴.۱
۶۶	داده‌های زمین آماری	۳.۴.۲
۶۸	توزیع پیشین پارامترهای مدل	۳.۴.۳
۶۹	ارزیابی کیفیت برآوردگرها	۳.۵
۷۱	برآورد کوچک ناحیه‌ای در صورت عدم اطلاع از برآوردگرهای مستقیم	۳.۶
۷۱	همبستگی اثرات تصادفی فضایی	۳.۶.۱
۷۲	مشخص‌سازی نرمال چندمتغیره	۳.۶.۲
۷۲	مشخص‌سازی اتورگرسیو شرطی	۳.۶.۳
۷۳	وام گرفتن اطلاعات در سطوح اجرایی بالاتر	۳.۶.۴

۷۵	کاربرد مدل فضایی بیزی برای برآورد کوچک ناحیه‌ای داده‌های بیمه استان گیلان	۴
۷۵	مقدمه	۴.۱
۷۶	معرفی داده‌ها	۴.۲
۷۷	معرفی مدل بیزی	۴.۳
۷۸	برازش مدل و تحلیل نتایج	۴.۴
۷۹	نتیجه‌گیری و آینده تحقیق	۴.۵
۸۱	مراجع	
۹۱	توزیع‌های شرطی کامل برای مشخص‌سازی CAR با مشاهدات گم‌شده	آ
۹۳	دستورات نرم‌افزار R و OPENBUGS	ب
۱۰۱	واژه‌نامه فارسی به انگلیسی	
۱۰۳	واژه‌نامه انگلیسی به فارسی	

فهرست تصاویر

۷ نقشه حوزه‌های استان گیلان	۱.۱
۲۱ نقشه محل سنگ‌های فیروزه در کشور هندوستان	۲.۱
۲۲ نمودار محل افراد بیمه شده در استان گیلان	۳.۱
۲۲ نمودار نقاط زلزله در سطح کشور	۴.۱
 انواع همبستگی فضایی، (شکل سمت راست) همبستگی فضایی مثبت، (شکل سمت چپ) همبستگی فضایی منفی	۵.۱
۲۵ همسایگی مرتبه ۱ تا ۳ برای ناحیه‌ی d	۶.۱
۲۶ نواحی مورد بررسی	۷.۱
۲۷ نمودار انواع برآوردهای کوچک ناحیه‌ای	۱.۲
۳۰ نقشه همسایگی شهرستان‌های استان گیلان	۱.۴
۷۷ نمودارهای اثر پارامترهای مدل بیزی داده‌های بیمه	۲.۴
۷۸		

فهرست جداول

۱.۴ نتایج برازش مدل بیزی بر روی داده‌های بیمه اجباری استان گیلان . . . ۷۹

فصل ۱

مفاهیم پایه برآورد کوچک ناحیه‌ای

با توجه به محدود بودن منابع، استفاده از داده‌ها در سطوح پایین جغرافیایی الزامی است، که در این صورت برآوردگرهای مستقیم^۱ با خطای استاندارد و ضریب تغییرات بزرگ همراه می‌شوند. لذا برای بالا بردن دقت برآوردگرها و کاهش واریانس آن‌ها از برآوردگرهای غیرمستقیم^۲ کوچک ناحیه‌ای استفاده می‌شود. بنابراین در برآورد نواحی کوچک^۳ روش‌هایی پیشنهاد شده‌اند. که هر کدام از آن‌ها حائز اهمیت هستند. در این فصل به تعاریف و مقدمات لازم برای ورود به فصل‌های بعدی خواهیم پرداخت. عمده از مطالب استفاده شده در این فصل برگرفته از کتاب برآورد کوچک ناحیه‌ای راثو (۲۰۰۳) و کتاب آمار فضایی محمدزاده (۱۳۹۴) است.

۱.۱ مقدمه

دانشواژه «کوچک ناحیه^۴» بیانگر هر زیر جامعه‌ای است که برای آن نتوان برآوردهای مستقیم را با دقت کافی تولید کرد (راثو، ۲۰۰۳). دیر زمانی است که تقاضا برای برآوردهای کوچک ناحیه‌ای به طور قابل توجهی افزایش یافته است. این توجه به دلیل افزایش تقاضای برآوردهای معتبر برای نواحی کوچک است. زیرا تصمیم‌های مربوط به کسب و کار، به‌خصوص تصمیم‌های

^۱ Direct estimators

^۲ Indirect estimators

^۳ Small area estimation

^۴ Small area

مربوط به بخش‌های خصوصی و دولتی که مرتبط با بنگاه‌های کوچک‌اند، به شدت متکی بر شرایط محلی هستند. در حال حاضر روش برآورد کوچک ناحیه‌ای برای بررسی مسائل اقتصادی آن‌هم در کشورهای اروپایی و کشورهای اتحاد شوروی اهمیت بسزایی دارد، زیرا این کشورها در حال دورشدن از تصمیم‌گیری‌های متمرکز می‌باشند.

عموما بررسی‌ها به گونه‌ای طراحی شده‌اند که برآورد پارامترهای مورد نظر، در سطوح ملی یا حوزه‌های کشوری حاصل می‌شوند. بنابراین این برآوردها در سطوح پایین‌تر نیز لازم هستند. از طرفی اندازه‌ی نمونه‌ها در این سطوح کوچک‌اند، پس، نمی‌توان این برآوردها را به‌طور قابل قبولی با داده‌های بررسی به‌دست آورد. اغلب داده‌های ثبتي در دسترس‌اند، اما تنها اطلاعات بسیار کمی را در رابطه با متغیرهای مورد نظر خواهند داشت.

در خصوص برآورد کوچک ناحیه‌ای، راثو (۲۰۰۳)، جیانگ و همکاران (۲۰۰۶)، پففرمن (۲۰۱۳) و راثو و مولینا (۲۰۱۵) تاریخچه‌ای را از شکل‌گیری آن را شرح داده‌اند و روش‌ها و مطالعات موجود در این زمینه را دسته‌بندی کرده‌اند (آرورا و لاهیری، ۱۹۹۷).

پرسش کلیدی که در برآورد نواحی کوچک مطرح خواهد شد این است که، وقتی اندازه‌ی نمونه شامل مشاهدات بسیار اندک و اطلاعات کم باشد، چگونه می‌توان برآوردهای بسیار معتبر و دقیقی را به‌دست آورد؟ برای پاسخ به این پرسش، دانستن این موضوع کافی است که زمانی که اندازه‌ی نمونه‌ی به‌دست آمده از یک ناحیه، کوچک باشد (حتی ممکن است گاهی نیز صفر باشد) برآوردهای مستقیم از یک انحراف معیار بزرگ و غیرقابل قبول برخوردار می‌شوند (سالواتی، ۲۰۰۴). به همین منظور برآوردهای به‌دست آمده از این داده‌ها نامعتبراند. لذا فنون برآورد نواحی کوچک را برای تحلیل و بررسی این مشکلات به‌کار می‌برند.

مشکل عمده‌ای که در برآورد نواحی کوچک وجود دارد، ناممکن بودن اندازه‌گیری متغیر هدف^۵ برای تک افراد حاضر در نواحی مورد نظر و استفاده کردن از یک نمونه است. حتی ممکن است، نمونه‌گیری از همه نواحی تحت مطالعه هزینه‌های مالی و نیز زمانی سنگینی را در پی داشته باشد (سارندال و همکاران، ۱۹۹۲).

کافی نبودن اندازه داده‌ها در آمارگیری از نواحی کوچک، موجب کم دقتی برآوردهای مستقیم در قسمت‌های مختلف نواحی می‌شود. لذا برای تولید آماره‌های معتبر و دقیق در نواحی کوچک، مطالعات زیادی انجام شده‌اند که با ارائه رهیافت‌های مناسب به دنبال حل این مشکلات هستند. یکی از این رهیافت‌ها را می‌توان به کارگیری منابع مختلف و اطلاعات کمکی در نظر گرفت تا با دادن قدرت قرضی به برآوردگر، دقت آن را افزایش دهند. به این نوع از برآوردگرها که باعث افزایش دادن دقت نواحی کوچک می‌شوند، برآوردگرهای غیرمستقیم می‌گویند. به‌طور کلی باید قدرت متغیر پاسخ‌مان را با استفاده از اطلاعات کمکی افزایش دهیم چرا که، موجب افزایش دقت برآوردگر مورد نظر می‌شود. اغلب از یک مدل رگرسیونی مناسب و برازش آن به نمونه منتخب و بهره‌گیری از اطلاعات کمکی اضافی، می‌توان نواحی خارج از نمونه‌گیری را برآورد کرد (بانرجی و همکاران، ۲۰۱۴).

^۵Target variable

در حال حاضر، با توجه به پیشرفت‌های محاسباتی و توان رایانه‌های امروزی، در دیدگاه کلاسیک و بیزی، به دست آوردن یک آماره معتبر در مدل‌های خطی کار پیچیده‌ای نیست. از این رو، در دیدگاه کلاسیک برای برازش مدل‌های رگرسیونی با اثرات تصادفی، در برآورد نواحی کوچک، معمولاً از روش بهترین پیشگوی ناریب خطی تجربی^۶ (EBLUP) استفاده می‌شود (رابینسون، ۱۹۹۱؛ جیانگ و لاهیری، ۲۰۰۶). اما هنگامی که متغیر هدف غیرنرمال باشد، برازش این مدل‌های مناسب و محاسبه EBLUP کار پیچیده‌ای خواهد شد. با توجه به مشکلات محاسباتی دیدگاه کلاسیک در مدل‌های نواحی کوچک، رهیافت مدل‌بندی بیزی آن‌ها، به دلیل وجود الگوریتم‌های نمونه‌گیری زنجیر مارکوف مونت کارلویی^۷ MCMC، طرفداران بیشتری پیدا کرده است (گوش و همکاران، ۱۹۹۴).

افزون بر این، دیدگاه بیزی مدل‌های نواحی کوچک مزایای ارزشمند دیگری نیز دارد؛ به‌طور جزئی‌تر، این دیدگاه یک چارچوب منسجم را فراهم می‌آورد که بر اساس آن می‌توان انواع متغیرهای هدف شامل پیوسته، دودویی، دسته‌ای و انواع ساختارهای وابستگی شامل مستقل، وابستگی فضایی^۸ و فضایی-زمانی^۹ نواحی فاقد اطلاعات مستقیم سرشماری شده را در نظر گرفت، تا با روش‌های محاسباتی یکسان، مدل را برازش داد.

با توجه به این که اشاره‌ای به وابستگی فضایی شده است، کافی است بدانید، نواحی جغرافیایی معمولاً از یکدیگر مستقل نیستند و ویژگی‌های مورد مطالعه در هر ناحیه مشابه نواحی همسایه هستند. به کارگیری این وابستگی بین نواحی، که به آن وابستگی فضایی گویند، می‌تواند به عنوان اطلاعات کمکی در برآوردهای نواحی کوچک موثر واقع شود تا برآوردهای دقیق و با کیفیت بهتری نتیجه شود. گوش و همکاران (۱۹۹۸)، راثو (۲۰۰۳)، پراتسی و سالواتی (۲۰۰۸) و جیانگ و لاهیری (۲۰۰۶) با استفاده از مدل‌های مختلف فضایی اثرات نواحی همسایه‌ها را در پیش‌گویی نواحی کوچک لحاظ کردند. یک ویژگی مطلوب اطلاعات مکانی آن است که همیشه و بدون تحمیل هزینه‌ای، در دسترس هستند. با توجه به این که استنباط آماری مدل‌ها در دیدگاه بیزی مبتنی بر توزیع پسین است و معمولاً، صورت تحلیلی خاصی ندارد، به کارگیری روش‌های MCMC پرکاربرد خواهد بود و برخلاف روش‌های کلاسیک که معمولاً بر گزاره‌های مجانبی متکی هستند، روش‌های بیزی با کمک الگوریتم‌های MCMC یک پاسخ کلی برای تقریب توزیع پسین فراهم می‌آورند. روش‌های برآوردیابی در رهیافت بیزی با وجود پیشرفت‌های محاسباتی از سختی محاسباتی کم‌تری برخوردار هستند. در این پایان‌نامه نواحی کوچک با استفاده از مدل‌های فضایی و غیرفضایی^{۱۰} و نیز به کارگیری از مدل‌های سلسه‌مراتبی بیزی^{۱۱} (HB) و الگوریتم‌های MCMC مد نظر قرار گرفته است. برخی از مقالات مروری مرتبط

^۶ Empirical best linear unbiased prediction

^۷ Markov chain monte carlo

^۸ Spatial dependence

^۹ Spatial-time

^{۱۰} Spatial and non-spatial models

^{۱۱} Bayesian hierarchical models

با برآورد کوچک ناحیه‌ای منتشر که کمک شایانی برای آشنایی بیشتر در این باره می‌کنند، عبارت‌اند از راثو (۱۹۸۶)، چاودری (۱۹۹۲)، گوش و راثو (۱۹۹۴)، مارکر (۱۹۹۹)، راثو (۲۰۰۱)، و پفرمن (۲۰۰۲). همچنین، کتاب‌هایی نیز درباره‌ی نظریه برآورد کوچک ناحیه‌ای منتشر شده‌اند، که می‌توان به راثو (۲۰۰۳) و موخوپادیای (۱۹۹۸) اشاره کرد (دیگل و همکاران، ۲۰۰۲).

۲.۱ علل تقاضا برای برآورد کوچک ناحیه‌ای

آمارهای کوچک ناحیه‌ای در انگلستان در سده‌ی یازدهم و در کانادا در سده‌ی هفدهم، به دو صورت مبتنی بر سرشماری یا اطلاعات ثبتی مورد استفاده قرار می‌گرفت (براکستون، ۱۹۸۷). در سال‌های اخیر آمارهای کوچک ناحیه‌ای به‌طور چشم‌گیری مورد توجه قرار گرفته‌اند. علت به‌کارگیری این آمارها نیز آن است، که در گذشته سیاست‌گذاران به آن دسته از برآوردهای تولیدشده‌ای که برای ناحیه‌های جغرافیایی بزرگ و وسیع بود، قانع بودند. اما در حال حاضر سیاست‌گذاران برای تصمیم‌گیری‌های خود نیازمند برآوردهایی در ناحیه‌های کوچکتر همانند استان و شهرستان هستند. به‌عنوان مثال، دولت اگر بخواهد برای رفع محرومیت در منطقه برنامه‌ریزی کند قبل از هر چیز باید برآوردهایی مانند نرخ بیکاری و تعداد افراد زیر خط فقر در آن منطقه را داشته باشد و صرفاً دانستن آمارها در سطح کشوری نمی‌تواند کافی باشد. به‌عنوان مثالی دیگر، وقتی یک بیماری در سطح کشور شیوع یافته است، برای کنترل سریع‌تر آن توجه به شرایط مناطق مختلف بسیار مهم است.

از دلایل دیگر نیاز به آمارهای کوچک ناحیه‌ای، می‌توان به فرمول‌بندی سیاست‌ها و برنامه‌ها در تخصیص منابع دولتی به کمک دولت، برنامه‌ریزی منطقه‌ای، قانون‌گذاری از طریق حکومت‌های ملی برای ناحیه‌های مختلف، تقاضا از جانب بخش خصوصی و تصمیم‌های بنگاه‌های کوچک اشاره کرد.

۳.۱ آشنایی با مفاهیم لازم در برآورد کوچک ناحیه‌ای

از دیرباز آمارگیری نمونه‌ای به‌عنوان وسیله‌ای کم هزینه برای دستیابی به برآوردهایی از کل جامعه یا زیربخش‌هایی از جامعه به کار می‌رفته است. در این راستا، دو نوع برآوردگر به نام‌های مستقیم و غیرمستقیم مورد توجه قرار گرفته‌اند. در آمارگیری نمونه‌ای، برآوردگری را مستقیم گویند، هرگاه بر داده‌های نمونه‌ای ناحیه مبتنی باشد. این نوع از برآوردهای مستقیم عموماً «طرح پایه^{۱۲}» هستند اما، با این حال می‌توانند از مدلی نشأت گرفته باشند و طبق همان مدل نیز توجیه شوند.

^{۱۲} Design based

برآوردگرهای طرح پایه از وزن‌های نمونه‌گیری استفاده می‌کنند، و استنباط‌هایی که انجام می‌دهند مبتنی بر توزیع‌های احتمالی هستند که توسط طرح نمونه‌گیری و ثابت فرض کردن مقادیر جامعه تولید می‌شوند. برآوردگرهای مستقیم «مدل یار^{۱۳}» که از مدل‌های «کاری^{۱۴}» استفاده می‌کنند نیز طرح پایه‌اند، و اهدافشان «استوار^{۱۵}» سازی استنباط‌ها در مقابل تعیین نادرست (بد مشخص‌سازی) مدل است. برآوردگرهای مدل یار مبنی بر نمونه‌ی به‌دست آمده از جامعه‌ای است که دارای مدل است که آن مدل از یک مدل کلی به‌دست آمده است. در نتیجه به آن «مدل یار» می‌گویند (برسلو، ۱۹۹۳).

در ادامه برخی از مفاهیم لازم در برآورد کوچک ناحیه‌ای را بیان می‌کنیم.

تعریف ۱.۳.۱. (برآوردگر غیرمستقیم). به برآوردگر، غیرمستقیمی گفته می‌شود که علاوه بر استفاده از داده‌های مشاهده شده در نمونه، از مقدارهای مورد مطالعه در ناحیه‌ها و زمان‌های مرتبط و داده‌های کمکی در سطح جامعه نیز استفاده می‌کند.

تعریف ۲.۳.۱. (ناحیه‌ی بزرگ^{۱۶}). به ناحیه‌ای اطلاق می‌شود که اندازه‌ی نمونه‌ی آن به اندازه‌ی کافی بزرگ باشد که برآوردگر مستقیم حاصل از آن از دقت کافی برخوردار باشد.

تعریف ۳.۳.۱. (ناحیه‌ی کوچک). به ناحیه‌ای اطلاق می‌شود که اندازه‌ی نمونه‌ی آن به اندازه‌ی کافی بزرگ نیست که برآوردگر مستقیم حاصل از آن از دقت کافی برخوردار باشد.

تعریف ۴.۳.۱. (برآوردگر مدل پایه^{۱۷}). برآوردگرهای غیرمستقیمی که مبتنی بر مدل‌های کوچک ناحیه‌ای هستند را برآوردگرهای مدل پایه می‌نامند (گلنن، ۲۰۰۶).

سه نوع برآوردگر «غیرمستقیم» معرفی شده‌اند که عبارت‌اند از «غیرمستقیم حوزه‌ای^{۱۸}»، «غیرمستقیم زمانی^{۱۹}» و «غیرمستقیم حوزه‌ای و زمانی^{۲۰}» (رائو، ۲۰۰۳).

تعریف ۵.۳.۱. (برآوردگر غیرمستقیم حوزه‌ای). برآوردگری است که از مقادیر پاسخ و سایر متغیرهای مربوط با حوزه‌ای دیگر، اما نه از دوره‌ی زمانی دیگر، بهره می‌گیرد.

تعریف ۶.۳.۱. (برآوردگر غیرمستقیم زمانی). برآوردگری است که از مقادیر پاسخ و سایر متغیرهای مربوط با دوره‌ی زمانی دیگر برای حوزه‌ی مورد نظر اما نه از حوزه‌ای دیگر، استفاده می‌کند.

^{۱۳} Model assisted

^{۱۴} Working

^{۱۵} Robust

^{۱۶} Large area

^{۱۷} Model based estimator

^{۱۸} Indirect domain

^{۱۹} Indirect time

^{۲۰} Indirect domain and time

تعریف ۷.۳.۱. (برآوردگر غیرمستقیم حوزه‌ای و زمانی). برآوردگری است که از مقادیر پاسخ و سایر متغیرهای مربوط با حوزه‌ای دیگر و نیز دوره‌ی زمانی دیگر بهره می‌گیرد (ابری، ۲۰۰۰).

بقیه مدل‌ها را می‌توان به دو دسته مدل‌های پیونددهنده‌ی ضمنی^{۲۱} و صریح^{۲۲} تقسیم نمود که در بخش ۱.۲ به‌طور کامل توضیح داده می‌شوند. همچنین دو نوع از برآوردگرهای غیرمستقیم سنتی^{۲۳} وجود دارند که از مدل‌های ضمنی استفاده می‌کنند. این برآوردگرها عموماً طرح پایه‌اند و واریانس‌های طرحی آن‌ها (یعنی، واریانس‌های متناظر با توزیع احتمال القا شده توسط طرح نمونه‌گیری) نسبت به واریانس‌های طرحی برآوردگرهای مستقیم کوچک می‌باشند. با این تفاوت که برآوردگرهای غیرمستقیم عموماً طرح اریب^{۲۴} اند و وقتی اندازه‌ی نمونه کلی افزایش یابد، مقدار طرح اریبی کاهش نخواهد یافت. اگر فرض بر این باشد که مدل پیونددهنده‌ی ضمنی به‌طور تقریبی درست باشد، آن‌گاه مقدار طرح اریبی کوچک می‌شود، و به مقداری از میانگین توان دوم خطای طرحی منجر می‌شود که به‌طور معنی‌داری در قیاس با میانگین توان دوم خطای برآوردگر مستقیم کوچک‌تر است. کاهش در میانگین تنها دلیل عمده برای بهره‌گیری از برآوردگرهای غیرمستقیم در میانگین است (رائو، ۲۰۰۳). برآوردگرهای غیرمستقیم مبتنی بر مدل‌های کوچک ناحیه‌ای را «برآوردگرهای مدل پایه» می‌نامند.

استنباط‌های ناشی از برآوردگرهای مدل پایه به توزیعی اشاره دارند که از مدل مفروض نتیجه می‌شود. بنابراین، مدل‌گزینی و اعتبارسنجی نقشی حیاتی در روش برآورد مدل پایه ایفا می‌کنند. اگر مدل‌های مفروض برازش خوبی را به داده‌ها فراهم نسازند، برآوردگرهای مدل پایه، مدل اریب خواهند بود، که به نوبه خود می‌توانند به استنباط‌های ناکارآمد منجر شوند. به دلیل برتری‌های استفاده از مدل‌های صریح کوچک ناحیه‌ای به‌طور کلی هرگاه بخواهیم برآوردگرهای غیرمستقیم را به‌دست آوریم، بایستی بر اساس مدل‌های صریح کوچک ناحیه‌ای باشند؛ زیرا

۱. تحت مدل مفروض، برآوردگرهای «بهینه^{۲۵}» را می‌توان به‌دست آورد.
۲. برخلاف معیارهای فراموضعی (که نسبت به ناحیه‌های کوچک گرفته شده‌اند) و اغلب با برآوردگرهای غیرمستقیم سنتی به‌کار می‌روند، معیارهای تغییرپذیری مختص ناحیه را می‌توان با هر برآوردگر همراه ساخت.
۳. اعتبار مدل‌ها را می‌توان از روی داده‌های نمونه‌ای سنجید.
۴. بسته به طبیعت متغیرهای پاسخ و پیچیدگی ساختارهای داده‌ها (مانند ساختارهای فضایی و سری زمانی)، انواع گوناگون مدل‌ها را می‌توان آزمود.

^{۲۱} Implicit

^{۲۲} explicit connective

^{۲۳} Indirect traditional estimators

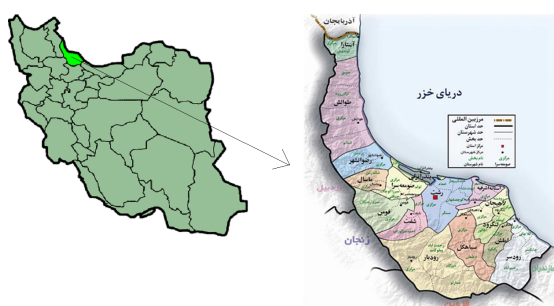
^{۲۴} Design bias

^{۲۵} Optimal

۱.۳.۱ برآورد کوچک ناحیه‌ای

اغلب آمارگیری‌های نمونه‌ای برای فراهم آوردن برآوردهای طرح پایه‌اند که از لحاظ آماری قابل اعتماد باشند، و همچنین از لحاظ جغرافیایی در سطح‌های ملی، استانی و ناحیه‌ای طراحی می‌شوند. به‌طور کلی این آمارگیری‌ها برای فراهم آوردن برآوردهایی برای انواع زیرجامعه‌ها و حوزه‌ها به‌کار برده می‌شوند.

برای بیان کردن حوزه‌ها همانند شکل ۱.۱ می‌توان این گونه تعریف کرد که، حوزه‌ها شامل ناحیه‌های جغرافیایی همانند ایالت یا استان، شهرستان، دهستان، یا ناحیه‌های شهرداری می‌باشند.



شکل ۱.۱: نقشه حوزه‌های استان گیلان

برای غلبه بر مشکل اندازه کم نمونه، روش‌های برآورد کوچک ناحیه‌ای مورد استفاده قرار می‌گیرد چرا که برآوردهایی تولید می‌کنند که بهتر از برآورد مستقیمی است که از نمونه در هر ناحیه کوچک به‌دست می‌آید. در حالت کلی به عنوان تعریفی از روش‌های برآورد کوچک ناحیه‌ای می‌توان گفت: ”روش‌های ایجاد برآوردهای قابل اعتماد و مناسب برای نواحی هستند، که در آن‌ها استفاده از روش‌های برآورد مستقیم، دقت مورد نیاز را فراهم نمی‌کند.“

استفاده از برآورد مستقیم، به معنی به‌کارگیری از روش‌های برآورد نمونه‌ای طرح پایه‌ی کلاسیک است و از واحدهای نمونه‌ی موجود در همان ناحیه استفاده می‌کند. بنابراین اگر نمونه ناحیه به اندازه‌ی کافی بزرگ نباشد که بتواند برآورد مستقیم و قابل اعتماد و همچنین با دقت کافی تولید کند آن ناحیه به عنوان ناحیه‌ی کوچک تلقی نمی‌شود. به این منظور، بسیاری از حوزه‌های مورد نظر (همانند شهرستان) ممکن است حتی دارای اندازه‌ی نمونه‌ی صفر باشند. به این دلیل از روش‌های برآورد کوچک ناحیه‌ای می‌توان برای تهیه‌ی آماره‌های قابل اعتماد برای ناحیه‌های کم نمونه یا بدون نمونه استفاده کرد.

تفاوتی که بین برآوردهای غیرمستقیم و مستقیم وجود دارد آن است که، برآوردهای غیرمستقیم برخلاف برآوردهای مستقیم تنها از واحدهای نمونه یا مقدار متغیر مورد نظر همان ناحیه استفاده نمی‌کنند، بلکه این برآوردها برای فراهم آوردن برآوردهایی با سطح دقت کافی از مقدارهای متغیر مورد نظر ناحیه‌ها یا دوره‌های زمانی مرتبط دیگر نیز استفاده

می‌کنند. لذا در این برآوردها از اطلاعات کمکی متغیر مورد نظر مانند اطلاعات آخرین سرشماری و اطلاعات ثبتي جاری نیز استفاده می‌شود. داده‌های متغیرهای مورد نظر ناحیه‌ها و دوره‌های زمانی مرتبط و اطلاعات کمکی از طریق یک مدل مناسب به فرآیند برآورد انتقال داده می‌شود. بنابراین در حالت کلی می‌توان گفت در دسترس بودن داده‌های کمکی خوب و تعیین مدل‌های پیونددهنده‌ی مناسب برای ساختن برآوردهای غیر مستقیم عامل‌های اساسی هستند. برای اطلاعات بیشتر راجع به برآوردهای کوچک ناحیه‌ای به پلاتنگ و سینگ (۱۹۸۶)، پلاتنگ و همکاران (۱۹۸۷)، چاودری (۱۹۹۴)، مارکر (۱۹۹۹)، راثو (۲۰۰۱)، پفرمن (۲۰۰۲) و راثو (۲۰۰۳) مراجعه کنید.

۴.۱ برآوردهای غیرمستقیم سنتی

این دسته از برآوردها، مبنی بر مدل‌های پیونددهنده‌ی ضمنی، به دو صورت برآوردهای هم‌گذاشتی^{۲۶} و ترکیبی^{۲۷} دسته‌بندی شده‌اند. حالت عمومی این برآوردها طرح پایه است و واریانس طرحی آن‌ها معمولاً نسبت به واریانس‌های طرحی برآوردهای مستقیم کوچک‌تر است. همچنین، برآوردهای غیرمستقیم سنتی در حالت کلی طرح اریب هستند و هنگامی که اندازه‌ی نمونه افزایش یابد مقدار طرح اریبی آن‌ها افزایش می‌یابد. بنابراین اگر مدل پیونددهنده‌ی ضمنی تقریباً درست باشد، مقدار طرح اریبی کوچک خواهد بود و این عمل منجر به مقداری از میانگین توان دوم خطای^{۲۸} (MSE) طرحی خواهد شد که در قیاس با MSE برآوردهای مستقیم کوچک‌تر است. علت اصلی به‌کارگیری برآوردهای غیرمستقیم، کاهش در MSE است.

۵.۱ روش برآورد حوزه‌ای مستقیم

داده‌های آمارگیری نمونه‌ای که از طرح‌های نمونه‌گیری مطرح در روش‌های نمونه‌گیری به‌دست می‌آیند، برای تولید برآوردهای مستقیم و قابل اعتماد از مجموع‌ها و میانگین‌های مرتبط به کل جامعه و ناحیه‌ها یا حوزه‌های بزرگ استفاده می‌کنند. برآوردهای مستقیم هر ناحیه، تنها مقدارهای متغیر مورد نظر مربوط به واحدهایی را که در آن ناحیه واقع‌اند، به کار می‌گیرد. برای فراگیری بیشتر از روش برآورد مستقیم می‌توان به کتاب‌های کوکران (۱۹۷۷)، هدایت و سینها (۱۹۹۱)، سارندال، سوئسون و رتمن (۱۹۹۲)، تامپسون (۱۹۹۷) و لهر (۱۹۹۹) مراجعه کنید.

^{۲۶} Synthetic

^{۲۷} combinational estimators

^{۲۸} Mean squared error

از طرفی دیگر، روش‌های مدل پایه برای تولید برآوردگرهای مستقیم و استنباط‌های مربوط به آن به کار می‌روند، و استنباط در این روش‌ها، بر اساس نمونه‌ی خاص استخراج‌شده صرف نظر از طرح نمونه‌گیری به دست می‌آیند (بروئر ۱۹۶۳، رویال ۱۹۷۰؛ والیانت و همکاران، ۲۰۰۱). اما متاسفانه، راهبردهای مدل پایه تحت شرایط تعیین نادرست مدل وقتی اندازه‌ی نمونه در ناحیه افزایش یابد، ضعیف عمل می‌کند. برای مثال، هسن و همکاران (۱۹۸۳) یک مورد از بد مشخص‌سازی مدل را معرفی کردند که از طریق آزمون‌های معنی‌داری در نمونه‌های تا اندازه‌ی ۴۰۰ قابل کشف نیست. آن‌ها سپس نشان دادند که احتمال‌های در برگرفتن میانگین، \bar{Y} ، توسط بازه‌های اطمینان مدل پایه در نمونه‌گیری‌های مکرر به‌طور قابل توجهی کمتر از سطح مطلوب است و وقتی اندازه‌ی نمونه افزایش یابد، این کم‌پوشانی بدتر می‌شود. این عملکرد ضعیف تا حد قابل توجهی به دلیل ناسازگاری مجانبی برآوردگر مدل پایه نسبت به طرح نمونه‌گیری تصادفی طبقه‌ای است که توسط هسن و همکاران (۱۹۸۳) به کار گرفته شد.

۱.۵.۱ رویکرد طرح پایه

ابتدا به بررسی روش طرح پایه به‌عنوان وسیله‌ی سنتی در روش‌های نمونه‌گیری برای به دست آوردن برآوردهایی از مجموع یا میانگین جامعه می‌پردازیم. برآوردهای مستقیم نواحی، با استفاده از نتایج مربوط به مجموع‌های جامعه، با تغییر دادن نمادها به سادگی به دست می‌آیند. فرض کنید جامعه‌ی تحت نمونه‌گیری U مبتنی بر N عنصر متفاوت است که با برچسب‌های $j = 1, \dots, N$ شناسایی می‌شوند. فرض کنید مشخصه مورد نظر مانند y مرتبط به عنصر j ، به‌طور دقیق اندازه‌گیری شده باشد. بنابراین فرض می‌شود که خطاهای اندازه‌گیری وجود ندارند. در این‌جا، پارامتر مورد نظر، مجموع جامعه $Y = \sum_{\{j: j \in U\}} y_j$ یا میانگین جامعه $\bar{Y} = Y/N$ است. برای استخراج نمونه s از U با احتمال $P(s)$ ، از طرح نمونه‌گیری استفاده می‌شود. احتمال استخراج مشاهده نمونه $P(s)$ می‌تواند متکی بر متغیرهای طرحی معلوم از قبیل متغیرهای نشانگر طبقه و معیارهای اندازه‌گیری خوشه‌ها باشد. در عمل، برای اجرای یک طرح نمونه‌گیری از یک برنامه‌ی نمونه‌گیری استفاده می‌شود. مثلاً، نمونه تصادفی ساده با اندازه‌ی n را می‌توان با بیرون آوردن عدد تصادفی از ۱ تا N بدون جایگذاری به دست آورد. عمومی‌ترین طرح‌های نمونه‌گیری عبارت‌اند از؛ نمونه‌گیری تصادفی ساده طبقه‌ای و نمونه‌گیری چند مرحله‌ای طبقه‌ای.

برای انجام استنباط‌هایی درباره‌ی مجموع Y ، مقادیر y مربوط به نمونه‌ی انتخاب‌شده‌ی s را مشاهده می‌کنیم. برای سادگی کار فرض می‌کنیم همه‌ی عناصر $j \in s$ را می‌توان مشاهده کرد و پاسخ کامل داشت. در رویکرد طرح پایه برآوردگر \hat{Y} از Y را طرح ناریب گویند اگر امید ریاضی طرحی \hat{Y} برابر Y باشد؛ یعنی،

$$E_P(\hat{Y}) = \sum P(s)\hat{Y}_s = Y \quad (1.1)$$

که در آن جمع بندی روی همه‌ی نمونه‌های ممکن s تحت طرح مشخص و \hat{Y}_s مقدار \hat{Y} برای نمونه‌ی s است. دقت کنید منظور از زیرنویس P امید ریاضی نسبت به توزیع احتمال ایجاد شده توسط طرح نمونه‌گیری است. واریانس طرحی \hat{Y} به صورت

$$V_P(\hat{Y}) = E_P[\hat{Y} - E_P(\hat{Y})]^2$$

نشان داده می‌شود. برآوردگری از $V_P(\hat{Y})$ را به صورت $v(\hat{Y}) = s^2(\hat{Y})$ نشان می‌دهند، و اگر $E_P[v(\hat{Y})] = V_P(\hat{Y})$ ، $v(\hat{Y})$ را P ناریب گویند. برآوردگری مانند \hat{Y} طرح سازگار است اگر \hat{Y} برآوردگری P ناریب باشد و وقتی اندازه‌ی نمونه افزایش یابد $V(\hat{Y})$ به صفر میل می‌کند. به‌طور کلی لازم است P سازگاری در ساختار دنباله‌ای از جامعه‌های U_v در نظر گرفته شود، به‌طوری که وقتی $v \rightarrow \infty$ ، هم اندازه‌ی نمونه n_v و هم اندازه‌ی جامعه N_v به ∞ میل کنند. ویژگی P سازگاری برآوردگر واریانس $v(\hat{Y})$ به‌طور مشابه تعریف می‌شود. اگر برآوردگر $v(\hat{Y})$ هر دو P سازگار باشند، آن‌گاه رویکرد طرح پایه، به دور از مقادیر جامعه‌ای، استنباط‌های معتبر درباره‌ی Y را به دست می‌دهند به این معنا که وقتی اندازه‌ی نمونه افزایش یابد متغیر محوری مکرر وقتی اندازه‌ی نمونه‌ی افزایش یابد، پس تقریباً $(1 - \alpha)100$ درصد از بازه‌های اطمینان $[\hat{Y} - z_{\alpha/2}s(\hat{Y}), \hat{Y} + z_{\alpha/2}s(\hat{Y})]$ مقدار واقعی Y را در بر می‌گیرند. وزن‌های طرحی $w_j(s)$ نقشی مهم در ایجاد برآوردگرهای طرح پایه ایفا می‌کنند. این وزن‌های پایه‌ای ممکن است هم به s و هم به عنصر $j (j \in s)$ وابسته باشند. یک انتخاب مهم $w_j(s) = 1/\pi_j$ است که در آن به‌ازای $j, j = 1, \dots, N$ ، کمیت‌های $\pi_j = \sum_{s:j \in s} P(s)$ همان احتمال‌های شمول^{۲۹} اند.

تعریف ۱.۵.۱. (احتمال‌های شمول). مرتبه‌ی اول و دوم متناظر با طرح نمونه‌گیری $P(\cdot)$ عبارت‌اند از:

$$\pi_i = \sum_{s: i \in s} P(s), \pi_{ij} = \sum_{s: i, j \in s} P(s)$$

که در آن $\sum_{s: i \in s}$ جمع بندی روی تمام نمونه‌های s در بردارنده‌ی عنصر i و $\sum_{s: i, j \in s}$ جمع بندی روی همه‌ی نمونه‌های s در بردارنده‌ی عنصر $i, j = 1, \dots, N$ را نشان می‌دهد (سمپت، ۲۰۰۵).

در نبود اطلاعات کمکی جامعه، برآوردگر گسترشی

$$\hat{Y} = \sum_{j: j \in s} w_j y_j \quad (۲.۱)$$

را به کار می‌برند. در این حالت، شرط (۱.۱) به شرط زیر تبدیل می‌شود:

$$\sum_{\{s: j \in s\}} P(s) w_j(s) = 1, \quad j = 1, \dots, N. \quad (۳.۱)$$

گزینه‌ی $w_j(s) = 1/\pi_j$ در شرط ناریبی (۲.۱) صدق می‌کند، بنابراین به برآوردگر مشهور هوروتیز-تامپسون (H-T)^{۳۰} می‌انجامد. و بر پایه‌ی هر طرح نمونه‌گیری، ناریب است که به

^{۲۹}Inclusion probabilities

^{۳۰}Horwitz-thompson estimator

اصطلاح طرح ناریب گفته می‌شود. برای جزئیات برآورد واریانس به کوکران (۱۹۷۷)، ولتر (۱۹۸۵) مراجعه کنید. راثو (۱۹۷۹) نشان داده است که برآوردگر ناریب نامنفی واریانس \hat{Y} لزوماً به صورت زیر است:

$$v(\hat{Y}) = v(y) = -\sum_{j < k} \sum_{j,k \in s} w_{jk}(s) b_j b_k \left(\frac{y_j}{b_j} - \frac{y_k}{b_k} \right)^2 \quad (4.1)$$

که در آن وزن‌های $w_{jk}(s)$ در شرط ناریبی صدق می‌کنند و ثابت‌های ناصفر b_j چنانند که وقتی به‌ازای همه‌ی j ها $y_j \propto b_j$ باشد، واریانس \hat{Y} صفر می‌شود. برای مثال، در حالت خاص برای $w = 1/\pi_j$ و طرح اندازه ثابت نمونه‌ای، داریم $b_j = \pi_j$ و $w_{jk}(s) = (\pi_{jk} - \pi_j \pi_k) / (\pi_{jk} \pi_j \pi_k)$ که در آن به‌ازای $j < k = 1, \dots, N$ کمیت‌های $\pi_{jk} = \sum_{s: (j,k) \in s} P(s)$ احتمال‌های شمول توام‌اند که فرض می‌شود مثبت هستند. برآوردگر واریانس (۴.۱) در این حالت به برآوردگر مشهور واریانس سن-ییتس-گروندی^{۳۱} (S-Y-G) تبدیل می‌شود. برای به‌دست آوردن برآوردگر حوزه‌ای کافی است، ابتدا تعریف می‌کنیم

$$y_{(ij)} = \begin{cases} y_j & j \in U_i \\ 0 & o.w \end{cases}$$

$$a_{(ij)} = \begin{cases} 1 & j \in U_i \\ 0 & o.w \end{cases}$$

بنابراین

$$Y(y_i) = \sum_{j \in U} a_{ij} = \sum_{j \in U} 1 = N_i \quad \text{و} \quad Y(y_i) = \sum_{j \in U} y_{ij} = \sum_{j \in U} y_j = Y_i$$

نظر و N_i اندازه‌ی آن است.

برآوردگرهای حوزه‌ای به راحتی به صورت $\hat{Y}_i = \hat{Y}(y_i) = \sum_{j \in s} w_j y_{ij} = \sum_{j \in s_i} w_j y_j$ به دست می‌آیند.

برای مطالعه‌ی بیشتر در مورد سایر برآوردگرها نظیر برآوردگر رگرسیون تعمیم یافته و برآوردگر نسبی به راثو (۲۰۰۳) مراجعه کنید.

انتقادی که به رویکرد طرح پایه (یا نمونه‌گیری احتمالاتی) وجود دارد، این است که هر چند استنباط‌های حاصل از آن آزاد فرض‌اند، به جای آن که تنها به نمونه‌ی خاص s استخراج شده اشاره کنند به نمونه‌گیری مکرر اشاره دارند. یافتن طرح «بهینه» نمونه‌گیری برای استفاده در برآوردگرهای مستقیم برای مجموع یا میانگین بزرگ ناحیه‌ها، طی سال‌های گذشته توجه زیادی را به خود معطوف کرده است. به‌ویژه، موضوع‌های مهم طراحی از قبیل تعداد طبقه‌ها، تشکیل طبقه‌ها، تخصیص نمونه، و احتمال‌های گزینش، مورد بررسی دقیق قرار گرفته است

^{۳۱}San-yates-groendi variance estimator

(کوکران، ۱۹۷۷). هدف در اینجا پیدا کردن طرحی «بهینه» است که MSE برآوردگر مستقیم را با توجه به هزینه‌ای ثابت می‌نیمم سازد.

این هدف در عمل به واسطه‌ی محدودیت‌های عملیاتی و عامل‌های دیگر به‌ندرت حاصل می‌شود. در نتیجه، طرحی «بینابینی» که «نزدیک» به طرح بهینه باشد به کار می‌رود. در عمل، پیش‌بینی و برنامه‌ریزی برای همه‌ی ناحیه‌ها و کاربست‌ها، در داده‌های آمارگیری میسر نیست، زیرا «مشتری همیشه بیش از آنچه در مرحله‌ی طراحی مشخص شده است، خواهد خواست» (فولر، ۱۹۷۳). در نتیجه، با توجه به تقاضای رو به رشد آمارهای کوچک ناحیه‌ای قابل اعتماد، برآوردگرهای غیرمستقیم همیشه در عمل لازم خواهند بود. همچنین موضوع‌هایی از طراحی نمونه‌گیری که تاثیری بر روش برآورد کوچک ناحیه‌ای دارند، به‌ویژه در برنامه‌ریزی و طراحی آمارگیری‌های نمونه‌ای بزرگ مقیاس را نیز می‌توان بررسی کرد که عبارت‌اند از: می‌نیمم‌سازی خوشه‌بندی، طبقه‌بندی، تخصیص نمونه، تلفیق آمارگیری‌ها، آمارگیری با چارچوب دوگانه و آمارگیری مکرر.

هرچند این روش‌ها برای به‌دست آوردن برآوردگرهای قابل اعتماد کاربرد دارند، اما در عمل هنگامی که اندازه‌ی نمونه کم باشد کمکی در به‌دست آوردن برآوردگرهای با دقت کافی برای کوچک ناحیه‌ها نمی‌کنند. برای مطالعه‌ی این روش‌ها به سینگ و همکاران (۱۹۹۴)، مارکر (۲۰۰۱) و راثو (۲۰۰۳) مراجعه کنید.

۶.۱ روش‌های برآورد حوزه‌ای غیرمستقیم

زمانی که اندازه‌ی نمونه به‌قدر کافی بزرگ باشد یا به‌طور کلی نمونه بهینه شده باشد خطاهای استاندارد قابل قبول خواهند بود. اما در ساختار برآورد کوچک ناحیه‌ای، برآوردهای مستقیم به‌دلیل اینکه اندازه نمونه‌ها به‌طور غیر قابل قبولی کوچک است، منجر به خطاهای استاندارد می‌شوند که به‌طور غیرقابل قبولی بزرگ‌اند؛ حتی گاهی اوقات ممکن است برای برخی ناحیه‌ها هیچ نمونه‌ای انتخاب نشده باشد. همچنین لازم است از برآوردگرهای غیرمستقیم که با بهره‌گیری از مدل‌ها و اطلاعات کمکی ساخته می‌شوند نیز استفاده شود چرا که خطای استاندارد کاهش می‌یابد. برآوردگرهای حوزه‌ای غیرمستقیم مبتنی بر مدل‌های پیوندی ضمنی را برآوردگرهای نامستقیم سنتی گویند، که عبارت‌اند از: برآوردگرهای هم‌گذاشتی، برآوردگرهای ترکیبی و برآوردگرهای جیمز-استاین^{۳۲}. با توجه به این که برآوردگرهای غیرمستقیم حاضر برخلاف برآوردگرهای غیرمستقیم جمعیت شناختی (راثو، ۲۰۰۳) تنها از داده‌های ثبتي و سرشماری استفاده می‌کنند و نمونه‌گیری در بین نیست، تا حد زیادی مبتنی بر داده‌های آمارگیری نمونه‌ای در ترکیب با داده‌های کمکی از جامعه هستند.

۷.۱ دیدگاه‌های استنباطی در برآورد کوچک ناحیه‌ای

در دیدگاه استنباطی معمولاً دو نوع دیدگاه وجود دارد، یکی از این دیدگاه‌ها، دیدگاه کلاسیک و دیگری بیزی است. در دیدگاه کلاسیک پارامتر مقدار ثابت ولی نامعلوم می‌باشد، اما هدف در این روش ارزیابی رفتار روش‌های استنباطی بر اساس اصل تکرارپذیری است. اما در دیدگاه بیزی مقدار پارامتر، تحقق از یک توزیع احتمالی می‌باشد که معروف به توزیع پیشین است. از طرفی هدف در این رویکرد ارزیابی عدم قطعیت موجود در روش‌های استنباطی است که توسط توزیع پسین تحقق می‌یابد. از این رو با توجه به پیشرفت‌هایی که در حیطه رایانه بوجود آمده است برازش مدل در دو دیدگاه کار سختی نیست، اما در دیدگاه کلاسیک برازش مدل رگرسیونی با اثرات تصادفی، در برآورد نواحی کوچک، از روش بهترین پیشگوی نارایب خطی تجربی میسر می‌باشد (رابینسون، ۱۹۹۱؛ جیانگ و لاهیری، ۲۰۰۶). اما این رویکرد همیشه عملی نخواهد بود زیرا اگر متغیر هدف غیرنرمال باشد به دست آوردن آن در روش^{۳۳} EBLUP عملی نخواهد بود (گوش و همکاران، ۱۹۹۸). با توجه به مشکلات این روش دیدگاه بیزی اهمیت پیدا می‌کند چرا که در این روش محاسبه مدل‌های نواحی کوچک به دلیل وجود الگوریتم MCMC قابل محاسبه خواهد بود.

۸.۱ روش‌های مبتنی بر نمونه‌گیری

روش‌های مبتنی بر نمونه‌گیری که به سه صورت زیر می‌باشد برای استنباط بیزی مورد استفاده قرار می‌گیرد.

۱.۸.۱ مونت کارلوی زنجیره مارکوفی MCMC

یکی از روش‌های معمول که به دلیل انتگرال‌های پیچیده موجود در تابع درستنمایی برای تعیین تحلیلی توزیع پسین به کار می‌رود MCMC است، چرا که شبیه‌سازی از توزیع پسین به کمک این روش آسان‌تر خواهد بود. این روش در سال ۱۹۴۹ با انتشار مقاله‌ای به نام روش مونت کارلو توسط دو ریاضیدان آمریکایی به نام‌های نیمن و یولام شناخته شد. پایه‌های نظری این روش از مدت‌ها قبل شناخته شده بود و در قرن نوزدهم و اوایل قرن بیستم گاهی اوقات مسائل آماری را با کمک شبیه‌سازی کمیت‌های تصادفی که همان روش مونت کارلو است، حل می‌کردند. اما تا پیش از اختراع رایانه‌های الکترونیکی از این روش به دلیل خسته کننده بودن شبیه‌سازی کمیت‌های تصادفی زیاد استفاده نمی‌شد. بنابراین می‌توان شروع استفاده از روش مونت کارلو را به عنوان یک تکنیک عددی عمومی، همزمان با اختراع رایانه‌های الکترونیکی دانست. نام مونت کارلو از شهر مونت کارلو واقع در موناکو گرفته شده است.

^{۳۳} Estimator best linear unbiased Predictor

به‌طور کلی، روش‌های MCMC بر اساس یک الگوریتم مکرر^{۳۴}، زنجیری تقلیل ناپذیر^{۳۵}، و نادره‌ای ساخته شده است و نیز از دنباله‌ای از پارامترهای مورد نظر تولید می‌شود، به گونه‌ای که توزیع مانای آن همان توزیع پسین است. در واقع با تکرار متوالی مراحل الگوریتم یک نمونه وابسته مارکوفی از توزیع پسین تولید می‌شود. دنباله آغازین تکرارهای زنجیر مارکوف قبل از زمان همگرایی، دوره داغیدن نامیده می‌شود. مشاهدات بعد از این دوره را می‌توان به عنوان مشاهداتی از توزیع پسین دانست. در روش‌های MCMC همگرایی زنجیر مساله‌ای اساسی به شمار می‌رود. از جمله مهم‌ترین این روش‌ها الگوریتم‌های گیبز^{۳۶} و متروپولیس-هستینگز^{۳۷} است. با به‌کارگیری از قضیه‌ی بیز، مدل (۱.۷) برای به‌دست آوردن توزیع پسین پارامترهای ناحیه‌ی کوچک استفاده می‌شود.

$$f(\mu, \lambda | y) = f(y, \mu | \lambda) f(\lambda) / f_1(y),$$

$$f_1(y) = \int f(y, \mu | \lambda) f(\lambda) d\mu d\lambda \quad (5.1)$$

$$f(\mu | y) = \int f(\mu | y, \lambda) d\lambda = \int f(\mu | y, \lambda) f(\lambda | y) dy$$

پس بنابر مدل بالا ارزیابی $f(\mu | y)$ و کمیت‌های پسینی مربوط مانند $E[h(\mu | y)]$ شامل انتگرال‌گیری‌های چند بعدی است. اما اغلب اجرای انتگرال‌گیری به‌طور تحلیلی تنها نسبت به برخی از مولفه‌های μ و y میسر است. اگر مسئله‌ی کاهش یافته شده تنها شامل انتگرال‌گیری یک یا دو بعدی باشد، انتگرال‌گیری عددی مستقیم می‌تواند به‌کار گرفته شود تا کمیت‌های پسینی مطلوب به‌دست آورده شوند. اما برای مسئله‌های پیچیده، لازم می‌آید که انتگرال‌های با بعد زیاد ارزیابی شوند. به‌نظر می‌رسد روش‌های تازه ابداع شده‌ی زنجیر مارکوفی MCMC می‌تواند تا حد زیادی مشکلات محاسباتی را برطرف کند. روش‌های زنجیره‌ی مارکوف مونت کارلویی MCMC شامل انواع الگوریتم‌های محاسباتی که برای شبیه‌سازی توزیع‌های پسین در مدل بیزی به‌کار می‌روند (گیلکس و همکاران، ۱۹۹۶). هنگام مطالعه مدل‌های پیچیده، تنها روش نمونه‌برداری از توزیع پسین، برآوردگرهای بیزی است. این روش نمونه‌هایی از توزیع پسینی تولید می‌کند، و سپس نمونه‌های شبیه‌سازی شده را برای تقریب کمیت‌های پسینی مطلوب به‌کار می‌برد. می‌توان برای تشخیص همگرایی و کارایی از بسته‌های نرم‌افزاری BUGS^{۳۸} و CODA روش MCMC استفاده کرد (گیلکس و همکاران، ۱۹۹۶).

با توجه به مطالب بیان شده فرض کنید $\eta = (\mu^T, \lambda^T)^T$ بردار پارامترهای کوچک ناحیه‌ای μ و پارامترهای مدل λ باشد. به‌طور کلی، به‌دلیل مخرج غیرقابل پیگیری $f_1(y)$ در مدل (۱.۷)، استخراج نمونه‌های مستقل از پسین توام $f(\mu, \lambda | y)$ میسر نیست. روش MCMC با

^{۳۴} Iterative

^{۳۵} Irreducibility

^{۳۶} Gibbs

^{۳۷} Metropolis-hastings

^{۳۸} Bayesian inference using gibbs sampling

تشکیل یک زنجیر مارکوفی $\{\eta^{(k)}, k = 0, 1, 2, \dots\}$ به گونه‌ای که توزیع $\eta^{(k)}$ به توزیع مانا (یا ناوردای^{۳۹}) یکتا برابر با $f(\eta|y)$ همگرا شود، از این مشکل پرهیز می‌کند؛ و این توزیع را با $\pi(\eta)$ نشان می‌دهند. بنابراین، بعد از یک دوره‌ی «تطبیق» به قدر کافی بزرگ مانند d می‌توانید $\eta^{(d+1)} \dots \eta^{(d+D)}$ را، صرف‌نظر از نقطه‌ی شروع $\eta^{(0)}$ ، همانند D نمونه‌ی مستقل از توزیع هدف تلقی $f(\eta|y)$ کنید.

برای ساختن یک زنجیر مارکوفی، نیاز به مشخص کردن احتمال انتقال یک گامی (یا هسته‌ی) $P(\eta^{(k+1)}|\eta^{(k)})$ دارید که تنها به «وضعیت» جاری زنجیر $\eta^{(k)}$ وابسته باشد. یعنی، توزیع شرطی $\eta^{(k+1)}$ به شرط $\eta^{(0)}, \dots, \eta^{(k)}$ به «سابقه‌ی» زنجیر $\{\eta^{(0)}, \dots, \eta^{(k+1)}\}$ وابسته نیست. هسته‌ی انتقال باید در شرط مانایی صدق کند:

$$\int \pi(\eta^{(k)}) P(\eta^{(k+1)}|\eta^{(k)}) d\eta^{(k)} = \pi(\eta^{(k+1)}) \quad (6.1)$$

معادله‌ی (۶.۱) نشان می‌دهد که اگر $\eta^{(k)}$ از توزیع $\pi(\cdot)$ باشد، $\eta^{(k+1)}$ نیز از توزیع $\pi(\cdot)$ خواهد بود. بنابراین اگر زنجیر «برگشت‌پذیر» باشد، مانایی برقرار خواهد بود:

$$\pi(\eta^{(k)}) P(\eta^{(k+1)}|\eta^{(k)}) d\eta^{(k)} = \pi(\eta^{(k+1)}) \quad (7.1)$$

درستی رابطه‌ی (۷.۱) این اطمینان را می‌دهد که توزیع مانای زنجیر ایجاد شده با $P(\cdot|\cdot)$ همان $\pi(\cdot)$ می‌باشد. همچنین لازم است این اطمینان حاصل شود که توزیع $\eta^{(k)}$ به شرط $\eta^{(0)}$ که با $P(\eta^{(k)}|\eta^{(0)})$ نشان داده می‌شود، صرف‌نظر از $\eta^{(0)}$ ، به $\pi(\eta^{(k)})$ همگرا می‌شود. برای برآورده شدن این شرط، لازم است زنجیر «نافروکاستنی^{۴۰}» و «نادوره‌ای» باشد. نافروکاستنی بدین معنی می‌باشد که زنجیر از همه‌ی نقطه‌های شروع $\eta^{(0)}$ با احتمال مثبت سرانجام به هر مجموعه‌ی غیرتهی در فضای وضعیت خواهد رسید. نادوره‌ای بودن به این معنی است که زنجیر بین مجموعه‌ها مجاز نیست و به گونه‌ای دوره‌ای نوسان می‌کند. در مورد یک زنجیر فروکاستنی و نادوره‌ای، قضیه‌ی ارگودیک^{۴۱} زیر صدق می‌کند: هنگامی که $D \rightarrow \infty$

$$\bar{h}_D = \frac{1}{D} \sum_{k=d+1}^{d+D} h(\eta^{(k)}) \xrightarrow{p} E[h(\eta|y)] \quad (8.1)$$

که در آن همگرایی در احتمال را نشان می‌دهد. بدین صورت، به‌ازای اندازه‌ی بزرگ از D ممکن است، بتوانید برآوردگر \bar{h}_D از $E[h(\eta|y)]$ را با دقت کافی به‌دست آورید. اما یافتن خطای استاندارد مونت کارلویی \bar{h}_D به‌دلیل وابستگی بین نمونه‌های شبیه‌سازی شده‌ی $\eta^{(d+1)}, \dots, \eta^{(d+D)}$ آسان نیست.

عملاً، روش‌های MCMC را می‌توان با مثالی از زنجیره مارکف که توسط تانر و وانگ (۱۹۸۷) مطرح کرده‌اند، نشان داد. فرض کنید که در توزیع پسین از $X = (z, y)$ استفاده کنید.

^{۳۹}Invariance

^{۴۰}irreducible

^{۴۱}Ergodic

چگالی‌های پسین حاشیه‌ای z و y را می‌توان همانند زیر مشتق‌گیری کرد:

$$p(\mathbf{z}) = \int p(\mathbf{z}|\mathbf{y})p(\mathbf{y})d\mathbf{y}$$

و

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

بنابراین $p(\mathbf{z}|\mathbf{y})$ و $p(\mathbf{y}|\mathbf{z})$ مشخص می‌باشد و برای $p(\mathbf{z})$ نیاز به حل کردن است. و از رویکرد مبتنی بر نمونه‌گیری استفاده می‌کنند، بنابراین می‌توان از $p_0(\mathbf{z})$ شروع کرد و $p(\mathbf{z})$ را برآورد کرد و در نهایت $\mathbf{Z}^{(o)} \sim p_0(\mathbf{z})$ را نوشت. و همین‌طور می‌توان نشان داد $\mathbf{Y}^{(l)} \sim p(\mathbf{y}|\mathbf{z}^{(o)})$. توزیع حاشیه‌ای $\mathbf{Y}^{(l)}$ برابر است با

$$p_{\setminus}(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{z})p_0(\mathbf{z})d\mathbf{z}$$

در مرحله‌ی بعد $\mathbf{Z}^{(l)} \sim p(\mathbf{z}|\mathbf{y}^{(l)})$ را با توزیع حاشیه‌ای زیر نشان می‌دهند:

$$\begin{aligned} p_{\setminus}(\mathbf{z}) &= \int p(\mathbf{z}|\mathbf{y})p_{\setminus}(\mathbf{y})d\mathbf{y} \\ &= \int p(\mathbf{z}|\mathbf{y}) \int p(\mathbf{y}|\mathbf{z}')p_0(\mathbf{z}')d\mathbf{z}'d\mathbf{y} \\ &= \int \left(\int p(\mathbf{z}|\mathbf{y})p(\mathbf{y}|\mathbf{z}')d\mathbf{y} \right) p_0(\mathbf{z}')d\mathbf{z}' \\ &= \int h(\mathbf{z}, \mathbf{z}')p_0(\mathbf{z}')d\mathbf{z}' = I_h(p_0(\mathbf{z})) \end{aligned}$$

تانر و وانگ (۱۹۸۷) نشان دادند که اگر این فرآیند را تکرار کنند، به تعداد کافی جفت $(\mathbf{Z}^{(i)}, \mathbf{Y}^{(i)})$ ، دنباله‌ی $\mathbf{Z}^{(i)}, \mathbf{Y}^{(i)}$ که همگرا به $p(\mathbf{z}, \mathbf{y})$ و نیز با توزیع حاشیه‌ای $p(\mathbf{z})$ و $p(\mathbf{y})$ همراه هستند، به‌دست آورده می‌شود. به عبارت دیگر برای بسندگی بزرگ i ، $\{\mathbf{Z}^{(i)}\}$ را می‌توان از یک نمونه در توزیع حاشیه‌ای $p(\mathbf{z})$ در نظر گرفت و $\{\mathbf{Y}^{(i)}\}$ را نیز می‌توان از یک نمونه با توزیع حاشیه‌ای $p(\mathbf{y})$ فرض کرد.

بنابراین، الگوریتم تکراری است و نیاز به یک دوره سوختن قبل از نمونه‌های ایجاد شده دارد و می‌توان فرض کرد از تراکم حاشیه‌ای شکل گرفته است. هنگامی که خروجی MCMC یک مدل تولید کند، می‌توان مقادیر پسین مورد نظر را محاسبه نمود و با توجه به نمونه‌هایی که برای هر جز بردار پارامترها به‌دست می‌آیند، آن‌ها را می‌توان به عنوان نمونه‌هایی از توزیع پسین حاشیه‌ای مربوط در نظر گرفت. ساده‌ترین روش برای به‌دست آوردن برآورد یک تابع از پارامتر خاص (برای مثال، $\phi = f(\theta_i)$) از طریق میانگین ارگودیک امکان‌پذیر می‌باشد و به‌صورت زیر است:

$$\hat{\phi} = \hat{f}(\theta_i) = \frac{1}{D} \sum_{k=d+1}^{d+D} \phi^{(k)},$$

d طول دوره‌ی سوختن و برای به‌دست آوردن همگرایی استفاده می‌شود. D اندازه‌ی نمونه باقی‌مانده و $\phi^{(k)} = f(\theta_{(i)}^{(k)})$ با $\theta_{(i)}^{(k)}$ از زنجیره‌ی k ام در نمونه‌ی MCMC برای θ_i می‌باشد. با توجه به برآورد دقت برآوردگرها، واریانس پسین را می‌توان طبق رابطه‌ی زیر محاسبه نمود:

$$\hat{Var}(\phi|y) = \frac{1}{D-1} \sum_{k=d+1}^{d+D} (\phi^{(k)} - \hat{\phi})^2.$$

اگر صورت ریاضی برای امیدریاضی شرطی ϕ به شرط داده‌ی y و تمامی پارامترها شناخته شده باشد، می‌توان از برآوردگر بهبود یافته راتو-بلک ولز^{۴۲} استفاده کرد (گلفند و اسمیت، ۱۹۹۱):

$$\phi^{(RB)} = \frac{1}{D} \sum_{k=d+1}^{d+D} E(\phi|y, \xi^{(k)}),$$

که $\{\xi^{(k)}\}$ توزیع پسین حاشیه‌ای از $y|\xi$ است. واریانس پسین همانند زیر محاسبه می‌شود:

$$\tilde{Var}(\phi|y) = \frac{1}{D} \sum_{k=d+1}^{d+D} Var(\phi|y, \xi^{(k)}) + \frac{1}{D-1} \left(\sum_{k=d+1}^{d+D} Var(\phi|y, \xi^{(k)} - \phi^{(RB)})^2 \right).$$

رابطه‌ی فوق معتبر خواهد بود، اگر مقادیر نمونه تولید شده تقریباً مستقل و یکسان (i.i.d) باشند. در برنامه‌های کاربردی، اغلب چندین زنجیره‌ی موازی با نقاط اولیه بیش‌پراکنش شده قابل اجرا خواهد بود.

فرض کنید زنجیره‌های موازی L قابل اجرا باشد، در این هنگام معادلات ارگودیک را می‌توان به‌صورت زیر تغییر داد:

$$\hat{\phi} = \frac{1}{D} \sum_{l=1}^L \sum_{k_l=d+1}^{d+D} \phi^{(k_l)}$$

d طول دوره‌ی سوختن و D اندازه‌ی نمونه‌ی باقی‌مانده (d و D فرض می‌شود برای همه زنجیره‌ها یکسان باشد)؛ و $\phi^{(k_l)} = f(\theta_i^{(k_l)})$ با داشتن $\theta_i^{(k_l)}$ در زنجیره‌ی k_l ام از نمونه (MCMC) برای θ_i در زنجیره‌ی l ام می‌باشد. واریانس پسین ϕ می‌تواند به‌صورت زیر برآورده شود:

$$\hat{Var}(\phi|y) = \frac{D-1}{D} W + \frac{1}{D} B$$

از این‌رو

- $\bar{\phi}_l = \frac{1}{D} \sum_{k_l=d+1}^{d+D} \phi^{(k_l)}$ - مقادیر میانگین در زنجیره‌ی l
- $\bar{\phi} = \frac{1}{L} \sum_{l=1}^L \bar{\phi}_l$ - میانگین هدف در میان تمام زنجیره‌ها
- $W = \frac{1}{L(D-1)} \sum_{l=1}^L \sum_{k_l=d+1}^{d+D} (\phi^{(k_l)} - \bar{\phi}_l)^2$ - واریانس درون زنجیره‌ای
- $B = \frac{D}{L-1} \sum_{l=1}^L (\bar{\phi}_l - \bar{\phi})^2$ - واریانس بین زنجیره‌ای.

۲.۸.۱ نمونه‌گیر گیبز

نمونه‌برداری گیبز به افتخار نام فیزیکدان جوسی ویلارد گیبز، نامگذاری شده است و اشاره به یک مقایسه بین الگوریتم نمونه‌برداری و فیزیک آماری دارد. الگوریتم شرح داده شده، توسط برادران استوارت و دونالد گمن در سال ۱۹۸۴، یعنی هشت دهه پس از مرگ گیبز تشریح شد. نمونه‌برداری گیبز در تجسم اولیه، یک حالت خاص از این الگوریتم متروپلیس-هستینگز است. نکته‌ای که در نمونه‌برداری گیبز وجود دارد این است که، برای یک توزیع چند متغیره، نمونه‌برداری از توزیع شرطی ساده‌تر است از محاسبه توزیع حاشیه‌ای، که با انتگرال‌گیری بر روی توزیع توام به‌دست می‌آید.

منظور از نمونه‌برداری گیبز یا نمونه‌بردار گیبز در مطالعات آماری، الگوریتمی است که بر مبنای تئوری زنجیره مارکوف مونت کارلو طراحی شده است. کاربرد این الگوریتم در تولید دنباله‌ای از مشاهدات از یک تابع توزیع احتمالاتی چند متغیره است که تولید نمونه از آن به صورت مستقیم دشوار است. این دنباله را می‌توان برای تخمین توزیع هم‌زمان (مثلاً برای تولید هیستوگرام توزیع)، تخمین توزیع حاشیه‌ای بر روی یک یا زیر مجموعه‌ای از متغیرهای توزیع (مانند پارامتر پنهان یا متغیرهای پنهان)، یا برای محاسبه یک انتگرال (مانند امید ریاضی متغیرها) استفاده نمود. اغلب برخی از متغیرها وابسته به مشاهدات هستند که مقدار آن‌ها مشخص است و بنابراین نیازی به نمونه‌برداری برای آن‌ها نیست.

نمونه‌برداری گیبز معمولاً به عنوان ابزاری برای استنباط آماری و به‌ویژه در استنباط بی‌زی استفاده می‌شود. این روش یک الگوریتم تصادفی (از آن جهت که با استفاده از اعداد تصادفی نمونه تولید می‌کند) که می‌تواند جایگزینی برای الگوریتم‌های قطعی در استنباط آماری باشد. نمونه‌برداری گیبز مانند دیگر الگوریتم‌های زنجیره مارکوف مونت کارلو، زنجیره مارکوفی را از نمونه‌ها تولید می‌کند، به‌طوری که هر نمونه وابسته به نمونه‌های نزدیک است؛ بنابراین اگر نمونه‌های مستقل مورد نظر است، باید نمونه‌برداری محتاطانه انجام پذیرد، این کار اغلب با نازک‌سازی زنجیره نمونه‌های حاصل شده انجام می‌شود، بدین شکل، تنها مقدار n ام زنجیره مثلاً ۱۰۰ام انتخاب می‌شود. علاوه بر این، نمونه‌های ابتدای زنجیره (تکرارهای سوخته) احتمالاً نمایانگر خوبی برای توزیع مورد نظر نخواهند بود.

با توجه به مطالب بالا می‌توان گفت که نمونه‌گیری گیبز یک الگوریتم MCMC با هسته تغییر وضعیت است که به کمک توزیع‌های شرطی کامل ساخته شده است.

بنابراین برای تولید نمونه‌های $\eta^{(k)}$ بردار η را به بلوک‌هایی مناسب η_1, \dots, η_r افراز می‌کنند. برخی از بلوک‌ها ممکن است تنها یک عنصر داشته باشند در صورتی که بلوک‌های دیگر ممکن است شامل بیش از یک عنصر باشند. برای مثال، مدل پایه‌ای در سطح ناحیه با $\mu = (\theta_1, \dots, \theta_m)^T = \theta$ و $\lambda = (\beta^T, \sigma_v^2)^T$ را در نظر بگیرید. در این حالت η را می‌توان به صورت $\eta_{m+2} = \sigma_v^2$ (برای $r = m + 2$) افراز کرد. ما به مجموعه‌ی

توزیع‌های شرطی گیزی زیر نیاز داریم:

$$f(\eta_1 | \eta_2, \dots, \eta_r, y), f(\eta_r | \eta_1, \dots, \eta_{r-1}, y), \dots, f(\eta_2 | \eta_1, \eta_3, \dots, \eta_r, y).$$

نمونه‌گیر گیزی از این توزیع‌های شرطی برای ساختن یک هسته‌ی انتقالی $P(\cdot|\cdot)$ ، به‌طوری که توزیع مانای زنجیر مارکوفی حاصل $\pi(\eta) = f(\eta|y)$ باشد، استفاده می‌شود. این نتیجه از این واقعیت بر می‌آید که $f(\eta|y)$ به‌وسیله‌ی مجموعه‌ی توزیع‌های شرطی گیزی به‌طور یکتا تعیین می‌شود.

اگر یک توزیع شرطی دارای صورتی استاندارد و بسته مثلا مانند نرمال یا گامای وارونه باشد، آن‌گاه نمونه‌ها می‌توانند مستقیما از آن توزیع شرطی تولید شوند. در غیر این صورت، الگوریتم‌های بدیل مانند نمونه‌گیری ردی متروپلیس-هستینگز می‌توانند به‌کار گرفته شوند تا نمونه‌هایی از توزیع شرطی تولید شوند. برخی مولفان به‌کارگیری M-H را برای حالت‌هایی به‌صورت بسته نیز پیشنهاد می‌کنند. اگر M-H تنها برای حالت‌های بدون صورت بسته به‌کار رود، آن‌گاه الگوریتم را M-H درون گیزی گویند. الگوریتم نمونه‌گیری گیزی شامل مراحل زیر است:

- مرحله‌ی اول. یک نقطه‌ی شروع $\eta^{(0)}$ با مولفه‌های $\eta_1^{(0)}, \dots, \eta_r^{(0)}$ را انتخاب کنید؛ قرار دهید $k = 0$. مثلا، می‌توان از برآوردهای بی‌تجربی (EB) μ یا پارامترهای مدل λ به‌عنوان مقدارهای شروع استفاده کرد. اما نقطه‌های شروع می‌توانند دلخواه باشند.
 - مرحله دوم. بردار $\eta^{(k+1)} = (\eta_1^{(k+1)}, \dots, \eta_r^{(k+1)})$ را به شرح زیر تولید کنید:
مولفه‌ی $\eta_1^{(k+1)}$ را از $f(\eta_1 | \eta_2^{(k)}, \dots, \eta_r^{(k)}, y)$ ، مولفه‌ی $\eta_2^{(k+1)}$ را از $f(\eta_2 | \eta_1^{(k+1)}, \eta_3^{(k)}, \dots, \eta_r^{(k)}, y)$ ، و $\eta_r^{(k+1)}$ را از $f(\eta_r | \eta_1^{(k+1)}, \dots, \eta_{r-1}^{(k+1)}, y)$ استخراج کنید.
 - مرحله‌ی سوم. قرار دهید $k = k + 1$ و به مرحله‌ی اول بروید.
- مراحل دوم و سوم یک چرخه برای هر k را تشکیل می‌دهند. دنباله‌ی $\{\eta^{(k)}\}$ که به‌وسیله‌ی نمونه‌گیر گیزی تولید شده یک زنجیر مارکوفی با توزیع مانای $\pi(\eta) = f(\eta|y)$ است؛ توجه کنید که هسته‌ی انتقالی تک گامی برابر حاصل ضرب r توزیع شرطی گیزی است (اسمیت و گلفند، ۱۹۹۰).

۳.۸.۱ الگوریتم متروپلیس-هستینگز درون گیزی

الگوریتم متروپلیس-هستینگز توسط متروپلیس و همکاران در سال (۱۹۵۳) و هستینگز (۱۹۷۰) معرفی شد. الگوریتم M-H برای هر یک از متغیرهای مورد نظر، با استفاده از تابع چگالی توام کامل و توزیع‌های پیشنهاد شده (مستقل)، نمونه‌ها را از توزیع احتمال شبیه‌سازی می‌کند.

این روش خاصا در مواردی به‌کار می‌رود که نمونه‌گیری مستقیم دشوار است. ترتیب به دست آمده را می‌توان برای تقریب‌گیری توزیع (تولید یک هیستوگرام) یا برای محاسبه

یک انتگرال استفاده کرد. متروپلیس-هستینگز و سایر الگوریتم‌های زنجیره مونت کارلویی معمولاً برای نمونه‌گیری از توزیع‌های چند بعدی در ابعاد زیاد، استفاده می‌شوند. همگرایی این الگوریتم نسبت به نمونه‌گیر گیبز منظم‌تر و کندتر می‌باشد. اما الگوریتم M-H انعطاف‌پذیری لازم را از طریق انتخاب چگالی پیشنهادی به‌دست می‌آورد. و از مشکلات الگوریتم M-H می‌توان تشخیص همگرایی زنجیر و زمان طولانی محاسبات دانست.

حال اگر همه‌ی توزیع‌های شرطی گیبزی صورت‌هایی بسته متعلق به خانواده‌های استاندارد داشته باشند، آن‌گاه تولید نمونه‌ها از توزیع‌های شرطی سر راست است. اگر توزیع شرطی صورتی بسته نداشته باشد، در این صورت برای تولید نمونه‌ها از آن توزیع شرطی چند روش وجود دارند. به‌ویژه، نمونه‌گیری ردی را برای توزیع‌های شرطی از نوع لگ کاو تک متغیره (گیلکز و وایلد، ۱۹۹۲) و به‌طور کلی‌تر M-H درون گیبزی را خواهیم داشت (متروپلیس و همکاران، ۱۹۵۳؛ هستینگز، ۱۹۷۰).

فرض کنید $f(\eta|\eta_{-i}^k, \mathbf{y})$ توزیع شرطی گیبزی پس از تکمیل $(i-1)$ استخراج اول از چرخه‌ی $(k+1)$ ام نمونه‌گیری گیبزی را نشان دهد، که در آن $\eta_{-i}^k = \{\eta_1^{k+1}, \dots, \eta_{i-1}^{k+1}, \eta_{i+1}^k, \dots, \eta_r^k\}$ الگوریتم M-H برای تولید یک نمونه‌ی $\eta_i^{(k+1)}$ از $f(\eta_i|\eta_{-i}^{(k)}, \mathbf{y})$ شامل مراحل زیر است.

- مرحله ی اول. توزیع $f(\eta_i|\eta_{-i}^{(k)}, \mathbf{y})$ را با یک چگالی پیشنهادی مانند $q_i(\eta_i|\eta_i^{(k)}, \eta_{-i}^{(k)})$ که نمونه‌گیری از آن آسان است از قبیل توزیع نرمال یا تی‌استودیت، تقریب بزنید. چگالی پیشنهادی ممکن است به مقدارهای جاری $\{\eta_i^{(k)}, \eta_{-i}^{(k)}\}$ وابسته باشد.
- مرحله‌ی دوم. برای η_i «نامزدی» مانند η_i^* از چگالی نامزد \mathbf{u} را از توزیع یکنواخت $(0, 1)$ تولید کنید.

- مرحله‌ی سوم. اگر $\mathbf{u} \leq a(\eta_{-i}^{(k)}, \eta_i^{(k)}, \eta_i^*)$ قرار دهید $\eta_i^{(k+1)} = \eta_i^*$ وگرنه $\eta_i^{(k+1)} = \eta_i^{(k)}$ ، که در آن احتمال پذیرش $a(\eta_{-i}^{(k)}, \eta_i^{(k)}, \eta_i^*)$ با رابطه‌ی زیر تعریف می‌شود:

$$a(\eta_{-i}^{(k)}, \eta_i^{(k)}, \eta_i^*) = \min\left\{\frac{f(\eta_i^*|\eta_{-i}^{(k)}, \mathbf{y})q_i(\eta_i^{(k)}|\eta_i^*, \eta_{-i}^{(k)})}{f(\eta_i^{(k)}|\eta_{-i}^{(k)}, \mathbf{y})q_i(\eta_i^*|\eta_i^{(k)}, \eta_{-i}^{(k)})}, 1\right\} \quad (9.1)$$

دقت کنید که احتمال پذیرش (۹.۱) به نسبت $f(\eta_i^*|\cdot)/f(\eta_i^{(k)}|\cdot)$ بستگی دارد، پس لازم است که $f(\eta_i|\cdot)$ را تنها تا حد ثابت تناسب بدانید، یعنی ثابت نرمال‌گر در (۹.۱) از صورت و مخرج حذف می‌شود. همچنین اگر چگالی پیشنهادی متقارن باشد، یعنی داشته باشیم $q_i(\eta_i^{(k)}|\eta_i^*, \cdot) = q_i(\eta_i^*|\eta_i^{(k)}, \cdot)$ ، آن‌گاه احتمال پذیرش به مقدار زیر فرو می‌کاهد (متروپلیس و همکاران، ۱۹۵۳).

$$a(\eta_{-i}^{(k)}, \eta_i^{(k)}, \eta_i^*) = \min\left\{\frac{f(\eta_i^*|\eta_{-i}^{(k)}, \mathbf{y})}{f(\eta_i^{(k)}|\eta_{-i}^{(k)}, \mathbf{y})}, 1\right\} \quad (10.1)$$

نمونه‌گیر گیبزی حالتی خاص از M-H به‌ازای $q_i(\eta_i|\eta_i^{(k)}, \eta_{-i}^{(k)}) = f(\eta_i|\eta_{-i}^{(k)}, \mathbf{y})$ و احتمال پذیرش متناظر برابر ۱ است به‌طوری که نامزد گیبزی η^* خودبه‌خود پذیرفته می‌شود. برای بهبود بخشیدن به همگرایی نمونه‌گیری MCMC، چند صورت مختلف از M-H پیشنهاد شده‌اند. برای مطالعه بیشتر می‌توانید به متروپلیس و همکاران، (۱۹۵۳) مراجعه کنید.

۹.۱ مقدمه‌ای بر آمار فضایی

گاهی اوقات در مطالعات محیطی با مشاهداتی سر و کار خواهیم داشت که نسبت به هم مستقل نمی‌باشند، و همبستگی بین آن‌ها ناشی از موقعیت و مکان قرار گرفتن مشاهدات در فضای مورد بررسی می‌باشد، بنابراین به این نوع از مشاهدات داده‌های فضایی^{۴۳} گویند (کرسی، ۱۹۹۳). به دلیل وجود همبستگی فضایی^{۴۴} بین آن‌ها روش‌های آماری برای تحلیل چنین داده‌هایی قابل به‌کارگیری نمی‌باشند و به همین منظور، لازم است ساختار همبستگی داده‌ها را در تحلیل‌شان به کار برد. داده‌های فضایی را می‌توان با توجه به انواع موقعیت‌ها، به سه گروه زیر تقسیم کرد:

داده‌های زمین آماری: این نوع از داده‌ها در موقعیت‌های ثابت و مشخص در ناحیه‌ای پیوسته مشاهده می‌شوند. حال ممکن است که متغیر مورد بررسی گسسته یا پیوسته باشد. به‌عنوان مثال، شکل ۲.۱ نشان‌دهنده‌ی میزان سنگ‌های فیروزه در یک معدن و یا نوعی از گیاه که در یک محدوده‌ی جنگلی است، پس، به آن‌ها داده‌های زمین آماری^{۴۵} گویند. هدف از تحلیل این نوع داده‌ها، پیش‌گویی پاسخ در مکان‌های مشاهده نشده است.



شکل ۲.۱: نقشه محل سنگ‌های فیروزه در کشور هندوستان

داده‌های شبکه‌ای: این نوع از داده‌ها مرتبط به مکان‌های ناحیه‌ای هستند یعنی بین مکان‌های مشاهده شده مکان دیگری وجود ندارد. این مکان‌ها ممکن است به دو حالت منظم و نامنظم باشند. معمولاً در عمل داده‌های شبکه‌ای^{۴۶} بر روی یک شبکه‌ی نامنظم جای دارند، مشابهی مجموعه‌ای از استان‌ها یا دیگر مرزهای منطقه‌ای. مثلاً، تعداد افراد بیمه شده در تمام مناطق استان گیلان می‌تواند نامنظم باشد و یا نرخ رشد اقتصادی هر استان

^{۴۳} Spatial data

^{۴۴} Spatial correlation

^{۴۵} Geostatistical data

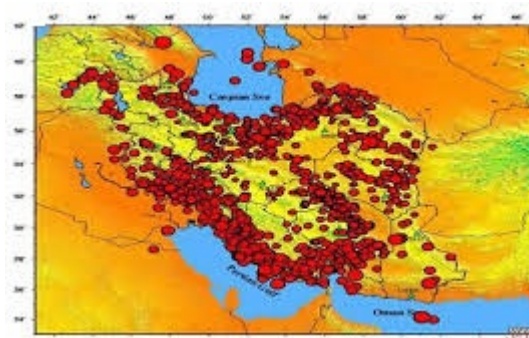
^{۴۶} Lattice data

مثال‌هایی از داده‌های شبکه‌ای می‌باشد. این نوع از داده‌ها به صورت گسسته و پیوسته تعریف شده‌اند. به طور معمول هدف از تحلیل داده‌های فضایی شبکه‌ای، مدل‌بندی احتمالاتی مشاهدات است، که در این پایان‌نامه از این الگو استفاده می‌شود. همانند شکل ۳.۱ خطوط آبی رنگ در نقشه مراکز بیمه استان گیلان را نشان می‌دهد.



شکل ۳.۱: نمودار محل افراد بیمه شده در استان گیلان

داده‌های الگو نقطه‌ای: در این حالت مکان یا موقعیت مشاهده شده خود متغیری تصادفی است. الگوهای نقطه‌ای ۴۷ شامل تعدادی متناهی از مکان‌ها در یک ناحیه‌اند که در آن‌ها یک صفت خاص اندازه‌گیری می‌شود. به طور معمول الگوهای نقطه‌ای به سه دسته به طور کامل تصادفی فضایی ۴۸ CSR، منظم ۴۹ و خوشه‌ای ۵۰ تقسیم می‌شوند. برای مثال، در شکل ۴.۱، موقعیت گونه‌ای از درختان را در یک ناحیه‌ی جنگلی یا موقعیت مراکز زلزله را نشان می‌دهد، که جز این نوع از داده‌های الگو نقطه‌ای به‌شمار می‌روند.



شکل ۴.۱: نمودار نقاط زلزله در سطح کشور

۴۷ Point data

۴۸ Complete spatial randomness

۴۹ Regular

۵۰ Cluster

۱.۹.۱ مدل بندی ساختار وابستگی فضایی

برای تحلیل داده‌های فضایی لازم است یک مدل آماری در نظر گرفته شود. در آمار فضایی به‌طور معمول یک میدان تصادفی به‌عنوان مدل آماری داده‌های فضایی در نظر گرفته می‌شود.

تعریف ۱.۹.۱. (میدان تصادفی)

میدان تصادفی مجموعه‌ای از متغیرهای تصادفی مانند $\{Z(s); s \in D\}$ است، که در آن مجموعه اندیس‌گذار D یک زیر مجموعه از فضای اقلیدسی d بعدی، $d \geq 1$ ، از \mathbb{R}^d است. پس خواهیم داشت:

تعریف ۲.۹.۱. (توابع میانگین، واریانس و کوواریانس میدان تصادفی)

در مورد میدان تصادفی $Z(\cdot)$ میانگین در موقعیت s و کوواریانس در موقعیت‌های s_1 و s_2 به‌ترتیب به‌صورت زیر:

$$\mu(s) = E[Z(s)], \quad s \in D$$

$$C(s_1, s_2) = Cov[Z(s_1), Z(s_2)] =$$

$$E[(Z(s_1) - \mu(s_1))(Z(s_2) - \mu(s_2))], \quad s_1, s_2 \in D$$

تعریف می‌شوند. برای $s = s_2, s = s_1$ ، واریانس میدان تصادفی $Z(\cdot)$ در مکان s به‌صورت

$$Var[Z(s)] = E[Z(s) - \mu(s)]^2 = C(s, s)$$

حاصل می‌شود. هر میدان تصادفی $Z(\cdot)$ را می‌توان به‌صورت

$$Z(s) = \mu(s) + \delta(s), \quad s \in D$$

تجزیه کرد، که در آن $\mu(s)$ تغییرات بزرگ‌مقیاس یا روند و $\delta(\cdot)$ ، فرایند خطا یا تغییرات کوچک‌مقیاس میدان تصادفی نامیده می‌شوند. تغییرات کوچک‌مقیاس ممکن است ناشی از خطای اندازه‌گیری یا تغییرپذیری در درون موقعیت مشاهده باشد و تغییرات بزرگ‌مقیاس ممکن است ناشی از تغییرپذیری بین موقعیت‌های مشاهده شده باشند. به‌طور معمول تحلیل داده‌های فضایی بر اساس مشاهدات نمونه دشوار است، اما برخی ویژگی‌های میدان تصادفی موجب ساده‌سازی مساله خواهند شد.

تعریف ۳.۹.۱. (مانای ذاتی). میدان تصادفی $Z(\cdot)$ را **مانای ذاتی**^{۵۱} گویند هرگاه:

۱. میانگین میدان تصادفی^{۵۲} ثابت یا مستقل از s باشد، به‌طوری که:

^{۵۱} Intrinsic stationary

^{۵۲} Random field

$$E(Z(s)) = \mu, \quad s \in D$$

۲. واریانس مربوط به عبارت $(Z(s_1) - Z(s_2))$ نیز تنها تابعی از فاصله بین دو موقعیت s_1 و s_2 باشد، بنابراین واریانس را به صورت زیر نشان داده می‌شود:

$$\text{Var}(Z(s_1) - Z(s_2)) = 2\gamma(s_1 - s_2), \quad s_1, s_2 \in D \subset \mathbb{R}^d$$

تعریف ۴.۹.۱. (مانای مرتبه دوم یا ضعیف). میدان تصادفی $Z(\cdot)$ را **مانای مرتبه دوم**^{۵۳} یا **ضعیف**^{۵۴} نامند هرگاه:

۱. میانگین میدان تصادفی $Z(\cdot)$ مستقل از s و ثابت باشد. به همین دلیل:

$$E(Z(s)) = \mu, \quad s \in D \subset \mathbb{R}^d$$

۲. کوواریانس $Z(s_1)$ و $Z(s_2)$ در صورتی که فاصله تابعی از موقعیت‌های s_1 و s_2 باشد بنابراین،

$$\text{Cov}(Z(s_1), Z(s_2)) = C(s_1 - s_2), \quad s_1, s_2 \in D \subset \mathbb{R}^d$$

که تحت این نوع از مانایی

$$\text{Var}(Z(s), Z(s)) = \text{Cov}(Z(s), Z(s)) = C(o) = \sigma^2$$

که بنابر مطالب فوق میدان تصادفی به موقعیت فضایی بستگی نخواهد داشت و تغییرپذیری میدان تصادفی در همه جا یکسان می‌باشد.

تعریف ۵.۹.۱. (مانای قوی). میدان تصادفی $Z(\cdot)$ را **مانای قوی** گویند، زمانی که برای تمامی موقعیت‌های s_1, \dots, s_n و تاخیر h ، توزیع توأم $Z(s_1), \dots, Z(s_n)$ نیز مانند توزیع توأم

$$Z(s_1 + h), \dots, Z(s_n + h)$$

باشد. بنابراین

$$(Z(s_1), \dots, Z(s_n)) \stackrel{D}{=} (Z(s_1 + h), \dots, Z(s_n + h))$$

در صورت انتقال موقعیت‌های s_1, \dots, s_n در راستای $\{h \in \mathbb{R}^d\}$ ، توزیع توأم $Z(s_1), \dots, Z(s_n)$ تغییر نخواهد کرد.

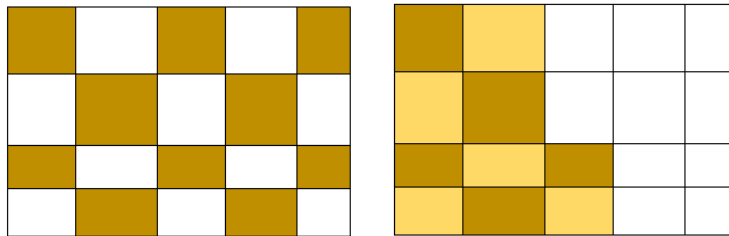
تذکر ۱.۹.۱. برای جمع‌بندی مطالب فوق می‌توان گفت که، مانای قوی شرط کافی برای مانای ضعیف و نیز مانای ضعیف شرطی کافی برای مانای ذاتی خواهد بود که عکس مطالب فوق صحت نخواهد داشت.

^{۵۳} Second

^{۵۴} stationary

۲.۹.۱ وارد کردن ساختار همبستگی فضایی در مدل

همان طوری که می‌دانید همبستگی فضایی یک خاصیت ذاتی است که این ویژگی در داده‌های ناحیه‌ای وجود دارد و باید این همبستگی در مطالعات لحاظ گردد. این نوع از همبستگی می‌تواند به دو صورت مثبت و منفی وجود داشته باشد. زمانی این همبستگی مثبت عمل می‌کند که نواحی نزدیک مشابه هم باشند و بالعکس زمانی این همبستگی منفی است که نواحی نزدیک به صورت غیرمشابه عمل کنند. برای مثال، در همبستگی فضایی مثبت، هنگامی که مقدار اندازه‌گیری شده‌ی متغیر مورد نظر در نواحی بزرگ باشد، نواحی مجاور، متمایل به مقدار بزرگ می‌شوند ولی برای همبستگی منفی این مقدار برای نواحی مجاور متمایل به مقدار کوچک خواهد بود. در شکل ۵.۱ این نوع از همبستگی‌ها را نشان می‌دهیم که نواحی رنگ زده شده یعنی همبستگی مثبت و نواحی به غیر این حالت یعنی همبستگی منفی. اولین گام



شکل ۵.۱: انواع همبستگی فضایی، (شکل سمت راست) همبستگی فضایی مثبت، (شکل سمت چپ) همبستگی فضایی منفی

برای مدل‌بندی فضایی بین نواحی، مشخص کردن روابط بین نواحی است، برای این عمل از همسایگی بهره می‌گیرند.

تعریف ۶.۹.۱. مجموعه نواحی که بیشترین اطلاع در خصوص ناحیه‌ی دلخواه، مانند d را ارائه دهد، مجموعه همسایه‌های ناحیه‌ی d می‌نامند و آن را با $N(d)$ نشان می‌دهند.

۳.۹.۱ ماتریس همسایگی

به صورت‌های مختلف می‌توان همسایگی را داده‌های نواحی بیان کرد. یکی از این رویکردها به این صورت است که ناحیه‌هایی که در فاصله‌ی مشخصی از ناحیه‌ی مورد بررسی قرار دارند، به‌عنوان همسایه‌های آن ناحیه‌ی در نظر گرفت. در واقع با مشخص کردن یک مقدار در آستانه‌ی d همه‌ی ناحیه‌هایی که در فاصله‌ی γ یا کمتر از آن قرار دارند، به‌عنوان همسایه‌ی ناحیه‌ی مورد بررسی در نظر گرفته می‌شود. رویکرد دیگری از همسایگی به صورت راس مشترک جغرافیایی تعریف می‌شود. یعنی، همه‌ی ناحیه‌هایی که دارای راس مشترک با ناحیه‌ی مورد نظر هستند، به‌عنوان همسایه‌های آن ناحیه در نظر گرفته می‌شوند. همچنین همسایگی

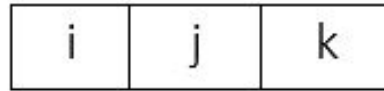
می‌تواند به صورت تبادلات بیمه‌ای بین شهرها نیز تبیین شود، پس آن دسته از شهرهایی که با شهر d از نظر بیمه‌ای در ارتباط اند به عنوان همسایه‌های شهر d در نظر گرفته می‌شوند. همسایگی برای داده‌های ناحیه‌ای به صورت مرز مشترک جغرافیایی نیز است، به طوری که برای هر ناحیه‌ی خاص ناحیه‌هایی که دارای مرز مشترک با ناحیه‌ی مورد نظر هستند، به عنوان همسایه‌ی مرتبه‌ی اول آن ناحیه می‌باشد، البته به غیر از همسایگی مرتبه‌ی اول، برای هر ناحیه‌ی خاص همسایگی مرتبه دوم، سوم و مراتب بالاتر نیز تعریف می‌شود (همانند شکل ۶.۱). حال بعد از اینکه همسایه‌های ناحیه‌ی مورد نظر را تعیین کردیم کافی است برای نشان دادن

			۳			
		۳	۲	۳		
	۳	۲	۱	۲	۳	
۳	۲	۱	d	۱	۲	۳
	۳	۲	۱	۲	۳	
		۳	۲	۳		
			۳			

شکل ۶.۱: همسایگی مرتبه ۱ تا ۳ برای ناحیه‌ی d

ارتباط بین آن‌ها از یک ماتریس همانند W با اعضای صفر و یک استفاده کنیم. برای مثال، اگر ناحیه‌ی مورد نظر d با ناحیه‌ی i همسایه باشد، بنابراین $w_{dj} = w_{jd} = 1$ و در غیر این صورت صفر خواهد بود. به علاوه، برای هر ناحیه‌ی d ، $w_{dd} = 0$ خواهد شد. برای نشان دادن مطالب گفته شده یک مثال را ارائه می‌دهیم (W را برای سه ناحیه در نظر می‌گیریم). ماتریس همسایگی طبق شکل ۷.۱ نشان داده شده است.

$$W = \begin{pmatrix} 0 & w_{ij} & w_{ik} \\ w_{ji} & 0 & w_{jk} \\ w_{ki} & w_{kj} & 0 \end{pmatrix} = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$$



شکل ۷.۱: نواحی مورد بررسی

ماتریس فوق را ماتریس همسایگی^{۵۵} می‌نامند. بنابراین نواحی که مجاور هم می‌باشند یعنی، مرز مشترک داشته باشند، به‌عنوان همسایه در نظر گرفته می‌شود. روش دیگری هم برای مشخص کردن ناحیه‌ی مورد نظر به کار گرفته‌اند که به‌صورت زیر خواهد بود:

$$w_{(dj)} = \begin{cases} 1 & \gamma_{dj} < \delta \\ 0 & o.w. \end{cases}$$

که γ_{dj} فاصله‌ی بین ناحیه‌ی d و همسایه‌ی j است. توجه کنید که همیشه نیازی نیست که از صفر و یک برای انتساب وزن استفاده شود، پس می‌توان از تابع توانی منفی فاصله‌ی بین نواحی نیز استفاده کرد و وزن را اینگونه برای اعضای ماتریس همسایگی لحاظ کرد، در نتیجه α یک مقدار معلوم مثبت و نیز مقدار ۱ یا ۲ را می‌پذیرد.

$$w_{(dj)} = \begin{cases} \gamma_{dj}^{-\alpha} & \alpha > 0 \\ 0 & o.w. \end{cases}$$

از طرفی نیز می‌توان از تابع نمایی منفی^{۵۶} به صورت زیر هم استفاده کرد.

$$w_{(dj)} = \begin{cases} \exp\{-\alpha\gamma_{dj}\} & \alpha > 0 \\ 0 & o.w. \end{cases}$$

یکی از مشکلاتی که روش بالا وجود دارد، دارا بودن وزن برای همه‌ی همسایگی‌ها است که یکسان در نظر می‌گیرند و این منطقی نیست، زیرا، بعضی از همسایگی‌ها احتمال دارند که اطلاعات بیشتری در مورد ناحیه‌ی مورد نظر دهند و نشان‌دهنده‌ی این است که وزن بیشتری را در اختیار دارند. به عبارت دیگر، نیازی نیست ماتریس همسایگی متقارن باشد همین که وزن یکسانی به همسایگی ماتریس W داده شود کافی است. برای انتساب وزن به هر یک از این همسایگی‌ها، می‌توان از عواملی همچون مرز مشترک نواحی، جمعیت ناحیه، مساحت نواحی، طول مشترک نواحی و غیره استفاده کرد.

^{۵۵} Neighborhood matrix

^{۵۶} Negative exponential function

جمع‌بندی

در سال‌های اخیر، روش برآورد کوچک ناحیه‌ای به دلیل نیاز به فراهم آوردن برآوردهای قابل اعتماد آن هم برای وضعیت‌های کم‌نمونه یا بدون نمونه در زیر جوامع، مورد توجه دولت‌ها، بخش‌های اقتصادی و محققان قرار گرفته است. در این فصل مفاهیم اولیه برآورد کوچک ناحیه‌ای و الگوریتم‌های در این موضوع را بررسی کردیم. در فصل آتی، به معرفی مدل‌های کوچک ناحیه‌ای خواهیم پرداخت.

فصل ۲

معرفی مدل‌های کوچک ناحیه‌ای

در ابتدای این فصل، اهمیت برآوردهای غیرمستقیم و روش‌های مبتنی بر مدل‌هایی که در این روش وجود دارند را معرفی می‌کنیم. همچنین آن دسته از مدل‌هایی را که در برآورد نواحی کوچک معروف می‌باشند، شرح می‌دهیم. از جمله این مدل‌ها می‌توان مدل فی-هریوت را نام برد. و در انتهای فصل مدل‌های آمیخته خطی^۱ (LMM) و خطی تعمیم‌یافته‌ی فضایی^۲ (SGLMM) را با ذکر مثال‌هایی بیان می‌کنیم. مطالب ذکر شده در این فصل عمدتاً برگرفته از کتاب برآورد کوچک ناحیه‌ای راتو (۲۰۰۳) و مقاله میلان کارگانیس (۲۰۰۹) است.

۱.۲ برآوردهای کوچک ناحیه‌ای

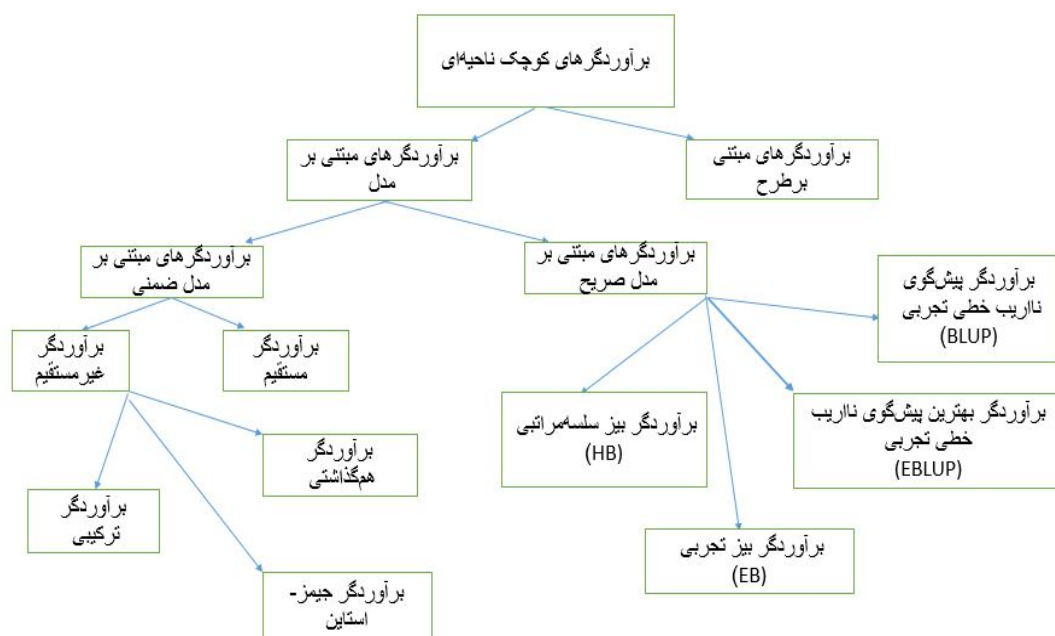
بنابر مطالب بیان شده در فصل اول، نمی‌توان از برآوردهای مستقیم برای برآورد پارامتر نواحی کوچک استفاده کرد. بنابراین، به برآوردهای غیرمستقیم توجه بیشتری می‌شود، زیرا این دسته از برآوردها دقت‌شان نسبت به برآوردهای مستقیم بیشتر است و برآوردهای قابل اطمینان و با کیفیت بهتری را ارائه می‌دهد. از طرفی برای برآورد پارامترهای نواحی کوچک می‌توان از طرح نمونه‌گیری استفاده کرد. اما اگر این طرح‌ها مبنی بر طرح برآوردهای مستقیم باشند، مستلزم هزینه‌های مالی و زمانی خواهند بود.

^۱Linear mixed models

^۲Spatial generalized Linear mixed models

به این منظور برای اینکه بتوان به برآوردگرهای نواحی کوچک دسترسی داشت برآوردگرهای مبتنی بر مدل را جایگزین برآوردگرهای مستقیم بر طرح کرده‌اند. در روش مبتنی بر مدل لزوماً از یک مدل آماری استفاده می‌کنند، بنابراین این مدل برای اینکه بتواند به برآوردگر یک ناحیه‌ی کوچک دسترسی داشته باشد، تنها باید از اطلاعات داده‌های نمونه‌گیری شده‌ی سایر نواحی کوچک و یا دوره‌های زمانی دیگر بهره بگیرد. بنابراین دو مورد از عواملی که کارایی برآوردگرهای ناحیه‌ی کوچک به آن بستگی خواهد داشت را می‌توان یکی متغیرهای پیشگو و دیگری انتخاب یک مدل مناسب برشمرد، زیرا در دسترس بودن داده‌های کمکی مناسب در سطح ناحیه یا در سطح واحد حائز اهمیت است. مدل‌های آماری بنابر ماهیت داده‌ها و داده‌های کمکی قابل دستیابی هستند.

با توجه به مطالب گفته شده، جدول زیر را برای آشنایی بیشتر نسبت به شاخص روش‌های برآورد کوچک ناحیه‌ای مطرح کرده‌ایم.



شکل ۱.۲: نمودار انواع برآوردگرهای کوچک ناحیه‌ای

سعی بر این است، مطالبی درباره‌ی برآوردگرهایی مبتنی بر مدل را ارائه دهیم.

برآوردگرهای مبتنی بر مدل

بنابر مطالب گفته شده، برآوردگرهای مبتنی بر طرح به دلیل دارا بودن واریانس بزرگ برای برآورد پارامتر نواحی کوچک مناسب نمی‌باشند. در روش مبتنی بر مدل، یک توزیع احتمال را روی متغیر تصادفی Y قرار می‌دهند. عملکرد این روش بر مبنای نمونه‌گیری جامعه متناهی از یک جامعه فرضی، متناهی است که با یک مدل تصادفی مشخص شده است. در برآوردگر مبتنی بر مدل برای اینکه بتوان دقت برآوردگر را در برآورد کوچک ناحیه‌ای افزایش داد و همچنین میزان واریانس برآوردگر کوچک ناحیه‌ای را کاهش داد کافی است از اطلاعات سایر نواحی و یا از اطلاعات سایر دوره‌های زمانی استفاده کرد.

در برآوردگرهای غیرمستقیم مبتنی بر مدل، اطلاعات مورد نیاز مربوط به مقادیر متغیرهای مورد مطالعه و کمکی نواحی و یا دوره‌های زمانی است و تمامی این مقادیر با به‌کارگیری از یک مدل ربط (ضمنی یا صریح) در برآورد پارامتر وارد می‌شود. مدل‌های ربط شامل دو نوع می‌باشند، یکی مدل‌های صریح و دیگری مدل ضمنی، برآوردگرهای مبتنی بر مدل شامل دو دسته‌ی کلی برآوردگرهای مبتنی بر مدل ضمنی و مبتنی بر مدل صریح هستند. در مدل ضمنی، به‌طور غیرصریح از ساختار موجود در جامعه استفاده می‌شود و تغییرپذیری نواحی کوچک را شامل نمی‌شود. با توجه به این مدل‌ها، فرض بر این است، که تمامی نواحی نسبت به متغیر مورد بررسی به‌طور مشابه رفتار می‌کند و تغییرپذیری در میان نواحی کوچک وجود ندارد. ولی در مدل‌های صریح از روابط بین متغیرهای مورد بررسی و اطلاعات کمکی و ترکیب آن‌ها با اثرهای تصادفی ناحیه‌ی مشخص، به‌صورت صریح به‌کار گرفته می‌شود.

برآوردگرهای مبتنی بر مدل ضمنی

این نوع از برآوردگرها را برآوردگرهای غیرمستقیم سنتی نیز می‌نامند. همان‌طوری که قبلاً اشاره شد، این نوع از برآوردگرها شامل سه دسته می‌باشند. برای جزئیات بیشتر می‌توانید به کتاب راثو (۲۰۰۳) مراجعه کنید.

۲.۲ معرفی مدل‌های کوچک ناحیه‌ای

به‌منظور اتخاذ اثرهای ناحیه‌ای و آگاهی از تغییرات بین ناحیه‌ای، مدل‌های ضمنی کمتر مورد استفاده قرار می‌گیرند و مدل‌های صریح جایگزین شده‌اند. معمولاً مدل‌های صریح با به‌کارگیری مدل‌های آمیخته‌ی خطی، اثرهای تصادفی ناحیه‌ای را به حساب می‌آورند. در واقع مدل‌های صریح کاراتر از مدل‌های ضمنی هستند زیرا، علاوه بر اتخاذ داده‌های کمکی تغییرات بین ناحیه‌ای را فراتر از آنچه که توسط اطلاعات کمکی موجود در مدل تبیین می‌شوند به حساب می‌آورند. مدل‌های صریح کوچک ناحیه‌ای را می‌توان به دو دسته‌ی زیر تقسیم کرد: «مدل در سطح واحد آماری^۳» و «مدل در سطح ناحیه^۴»

^۳ Unit level model

^۴ Area level model

● مدل در سطح واحد آماری

این مدل‌ها مقادیر واحد متغیر تحت مطالعه را با متغیرهای کمکی خاص آن واحد معین ربط می‌دهند. در این مدل‌ها، واحدهای قابل مشاهده، یک واحد انفرادی در ناحیه هستند و نواحی اغلب اثر خوشه‌ای دارند یعنی مشاهدات واحدهای انفرادی درون یک ناحیه، نسبت به کل مشاهدات به یکدیگر نزدیک‌تراند. این مدل‌ها نیز انواع گوناگونی دارند. هنگامی که مقادیر متغیرهای کمکی واحد جامعه در دسترس باشند و بتوان روی آن‌ها مدل‌سازی کرد، این مدل‌ها به کار گرفته می‌شوند. این نوع مدل‌ها مقادیر واحد انفرادی از متغیر تحت مطالعه را با متغیرهای کمکی آن واحد معین پیوند می‌دهند.

● مدل در سطح ناحیه

مدل‌هایی که برآوردگرهای مستقیم ناحیه‌ی کوچک را با متغیرهای کمکی ناحیه‌ی معین ربط می‌دهند، مدل در سطح ناحیه می‌نامند. زمانی که داده‌های مربوط به سطح واحد در دسترس نباشند، استفاده از این مدل‌ها ضرورت پیدا می‌کند. این مدل‌ها برآوردگرهای مستقیم ناحیه‌ی کوچک را با متغیرهای کمکی سطح ناحیه‌ی مورد نظر پیوند می‌دهند. در این مدل‌ها واحدهای قابل مشاهده برای متغیرهای پاسخ و کمکی ناحیه‌ها می‌باشند و اثرهای تصادفی ناحیه نیز به عنوان مانده‌های مدل خطی معین شده برای متغیر پاسخ هستند. به‌طور کلی یک مدل در سطح ناحیه دو قسمت دارد، یک قسمت مربوط به مدل نمونه‌گیری برای برآوردگرهای مستقیم که مبتنی بر طرح است و قسمت دوم مدل ربطی است که پارامترهای ناحیه‌ی کوچک را به متغیرهای کمکی در سطح ناحیه ربط می‌دهد که معمولاً از نوع رگرسیونی است. البته لازم به ذکر است در صورت وجود اطلاعات کمکی در سطح واحد نیز می‌توان مدل در سطح ناحیه را به کار برد.

مثال ۱.۲.۲. برای مثال، فرض کنید در برآورد نرخ بیکاری ناحیه‌ی کوچک اطلاعات مربوط به تعداد افراد باسواد خانوار به عنوان داده‌ی کمکی در سطح خانوارها (واحد آماری) در ناحیه‌ی مورد نظر در دسترس است. لذا با میانگین‌گیری می‌توان داده‌ی کمکی را در سطح ناحیه‌ی مورد نظر بوجود آورد.

روش‌های سنتی غیرمستقیم بر پایه‌ی مدل‌های ضمنی استوارند که از طریق داده‌های تکمیلی، پیوندی با ناحیه‌های کوچک مرتبط با ناحیه‌ی مربوطه برقرار می‌سازند. حال اگر به مدل‌های کوچک ناحیه‌ای صریح رجوع کنید می‌توانید به این نتیجه برسید که برای تغییرات بین ناحیه‌ای فرصتی خاص ایجاد شده است. کاربست مدل‌های صریح چند مزیت خواهد داشت که عبارت‌اند از

۱. برای پی بردن مدل (هایی) مناسب که به داده‌ها خوب برازش دهد، می‌توان روش‌های تشخیصی مدل را به کار گرفت. این نوع از روش‌های تشخیصی شامل تحلیل مانده‌ها برای کشف انحراف‌ها از مدل مفروض، گزینش متغیرهای کمکی برای مدل، و مباحث تشخیصی حذف برخی مشاهده‌ها برای کشف مشاهده‌ها موثر است.

۲. می‌توان معیارهای دقت ناحیه ویژه را برای هر برآورد کوچک ناحیه‌ای به دست آورد. این مطلب برتری نسبتاً خوبی را به معیارهای کلی که روی ناحیه‌های کوچک و متوسط قرار دارند، و اغلب اوقات با برآوردهای هم‌گذاشتی به کار می‌روند.

۳. مدل‌های آمیخته‌ی خطی و نیز مدل‌های غیرخطی، همچون مدل‌های رگرسیون لجستیک^۵ و مدل‌های خطی تعمیم‌یافته^۶ شامل اثرهای تصادفی ناحیه‌ای هستند که می‌توانند مورد بررسی قرار گیرند. با این مدل‌ها می‌توان ساختار داده‌های پیچیده را همانند وابستگی فضایی را که یکی از اهداف اصلی این پایان‌نامه است، وارد مدل بندی کرد.

۴. پیشرفت‌های روش‌شناختی اخیر برای مدل‌هایی که شامل اثرهای تصادفی هستند راهکاری برای استنباط‌های دقیق کوچک ناحیه‌ای می‌باشد.

دقت کنید که می‌توان مدل‌های گوناگونی برای برآورد کوچک ناحیه‌ای انتخاب کرد. اما مدل باید با مشورت متخصصان موضوع مورد مطالعه و کاربران نهایی این حیطة انتخاب شود. آن‌ها نقش موثری در انتخاب متغیرهای کمکی نیز دارند. لذا پرکاربرد بودن هر روش مدل پایه به در دسترس بودن داده‌های کمکی خوب بستگی دارد. بنابراین باید به جمع‌آوری آن دسته از متغیرهای کمکی که پیش‌گویی‌های خوبی از متغیرهای مورد بررسی دارند، توجه بیشتری شود. اکنون به شرح مختصری از مدل‌های کوچک ناحیه‌ای می‌پردازیم.

۳.۲ مدل پایه‌ای در سطح ناحیه

فرض کنید که $\theta_i = g(\bar{Y}_i)$ به‌ازای تابعی مشخص مانند $g(\cdot)$ با داده‌های کمکی ناحیه ویژه‌ی $\mathbf{z}_i = (z_{1i}, \dots, z_{pi})^T$ از طریق مدل خطی زیر مرتبط است:

$$\theta_i = \mathbf{z}_i^T \boldsymbol{\beta} + b_i v_i, \quad i = 1, \dots, m \quad (1.2)$$

که در آن b_i ها مقادیر ثابت مثبت و معلوم‌اند و $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$ برداری با ابعاد $p \times 1$ از ضریب‌های رگرسیونی است. به‌علاوه، که v_i ها اثرتصادفی ناحیه ویژه هستند، که فرض می‌شوند مستقل و هم‌توزیع (i.i.d) با

$$E_m(v_i) = 0, \quad V_m(v_i) = \sigma_v^2 (\geq 0) \quad (2.2)$$

می‌باشند، که در آن E_m امید ریاضی مدل و V_m واریانس مدل را نشان می‌دهند. این فرض به‌صورت $v_i \stackrel{iid}{\sim} N(0, \sigma_v^2)$ است. نرمال بودن اثرهای تصادفی v_i نیز اغلب اوقات در نظر گرفته می‌شود. پارامتر σ_v^2 معیاری از همگنی ناحیه‌ها پس از به‌حساب آوردن متغیرهای کمکی \mathbf{z}_i است.

^۵ Logistic regression models

^۶ Generalized linear mixed models

در بعضی کاربردها، همه‌ی ناحیه‌ها به عنوان عضو نمونه انتخاب نمی‌شوند. فرض کنید M ناحیه در جامعه داریم و تنها m ناحیه در نمونه انتخاب می‌شوند، $m < M$ مدلی را به مدل (۱.۲) برای جامعه فرض می‌کند. به علاوه فرض کنید که ناحیه‌های نمونه از مدل جامعه پیروی می‌کنند، یعنی آریبی در گزینش ناحیه‌های نمونه‌ای وجود ندارد به طوری که در (۱.۲) برای ناحیه‌های نمونه‌گیری شده صادق است.

برای انجام استنباط‌هایی درباره‌ی میانگین‌های کوچک ناحیه‌ای \bar{Y}_i تحت مدل (۱.۲)، فرض کنید که برآوردهای مستقیم \hat{Y}_i قابل دسترس باشد. حال فرض بر این است که

$$\hat{\theta}_i = g(\hat{Y}_i) = \theta_i + e_i, \quad i = 1, \dots, m \quad (3.2)$$

در آن e_i ها خطاهای نمونه‌گیری مستقل هستند و

$$E_p(e_i|\theta_i) = 0, \quad V_p(e_i|\theta_i) = \Psi_i \quad (4.2)$$

همچنین مرسوم می‌باشد. پس، فرض می‌شود که واریانس‌های نمونه‌گیری، Ψ_i ، معلوم‌اند. فرض‌های بالا ممکن است در برخی کاربردها کاملاً محدود کننده باشد. مثلاً، اگر $g(\cdot)$ تابع غیرخطی و اندازه‌ی نمونه‌ی ناحیه n_i کوچک باشد، برآوردهای مستقیم $\hat{\theta}_i$ ممکن است آریب باشد. اگر ناحیه‌های کوچک طرح نمونه‌گیری را قطع کنند، خطاهای نمونه‌گیری ممکن است وابسته باشد. فرض واریانس معلوم Ψ_i را می‌توان با برآورد Ψ_i از داده‌های نمونه‌ای در سطح واحد آماری و سپس هموارسازی واریانس‌های برآورد شده‌ی $\hat{\Psi}_i$ برای به دست آوردن برآوردی پایدارتر از Ψ_i ، تخفیف داد. نرمال بودن $\hat{\theta}_i$ نیز اغلب فرض می‌شود، اما به واسطه‌ی تاثیر قضیه‌ی حد مرکزی بر $\hat{\theta}_i$ ، این فرض ممکن است به اندازه‌ی فرض نرمال بودن اثرهای تصادفی محدود کننده نباشد. مدل‌های یقینی برای θ_i با قرار دادن $\sigma_v^2 = 0$ به دست می‌آیند، یعنی $\theta_i = \mathbf{z}_i^T \beta$. چنین مدل‌هایی به برآوردهای هم‌گذاشتی می‌انجامند که تغییرات محلی را به غیر از تغییرات بازتاب یافته در متغیرهای کمکی \mathbf{z}_i ، در نظر نمی‌گیرند. با استفاده از ترکیب دو مدل (۱.۲) و (۳.۲) مدل زیر به دست می‌آید:

$$\hat{\theta}_i = \mathbf{z}_i^T \beta + b_i v_i + e_i, \quad i = 1, \dots, m \quad (5.2)$$

با توجه به این که مدل (۵.۲) شامل خطاهای ایجاد شده از طرح e_i و نیز خطاهای مدل v_i است. فرض کنید e_i و v_i مستقل‌اند. مدل (۵.۲) حالت خاصی از مدل آمیخته‌ی خطی است. اگر اندازه‌ی نمونه‌ای n_i در ناحیه‌ی i ام کوچک باشد و θ_i تابعی غیرخطی از مجموع Y_i باشد، حتی اگر برآوردهای مستقیم \hat{Y}_i طرح نااریب باشد، فرض $E_p(e_i|\theta_i) = 0$ در مدل نمونه‌گیری (۳.۲) ممکن است معتبر نباشد. مدلی واقعی‌تر با فرمول

$$\hat{Y}_i = Y_i + e_i^*, \quad i = 1, \dots, m \quad (6.2)$$

^Y Small area means

و $EP(e_i^*|Y_i) = 0$ مطرح می‌شود، یعنی \hat{Y}_i برای مجموع Y_i طرح نارایب است.

در این حالت، مدل‌های نمونه‌گیری و پیونددهنده با هم جور نیستند. در نتیجه نمی‌توان مدل (۶.۲) را با مدل پیونددهنده‌ی (۱.۲) ترکیب کرد تا مدل آمیخته‌ی خطی به شکل مدل (۳.۲) به دست آید. ضمن این که نتایج استاندارد نظریه‌ی مدل‌های آمیخته‌ی خطی، کاربرد ندارد (رائو، ۲۰۰۳). روش پیشنهادی آماردانان‌ها برای حل این مشکل، با استفاده از رویکرد بیزی سلسله‌مراتبی^۸ HB برای اداره کردن مدل‌های ناجور نمونه‌گیری و پیونددهنده است.

مثال ۱.۳.۲. (شمار افراد فقیر). مدل پایه‌ای در سطح ناحیه (۵.۲) اخیراً برای تولید برآوردهای شهرستانی مدل پایه از کودکان فقیر در سن تحصیل در ایالت‌های متحده به کار رفته است (شورای پژوهش ملی، ۲۰۰۰). با استفاده از این برآوردها، وزارت آموزش و پرورش ایالت‌های متحده سالانه بیش از ۷ میلیارد دلار از بودجه‌ی عمومی را به شهرستان‌ها تخصیص می‌دهد، و سپس ایالت‌ها این بودجه را در بین منطقه‌های مدرسه‌ای توزیع می‌کنند. در گذشته، بودجه‌ها بر اساس شمارهای برآوردشده از سرشماری قبلی تخصیص می‌یافتند، اما شمار افراد فقیر در طی زمان به‌طور معنی داری تغییر کرده است.

در این کاربرد، $\theta_i = \log(Y_i)$ ، که در آن Y_i شمار واقعی افراد فقیر شهرستان (کوچک ناحیه‌ی) i ام است. برآوردهای مستقیم \hat{Y}_i به‌صورت میانگین موزون سه‌ساله‌ی کودکان فقیر در سن تحصیل (زیر ۱۸ سال) از نشریه‌ی تکمیلی ماه مارس آمارگیری جمعیت جاری (CPS) به دست آمدند. متغیرهای پیشگو در سطح ناحیه، z_i ، از داده‌ی ثبتي به دست آمده است. شهرستان‌های با نمونه‌ی (CPS) ولی، بدون کودکان فقیر در سن تحصیل کنار گذاشته شده است، زیرا $\hat{Y}_i = 0$ و $\log \hat{Y}_i = -\infty$ می‌باشد. برای شهرستان‌هایی که در نمونه‌ی (CPS) انتخاب نشدند، تنها متغیرهای کمکی z_i موجودند.

۱.۳.۲ گسترش‌های مدل در سطح ناحیه (A)

شکل (۱.۲) بعضی از مدل‌های معروف آماری را تا حدودی نشان می‌دهد. اما این مدل‌ها تنها بخشی از کل مطالبی است که معرفی نموده‌ایم، برای آشنایی بیشتر می‌توانید به کتاب رائو (۲۰۰۳) مراجعه کنید. گسترش مدل در سطح ناحیه شامل پنج دسته کلی که عبارت‌اند از؛ مدل فی-هریوت، مدل فی-هریوت چندمتغیره، مدل با خطاهای نمونه‌گیری همبسته، مدل‌های سری زمانی و مقطعی و مدل‌های فضایی است که در بخش‌های زیر به‌طور مفصل شرح خواهیم داد.

مدل فی-هریوت

مدل فی-هریوت یکی از مدل‌های موجود در سطح ناحیه است، این مدل نخستین بار توسط فی و هریوت در سال (۱۹۷۹) مطرح شد. ترکیب این مدل شامل مدل پیونددهنده و مدل

^۸Hierarchical bayes

نمونه‌گیری می‌باشد که در آن مدل نمونه‌گیری نشان‌دهنده‌ی ارتباط بین پارامتر مورد نظر کوچک ناحیه‌ای و برآوردگر مستقیم مرتبط به آن است، همچنین مدل پیونددهنده به کمک یک مدل رگرسیونی در بین پارامتر کوچک ناحیه‌ای و متغیرهای کمکی در سطح ناحیه پیوند برقرار می‌کند. حال اگر نواحی بزرگ مورد نظر دارای طرح نمونه‌گیری متناسب باشد و شامل M نواحی کوچک باشد، مدل پیونددهنده درباره‌ی هریک از این M نواحی صدق می‌کند (پراتسی و سالواتی، ۲۰۰۶). بنابراین تنها m نواحی کوچک در میان سایر نواحی نمونه‌گیری قرار می‌گیرد، در این صورت فرض بر این است که اریبی انتخاب نمونه وجود ندارد و این m نواحی کوچک از یک مدل پیونددهنده‌ی موجود در جامعه پیروی می‌کند. و برای آشنایی با این مدل می‌توان آن را به‌صورت زیر نمایش داد:

$$\theta_i = \mathbf{z}_i^T \boldsymbol{\beta} + b_i v_i; \quad i = 1, \dots, m \quad (7.2)$$

که θ_i پارامتر مورد نظر ناحیه‌ی i ام و $\mathbf{z}_i = (z_{i1}, z_{i2}, \dots, z_{ip})$ بردار p بعدی از متغیرهای کمکی ناحیه‌ی i ام است و b_i یک مقدار ثابت معلوم و $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)$ بردار $1 \times p$ از ضرایب رگرسیونی می‌باشد و عناصر آن نیز نامعلوم و ثابت هستند. اثر تصادفی ناحیه‌ی v_i است که مستقل و هم‌توزیع فرض می‌شود (i.i.d) بنابراین،

$$E(v_i) = 0, \quad V(v_i) = \delta_v^2.$$

توزیع v_i ها را نرمال در نظر می‌گیرند، پس $v_i \stackrel{iid}{\sim} N(0, \delta_v^2)$ و v_i ها و e_i ها نیز مستقل از هم هستند. مدل نمونه‌گیری هم به صورت زیر است:

$$\hat{\theta}_i = g(\hat{Y}_i) = \theta_i + e_i; \quad i = 1, \dots, m. \quad (8.2)$$

در مورد مدل فوق می‌توان گفت که $\hat{\theta}_i$ برآوردگر مستقیم مرتبط به θ_i و e_i خطای نمونه‌گیری و از هم مستقل اند لذا بنابر $E_p(e_i | \theta_i) = 0$ و $V_p(e_i | \theta_i) = \Psi_i$ می‌توان گفت، برآوردگر مستقیم θ_i طرح نارایب است و به‌صورت تابعی از میانگین طرحی ناحیه‌ی کوچک i ام یعنی \hat{Y}_i در نظر گرفته می‌شود. برای میانگین طرحی کوچک ناحیه‌ای در ناحیه‌ی i ام معمولاً از برآوردگر هورویتز-تامسون استفاده می‌کنند که این نوع از برآوردگر دارای ویژگی طرح ناراییبی است.

Ψ_i ها در واقع همان واریانس‌های نمونه‌گیری اند که معلوم فرض می‌شوند. بنابر شرایط مطرح شده، محدودیت در بعضی مسائل کاربردی وجود دارد. زمانی $\hat{\theta}_i$ ها برای θ_i ها می‌تواند طرح اریب باشد، که به عنوان مثال، اگر g یک تابع غیرخطی باشد، آن‌گاه n_i کوچک می‌شود. اگر خطاهای نمونه‌گیری (منظور همان e_i ها) نسبت به هم وابسته باشند و نواحی کوچک، طرح نمونه‌گیری را قطع کنند، بایستی از سایر مدل‌ها استفاده کرد. اشاره به این نکته لازم است که با فرض معلوم بودن Ψ_i ها و نرمال بودن برآوردگر مستقیم، $\hat{\theta}_i$ ها محدودکننده نمی‌باشند، زیرا به جای e_i ها می‌توان از برآورد هموارشده‌ی آن‌ها که از طریق مشاهدات نمونه‌ای حاصل می‌شوند، استفاده کرد و علاوه بر این برآوردگرهای $\hat{\theta}_i$ ها یک ترکیب خطی از \hat{Y}_i ها می‌باشند.

بنابر قضیه‌ی حد مرکزی می‌توان نتیجه گرفت که \hat{Y}_i ها به دنبال $\hat{\theta}_i$ ها هستند و از توزیع نرمال پیروی می‌کنند که هر کدام از $\hat{\theta}_i$ ها به صورت حاشیه‌ای دارای توزیع نرمال خواهند بود. در انتها با ترکیب مدل‌های (۷.۲) و (۸.۲) می‌توان مدل نهایی را به صورت زیر نمایش داد:

$$\hat{\theta}_i = \mathbf{z}_i^T \boldsymbol{\beta} + b_i v_i + e_i, \quad i = 1, \dots, m. \quad (9.2)$$

مدل گفته شده در بالا، یک صورت خاص از یک مدل آمیخته خطی می‌باشد. اگر $E_p(e_i|\theta_i) = 0$ در مدل (۸.۲) برقرار نباشد، برآوردگر مستقیم $\hat{\theta}_i$ برای θ_i طرح اریب می‌شود، همانند زمانی که g یک تابع غیرخطی از یک مجموعه در ناحیه‌ی کوچک i ام باشد و n_i کوچک باشد، در این صورت می‌توان گفت که اگر برآوردگر مستقیم \hat{Y}_i برای Y_i طرح نارایب باشد آن‌گاه $\hat{\theta}_i$ برای θ_i طرح اریب خواهد بود. و به همین دلیل مدل نمونه‌گیری زیر را به جای مدل (۸.۲) به کار می‌برند. یعنی:

$$\hat{Y}_i = Y_i + e_i^* \quad (10.2)$$

اگر $E_p(e_i^*|\theta_i) = 0$ برقرار باشد، برآوردگر مستقیم \hat{Y}_i برای Y_i یک برآوردگر طرح نارایب خواهد شد. به همین دلیل مدل (۱۰.۲) یک مدل نمونه‌گیری است و با مدل (۷.۲) که یک مدل پیونددهنده است قادر به ترکیب شدن نمی‌باشند پس، از بوجد آوردن مدل ترکیبی (۹.۲) معذور خواهیم بود. در نتیجه یک مدل آمیخته‌ی خطی و نیز امکان دسترسی به پیشگوی EBLUP را نخواهیم داشت و این به منزله استفاده از مدل سلسله مراتبی HB خواهد بود.

مدل فی-هریوت چندمتغیره

در نظر بگیرید که بردار $1 \times r$ از برآوردگرهای آمارگیری $\hat{\theta}_i = (\hat{\theta}_{i1}, \dots, \hat{\theta}_{ir})^T$ باشد، پس خواهیم داشت:

$$\hat{\theta}_i = \boldsymbol{\theta}_i + \mathbf{e}_i, \quad i = 1, \dots, m \quad (11.2)$$

و $\boldsymbol{\theta}_i = (\theta_{i1}, \dots, \theta_{ir})^T$ است که $\theta_{ij} = g_i(\bar{Y}_{ij})$ به‌ازای $j = 1, \dots, r$ برقرار می‌باشد. خطاهای نمونه‌گیری $\mathbf{e}_i = (e_{i1}, \dots, e_{ir})^T$ برداری با متغیرهای مستقل نرمال r بعدی و دارای توزیع $(0, \boldsymbol{\Psi}_i)$ با میانگین ۰ و ماتریس‌های کوواریانس معلوم $\boldsymbol{\Psi}_i$ به شرط $\boldsymbol{\theta}_i$ است. در اینجا $N_r(0, \boldsymbol{\Psi}_i)$ بردار r بعدی و \bar{Y}_{ij} میانگین ناحیه‌ی کوچک i ام برای مشخصه‌ی j ام است. همچنین فرض کنید که $\boldsymbol{\theta}_i$ با داده‌های کمکی ناحیه ویژه‌ی \mathbf{z}_{ij} از طریق مدل خطی

$$\boldsymbol{\theta}_i = \mathbf{Z}_i \boldsymbol{\beta} + \mathbf{v}_i, \quad i = 1, \dots, m \quad (12.2)$$

ارتباط دارد که در آن اثرهای تصادفی ناحیه ویژه‌ی \mathbf{v}_i متغیرهای $N_r(0, \boldsymbol{\Sigma}_v)$ و مستقل‌اند، ماتریس \mathbf{Z}_i از مرتبه‌ی $r \times rp$ است که سطر i ام آن با $(0^T, \dots, 0^T, \mathbf{z}_{ij}^T, 0^T, \dots, 0^T)$ و $\boldsymbol{\beta}$ بردار rp

بعدی ضریب‌های رگرسیونی است. لذا بردار صفر دارای ابعاد $1 \times p$ و Z_{ij} در مکان j ام از بردار سطری (سطر j ام) واقع است. با ترکیب مدل‌های (۱۱.۲) و (۱۲.۲)، مدل آمیخته‌ی خطی چند متغیره‌ی

$$\hat{\theta}_i = Z_i\beta + v_i + e_i \quad (13.2)$$

به دست می‌آید. مدل (۱۳.۲) نوعی از مدل فی-هریوت (۵.۲) است زمانی که $b_i = 1$ باشد. فی (۱۹۸۷) و داتا و همکاران (۱۹۹۱) گسترش چند متغیره‌ی (۷.۲) را پیشنهاد کردند و نشان دادند این مدل می‌تواند به برآوردگرهای کاراتر از میانگین‌های کوچک ناحیه‌ای \bar{Y}_{ij} منجر شود زیرا، برخلاف مدل تک متغیره‌ی (۵.۲)، از همبستگی بین مولفه‌های $\hat{\theta}_i$ بهره می‌گیرد. همچنین داتا و همکاران (۱۹۹۱) مدل چندمتغیره‌ی (۱۳.۲) را برای برآورد میانه‌ی درآمد جاری مربوط به خانوارهای چهارنفری در هر یک از ایالت‌های متحده به کار بستند.

مدل با خطاهای نمونه‌گیری همبسته

یک گسترش طبیعی از مدل فی-هریوت (۵.۲) عبارت است از مدلی با خطاهای نمونه‌گیری همبسته e_i . فرض کنید $\theta = (\theta_1, \dots, \theta_m)^T$ ، $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_m)^T$ و $e = (e_1, \dots, e_m)^T$ برقرار باشد، بنابراین مدل فی-هریوت زیر

$$\hat{\theta} = \theta + e, \quad (14.2)$$

دارای توزیعی با $e|\theta \sim N_m(0, \Psi)$ ، که در آن ماتریس کوواریانس دارای خطاهای نمونه‌گیری $\Psi = (\psi_{ij})$ است. بعد از ترکیب شدن مدل (۱۴.۲) با مدل (۱.۲) برای θ_i ها تعمیمی از مدل فی-هریوت (۵.۲) را به دست می‌آورند. لذا اگر به ازای تمامی i ها در مدل (۵.۲) $b_i = 1$ ، جای گذاری شود، بنابراین مدل ترکیبی به دست می‌آید که می‌توان آن را به صورت زیر نوشت:

$$\hat{\theta} = X\beta + v + e \quad (15.2)$$

با توجه به مدل فوق $v = (v_1, \dots, v_m)^T$ و x ماتریسی با بعد $m \times p$ با سطر i ام در x_i^T است. اما در عمل، به جای Ψ یک برآوردگر آمارگیری $\hat{\Psi}$ یا برآوردگری هموار شده جانشین می‌شود، اما در بعضی موارد تغییرپذیری مربوط به این برآوردگر اغلب نادیده گرفته می‌شود.

مثال ۲.۳.۲. (کم‌شماری سرشماری). در ساختار برآورد مقدار کم‌شماری در سرشماری ۱۹۹۰ ایالت‌های متحده، جامعه به ۳۵۷ پساتبقه مرکب از ۵۱ گروه پساتبقه‌ای تقسیم شد و هر یک از آن‌ها به ۷ دسته‌ی سنی-جنسی تقسیم شدند. تعداد ۵۱ گروه پساتبقه‌ای بر اساس نژاد یا قومیت، تصرف مسکن (مالک، مستاجر)، نوع ناحیه و منطقه تعریف شدند. برآوردهای

دوگان سامانه $\hat{\theta}_i$ برای هر پساطبقه $i = 1, \dots, 357$ با استفاده از داده‌های حاصل از آمارگیری پساشمارشی (PES) سال ۱۹۹۰ به دست آمدند. در اینجا $\hat{\theta}_i$ ضریب تعدیل برآورد شده برای پساطبقه‌ی i است. مدلی به شکل (۱۵.۲) برای به دست آوردن برآوردهای هموارشده‌ی ضریب‌های تعدیل θ_i به کار گرفته شد (ایساک‌ی و همکاران، ۲۰۰۰).

مدل‌های سری زمانی و مقطعی

سری زمانی به مجموعه‌ای از مشاهدات کمی گفته می‌شود که در فواصل زمانی و به صورت متوالی اندازه‌گیری می‌شود. به بیان دیگر، سری زمانی مجموعه‌ای از مشاهدات یک متغیر است که در نقاط گسسته‌ای از زمان، که معمولاً فاصله‌های مساوی دارند، اندازه‌گیری و بر حسب زمان مرتب شده‌اند. بنابراین یک سری زمانی از مشاهده یک پدیده در طول زمان به دست می‌آید. بنابراین، بسیاری از آمارگیری‌های نمونه‌ای با جانشین کردن پاره‌ای از عنصرهای نمونه‌ای، در طی زمان تکرار می‌شوند. مثلاً، در آمارگیری جمعیت (CPS) ماهانه‌ی ایالت‌های متحده یک خانوار انفرادی به مدت چهار ماه پیاپی در نمونه باقی می‌ماند. سپس برای هشت ماه بعدی از نمونه خارج می‌شود، و برای چهار ماه بعدی دیگر به نمونه باز می‌گردد. در آمارگیری نیروی کار (LFS) ماهانه‌ی کانادا، هر خانوار انفرادی به مدت شش ماه پیاپی در نمونه می‌ماند و سپس از نمونه بیرون می‌رود. در مورد چنین آمارگیری‌های مکرر با وام گرفتن قرضی هم از ناحیه‌های کوچک و هم از دوره‌های زمانی، می‌توان افزایش قابل ملاحظه‌ای در کارایی برآوردها به دست آورد.

رائو و یو (۱۹۹۲، ۱۹۹۴) گسترشی از مدل پایه‌ای فی-هریوت (۵.۲) را برای اداره‌ی داده‌های سری زمانی و مقطعی پیشنهاد کردند. مدل آن‌ها شامل یک مدل خطای نمونه‌گیری

$$\hat{\theta}_{it} = \theta_{it} + e_{it}, \quad t = 1, \dots, T; i = 1, \dots, m \quad (16.2)$$

و مدل پیوند دهنده‌ی

$$\hat{\theta}_{it} = \mathbf{z}_{it}^T \beta + v_i + u_{it} \quad (17.2)$$

بود. در اینجا $\hat{\theta}_{it}$ برآوردها برای کوچک ناحیه‌ای i در زمان t است، پس، $\theta_{it} = g(\bar{Y}_{it})$ تابعی از میانگین کوچک ناحیه‌ای \bar{Y}_{it} است، e_{it} ‌ها عبارت‌اند از، خطاهای نمونه‌گیری با توزیع شرطی نرمال به شرط θ_{it} ‌ها که دارای میانگین‌های صفر و یک ماتریس کوواریانس قطری بلوکی معلوم Ψ با بلوک‌های Ψ_i ‌اند، و برداری از متغیرهای کمکی ناحیه ویژه است که برخی از آن‌ها، مثلاً داده‌های ثبتي، ممکن است با زمان تغییر کنند. به علاوه، $v_i \stackrel{iid}{\sim} N(0, \sigma_v^2)$ و فرض می‌شود که u_{it} از یک فرآیند اتورگرسیو مرتبه‌ی اول مشترک برای هر i پیروی می‌کند. یعنی،

$$u_{it} = \rho u_{i,t-1} + \varepsilon_{it}, \quad |\rho| < 1 \quad (18.2)$$

که در آن $\varepsilon_{it} \stackrel{iid}{\sim} N(0, \sigma^2)$. فرض می‌شود که خطاهای $\{e_{it}\}$ ، $\{v_i\}$ و $\{\varepsilon_{it}\}$ نیز مستقل از یکدیگر هستند. مدل‌هایی که به شکل (۱۷.۲) و (۱۸.۲) با نادیده گرفتن خطای نمونه‌گیری، به‌طور گسترده مقالات اقتصادسنجی نیز به‌کار رفته‌اند (اندرسون و هسیاوا، ۱۹۸۱).

مدل (۱۷.۲) برای θ_{it} ‌ها هم به اثرهای ناحیه ویژه‌ی v_i و هم به اثرهای (ناحیه \times زمان) ویژه‌ی u_{it} که به ازای هر i در طول زمان همبسته‌اند، و وابستگی دارند. همچنین می‌توانید (۱۷.۲) را به‌صورت مدلی با تاخیر توزیع‌شده بیان کنید:

$$\theta_{it} = \rho\theta_{i,t-1} + (\mathbf{z}_{it} - \rho\mathbf{z}_{i,t-1})^T \boldsymbol{\beta} + (1 - \rho)v_i + \varepsilon_{it} \quad (19.2)$$

صورت بدیل (۱۹.۲) متغیر θ_{it} را به میانگین دوره‌ی قبلی $\theta_{i,t-1}$ ، مقدارهای متغیرهای کمکی در نقطه‌های زمانی t و $(t-1)$ ، اثرهای تصادفی کوچک ناحیه‌ای v_i ، و اثرهای (ناحیه \times زمان) ε_{it} ارتباط می‌دهد. مدل‌های پیچیده‌تر از (۱۸.۲) برای u_{it} را می‌توان با فرض یک فرآیند میانگین متحرک اتورگرسیو (ARIMA) فرمول‌بندی کرد، اما مقدار افزایش حاصل در کارایی نسبت به (۱۸.۲) نامحتمل است که معنی‌دار باشد.

گوش و نانجیا (۱۹۹۳) و گوش، نانجیا و کیم (۱۹۹۶) نیز یک مدل سری زمانی مقطعی برای برآورد کوچک ناحیه‌ای پیشنهاد کردند. مدل آن‌ها به‌صورت زیر است:

$$\hat{\theta}_{it} | \theta_{it} \stackrel{ind}{\sim} N(\theta_{it}, \psi_{it}) \quad (20.2)$$

$$\theta_{it} | \boldsymbol{\alpha}_t \stackrel{ind}{\sim} N(F_{it}\boldsymbol{\beta} + \mathbf{w}_{it}^T \boldsymbol{\alpha}_t, \sigma_t^2) \quad (21.2)$$

و

$$\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1} \stackrel{ind}{\sim} N_r(\mathbf{H}_t \boldsymbol{\alpha}_{t-1}, \Delta) \quad (22.2)$$

در اینجا \mathbf{z}_{it} و \mathbf{w}_{it} بردارهای متغیرهای کمکی ناحیه ویژه‌اند، ψ_{it} ‌ها واریانس نمونه‌گیری‌اند که معلوم فرض می‌شوند، برداری $\boldsymbol{\alpha}_t$ از اثرهای تصادفی زمان ویژه، و \mathbf{H}_t یک ماتریس $r \times r$ معلوم است. فضای وضعیت (۲۲.۲) در حالت تک متغیره ($r = 1$) با $H_t = 1$ به مدل مشهور گام‌زدن تصادفی فرو می‌کاهد. مدل بالا مبتلا به دو محدودیت است: (یک) فرض می‌شود که برآوردهای مستقیم $\hat{\theta}_{it}$ برای هر i در طول زمان مستقل باشند. این فرض در ساختار آمارگیری مکرر با نمونه‌های هم‌پوش، از قبیل CPS و LFS، واقع‌گرایانه نیست. (دو) اثرهای تصادفی ناحیه ویژه در مدل گنجانده نشده‌اند که منجر به ترنجش زیاد برآوردهای کوچک ناحیه‌ای همانند برآوردهای هم‌گذاشتی می‌شود.

داتا و همکاران (۲۰۰۲) و یو (۱۹۹۹) مدل‌های نمونه‌گیری و پیونددهنده‌ی رآئو-یو (۱۶.۲) و (۱۷.۲) را به‌کار بستند، اما به‌جای AR(1) مدل (۱۸.۲)، برای u_{it} یک مدل گام‌زدن تصادفی $u_{it} = u_{i,t-1} + \varepsilon_{it}$ را گذاشتند که با رابطه‌ی (۱۸.۲) به ازای $\rho = 1$ داده می‌شود. داتا و همکاران (۱۹۹۹) مدلی مشابه را در نظر گرفتند اما جمله‌هایی دیگر به مدل پیونددهنده افزودند تا تغییرات فصلی را در کاربرد آن‌ها بازتاب دهد.

پففرمن و بورک (۱۹۹۰) مدلی کلی شامل اثرهای تصادفی (ناحیه \times زمان) ویژه را پیشنهاد کردند. مدل آن‌ها به صورت زیر است:

$$\hat{\theta}_{it} = \theta_{it} + e_{it} \quad (23.2)$$

$$\theta_{it} = \mathbf{z}_{it}^T \beta_{it} \quad (24.2)$$

که در آن ضریب‌های $\beta_{it} = (\beta_{it}, \dots, \beta_{itp})^T$ به طور مقطعی و در طی زمان مجاز به تغییر هستند. فرض می‌شود که خطاهای e_{it} برای هر ناحیه‌ی i به طور پیاپی ناهمبسته با میانگین صفر و واریانس ψ_{it} باشند. تغییرات β_{it} در طی زمان با مدل فضای وضعیت زیر مشخص می‌شود:

$$\begin{bmatrix} \beta_{itj} \\ \beta_{it} \end{bmatrix} = \mathbf{T}_j \begin{bmatrix} \beta_{i,t-1,j} \\ \beta_{ij} \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} v_{itj}, \quad j = 0, 1, \dots, p \quad (25.2)$$

در اینجا β_{ij} ‌ها ضریب‌های ثابت‌اند، \mathbf{T}_j یک ماتریس معلوم 2×2 با سطر دوم $(0, 1)$ است، و خطاهای مدل $\{v_{itj}\}$ به ازای هر i در طی زمان ناهمبسته‌اند و دارای میانگین صفر و واریانس‌های $E_m(v_{itj}, v_{itl}) = \sigma_{vj}$ ، به ازای $j, l = 0, 1, \dots, p$ هستند.

فرمول‌بندی (۲۵.۲) چند مدل سودمند را در بردارد. اول، گزینه‌ی $\mathbf{T}_j = \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix}$ مدل

مشهور رگرسیونی با ضریب تصادفی $\beta_{itj} = \beta_{ij} + v_{itj}$ را می‌دهد. مدل آشنای گام‌زدن تصادفی $\beta_{itj} = \beta_{i,t-1,j} + v_{itj}$ از انتخاب $\mathbf{T}_j = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ حاصل می‌شود. در این حالت ضریب β_{ij} اضافی

است و باید حذف شود به طوری که $\mathbf{T}_j = 1$. گزینه‌ی $\mathbf{T}_j = \begin{bmatrix} \rho & 1 - \rho \\ 0 & 1 \end{bmatrix}$ یک مدل AR(1) را

می‌دهد: $\beta_{itj} - \beta_{ij} = \rho(\beta_{i,t-1,j} - \beta_{ij} + v_{itj})$. مدل فضای وضعیت (۲۵.۲) کاملاً کلی است، اما فرض خطاهای نمونه‌گیری به طور پیاپی ناهمبسته‌ی e_{it} در (۲۳.۲) در ساختار آمارگیری‌های مکرر با نمونه‌های هم‌پوش، محدود کننده است.

مثال ۳.۳.۲. (نرخ‌های بیکاری ایالت‌های متحده). داتا و همکاران (۱۹۹۹) مدل خود را برای برآورد نرخ‌های بیکاری ماهانه‌ی مربوط به ۴۹ ایالت آمریکا (با کنار گذاشتن ایالت نیویورک) و منطقه کلمبیا (شهر واشنگتن) ($m = 50$) که توسط انواع سازمان‌های فدرال به منظور تخصیص بودجه‌ها و تنظیم سیاست‌ها مورد استفاده قرار می‌گیرند، به کار بستند. آنان دوره‌ی زمانی ژانویه‌ی ۱۹۸۵ تا دسامبر ۱۹۸۸ را در نظر گرفتند و برآوردهای CPS را به عنوان $\hat{\theta}_{it}$ و نرخ ادعای بیمه‌ی بیکاری UI (درصد کارکنان بیکاری که از بین شاغلان غیر کشاورزی ادعای مزایای بیکاری می‌کنند) را به عنوان داده‌های کمکی، \mathbf{z}_{it} به کار گرفتند. تغییرات فصلی در نرخ‌های بیکاری ماهانه از راه وارد کردن اثرهای تصادفی ماه و سال در مدل به حساب آورده شدند.

مدل فضایی

مدل پایه‌ای فی-هریوت (۵.۲) اثرهای کوچک ناحیه‌ای v_i را به صورت (i.i.d) فرض می‌کند، اما در برخی کاربردها شاید امتحان مدل‌هایی که همبستگی بین v_i ها را مجاز می‌شمارند واقع‌گرایانه‌تر باشد. مدل‌های فضایی برای v_i ها وقتی به کار می‌روند که بتوان ناحیه‌های «همسایه»ی هر ناحیه را تعریف کرد. چنین مدل‌هایی همبستگی‌های بین v_i ها، مثلاً، همبستگی‌هایی را که در ساختار برآورد نرخ‌های محلی بیماری و مرگ‌ومیر، به هم‌جواری جغرافیایی وابسته‌اند، القا می‌کنند. کرسی (۱۹۹۱) از یک مدل فضایی برای برآورد کوچک ناحیه‌ای در ساختار سرشماری استفاده کرد.

۴.۲ مدل پایه‌ای در سطح واحد آماری

فرض کنید که داده‌های کمکی واحد ویژه $\mathbf{x}_{it} = \{x_{ij1}, \dots, x_{ijp}\}$ برای هر عنصر j ی در جامعه در هر کوچک ناحیه‌ی i در دسترس‌اند. اغلب کافی است فرض کنید که تنها میانگین‌های جامعه‌ای \bar{X}_i معلوم‌اند. این اطلاعات از سرشماری و اطلاعات ثبتي به دست می‌آیند. به علاوه فرض می‌شود متغیر مورد نظر، y_{ij} از طریق یک مدل رگرسیون خطی با خطاهای آشيانی یک لایه‌ای e_{ij} به x_{ij} به صورت زیر مربوط می‌شود:

$$y_{ij} = \mathbf{x}_i^T \boldsymbol{\beta} + v_i + e_{ij}, \quad j = 1, \dots, N_i, i = 1, \dots, m \quad (26.2)$$

که در اینجا فرض بر این است که اثرهای ناحیه‌ی ویژه v_i متغیرهای تصادفی (i.i.d) هستند که در شرایط مدل زیر صدق می‌کنند،

$$E_m(v_i) = 0, \quad V_m(v_i) = \sigma_v^2 (\geq 0)$$

آن‌گاه $e_{ij} = k_{ij} \tilde{e}_{ij}$ با ثابت‌های معلوم k_{ij} و \tilde{e}_{ij} متغیرهای تصادفی (i.i.d) مستقل از v_i ها هستند و

$$E_m(\tilde{e}_{ij}) = 0, \quad V_m(\tilde{e}_{ij}) = \sigma_e^2. \quad (27.2)$$

افزون بر آن، نرمال بودن v_i ها و e_{ij} ها اغلب فرض می‌شود. پارامترهای مورد نظر میانگین‌های \bar{Y}_i یا مجموع‌های Y_i کوچک ناحیه‌ای‌اند. مدل‌های رگرسیونی استاندارد با قرار دادن $\sigma_v^2 = 0$ یا به‌طور هم‌ارز $v_i = 0$ در (۲۶.۲) به دست می‌آیند. این نوع مدل‌ها به برآوردگرهایی از نوع هم‌گذاشتی می‌انجامند.

فرض کنید که نمونه‌ی s_i با اندازه‌ی n_i از N_i واحد درون ناحیه i ام ($i = 1, \dots, m$) گرفته می‌شود و مقدارهای نمونه‌ای نیز از مدل (۲۶.۲) پیروی می‌کنند. فرض اخیر تحت نمونه‌گیری تصادفی ساده از هر ناحیه یا به‌طور کلی‌تر برای طرح‌های نمونه‌گیری بی‌که اطلاعات کمکی x_{ij}

را در گزینش نمونه‌های s_i به کار می‌گیرند، صادق است. برای درک بهتر این نکته، (۲۶.۲)

$$\mathbf{y}_i^P = \mathbf{X}_i^P \boldsymbol{\beta} + v_i \mathbf{1}_i^P + \mathbf{e}_i^P, \quad i = 1, \dots, m \quad (28.2)$$

که در آن \mathbf{X}_i^P از مرتبه‌ی $N_i \times p$ ، بردارهای $\mathbf{1}_i^P$ و \mathbf{e}_i^P از مرتبه‌ی $N_i \times 1$ اند، و $\mathbf{1}_i^P = (1, \dots, 1)^T$. سپس (۲۸.۲) را به دو بخش نمونه‌ای و غیرنمونه‌ای افراز می‌کنیم:

$$\mathbf{y}_i^P = \begin{bmatrix} \mathbf{y}_i \\ \mathbf{y}_i^* \end{bmatrix} = \begin{bmatrix} \mathbf{X}_i \\ \mathbf{X}_i^* \end{bmatrix} \boldsymbol{\beta} + v_i \begin{bmatrix} \mathbf{1}_i \\ \mathbf{1}_i^* \end{bmatrix} + \begin{bmatrix} \mathbf{e}_i \\ \mathbf{e}_i^* \end{bmatrix} \quad (29.2)$$

که در آن بالا نویس * واحدهای نمونه‌گیری نشده را نشان می‌دهد. اگر مدل برای نمونه صادق باشد، یعنی، اگر اریبی گزینش وجود نداشته باشد، آن‌گاه استنباط‌ها درباره‌ی پارامترهای زیر

$$\boldsymbol{\psi} = (\boldsymbol{\beta}^T, \sigma_v^2, \sigma_e^2)^T$$

بر پایه‌ی توزیع احتمال زیر استوارند:

$$f(\mathbf{y}_i | \mathbf{X}_i^P, \boldsymbol{\psi}) = \int f(\mathbf{y}_i, \mathbf{y}_i^* | \mathbf{X}_i^P, \boldsymbol{\psi}) d\mathbf{y}_i^*, \quad i = 1, \dots, m \quad (30.2)$$

که در آن $f(\mathbf{y}_i, \mathbf{y}_i^* | \mathbf{X}_i^P, \boldsymbol{\psi})$ توزیع توام مفروض \mathbf{y}_i و \mathbf{y}_i^* است. از طرفی با در نظر گرفتن فرض $\mathbf{a}_i = (a_{i1}, \dots, a_{iN_i})$ که در آن $a_{ij} = 1$ اگر $j \in s_i$ و در غیر این صورت $a_{ij} = 0$ ، توزیع داده‌های نمونه‌ای $(\mathbf{y}_i, \mathbf{a}_i)$ با فرمول زیر داده می‌شود.

$$\begin{aligned} f(\mathbf{y}_i, \mathbf{a}_i | \mathbf{X}_i^P, \boldsymbol{\psi}) &= \\ & \int f(\mathbf{y}_i, \mathbf{y}_i^* | \mathbf{X}_i^P, \boldsymbol{\psi}) f(\mathbf{a}_i | \mathbf{y}_i, \mathbf{y}_i^*, \mathbf{X}_i^P) d\mathbf{y}_i^* = \\ & \left[\int f(\mathbf{y}_i, \mathbf{y}_i^* | \mathbf{X}_i^P, \boldsymbol{\psi}) d\mathbf{y}_i^* \right] f(\mathbf{a}_i | \mathbf{X}_i^P) \end{aligned}$$

به شرط آنکه رابطه‌ی زیر برقرار باشد:

$$f(\mathbf{a}_i | \mathbf{y}_i, \mathbf{y}_i^*, \mathbf{X}_i^P) = f(\mathbf{a}_i | \mathbf{X}_i^P)$$

یعنی احتمال‌های گزینش نمونه به \mathbf{y}_i^P وابسته نیستند اما ممکن است به \mathbf{X}_i^P وابسته باشند. در این حالت، اریبی گزینش وجود ندارد و می‌توانید فرض کنید که مقدارهای نمونه‌ای نیز از مدل مفروض پیروی می‌کنند، یعنی، از $f(\mathbf{y}_i | \mathbf{X}_i^P, \boldsymbol{\psi})$ برای استنباط‌هایی درباره $\boldsymbol{\psi}$ می‌توان استفاده کرد (اسمیت، ۱۹۸۳).

اگر احتمال‌های گزینش نمونه‌ای به متغیری کمکی، مثلاً \mathbf{z}_i^P ، وابسته باشد که در \mathbf{X}_i^P گنجانده نشده است، آن‌گاه توزیع داده‌های نمونه‌ای $(\mathbf{y}_i, \mathbf{a}_i)$ عبارت‌اند از

$$f(\mathbf{y}_i, \mathbf{a}_i | \mathbf{X}_i^P, \mathbf{z}_i^P, \boldsymbol{\psi}) = \left[\int f(\mathbf{y}_i, \mathbf{y}_i^* | \mathbf{X}_i^P, \mathbf{z}_i^P, \boldsymbol{\psi}) d\mathbf{y}_i^* \right] f(\mathbf{a}_i | \mathbf{z}_i^P, \mathbf{X}_i^P)$$

در این حالت، استنباط درباره‌ی ψ مبتنی بر $f(y_i | \mathbf{X}_i^p, \mathbf{z}_i^p, \psi)$ است که با (۳۰.۲) تفاوت دارد، مگر آنکه، \mathbf{z}_i^p به شرط \mathbf{X}_i^p با y_i^p ارتباط نداشته باشد. در این حالت، اریبی‌گزینش خواهید داشت و بنابراین نمی‌توانید فرض کنید که مدل (۲۸.۲) برای مقادیرهای نمونه‌ای برقرار است. شاید بتوانید مدل (۲۸.۲) را با گنجاندن \mathbf{z}_i^p و سپس آزمون فرض معنی‌دار بودن ضریب رگرسیونی مربوط با استفاده از داده‌های نمونه‌ای، گسترش دهید. اگر فرض صفر رد نشود، آن‌گاه می‌توانید فرض کنید که مدل اصلی (۲۸.۲) برای مقادیرهای نمونه‌ای نیز برقرار است (اسکینر، ۱۹۹۴).

مدل (۲۸.۲) تحت نمونه‌گیری خوشه‌ای دومرحله‌ای درون ناحیه‌های کوچک نیز مناسب نیست، زیرا اثرهای تصادفی خوشه‌ها در آن وارد نشده‌اند. اما می‌توانید مدل را برای در نظر گرفتن چنین ویژگی‌هایی گسترش دهید.

هدف یافتن برآورد پارامترهای میانگین‌های کوچک ناحیه‌ای $\theta_i = \bar{Y}_i = \mathbf{x}_i^T \beta + v_i$ یا مجموع Y_i کوچک ناحیه‌ای است. در مدل (۲۸.۲) اگر قرار دهیم $\sigma_v^2 = 0$ یا به‌طور هم‌ارز $v_i = 0$ ، آن‌گاه مدل‌های رگرسیونی استاندارد به‌دست می‌آید که نوعی از برآوردگرهای هم‌گذاشتی است. میانگین کوچک ناحیه‌ای \bar{Y}_i را به‌صورت زیر می‌نویسند:

$$\bar{Y}_i = f_i \bar{y}_i + (1 - f_i) \bar{Y}_i^* \quad (۳۱.۲)$$

که $f_i = n_i / N_i$ و \bar{y}_i و \bar{Y}_i^* به ترتیب میانگین‌های واحدهای نمونه گرفته شده و نمونه گرفته نشده را نشان می‌دهند. از مدل (۳۱.۲) داریم برآورد میانگین کوچک ناحیه‌ای \bar{Y}_i معادل با برآورد مقدار تحقق‌یافته متغیر تصادفی \bar{Y}_i^* به شرط داده‌های نمونه‌ای $\{y_i\}$ و داده‌های کمکی $\{\mathbf{X}_i^p\}$ است. اگر اندازه‌ی جامعه N_i بزرگ باشد، با توجه به اینکه $\bar{Y}_i = \bar{\mathbf{x}}_i^T \beta + v_i + \bar{E}_i$ و $\bar{E}_i \approx 0$ که در آن میانگین N_i خطای e_{ij} و $\bar{\mathbf{X}}_i$ میانگین معلوم \mathbf{X}_i^p است، می‌توانید میانگین‌های کوچک ناحیه‌ای را به‌صورت

$$\bar{Y}_i = \bar{\mathbf{X}}_i^T \beta + v_i \quad (۳۲.۲)$$

بنویسید. از مدل (۳۲.۲) چنین بر می‌آید که برآورد کردن \bar{Y}_i هم‌ارز با برآورد کردن ترکیب خطی از β و مقدار تحقق‌یافته متغیر تصادفی v_i است.

مثال ۱.۴.۲. (دستمزدها و حقوق‌ها). راثو و چاودری (۱۹۹۵) جامعه مالیات پردازان غیرشرکتی از استان نووا اسکوشیای کانادا را بررسی کردند. آن‌ها مدل

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + v_i + x_{ij}^{1/2} \tilde{e}_{ij}$$

را که حالت خاصی از (۲۶.۲) با $k_{ij} = x_{ij}^{1/2}$ است، پیشنهاد کردند. در اینجا y_{ij} و x_{ij} کل دستمزد و حقوق‌ها و درآمد ناخالص را برای بنگاه زام در ناحیه‌ی i ام نشان می‌دهند. برای برآورد مجموع‌های کوچک ناحیه‌ای Y_i یا میانگین‌های \bar{Y}_i از نمونه‌گیری تصادفی ساده از جامعه‌ی سراسری استفاده کردند.

۱.۴.۲ گسترش‌های مدل در سطح واحد آماری (B)

گسترش مدل در سطح واحد آماری انواع گوناگونی دارند که عبارت‌اند از: مدل رگرسیونی با خطای آشیانی چندمتغیره، مدل خطی با واریانس خطای تصادفی، مدل رگرسیونی با خطای آشیانی دولایه، مدل دوسطحی و مدل آمیخته‌ی خطی کلی.

تذکر ۱.۴.۲. سازگاری ضرایب برآوردگرهای متغیرهای کمکی بین ناحیه و مدل‌های سطح واحد همیشه کامل نخواهد بود (گرینلند و همکاران، ۱۹۹۴). انتظار می‌رود که داشتن مقادیر مشابه در هر دو مورد مناسب باشد (برازش داده شده) اما هنگامی که از داده جمع‌آوری شده استفاده شود، ممکن است تعداد آریبی در برآوردگرها مشاهده شود که به آریبی بوم شناختی معرفی شده است. این حالت ممکن است اتفاق بیفتد، برای مثال هنگامی که اختلالی در متغیرهای کمکی سطح فردی و متغیرهای کمکی جمع‌آوری شده وجود داشته باشد نمی‌توان آن را شرح داد. از این رو، دقت کنید که به‌طور معمول ضرایب را همانند اندازه‌گیری اثرات سطح فردی تفسیر کنید. انتخاب مدل نیز وابسته به نوع داده خواهد بود. داده جمع‌آوری شده آسان‌تر به دست می‌آید، همانند موسسات و دفاتر مختلف آماری که داده‌ی آماری سطح ناحیه را با حجم‌های زیادی ارائه می‌دهند.

۵.۲ مدل رائو-یو

برای نشان دادن کارایی این مدل مثالی را مطرح کرده‌ایم. در نظرسنجی نیروی کار در کانادا، خانوارهای مورد نظر را طی ۶ ماه متوالی و در هر ماه مصاحبه می‌کنند. همچنین ارزیابی کودکان و نوجوانان سنین ۵ تا ۱۷ سال برای بررسی این که چه تعداد از این افراد در زیر خط فقر جای دارند یا خیر. از طرفی زمان سرشماری برای بررسی این موضوع هر دو سال و یا هر سال انجام می‌شود. زیرا دولت برای تخصیص دهی بودجه نیازمند داشتن چنین اطلاعاتی است که کدام مناطق نیازمند بودجه‌ی بیشتری است. اما همیشه سرشماری کردن در بازه‌های زمانی طولانی کارا نیست لذا در هر نظرسنجی ترکیب کردن داده‌ها در طول زمان متوالی منجر به افزایش اندازه نمونه و ... می‌شود، پس این روش دقت در برآوردهای نواحی کوچک را افزایش می‌دهد. لذا رائو و یو^۹ (۱۹۹۲، ۱۹۹۴) مدل فی-هریوت را به دو قسمت تقسیم کردند :

۱. مدل نمونه‌گیری

۲. مدل ربط (پیوند)

با این عمل توانستند خیلی از نواقص مدل فی-هریوت را برطرف کنند. (مانند نادیده گرفتن ویژگی مدل) بنابراین با توجه به مطالب گفته شده مدل رائو-یو را شرح خواهیم داد.

^۹Rao-yo model

- X_i برداری از کاراکترهای ناحیه
- θ_i کاراکترهای ناحیه غیرشناخته شده از عواملی همچون میانگین، واریانس، مد و ...
- $\tilde{\theta}_i$ یک برآوردگر مستقیم از θ_i (میانگین نمونه، برآوردگرهای مبتنی بر طرح و ...)
- m تعدادی از نواحی با داده‌های نمونه
- n_i اندازه نمونه در ناحیه‌ی i ام و N_i اندازه ناحیه

تذکر ۱.۵.۲. رابطه‌ی زیر صورت کلی مدل را در حالت نمونه‌گیری (یعنی در این مدل از خود مدل استفاده نمی‌شود، به‌طور حتم اگر یک برآوردگر مستقیم داشته باشیم این برآوردگر با واقعیت خود می‌تواند فرقی داشته باشد که این فرق ناشی از خطای نمونه‌گیری می‌باشد، یعنی e_i) نشان می‌دهد.

$$\tilde{\theta}_i = \theta_i + e_i$$

تذکر ۲.۵.۲. همچنین رابطه‌ی زیر هم صورت کلی مدل را در حالت ربط (در واقع θ شما خود از یک مدل تولید می‌شود و در این مدل از ویژگی‌های حالت نمونه‌گیری برای پیدا کردن برآوردگر استفاده نخواهیم کرد و تنها از ویژگی‌های خود مدل بهره‌مند می‌شویم تا برآوردگر را بهتر شناسایی نماییم) نشان می‌دهد.

$$\theta_i = \mathbf{X}_i^T \beta + u_i$$

X_i بردار مقادیر متغیر کمکی و با ابعاد $(p \times 1)$ و با استفاده از اطلاعاتی که در سطح ناحیه قرار دارد به دست آورده می‌شود و β نیز برداری با ابعاد $(p \times 1)$ و نشان‌دهنده‌ی ضرایب رگرسیونی است. لذا e_i خطای نمونه‌گیری و مستقل همچنین u_i نشان‌دهنده‌ی اثرات تصادفی و دارای توزیع زیر می‌باشند:

$$e_i \stackrel{ind}{\sim} N(0, \sigma_{D_i}^2) \text{ و } u_i \stackrel{iid}{\sim} N(0, \sigma_u^2)$$

۶.۲ مدل‌های آمیخته‌ی خطی کلی

۱.۶.۲ مقدمه

در برآورد کوچک ناحیه‌ای استفاده کردن از مدل‌های آمیخته خطی مناسب است، دلیل استفاده کردن از این نوع مدل‌ها در ترکیب کردن اطلاعات از منابع مختلف و در نظر گرفتن منابع مختلف تولید کننده‌ی خطا، که به‌طور ویژه موثر می‌باشند. مدل‌های آمیخته‌ی خطی نوعی از مدل‌های آماری هستند که از اثرهای ثابت و تصادفی تشکیل شده‌اند. بنابراین این نوع از مدل‌ها تفاوتی که بین نواحی وجود دارد را نمی‌تواند به کمک اثرهای ثابت به دست آورد، پس به کمک اثرهای تصادفی محاسبه می‌شوند، از طرفی مدل‌های خطی رگرسیونی وجود دارد که

تفاوت بین نواحی را به کمک متغیرهای کمکی نشان می‌دهد. به‌طور حتم در مفهوم کوچک ناحیه‌ای، این دسته از مدل‌ها نوعاً به دو گروه کلی تقسیم می‌شوند یکی در سطح ناحیه و دیگری در سطح واحد. به منظور آشنایی بیشتر این موضوع مطالب مختصری از مدل‌های آمیخته‌ی خطی و تعمیم‌یافته و تعمیم‌یافته فضایی را شرح می‌دهیم.

۲.۶.۲ مدل آمیخته‌ی خطی LMM

همان‌طوری که در مقدمه اشاره کردیم، مدل آمیخته‌ی خطی یک مدل آماری است که شامل اثرهای ثابت و تصادفی می‌باشد. کاربرد این مدل را می‌توان در شاخه‌های علوم زیستی و اجتماعی دید. از طرفی مدل آمیخته‌ی خطی در برآورد کوچک ناحیه‌ای پرکاربرد است، زیرا استفاده از این مدل می‌تواند تفاوت بین نواحی را که توسط اثرهای ثابت مدل محاسبه نمی‌شوند، با به‌کارگیری از اثرهای تصادفی انجام داد و همین‌طور مدل‌های رگرسیونی خطی وجود دارد که تفاوت بین نواحی را فقط از طریق اطلاعات کمکی نشان می‌دهد. مدل‌های آمیخته‌ی خطی را می‌توان بر سه نوع برشمرد. مدل‌هایی با اثرهای ثابت، مدل‌هایی با اثرهای تصادفی و مدل‌هایی با اثرهای آمیخته.

بهترین راه حلی که به ما کمک می‌کند تا بتوان مدل آمیخته‌ی خطی را راحت‌تر درک کرد، ارائه یک مدل رگرسیونی است، این مدل برای اینکه ارتباط بین متغیر پاسخ و متغیرهای کمکی را ایجاد کند، استفاده می‌شود و آن را به‌صورت زیر نشان می‌دهیم:

$$y = X\beta + e, \quad (۳۳.۲)$$

باتوجه به معادله‌ی بالا y بردار مشاهدات، X ماتریس طرح معلوم، β بردار نامعلوم از ضرایب رگرسیونی است، که اغلب اثرهای ثابت نامیده می‌شوند و e یک بردار خطاهای تصادفی غیرقابل مشاهده است. به‌طور کل فرض نرمال بودن خطاها را در نظر بگیرید. خطاها دارای توزیع نرمال با میانگین صفر و واریانس σ^2 می‌باشند، یعنی $e \stackrel{iid}{\sim} N(0, \sigma^2 I)$. بنابراین توزیع مقدارهای مشاهده شده نیز نرمال خواهد شد، و در نتیجه داریم $y \stackrel{iid}{\sim} N(X\beta, \sigma^2 I)$ و صورت ماتریسی این مدل به شکل زیر است:

$$y = X\beta + Zv + e, \quad (۳۴.۲)$$

که $y_{n \times 1}$ بردار مشاهدات، $\beta_{p \times 1}$ بردار اثرهای ثابت، $v_{q \times 1}$ بردار اثرهای تصادفی، $e_{n \times 1}$ بردار خطاهای تصادفی، $X_{n \times p}$ ماتریس طرح برای اثرهای ثابت و $Z_{q \times 1}$ ماتریس طرح برای اثرهای تصادفی که هر دو آن‌ها معلوم هستند. از طرفی توجه به اینکه e و v متغیرهای تصادفی غیرهمبسته با میانگین‌های صفر و دارای ماتریس کوواریانس می‌باشد و می‌توان به‌صورت زیر آن را نشان داد:

$$Var(v) = G, Var(e) = R, Cov(v, e) = 0$$

لذا ماتریس‌های کوواریانس شامل بعضی از پارامترهای پراکندگی نامعلوم می‌باشند. اگر فرض نرمال بودن اثرهای تصادفی و خطاها آورده شود، بنابراین خواهیم داشت $v \sim N(0, G)$ و $e \sim N(0, R)$ ، که این مدل را مدل آمیخته‌ی خطی گاوسی^{۱۰} می‌نامند. همچنین فرض می‌شود که توزیع حاشیه‌ای y نرمال چند متغیره‌ای است که مشمول میانگین $X\beta$ و ماتریس کوواریانس $V = ZGZ^T + R$ که تابعی از بردار v می‌باشد پس $V = V(v)$ و در نتیجه $y \sim N(X\beta, V)$ را خواهیم داشت (جیانگ و لاهیری، ۲۰۰۶). همان طوری که در مدل آمیخته‌ی خطی (۲.۳۴) مشخص است، این مدل‌ها سه مولفه دارند (رائو، ۲۰۰۳):

- مولفه تصادفی: این مولفه متغیر پاسخ y را با یک توزیع احتمال مناسب تعریف می‌کند. توزیع احتمال متغیر پاسخ می‌تواند توزیع‌های مختلفی باشد که معمولاً برای داده‌های گسسته توزیع دوجمله‌ای و یا پواسون و برای پیوسته توزیع نرمال مورد استفاده قرار می‌گیرند.
- مولفه سیستماتیک: این مولفه متغیرهای کمکی استفاده شده یعنی، متغیرهایی که نقش x_j و z_j را در مدل دارند، به عنوان پیشگوهای مدل تعیین شده و به صورت ترکیب خطی از متغیرهای پیشگو یعنی $\eta = x\beta + zv$ تحت عنوان پیشگوی خطی در نظر می‌گیرند.
- مولفه ربط: تابع یا پیوندی بین مولفه سیستماتیک و امیدریاضی (میانگین) مولفه تصادفی است که در آن نحوه ارتباط $\mu = E(y|x, z)$ با متغیرهای کمکی در پیشگوی خطی نشان داده می‌شود و می‌توان میانگین μ را به طور مستقیم و یا به صورت تابع $\eta = g(\mu)$ و به صورت $\eta = g(\mu) = x\beta + zv$ بیان کرد، بنابراین، $g(\cdot)$ را تابع ربط می‌نامند. در تابع ربط، میانگین متغیر پاسخ با پیشگوی خطی ربط داده می‌شود و این مساله از تبدیل بر روی متغیر پاسخ متفاوت است. همچنین نکته مهم در انتخاب تابع ربط، این است که تابع انتخاب شده مقادیر پیشگویی قابل قبول از متغیر پاسخ را نتیجه می‌دهد (جیانگ و لاهیری، ۲۰۰۶؛ رائو، ۲۰۰۳).

باتوجه به مطالب گفته شده، برای وضوح بهتر مباحثی را در پیرامون مدل‌های آمیخته‌ی خطی تعمیم‌یافته GLMM شرح می‌دهیم.

۳.۶.۲ مدل آمیخته‌ی خطی تعمیم‌یافته GLMM

بنابر مطالبی که در بخش قبلی گفته شد، مدل‌های رگرسیونی را در مواقعی می‌توان استفاده کرد که متغیر پاسخ، متغیری پیوسته باشد اما با این حال شروطی نیز لازم است:

۱. مشاهدات از هم مستقل در نظر گرفته شود.
۲. می‌توان روابط بین متغیر پاسخ و اثرهای ثابت را به کمک یک تابع خطی مدل‌سازی کرد.
۳. از طرفی واریانس متغیر پاسخ یک مقدار ثابت باشد.

^{۱۰} Gaussian linear mixed model

۴. خطاها معمولاً از توزیع نرمال پیروی می‌کنند.

بنابر شروط گفته شده، در بعضی کاربردها می‌توان آن‌ها را نقض کرد، بنابراین استفاده از مدل‌های LMM و GLMM راه حلی برای برطرف کردن این نوع نقص‌ها می‌باشد. در مدل آمیخته‌ی خطی فرض بر این است که، متغیرهای پاسخ از هم مستقل‌اند و از یک توزیع خاص که جزء خانواده‌های نمایی است پیروی می‌کند که در بیشتر مواقع این توزیع‌ها شامل توزیع‌های نرمال، پواسون و برنولی می‌باشد. لذا میانگین متغیر پاسخ μ به ترکیب خطی از متغیرهای کمکی $\xi_1 = \mathbf{X}\beta$ به کمک معکوس تابع پیوند مرتبط می‌شوند، پس خواهیم داشت:

$$E(\mathbf{y}) = \boldsymbol{\mu} = g^{-1}(\mathbf{X}\boldsymbol{\beta}) = g^{-1}(\xi_1)$$

β برداری از ضرایب رگرسیونی نامعلوم و همچنین g یک تابع پیوند است. اما در حالت کلی مدل رگرسیونی یک حالت خاص از مدل GLM می‌باشد که تابع پیوند آن‌ها تابع همانی است. لذا اگر مشاهدات نسبت به هم مستقل نباشند، مدل مطرح شده همان GLM است در اغلب موارد مناسب نیستند به همین دلیل، لازم است از مدل دیگری استفاده کرد و آن هم مدل‌های آمیخته‌ی خطی تعمیم‌یافته می‌باشد.

در مدل‌های آمیخته‌ی خطی تعمیم‌یافته پیشگوی خطی همانند مدل‌های آمیخته‌ی خطی شامل اثرهای ثابت و تصادفی می‌باشد، یعنی:

$$\xi_2 = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}$$

همچنین در این دسته از مدل‌ها،

$$E(\mathbf{y}|\mathbf{u}) = g^{-1}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) = g^{-1}(\xi_2)$$

$\mathbf{X}_{n \times p}$ ماتریس طرح برای اثرهای ثابت $\beta_{p \times 1}$ ، $\mathbf{Z}_{n \times q}$ یک ماتریس طرح برای اثرهای تصادفی $\mathbf{u}_{q \times 1}$ و g یک تابع پیوند است. حال فرض کنید که مدل آمیخته‌ی خطی را بسط دهیم آن‌گاه مشاهدات y_1, y_2, \dots, y_n به صورت شرطی مستقل هستند، یعنی $y_i | \mathbf{u} \sim N(\mathbf{x}_i \boldsymbol{\beta} + \mathbf{z}_i \mathbf{u}, \tau^2)$. با توجه به این می‌توان گفت که x_i و z_i سطرهای i ام ماتریس‌های \mathbf{X} و \mathbf{Z} می‌باشند و پارامتر τ^2 یک واریانس نامعلوم است و با داشتن فرض نرمال و $\mathbf{R} = \tau^2 \mathbf{I}$ منجر به مدل $\boldsymbol{\theta}_i = \mathbf{X}_i \boldsymbol{\beta} + v_i$ می‌شود. اکنون اگر فرض بر این باشد که، بردارهای تصادفی \mathbf{u} برقرار باشند و y_1, y_2, \dots, y_n به‌طور شرطی مستقل باشند، آن‌گاه توزیع شرطی $y_i | \mathbf{u}$ یکی از اعضای خانواده توزیع‌های نمایی با تابع چگالی احتمال زیر می‌شود:

$$f_i(y_i | \mathbf{u}) = \exp\left[\frac{y_i \xi_i - b(\xi_i)}{a_i(\phi)} + c_i(y_i, \phi)\right],$$

که ϕ پارامتر پراکندگی و $a_i(\cdot), b_i(\cdot), c_i(\cdot, \cdot)$ توابعی معلوم هستند و اگرچه در برخی موارد اینکه نامعلوم هستند نیز امکان دارد. جیانگ و لاهیری (۲۰۰۶) ژانگ و لی (۲۰۱۱)، بعد از انتخاب مدل آماری مناسب با توجه به شرایط مساله روش‌های مرسوم که برای برآورد کوچک ناحیه‌ای شدنی می‌باشد را تبیین کردند.

۴.۶.۲ مدل‌های آمیخته خطی تعمیم‌یافته فضایی SGLMM

فرض کنید G ناحیه فضایی مشبکه‌ای و $Y(s_i)$ و $i = 1, \dots, n$ متغیر پاسخ شمارشی برای ناحیه s_i باشد. همچنین فرض کنید در هر ناحیه مقادیر بردار p بعدی متغیرهای تبیینی $\mathbf{x}(s_i)$ مشاهده شده باشند. یک SGLM به صورت زیر تعریف می‌شود:

$$E(Y(s_i)|v(s_i)) = g^{-1}(\mathbf{x}^T(s_i)\boldsymbol{\beta} + v(s_i) + \epsilon(s_i)), \quad i = 1, \dots, n \quad (۳۵.۲)$$

که در آن $g(\cdot)$ یک تابع پیوند مشتق‌پذیر و معکوس‌پذیر است، $\boldsymbol{\beta}$ بردار p بعدی پارامترهای ثابت مدل است و $\epsilon(s_i)$ جمله خطای تصادفی غیرساختارمند با توزیع نرمال $N(0, \sigma_\epsilon^2)$ است که اغلب برای خطای اندازه‌گیری در نظر گرفته می‌شود. اثر تصادفی فضایی $\mathbf{v} = (v(s_1), \dots, v(s_n))'$ از یک مدل اتورگرسیو شرطی^{۱۱} (CAR) همانند مدل زیر پیروی می‌کند.

$$v_i | \mathbf{v}_{-i}, \sigma_v^2 \sim N\left(\sum_{j \in \delta_i} \frac{d_{ij}}{|\delta_i|} v_j, \frac{\sigma_v^2}{|\delta_i|}\right)$$

و توزیع متغیر پاسخ $Y(s_i)$ به شرط $v(s_i)$ ، از خانواده توزیع‌های نمایی با تابع چگالی

$$f(y(s_i)|v(s_i); \boldsymbol{\beta}, \phi) = \exp\left[\frac{1}{\phi} (a(\mu_i)y(s_i) - b(\mu_i)) + c(y(s_i), \phi)\right]$$

پیروی می‌کند، که در آن $\mu_i = \mu(s_i) = E(Y(s_i)|v(s_i))$ پارامتر پراکندگی معلوم و توابع $a(\cdot)$ ، $b(\cdot)$ و $c(\cdot, \cdot)$ نیز معلوم هستند.

تابع درست‌نمایی کناری مدل (۳۵.۲)، برای بردار مشاهدات \mathbf{y} ، به صورت

$$L(\boldsymbol{\beta}, \sigma_v^2, \sigma_\epsilon^2 | \mathbf{y}) \propto \int \prod_{i=1}^n f(y(s_i)|v(s_i); \boldsymbol{\beta}) \phi_n(\mathbf{v}; 0, \sigma_v^2 A^{-1}) d\mathbf{v} \quad (۳۶.۲)$$

است، که در آن $\phi_n(\mathbf{v}; 0, \sigma_v^2 A^{-1})$ تابع چگالی نرمال n متغیره با بردار میانگین صفر و ماتریس کوواریانس $\sigma_v^2 A^{-1}$ است. محاسبه تابع درست‌نمایی (۳۶.۲) مستلزم حل انتگرالی با بعد برابر تعداد مشاهدات است و برای داده‌های (تعداد نواحی) حجیم، محاسبه آن کاری طاقت‌فرسا و مشکل‌ساز است. هر چند روش‌های تقریبی مختلفی در دیدگاه مبتنی بر درست‌نمایی برای محاسبه این درست‌نمایی پیشنهاد شده‌اند (بریسلو و کلیتون، ۱۹۹۳). اما استفاده از چارچوب بیزی برای رده SGLMMs، به دلیل وجود الگوریتم‌های MCMC، انتخاب اول است.

کاربرد مدل GLM در برآورد کوچک ناحیه‌ای

۱. مدل‌های اثرات آمیخته در سطح ناحیه

مدل‌های اثرات آمیخته نوع A (یعنی همان مدل در سطح ناحیه) را می‌توانید با استفاده از مدلی که به فرم زیر، نمایش دهید:

$$\hat{\mathbf{Y}} = \bar{\mathbf{X}}\boldsymbol{\beta} + \mathbf{Z}u + \boldsymbol{\epsilon}$$

^{۱۱}Conditional autoregressive

که در آن \hat{Y}_D برداری از برآوردهای غیرمستقیم، Z نیز به توصیفی از ساختار اثرات تصادفی u می‌پردازد و ε جمله خطا می‌باشد. در اصل می‌خواهیم فرض کنیم که Z یک ماتریس همبندی است، اما یک ساختار پیچیده‌تر، زمانی می‌تواند مورد استفاده قرار بگیرد که مقدار آن در یک ناحیه به چندین اثر تصادفی بستگی داشته باشد. فرض می‌کنیم که اثرات تصادفی u_i دارای توزیع نرمال با میانگین صفر و واریانس σ_u^2 است. خطای نمونه‌گیری یا تغییرات فردی ε ، مستقل با میانگین صفر و واریانس قطری $(\hat{\sigma}_i^2)$ در نظر گرفته می‌شود. علاوه بر این D و V ماتریس واریانس u و ε اند. ناشرطی بودن روی اثرات تصادفی، برای واریانس \hat{Y}_D می‌تواند به این صورت $G = \hat{Z}DZ + V$ باشد. با توجه اینکه واریانس شناخته شده است، پارامتر β می‌تواند توسط میانگین‌هایی از کمترین توان‌های دوم تعمیم یافته شده برآورد شود، هر چند که در این مورد با ماکسیم درست‌نمایی معادل است. اثرات تصادفی را می‌توان به چندین روش به شرح زیر برآورد کرد. از این رو، برآوردها

$$\hat{Y} = \bar{X}_i \hat{\beta} + Z \hat{u}$$

به طوری که برآوردهای غیرمستقیم برای برخی نواحی از دست رفته است اما با این حال متغیرهای کمکی در سطح ناحیه وجود دارد، و مقدار \hat{u} باید به صفر برسد و آن‌گاه برآوردها ترکیبی محاسبه می‌شود.

$$\hat{Y} = \bar{X}_i \hat{\beta}$$

اگر اثرات تصادفی یک ساختار داده شود برآوردها قبلی می‌تواند برای برآورد اثرات تصادفی حساب شود.

۲. مدل‌های اثرات آمیخته در سطح واحد

برای مدل‌های نوع B (یعنی همان مدل در سطح واحد) را می‌توانید مشابه فرمولی زیر داشته باشید:

$$y_{ij} = x_{ij} \beta + \sum_l z_{li} u_i + e_{ij}$$

که در این صورت e_{ij} تغییرات تصادفی فردی یا واحد j ام در ناحیه‌ی i ام است. این مدل‌ها همچنین می‌توانند با تقسیم داده‌ها به واحدهای نمونه‌ی (y_1) و غیرنمونه‌ی (y_2) بیان شوند:

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \beta + \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} u + \begin{bmatrix} e_1 \\ e_2 \end{bmatrix}$$

توجه داشته باشید که طول از y_1 و y_2 را به صورت $\sum_{i=1}^n (N_i - n_i)$ و $\sum_{i=1}^n n_i$ است و تمامی ماتریس‌های دیگر نیز بر اساس اندازه‌ها می‌باشند. برآورد از مدلی که مبنی بر معادله داده مشاهده شده

$$y_1 = x_1 \beta + Z_1 u + e_1$$

و روش‌های متفاوتی که می‌توان برای ارایه برآوردهای پارامتر از مدل استفاده کرد امکان پذیر خواهد بود. برآورد از پارامتر β توسط ماکسیم درستنمایی^{۱۲} (ML)، یا شبه احتمال تاوانیده^{۱۳} نیز انجام می‌شود، ضمن اینکه ماکسیم درستنمایی مقید^{۱۴} (REML) می‌تواند برای به‌دست آوردن غیراریبی برآورد واریانس استفاده شود. مک کولا و سیرل (۲۰۰۱)، استراتژی‌های دیگری را برای به‌دست آوردن برآوردهایی از پارامترهای مختلف مطرح کرده‌اند. علاوه بر این ایراریا کانسرسیوم (۲۰۰۴) الگوریتم‌های دقیق و ترفندهای محاسباتی را برای برآورد پارامترهایی در نوع‌های متفاوت آن هم از مدل‌های آمیخته ارایه داده است. سپس برآورد کوچک ناحیه‌ای با جمع‌آوری مشاهدات و مقادیر برآورد شده در یک روش به‌دست می‌آید:

$$\hat{Y}_i = \frac{n_i}{N_i} \sum_{j=1}^{n_i} y_{ij} + \frac{N_i - n_i}{N_i} \sum_{j=1}^{N_i - n_i} \hat{y}_{ij} = \frac{n_i}{N_i} \sum_{j=1}^{n_i} y_{ij} + \frac{N_i - n_i}{N_i} \sum_{j=1}^{N_i - n_i} x_{ij} \hat{\beta} + \hat{u}_i \quad (37.2)$$

ضمن اینکه \hat{y}_{ij} برآوردی برای واحد j ام در ناحیه‌ی i ام می‌باشد. هنگامی که اندازه نمونه‌ی n_i نسبت به N_i خیلی کوچک باشد برآوردگر مورد استفاده به شرح زیر است:

$$\hat{Y}_i = \frac{1}{N_i} \sum_{j=1}^{N_i} x_{ij} \hat{\beta} + \hat{u}_i = \bar{X}_i \hat{\beta} + \hat{u}_i$$

۷.۲ برآوردهای غیرمستقیم و روش‌های مبتنی بر مدل

۱.۷.۲ برآوردگر ترکیبی

هنگامی که متغیری مدام اندازه‌گیری شود، مدل خطی می‌تواند با ایجاد یک ارتباط بین متغیر هدف و متغیرهای کمکی استفاده شود. فی-هریوت مدل زیر را که ترکیبی از برآوردگر مستقیم و رگرسیونی می‌باشد پیشنهاد داده است:

$$\bar{Y}_i = \beta \bar{X}_i + \epsilon_i; \quad \epsilon_i \sim N(0, V_i^2) \quad (38.2)$$

بنابراین V_i^2 واریانس طرحی است. برآوردگر ترکیبی می‌تواند با استفاده از برآورد β در مدل فوق به‌دست آورده شود، پس خواهیم داشت:

$$\hat{Y}_{SA,i} = \hat{\beta} \bar{X}_i$$

$\hat{Y}_{SA,i}$ یک برآوردگر ترکیبی در ناحیه‌ی i ام است. همانند روش‌های قبلی متغیرهای کمکی برای همه نواحی کوچک در دسترس‌اند. این رویکرد می‌تواند پاسخ‌های غیرنرمال را از طریق

^{۱۲} Maximum likelihood

^{۱۳} quasi-probability penalized

^{۱۴} Restricted maximum likelihood

میانگین‌های مدل آمیخته خطی تعمیم‌یافته در پاسخ و پیش‌بینی خطی با یک تابع متناسب که مرتبط می‌باشد، گسترش دهد. برآورد واریانس را می‌توان به‌شيوه‌ی زیر انجام داد:

$$\text{var}[\hat{\mathbf{Y}}_{\text{SA},i}] = E[(\hat{\mathbf{Y}}_{\text{SA},i} - \hat{\beta}\bar{\mathbf{X}}_i - \epsilon_i)^2] = \bar{\mathbf{X}}_i^T \text{var}[\hat{\beta}]\bar{\mathbf{X}}_i + \mathbf{V}_i^2$$

معادله مدل در سطح واحد در این برآوردگر مشمول نمونه‌های در دسترس برای هر ناحیه و تحت مدل‌سازی زیر پیشنهاد شده‌است:

$$y_{ij} = \mathbf{x}_{ij}\beta + e_{ij}$$

e_{ij} خطای فردی و دارای توزیع نرمال با میانگین ۰ و واریانس σ_i^2 است. فرض کنیم تغییرات در داخل ناحیه داخلی می‌تواند از ناحیه‌ای به ناحیه‌ی دیگر متفاوت باشد. اگر نسبت n_i/N_i خیلی کوچک باشد همبستگی جمعیت کل را می‌توان چشم‌پوشی کرد و برآوردگر ترکیبی مدل را در سطح واحد به‌صورت زیر نشان داد:

$$\hat{\mathbf{Y}}_{\text{SB},i} = \frac{1}{N_i} \sum_{j=1}^{N_i} \mathbf{x}_{ij}\hat{\beta} = \bar{\mathbf{X}}_i\hat{\beta}$$

هرچند این نسبت برای مدل‌های غیرخطی یا غیرنرمال صورت نمی‌پذیرد. این نوع از برآوردگر ترکیبی به‌طور گسترده توسط اداره آمار ملی در انگلیس برای ایجاد گزارش‌های متفاوت ناحیه‌ای مورد استفاده قرار می‌گیرد. یار و همکاران (۲۰۰۲) و هدی و همکاران (۲۰۰۳) متکی به نوعی از مدل در سطح واحد می‌باشند، اما با توجه به اینکه هیچ متغیر کمکی در دسترس در سطح فردی نیست آن را به‌صورت $\mathbf{x}_{ij} = \mathbf{X}_i$ در نظر می‌گیرند. واریانس این برآوردگر به‌صورت زیر خواهد بود:

$$\text{var}[\hat{\mathbf{Y}}_i] = \bar{\mathbf{X}}_i^T \text{var}\hat{\beta}\bar{\mathbf{X}}_i + \sigma_i^2/N_i$$

همچنین توجه کنید، برآورد ترکیبی مدل در سطح ناحیه برآوردگری برای ناحیه را فراهم می‌کند، برآوردگر مدل در سطح واحد مقادیر متغیر مورد نظر را برای افرادی که نمونه‌برداری نشده‌اند ($\forall j$ در ناحیه i) $(\hat{y}_{ij} = \hat{\beta}x_{ij})$ و برآوردگرهایی که باید بعد از آن به‌طور متوسط بر روی تمام افراد در ناحیه محاسبه شوند را در نظر می‌گیرد.

۸.۲ برآوردگرهای مبتنی بر مدل صریح

بعد از انتخاب مدل مناسب و گزینش متغیرهای کمکی خوب که با متغیر تحت بررسی وابستگی نسبتاً خوبی را داشته باشد، نوبت به انتخاب برآوردگر مناسب برای نواحی کوچک است. به‌طور حتم دو رویکرد متفاوت در یافتن برآورد وجود دارد که یکی، رویکرد فراوانی‌گرای کلاسیک و دیگری روش بیزی می‌باشد. از طرفی می‌دانید که برآوردگر غیرمستقیم مبتنی بر مدل صریح را، برآوردگر غیرمستقیم مدرن گویند. این نوع از برآوردگرها، برآوردگرهای مبتنی بر مدل صریح

هستند که به ۴ دسته تقسیم می‌شوند و همین طور این ۴ دسته زیر گروه همان روش‌های فراوان‌گرای کلاسیک و روش بیزی‌اند پس، عبارت‌اند از:

۱. برآوردگرهای پیشگوی نارایب خطی تجربی BLUP

۲. برآوردگرهای بهترین پیشگوی نارایب خطی تجربی EBLUP

۳. برآوردگرهای بیز تجربی EB

۴. برآوردگرهای بیزی سلسله مراتبی HB

در ادامه درباره‌ی این برآوردگرها مطالبی را ارائه خواهیم داد (گوش و رائو، ۱۹۹۴).

۱.۸.۲ برآوردگرهای پیشگوی نارایب خطی تجربی BLUP

میانگین و مجموع نواحی کوچک را می‌توان به صورت ترکیب‌های خطی اثرهای ثابت و تصادفی بیان کرد. بنابراین روش برآوردگرهای پیشگوی نارایب خطی برای این دسته از پارامترها را می‌توان در چارچوب بسامدگرای کلاسیک، با مرتبط قرار دادن به نتایج کلی در مورد روش برآورد BLUP به دست آورد. این بهترین در این جمله یعنی کمترین میانگین توان دوم خطا را داشته باشد. این نوع از برآوردگرها در رده‌ی برآوردگرهای نارایب خطی MSE را مینیمم می‌سازد و همچنین به نرمال بودن اثرهای تصادفی وابسته نیست (مک کولا و سیرل، ۲۰۰۱). این دسته از برآوردگرها، واریانس‌ها و کوواریانس‌های اثرهای تصادفی را به کمک روش برازش ثابت‌ها یا روش گشتاوری برآورد می‌کند. یکی دیگر از روش‌های که در دسترس است و به وسیله‌ی آن می‌توان برآورد مولفه‌های واریانس و کوواریانس را به دست آورد و نیز با لحاظ فرض نرمال بودن اثرها، روش‌های ماکسیمم درست‌نمایی و یا ماکسیمم درست‌نمایی مقید می‌باشد. با به کارگیری این مولفه‌های برآوردشده در برآوردگر BLUP، برآوردگری دو مرحله‌ای به دست می‌آید که در مقایسه با برآوردگر بیز تجربی بهتر عمل می‌کند و از آن به عنوان برآوردگر BLUP تجربی یا EBLUP یاد می‌شود (هارویل، ۱۹۹۱).

بنابراین با توجه به مدلی که قبلاً ارائه شده بود، برآوردگر را می‌توان به شرح زیر نشان داد:

$$E[u|y] = \hat{u} = \mathbf{DZ}^T \Sigma_{11}^{-1} (\mathbf{y} - \mathbf{X}_1 \beta)$$

هنگامی که اثرات تصادفی در این روش برآورده شوند؛ معادله‌ی (۳۷.۲) به یک برآوردگر BLUP تبدیل می‌شود و مجموع برآوردگر مستقیم مبتنی بر نمونه در ناحیه‌ی i ام می‌شود و قسمت ثابت مبتنی بر متغیرهای کمکی و برآوردگری از اثر تصادفی است.

توجه کنید اگر در ناحیه‌ی i ام هیچ نمونه‌ای در دسترس نباشد برآوردگر BLUP به برآوردگر ترکیبی کاهش می‌یابد، زیرا، برآوردگر مستقیم وجود ندارد و $\hat{u} = 0$ است. هنگامی که واریانس اثرات تصادفی ناشناخته باشد و از داده برآورد شود، برآوردگرهای مشابه می‌توانند توسعه یابند.

بنابراین موارد بیان شده هم، بهترین برآوردگر ناریب خطی تجربی BLUP نامیده می‌شود. راثو (۲۰۰۳) در مورد استفاده از برآوردگرهای BLUP و EBLUP در برآورد کوچک ناحیه‌ای پیشنهادهایی را مطرح نمود.

برای مدل در سطح ناحیه برآورد میانگین توان دوم خطا MSE اغلب با تقریب ۳ مولفه محاسبه می‌شود که به ترتیب عبارت‌اند از: (G_1) نشان‌دهنده‌ی عدم اطمینان در برآوردگر، (G_2) برآورد β و (G_3) برآورد واریانس σ^2 می‌باشند. که اغلب نیز به صورت زیر نشان داده می‌شود:

$$MSE[\hat{Y}_i] \approx G_1 + G_2 + G_3$$

از طرفی MSE برای مدل در سطح واحد نیز مشابه مدل در سطح ناحیه انجام می‌گیرد. برای توضیح بیشتر به روش‌های EBLUP می‌پردازیم.

۲.۸.۲ برآوردگرهای بهترین پیشگوی ناریب خطی تجربی EBLUP

برآوردگر مطرح شده در دو مرحله به دست می‌آید. همان طوری که از نام این برآوردگر مشخص است در مرحله‌ی اول، بهترین پیشگوی ناریب خطی در رده‌ی برآوردگرهای ناریب خطی، زمانی به دست می‌آید که مرتبط به کمیت مورد بررسی باشد و میانگین توان دوم خطا را مینیمم سازد. این دسته از برآوردگرها به واریانس و کوواریانس اثر تصادفی در مدل بستگی دارند. در مرحله‌ی بعدی، برآوردگر بهترین پیشگوی ناریب خطی، زمانی به دست می‌آید که برآوردگرهای مناسب را برای پارامترهای واریانس و کوواریانس، در بهترین پیشگوی ناریب خطی در مرحله‌ی اول جایگزین نماییم. همچنین برای به دست آوردن برآورد واریانس‌ها معمولاً از روش گشتاوری یا روش لاگرانژ استفاده می‌شود. و گاهی نیز با فرض نرمال بودن، می‌توان از روش دیگری برای برآورد مولفه‌های واریانس استفاده کرد که برای این عمل روش‌های ماکسیمم درست‌نمایی یا ماکسیمم مقید معرفی شده است.

با توجه به مطالب بیان شده اکنون به توضیح روش برآوردگر بهترین پیشگوی ناریب خطی تجربی پرداخته می‌شود. فرض کنید داده‌های نمونه‌ای از یک مدل آمیخته‌ی خطی به صورت زیر پیروی می‌کند، لذا خواهیم داشت:

$$y = X\beta + Zv + e. \quad (39.2)$$

• y برداری از مشاهدات $(y_1^T, \dots, y_n^T)^T$ و y_i برداری با ابعاد $1 \times n$

• $\mathbf{X} = \text{col}_{1 \leq i \leq m}(\mathbf{X}_i)$ و \mathbf{X}_i ماتریسی با ابعاد $(n_i \times p)$ ،

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_m \end{pmatrix}$$

• $\mathbf{Z} = \text{diag}_{1 \leq i \leq m}(\mathbf{Z}_i)$ و \mathbf{Z}_i ماتریسی با ابعاد $(n_i \times h_i)$

$$\mathbf{Z} = \begin{pmatrix} \mathbf{Z}_1 & 0 & \dots & 0 \\ 0 & \mathbf{Z}_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \mathbf{Z}_m \end{pmatrix}$$

همان‌طوری که همانند بالا اشاره داشته‌ایم، \mathbf{y} بردار مشاهدات نمونه‌ای با بعد $n \times 1$ ، \mathbf{X} و \mathbf{Z} ماتریس‌های معلومی هستند که به ترتیب با ابعاد $n_i \times p_i$ و $n_i \times h_i$ نمایش داده می‌شوند. $\beta = (\beta_1, \dots, \beta_p)^T$ برداری از ضرایب رگرسیونی و همچنین \mathbf{v} بردار اثرات تصادفی که با ابعاد $(h_i \times 1)$ می‌باشند. از طرفی \mathbf{v} و \mathbf{e} نیز دارای توزیعی‌هایی با بردار میانگین $\mathbf{0}$ و ماتریس کوواریانس Σ_e و Σ_v اند، در نتیجه مستقل نیز می‌باشند. ماتریس کوواریانس‌های اشاره شده، دارای مولفه‌های واریانسی می‌باشند که این مولفه‌ها در بردار $\delta = (\delta_1, \dots, \delta_q)^T$ آورده می‌شوند. ماتریس کوواریانس برای مدل (۳۹.۲) به صورت زیر به دست می‌آید:

$$\text{Var}(\mathbf{y}) = \mathbf{V} = \mathbf{V}(\delta) = \Sigma_e + \mathbf{Z}\Sigma_v\mathbf{Z}^T,$$

که برای تمامی δ ها، نامنفرد است.

ترکیب خطی زیر که از ضرایب رگرسیونی β و مقادیر تحقق‌یافته‌ی \mathbf{v} است را برای بردارهای ثابت و معین η و \mathbf{m} در نظر بگیرید.

$$\mu = \eta^T \beta + \mathbf{m}^T \tilde{\mathbf{v}}$$

در روش بهترین پیشگوی ناریب خطی تجربی در پی دستیابی برآوردگر μ یعنی همان $\hat{\mu}$ هستیم، که این نوع از برآوردگر برای μ ناریب می‌باشد. و اگر فرض بر این باشد که، پارامترهای δ معلوم باشند، بنابراین پیشگوی ناریب خطی به صورت زیر به دست می‌آید:

$$\tilde{\mu}^H = t(\delta, \mathbf{y}) = \eta^T \tilde{\beta} + \mathbf{m}^T \tilde{\mathbf{v}}, \quad (40.2)$$

بالانویس H در $\tilde{\mu}^H$ نشان‌دهنده‌ی نام هندرسون^{۱۵} است. و $\tilde{\beta}$ بهترین برآوردگر ناریب خطی β است و آن را به صورت زیر نشان می‌دهیم:

$$\tilde{\beta} = \tilde{\beta}(\delta) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{Y}$$

^{۱۵}Henderson

و \tilde{v} برابر است با

$$\tilde{v} = \tilde{v}(\delta) = \Sigma_v \mathbf{Z}^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\tilde{\beta}).$$

بنابراین $\tilde{\mu}^H \tilde{\mu}$ برآوردگر پیشگوی ناریب خطی است که در رده‌ی برآوردگرهای خطی ناریب قرار دارد که بهترین برآوردگر می‌باشد. برآوردگر δ که مرتبط با عناصر Σ_e و Σ_v می‌باشد، فرض می‌شود که معلوم باشد در حالی که در عمل این‌گونه عمل نمی‌باشد و عبارت δ را باید به روش‌های مختلفی همانند روش گشتاوری، حداقل مربعات تعمیم‌یافته، ماکسیمم درست‌نمایی و ماکسیمم درست‌نمایی مقید برآورده شود و در صورت برآورده شدن جایگزین می‌شود که بهترین برآوردگر پیشگوی ناریب خطی تجربی به‌دست آید.

۳.۸.۲ برآوردگر بیز تجربی EB

برای دستیابی به این نوع از برآوردگر، کافی است از سه مرحله‌ی زیر استفاده کنید:

مرحله‌ی اول: برای به‌دست آوردن توزیع پسین پارامترهای مورد نظر μ بنابر شرط y یعنی $f(\mu|y, \lambda)$ نیاز است، لذا برای به‌دست آوردن آن از دو توزیع شرطی $f(\mu|y, \lambda_1)$ و $f(\mu|y, \lambda_2)$ که در آن بردار پارامتر مدل، مربوط به عبارت $\lambda = (\lambda_1, \lambda_2)$ می‌باشد، استفاده می‌کنند.

مرحله‌ی دوم: بنابر پارامترهای مدل گفته شده یعنی همان λ ها، از تابع چگالی حاشیه‌ای $f(y|\lambda)$ و به کمک روش‌های ماکسیمم درست‌نمایی، گشتاوری و غیره برآورد می‌شوند.

مرحله‌ی سوم: برای به‌دست آوردن پارامتر μ از تابع چگالی پسین برآورد شده‌ی $f(\mu|y, \hat{\lambda})$ استفاده می‌کنند و لذا برآوردگر $\hat{\lambda}$ مرتبط به برآوردگر λ می‌باشد.

به‌طور کلی بهترین پیشگو برای یک کمیت، زمانی به‌دست می‌آید که در روش بیزی با در نظر گرفتن تابع زیان دوم خطاها، امید شرطی کمیت تحت مطالعه برای ناحیه‌ی کوچک بنابر اینکه داده‌ها و پارامترهای مدل داده شده باشند، به‌دست آورده می‌شود و همچنین برای این محاسبات بعضی فرض‌های توزیعی مورد نیاز است.

برآوردگر بیز تجربی با استفاده از برآوردگرهای مناسب پارامترهای مدل و با به‌کارگیری از روش‌های ماکسیمم درست‌نمایی و روش گشتاوری به‌دست می‌آید. بنابراین همانند برآوردگر بهترین پیشگوی خطی ناریب تجربی، مدل (۳۹.۲) را در نظر بگیرید، که دارای توزیع‌های زیر می‌باشند:

$$y_i | v_i \stackrel{iid}{\sim} N(\mathbf{X}_i \beta + \mathbf{Z}_i v_i, \Sigma_{e_i}),$$

$$v_i \stackrel{iid}{\sim} N(0, \Sigma_{v_i}), \quad i = 1, \dots, m.$$

از این رو، Σ_{e_i} و Σ_{v_i} پارامترهای واریانسی می‌باشند که با δ نشان داده می‌شوند، و در معادله بالا $i = 1, \dots, m$ نیز نشان دهنده‌ی نواحی کوچک هستند. اما برای بردارهای ثابت و معین η

و m بنا بر فرض اگر رابطه‌ی زیر برقرار باشد،

$$\mu_i = \eta_i^T \beta + m_i^T v_i$$

آن‌گاه برآورد پیشگوی بیزی به صورت امید شرطی μ_i مشروط بر y_i ، β و δ است که به صورت زیر نشان داده می‌شود:

$$\mu_i^B = \hat{\mu}_i^B(\beta, \delta) = E(\mu_i | y_i, \beta, \delta) = \eta_i^T \beta + m_i^T \hat{v}_i^B; \quad (41.2)$$

و در این برآوردگر

$$v_i^B = E(v_i | y_i, \beta, \delta) = \Sigma_{v_i} Z_i^T V_i^{-1} (y_i - X_i \beta),$$

$$V_i = \Sigma_{e_i} + Z_i \Sigma_{v_i} Z_i^T,$$

بنابراین $\hat{\mu}_i^B$ از توزیع پسین زیر به دست می‌آید:

$$\mu_i | y_i, \beta, \delta \stackrel{iid}{\sim} N(\hat{\mu}_i^B, g_i(\delta))$$

پس

$$g_i(\delta) = m_i^T (\Sigma_{v_i} - \Sigma_{v_i} Z_i^T V_i^{-1} \Sigma_{v_i}) m_i.$$

در رابطه‌ی بالا یعنی مدل (41.2)، دیده می‌شود که این برآوردگر به پارامترهای β و δ وابسته است، به همین منظور لازم می‌شماریم که این پارامترها برآورد شوند. در مرحله‌ی بعدی، این پارامترها برای $i = 1, \dots, m$ از توزیع حاشیه‌ای $y_i \stackrel{iid}{\sim} N(X_i \beta, V_i)$ با به کارگیری از روش‌های ماکسیمم درست‌نمایی یا ماکسیمم درست‌نمایی مقید برآورده می‌شوند و به صورت $\hat{\beta}$ و $\hat{\delta}$ نشان داده می‌شوند و بعد در پارامتر $\hat{\mu}_i^B$ قرار می‌گیرند که بدین صورت برآوردگر بیز تجربی به صورت زیر خواهد بود:

$$\hat{\mu}_i^{EB} = \hat{\mu}_i^{EB}(\hat{\beta}, \hat{\delta}) = \eta_i^T \hat{\beta} + m_i^T \hat{v}_i^B(\hat{\beta}, \hat{\delta}).$$

پس، $\hat{\mu}_i^{EB}$ میانگین برآورد شده‌ی توزیع پسین عبارت $f(\mu_i | y_i, \beta, \delta)$ به شکل $N(\hat{\mu}_i^{EB}, g_i(\hat{\delta}))$ می‌باشد.

۴.۸.۲ برآوردگر بیز سلسله‌مراتبی HB

اغلب اوقات روش بیز سلسله‌مراتبی را به عنوان روش کاملاً بیزی می‌نامند، چرا که توزیع پیشین را برای پارامترهای مدل فرض می‌کنند و سپس نتیجه‌گیری را براساس توزیع پسین انجام می‌دهند. به طور کلی، تحت میانگین مربع خطا، پارامتر نواحی کوچک توسط میانگین‌های پیشین برآورد می‌شوند و عدم حتمیت توسط واریانس پسین اندازه‌گیری می‌شود (گلفند و اسمیت، ۱۹۹۰). روش بیز سلسله‌مراتبی، روش کاملاً ساده و قابل درک می‌باشد. هر چند

نیازمند محاسبات پیچیده همراه با انتگرال‌گیری‌های چند بعدی است. پیشرفت‌های اخیر در ابعاد محاسباتی، به‌ویژه در نمونه‌گیری گیز در رویکرد بیزی قضیه بیز را بسیار پرکاربرد و ساده کرده است. در روش سلسله‌مراتبی برای اینکه بتوان برآوردگری را به‌دست آورد، در ابتدا کافی است توزیع پسین $f(\mu|y)$ را محاسبه کرد. لذا برآوردگر پارامتر مورد نظر $\phi = h(\mu)$ از طریق برآورد میانگین پسین $\hat{\phi}^{HB} = E(h(\mu)|y)$ به‌دست می‌آید و برای دقت نیز از عبارت $V(h(\mu)|y)$ استفاده می‌شود.

همان‌طوری که می‌دانید برای پیدا کردن توزیع پسین $f(\mu|y)$ به جای برآورد λ ، از توزیع پیشین $f(\lambda)$ استفاده می‌کنند که روی مدل تعریف شده است، چرا که با این عمل توزیع پسین مورد نظر $f(\mu|y)$ برای پارامترهای کوچک ناحیه‌ای برآورده می‌شوند. برای به‌دست آوردن توزیع پسین با استفاده از قضیه‌ی بیز، کافی است مدل را در دو مرحله برای داده‌ها و پارامترها با استفاده از توزیع‌های $f(\mu|\lambda_2)f(y|\mu, \lambda_1)$ که $\lambda = (\lambda_1^T, \lambda_2^T)^T$ بردار پارامترهای مدل می‌باشد، صورت گیرد. با بهره‌گیری از قضیه‌ی بیز، برای به‌دست آوردن توزیع پسین پارامترهای ناحیه‌ی کوچک همانند رابطه‌ی (۱.۵) به‌عمل می‌آید. برآوردگر بیز سلسله‌مراتبی را می‌توان با استفاده از امیدریاضی بر روی توزیع پسین، برای پارامترهای کوچک ناحیه‌ای به‌دست آورد. همان‌طوری که در به‌دست آوردن این امیدریاضی، اغلب اوقات با انتگرال‌های پیچیده و با ابعاد بالا مواجه می‌شویم پس، محاسبات آن‌ها سختی نیز به همراه خواهد داشت. برای اینکه بتوان از سختی محاسبات کم کرد، بهتر است از روش‌های دیگر استفاده کرد و به همین دلیل به‌کاربردن از انواع روش‌های مونت کارلوی زنجیره‌ی مارکوفی که شامل نمونه‌گیری گیز است، ملزوم می‌باشد. اساس عملکرد این روش‌ها، ساختن نمونه‌های شبیه‌سازی برای تمامی پارامترهای مدل $\eta = (\lambda_1, \lambda_2, \mu)$ است، به‌طوری که این نمونه‌ها یک زنجیره‌ی مارکوفی $\eta^{(k)}$ ، $K = 0, 1, 2, \dots$ که در آن توزیع $\eta^{(k)}$ به توزیع پسین یکتای $\pi(\eta) = f(\eta|y)$ همگرا می‌شود، را تشکیل می‌دهند. در الگوریتم به‌کار گرفته شده، برای ایجاد نمونه‌ها ابتدا از نقطه‌ی اولیه $\eta^{(0)}$ بعد از انتخاب f نمونه‌ی اولیه، F نمونه‌ی $\eta^{(f+1)}, \dots, \eta^{(f+F)}$ را از توزیع هدف $f(\eta|y)$ به‌دست می‌آورند. و در انتها میانگین نمونه‌های به‌دست آمده، برآورد پارامتر را نتیجه می‌دهند. چند تن از افرادی که در این زمینه کار کرده‌اند به شرح زیر می‌باشد:

بل (۱۹۹۹) با به‌کارگیری از روش بیز سلسله‌مراتبی در مدل سطح ناحیه، برآورد درآمد متوسط و نیز تعداد کودکان زیر خط فقر را در همه‌ی ایالت‌های متحده آمریکا به‌دست آورد. داتا و همکاران (۱۹۹۶) درآمد متوسط خانوارهای ۴ نفره ایالت متحده را با استفاده از مدل‌های HB انجام دادند. و همچنین داتا و همکاران در سال (۱۹۹۹) با این روش بیان شده توانستند نرخ بیکاری را برای همه ایالت‌های متحده به‌دست آورند.

جمع‌بندی

در این فصل سعی بر آن بوده است که مدل‌های برآورد کوچک ناحیه‌ای را معرفی کنیم و همچنین اهمیت هر کدام از این مدل‌ها را با ذکر مثال‌هایی تشریح نمودیم. لذا تمامی مدل‌های بیان شده در این فصل اهمیت دارد ولی، آنچه که برایمان در این پایان‌نامه مهم می‌باشد، به‌کارگیری از مدل‌های فضایی است، چرا که این مدل‌ها موجب افزایش دقت برآوردگر مورد نظر می‌شوند. به همین منظور، در فصل بعدی مدل بیزی برآورد کوچک ناحیه‌ای را مطرح می‌نماییم.

فصل ۳

مدل بندی بیزی برآورد کوچک ناحیه‌ای

برای بالا بردن دقت برآوردگرها و کاهش واریانس از برآوردگرهای غیرمستقیم در ناحیه‌ی کوچک استفاده می‌شود و به کارگیری از روش‌های بیزی نیز راه حل مناسبی برای موارد فوق خواهد بود. در این فصل به بیان اهمیت روش‌های بیزی در برآورد کوچک ناحیه‌ای خواهیم پرداخت. لذا استفاده کردن از رهیافت بیزی، توزیع پیشین مناسبی را برای پارامترهای مدل ارائه خواهد داد. از طرفی با اعمال محدودیت روی پارامتر فضایی و تعریف مناسب ماتریس همسایگی، مدل CAR برای مدل پیوند مطرح می‌شود. بنابراین برای تعیین برآوردهای بیزی، با در نظر گرفتن پیشین مناسب و استفاده از مدل‌های مطرح شده، توزیع پسین پارامترها به دست می‌آید و پارامتر مورد علاقه به روش بیزی برآورد می‌شود. مطالب مطرح شده در این فصل، از مقاله استنباط بیزی برآورد کوچک ناحیه‌ای ریچاردسون و لی در سال (۲۰۱۰) گرفته شده است.

۱.۳ روش‌های بیزی

در دیدگاه کلاسیک، مدل بندی بر اساس تابع درست‌نمایی است. ولی در دیدگاه بیزی، مقدار واقعی پارامتر تحقیقی از یک توزیع تصادفی می‌باشد. در استنباط بیزی، علاوه بر مشخص کردن مدل برای داده‌های مشاهده شده‌ی بردار y به شرط یک بردار از پارامترهای مجهول θ به صورت توزیع $f(y|\theta)$ ، فرض می‌شود که θ یک کمیت تصادفی با توزیع پیشین $\pi(\theta|\eta)$ و η

نیز یک بردار ابر پارامتر می‌باشد. استنباط بیزی بر اساس توزیع پسینی به صورت

$$\pi(\theta|y, \eta) = \frac{\pi(\theta, y|\eta)}{f(y|\eta)} = \frac{\pi(\theta, y|\eta)}{\int f(y, u|\eta) du} = \frac{f(y|\theta)\pi(\theta|\eta)}{\int f(y|u)\pi(u|\eta) du},$$

محاسبه می‌شود که این قضیه به قضیه‌ی بیز بر می‌گردد. اگر η معلوم باشد آن‌گاه در نوشتن $\pi(\theta|y, \eta)$ ، η نوشته نمی‌شود. در دیدگاه بیزی، همه‌ی استنباط‌ها مبنی بر توزیع پسین به دست می‌آید.

انواع توزیع‌های پیشینی

به طور کلی توزیع‌های پیشینی دارای نوع‌های مختلفی می‌باشند، اما تنها ۴ نوع از آن‌ها را شرح می‌دهیم:

۱. توزیع پیشین سره^۱: به توزیعی گویند که در آن

$$\int_{\Theta} \pi(\theta) d\theta = 1$$

باشد.

۲. توزیع پیشین ناسره^۲: به توزیعی گویند که در آن

$$\int_{\Theta} \pi(\theta) d\theta = +\infty$$

باشد.

۳. توزیع پیشین ناآگاهی‌بخش^۳: یکی از مهمترین پیشین‌ها، پیشین ناآگاهی‌بخش می‌باشد و در مواردی که اطلاع چندانی در مورد پارامتر موجود نباشد، مورد استفاده قرار می‌گیرد. در این حالت محقق احتمال یکسانی را به همه پیشامدهای ممکن تخصیص می‌دهد؛ یعنی $\pi(\theta) = k$.

۴. توزیع پیشین مبهم^۴: توزیع پیشین سره‌ای که واریانس خیلی بزرگ دارد، به طوری که به یک پیشین ناآگاهی‌بخش نزدیک می‌شود.

پیشین‌های متفاوتی را می‌توان در تحلیل بیزی در نظر گرفت. انتخاب توزیع پیشین یکی از مسائل مهم در آمار بیزی می‌باشد، زیرا این توزیع تاثیر نسبتاً زیاد و قابل توجهی در نتیجه‌ی استنباط دارد. پس، این انتخاب باید به بهترین نحو انجام شود. بنابر مطالب فوق، تحلیل بیزی، مبنی بر معرفی پیشین‌ها، توزیع پسینی و توزیع‌های شرطی کامل است؛ که مهم‌ترین روش در این پایان‌نامه مدل‌های سلسله‌مراتبی بیزی می‌باشد.

^۱ Proper priori

^۲ ImProper priori

^۳ informative priori

^۴ vague priori

۲.۳ ساختار سلسله‌مراتبی مدل‌بندی بیزی برای نواحی کوچک

فرض کنید $y_{(n)} = (y_1, \dots, y_n)$ بردار نمونه‌گیری شده با اندازه‌ی n است. ساختار مدل‌بندی مدل سلسله‌مراتبی کلی زیر را در نظر می‌گیریم:

$$\text{مرتبه‌ی اول: } y_{(n)} | \mathbf{X} = \mathbf{x} \sim h(y_{(n)} | \mathbf{X} = \mathbf{x}, \theta_1)$$

$$\text{مرتبه‌ی دوم: } \mathbf{X} \sim g(\mathbf{x} | \theta_2)$$

در اینجا $y_{(n)}$ قابل مشاهده اما \mathbf{x} (پنهان) غیرقابل مشاهده است. پارامترهای مورد علاقه برای استنباط، پارامترهای $\theta = (\theta_1, \theta_2)$ هستند. هدف تحلیل برآورد پارامترهای θ و پیشگویی وضعیت پنهان \mathbf{x} است. تابع درست‌نمایی ساختار مدل‌بندی سلسله‌مراتبی به صورت زیر است:

$$L(\theta, y_{(n)}) = \int h(y_{(n)} | \mathbf{x}, \theta_1) g(\mathbf{x}, \theta_2) d\mathbf{x}$$

مشکلات مرتبط در استفاده از این تابع برای استنباط بیشتر، محاسباتی هستند، به این معنا که:

۱. به‌طور کلی محاسبه‌ی تابع درست‌نمایی شامل انتگرال‌هایی با بعد بسیار زیاد است.
۲. به‌دست آوردن مقداری از θ که $L(\theta; y_{(n)})$ را ماکسیمم سازد، با به‌کار بردن روش جستجوی عددی به دلیل طبیعت تصادفی درست‌نمایی برآورد شده بسیار مشکل‌ساز است؛ دقت کنید که طبیعت تصادفی تابع درست‌نمایی در حقیقت ناشی از متفاوت بودن مقادیر نمونه در هر نمونه‌گیری مجدد است.
۳. و در آخر محاسبه‌ی خطاهای معیار به‌دست آمده از پارامترهای برآورد شده در روش محاسبه‌ی عددی مشتقات مرتبه‌ی دوم تابع لگ درست‌نمایی، شامل مشکلات و سختی‌های بیشتری می‌باشد.

۳.۳ بررسی رویکرد بیزی

فرض کنید مشاهدات $y_{(n)} = (y_1, \dots, y_n)$ یک نمونه‌ی n تایی از مدل آماری سلسله‌مراتبی زیر باشد:

$$y_{(n)} \sim f(y_{(n)} | \mathbf{X}, \theta_1)$$

$$\mathbf{X} \sim g(\mathbf{x} | \theta_2)$$

که در آن f و g توابع چگالی توأم، \mathbf{X} بردار کمیت‌های تصادفی یا فرآیندهای موثر بر مشاهدات، و $\theta_1 = (\theta_{11}, \dots, \theta_{1q})$ بردار پارامترهای ثابت و نامعلوم موثر بر مشاهدات، و $\theta_2 = (\theta_{21}, \dots, \theta_{2p})$

بردار پارامترهای ثابت و نامعلوم مرتبط با فرآیند X است. دقت کنید که در ادبیات مدل‌های سلسه‌مراتبی، کمیت‌های تصادفی درون X را اثرات تصادفی، متغیر پنهان یا وضعیت سیستم گویند.

تابع درست‌نمایی برای مدل کلی سلسه‌مراتبی معرفی شده در بالا به صورت زیر است:

$$L(\theta; y_{(n)}) = \int h(y_{(n)}|x, \theta_1)g(x, \theta_2)dx \quad (1.3)$$

که در آن y بردار مشاهدات است (داده‌های مشاهده شده). برآورد ماکسیمم درست‌نمایی پارامترها که با $(\hat{\theta}_1, \hat{\theta}_2) = (\hat{\theta}_{11}, \dots, \hat{\theta}_{1q}; \hat{\theta}_{21}, \dots, \hat{\theta}_{2p})$ نشان داده می‌شود، مقادیر پارامترهای $(\theta_1, \theta_2) = (\theta_{11}, \dots, \theta_{1q}; \theta_{21}, \dots, \theta_{2p})$ هستند که توأم تابع درست‌نمایی را ماکسیمم می‌کنند. واضح است که محاسبه‌ی تابع درست‌نمایی و برآوردهای ماکسیمم درست‌نمایی مشمول کار محاسبه‌ای بسیار سخت و پیچیده با انتگرال‌گیری‌هایی با بعد بالا است. به طور کامل رویکرد بیزی از عهده‌ی انتگرال‌هایی با بعد زیاد بر می‌آید. رویکرد بیزی با این فرض که پارامترهای مدل، متغیرهای تصادفی هستند، شروع می‌شود. توزیع توأم آماری این پارامترها را توزیع پیشینی می‌گویند. توزیع پیشینی باور شخص نسبت به مقادیر پارامترها قبل از جمع‌آوری داده‌ها را نشان می‌دهد. سپس توزیع پیشینی را با تابع درست‌نمایی و با استفاده از قضیه‌ی بیز توأم می‌کنند که در نتیجه‌ی آن را توزیع پسینی می‌نامند. این توزیع باور شخص می‌باشد که پس از جمع‌آوری داده‌ها نسبت به پارامترهای جامعه است.

توجه کنید که در مدل بندی مدل سلسه‌مراتبی، متغیرهای X نامعلوم و تصادفی هستند. توزیع پیشین را اغلب با نماد $\pi(\theta_1, \theta_2)$ نشان می‌دهند. بنابر قضیه‌ی بیز، توزیع پسین توأم کمیت‌های نامعلوم (θ_1, θ_2, X) به شرط مقادیر مشاهده شده‌ی $y_{(n)} = (y_1, \dots, y_n)$ به صورت زیر ارائه می‌دهد:

$$h(\theta_1, \theta_2, X|y_{(n)}) = \frac{f(y_{(n)}|x, \theta_1)g(x|\theta_2)\pi(\theta_1, \theta_2)}{\int f(y_{(n)}|x, \theta_1)g(x|\theta_2)\pi(\theta_1, \theta_2)dXd\theta_1d\theta_2} \quad (2.3)$$

توزیع پسینی حاشیه‌ای پارامترها که با $\pi(\theta_1, \theta_2, X|y_{(n)})$ نشان داده می‌شود به سادگی با انتگرال‌گیری از تابع (۲.۳) نسبت به X به دست می‌آید. در ابتدا به نظر می‌رسد که مدل (۲.۳) مشکل انتگرال‌گیری‌های با بعد زیاد را برای محاسبات درست‌نمایی به یک شکل بزرگتر انتگرال‌گیری‌های با بعد بیشتر در مخرج معادله‌ی (۲.۳) و سپس انتگرال‌گیری از تابع توزیع پسینی حاشیه‌ای را جایگزین می‌کند. اما همان گونه که در زیر توضیح می‌دهیم این گونه نیست. الگوریتم زنجیر مارکوف مونت کارلویی ابزار محاسباتی است که اعداد تصادفی را از توزیع پسینی (۲.۳) تنها با استفاده از صورت کسر بدون این که نیازی به انتگرال‌گیری مخرج وجود داشته باشد را تولید می‌کند. برای اطلاعات بیشتر در این باره به کسلا و جرج (۱۹۹۲)، چیب و گرینبرگ (۱۹۹۵)، گیلکس و همکاران (۱۹۹۵) و رابرت و کسلا (۲۰۱۰) مراجعه کنید. توجه کنید که صورت کسر نیاز به انتگرال‌گیری ندارد. فرض کنید که اعداد تصادفی تولید شده به وسیله‌ی الگوریتم مونت کارلویی زنجیره‌ی مارکوف به صورت $(\theta_1, \theta_2, X)_j$ به ازای $j = 1, \dots, B$

نشان داده شود. در اینجا B تعداد مشاهدات استخراج شده از تابع توزیع (۲.۳) است. این مقدار به اندازه‌ی کافی بزرگ است تا بتواند برآورد مناسبی از (۲.۳) را به دست آورد. همان طوری که می‌دانید با به دست آوردن توزیع (۲.۳) محاسبه‌ی توزیع پسینی $\pi(\theta_1, \theta_2 | y_{(n)})$ نیاز به انتگرال‌گیری نسبت به X در فرمول (۲.۳) دارد. بنابراین چنین انتگرال‌گیری‌هایی لازم نمی‌باشد. توزیع حاشیه‌ای (θ_1, θ_2) به راحتی با حذف مولفه‌ی X از اعداد تصادفی $(\theta_1, \theta_2, X)_j$ به صورت $(\theta_1, \theta_2)_j$ به ازای $j = 1, \dots, B$ به دست می‌آید. به طور مشابه، مقادیر میانگین و واریانس اعداد تصادفی ایجاد شده‌ی $(\theta_1, \theta_2)_j$ به ازای $j = 1, \dots, B$ است. در نتیجه، فرآیند شبیه‌سازی توزیع پسینی حاشیه‌ای نیازی به هیچ انتگرال‌گیری‌ای ندارد.

۴.۳ مدل سلسله‌مراتبی بیزی فضایی

فرض کنید یک مورد از داده‌ی سطح فردی بر روی متغیر هدف و متغیرهای کمکی نمونه بررسی شده در همه‌ی نواحی وجود داشته باشد. توجه کنید، مدل‌ها با داده جمع‌آوری شده، در نظر گرفته می‌شود. برای متغیرهای گسسته، پاسخ می‌تواند با استفاده از توزیع پواسون مدل‌سازی شود (رائو، ۲۰۰۳):

$$y_{ij} | \mu_{ij}, \sigma_e^2 \sim P(\mu_{ij}, \sigma_e^2) \quad (۴.۳)$$

بنابراین μ_{ij} مقدار واقعی (درست) متغیر هدف از فرد j ام در نمونه‌ای از ناحیه‌ی i ام است، σ_e^2 نشان‌دهنده‌ی تغییرات نمونه‌گیری فردی می‌باشد و فرض می‌شود در همه نواحی شدنی است ($\sigma_i^2 = \sigma_e^2$). اگرچه به طور معمول متغیرهای کمکی، به صورت وابستگی مدل بین میانگین و تعداد عوامل توضیحی، تعریف شده است. و این احتمال هم وجود دارد، برخی واریانس باقی‌مانده توسط متغیرهای کمکی غیرقابل توضیح باقی بماند. به طور حتم می‌توان با توجه به اثرات تصادفی در مدل مورد توجه قرار بگیرد. این اثرات تصادفی، الگوهای ناشناخته‌ای را همانند وابستگی فضایی و تغییرات بین نواحی در نظر می‌گیرند. بنابراین، مدل رگرسیون اثرات تصادفی زیر را برای میانگین‌های سطح واحد در نظر می‌گیرد:

$$\mu_{ij} = \alpha + X_{ij}\beta + u_i + v_i \quad (۴.۳)$$

که منجر به مدل میانگین سطح ناحیه می‌شود:

$$\mu_i = \sum_j \frac{\mu_{ij}}{N_i} = \alpha + \bar{X}_i\beta + u_i + v_i \quad (۵.۳)$$

α عرض از مبدا مدل و β برداری از ضرایب متغیرهای کمکی X_{ij} ، u_i اثرات تصادفی که برای تغییرات سطح ناحیه و توزیع مستقل شده زیر به شمار می‌رود:

$$u_i | \sigma_u^2 \sim P(0, \sigma_u^2)$$

لذا u_i نشان‌دهنده‌ی همبستگی اثرات تصادفی فضایی است.

۱.۴.۳ داده‌های شبکه‌ای

یکی از معمول‌ترین مدلی که در برآورد نواحی کوچک استفاده می‌شود، مدل اتورگرسیوشرطی CAR است. این مدل در سال (۱۹۹۱) توسط بی‌سگ و همکارانش معرفی شده است. این مدل به دلیل دارا بودن ساختار ماتریس مجاورت بسیار پرکاربرد خواهد بود. از طرفی این مدل‌ها توسط توزیع‌های احتمال شرطی ساخته می‌شوند. همچنین بسته به موقعیت قرار گرفتن داده‌ها در مدل‌بندی توزیع احتمال نقش مهمی را ایفا می‌کنند. پس ملزوم می‌باشد که نحوه‌ی تاثیر قرار گرفتن این نوع از موقعیت‌ها در مدل‌بندی مورد بررسی قرار گیرد. بیشتر داده‌های نواحی کوچک فضایی شبکه‌ای هستند. از طرفی یکی از پرکاربردترین مدل‌های فضایی برای داده‌های شبکه‌ای و وارد کردن وابستگی فضایی از طریق ساختار همسایگی نواحی، مدل اتورگرسیوشرطی است.

به دلیل به کارگیری این مدل در زمینه‌های مختلف از جمله روش MCMC، که برای برازش رده مشخصی از مدل‌های سلسله‌مراتبی به کار می‌رود، اهمیت می‌یابد. از این‌رو، در مدل‌های فضایی، فرضیه‌ی مستقل بودن برداشته می‌شود و اثرات تصادفی ناحیه v_i به صورت وابسته و تحت تاثیر موقعیت جغرافیایی و همسایگی نواحی قرار می‌گیرد. توزیع‌ها و مدل‌های مختلف فضایی اغلب، از مدل اتورگرسیوشرطی استفاده می‌کنند. تحت این مشخص‌سازی، توزیع شرطی از u_i به شرط مقادیر v_i در تمام نواحی باقی مانده تنها مستلزم همسایگی نواحی می‌باشد.

$$v_i | v_{-i}, \sigma_v^2 \sim N\left(\sum_{j \in \delta_i} \frac{d_{ij}}{|\delta_i|} v_j, \frac{\sigma_v^2}{|\delta_i|}\right) \quad (۶.۳)$$

δ_i مجموعه‌ای از همسایگی‌های ناحیه‌ی i ام و $|\delta_i|$ تعداد همسایگی‌ها در نظر گرفته می‌شود. علاوه بر این، قیدی را اضافه می‌کنند تا مجموع مقادیر همه‌ی اثرات تصادفی v_i برابر صفر شود تا عرض از مبدا و اثرات تصادفی را شناسایی کند (بانرجی و همکاران، ۲۰۰۴).

۲.۴.۳ داده‌های زمین آماری

در داده‌های زمین آماری، فرض می‌شود که $w_i = w(s_i)$ تحقق‌ی از یک میدان تصادفی گاوسی^۵ (GRF) به صورت $\{w(s), s \in S\}$ است، و برداری از یک توزیع نرمال چندمتغیره^۶ $w = (w_1, \dots, w_n)$ است. بنابراین، به عنوان جایگزین برای مشخصات شرطی می‌توانید میانگین μ_{ij} که شامل اثرات تصادفی فضایی w_i و همچنین همبسته است را انتخاب کنید، و از طرفی، مطابق فاصله‌ی d_{kl} که بین دو ناحیه‌ی k و l قرار می‌گیرد، مطابق زیر عمل کنید (دیاگ و همکاران، ۱۹۹۸):

$$\mu_{ij} = \alpha + x_{ij}\beta + w_i \quad (۷.۳)$$

^۵Gaussian random field

^۶Multivariate normality

به‌عنوان نرمال چند متغیره توزیع شده است.

$$w_i | \Sigma \sim MVN(0, \Sigma); \Sigma_{kl} = \sigma_w^2 \exp\{-(\phi d_{kl})\} \quad (۸.۳)$$

σ_w^2 واریانس در هر نقطه و ϕ پارامتر هموارسازی که مقیاس بین نواحی را کنترل می‌کند. برخلاف مدل (۵.۳) اثرات تصادفی مستقل u_i را در مدل (۷.۳) وارد نمی‌کنند. هدف برای انجام مدل (۵.۳) مبنی بر این است که وابستگی فضایی اتورگرسیو شرطی ذاتی، اثرات تصادفی مدل (۶.۳) را به‌وسیله‌ی ساختار همسایگی که از پیش تعیین شده است، تشریح کند. از این رو، اثرات غیرساختاری، همیشه این اجازه را برای یادگیری بیزی درباره‌ی وابستگی فضایی قوی در داده از طریق همبستگی نسبی u_i و w_i پسینی خواهد داد (بی‌سگ و همکاران، ۱۹۹۱ و کارلین و همکاران، ۲۰۰۰). درباره‌ی مدل (۷.۳) یادگیری بیزی درباره‌ی وابستگی فضایی قوی اثرات تصادفی w_i به‌طور مستقیم از طریق برآوردگر پسین پارامتر همبستگی ϕ در مدل (۸.۳) انجام می‌گیرد (مستلزم هیچ همبستگی فضایی نیست، $\phi \rightarrow 0$). در این روش ممکن است اثر تصادفی مستقل جداگانه در مدل (۷.۳) باشد. اما در عمل، می‌تواند منجر به ضعیف شدن توزیع‌های پسین شود (دیاگ و همکاران، ۲۰۰۲). بنابراین برای همه مدل‌ها برآورد سطح ناحیه معقول می‌باشد.

$$\hat{\mu}_{b,i} = E_{\circ|y}[\alpha + \bar{X}_i \beta + z_i] = \hat{\alpha} + \bar{X}_i \hat{\beta} + \hat{z}_i$$

$E_{\circ|y}[\cdot]$ نشان‌دهنده‌ی امیدریاضی پسین و z_i نشان‌دهنده‌ی اثرات تصادفی است که هر دو با $z_i = u_i + v_i$ یا $z_i = w_i$ مشخص شده‌اند. در این مورد به‌ترتیب میانگین‌های پسین $\hat{\alpha}$ ، $\hat{\beta}$ و \hat{z}_i و α ، β و z_i محاسبه می‌شوند. از طرفی هم، فرض کنید که میانگین‌های سطح ناحیه متغیرهای کمکی \bar{X}_i در دسترس‌اند. اساساً مدل (۸.۳) توسط بست و همکاران در سال (۱۹۸۸) مطرح شده است که شامل انواع مختلفی از اثرات تصادفی می‌باشد. واریانس σ_e^2 ناحیه‌ی داخلی را برای همه‌ی نواحی در نظر می‌گیرد که معمولاً غیر واقعی است زیرا، تغییرات فردی بین نواحی متفاوت است. به‌طور کلی، گسترش این مدل در هر ناحیه، واریانس σ_i^2 متفاوتی را در نظر می‌گیرد.

$$y_{ij} | \mu_{ij}, \sigma_i^2 \sim P(\mu_{ij}, \sigma_i^2) \quad (۹.۳)$$

$$\sigma_i^2 \sim \pi(\sigma_i^2),$$

$$\text{Log}(\mu_{ij}) = \alpha + x_{ij} \beta + w_i$$

$\pi(\sigma_i^2)$ یک توزیع پیشین مبهم است. در این مورد، هر واریانس تنها با استفاده از اطلاعات نمونه در ناحیه‌ی i ام برآورد می‌شود. هنگامی که داده بررسی در داخل هر ناحیه به‌طور جداگانه باشد، می‌تواند منجر به برآوردگرهای ضعیف‌تری از σ_i^2 شود. یک جایگزینی می‌تواند این باشد که، با به‌کارگیری از ساختار سلسله‌مراتبی بر روی واریانس‌ها و وام گرفتن از اطلاعات

نواحی بتوان برآوردگرهای استوار را به دست آورد. به طور ویژه، می‌توانید لگاریتم واریانس‌ها را به صورت زیر مدل سازی کنید:

$$y_{ij} | \mu_{ij}, \sigma_i^2 \sim P(\mu_{ij}, \sigma_i^2) \quad (10.3)$$

$$\text{Log}(\sigma_i^2) | \sigma^2 \sim P(0, \sigma^2)$$

$$\sigma_i^2 \sim \pi(\sigma^2)$$

مدل آخر در استفاده از توابع واریانس تعمیم یافته در هموارسازی واریانس‌های سطح ناحیه مشابه است، و به طور کامل مدل پایه می‌باشد. بنابراین، عدم اطمینان درباره‌ی برآوردگرهای واریانس منجر به واریانس پسین در برآوردگرهای کوچک ناحیه‌ای می‌شود. مدل (۹.۳) اساساً همان مدلی است که توسط آرورا و لاهیری در سال (۱۹۹۷) با اثرات تصادفی پیشنهاد شده است. آرورا و لاهیری (۱۹۹۷) با به کارگیری از مدل سطح ناحیه مشابه با اثرات تصادفی مستقل و با هدف رویکرد بیز تجربی برای برآورد اثرات تصادفی پیشنهاداتی را مطرح کرده‌اند. راتو و همکاران (۱۹۹۲) تقریبی از MSE را در مدل سطح واحد با واریانس‌های سطح ناحیه که متفاوت با توزیع پیشین مشترک است، مطرح کرده‌اند.

۳.۴.۳ توزیع پیشین پارامترهای مدل

برای عرض از مبدا α و ضرایب متغیرهای کمکی β پیشین‌های تخت ناسره به کار گرفته می‌شود اما، پسین‌های سره به وجود آمده است. گامای وارونه را همانند توزیع پسین برای هر یک از واریانس‌های $\sigma_u^2, \sigma_v^2, \sigma_w^2, \sigma_i^2$ استفاده می‌کنند. به منظور اطلاعات پیشین مبهم و یادگیری مدل از داده، مقادیر کوچک را برای پارامترهای توزیع گاما استفاده می‌کنند. پیشین برای پارامتر ϕ وابسته به دامنه و مقیاس اندازه‌گیری فاصله‌های بین نواحی کوچک است. گلמן (۲۰۰۶) نشان داد که پیشین‌های معکوس گاما با پارامترهای مقیاس و شکل برای واریانس‌های اثرات تصادفی ممکن است، انقباض دروغین ایجاد کند. به خصوص زمانی که تعداد گروه‌ها کوچک باشد، تعداد اندکی مشاهدات در هر گروه وجود دارد، و از این رو، چندین جایگزین پیشنهاد شده است. همیشه از یک توزیع نیمه‌کوشی بر روی انحراف استاندارد از اثرات تصادفی استفاده می‌کنند. اما در برآوردگرهای کوچک ناحیه‌ای اختلاف معنی‌داری را نخواهد داشت.

مثال ۱.۴.۳. برای داده سطح ناحیه می‌توان مدل را در معادله (۸.۳) نشان داد. پس، به دلیل داشتن اثرات تصادفی و متغیرهای کمکی می‌توان آن را گسترش داد. برای مثال:

$$\hat{Y}_i | \mu_i, \hat{V}_i^2 \sim P(\mu_i, \hat{V}_i^2) \quad (11.3)$$

$$\mu_i = \alpha^* + \bar{X}_i \beta^* + u_i^* + v_i^*$$

u_i^* و v_i^* به ترتیب نشان دهنده ی پواسون و توزیع های CAR می باشند. واریانس های سطح ناحیه در برآوردگرهای میانگین ناحیه به طور فرض شناخته شده اند، لذا با یک مربع نشان داده می شوند. برآورد از میانگین ناحیه به صورت زیر ارائه می شود:

$$\hat{\mu}_{B,i} = E_{o|y}[\alpha^* + \bar{X}_i\beta^* + u_i^* + v_i^*] = \hat{\alpha}^* + \bar{X}_i\hat{\beta}^* + \hat{u}_i^* + \hat{v}_i^*$$

همانند قبل نماد $\hat{\mu}$ نشان دهنده ی میانگین پسین از پارامتر مربوط است. برای مدل فوق $\hat{\alpha}^*$ و $\hat{\beta}^*$ به جای α و β به منظور نشان دادن این عملکرد بیان شده است تا مدل های سطح ناحیه، برآوردگرهای متفاوتی را در مدل های سطح واحد ارائه دهند. به طور مشابه، برآوردگرهای اثرات تصادفی u_i^* و v_i^* مشابه اند ولی متفاوت از u_i و v_i خواهند بود.

۵.۳ ارزیابی کیفیت برآوردگرها

ارزیابی کیفیت برآوردگرهای کوچک ناحیه ای که با مدل های سطح ناحیه و واحد به دست آمده است، می تواند در عمل متفاوت باشد. معمولاً به دست آوردن واریانس ها و میانگین توان دوم خطاهای پیش بینی^۷ مهم می باشد اما، به دست آوردن برآوردگرهای خوب از (MSPE) معمولاً متفاوت بر روش های SAE است. از این رو، دارای فرم بسته ای می باشند که منجر به برآورد پارامترهای مدل نمی شوند. فرمول تقریبی مختلفی پیشنهاد شده است که عبارت اند از: برآوردگرهای بوت استرپ، جک نایف و ... (جیانگ و لاهیری، ۲۰۰۶ و راثو، ۲۰۰۳).

از سوی دیگر، اندازه گیری بیز طبیعی دقیق و واریانس پسین برآوردگرهای کوچک ناحیه ای به طور خودکار از پسین به دست آمده است و به طور کامل عدم اطمینان برای همه ی پارامترهای مدل محاسبه می شود. یک معیار که می تواند برای مقایسه مدل در آمار بیزی به کار رود، معیار انحراف اطلاعات^۸ (DIC) است (اسپیگل هالتر و همکاران، ۲۰۰۲). بر این اساس، انحراف از مدل تاوانیده برای پیچیدگی مدل و تفسیر آن مشابه^۹ (AIC) است. پس، مدل های ترجیح داده می شود که میزان DIC کمتری داشته باشند.

برای مطالعات شبیه سازی زمانی که مقدار میانگین واقعی (درست) \bar{Y}_i شناخته شده باشد، می توانید ارزیابی نسبی^{۱۰} (RB) و ریشه ی میانگین توان دوم خطای نسبی^{۱۱} RRMSE را برای ارزیابی و درستی برآوردگرهای کوچک ناحیه ای محاسبه کنید. توجه کنید RRMSE یعنی،

$$RRMSE = \sqrt{MSE} / (\text{مقدار واقعی})$$

^۷ Mean squared error prediction

^۸ Information deviation criterion

^۹ Akaike information criterion

^{۱۰} Relative bias

^{۱۱} Relative root mean squared error

لذا همانند زیر تعریف شده است:

$$RB_i = \frac{1}{K} \frac{\sum_{k=1}^K (\hat{Y}_i^{(k)} - \bar{Y}_i)}{\bar{Y}_i},$$

$$RRMSE_i = \frac{\sum_{k=1}^K (\hat{Y}_i^{(k)} - \bar{Y}_i)^2}{\bar{Y}_i}$$

k شاخص بررسی نمونه‌ها است.

همانند اندازه‌گیری فراموضعی، می‌توان میانگین مطلق اریبی نسبی و میانگین ریشه‌ی میانگین مربعات توان دوم خطای نسبی^{۱۲} (MRRMSE) را در نظر گرفت:

$$MARB = \frac{1}{m} \sum_{i=1}^m |RB_i|, \quad (۱۲.۳)$$

$$MRRMSE = \frac{1}{m} \sum_{i=1}^m RRMSE_i$$

بنابراین m آخرین شماره از نواحی است. در این روش، به ترتیب می‌توانید ارزیابی کنید که آیا برآوردگرها اریب هستند یا خیر. هر یک از این دو معیار را می‌توانید برای تصمیم‌گیری بر روی بهترین در مدل استفاده کنید. پس، برآوردگری بهتر خواهد بود که مقادیر کوچکتری از میانگین قدرمطلق اریبی نسبی^{۱۳} MARB و میانگین نسبی ریشه میانگین توان دوم خطا MRRMSE را ارائه دهد.

همچنین، اریبی نسبی RBvar و ریشه نسبی میانگین توان دوم خطا RRMSEvar برآوردگرهای واریانس را در اطلاعات شبیه‌سازی محاسبه می‌کنند. علاوه بر این، به بررسی خطای واقعی برآوردگر کوچک ناحیه‌ای که چه میزانی است هم می‌پردازد. لذا همانند زیر تعریف شده است:

$$RBvar_i = \frac{1}{K} \frac{\sum_{k=1}^K (v\hat{ar}(\bar{Y}_i^{(k)}) - EMSE_i)}{EMSE_i}$$

$$RRMSEvar_i = \frac{1}{K} \frac{\sum_{k=1}^K (v\hat{ar}(\bar{Y}_i^{(k)}) - EMSE_i)^2}{EMSE_i}$$

و ($EMSE_i$) خطای تجربی واقعی است. لذا خواهیم داشت:

$$EMSE_i = \frac{1}{K} \sum_{k=1}^K (v\hat{ar}(\bar{Y}_i^{(k)}) - \bar{Y}_i)^2$$

میانگین مطلق اریبی نسبی MARBvar و میانگین نسبی ریشه میانگین توان دوم خطا MRRMSEvar برآوردگرهای واریانس نیز همانند مدل (۱۲.۳) تعریف می‌شود.

^{۱۲} Mean relative root mean squared error

^{۱۳} Mean absolute relative bias

۶.۳ برآورد کوچک ناحیه‌ای در صورت عدم اطلاع از برآوردگرهای مستقیم

به‌منظور کاهش هزینه‌ها، اغلب بررسی‌ها در یک زیر مجموعه از نواحی با نمونه‌گیری از جمعیت انجام می‌شود که نماینده کل نواحی مورد مطالعه است. این بدان معنی است که برآوردگرهای مستقیم می‌تواند تنها برای تعداد کمی نواحی ارائه شود و برآوردگرهایی که برای آن‌ها نمونه در نواحی وجود ندارد باید از روش‌های دیگری به‌دست آورده شوند. بنابراین می‌توانید نواحی را به دو قسمت تقسیم کنید: ۱. نواحی داخل نمونه ۲. نواحی خارج نمونه. در مدل‌های سطح ناحیه، مقادیر از دست رفته \hat{Y}_i و \hat{x}_i برای نواحی خارج نمونه وجود خواهد داشت ضمن اینکه، در مدل‌های سطح واحد مقادیر از دست رفته y_{ij} و x_{ij} در نواحی خارج نمونه وجود دارد. از طرفی فرض می‌شود، متغیرهای کمکی سطح ناحیه \bar{x}_i برای همه‌ی نواحی در دسترس‌اند و لذا، در بررسی‌ها از منابع متفاوتی به‌دست آورده می‌شوند. این یک فرض کلیدی است تا بتوانید برآوردهای نواحی کوچک را برای همه نواحی ارائه کنید. توجه کنید در اینجا فقط به مشکلات داده‌ای که در طرح بررسی از دست رفته است، پرداخته می‌شود.

برای مثال، عدم پاسخ در بررسی‌ها، یکی دیگر از منابع مشترک داده از دست رفته در برآوردهای کوچک ناحیه‌ای است اما در اینجا این مسئله برطرف نمی‌شود. یک رویکرد ساده در برخورد با مشکل عدم مشاهدات مستقیم در نواحی مختلف، به‌کارگیری از یک مدل رگرسیونی بر روی برخی متغیرهای کمکی است که از داده بررسی استفاده می‌شود. برآوردها برای نواحی خارج نمونه با تکیه بر مدل برآورده شده‌اند و اطلاعات اضافی محاسبه می‌شوند (همانند، متغیرهای کمکی سطح ناحیه). مشکل اصلی این روش این است که مقادیر جانپی شده برای عدم اطمینان در برآورد ضرایب رگرسیونی یا همبستگی فضایی بین متغیر هدف در نواحی مختلف محاسبه نمی‌شوند.

۱.۶.۳ همبستگی اثرات تصادفی فضایی

اگر هر یک از مدل‌های بخش (۴.۳) و مدل‌های مثال (۱.۴.۳) را در نظر بگیرید، اثرات تصادفی همبستگی فضایی می‌تواند علاوه بر متغیرهای کمکی هنگام پیش‌گویی برآوردهای کوچک ناحیه‌ای در نواحی خارج نمونه مورد محاسبه قرار بگیرد. این رویکرد هنگامی اهمیت پیدا می‌کند که، وام گرفتن اطلاعات نواحی همسایه‌ها به کمک برآوردگرهای مستقیم از دست‌رفته باشد.

روش اطلاعات برای نواحی که با مشاهدات غیرمستقیم از سایر نواحی وام گرفته است، به شرح زیر می‌باشد. همچنین اگر می‌خواهید برآورد در نواحی خارج نمونه را داشته باشید از مدل‌های سطح ناحیه که در معادله‌ی (۱۱.۳) مطرح شده است، استفاده کنید، بنابراین

می‌تواند همانند مدل زیر بیان شود:

$$\begin{bmatrix} \hat{Y}_s \\ \mu_s \end{bmatrix} = \alpha + \begin{bmatrix} \bar{X}_s \\ \bar{X}_s \end{bmatrix} \beta + \begin{bmatrix} z_s \\ z_s \end{bmatrix} + \begin{bmatrix} e_s \\ 0 \end{bmatrix} \quad (13.3)$$

z نشان دهنده‌ی همبستگی اثرات تصادفی فضایی و زیر نویس s نشان دهنده‌ی مشاهدات (داخل نمونه) نواحی و \underline{s} نواحی مشاهده نشده است. مقدار z_s را می‌توان توسط همبستگی (فضایی) با z_s برآورد کرد.

۲.۶.۳ مشخص سازی نرمال چندمتغیره

هنگامی که بردار کاملی از اثرات تصادفی فضایی $z = w$ دارای توزیعی باشد، آن گاه همانند مدل (۸.۳) نشان داده می‌شود. توزیع شرطی از $w_s | w_{\underline{s}}$ را طبق رابطه‌ی زیر نشان می‌دهند (دیاگ و همکاران، ۱۹۹۸):

$$MVN(\Sigma_{ss} \Sigma_{ss}^{-1} w_s, \Sigma_{ss} - \Sigma_{ss}^T \Sigma_{ss} \Sigma_{ss})$$

برآورد از \hat{w}_s هنگامی است که امید ریاضی پسین میانگینی از MVN شرطی باشد، $E_{\cdot|y}[\Sigma_{ss} \Sigma_{ss}^{-1} w_s]$. پس، برآوردگر نهایی برای مجموعه‌ای از نواحی در \underline{s} همانند مدل زیر است:

$$\hat{Y}_s = \hat{\alpha} + \hat{\beta} \bar{X}_s + \hat{w}_s$$

این برآوردگر می‌تواند به عنوان نسخه‌ی بهبود یافته‌ی برآوردگر ترکیبی و احتمال اینکه اریبی در این برآوردگر کاهش یابد، در نظر گرفته شود.

۳.۶.۳ مشخص سازی اتورگرسیو شرطی

پیش گویی اثرات تصادفی، w_s ، در نواحی خارج نمونه با استفاده از مشخص سازی مدل نرمال چندمتغیره (۸.۳) غیرمبهم است. می‌توان به این علت بیان کرد که توزیع پیوند (ربط) از $w = (w_s, w_{\underline{s}})$ و پیش گویی از $w_s | w_{\underline{s}}$ منحصر به فرد تعیین شده است. در مقابل، پیش گویی اثرات تصادفی v_s CAR است، بنابراین توزیع پیوند برای بردار کامل $v = (v_s, v_{\underline{s}})$ تعریف نشده است (بانرجی و همکاران، ۲۰۰۴). در عوض، پیش گویی فرآیندهای که به طور مستقیم با توزیع‌های شرطی $v_s | v_{\underline{s}}$ انجام شده است، به خوبی نیز تعریف می‌شود، اما با این تفاوت که در عمل منحصر به فرد نمی‌باشد. این نکته توسط بانرجی و همکاران (۲۰۰۴) مطرح شد. برای این مورد یک مدل اتورگرسیو شرطی بر روی نقطه‌ای از سطح داده برآزش داده شده است.

پیش گویی در یک مکان جدید می‌تواند در ساختار یک مدل CAR برای مجموعه‌ای از مکان‌های مشاهده شده انجام شود. و به طور جداگانه، توزیع شرطی تعیین شده از مکان جدید با توجه به مشاهدات یا با ایجاد یک مدل CAR برای مجموعه کاملی از مکان‌های مشاهده شده

جدید نیز انجام گیرد. هر دو روش معتبر است، اما پیش‌گویی توزیع‌های مختلف مقدم‌تر است. در مورد داده سطح ناحیه، منطقی نیست که رویکرد قبلی را در نظر گرفت، بنابراین مشخص نیست که چطور ساختار مجاورت فقط نواحی خارج نمونه را در نواحی چشم‌پوشی می‌کند. بنابراین مدل CAR (۶.۳) برای مجموعه کاملی از اثرات تصادفی فضایی که در نواحی داخل نمونه و در خارج نمونه که $v = (v_s; v_g)$ است، را مشخص می‌کنند. و به سادگی پاسخ داده که در نواحی خارج نمونه از دست رفته باشد را مورد بررسی قرار می‌دهد. این عمل منجر به مجموعه تغییر یافته از توزیع‌های شرطی کامل برای اثرات تصادفی فضایی در نواحی خارج نمونه در طرح MCMC و نیز با به‌کارگیری از برآوردگر توزیع پسین انجام می‌شود. اگرچه مشخص‌سازی برای پیش‌گویی اثرات تصادفی خارج نمونه مشخص است ولی زمانی که مشاهدات مستقیم از چندین نواحی باشند، اغلب با مشکلات زیر مواجه می‌شوید:

۱. برآزش v_i در نواحی داخل نمونه ممکن است مشکل‌ساز باشد، اگر بیشتر یا همه‌ی همسایگی نواحی خارج نمونه باشد.

۲. به‌طور مشابه هنگامی که v_i در نواحی خارج نمونه پیش‌گویی می‌شود، ممکن است نواحی با مقدار کم یا همسایگی‌های داخل نمونه وجود نداشته باشند.

در اصل مدل می‌تواند برآزش داده شود ولی، هنگامی اتفاق می‌افتد که در آن‌جا نواحی با همسایگی از دست رفته باشد، اما اگر تعدادی از آن‌ها از دست رفته باشد شیوه‌ی برآورد خیلی پایدار نخواهد بود. اثرات تصادفی غیرفضایی u هم می‌تواند در مدل اثر گذارد، اما باید تمامی اعضای u_g به صفر برسد تا از مشکلات شناسایی‌پذیری با v_g دوری کند. به هر حال، اختلاف (تفاوت) اندازه‌گیری اثرات تصادفی بین پیش‌گویی اثرات داده و پاسخ در مدل غیرساختاری است. اگر پاسخ نقض شود، اثرات تصادفی u نمی‌تواند همبستگی بین نواحی را داشته باشد. پس برآورد مقیاس مقادیر در u_g غیرممکن می‌شود.

۴.۶.۳ وام گرفتن اطلاعات در سطوح اجرایی بالاتر

همان‌طوری که قبلاً اشاره شد، هنگامی که تعدادی داده در داخل نمونه نباشد، ممکن است پیش‌گویی اثرات تصادفی فضایی در سطح بالا با استفاده از یک مدل CAR مشکل‌ساز باشد زیرا، اطلاعات از همسایگی نواحی وام گرفته است و به اندازه‌ی کافی این احتمال وجود دارد که همسایگی داخل نمونه در دسترس برآورد طرح فضایی نباشد. علاوه بر این مشکل، مشخص‌سازی مدل CAR در سطح جغرافیایی بالاتر می‌تواند با تنگ بودن در داده استفاده شود. فرض می‌شود، همه‌ی نواحی از یک قسمت از M سطوح اجرایی بالاتر و همچنین هر کدام از آن‌ها از چند نواحی ساخته شده باشد. از طرفی، اثر تصادفی r_k ، $k = 1, \dots, M$ را به‌جای اثر فضایی v_i اضافه می‌کنند، بنابراین میانگین در مدل‌های سطح ناحیه همانند زیر تجزیه

می‌شود، پس خواهیم داشت:

$$\mu_i = \alpha + \beta \bar{X}_i + u_i + r_{k(i)} \quad (14.3)$$

شاخص $k(i)$ نشان‌دهنده‌ی سطوح اجرایی بالاتر در ناحیه‌ی i ام است. یک رابطه مشابه را می‌توان برای ساختن μ_{ij} در مدل سطح ناحیه پیشنهاد کرد. ساختار فضایی با ویژگی تغییرات فضایی بزرگ مقیاس را در بر می‌گیرد. یک پیشین CAR مشابه آنچه که برای v_i ، تعیین کرده‌ایم را فرض می‌کنیم، اما مجاورت، مطابق سطح اجرایی بالاتر تعریف می‌شود و فرض می‌شود حداقل یک ناحیه در هر منطقه نمونه‌برداری می‌شود. حال اگر الگوی فضایی در این سطح بالاتر، خیلی ضعیف یا اصلاً وجود نداشته باشد اثرات تصادفی $r_{k(i)}$ ممکن است، یک توزیع غیرفضایی همانند $P(0, \sigma_r^2)$ را تعیین کنند.

جمع بندی

به‌طور خلاصه می‌توان گفت در این فصل، بنابر اهمیت مدل‌های بیزی و ساختار سلسه‌مراتبی، می‌توان آن‌ها را در برآورد کوچک ناحیه‌ای لحاظ نمود. با توجه به اینکه، برآورد کوچک ناحیه‌ای نیز، به‌دلیل وابستگی بین نواحی حالت فضایی بودن را در خودش دارد این موضوع کمکی بر آن خواهد داشت که از داده‌های شبکه‌ای استفاده کنیم. لذا مشخصاتی که در این مدل وجود دارد را به‌طور کامل بیان نمودیم و به کمک مدل اتورگرسیو شرطی توانسته‌ایم در فصل بعدی داده‌های واقعی خود را به کمک مدل سلسه‌مراتبی برازش دهیم. نتایج حاصل از داده‌ها نیز در فصل آتی بیان می‌کنیم.

فصل ۴

کاربرد مدل فضایی بیزی برای برآورد کوچک ناحیه‌ای داده‌های بیمه استان گیلان

۱.۴ مقدمه

هدف از این فصل، تحلیل تعداد بیمه‌شدگان اجباری در شهرستان‌های استان گیلان است. با توجه به ماهیت شمارشی پاسخ در این مجموعه داده و وجود اطلاعات مکانی همسایگی شهرستان‌ها، از رده مدل‌های آمیخته خطی تعمیم‌یافته فضایی استفاده می‌کنیم. آماردان‌های مختلفی از این مدل‌ها در بحث نواحی کوچک استفاده کرده‌اند. به‌عنوان دو نمونه می‌توان به کلینشمیت و همکاران (۲۰۰۰) و ژانگ و لی (۲۰۱۱) اشاره کرد. استنباط مبتنی بر درست‌نمایی در رده SGLMMs، به‌جز نسخه‌های خطی آن، به دلیل وجود انتگرال‌های پیچیده و معمولاً با بعد بالا برای محاسبه تابع درست‌نمایی مدل، بسیار مشکل و چالش‌برانگیز است (بریسلو و کلیتون، ۱۹۹۳؛ مک کالاک، ۱۹۹۷). به همین دلیل، معمولاً برای تحلیل و استنباط در این مدل‌ها از دیدگاه بیزی استفاده می‌شود (دیگل و همکاران، ۱۹۹۸). دلیل معمول برای انتخاب دیدگاه بیزی، برآزش ساده مدل به خاطر وجود الگوریتم‌های مونت کارلوی زنجیر مارکوفی MCMC است (رابرت و کسلا، ۲۰۰۵). در این فصل، از الگوریتم‌های MCMC برای

برازش یک مدل SGLM در نواحی کوچک (شهرستان‌های استان گیلان) برای تحلیل تعداد بیمه‌شدگان اجباری، استفاده می‌کنیم. در ادامه، در بخش (۲.۴) به معرفی داده‌ها و در بخش بعدی نیز به معرفی مدل بیزی و در آخر نیز به برازش و تحلیل مدل پیشنهادی می‌پردازیم.

۲.۴ معرفی داده‌ها

در صنعت بیمه، بنا بر مبنای قانونی، اهداف و مقررات حاکم، دو دسته بیمه به نام‌های، بیمه‌های اجتماعی و بیمه‌های بازرگانی وجود دارند.

تعریف ۱.۲.۴. (بیمه‌های اجتماعی یا بیمه‌های اجباری) این نوع بیمه، به بیمه‌های ناشی از قانون نیز معروف هستند که بیشتر به کارگران و اقشار کم‌درآمد جامعه مربوط می‌شوند. قانون‌گذار برای اقشاری از جامعه که از یک سو نیروی تولیدی جامعه محسوب می‌شوند و از سوی دیگر خود به فکر آینده و معیشت خود نیستند یا برای آینده خود نمی‌توانند برنامه مناسبی داشته باشند، دولت را موظف کرده است که برای حمایت از آن‌ها، بیمه‌های اجتماعی را تعریف کند و این اقشار را زیر چتر این نوع بیمه‌ها قرار دهد.

تعریف ۲.۲.۴. (بیمه‌های بازرگانی) این نوع بیمه، به بیمه‌های اختیاری نیز معروف هستند، به میل و اراده خود بیمه‌گذار تعیین می‌شوند. در این نوع بیمه‌ها تعهد دوطرفه است؛ یعنی بیمه‌گر به ازای دریافت حق بیمه از بیمه‌گذار، تأمین پوشش‌های بیمه‌ای مشخصی را در اختیار وی قرار می‌دهد.

بنا به اهمیت نوع بیمه اجباری، در این پایان‌نامه داده‌های بیمه اجباری استان گیلان را بر اساس یک مدل بیزی نواحی کوچک فضایی تحلیل می‌کنیم. در این مجموعه داده، متغیر پاسخ تعداد بیمه‌شدگان اجباری در ۱۶ شهرستان استان گیلان است. همچنین متغیرهای تبیینی شامل تعداد خانوارها، جمعیت و مستمری تبعی در هر شهرستان هستند. نمایی از ناحیه فضایی تحت مطالعه و نقشه همسایگی شهرستان‌های استان گیلان در شکل ۱.۴ مشاهده می‌شود.



شکل ۱.۴: نقشه همسایگی شهرستان‌های استان گیلان

۳.۴ معرفی مدل بیزی

با توجه به شمارشی بودن متغیر پاسخ، از توزیع پواسون با تابع پیوند لگاریتمی برای مدل بندی آن استفاده کردیم. چرا که همان طوری که اشاره داشتیم، داده‌های بیمه نیز ماهیت شمارشی دارند. بنابراین با استفاده از مدل (۳۵.۲) برای داده‌های مورد نظر می‌توان مدل پیشنهادی زیر را در نظر گرفت.

$$\log(\mu(s_i)) = \eta(s_i) = \beta_0 + \beta_1 x_1(s_i) + \beta_2 x_2(s_i) + \beta_3 x_3(s_i) + \mathbf{v}(s_i) + \epsilon(s_i), \quad i = 1, \dots, 16 \quad (1.4)$$

و در آن x_1 تعداد خانوار، x_2 تعداد جمعیت و x_3 تعداد مستمری تبعی را نشان می‌دهد. برای برازش مدل همه متغیرهای تبیینی را استاندارد کردیم. برای تشکیل مدل بیزی، باید توزیع‌های پیشین پارامترها را مشخص کنیم. در مدل (۱.۴) برای همه پارامترهای رگرسیونی توزیع پیشین نرمال $N(0, 100)$ را در نظر گرفتیم و برای دو پارامتر رگرسیونی σ_v^2 و σ_ϵ^2 توزیع گامای معکوس با مقادیر ابرپارامتر 0.001 و 0.001 را تعیین کردیم. بنابراین مدل بیزی کامل

برای داده‌های بیمه را می‌توان به صورت مدل سلسله‌مراتبی زیر نوشت:

$$y(s_i) \sim Pois(\mu(s_i))$$

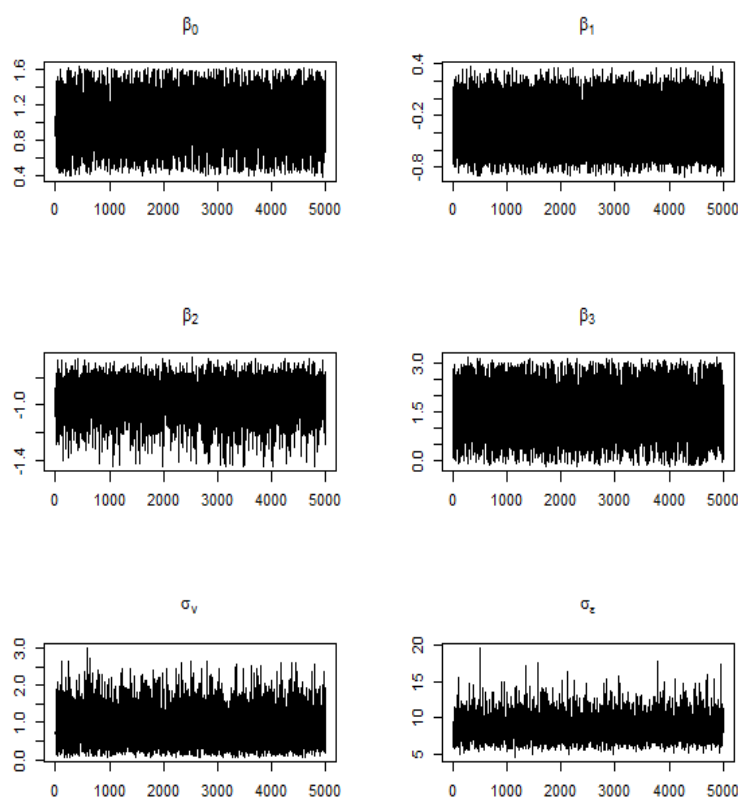
$$\log(\mu(s_i)) = \eta(s_i)$$

$$v \sim CAR(0, \sigma_v^2 A^{-1})$$

$$\epsilon(s_i) \sim N(0, \sigma_\epsilon^2)$$

۴.۴ برآزش مدل و تحلیل نتایج

برنامه BUGS در سال (۱۹۸۹) توسط اندرو توماس و دیوید اشپیگل‌هالتر معرفی شد. برای برآزش مدل، از یک الگوریتم MCMCs استفاده کردیم که در نرم‌افزار Openbugs اجرا شد. یک نمونه مونت کارلویی به حجم ۱۰۰۰۰ تولید کردیم که دوره سوزاندن الگوریتم را ۵۰۰۰ انتخاب کردیم. همچنین طول گام انتخاب نمونه‌ها را ۱۰ در نظر گرفتیم که از وابستگی نمونه‌ها بکاهیم. بنابراین همه نتایج بر اساس ۱۰۰۰ نمونه نهایی گزارش شده‌اند. نمودارهای اثر نمونه‌های تولید شده برای پارامترهای توزیع پسین در شکل ۲.۴ نمایش داده شده‌اند. به سادگی می‌توان دریافت رفتار آمیختگی زنجیر برای همه پارامترهای مدل خوب است.



شکل ۲.۴: نمودارهای اثر پارامترهای مدل بیزی داده‌های بیمه

برآورد پارامترهای مدل به همراه نواحی اعتبار (CI) ۹۵ درصد متناظر هر کدام در جدول ۱.۴ گزارش شده‌اند. با توجه به اعداد جدول، می‌توان گفت متغیر تعداد خانوار بر پاسخ تاثیر معنی‌داری ندارد. اما در مقابل جمعیت هر شهرستان اثر منفی معنی‌دار بر تعداد بیمه‌های اجباری دارد و تعداد مستمری تبعی اثر مثبت معنی‌دار دارد. از طرفی، حضور اثر فضایی در مدل لازم است. اما واریانس بزرگ جمله خطا در مقابل اثر فضایی، حاکی از آن است که اطلاعات مهم دیگری هستند که بر تعداد بیمه‌های اجباری اثر دارند ولی نادیده گرفته شده‌اند. بنابراین، تحلیل این داده‌ها و برآوردهای نواحی کوچک برای تعداد کل بیمه‌شدگان در سطح شهرستان‌ها نیازمند دقت بیشتر و در نظر گرفتن سایر اطلاعات مهم است.

جدول ۱.۴: نتایج برازش مدل بیزی بر روی داده‌های بیمه اجباری استان گیلان

اثر	برآورد	انحراف معیار	کران پایین CI	کران بالا CI
β_0	۱/۰	۰/۳	۰/۵	۱/۶
β_1	-۰/۴	۰/۳	-۰/۸	۰/۲
β_2	-۰/۹	۰/۱	-۱/۳	-۰/۷
β_3	۱/۳	۰/۸	۰/۰	۲/۹
σ_v^2	۰/۹	۰/۵	۰/۱	۱/۹
σ_e^2	۸/۶	۱/۶	۶/۰	۱۲/۴

۵.۴ نتیجه‌گیری و آینده تحقیق

در این پایان‌نامه، از یک مدل CAR برای مدل‌بندی داده‌های فضایی شبکه‌ای استفاده کردیم. برازش مدل بیزی مطرح شده نیز به کمک الگوریتم‌های MCMC انجام شد. اهمیت این نوع مدل‌ها برای مدل‌بندی داده‌های فضایی نواحی کوچک، به دلیل عملی و مناسب بودن برازش و تفسیر ساده مدل، در کاربردهای واقعی قابل توجه است. با توجه به مطالبی که در این پایان‌نامه مطرح شد، برای آینده تحقیق می‌توان استنباط بیزی برآورد کوچک ناحیه‌ای را با مدل‌های ICAR و MCAR انجام داد که به نوبه خود می‌تواند جالب توجه و اهمیت داشته باشد. و همچنین استفاده از مدل‌های شمارشی مانند مدل‌های دوجمله‌ای، بتا-دوجمله‌ای، برنولی، دوجمله‌ای منفی و وایبل گسسته و نیز به‌کارگیری از الگوریتم INLA می‌تواند مفید و پرکاربرد باشد.

مراجع

- [۱] محمدزاده م، (۱۳۹۴)، «آمارفضایی» چاپ اول، مرکز نشر آثار علمی دانشگاه تربیت مدرس، تهران.
- [2] Anderson, T.W., and Hsiao, C. (1981), Formulation and estimation of dynamic models using panel data, *Journal of Econometrics* 18, 67–82.
- [3] Arora, V., and Lahiri, P. (1997), On the superiority of the bayesian method over the BLUP in small area estimation problems, *Statistica Sinica*, 7, 1053–1063.
- [4] Arora, V., Lahiri, P. and Mukherjee, K. (1997), Empirical bayes estimation of finite population means from complex surveys, *Journal of the American Statistical Association*, 92(440), 1555–1562.
- [5] Brewer, K.R.W. (1963), Ratio estimation and finite populations: some results deducible from the assumption of an underlying stochastic process, *Australian Journal of Statistics*, 5, 5–13.
- [6] Brackstone, G.J. (1987), Small area data: Policy issues and technical challenges, in R. Platek, J.N.K. Rao, C.E., Sarndal and M.P., Singh (Eds.), *Small Area Statistics*, New York: Wiley, pp. 3–20.
- [7] Battese, G.E., Harter, R.M., and Fuller, W.A. (1988), An error-components model for prediction of county crop areas using survey and satellite data, *Journal of the American Statistical Association*, 83, 28–36.
- [8] Besag, J., J.C., York, and A., Molline (1991), Bayesian image restoration, with two applications in spatial statistics, *Annals of the Institute of Statistical Mathematics* 43, 1-59.
- [9] Breslow, N., and Clayton, D. (1993), Approximate inference in generalized linear mixed models, *Journal of the American Statistical Association*, 88, 9–25.
- [10] Bell, W.R., (1999), Accounting for uncertainty about variances in small area estimation, *Bulletin of The International Statistical Institute*,

- [11] Banerjee, S., B.P., Carlin, and A.E., Gelfand (2004), Hierarchical Modeling and Analysis for Spatial Data, Chapman Hall-CRC, Boca Raton, Florida.
- [12] Best, N., Richardson, S., and Thomson, A., (2005), A comparison of bayesian spatial models for disease mapping, *Statistical Methods in Medical Research*, 14, 35–59.
- [13] Banerjee, S., Carlin, B.P., and Gelfand, A.E., (2014), Hierarchical Modeling and Analysis for Spatial Data. CRC Press, New york.
- [14] Cressie, N. (1991), Small-area prediction of undercount using the general linear model, *Proceedings of Statistics Symposium 90: Measurement and Improvement of Data Quality*, Ottawa: Statistics Canada, pp. 93–105.
- [15] Casella, G., and George, E.I., (1992), Explaining the gibbs sampler, *The American Statistician* 46(3), 167-174, 39, 41.
- [16] Cressie, N. (1993), *Statistics for Spatial Data*, New York: Wiley.
- [17] Chaudhuri, A. (1993), Small domain statistics: A review, *Statistica Neerlandica* ,48, 215–236.
- [18] Chib, S., and Greenberg, E. (1995), Understanding the metropolis-hastings algorithm, *American Statistician*, 49, 327–335.
- [19] Cochran, W.G., (1977), *Sampling techniques* (3rd Eds.), John Wiley, Sons, Inc., New York.
- [20] Carlin, B.P., and Louis, T.A., (2000), *Bayes and Empirical Bayes Methods for Data Analysis* (2nd Eds.), London: Chapman and Hall.
- [21] Datta, G.S., and Ghosh, M., (1991), Bayesian prediction in linear models: Applications to small area estimation, *The Annals of Statistics*, 17, 48–1770.
- [22] Datta, G.S., Ghosh, M., Nangia, N., and Natarajan, K., (1996), Estimation of median income of four-person families: A bayesian approach, in W.A. Berry, K.M., Chaloner, and J.K., Geweke (Eds.), *Bayesian Analysis in Statistics and Econometrics*1, New York: Wiley, pp. 129–140.
- [23] Diggle, P. J., Tawn, J. A., and Moyeed, R. A. (1998). Model-based geostatistics, *Journal of the Royal Statistical Society: Series C Applied Statistics*, 47(3), 299–350.

- [24] Datta, G.S., Lahiri, P., Maiti, T., and Lu, K.L. (1999), Hierarchical bayes estimation of unemployment rates for the u.s., states, *Journal of the American Statistical Association*, 94, 1074–1082.
- [25] Datta, G.S., Kubakawa, T., and Rae, J.N.K., (2002), Estimation of MSE in small area estimation, Technical Report, Department of Statistics, University of Georgia, Athens.
- [26] Diggle, P., Moyeed, R., Rowlingson, B., and Thomson, M., (2002), Childhood malaria in the gambia: A case-study in model-based geostatistics, *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 51(4), 493–506.
- [27] Eberly, L.E., and Carlin, B.P., (2000). Identifiability and convergence issues for markov chain monte carlo fitting of spatial models, *Statistics in Medicine*, 19(17–18), 2279–2294.
- [28] EURAREA Consortium (2004), Project reference volume. Technical Report, EURAREA Consortium .
- [29] Fuller, W.A., and Battese, G.E., (1973), Transformations for estimation of linear models with nested-error structure, *Journal of the American Statistical Association*, 68, 626–632.
- [30] Fay, R.E., and Herriot, R.A., (1979), Estimates of income for small places: An application of james-stein procedures to census data, *Journal of the American Statistical Association*, 74(366a), 269–277.
- [31] Fuller, W.A., (1989), Prediction of true values for the measurement error model, in conference on statistical analysis of measurement error models and applications, Humboldt State University.
- [32] Fuller, W.A., (1999), Environmental surveys over time, *Journal of Agricultural, Biological and Environmental Statistics*, 4, 331–345.
- [33] Gonzalez, J.F., Placek, P.J., and Scott, C., (1966), Synthetic estimation of followback surveys at the national center for health statistics, in W.L., Schaible (Eds.), *In direct Estimators in U.S. Federal Programs*, Springer-Verlag: New York, pp. 16–27.
- [34] Gonzalez, M.E., and Wakesberg, J., (1973), Estimation of the error of synthetic estimates, Paper presented at the first meeting of the International Association of Survey Statisticians, Vienna, Austria.
- [35] Gonzalez, M.E., (1973), Use and evaluation of synthetic estimates, proceedings of the social statistics section, American Statistical Association, pp. 33–36.

- [36] Gonzalez, M.E., and Hoza, C., (1978), Small-area estimation with application to unemployment and housing estimates, *Journal of the American Statistical Association*, 73, 7–15.
- [37] Gelfand, A.E., and Smith, A.F.M., (1990), Sample-based approaches to calculating marginal densities, *Journal of the American Statistical Association*, , 85, 972–985.
- [38] Gelfand, A.E., and Smith, A.F.M., (1991), Gibbs sampling for marginal posterior expectations, *Communications in Statistics-Theory and Methods*, 20, 1747–1766.
- [39] Gilks, W.R., and Wild, P., (1992), Adaptive rejection sampling for gibbs sampling, *Applied Statistics*, 41, 337–348.
- [40] Ghosh, M., (1992), Constrained bayes estimation with applications, *Journal of the American Statistical Association*, 87(418), 533–540.
- [41] Ghosh, M., and Nangia, N., (1993), Estimation of median income of fourpers on families: A bayesian time series approach, technic report, depar, tment of Statistics, University of Florida, Gainesville.
- [42] Ghosh, M., and J.N.K., Rao (1994), Small area estimation: An Appraisal, In: *Statistical Science* 9 (1), pp. 55–93.
- [43] Greenland, S. and J., Robins (1994), Ecologic studies biases, misconceptions, and counterexamples, *American Journal of Epidemiology*, 139 (8), 747–760.
- [44] Gilks, W.R., Richardson, S., and Spiegelhalter, D. (Eds.) (1995), *Markov chain Monte Carlo in practice*, CRC press.
- [45] Gilks, W.R., Richardson, S., and Spiegelhalter, D.J., (Eds.) (1996), *Markov chain monte carlo in practice*, London: Chapman and Hall.
- [46] Ghosh, M., Nangia, N., and Kim, D. (1996), Estimation of median income of four-person families: A bayesian time series approach, *Journal of the American Statistical Association*, 91, 1423–1431.
- [47] Ghosh, M., Natarajan, K., Stroud, T.W.F., and Carlin, B.P., (1998), Generalized linear models for small-area estimation, *Journal of the American Statistical Association*, 93(441), 273–282.
- [48] Gelman, A., Carlin, J.B., Stern, H.S., and Rubin, D.B., (2003), *Bayesian data analysis*, Chapman and Hall/CRC Texts in Statistical Science.

- [49] Gelman, A. (2006), Prior distributions for variance parameters in hierarchical models comment on article by browne and draper. *Bayesian Analysis*, 1(3), 515–534.
- [50] Hastings, W.K., (1970), Monte carlo sampling methods using markov chains and their applications, *Biometrika*, 57, 97–109.
- [51] Hansen, M.H., Madow, W.G., and Tepping, B.J., (1983), An evaluation of model-dependent and probability sampling inferences in sample surveys, *Journal of the American Statistical Association*, 78, 776–793.
- [52] Hedayat, A.S., and Sinha, B.K., (1991), *Design and inference in finite population sampling*, New York: Wiley.
- [53] Harville, D.A., (1991), Comment, *Statistical Science*, 6, 35-39.
- [54] Hodges J, Sargent D., (2001), Counting degrees of freedom in hierarchical and other richly parameterized models. *Biometrika*, 88, 367–379.
- [55] Isaki, C.T., Tsay, J.H., and Fuller, W.A., (2000), Estimation of census adjustment factors, *Survey Methodology*, 26, 31–42.
- [56] Jiang, J. and Lahiri, P., (2006)., *Mixed Model Prediction and Small Area Estimation*. *Test*, 15, 1–96.
- [57] Jackson, C., Best, N., and Richardson, S., (2006)., Improving ecological inference using individual-level data. *Statistics in medicine*, 25(12), 2136–2159.
- [58] Jiang, J., and P., Lahiri (2006)., *Mixed Model Prediction and Small Area Estimation*, In: *Test* 15 (1), pp. 1–96. *Annals of Statistics* , 30, 1782–1810.
- [59] Kleffe, J., and Rao, J.N.K., (1992), Estimation of mean square error of empirical best linear unbiased predictors under a random error variance linear model. *Journal of Multivariate Analysis*, 43(1), 1–15.
- [60] Kleinschmidt I., Bagayako M., Clarke GPY, et al(2000), A Spatial Statistical Approach to Malaria Mapping, *International Journal of Epidemiology*, 29, 355–361.
- [61] Levy, P.S., (1971), The Use of Mortality data in evaluating synthetic estimates, Proceedings of the social statistics section, American Statistical Association, pp. 328–331.
- [62] Louis, T.A., (1984), Estimating a population of parameter values using bayes and empirical bayes methods, *Journal of the American Statistical Association*, 79(386), 393–398.

- [63] Lahiri, P. (1990), Adjusted bayes and empirical bayes estimation in finite population sampling, *Sankhyā: The Indian Journal of Statistics, Series B*, 50–66.
- [64] Lohr, S.L., (1999), *Sampling: Design and Analysis*, Pacific Grove; CA:Duxbury.
- [65] LeSage, J.P., and Pace, R.K., (2004), Models for spatially dependent missing data, *The Journal of Real Estate Finance and Economics*, 29(2), 233–254.
- [66] Lin, R., Louis, T.A., Paddock, S.M., and Ridgeway, G., (2009), Ranking USRDS provider specific SMRs from 1998–2001, *Health Services and Outcomes Research Methodology*, 9(1), 22–38.
- [67] Metropolis, N., Rosenbluth, A.W., Rosenbluth, M.N., Teller, A.H., and Teller, E., (1953), Equations of state calculations by fast computing machines, *Journal of Chemical Physics*, 21, 1087–1091.
- [68] McCulloch, C.E., (1997), Maximum likelihood algorithms for generalized linear mixed models, *Journal of the American Statistical Association*, 92, 162–170.
- [69] Mukhopadhyay, P. (1998), *Small area estimation in survey sampling*, New Delhi: Narosa Publishing House.
- [70] Marker, D.A., (1999), Organization of small area estimators using a generalized linear regression framework, *Journal of Official Statistics*, 15, 1–24.
- [71] McCulloch CE, Searle SR (2001). *Generalized, Linear, and Mixed Models*. John Wiley and Sons, Inc., New York.
- [72] Moura, F.A.S., and Meung, D. (2002), Small area estimation using multilevel models, *Survey Methodology*, 25, 73–80.
- [73] Milana Karaganis, (2009), *A hierarchical bayes approach*, department of statistics, The University of Manitoba.
- [74] Normand, S.L.T., Glickman, M.E., and Gatsonis, C.A., (1997), Statistical methods for profiling providers of medical care: issues and applications, *Journal of the American Statistical Association*, 92(439), 803–814.
- [75] Purcell, N.J., and Kish, L., (1980), Postcensal estimates for local areas (or Domains), *International Statistical Review*, 48, 3–18.

- [76] Platek, R. and Singh, M.P., (Eds.,) (1986), Small area statistics: Contributed papers, Laboratory for research in statistics and probability, Carleton university, Ottawa, Canada.
- [77] Platek, R., Rao, J.N.K., Sarndal, C.E., and Singh, M.P., (Eds.,) (1987), Small Area Statistics, New York: Wiley.
- [78] Pfeffermann, D., and Burck, L., (1990), Robust small area estimation combining time series and cross-sectional data, *Survey Methodology*, 16, 217–237.
- [79] Pfeffermann, D. (2002), Small area estimation-new developments and directions, In: *International Statistical Review*, 70 (1), pp. 125–143.
- [80] Petrucci, A., and Salvati, N. (2006), Small area estimation for spatial correlation in watershed erosion assessment, *Journal of Agricultural, Biological, and Environmental Statistics*, 11, 2–169.
- [81] Paddock, S.M., Ridgeway, G., Lin, R., and Louis, T.A., (2006), Flexible distributions for triple-goal estimates in two-stage hierarchical models, *Computational Statistics and Data Analysis*, 50(11), 3243–3262.
- [82] Pratesi, M. Salvati, N. (2008), Small area estimation: The EBLUP estimator based on spatially correlated random area effects, *Statistical Methods and Applications*, 17, 113–141.
- [83] Pfeffermann, D. (2013), New important developments in small area estimation, *Statistical Science*, 28, 40–68.
- [84] Royal, R.M., (1970), On finite population sampling theory under certain linear regression, *Biometrika*, 57, 377–387.
- [85] Rao, J.N.K., (1979), On deriving mean square errors and their non-negative unbiased estimators in finite population sampling, *Journal of the Indian Statistical Association*, 17, 125–136.
- [86] Rao, J.N.K., (1986), Synthetic estimators, SPREE and best model based predictors, *Proceedings of the conference on survey research methods in agriculture, US*, Department of Agriculture, Washington, DC, pp. 1–16.
- [87] Robinson, G.K., (1991), That BLUP is a good thing: The estimation of random effects, *Statistical science*, 6(1), 15–32.
- [88] Rao, J.N.K., and Yu, M. (1992), Small area estimation by combining time series and cross-sectional data, *Proceedings of the section on survey research method*, American Statistical Association, pp. 1–9.

- [89] Rao, J.N.K., (1992), Estimating totals and distribution functions using auxiliary information at the estimation stage, Proceedings of the workshop on uses of auxiliary information in surveys, Statistics Sweden.
- [90] Rao, J.N.K., and Yu, M. (1994), Small area estimation by combining time series and cross-sectional data, *Canadian Journal of Statistics*, 22, 511–528.
- [91] Rao, J.N.K., and Choudhry, G.H., (1995), Small area estimation: Overview and empirical study, in B.G., Cox, D.A., Binder, B.N., Chinnappa, A. Christianson, M.J., Colledge, and P.S., Kott (Eds.), *Business survey methods*, New York: Wiley, pp. 527–542.
- [92] Rao, J.N.K., (2001), EB and EBLUP in small area estimation, in S.E., Ahmed and N. Reid (Eds.), *Empirical bayes and likelihood inference*, Lecture notes in statistics 148, New York: Springer, pp. 33–43.
- [93] Rao, J.N.K., (2003), *Small Area Estimation*, Wiley, Hoboken, New Jersey.
- [94] Robert, C.P., and Casella, G., (2005), *Monte carlo statistical methods*, Second Ed., Springer-Verlag, New York.
- [95] S. Richardson, G. Li., (2010), *Bayesian statistics for small area estimation*, Department of epidemiology and public health imperial college london st, Mary's campus, Norfolk place w2 1PG London-United Kingdom.
- [96] Rao, J.N.K., and Molina, I., (2015), *Small Area Estimation*, 2nd Edition, Wiley, New York.
- [97] Smith, S.K., and Lewis, B.B., (1980), Some New Techniques for Applying the Housing Unit Method of Local Population Estimations, *Demography*, 17, 349–340.
- [98] Smith, T.M.F., (1983), On the validity of inferences from non-random samples, *Journal of the Royal Statistical Society, Series A*, 146, 394–403.
- [99] Särndal, C.E., B. Swensson and J.H., Wretman (1992), *Model Assisted Survey Sampling*, New York: Springer.
- [100] Sardnal, C.E., Swensson, B., and Wretman, J.H., (1992), *Model assisted survey sampling*, New York Springer-Verlag.
- [101] Singh, A.C., Mantel, J.H., and Thomas, B.W., (1994), Time series EBLUPs for small areas using survey data, *Survey Methodology*, 20, 33–43.

- [102] Skinner, C.J., (1994), Sample models and weights, Proceedings of the section on survey research methods, American statistical association, Washington, DC, pp. 133–142.
- [103] Spiegelhalter, D.J., N.G., Best, B.P., Carlin, and A. van Der Linde (2002), Bayesian measures of model complexity and
t with discussion, J.R., Statist. Soc. B., 64, 583–639.
- [104] Salvati, N. (2004), Small area estimation by spatial models: the spatial empirical best linear unbiased prediction spatial EBLUP, Working Paper 2004/03, University of Florence.
- [105] Singh, B. B., G.K., Shukla, and D. Kundu (2005), Spatio-Temporal models in small area estimation, Survey Methodology 31 (2), 183–196.
- [106] Tanner, M.A., and Wong, W.H., (1987), The calculation of posterior distributions by data augmentation with discussion, *Journal of the American Statistical Association*, 82, 528–550.
- [107] Thompson, M.E., (1997), Theory of sample survey, London: Chapman and Hall.
- [108] You, Y. (1999), Hierarchical bayes and related methods for model-based small area estimation, Unpublished Ph.D. Thesis, Carleton University, Ottawa, Canada.
- [109] Yar M, Hennell S, Clarke P, Meltzer H, Gatward R., (2002), Childhood mental disorder in england: Ward estimates, Technical report, Office for National Statistics, United Kingdom.
- [110] Valliant, R., Dorfmap, A.H., and Royall, R.M., (2001), Finite Population Sampling and Inference: A Prediction Approach, New York: Wiley.
- [111] Zhang W. Li, N., (2011), Prevalence, risk factors, and management of prehypertension, *International Journal of Hypertension*, 60, 53–59.

پیوست آ

توزیع‌های شرطی کامل برای مشخص‌سازی CAR با مشاهدات گم‌شده

فرض کنید یک مدل سطح ناحیه دارید و داده از نواحی l ام اول وجود دارد بنابراین، مقادیر برآوردهای مستقیم و واریانس‌های نمونه‌گیری برای این نواحی نیز وجود داشته باشد. در این زمینه، $s = 1, \dots, l$ و $s = l + 1, \dots, m$ است. توزیع‌های شرطی کامل برای این مدل با استفاده از نمونه‌گیری گیبز به صورت زیر است:

$$\Pi(\alpha, \beta | \dots) = N((\bar{X}^T V^{-1} \bar{X})^{-1} \bar{X}^T V^{-1} (\hat{Y} - v), (\bar{X}^T V^{-1} \bar{X})^{-1})$$

پس،

$$\hat{Y} = (\hat{Y}_1, \dots, \hat{Y}_l)^T, v = (v_1, \dots, v_l)^T, V = \text{diag}(\hat{V}_1^2, \dots, \hat{V}_l^2)$$

9

$$\bar{X}^T = \begin{pmatrix} 1 & \dots & 1 \\ \bar{X}_1 & \dots & \bar{X}_l \end{pmatrix}$$

$$\Pi(v_i | \dots) \propto \exp\left\{ \frac{-1}{\hat{\sigma}_i^2} (\bar{Y}_i - \alpha - \beta \bar{X}_i - v_i)^2 + \frac{1}{\sigma_v^2} n_i (v_i - \bar{v}_i^2) \right\}, \quad i \in s$$

$$\Pi(v_j|\dots) \propto \exp\left\{-\frac{1}{\sigma_v^2} n_j (v_j - \bar{v}_j)^2\right\}, \quad j \in \underline{s}$$

لذا به ترتیب n_i و n_j نشان‌دهنده‌ی تعدادی از همسایگی‌های نواحی i و j می‌باشند.

$$\Pi(\sigma_v^2|\dots) = Ga^{-1}\left(\epsilon + \frac{n}{\sigma_v^2}, \epsilon + \frac{1}{\sigma_v^2} \sum_{k \sim l} (v_k - v_l)^2\right)$$

از این رو، شرطی کامل برای هر v_i مستقیماً وابسته بر مشاهدات داده و همسایگی‌های آن است، ضمن اینکه v_j تنها وابسته بر روی همسایگی‌های \bar{v}_j و غیرمستقیم بر روی مشاهدات داده می‌باشد. علاوه بر این، همان‌گونه که انتظار می‌رفت α و β تنها اطلاعاتی هستند که بوسیله داده از نواحی تحت بررسی به دست آورده می‌شوند.

پیوست ب

دستورات نرم افزار R و OPENBUGS

در این قسمت دستورات مربوط به کدنویسی داده‌های بیمه استان گیلان ارائه شده است.

```
model
{
for(i in 1:N)
{
y[i] ~ dpois(mu[i])
log(mu[i]) <- alpha + beta1*x1[i] + beta2*x2[i] +beta3*x3[i] + u[i] + v[i]
u[i] ~ dnorm(0, precu)
}
v[1:N] ~ car.normal(adj[], weights[], num[], precv)
for (k in 1:sumNumNeigh){ weights[k] <- 1}
alpha ~ dflat()
beta1 ~ dnorm(0,1.0E-5)
beta2 ~ dnorm(0,1.0E-5)
beta3 ~ dnorm(0,1.0E-5)
```

```

precu ~ dgamma(0.001, 0.001)
precv ~ dgamma(0.1, 0.1)
sigmau<-1/precu
sigmav<-1/precv
}

#liatticels
#Data

list(N=16,
num=c(1 ,2, 3, 2, 6, 4, 3, 4 ,3, 3, 6, 3,
4, 4, 7, 3),
adj=c(4 ,
15, 8,
13, 11, 15,
13, 1,
15, 10, 7, 12, 14, 8,
10, 9, 11, 15,
5, 12, 14,
2, 15, 5, 14,
16, 11, 6,
6, 15, 5,
15, 3, 6, 9, 16, 13,
7, 5, 14,
4, 16, 11, 3,
8, 5, 12, 7,
3, 11, 6, 10, 5, 8, 2,
13, 11, 9
),
y=c(5547,4822, 17454,
6317,
3538,
4112,
5195 ,

```

```
14081,  
5,573,  
12170,  
7422,  
3138,  
5410,  
7977,  
115632,  
2283 ),
```

```
x1=c(28742, 38824, 48193, 61055, 16351, 18416, 51586, 58378,
```

```
31209, 13220, 41975, 15306, 22246, 49351, 179456, 16901),
```

```
x2=c(91257, 108130, 139016, 200649, 46975, 54226, 147399, 167544,
```

```
92310, 42408, 125074, 43225, 69865, 140686, 523749, 52649),
```

```
x3=c(3286, 2417,9657,4272, 1094,1647, 3,069,6538,2606, 8985, 3591,
```

```
1184, 4720, 4578, 55630, 1129),
```

```
sumNumNeigh = 58)
```

```
list(alpha=0, beta1=0.5, beta2=0.5, beta3=0.5, precv=0.01, precu=0.01),
```

```
precu ~ dgamma(0.001, 0.001)
```

```
precv ~ dgamma(0.1, 0.1)
```

```
sigmau<-1/precu
```

```
rm(list=ls())
```

```
library(R2OpenBUGS)
```

```
l.model <- function(){
```

```
## Likelihood
```

```
for(i in 1:N){
```

```

y[i] ~ dpois(mu[i])
log(mu[i]) <- beta[1] + beta[2]*x1[i] + beta[3]*x2[i] +
beta[4]*x3[i] + v[i] + u[i]
u[i] ~ dnorm(0, precu)
# Area-specific relative risk (for maps)
# RR[i] <- exp(beta0 + beta1 * x1[i] +
# beta2 * x2[i] + beta3 * x3[i] + v[i])
}
# CAR prior distribution for random effects:
v[1:N] ~ car.normal(adj[], weights[], num[], precv)
for (k in 1:sumNumNeigh){
weights[k] <- 1
}
precu ~ dgamma(0.01, 0.01)
precv ~ dgamma(0.01, 0.01)
sigmau <- sqrt(1/precu)
sigmav <- sqrt(1/precv)
#
for (j in 1:4){
beta[j] ~ dnorm(0,.01)
}
}

# Loading data

l.data <- list(N = 16,
num=c(1 ,2, 3, 2, 6, 4, 3, 4 ,3, 3, 6, 3, 4, 4, 7, 3),
adj=c(4 ,
15, 8,
13, 11, 15,
13, 1,
15, 10, 7, 12, 14, 8,
10, 9, 11, 15,
5, 12, 14,

```

2, 15, 5, 14,
16, 11, 6,
6, 15, 5,
15, 3, 6, 9, 16, 13,
7, 5, 14,
4, 16, 11, 3,
8, 5, 12, 7,
3, 11, 6, 10, 5, 8, 2,
13, 11, 9
)

y=c(5547,4822, 17454
, 6317,
3538,
4112,
5195 ,
14081,
5,573 ,
12170 ,
7422 ,
3138,
5410,
7977,
115632,
2283),

x1=c(-0.363253508829662, -0.109955722412744, 0.125428820429955,
0.448570670151676, -0.674562065811612, -0.622681493608584, 0.210673750931106,
0.381314354516079, -0.301273183587593, -0.753224570202742, -0.0307907427460408,
-0.70081639896036, -0.526457478336138, 0.154522138885697, 3.42324942531429,
-0.660743995733324),

x2=c(-0.314534366653528, -0.16939426940249, 0.0962844506648097,
0.626446224734223, -0.695444328513923, -0.633071849099043, 0.168394296883944,

```
0.341679857194627, -0.305476551066179, -0.734729270382322, -0.0236434362476521,
-0.727701506531549, -0.498546513894339, 0.110649647275325, 3.40572468273282,
-0.646637067694722),
```

```
x3=c(-0.252182437726062, -0.319255229612263, 0.239556062972422,
-0.176079131857485, -0.421369503105293, -0.378686817359529, -0.505577116393952,
-0.500482980301329, -0.00118045934410152, -0.304667476256115,
0.187688495483899, -0.228641354267729, -0.41442295388808, -0.141500753531803,
-0.152460864518961, 3.78793058700498, -0.418668067298599),
sumNumNeigh = 58)
```

```
l.inits <- function() {
list(beta = c(0, 0, 0, 0), precv=0.1, precu=0.1,
v= c(0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0),
u= c(0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0))
}
```

```
### Fitting model
```

```
l.out <- bugs(data = l.data, inits = l.inits,
parameters.to.save = c("beta", "sigmav", "sigmau"),
model.file = l.model,
n.chains = 1, n.thin=10, n.iter = 10000, debug=FALSE)
```

```
l.out
```

```
l.out$summary
```

```
par(mfrow=c(2,3))
```

```
#Intercept
```

```
plot(density(l.out$sims.list$beta[,1]), main=expression(alpha), xlab="")
```

```
#abline(v = seq(-0.5, 1), col = "grey")
```

```
abline(v = mean(l.out$sims.list$beta[,1]), col = "red", lwd = 1)
```

```
abline(v = median(l.out$sims.list$beta[,1]), col = "red", lwd = 1, lty = 2)
```

```
plot(density(l.out$sims.list$beta[,2]), main=expression(beta[1]), xlab="")
abline(v = mean(l.out$sims.list$beta[,2]), col = "red", lwd = 1)
abline(v = median(l.out$sims.list$beta[,2]), col = "red", lwd = 1, lty = 2)

plot(density(l.out$sims.list$beta[,3]), main=expression(beta[2]), xlab="")
abline(v = mean(l.out$sims.list$beta[,3]), col = "red", lwd = 1)
abline(v = median(l.out$sims.list$beta[,3]), col = "red", lwd = 1, lty = 2)

plot(density(l.out$sims.list$beta[,4]), main=expression(beta[3]), xlab="")
abline(v = mean(l.out$sims.list$beta[,4]), col = "red", lwd = 1)
abline(v = median(l.out$sims.list$beta[,4]), col = "red", lwd = 1, lty = 2)

plot(density(l.out$sims.list$sigmav), main=expression(sigma[v]), xlab="")
abline(v = mean(l.out$sims.list$sigmav), col = "red", lwd = 1)
abline(v = median(l.out$sims.list$sigmav), col = "red", lwd = 1, lty = 2)

plot(density(l.out$sims.list$sigmau), main=expression(sigma[u]), xlab="")
abline(v = mean(l.out$sims.list$sigmau), col = "red", lwd = 1)
abline(v = median(l.out$sims.list$sigmau), col = "red", lwd = 1, lty = 2)

par(mfrow=c(3,2))
acf(l.out$sims.list$beta[,1], main=expression(alpha), xlab="")
acf(l.out$sims.list$beta[,2], main=expression(beta[1]), xlab="")
acf(l.out$sims.list$beta[,3], main=expression(beta[2]), xlab="")
acf(l.out$sims.list$beta[,4], main=expression(beta[3]), xlab="")
acf(l.out$sims.list$sigmav, main=expression(sigma[v]), xlab="")
acf(l.out$sims.list$sigmau, main=expression(sigma[u]), xlab="")

par(mfrow=c(3,2))
ts.plot(l.out$sims.list$beta[,1], main=expression(alpha), xlab="")
ts.plot(l.out$sims.list$beta[,2], main=expression(beta[1]), xlab="")
ts.plot(l.out$sims.list$beta[,3], main=expression(beta[2]), xlab="")
ts.plot(l.out$sims.list$beta[,4], main=expression(beta[3]), xlab="")
```

```
ts.plot(1.out$sims.list$sigmav, main=expression(sigma[v]), xlab="")  
ts.plot(1.out$sims.list$sigmau, main=expression(sigma[u]), xlab="")
```

واژه‌نامه فارسی به انگلیسی

Conditional autoregressive	اتورگرسیو شرطی
Small area estimation	برآوردگر کوچک ناحیه‌ای
Direct estimates	برآوردگر مستقیم
Indirect estimators	برآوردگر غیرمستقیم
Point data	داده الگو نقطه‌ای
Logistic regression	رگرسیون لجستیک
Neighborhood matrix	ماتریس همسایگی
Unit level model	مدل سطح واحد
Mean squared error	میانگین مربع خطا
Random field	میدان تصادفی
Horwitz-thompson	هوروتیز-تامپسون

واژه‌نامه انگلیسی به فارسی

Conditional autoregressive	اتورگرسیو شرطی
Direct estimates	برآوردگر مستقیم
Horwitz-thompson	هوروتیز-تامپسون
Indirect estimators	برآوردگر غیرمستقیم
Logistic regression	رگرسیون لجستیک
Mean squared error	میانگین مربع خطا
Neighborhood matrix	ماتریس همسایگی
Point data	داده الگو نقطه‌ای
Random field	میدان تصادفی
Small area estimation	برآوردگر کوچک ناحیه‌ای
Unit level model	مدل سطح واحد

Abstract

In recent years, the problem of small-scale regionalization has to be considered because of the need reliable statistics have been very much considered. The major problem here is the impossibility of measuring the target variable for each individual in the area in question. Even sampling from all areas under study could lead to high financial and time costs. Several statistical models have been proposed to achieve this goal. The main purpose of these models is the use of auxiliary information to upgrade direct estimates. To this end, the importance of using spatial information of areas within the framework of spatial models plays an effective role in analyzing small areas. The similarity of space between adjacent areas is also useful information that small space-based models have been proposed to use for this information. In view of the fact that in many analyzes of small regions with nonnormal responses such as numeric responses, we classify the generalized linear space mixing models in the framework of Bayesian inference for the analysis of small area data. By studying the simulation and analysis of the insurance data package in Gilan province, we show the performance of these models.

Keywords: Small area estimation, Spatial models, Generalized linear Mixed models, Bayesian inference.



Shahrood University of Technology

Faculty of Mathematical Sciences

MSc Thesis in: Mathematical Statistics

**Bayesian Spatial and non-Spatial Random
Effects Models for Small Area Estimation**

By: Motahareh Yousefi

Supervisors

Dr. Mohammad Reza Rabiei

Dr. Hossein Baghishani

September, 2018