

حاشا
البربر
البربر



دانشکده علوم ریاضی

رشته آمار، گرایش آمار ریاضی

پایان نامه کارشناسی ارشد

ناهمگنی فضایی در مدل‌های تعیین قیمت محصول با استفاده از یک رهیافت رگرسیون موزون موضعی تکراری

نگارنده: رویا پیرواولیاء

استاد راهنما

دکتر حسین باغیشنی

بهمن ۱۳۹۷

تقدیم بہ دستان مادر و خاطرات پدرم

سپاس گزارمی...

سپاس خدای بزرگ را که مرایاری رساند تا بتوانم این مقطع تحصیلی را به پایان رسانده و گامی در راستای اعتلای علم بردارم. از استاد راهنمای گران قدرم جناب آقای دکتر حسین باغشینی که وجودشان همیشه قوتی برای انجام کارهایم بوده است و بدون شک انجام این پایان نامه بدون کمک و راهنمایی های ارزنده ایشان امکان پذیر نبود، کمال تشکر را دارم. از اساتید گرامی جناب آقای دکتر محمد رضاری و سرکار خانم دکتر نگار اقبال که زحمات داورسی این پایان نامه را داشتند نیز سپاس گزارم. تشکر از اساتید گروه آمار دانشگاه صنعتی شاهرود، آقایان دکتر داوود شاهسونی، دکتر احمد نزلتی و دکتر محمد آرشی، که توفیق دانشجویی در محضرشان را داشتم. در پایان از تمامی عزیزانی که در طول انجام این پروژه مرایاری کرده اند کمال تشکر و قدردانی را ابراز می نمایم.

رویا پیرو اولیاء
بهمن ۱۳۹۷

تعهد نامه

این جانب رویا پیرواولیاء دانشجوی کارشناسی ارشد رشته آمار دانشگاه صنعتی شاهرود، نویسنده پایان نامه با عنوان **ناهمگنی فضایی در مدل های تعیین قیمت محصول با استفاده از یک رهیافت رگرسیون موزون موضعی تکراری**، تحت راهنمایی دکتر حسین باغیشنی 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 خواهد 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 آن ها) استفاده شده است، ضوابط و اصول اخلاقی رعایت شده اند.
- 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.

رویا پیرواولیاء

بهمن 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 است. 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

بسیاری از مطالعات تجربی، مستلزم استفاده از متغیرهایی است که تحت تأثیر موقعیت مکانی مشاهدات هستند. در این موارد میزان و نحوه تأثیر ساختار فضایی داده‌ها دارای اهمیت است و نادیده گرفتن آن موجب خطا در برآورد و از دست رفتن اطلاعات مهمی می‌شود. در برخورد با چنین داده‌هایی مدل محبوب و پرکاربرد رگرسیون خطی پاسخگو نیست؛ زیرا عامل موقعیت فضایی موجب نقض پذیره‌های همگنی و ناهمبستگی این مدل می‌شود. بنابراین در چنین مواردی باید از مدل‌های متناسب با ساختار فضایی استفاده شود. برای داشتن یک مدل فضایی منعطف، لازم است وابستگی و ناهمگنی فضایی با هم در نظر گرفته شوند. این در حالی است که ساختار وابستگی و ناهمگنی موجود در داده‌ها نامعلوم است و اغلب هیچ‌گونه اطلاع پیشینی هم در مورد این ساختار در دسترس نیست. در این پایان‌نامه از یک رهیافت رگرسیون موزون تکراری برای تعیین ناحیه‌های همگن تحت عنوان رژیم‌های فضایی استفاده شده است. سپس این رژیم‌ها در دسته‌ای از مدل‌های اتورگرسیون فضایی معروف به مدل‌های اقتصادسنجی فضایی، اعمال و نتیجه با عنوان مدل‌های اقتصادسنجی با رژیم فضایی درونی به داده‌ها برازش داده شده‌اند. این مدل‌ها بر اساس معیار اطلاع آکاییک با مدل رگرسیون خطی و مدل‌های اقتصادسنجی فضایی مقایسه شده‌اند؛ نتایج بیانگر عملکرد بهتر مدل‌های اقتصادسنجی با رژیم فضایی درونی هستند.

کلمات کلیدی: ناهمگنی فضایی، وابستگی فضایی، رژیم فضایی، مدل‌های اقتصادسنجی، مدل‌های اقتصادسنجی با رژیم فضایی درونی.

پیش‌گفتار

تحلیل داده‌های فضایی، بر اساس یک نمونه مشاهده‌شده در ناحیه جغرافیایی تحت مطالعه در عمل می‌تواند دشوار باشد؛ چرا که شکل وابستگی و ساختار همگنی موجود در داده‌ها نامعلوم هستند. یکی از روش‌های مدل‌بندی داده‌های فضایی، استفاده از مدل‌های اتورگرسیو فضایی است که ساختار وابستگی داده‌ها را از طریق ضریب اتورگرسیو فضایی در مدل رگرسیونی لحاظ می‌کنند. این مدل‌ها در مطالعات اقتصادی نظیر علوم منطقه‌ای، اقتصاد شهری و اقتصاد کشاورزی کاربرد دارند و با عنوان مدل‌های اقتصادسنجی فضایی نیز شناخته می‌شوند.

اولین بار ویتل (۱۹۵۴) با فرض این که خطاهای هر موقعیت فقط به موقعیت‌های مجاور وابسته است، از این روش برای مدل‌بندی وابستگی فضایی استفاده کرد. بسج (۱۹۷۴)، اورد (۱۹۷۵)، انسلین (۱۹۸۸a)، هاینینگ (۱۹۹۰)، کرسی (۱۹۹۲)، لی‌سیج و پیس (۲۰۰۹) و اسمیت (۲۰۱۶) نیز به تحلیل و مدل‌بندی داده‌های فضایی با استفاده از مدل‌های اتورگرسیو فضایی پرداخته‌اند. این مدل‌ها بدون توجه به ناهمگنی فضایی برای کل فضای مورد مطالعه، یک ساختار وابستگی واحد ارائه می‌دهند.

یکی دیگر از روش‌های مدل‌بندی ساختار فضایی، که ناهمگنی و وابستگی فضایی را با ارائه یک مدل موضعی برای هر موقعیت فضایی توأم در نظر می‌گیرد، مدل رگرسیون موزون جغرافیایی (فضایی) است. از این مدل در برخی متون (غیرفضایی) با عنوان رگرسیون موزون موضعی نیز یاد می‌شود. مدل رگرسیون موزون جغرافیایی به‌طور کامل توسط فودرینگ‌هام و همکاران (۲۰۰۳) تشریح شده است. آندرانو و همکاران (۲۰۱۶) و بیل و همکاران (۲۰۱۷) دو روش مشابه برای مدل‌بندی ساختار فضایی داده‌ها ارائه کردند. در این روش‌ها وابستگی و ناهمگنی فضایی در قالب رژیم‌های فضایی مختلف مدل‌بندی می‌شوند.

با این مقدمه، با هدف معرفی یک مدل ایده‌آل برای مدل‌بندی ساختار داده‌های فضایی، پایان‌نامه خود را در چهار فصل و یک پیوست ارائه می‌کنم.

۱. فصل اول به مروری بر رگرسیون خطی، پذیره‌ها و راه‌کارهای موجود برای برخورد با عدم برقراری پذیره‌های ناهمبستگی و همگنی در آن، اختصاص دارد. همچنین به دلیل اهمیت اثر عامل موقعیت فضایی در نقض این پذیره‌ها به مقدمه‌ای از آمار فضایی و همچنین ارائه آزمون‌هایی برای بررسی حضور ناهمگنی و وابستگی فضایی، پرداخته شده است.

۲. در فصل دوم، چند روش معمول مدل‌بندی ساختار وابستگی و ناهمگنی فضایی، روش‌های برآورد پارامترها و معایب هر کدام از آن‌ها بیان شده‌اند.

۳. در فصل سوم، مدل پیشنهادی بیل و همکاران (۲۰۱۷) برای مدل‌بندی هم‌زمان ساختار وابستگی و ناهمگنی فضایی ارائه و تشریح شده است.

۴. در فصل چهارم، کارایی بهتر مدل پیشنهادی فصل سوم نسبت به مدل‌های معرفی شده در فصل دوم در قالب دو مثال از داده‌های واقعی تعیین قیمت مسکن نمایش داده شده است.

۵. در پیوست نیز بسته‌افزارها و کدهای نرم‌افزار R برای تحلیل داده‌ها و برازش مدل‌های پیشنهادی، به صورت مرحله به مرحله، ارائه شده‌اند.

فهرست مطالب

ش	فهرست تصاویر
ث	فهرست جداول
۱	۱ مفاهیم و مقدمات
۱	۱.۱ مقدمه
۲	۲.۱ مدل رگرسیون خطی
۳	۱.۲.۱ روش کمترین توان‌های دوم معمولی
۳	۲.۲.۱ راه‌های برخورد با ناهمگنی در مدل رگرسیون خطی کلاسیک
۶	۳.۱ مقدمه‌ای بر آمار فضایی
۷	۱.۳.۱ داده‌های فضایی
۸	۲.۳.۱ مدل آماری
۹	۳.۳.۱ ساختار همبستگی فضایی
۱۲	۴.۳.۱ ناهمگنی فضایی
۱۴	۴.۱ سایر مفاهیم و مقدمات
۱۴	۱.۴.۱ مدل خودهمبستگی
۱۵	۲.۴.۱ ماتریس وزن فضایی
۱۶	۳.۴.۱ آماره آزمون موران نوع I
۱۷	۴.۴.۱ آزمون‌های مبتنی بر درست‌نمایی
۲۰	۵.۴.۱ آزمون ناهمگنی در حضور خودهمبستگی
۲۱	۶.۴.۱ ضرب هادامارد
۲۳	۲ مدل‌بندی ساختار وابستگی فضایی
۲۳	۱.۲ مقدمه
۲۴	۲.۲ وارد کردن ساختار وابستگی از طریق ماتریس کواریانس
۲۶	۳.۲ مدل‌های اقتصادسنجی

۲۶	مدل خطای فضایی	۱.۳.۲
۲۹	مدل تاخیر فضایی	۲.۳.۲
۳۱	مدل ترکیبی فضایی	۳.۳.۲
۳۲	مدل دوربین فضایی	۴.۳.۲
۳۴	رگرسیون موزون جغرافیایی	۴.۲
۳۵	استنباط مدل	۱.۴.۲
۳۷	روش‌های انتساب وزن	۲.۴.۲
۴۱	انتخاب بهترین همسایگی	۳.۴.۲
۴۵	مدل رگرسیون موزون موضعی تکراری	۳
۴۵	مقدمه	۱.۳
۴۶	مرحله اول: کنترل ناهمگنی فضایی	۲.۳
۴۷	اختصاص وزن‌های اولیه و برازش مدل GWR	۱.۲.۳
۴۸	شرط توقف الگوریتم تکرار	۲.۲.۳
۴۸	آماره آزمون همگنی ضرایب	۳.۲.۳
۴۹	آزمون همگنی ضرایب	۴.۲.۳
۵۰	الگوریتم به‌روزرسانی وزن‌ها	۵.۲.۳
۵۲	مرحله دوم: مدل‌های اقتصادسنجی با رژیم فضایی درونی	۳.۳
۵۳	مدل ترکیبی فضایی با رژیم درونی (ESR-SAC)	۱.۳.۳
۵۴	مدل دوربین فضایی با رژیم درونی (ESR-SDM)	۲.۳.۳
۵۵	ارزیابی مدل پیشنهادی	۴
۵۵	مقدمه	۱.۴
۵۵	مثال کاربردی اول	۲.۴
۵۷	انجام آزمون‌های آماری و برازش مدل‌های فراموضعی فضایی	۱.۲.۴
۵۹	ناهمگنی در داده‌ها	۲.۲.۴
۶۰	برازش مدل پیشنهادی	۳.۲.۴
۶۴	مثال کاربردی دوم	۳.۴
۶۶	انجام آزمون‌های آماری و برازش مدل‌های فراموضعی	۱.۳.۴
۶۷	ناهمگنی در داده‌ها	۲.۳.۴
۶۸	برازش مدل پیشنهادی	۳.۳.۴
۷۲	نتیجه و آینده تحقیق	۴.۴
۷۵	برنامه‌های رایانه‌ای با نرم‌افزار R	آ
۷۵	برنامه‌های داده بالتیمور	۱.آ

۸۵	۲.آ برنامه‌های داده مسکن
۹۷		مراجع
۱۰۱		واژه‌نامه فارسی به انگلیسی
۱۰۵		واژه‌نامه انگلیسی به فارسی

فهرست تصاویر

۳۸	رگرسیون موزون فضایی با هسته فضایی ثابت	۱.۲
۳۹	رگرسیون موزون فضایی با هسته فضایی سازوار	۲.۲
۴۱	مبادله واریانس-اریبی	۳.۲
۵۶	موقعیت شهر بالتیمور	۱.۴
		رژیم‌های فضایی مجموعه داده بالتیمور، نقاط * رژیم فضایی اول و نقاط	۲.۴
۶۱	○ رژیم فضایی دوم را نشان می‌دهند	
۶۵	موقعیت شهرستان لوکاس	۳.۴
۶۵	داده‌های انتخاب شده از مجموعه داده مسکن	۴.۴
		رژیم‌های فضایی مجموعه داده مسکن لوکاس، نقاط Δ رژیم فضایی اول، نقاط • رژیم فضایی دوم، نقاط × رژیم فضایی سوم و نقاط + رژیم فضایی	۵.۴
۶۹	چهارم را نشان می‌دهند	
۷۰	تابع تغییرات وزن‌های روش IGWR	۶.۴

فهرست جداول

۵۶	متغیرهای مجموعه داده بالتیمور	۱.۴
۵۷	نتایج آزمون‌های خودهمبستگی فضایی برای داده‌های بالتیمور	۲.۴
۵۸	برآورد پارامترهای مدل‌های فراموضعی برای مجموعه داده بالتیمور	۳.۴
۶۰	نتایج آزمون‌های ناهمگنی فضایی داده‌های بالتیمور	۴.۴
۶۰	نتایج برازش مدل GWR برای داده‌های بالتیمور	۵.۴
۶۲	برآورد پارامترهای مدل‌های با رژیم فضایی درونی برای مجموعه داده بالتیمور	۶.۴
۶۶	متغیرهای داده لوکاس	۷.۴
۶۶	نتایج آزمون‌های خودهمبستگی فضایی برای داده‌های لوکاس	۸.۴
۶۷	برآورد پارامترهای مدل‌های فراموضعی برای مجموعه داده مسکن لوکاس	۹.۴
۶۸	نتایج آزمون‌های ناهمگنی فضایی داده‌های مسکن لوکاس	۱۰.۴
۶۸	نتایج برازش مدل GWR برای داده‌های مسکن لوکاس	۱۱.۴
	برآورد پارامترهای مدل‌های با رژیم فضایی درونی برای مجموعه داده مسکن شهر	۱۲.۴
۷۰	لوکاس	

فصل ۱

مفاهیم و مقدمات

۱.۱ مقدمه

تجربه نشان داده است که تحلیل مسائل اقتصادی اغلب با چالش‌هایی مواجه است. چرا که داده‌ها به شکل‌های ناشناخته‌ای ناهمگن^۱ و به یکدیگر وابسته‌اند. در موارد متعددی این ناهمگنی و وابستگی داده‌ها حاصل از فضا و موقعیت قرارگیری آن‌ها نسبت به یکدیگر است. به این نوع داده‌ها، داده‌های فضایی^۲ می‌گویند (انسلین، ۱۹۸۸a). دیدگاه‌های مختلفی برای مدل‌بندی داده‌های فضایی پیشنهاد شده‌اند. یکی از دیدگاه‌ها مدل‌های اقتصادسنجی فضایی^۳ است که به‌طور گسترده برای مدل‌بندی ساختار وابستگی داده‌ها معرفی و استفاده می‌شوند. اما برخورد مستقیم با مشکل ناهمگنی^۴ داده‌های فضایی معمولاً کنار گذاشته و نادیده گرفته می‌شود. در این پایان‌نامه، قصد داریم که هر دو این ویژگی‌های واقعی داده‌ها را در یک مدل رگرسیونی لحاظ کنیم. این کار مستلزم معرفی برخی مفاهیم و تعاریف اولیه است که در این فصل ارائه می‌کنیم.

¹Heterogeneous

²Spatial data

³Spatial econometric models

⁴Heterogeneity

۲.۱ مدل رگرسیون خطی

ساده‌ترین و کاربردی‌ترین مدل رگرسیونی، رگرسیون خطی است که معادله آن در حالت کلی به صورت

$$y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip} + \varepsilon_i \quad i = 1, 2, \dots, n \quad (1.1)$$

است (مونت‌گومری و همکاران، ۲۰۱۲). همچنین می‌توان معادله (۱.۱) را به صورت ماتریسی زیر بازنویسی کرد:

$$y = X\beta + \varepsilon, \quad (2.1)$$

که در آن $y = (y_1, \dots, y_n)^T$ بردار متغیرهای پاسخ، X ماتریس متغیرهای تبیینی با بعد $n \times (p+1)$ ، و β بردار ضرایب نامعلوم رگرسیونی با بعد $(p+1) \times 1$ است که بیان‌گر چگونگی رابطه خطی بین متغیرهای تبیینی و متغیر پاسخ است. همچنین ε بردار $n \times 1$ خطاهای تصادفی است.

در رگرسیون خطی، برای اعتبار استنباط‌های آماری انجام‌شده، پذیره‌هایی را در نظر می‌گیرند. این پذیره‌ها شامل موارد زیر هستند:

۱- رابطه بین X و y خطی بوده و متغیرهای تبیینی مستقل از یکدیگرند.

۲- میانگین جمله خطا برابر صفر است، یعنی $E(\varepsilon) = 0$.

۳- واریانس جمله خطا ثابت است، یعنی $\text{Var}(\varepsilon) = \sigma^2 I$. این پذیره بیان می‌کند که واریانس ε یا y به مقادیر X بستگی ندارد. این پذیره به واریانس متجانس^۵ (همگن^۶) معروف است.

۴- جملات خطا ناهم‌بسته‌اند.

۵- جملات خطا دارای توزیع نرمال هستند.

تحت پذیره‌های شماره ۴ و ۵، جمله‌های خطای $\varepsilon_1, \dots, \varepsilon_n$ مستقل هستند. به‌منظور برآورد ضرایب رگرسیونی، یک روش معمول استفاده از روش کمترین توان‌های دوم معمولی^۷ (OLS) است (مونت‌گومری و همکاران، ۲۰۱۲)، که در ادامه آن را بازگو می‌کنیم.

⁵Homoscedastic

⁶Homogenous

⁷Ordinary least square

۱.۲.۱ روش کمترین توان‌های دوم معمولی

در این روش، ضرایب رگرسیونی در مدل (۲.۱) از طریق کمینه کردن مجموع توان‌های دوم خطا برآورد می‌شوند. برای این منظور فرض می‌شود که همواره $n > p$ و همه خطاها ناهم‌بسته و دارای میانگین صفر و واریانس ثابت هستند. تابع هدف کمترین توان‌های دوم به صورت زیر است:

$$\begin{aligned} S(\beta) &= \varepsilon^T \varepsilon \\ &= (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \\ &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\beta - \beta^T \mathbf{X}^T \mathbf{y} + \beta^T \mathbf{X}^T \mathbf{X}\beta \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\beta + \beta^T \mathbf{X}^T \mathbf{X}\beta \end{aligned}$$

و برآورد کمترین توان‌های دوم ضرایب رگرسیونی از حل معادله زیر نتیجه می‌شود:

$$\frac{\partial S(\beta)}{\partial \beta} = -2\mathbf{y}^T \mathbf{X} + 2\beta^T \mathbf{X}^T \mathbf{X} = 0.$$

بنابراین

$$\begin{aligned} \mathbf{y}^T \mathbf{X} &= \beta^T \mathbf{X}^T \mathbf{X} \\ \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} &= \beta^T \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \end{aligned}$$

و در نهایت

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (3.1)$$

همان‌طور که گفته شد، برای اعتبار نتایج استنباط در مدل رگرسیون خطی باید چند پذیره پایه‌ای برای جمله خطای مدل برقرار باشند. اما در عمل و کاربرد ممکن است که هر یک از این پذیره‌ها برقرار نباشند. به‌طور ویژه ممکن است پذیره همگنی جمله‌های خطا معتبر نباشد. در این صورت برای مثلاً برآورد ضرایب، دیگر مجاز به استفاده از روش OLS نخواهیم بود و در صورت عدم توجه به این موضوع از کارایی و صحت مدل کاسته می‌شود. بنابراین ناهمگنی جمله‌های خطا باید در مدل لحاظ شود و این کار در رگرسیون از طریق انجام تبدیلات مناسب یا تغییر روش برآورد پارامترها امکان‌پذیر است. در ادامه روش‌های معمول برخورد با ناهمگنی واریانس جمله خطا را مطرح می‌کنیم.

۲.۲.۱ راه‌های برخورد با ناهمگنی در مدل رگرسیون خطی کلاسیک

در این بخش به بیان متداول‌ترین روش‌های برخورد با ناهمگنی واریانس در مدل‌های رگرسیون خطی کلاسیک می‌پردازیم.

تبدیل باکس-کاکس

در صورتی که جمله خطا و در نتیجه y ، نانرمال یا دارای واریانس ناهمگن باشد، تعدیل توانی y^λ یک راه مناسب برای رفع آن‌ها است. باکس و کاکس (۱۹۶۴) نشان دادند که پارامترهای رگرسیون و λ می‌توانند هم‌زمان، با استفاده از روش درست‌نمایی ماکسیمم برآورد شوند. اما این تبدیل در حالتی که λ برابر صفر شود، دچار مشکل می‌شود. لذا تبدیل مناسب برای استفاده را به صورت

$$y^{(\lambda)} = \begin{cases} \frac{y^\lambda - 1}{\lambda y^{\lambda-1}} & \lambda \neq 0 \\ y \ln y & \lambda = 0 \end{cases}$$

معرفی کردند، که در آن $\hat{y} = \ln^{-1}[1/n \sum \ln y_i]$ میانگین هندسی^۸ مشاهدات است و مدل

$$y^{(\lambda)} = \mathbf{X}\beta + \varepsilon$$

با روش کمترین توان‌های دوم یا درست‌نمایی ماکسیمم برآورد می‌شود.

تبدیل لگاریتمی

این تبدیل حالت خاصی از تبدیل باکس-کاکس ($\lambda = 0$) است. با تبدیل لگاریتمی غالباً ناهمسانی واریانس کاهش می‌یابد، زیرا مقیاس‌های اندازه‌گیری متغیرها را تحت تاثیر قرار می‌دهد. اما این تبدیل در مواردی که بعضی از مقادیر X و y صفر یا منفی باشند، کاربرد ندارد. از طرفی این تبدیل می‌تواند منجر به تولید همبستگی ساختگی شود.

کمترین توان‌های دوم تعمیم‌یافته

زمانی که جمله خطا دارای خودهمبستگی و واریانس ناهمگن است، یعنی $E(\varepsilon) = 0$ و

$$\text{Var}(\varepsilon) = \sigma^2 \mathbf{V}$$

که در آن \mathbf{V} یک ماتریس متقارن غیرقطری است، می‌توان از روش کمترین توان‌های دوم تعمیم‌یافته^۹ (GLS) برای برآورد ضرایب رگرسیونی استفاده کرد. در این روش با تبدیل داده‌ها به یک مجموعه داده جدید که دارای پذیره‌های لازم برای به کارگیری روش OLS است، خودهمبستگی و ناهمگنی در برآورد وارد می‌شود. البته چگونگی ساختار ماتریس \mathbf{V} موجب ایجاد مشکل در این روش می‌شود، اما با فرض این که \mathbf{V} نامنفرد^{۱۰} و مثبت باشد، ماتریس نامنفرد و متقارن \mathbf{K} را می‌توان طوری تعیین کرد که $\mathbf{K}\mathbf{K} = \mathbf{K}^T\mathbf{K} = \mathbf{V}$ و مدل جدید به صورت زیر تعریف می‌شود:

$$\mathbf{z} = \mathbf{B}\beta + \mathbf{g}$$

^۸Geometric mean

^۹Generalized least squares

^{۱۰}Nonsingular

که در آن

$$\mathbf{z} = \mathbf{K}^{-1}\mathbf{y}, \quad \mathbf{B} = \mathbf{K}^{-1}\mathbf{X}, \quad \mathbf{g} = \mathbf{K}^{-1}\boldsymbol{\varepsilon}.$$

در این مدل جدید پذیره‌های پایه روش OLS برقرار هستند. یعنی

$$E(\mathbf{g}) = E(\mathbf{K}^{-1}\boldsymbol{\varepsilon}) = \mathbf{K}^{-1}E(\boldsymbol{\varepsilon}) = \mathbf{0}$$

و

$$\begin{aligned} \text{Var}(\mathbf{g}) &= \text{Var}(\mathbf{K}^{-1}\boldsymbol{\varepsilon}) = \mathbf{K}^{-1}\text{Var}(\boldsymbol{\varepsilon})(\mathbf{K}^{-1})^T \\ &= \mathbf{K}^{-1}\sigma^2\mathbf{V}(\mathbf{K}^T)^{-1} = \mathbf{K}^{-1}\sigma^2\mathbf{V}(\mathbf{K})^{-1} \\ &= \sigma^2\mathbf{K}^{-1}\mathbf{K}\mathbf{K}^{-1} = \sigma^2\mathbf{I}. \end{aligned}$$

بنابراین فرآیند برآورد همانند OLS انجام می‌شود. به عبارت دقیق‌تر

$$\begin{aligned} S(\boldsymbol{\beta}) &= \mathbf{g}^T\mathbf{g} \\ &= \boldsymbol{\varepsilon}^T\mathbf{V}^{-1}\boldsymbol{\varepsilon} \\ &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \end{aligned}$$

و بنابراین

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{y}. \quad (4.1)$$

کمترین توان‌های دوم موزون

مدل رگرسیون خطی با واریانس خطای ناهمگن را می‌توان با روش کمترین توان‌های دوم موزون^{۱۱} (WLS) برازش داد. در این روش برآورد، تفاضل بین مقدار واقعی y_i و مقدار برآوردشده آن در وزن w_i ضرب می‌شود. این روش حالت خاصی از روش GLS است که در آن درایه‌های غیر قطر اصلی صفر هستند. تابع هدف روش WLS به صورت زیر نوشته می‌شود:

$$\begin{aligned} S(\boldsymbol{\beta}) &= \boldsymbol{\varepsilon}^T\mathbf{W}\boldsymbol{\varepsilon} \\ &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T\mathbf{W}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \mathbf{y}^T\mathbf{W}\mathbf{y} - \mathbf{y}^T\mathbf{W}\mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}^T\mathbf{X}^T\mathbf{W}\mathbf{y} + \boldsymbol{\beta}^T\mathbf{X}^T\mathbf{W}\mathbf{X}\boldsymbol{\beta} \\ &= \mathbf{y}^T\mathbf{W}\mathbf{y} - 2\mathbf{y}^T\mathbf{W}\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^T\mathbf{X}^T\mathbf{W}\mathbf{X}\boldsymbol{\beta} \end{aligned}$$

و برآورد کمترین توان‌های دوم موزون از حل معادله زیر نتیجه می‌شود:

$$\frac{\partial S(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = -2\mathbf{y}^T\mathbf{W}\mathbf{X} + 2\boldsymbol{\beta}^T\mathbf{X}^T\mathbf{W}\mathbf{X} = \mathbf{0}.$$

¹¹Weighted least squares

بنابراین

$$y^T \mathbf{W} \mathbf{X} (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} = \beta^T \mathbf{X}^T \mathbf{W} \mathbf{X} (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1}$$

و در نهایت

$$\hat{\beta} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{y}. \quad (5.1)$$

برای استفاده از این روش، ماتریس وزن \mathbf{W} باید مشخص شود. معمولاً وزن‌ها متناسب با معکوس واریانس جمله خطا تعریف می‌شوند؛ یعنی به مشاهدات با واریانس بزرگ وزن کوچک‌تر و به مشاهدات با واریانس کوچک وزن بزرگ‌تر تخصیص می‌یابد. برای مثال، اگر $\text{Var}(\varepsilon) = \sigma^2 \mathbf{V}$ باشد که در آن \mathbf{V} یک ماتریس قطری با بعد n است، آن‌گاه $\mathbf{W} = \mathbf{V}^{-1}$.

یکی دیگر از پذیره‌های پایه‌ای رگرسیون خطی، ناهمبسته بودن جمله خطاست. این پذیره نیز در بسیاری موارد در برخورد با داده‌های واقعی، برقرار نیست؛ یعنی جمله‌های خطا به شکلی به هم وابسته‌اند. این وابستگی در داده‌های سری زمانی (داده‌هایی که در طول زمان جمع‌آوری شده‌اند)، داده‌های طولی و داده‌های فضایی وجود دارد. در این پایان‌نامه، تمرکز بر داده‌های فضایی است؛ لذا برای درک بهتر ساختار وابستگی در این داده‌ها، در بخش بعد به بیان مفاهیم و مقدمات آمار فضایی می‌پردازیم.

۳.۱ مقدمه‌ای بر آمار فضایی

اغلب روش‌های معمولی آمار شامل رگرسیون خطی، OLS و WLS که در بخش قبل مطرح شدند مبتنی بر استقلال مشاهدات است، اما این پذیره در عمل در بسیاری از موارد برقرار نیست. کاربردهای مختلفی وجود دارند که در آن‌ها مشاهدات به یکدیگر وابسته بوده و نادیده گرفتن این وابستگی، سبب از بین رفتن اطلاعات مفیدی می‌شود و در نتیجه، تحلیل‌های آماری دور از واقعیت خواهند بود. در چنین مواقعی روش‌های معمول کارایی خود را از دست داده و باید روش دیگری را برای تحلیل مناسب داده‌ها انتخاب کرد.

این مسأله که مشاهده مربوط به یک موقعیت مکانی به مشاهدات موقعیت‌های مجاور خود وابسته باشد، کاملاً منطقی به نظر می‌رسد. این امر در داده‌های سری زمانی نیز به نوعی دیده می‌شود، با این تفاوت که در داده‌های زمانی به دلیل تک بعدی بودن زمان، مشاهدات در طول محور زمان، به‌طور متوالی و به ترتیب جمع‌آوری می‌شوند و در نتیجه وابستگی در مشاهدات پیاپی ظاهر می‌شود. اما گاهی وابستگی داده‌ها به وضعیت قرار گرفتن آن‌ها در فضای جغرافیایی تحت مطالعه بستگی دارد و در واقع وابستگی آن‌ها تابعی از فاصله مکانی آن‌ها از هم است. در برخورد با چنین داده‌هایی با دو دسته از عوامل روبه‌رو هستیم، که به

مکان استقرار متغیرها مربوط می‌شوند و می‌توان آن‌ها را تحت عناوین وابستگی فضایی^{۱۲} یا به بیان ضعیف‌تر خودهمبستگی فضایی^{۱۳} و ناهمگنی فضایی^{۱۴} توضیح داد. به شاخه‌ای از علم آمار که به بررسی و تحلیل این نوع داده‌ها می‌پردازد، آمار فضایی گفته می‌شود. در واقع در آمار فضایی سعی بر این است که وابستگی بین مقادیر مختلف یک متغیر، از حیث فاصله و جهت قرار گرفتن آن‌ها نسبت به هم، در مدل‌بندی لحاظ شود. این ارتباط فضایی که معمولاً در قالب روابط ریاضی بیان می‌شود، ساختار فضایی نام دارد.

در ادامه تعاریف و مفاهیم اولیه آمار فضایی را مطرح می‌کنیم. این مطالب از محمدزاده (۱۳۹۱)، انسلین (۱۹۸۸a) و لی‌سیج و پیس (۲۰۰۹) گرفته شده‌اند.

۱.۳.۱ داده‌های فضایی

در آمار فضایی متغیر مورد اندازه‌گیری ممکن است گسسته یا پیوسته و موقعیت مشاهدات نیز ممکن است گسسته یا پیوسته، نقطه‌ای یا ناحیه‌ای، منظم یا نامنظم باشد. مشاهدات فضایی با توجه به انواع موقعیت‌ها به سه گروه تقسیم می‌شوند:

زمین‌آماری^{۱۵}: این نوع داده‌ها در موقعیت‌های ثابت و مشخص در ناحیه‌ای پیوسته مشاهده می‌شوند. متغیر مورد بررسی ممکن است گسسته یا پیوسته باشد. در داده‌های زمین‌آماری معمولاً پیش‌گویی مقدار متغیر در یک موقعیت جدید مد نظر است.

مشبکه‌ای^{۱۶}: این نوع داده‌ها مربوط به مکان‌های ناحیه‌ای هستند، که این مکان‌ها ممکن است منظم یا نامنظم باشند. نواحی می‌توانند شامل شهرهای یک استان، مناطق شهرداری یک شهر و موارد مشابه باشند. به‌طور معمول هدف از تحلیل داده‌های مشبکه‌ای مدل‌بندی احتمالاتی مشاهدات است.

الگونقطه‌ای^{۱۷}: در این حالت مکان یا موقعیت مشاهده‌شده خود متغیری تصادفی است. الگوهای نقطه‌ای شامل تعدادی متناهی از مکان‌ها در یک ناحیه‌اند که در آن‌ها یک صفت خاص اندازه‌گیری می‌شود. این داده‌ها به سه دسته به‌طور کامل تصادفی فضایی، منظم و خوشه‌ای تقسیم و به مدل‌بندی آن‌ها اقدام می‌شود.

¹²Spatial dependence

¹³Spatial autocorrelation

¹⁴Spatial heterogeneity

¹⁵Geostatistical

¹⁶Lattice

¹⁷Point pattern

۲.۳.۱ مدل آماری

برای تحلیل داده‌های فضایی نیازمند یک مدل آماری هستیم. در آمار فضایی معمولاً از یک میدان تصادفی به صورت

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

برای مدل‌بندی داده‌ها استفاده می‌شود، که در آن D یک زیرمجموعه از فضای اقلیدسی d بعدی، $d \geq 1$ ، از \mathbb{R}^d است. مولفه $\mu(\cdot)$ تغییرات بزرگ‌مقیاس^{۱۸} یا روند^{۱۹} و $\delta(\cdot)$ فرایند خطا یا تغییرات کوچک‌مقیاس^{۲۰} نامیده می‌شوند. تغییرات کوچک‌مقیاس ممکن است ناشی از خطای اندازه‌گیری یا تغییر در درون موقعیت مشاهده‌شده و تغییرات بزرگ‌مقیاس ممکن است ناشی از تغییرپذیری بین موقعیت‌های مشاهده‌شده باشد.

در ادامه به معرفی برخی ویژگی‌های میدان تصادفی که موجب ساده‌سازی مسأله می‌شود، می‌پردازیم. باید توجه داشت که با وجود این‌که تمامی این پذیره‌ها موجب ساده‌سازی مسأله می‌شود، اما در صورت عدم برقراری، استفاده از مدل‌های فضایی بر پایه این پذیره‌ها موجب نامناسب شدن برآورد ساختار فضایی میدان و کاهش دقت نتایج حاصل از تحلیل داده‌های فضایی از جمله رگرسیون می‌شود.

تعریف ۱.۳.۱. میدان تصادفی $Z(\cdot)$ ایستای ذاتی^{۲۱} نامیده می‌شود، هرگاه میانگین میدان تصادفی مستقل از s باشد، یعنی $E(Z(s)) = \mu$ و واریانس تفاضل مقادیر $Z(s_1)$ و $Z(s_2)$ فقط تابعی از فاصله موقعیت‌های s_1 و s_2 باشد، یعنی

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

تعریف ۲.۳.۱. میدان تصادفی $Z(\cdot)$ ، ایستای مرتبه دوم^{۲۲} نامیده می‌شود، هرگاه میانگین آن ثابت و کواریانس $Z(s_1)$ و $Z(s_2)$ ، فقط تابعی از فاصله موقعیت‌ها باشد. تحت ایستایی مرتبه دوم، واریانس میدان تصادفی به موقعیت فضایی بستگی ندارد و تغییرپذیری میدان در همه جا یکسان است. یعنی

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

تعریف ۳.۳.۱. در صورتی که توزیع میدان تصادفی تحت هر فاصله h ثابت بماند، آن میدان را ایستای قوی^{۲۳} می‌گوییم. به عبارتی

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

¹⁸Large scale

¹⁹Trend

²⁰Small scale

²¹Intrinsic stationary

²²Second order stationary

²³Strong stationary

که منظور از $\stackrel{D}{=}$ هم‌توزیعی بردار متغیرهای تصادفی است.

تعریف ۴.۳.۱. در صورتی که توزیع توام هر تعداد متناهی از متغیرهای میدان تصادفی $Z(\cdot)$ نرمال چندمتغیره باشد، آن را میدان تصادفی گاوسی^{۲۴} گویند.

۳.۳.۱ ساختار همبستگی فضایی

در آمار فضایی مفاهیمی مشابه واریانس و کواریانس در آمار کلاسیک، برای بیان ساختار همبستگی فضایی وجود دارند.

تعریف ۵.۳.۱. واریانس تفاضل مقادیر میدان تصادفی در دو موقعیت s و $s+h$ را تغییرنگار^{۲۵} می‌نامند و به صورت زیر نمادگذاری می‌شود:

$$\gamma(h) = \text{Var}(Z(s+h) - Z(s)).$$

تغییرنگار معیاری برای نشان دادن میزان اثرات متقابل یک موقعیت با سایر موقعیت‌های فضای مورد مطالعه است. بنابراین کوچک بودن تغییرنگار، نشان‌گر وابستگی زیاد و بزرگ بودن آن نشانه وابستگی کم میدان تصادفی است. شیب نمودار تغییرنگار معرف همبستگی فضایی است، یعنی اگر برای تمام مقادیر h تغییرنگار ثابت باشد همبستگی فضایی داده‌ها بسیار ضعیف است.

تغییرنگار دارای سه پارامتر دامنه^{۲۶}، ازاره^{۲۷} و اثرقطعه‌ای^{۲۸} است.

دامنه: بازه تغییرات تغییرنگار که خارج از آن به حالت افقی درآمده و ثابت می‌ماند، دامنه نامیده می‌شود. خارج از دامنه، مشاهدات تقریباً اثری بر هم ندارند، یعنی تقریباً ناهمبسته هستند. از این‌رو داده‌های خارج از دامنه با روش‌های کلاسیک آماری قابل تحلیل‌اند.

إزاره: به‌طور معمول تغییرنگار تابعی افزایشی از تاخیر $\|h\|$ است و ممکن است به کران بالایی منتهی شود. منظور از $\|h\|$ ، فاصله اقلیدسی بین دو مشاهده فضایی با فاصله h در یک جهت مشخص است. چنین کرانی ازاره تغییرنگار نامیده شده و عبارت است از

$$\lim_{\|h\| \rightarrow \infty} \gamma(h) = c.$$

با توجه به طبیعت داده‌های فضایی که انتظار می‌رود در آن‌ها با افزایش فاصله مکانی، همبستگی داده‌ها کاهش یابد و در نتیجه تغییرنگار ثابت بماند، سعی می‌شود که بیش‌تر

²⁴Gaussian random field

²⁵Variogram

²⁶Range

²⁷Sill

²⁸Nugget effect

از مدل‌های کران‌دار استفاده شود. ازاره به عنوان یکی از پارامترهای مهم تابع تغییرنگار، از روی مشاهدات برآورد می‌شود. تعیین کران‌دار بودن تغییرنگار از روی مشاهدات مشکل است و در عمل با رسم نمودار تغییرنگار تجربی مقدار c برآورد می‌شود. مقداری که تابع تغییرنگار تجربی بعد از رسیدن به دامنه، حول آن نوسان می‌کند، ازاره را نمایش می‌دهد.

اثر قطعه‌ای: مقدار تغییرنگار در مبدا مختصات ($h = 0$)، اثر قطعه‌ای نامیده شده و به صورت

$$\lim_{\|h\| \rightarrow 0} \gamma(h) = c_0$$

تعریف می‌شود. به لحاظ نظری اثر قطعه‌ای باید صفر باشد، زیرا دو نمونه از موقعیت واحد، باید کمیت یکسانی داشته باشند. ولی در عمل تغییرنگارهای تجربی پیرو این وضعیت نیستند. دلایل این موضوع را می‌توان خطای نمونه‌گیری، خطای اندازه‌گیری یا تغییرات شدید کمیت مورد بررسی در موقعیت‌های نزدیک به هم دانست. استفاده از اثر قطعه‌ای می‌تواند معیاری برای تشخیص وجود ساختار فضایی در داده‌ها باشد. هر اندازه نسبت $\frac{c}{c+c_0}$ کوچکتر از ۰/۵ باشد، ساختار فضایی داده‌ها ضعیف‌تر است.

یکی از ویژگی‌های تغییرنگار، همیشه منفی شرطی^{۲۹} بودن آن است. همچنین ماترون (۱۹۷۱) ثابت کرد که اگر $\|h\| \rightarrow \infty$ آن گاه $\frac{\gamma(h)}{\|h\|^2} \rightarrow 0$.

مدل‌های معتبر تغییرنگار

مدل‌های معتبر تغییرنگار توابعی هستند که در شرط همیشه منفی شرطی صدق می‌کنند. جورنل و هویج‌برگتس (۱۹۷۸) مدل‌های پارامتری معتبر مختلفی را معرفی کردند. مدل‌های تغییرنگار شامل دو گروه با ازاره و فاقد ازاره هستند. مدل‌های تغییرنگار همسان‌گرد^{۳۰} دارای ازاره، به‌طور معمول بعد از دامنه تخت شده و با افزایش فاصله موقعیت‌ها تغییر معنی‌داری در آن‌ها صورت نمی‌پذیرد. همسان‌گردی را در ادامه تعریف خواهیم کرد. انواع مدل‌های نیم‌تغییرنگار معتبر همسان‌گرد دارای ازاره عبارتند از:

۱. **مدل گاوسی:** تابع نیم‌تغییرنگار گاوسی به‌صورت

$$\gamma(h) = c_0 + c(1 - e^{-\frac{\|h\|^2}{a^2}}), \quad h \in \mathbb{R}^d, \quad d \geq 1$$

است و با سه پارامتر a ، c و c_0 مشخص می‌شود، که در آن a دامنه، c_0 اثر قطعه‌ای و $c + c_0$ ازاره است.

²⁹Conditional negative definite

³⁰Isotropic

۲. مدل نمایی: تابع نیم‌تغییرنگار نمایی

$$\gamma(h) = c_0 + c(1 - e^{-\frac{\|h\|}{a}}), \quad h \in \mathbb{R}^d, \quad d \geq 1$$

دارای سه پارامتر a ، c و c_0 است. دامنه آن به‌طور معمول بسیار بزرگ است. وقتی داده‌ها در محدوده مورد بررسی دارای روند باشند یا نسبت به ابعاد محدوده تحت پوشش نمونه‌برداری دارای دامنه بزرگی باشند، تغییرنگار آن‌ها از مدل نمایی پیروی می‌کند.

۳. مدل کروی: ضابطه تابع نیم‌تغییرنگار کروی به‌صورت

$$\gamma(h) = \begin{cases} c_0 + c \left(\frac{3}{4} \frac{\|h\|}{a} - \frac{1}{4} \frac{\|h\|^2}{a^2} \right) & 0 < \|h\| \leq a, h \in \mathbb{R}^d, \quad d = 1, 2, 3 \\ c_0 + c & \|h\| > a \end{cases}$$

است.

به‌طور معمول مقدار مدل‌های نیم‌تغییرنگار بدون اِزاره با افزایش فاصله هم‌چنان افزایش پیدا می‌کند. انواع این مدل‌ها عبارت‌اند از:

۴. مدل توانی: تابع نیم‌تغییرنگار توانی با سه پارامتر λ ، a و c_0 به‌صورت

$$\gamma(h) = c_0 + a\|h\|^\lambda \quad 0 < \lambda < 2, \quad h \in \mathbb{R}^d, \quad d \geq 1$$

تعریف می‌شود. به پارامتر a شیب مدل گفته می‌شود. این مدل در حالت $\lambda = 1$ به مدل خطی منتهی می‌شود.

۵. مدل لگاریتمی: تابع نیم‌تغییرنگار لگاریتمی به‌صورت

$$\gamma(h) = 3a \ln(\|h\|)$$

است. نمودار این مدل دارای شیب $3a$ است که a پراکندگی مطلق نامیده می‌شود. چون برای $\|h\| < 1$ مقدار $\gamma(h)$ منفی است، این مدل برای نمونه‌گیری‌های با ابعاد کوچک قابل استفاده نیست.

۶. مدل موجی: تابع نیم‌تغییرنگار موجی به‌صورت

$$\gamma(h) = c_0 + c_w(1 - a_w \sin(-\|h\|/a_w)/\|h\|), \quad h \in \mathbb{R}^d, \quad d = 1, 2, 3$$

است و با سه پارامتر $a_w \geq 0$ ، $c_w \geq 0$ و $c_0 \geq 0$ مشخص می‌شود.

۷. مدل سهمی‌گون: تابع نیم‌تغییرنگار سهمی‌گون به‌صورت یک سهمی درجه دوم به‌صورت

$$\gamma(h) = \frac{1}{4} a^2 \|h\|^2$$

است. اگر میدان تصادفی $Z(\cdot)$ دارای روند خطی باشد، یعنی

$$Z(s+h) = Z(s) + ah$$

آن گاه مدل تغییرنگار آن سهمی گون خواهد بود، زیرا

$$\gamma(h) = E[(Z(s+h) - Z(s))^2] = a^2 h^2.$$

تعریف ۶.۳.۱. کواریانس دو متغیر $Z(s)$ و $Z(s+h)$ را هم‌تغییرنگار^{۳۱} می‌نامند و به صورت زیر نمادگذاری می‌شود:

$$C(h) = \text{Cov}(Z(s), Z(s+h)).$$

هم‌تغییرنگار میزان شباهت تغییرپذیری دو متغیر $Z(s)$ و $Z(s+h)$ را نشان می‌دهد. بنابراین در صورت وجود همبستگی فضایی، نقاط نزدیک به هم تشابه تغییرپذیری بیشتری دارند و مقدار هم‌تغییرنگار بزرگتر است.

تعریف ۷.۳.۱. ضریب همبستگی بین $Z(s)$ و $Z(s+h)$ همبستگی‌نگار^{۳۲} نامیده می‌شود و با شرط $C(0) > 0$ به صورت

$$\rho(s, s+h) = \rho(h) = \frac{C(h)}{C(0)} = 1 - \frac{\gamma(h)}{C(0)}$$

تعریف می‌شود. همبستگی‌نگار هم مانند تغییرنگار تابعی از فاصله بین موقعیت‌ها است و با افزایش فاصله، مقدار آن کم و با کاهش فاصله مقدار همبستگی‌نگار افزایش می‌یابد.

تعریف ۸.۳.۱. میدان تصادفی $Z(\cdot)$ را همسان‌گرد گویند هرگاه تغییرنگار، هم‌تغییرنگار و همبستگی‌نگار آن همسان‌گرد باشند، یعنی تنها تابعی از $\|h\|$ باشند و به جهت h بستگی نداشته باشند.

تعریف ۹.۳.۱. میدان تصادفی $Z(\cdot)$ که ایستای مرتبه دوم و همسان‌گرد باشد را همگن گویند. در غیر این صورت میدان ناهمگن است.

۴.۳.۱ ناهمگنی فضایی

ناهمگنی فضایی یکی دیگر از اثرات فضایی و صرفاً نایستایی ساختاری است. این نایستایی به دو شکل بی‌ثباتی در واریانس (واریانس نامتجانس^{۳۳}) یا ضرایب مدل (ضرایب متغیر^{۳۴})، رژیم‌های فضایی^{۳۵})، بروز پیدا می‌کند. این دو جنبه از ناهمگنی کاملاً مجزا از یکدیگر هستند

³¹Covariogram

³²Correlogram

³³Heteroscedasticity

³⁴Variable coefficients

³⁵Spatial regimes

(انسلین، ۱۹۸۸a). اولی، نتیجه حذف متغیرهای تبیینی لازم یا سایر شکل‌های مدل‌بندی نادرست، که منجر به بی‌ثباتی واریانس جمله خطا می‌شود، و دومی ناشی از ثابت نبودن ارتباط بین متغیر پاسخ و متغیرهای تبیینی در کل فضای مورد مطالعه است. همان‌طور که در بخش قبل گفته شد، می‌توان این ناهمگنی را با ابزار و تبدیلات مختلف استانداردسازی، سنجید. اما سه دلیل وجود دارند که تایید می‌کنند ناهمگنی را به صراحت در مدل در نظر بگیریم.

اول این که ساختار این نایستایی از جنس فضایی است. به این معنی که مکان مشاهدات در تعیین شکل نایستایی مهم است. به‌عنوان مثال، برای مدل‌بندی واریانس نامتجانس با طرح گروهی، می‌توان ناحیه مورد مطالعه را به زیرناحیه‌هایی با واریانس متفاوت تقسیم کرد. برای تفهیم بیشتر، یک مجموعه S از N واحد جغرافیایی (مثل استان‌ها، شهرها و سرشماری تراکم) را در نظر بگیرید. این مجموعه را به R زیرمجموعه S_r به هم متصل که هم‌پوشانی ندارند، تقسیم می‌کنیم به طوری که برای هر $r \neq s$ $S_r \cap S_s = \emptyset$ و $\bigcup_{r=1, \dots, R} S_r = S$. با این روش از این پس واریانس نامتجانس گروهی، به صورت واریانس خطای خوشه‌بندی فضایی برای هر مشاهده دنبال می‌شود. یعنی برای هر $i \in S_r$ $\text{Var}(\varepsilon_i) = \sigma^2$. به‌طور مشابه، تنوع در ضرایب می‌تواند با رژیم‌های فضایی یا زیرمجموعه‌های جغرافیایی داده‌ها که در آن‌ها برای هر $i \in S_r$ $\beta_i = \beta_r$ مشخص شود. به بیان دیگر پارامترهای مدل در زیرمجموعه‌هایی از ناحیه تحت مطالعه مقادیر مختلف اختیار می‌کنند؛ یعنی رابطه بین متغیر تبیینی و متغیر پاسخ در کل فضای مطالعه پایا نیست.

دوم این که به دلیل ساختار فضایی، ناهمگنی فضایی و وابستگی فضایی به صورت مشترک رخ می‌دهند و جداسازی آن‌ها مشکل است. این مسأله در متون تخصصی به‌عنوان مشکل معکوس^{۳۶} شناخته می‌شود. همچنین به نشدنی بودن تمایز بین آلودگی^{۳۷} واقعی و ظاهری مربوط است (انسلین، ۲۰۱۰).

سوم این که در یک بخش‌بندی کلی، خودهمبستگی فضایی و ناهمگنی فضایی ممکن است معادل در نظر گرفته شوند. یعنی مثلاً یک خوشه فضایی از باقی‌مانده‌های کرانگین^{۳۸} ممکن است به علت وجود ناهمگنی فضایی (مثلاً واریانس نامتجانس گروهی) یا خودهمبستگی فضایی تفسیر شود. لذا لازم است که هر دو جنبه ناهمگنی و وابستگی فضایی با هم در نظر گرفته شوند.

اما اکثراً در تحقیقات و مدل‌های فضایی نظیر مدل میدان تصادفی گاوسی، مدل‌های اتورگرسیو و مدل‌های اقتصادسنجی فضایی، تمرکز بر وابستگی فضایی است و اثر ناهمگنی فضایی به‌انزوا رفته است (انسلین، ۱۹۹۹). به‌عنوان مثال، نوسانات قیمت مسکن عموماً از کلان‌شهرها شروع شده و در مراحل بعد به سایر شهرها سرریز می‌شود. لذا بحث مجاورت

³⁶Inverse problem

³⁷Contagion

³⁸Extreme

و وابستگی فضایی مطرح می‌شود و در تحلیل‌های فضایی وارد می‌شود؛ اما این موضوع نیز وجود دارد که متغیرهای اثرگذار بر قیمت مسکن در شهرهای مختلف دارای تاثیرات متفاوت هستند و از مکانی به مکان دیگر تغییر می‌کنند. هنگامی که این تنوع بیش از حد رخ دهد، نشان از ناهمگنی فضایی است و باید در مدل لحاظ شود. در فصل دوم، مدل رگرسیون موزون جغرافیایی^{۳۹} (GWR) را که در آن به این امر توجه شده است و ناهمگنی فضایی نیز در کنار وابستگی فضایی در مدل لحاظ می‌شود، معرفی می‌کنیم.

۴.۱ سایر مفاهیم و مقدمات

در این بخش یک روش مدل‌بندی ساختار وابستگی فضایی و ماتریس وزن فضایی به‌عنوان ابزاری برای این مدل‌بندی مطرح می‌شود. همچنین آزمون‌های کاربردی معمول، برای تشخیص وجود اثرات فضایی در داده‌ها معرفی می‌شوند.

۱.۴.۱ مدل خودهمبستگی

یکی از شاخه‌های جالب و در حال رشد آمار فضایی مربوط به خودهمبستگی فضایی است. خودهمبستگی فضایی مفهومی نسبتاً ساده و در حقیقت بسط همین مفهوم در آمار متعارف است. خودهمبستگی، خواه از نوع فضایی باشد یا نه، یک معیار برای بیان تشابه و همبستگی بین مشاهدات همسایه است. مدل

$$\mathbf{u} = \rho \mathbf{W}\mathbf{u} + \varepsilon \quad (۶.۱)$$

برای بیان خودهمبستگی به کار می‌رود که در آن همسایگی‌ها از طریق ماتریس وزن \mathbf{W} معرفی می‌شوند. در این مدل از میانگین وزنی همسایگی‌های مشاهده i ، یعنی

$$u_i = w_{i1}u_1 + w_{i2}u_2 + \dots + w_{iN}u_N$$

برای برآورد مقدار آن استفاده می‌شود. پارامتر ρ ضریب خودهمبستگی ترکیب خطی مشاهده u_i با همسایگی‌های خود است و متوسط تاثیر مشاهدات همسایه بر مقدار مشاهده i را اندازه‌گیری می‌کند. این موضوع همان‌طور که یک چالش محسوب می‌شود و موجب پیچیدگی آزمون‌های آماری می‌شود، اما با اطلاعات فضایی نیز مطابقت دارد.

از این رو می‌توان مدل خودهمبستگی را با رگرسیون خطی ترکیب و ساختار فضایی را مدل‌بندی کرد. این مدل‌ها با عنوان مدل‌های اتورگرسیون فضایی^{۴۰} شناخته می‌شوند. در این مدل‌بندی جمله خودهمبستگی، که در آن ماتریس \mathbf{W} وزن همسایگی‌های مکانی را نشان

³⁹Geographical weighted regression

⁴⁰Spatial autoregressive model

می‌دهد، با توجه به شرایط داده‌ها در متغیر پاسخ، متغیرهای تبیینی، جمله خطا یا ترکیبی از این‌ها قرار می‌گیرد و مدل‌های رگرسیون فضایی مختلفی را تولید می‌کند که در فصل دوم تحت عنوان مدل‌های اتورگرسیو شرطی، اتورگرسیو هم‌زمان و مدل‌های اقتصادسنجی به تشریح آن‌ها می‌پردازیم.

۲.۴.۱ ماتریس وزن فضایی

ماتریس وزن فضایی، ابزاری است که شدت رابطه بین مشاهدات در یک همسایگی را به صورت کمی نشان می‌دهد. این ماتریس طوری تعریف می‌شود که موافق با ساختار همبستگی فضایی است؛ یعنی به همسایه‌های نزدیک وزن بیشتر و به همسایه‌های دورتر وزن کمتری اختصاص داده می‌شود. ماتریس وزن فضایی یک ماتریس مربعی از مرتبه n (تعداد مشاهدات) است که مولفه‌های آن، یعنی w_{ij} ، دارای خواص زیر هستند:

۱. نامنفی هستند.

۲. به ازای $i = 1, \dots, n$ ، $w_{ii} = 0$ و $w_{ij} = w_{ji}$.

به منظور تعیین w_{ij} دو منبع اطلاعاتی موجود است. یکی موقعیت مشاهدات در صفحه مختصات که از طریق طول و عرض جغرافیایی، فاصله هر نقطه در فضا را نسبت به نقاط دیگر نشان می‌دهد. دومی مجاورت و همسایگی است که اغلب در مورد داده‌های شبکه‌ای مورد استفاده قرار می‌گیرد و بر این اساس که مناطق i و j حاشیه یا مرز مشترک دارند یا خیر وزن w_{ij} را تعیین می‌کند. در ادامه تعدادی از روش‌های وزن‌دهی بر پایه این دو منبع اطلاع را بیان می‌کنیم.

ماتریس فاصله^{۴۱}: در این ماتریس میزان تاثیر و ارتباط هر جفت مشاهده تابعی از فاصله آن‌ها است. یعنی

$$w_{ij} = \begin{cases} d_{ij} & i \neq j \\ 0 & i = j \end{cases} \quad (7.1)$$

ماتریس معکوس-فاصله^{۴۲}: در این حالت شدت رابطه بین دو مشاهده با فاصله آن‌ها از هم نسبت عکس دارد. یعنی

$$w_{ij} = \begin{cases} \frac{1}{d_{ij}} & i \neq j \\ 0 & i = j \end{cases} \quad (8.1)$$

⁴¹Distance matrix

⁴²Inverse-distance matrix

ماتریس مجاورت ضلعی^{۴۳}: این ماتریس وزن در خصوص داده‌های شبکه‌ای منظم مورد استفاده قرار می‌گیرد و هرگاه نواحی i و j دارای ضلع مشترک باشند، $w_{ij} = 1$ و در غیر این صورت $w_{ij} = 0$.

ماتریس مجاورت رأسی^{۴۴}: در این روش وزن‌دهی، دو ناحیه همسایه یکدیگر محسوب می‌شوند و w_{ij} مقدار ۱ می‌گیرد، اگر دارای یک رأس مشترک باشند. این روش نیز در خصوص داده‌های شبکه‌ای منظم به کار می‌رود.

ماتریس مجاورت شعاعی^{۴۵}: در این روش که برای داده‌های شبکه‌ای منظم به کار می‌رود، دایره‌ای به مرکز ناحیه i و شعاع طول ضلع نواحی رسم می‌شود. به نواحی‌ای که دایره از آن‌ها عبور می‌کند یعنی نواحی‌ای که دارای ضلع و/یا رأس مشترک با ناحیه i هستند، وزن ۱ و در غیر این صورت وزن صفر تخصیص داده می‌شود.

در خصوص داده‌های شبکه‌ای نامنظم، نواحی‌ای که دارای مرز مشترک هستند معمولاً همسایه محسوب می‌شوند.

۳.۴.۱ آماره آزمون موران نوع I

متداول‌ترین آزمون حضور خودهمبستگی فضایی توسط یک آماره که مشابه آن اولین بار توسط موران (۱۹۴۸) برای بررسی همبستگی سری‌های زمانی معرفی شد، انجام می‌شود. این آماره که به تحلیل الگوهای پراکنش و توزیع خطاها در فضا می‌پردازد، به آماره موران I^{۴۶} معروف است. صورت ماتریسی آماره موران به صورت:

$$I = (n/S_0)(e^T W e / e^T e)$$

است، که در آن e بردار مقادیر باقی‌مانده حاصل از رگرسیون OLS است و $S_0 = \sum_i \sum_j w_{ij}$. مقدار امید و واریانس این آماره تحت فرضیه نبود خودهمبستگی فضایی (فرضیه صفر) به صورت زیر است:

$$E(I) = \frac{tr(MW)}{n-p} \quad (9.1)$$

$$Var(I) = \frac{tr(MWMW^T) + tr(MWMW) + [tr(MW)]^2}{(n-p)(n-p+2)} - [E(I)]^2 \quad (10.1)$$

که در آن ماتریس M برابر است با

$$M = I - X(X^T X)^{-1} X^T.$$

⁴³Edge contiguity matrix

⁴⁴Vertex contiguity matrix

⁴⁵Radius contiguity matrix

⁴⁶Moran I statistic

این آماره نشان می‌دهد که پراکندگی مقادیر خطا با در نظر گرفتن متغیرهای مورد مطالعه از یک الگوی خوشه‌ای پیروی می‌کند یا نه. اگر باقی‌مانده‌ها یا مقادیر متغیرهای مربوط به آن‌ها به‌طور تصادفی در ناحیه تحت مطالعه توزیع شده باشند، نباید بین آن‌ها ارتباطی وجود داشته باشد. به‌طور کلی مقادیر مثبت یا منفی آماره موران نشان از خودهمبستگی مثبت یا منفی و همچنین مقادیر نزدیک به مقدار مورد انتظار (۹.۱)، بیانگر نبود خودهمبستگی یعنی الگوی فضایی تصادفی است. محدوده تغییرات این آماره برخلاف ضریب همبستگی، بازه $[-1, 1]$ نیست.

همچنین آماره

$$Z = \frac{I - E(I)}{\sqrt{\text{Var}(I)}} \quad (11.1)$$

برای نمونه‌های به اندازه کافی بزرگ دارای توزیع نرمال استاندارد است. اگر $|Z| < z_{\alpha/2}$ آن‌گاه آزمون خودهمبستگی فضایی در سطح α معنی‌دار نبوده و بین جمله‌های خطا در مدل رگرسیونی، خودهمبستگی فضایی وجود ندارد.

۴.۴.۱ آزمون‌های مبتنی بر درستنمایی

با توجه به استفاده گسترده از روش درستنمایی ماکسیمم در برآورد مدل‌های فضایی، اکثر آزمون‌های فرضیه برای پارامترهای این مدل‌ها مبتنی بر ملاحظات جانبی است. دو نمونه از این آزمون‌های جانبی بر پایه درستنمایی ماکسیمم، آزمون نسبت درستنمایی^{۴۷} (LR) و آزمون ضریب لاگرانژ^{۴۸} (LM) است.

آزمون نسبت درستنمایی (LR)

آزمون نسبت درستنمایی، فرضیه‌های

$$\begin{cases} H_0 : \lambda = 0 \\ H_1 : \lambda \neq 0 \end{cases}$$

که در آن λ ضریب خودهمبستگی جمله خطا است را مورد آزمون قرار می‌دهد. آماره این آزمون بر پایه تفاضل تابع درستنمایی مدل رگرسیونی فضایی با حضور λ و تابع درستنمایی آن تحت فرضیه صفر بنا شده است و به‌صورت زیر معرفی می‌شود:

$$LR = n \left(\ln \sigma_0^2 - \ln \sigma^2 \right) + 2 \ln |\mathbf{I} - \lambda \mathbf{W}|. \quad (12.1)$$

این آماره دارای توزیع کای-دو با یک درجه آزادی است که در آن n تعداد مشاهدات و σ_0^2 و σ^2 به‌ترتیب برآورد درستنمایی ماکسیمم واریانس خطا در مدل فضایی و برآورد آن تحت فرضیه

⁴⁷Likelihood ratio

⁴⁸Lagrange multiplier

صفر است. بزرگی آماره LR نشان از اختلاف بین این دو مدل است و موجب رد فرضیه صفر می‌شود؛ یعنی خودهمبستگی فضایی در جمله‌های خطا موجود است.

آزمون ضریب لاگرانژ (LM)

آزمون ضریب لاگرانژ به دو شکل مختلف برای تشخیص حضور خودهمبستگی در جمله خطا و متغیر وابسته وجود دارد، که به ترتیب آزمون ضریب لاگرانژ خطای فضایی^{۴۹} (LM_λ) و آزمون ضریب لاگرانژ تاخیر فضایی^{۵۰} (LM_ρ) نامیده می‌شوند.

الف) آزمون ضریب لاگرانژ خطای فضایی

آماره آزمون ضریب لاگرانژ خطا به صورت زیر تعریف می‌شود:

$$LM_\lambda = \frac{1}{T} \left(\frac{\mathbf{e}^T \mathbf{W} \mathbf{e}}{S^2} \right)^2. \quad (13.1)$$

این آماره در مدل رگرسیون

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{v} \quad \mathbf{v} = \lambda \mathbf{W}\mathbf{v} + \varepsilon \quad (14.1)$$

که در آن $\varepsilon \sim MVN(0, \sigma^2 \mathbf{I})$ است، فرضیه‌های زیر را مورد آزمون قرار می‌دهد:

$$\begin{cases} H_0 : \lambda = 0 \\ H_1 : \lambda \neq 0 \end{cases}$$

یعنی آزمون می‌کند که آیا جمله‌های خطا در بردار \mathbf{v} ناهمبسته هستند یا خیر. این مدل رگرسیونی (مدل خطای فضایی^{۵۱} (SEM)) در فصل بعد به تفصیل بیان می‌شود. در این آماره S^2 برآورد درست‌نمایی ماکسیمم پارامتر σ^2 است که به صورت زیر تعریف می‌شود:

$$S^2 = \frac{\mathbf{e}^T \mathbf{e}}{n}$$

و همچنین

$$T = \text{tr}(\mathbf{W}^T \mathbf{W} + \mathbf{W}^2).$$

آماره آزمون (۱۳.۱) جدا از جمله مقیاس‌گذاری $\text{tr}(\mathbf{W}^T \mathbf{W} + \mathbf{W}^2)$ برابر توان دوم آماره موران است و تحت درستی فرضیه H_0 ، دارای توزیع تقریبی کای-دو با ۱ درجه آزادی است. تحت فرضیه H_1 ، لگاریتم تابع درست‌نمایی بردار \mathbf{y} با توجه به (۱۴.۱) به صورت زیر است:

$$\ln L = -\frac{n}{2} \ln(2\pi\sigma^2) + \ln |\mathbf{I}_n - \lambda \mathbf{W}| - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{I}_n - \lambda \mathbf{W}) (\mathbf{y} - \mathbf{X}\beta).$$

⁴⁹Lagrange multiplier error test

⁵⁰Lagrange multiplier lag test

⁵¹Spatial error model

برای به دست آوردن برآوردهای درست‌نمایی ماکسیمم پارامترها تحت فرضیه H_0 ، تابع لاگرانژ را به صورت زیر تعریف می‌کنیم:

$$L_R = \ln L - t\lambda.$$

اکنون مشتق تابع لاگرانژ را نسبت به تک تک پارامترها به دست می‌آوریم:

$$\begin{aligned} \frac{\partial L_R}{\partial \beta^T} &= \frac{1}{\sigma^2} \mathbf{X}^T (\mathbf{I}_n - \lambda \mathbf{W})^T (\mathbf{I}_n - \lambda \mathbf{W}) (\mathbf{y} - \mathbf{X}\beta) = 0 \\ \frac{\partial L_R}{\partial \sigma^2} &= \frac{1}{2\sigma^2} \left\{ n + \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{I}_n - \lambda \mathbf{W})^T (\mathbf{I}_n - \lambda \mathbf{W}) (\mathbf{y} - \mathbf{X}\beta) \right\} = 0 \\ \frac{\partial L_R}{\partial \lambda} &= -t - \sum_{j=1}^n \delta_j (1 - \lambda \delta_j)^{-1} + \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^T \{ (\mathbf{I}_n - \lambda \mathbf{W})^T \mathbf{W} \\ &\quad + \mathbf{W}^T (\mathbf{I}_n - \lambda \mathbf{W}) \} (\mathbf{y} - \mathbf{X}\beta) = 0 \\ \frac{\partial L_R}{\partial t} &= -\lambda = 0 \end{aligned} \quad (15.1)$$

به طوری که δ_j ها مقادیر ویژه ماتریس \mathbf{W} هستند.

با حل معادلات (15.1)، برآورد ماکسیمم درست‌نمایی پارامتر t به صورت زیر به دست می‌آید:

$$\hat{t} = -\text{tr}(\mathbf{W}) + \frac{1}{2\hat{\sigma}^2} (\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{W} + \mathbf{W}^T) (\mathbf{y} - \mathbf{X}\hat{\beta})$$

به طوری که $\hat{\beta}$ برآورد OLS پارامتر β است. هم‌چنین

$$s^2 = \hat{\sigma}^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{y} - \mathbf{X}\hat{\beta})}{n} = \frac{\mathbf{e}^T \mathbf{e}}{n}.$$

ار آن‌جا که $\text{tr}(\mathbf{W}) = 0$ ، پس

$$\hat{t} = \frac{n\mathbf{e}^T (\mathbf{W} + \mathbf{W}^T) \mathbf{e}}{2\mathbf{e}^T \mathbf{e}}.$$

تحت فرضیه H_0 ، تابع اطلاع پارامتر λ به صورت زیر به دست می‌آید:

$$T = -E_{H_0} \left(\frac{\partial^2 \ln L}{\partial \lambda^2} \right) = \text{tr}(\mathbf{W}^2 + \mathbf{W}^T \mathbf{W}).$$

بنابراین آماره آزمون، بر اساس رفتار $\hat{t} [\text{tr}(\mathbf{W}^2 + \mathbf{W}^T \mathbf{W})]^{-\frac{1}{2}}$ ساخته می‌شود. برای اطلاعات بیشتر درباره اثبات بالا به بوریچ (۱۹۸۰) مراجعه کنید.

(ب) آزمون ضریب لاگرانژ تأخیر فضایی

این‌بار فرضیه مورد بررسی، وجود خودهمبستگی فضایی در متغیر وابسته مدل است. در صورت رد فرضیه صفر، مدل زیر، که به مدل تأخیر فضایی δ^2 (SLM) معروف است و در فصل بعد به تفصیل معرفی خواهد شد، به داده‌ها برازش داده می‌شود:

$$\mathbf{y} = \rho \mathbf{W}\mathbf{y} + \mathbf{X}\beta + \varepsilon. \quad (16.1)$$

آزمون فرضیه مورد نظر به صورت زیر تعریف می‌شود:

$$\begin{cases} H_0 : \rho = 0 \\ H_1 : \rho \neq 0 \end{cases} \quad (17.1)$$

آماره آزمون ضریب لاگرانژ برای فرضیه‌های (۱۷.۱)، به صورت زیر است:

$$LM_\rho = \frac{1}{J} \left(\frac{\mathbf{e}^T \mathbf{W} \mathbf{y}}{S^2} \right)^2 \quad (18.1)$$

به طوری که

$$J = \frac{[(\mathbf{W}\mathbf{X}\hat{\beta})^T \mathbf{M}(\mathbf{W}\mathbf{X}\hat{\beta}) + TS^2]}{S^2}$$

برای محاسبه آماره آزمون (۱۸.۱) ابتدا مدل (۱۶.۱) که در آن $\varepsilon \sim MVN(0, \sigma^2 \mathbf{I})$ است تحت فرضیه صفر به کمک روش OLS برآورد می‌شود، به طوری که $\hat{\beta}$ برآورد حداقل مربعات معمولی بردار پارامتر β است. آماره آزمون (۱۸.۱) نیز، تحت درستی فرضیه H_0 ، دارای توزیع تقریبی کای-دو با ۱ درجه آزادی است. اگر آزمون (۱۷.۱) معنی‌دار باشد، به این معناست که باید جمله تأخیری $\rho \mathbf{W} \mathbf{y}$ نیز در مدل رگرسیونی لحاظ شود. در غیر این صورت، یعنی اگر معنی‌دار نباشد، نیازی به در نظر گرفتن این جمله در مدل نیست. چگونگی به دست آمدن آماره آزمون (۱۸.۱) در بوریچ (۱۹۸۰) آمده است.

۵.۴.۱ آزمون ناهمگنی در حضور خودهمبستگی

زمانی که جمله‌های خطا دارای وابستگی باشند، توابع پارامتری معمول برای آزمون ناهمگنی معتبر نیستند. چرا که توزیع مجانبی آماره این آزمون‌ها بر اساس پذیره استقلال توزیع خطاها است و در غیاب استقلال این نتایج حفظ نخواهند شد. بنابراین با اثبات حضور خودهمبستگی فضایی، نیازمند روشی هستیم که ناهمگنی را با وجود خودهمبستگی آزمون کند.

انسلین (۱۹۸۸a) برای آزمون وجود ناهمگنی در حضور خودهمبستگی فضایی آماره‌ای را معرفی کرده است. این آماره با ایجاد تغییراتی در ضابطه اصلی آماره بریچ-پیگان^{۵۳} (بریچ و پیگان، ۱۹۷۹) که برای آزمون ناهمگنی در شرایط استقلال به کار می‌رود، به شکل زیر معرفی شده است:

$$LM = (1/4\sigma^4) \mathbf{f}^T \mathbf{Z} \mathbf{I}^{-1} \mathbf{Z}^T \mathbf{f} \quad (19.1)$$

که در آن \mathbf{Z} یک ماتریس $n \times (p+1)$ از بردارهای \mathbf{z} برای هر مشاهده است، $f_i = (\sigma^{-1} e_i)^2 - 1$ و e_i ها درایه‌های بردار $\mathbf{e} = (1-\lambda)(\mathbf{y} - \mathbf{X}\beta)$ ، یعنی باقی مانده‌های درستی‌مابی ماکسیمم رگرسیون فضایی با خطاهای خودهمبسته هستند. σ^2 نیز برآورد ML واریانس خطای رگرسیون فضایی است. برای اطلاع بیشتر در مورد این آماره به انسلین (۱۹۸۸a) مراجعه کنید.

⁵³Breusch-Pagan

۶.۴.۱ ضرب هادامارد

تعاریف متعددی برای ضرب ماتریس‌ها وجود دارند. یکی از این تعاریف، ضرب هادامارد^{۵۴} یا ضرب درایه‌ای است (میلین، ۲۰۰۷). در این روش دو ماتریس باید دارای ابعاد یکسان باشند و ماتریس حاصل ماتریسی با همان ابعاد است. ضرب هادامارد دو ماتریس A و B به صورت زیر تعریف می‌شود:

$$[A \odot B]_{ij} = [A]_{ij}[B]_{ij} \quad 1 \leq i \leq m, 1 \leq j \leq n.$$

روش هادامارد در الگوریتم‌های فشرده‌سازی (الگوریتم‌های رایج برای فشرده‌سازی داده‌های چندرسانه‌ای مثل الگوریتم JPEG) استفاده می‌شود.

⁵⁴Hadamard

فصل ۲

مدل بندی ساختار وابستگی فضایی

۱.۲ مقدمه

به طور کلی در تحلیل مدل های رگرسیونی خطاهای مدل ناهمبسته در نظر گرفته می شوند. اما این پذیره در بسیاری از موقعیت های کاربردی برقرار نیست. به ویژه در تحلیل داده های فضایی که مشاهدات برحسب موقعیت به یکدیگر وابسته اند و مشاهدات نزدیک به هم همبستگی بیشتری دارند، این پذیره اشتباه است. بنابراین لازم است که این امر در مدل بندی داده های فضایی مورد توجه قرار گیرد.

یک راه مرسوم برای لحاظ کردن ساختار وابستگی فضایی، استفاده از مدل های اتورگرسیو فضایی است. در طول سال های اخیر، مدل های مختلفی بر این مبنا معرفی شده اند و پیوسته به روز می شوند؛ به عنوان مثال می توان به دو مدل متداول و بسیار محبوب داده های شبکه ای یعنی اتورگرسیو شرطی^۱ (CAR) و اتورگرسیو هم زمان^۲ (SAR) اشاره کرد. هر دو مدل وابستگی فضایی را در ساختار ماتریس کواریانس، به عنوان یک تابع از ماتریس وزن W و اغلب یک پارامتر همبستگی فضایی ناشناخته ثابت، تولید می کنند. دسته ای دیگر از مدل های اتورگرسیو نیز وجود دارند که به دلیل کاربرد بسیار آنها در مسائل اقتصادی (نظیر مطالعات بازرگانی، جمعیت شناسی و مالی) به مدل های اقتصادسنجی معروف شده اند. در این فصل به توضیح هر

^۱Conditional autoregressive model

^۲Simultaneously autoregressive model

یک از این مدل‌ها می‌پردازیم و به دلیل ناتوانی آن‌ها در مدل‌بندی و تفسیر ناهمگنی فضایی، مدل رگرسیون موزون جغرافیایی را معرفی می‌کنیم.

۲.۲ وارد کردن ساختار وابستگی از طریق ماتریس کواریانس

همان‌گونه که گفته شد، یک راه لحاظ کردن وابستگی فضایی در مدل رگرسیونی، وارد کردن آن از طریق ساختار وابستگی مقادیر خطاست. در مدل رگرسیونی پایه (۲.۱) جمله خطا یک متغیر تصادفی مستقل نرمال با میانگین صفر و واریانس ثابت معرفی شد، یعنی

$$y \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}). \quad (1.2)$$

اما در مورد داده‌های دارای ساختار فضایی، این پذیره ناهمبستگی و استقلال قابل قبول نیست زیرا مشاهدات حداقل از طریق مکان قرارگیری به یکدیگر وابسته هستند. در این موارد بهتر است که مدل (۱.۲) به صورت $y \sim N(\mu, \Lambda)$ نوشته شود، که در آن $\mu = \mathbf{X}\beta$ و Λ یک ماتریس واریانس متقارن مثبت است. در این ماتریس به دلیل ساختار فضایی موقعیت‌هایی که به یکدیگر نزدیک‌ترند، دارای کواریانس بزرگ‌ترند.

بنابراین مدل رگرسیون فضایی به صورت عمومی زیر نوشته می‌شود:

$$y = \mathbf{X}\beta + \varepsilon \quad \varepsilon \sim N(0, \Lambda) \quad (2.2)$$

به دلیل ساختار پیچیده معمول ماتریس Λ ، برای ساده‌سازی مدل‌بندی، پذیره‌هایی به مدل اضافه می‌شوند. در ادامه به معرفی این روش‌ها که مورد توجه این پایان‌نامه است، می‌پردازیم.

ساختار وابستگی داده‌های شبکه‌ای

فرض کنید $\{Y(A_i) : A_i \in (A_1, \dots, A_n)\}$ یک میدان تصادفی نرمال است که $\{A_1, \dots, A_n\}$ نواحی شبکه‌ای هستند یعنی $A_1 \cup A_2 \cup \dots \cup A_n = D$ و $A_i \cap A_j = \emptyset \quad \forall i \neq j$. در این بخش دو مدل متداول اتورگرسیو شرطی و اتورگرسیو هم‌زمان برای مدل‌بندی این ساختار معرفی می‌شوند.

مدل اتورگرسیو شرطی

یک راه مدل‌بندی $\{Y(A_i) : A_i \in (A_1, \dots, A_n)\}$ با مشاهدات y_i ، مدل اتورگرسیو شرطی است که به صورت

$$y_i | \{y_j, j \neq i\} \sim N(\mu_i + \sum_{j=1}^n c_{ij}(y_j - \mu_j), \tau^2) \quad (3.2)$$

معرفی می‌شود. از این پس $y(A_i)$ به اختصار y_i نوشته می‌شود. در مدل‌های CAR، علاوه بر وابستگی متغیر پاسخ y به متغیرهای تبیینی X ، y_i ها نیز وابسته به یکدیگر در نظر گرفته می‌شوند. توزیع y_i به شرط سایر مقادیر بردار y نیز نرمال است و به صورت $y_j - \mu_j$ ، یعنی اختلاف بین مقدار واقعی y_j و مقدار مورد انتظار آن، بیان شده است. برای c_{ij} ها که مقادیر ثابت معلوم یا نامعلوم هستند، محدودیت‌هایی قرار دارند. از جمله $c_{ij} = c_{ji}$ و $c_{ii} = 0$. محدودیت اول به سادگی به این معناست که، توزیع شرطی y_i نمی‌تواند به خودش وابسته باشد و محدودیت دوم به لزوم مقارن بودن وزن‌ها تاکید می‌کند.

مدل (۳.۲) را می‌توان به شکل ماتریسی

$$y \sim N(\mu, (I - C)\tau^2) \quad (4.2)$$

نوشت، که در آن ماتریس مقادیر C برای بازتاب ساختار فضایی داده‌ها به صورت‌های مختلفی تعریف می‌شود. یک راه معمول ساخت C معرفی آن توسط یک پارامتر وابستگی فضایی (ρ_c) و ماتریس وزن W است، که ساختار همسایگی‌های فضایی را وارد مدل می‌کند. تعاریف مختلفی برای ماتریس وزن W ارائه شده‌اند؛ تعدادی از آن‌ها را در فصل یک معرفی کردیم. برای مطالعه بیشتر در خصوص نحوه انتساب وزن‌ها در این مدل، به کلایتون و برناردینی (۱۹۹۲) مراجعه کنید.

از مدل CAR به‌طور گسترده برای مدل‌بندی داده‌های شبکه‌ای استفاده شده است. به‌عنوان مثال، در متون پزشکی برای تصویرسازی جغرافیایی بیماری‌ها توسط کلایتون و کالدر (۱۹۸۷) و مولی و ریچاردسان (۱۹۹۱) مورد استفاده قرار گرفته است. همچنین در متون اقتصادی می‌توان از انسلین و فلوراکس (۱۹۹۵) و کلیجان و پروچا (۱۹۹۹) نام برد. همچنین فودرینگ‌هام و همکاران (۲۰۰۳) با توجه به این که علاوه بر ویژگی‌های یک مسکن (مانند زیربنا و سن ساختمان)، ارزش بناهای مجاور نیز در تعیین قیمت مسکن موثر است، از مدل CAR برای تحلیل قیمت مسکن استفاده کردند.

مدل اتورگرسیو هم‌زمان

یک مدل اتورگرسیو هم‌زمان به صورت زیر تعریف می‌شود:

$$y_i \sim N(\mu_i + \sum_{j=1}^n b_{ij}(y_j - \mu_j), \lambda^2) \quad (5.2)$$

در این مدل بر خلاف مدل CAR توزیع y_i ها شرطی نیست؛ بلکه برای هر y_i به‌عنوان یک مجموعه از معادلات هم‌زمان تعریف شده است. به این معنا که مقدار واقعی y_i با مقدار مورد انتظار y_j همبستگی دارد. محدودیت $b_{ii} = 0$ برای مدل SAR نیز برقرار است اما شرط تقارن وجود ندارد. معادله (۵.۲) را می‌توان به صورت زیر بازنویسی کرد:

$$y \sim N(\mu, \lambda^2(I - B)^{-1}(I - B^T)^{-1}). \quad (6.2)$$

در این مدل نیز همانند مدل CAR ماتریس B اغلب به صورت $B = \rho_s W$ بیان می‌شود و روش انتخاب وزن‌ها نیز همانند CAR است.

۳.۲ مدل‌های اقتصادسنجی

یکی از تحولات و پیشرفت‌های ایجادشده در به‌کارگیری روش‌های کمی در علوم رفتاری به ویژه اقتصاد، تکامل شاخه اقتصادسنجی به اقتصادسنجی فضایی است. تفاوت اقتصادسنجی فضایی با اقتصادسنجی مرسوم، در توانایی و کاربرد آن در مواجهه با داده‌هایی است که دارای اطلاعات مکانی هستند. اقتصادسنجی فضایی، اثرات فضایی را به مدل‌های رگرسیونی اضافه می‌کند. با توجه به این که اقتصادسنجی فضایی به تحلیل داده‌های فضایی می‌پردازد، نخستین کاربردهای آن در علم اقتصاد در علوم منطقه‌ای، اقتصاد شهری و املاک و مستغلات (از جمله بازار مسکن) و جغرافیای اقتصادی بروز می‌یابد، اما به تدریج در شاخه‌های دیگر علم اقتصاد نظیر مطالعات تقاضا، اقتصاد بین‌الملل، اقتصاد نیروی کار، اقتصاد بخش عمومی و اقتصاد کشاورزی و محیط زیستی نیز از روش اقتصادسنجی فضایی استفاده شده است. این شاخه از اقتصادسنجی توسط پالینک و کلاسن (۱۹۷۹) و انسلین (۱۹۸۰) گسترش پیدا کرده است. کارهای متعددی توسط محققین مختلف درباره روش‌های مختلف برآورد و آزمون‌های تخصیص مدل‌های فضایی صورت گرفته‌اند. از جمله می‌توان به انسلین (۱۹۸۸a)، هاینینگ (۱۹۹۰)، کرسی (۱۹۹۳) و لی‌سیچ و پیس (۲۰۰۹) اشاره کرد. همان‌طور که گفته شد، مدل‌های بسیاری در ادبیات اقتصادسنجی فضایی مطرح شده‌اند. اما علی‌رغم این گستردگی تنها تعدادی از این مدل‌ها در مطالعات تجربی مورد استفاده قرار می‌گیرند. دلیل این امر وجود چالش‌های مدل‌بندی و محاسباتی در تعمیم و برازش اقتصادسنجی سایر مدل‌های فضایی است. در ادامه به معرفی این مدل‌ها و نحوه برآورد پارامترهای آن‌ها می‌پردازیم.

۱.۳.۲ مدل خطای فضایی

برای معرفی این مدل از مدل عمومی رگرسیون فضایی (۲.۲) شروع می‌کنیم. برای مشاهده i ام فرض می‌کنیم

$$y_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + v_i, \quad i = 1, \dots, n$$

که در آن جمله خطا برای ناحیه i ام، یعنی v_i ، خودهمبسته فضایی و دارای توزیع نرمال است. به‌طور دقیق‌تر

$$v_i = \lambda \sum_j w_{ij} v_j + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2), \quad i = 1, \dots, n.$$

همان‌طور که گفته شد، λ پارامتر وابستگی فضایی و w_{ij} ‌ها وزن‌های فضایی با شرط $w_{ii} = 0$ هستند. در حقیقت در این مدل (همان‌طور که از اسم آن مشخص است) فرض شده است

که کل ساختار وابستگی‌های فضایی در جمله خطای مدل وجود دارد. مدل خطای فضایی را می‌توان به صورت ماتریسی زیر نوشت:

$$y = X\beta + v, \quad v = \lambda Wv + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n). \quad (7.2)$$

مدل (۷.۲) برای برآورد ضرایب، باید یکپارچه شده و به صورت یک مدل رگرسیون خطی تعمیم یافته نوشته شود. می‌دانیم

$$\begin{aligned} v &= \lambda Wv + \varepsilon \\ &= (I_n - \lambda W)^{-1} \varepsilon = B_\lambda^{-1} \varepsilon. \end{aligned} \quad (8.2)$$

همچنین

$$\begin{aligned} \text{Cov}(v) &= \text{Cov}(B_\lambda^{-1} \varepsilon) = B_\lambda^{-1} \text{Cov}(\varepsilon) (B_\lambda^{-1})^T \\ &= B_\lambda^{-1} (\sigma^2 I_n) (B_\lambda^{-1})^T = \sigma^2 B_\lambda^{-1} (B_\lambda^T)^{-1} \\ &= \sigma^2 (B_\lambda^T B_\lambda)^{-1} = \sigma^2 V_\lambda. \end{aligned} \quad (9.2)$$

بنابراین مدل (۷.۲) را می‌توان به شکل زیر بازنویسی کرد:

$$y = X\beta + v, \quad v \sim N(0, \sigma^2 V_\lambda). \quad (10.2)$$

در صورت رد فرضیه صفر در آزمون موران (۱۱.۱)، نسبت درست‌نمایی (۱۲.۱) یا ضریب لاگرانژ خطا (۱۳.۱) از این مدل استفاده می‌شود. در ادامه دو روش برآورد پارامترهای این مدل فضایی را بیان خواهیم کرد.

روش GLS

با توجه به ساختار غیرقطری ماتریس واریانس جمله خطا در معادله (۱۰.۲) نمی‌توان از روش OLS برای برآورد ضرایب استفاده کرد. در مقابل می‌توان از روش GLS استفاده کرد. همان‌طور که در فصل ۱ بیان شد، این روش در صورت وجود ماتریس متقارن K به طوری که

$$KK = K^T K = V = \text{Var}(v) / \sigma^2$$

با اعمال تغییراتی در داده‌ها شرایط برآورد OLS را فراهم می‌آورد. با توجه به (۹.۲) در مدل SEM این شرط برقرار است، زیرا

$$\text{Var}(v) = \sigma^2 V = \sigma^2 (B^{-1})^T B^{-1}.$$

با تغییر متغیرهای

$$\tilde{Y} = By, \quad \tilde{X} = BX, \quad \tilde{v} = Bv$$

خواهیم داشت

$$E(\tilde{v}) = 0$$

و

$$\text{Var}(\tilde{v}) = \text{Var}(\mathbf{B}v) = \mathbf{B}\text{Var}(v)\mathbf{B}^T = \mathbf{B}\sigma^2(\mathbf{B}^{-1})^T\mathbf{B}^{-1}\mathbf{B}^T = \sigma^2\mathbf{B}\mathbf{B}^{-1}(\mathbf{B}^T)^{-1}\mathbf{B}^T = \sigma^2\mathbf{I}.$$

بنابراین

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

برآوردگر OLS متغیرهای جدید است. در نتیجه، برآوردگر GLS پارامترهای مدل SEM به صورت زیر خواهد بود:

$$(11.2)$$

$$\hat{\beta} = ((\mathbf{B}\mathbf{X})^T(\mathbf{B}\mathbf{X}))^{-1}(\mathbf{B}\mathbf{X})^T(\mathbf{B}\mathbf{y}) = (\mathbf{X}^T\mathbf{B}^T\mathbf{B}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{B}^T\mathbf{B}\mathbf{y} = (\mathbf{X}^T\mathbf{V}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}\mathbf{y}.$$

روش درست‌نمایی ماکسیمم

از آن جا که بردار تصادفی \mathbf{y} ترکیبی خطی از بردار تصادفی v می باشد (رابطه (10.2) را ببینید)، بنابراین \mathbf{y} نیز دارای توزیع نرمال n متغیره به شکل زیر است:

$$\mathbf{y} \sim \text{MVN}(\mathbf{X}\beta, \sigma^2\mathbf{V}_\lambda).$$

بنابراین تابع درست‌نمایی برای مشاهداتی از این بردار تصادفی به صورت

$$L(\theta) = (2\pi)^{-\frac{n}{2}}|\sigma^2\mathbf{V}_\lambda|^{-\frac{1}{2}} \times \exp\left\{-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\beta)^T\mathbf{V}_\lambda^{-1}(\mathbf{y} - \mathbf{X}\beta)\right\}$$

است، که در آن $\theta = [\lambda, \sigma^2, \beta^T]$. از ویژگی‌های دترمینان و هم‌چنین تقارن B نتیجه می شود

$$|\sigma^2\mathbf{V}_\lambda| = |\sigma^2(\mathbf{B}_\lambda^T\mathbf{B}_\lambda)^{-1}| = (\sigma^2)^n|\mathbf{B}_\lambda|^{-2}.$$

هم‌چنین از (8.2) و (10.2) نتیجه می شود

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\beta + \mathbf{B}_\lambda^{-1}\varepsilon \\ \varepsilon &= \mathbf{B}_\lambda(\mathbf{y} - \mathbf{X}\beta) = \mathbf{B}_\lambda\mathbf{y} - \mathbf{B}_\lambda\mathbf{X}\beta. \end{aligned} \quad (12.2)$$

بنابراین تابع لگاریتم درست‌نمایی مدل به شکل زیر به دست می آید:

$$\ln L = -\frac{n}{2}\ln(2\pi) - \frac{n}{2}\ln(\sigma^2) + \ln|\mathbf{B}_\lambda| - \frac{1}{2\sigma^2}\varepsilon^T\varepsilon. \quad (13.2)$$

ماکسیمم کردن تابع (۱۳.۲) نسبت به β ، معادل مینیمم کردن جمله آخر آن نسبت به β است. یعنی برآوردگر OLS پارامتر β در مدل (۱۲.۲) همان برآوردگر درست‌نمایی ماکسیمم (ML) β است. یعنی

$$\hat{\beta}_{ML} = \left((\mathbf{B}_\lambda \mathbf{X})^T (\mathbf{B}_\lambda \mathbf{X}) \right)^{-1} (\mathbf{B}_\lambda \mathbf{X})^T (\mathbf{B}_\lambda \mathbf{y}). \quad (14.2)$$

همچنین برآوردگر درست‌نمایی ماکسیمم برای σ^2 با حل معادله زیر محاسبه می‌شود:

$$\begin{aligned} \frac{\partial \ln L}{\partial \sigma^2} &= \frac{-n}{2\sigma^2} + \frac{2\mathbf{e}^T \mathbf{e}}{4\sigma^4} \\ &= \frac{-n\sigma^2 + \mathbf{e}^T \mathbf{e}}{2\sigma^4} = 0 \end{aligned} \quad (15.2)$$

که در آن $\mathbf{e} = \mathbf{B}_\lambda \mathbf{y} - \mathbf{B}_\lambda \mathbf{X} \hat{\beta}_{ML}$. نتیجه حل معادله (۱۵.۲) به برآوردگر زیر منتهی می‌شود:

$$\hat{\sigma}_{ML}^2 = \frac{\mathbf{e}^T \mathbf{e}}{n}.$$

۲.۳.۲ مدل تاخیر فضایی

یک مدل خطی دیگر بر پایه خودهمبستگی فضایی، مدل تاخیر فضایی است. در این مدل فرض شده است که روابط خودهمبستگی در متغیرهای وابسته وجود دارند. اگر مشابه مدل قبل، فرض کنیم که این روابط فضایی از طریق ماتریس وزن $\mathbf{W} = (w_{ij} : i, j = 1, \dots, n)$ با شرط $w_{ii} = 0$ لحاظ می‌شوند، آن‌گاه یک راه ساده برای نوشتن مدل به شکل زیر است:

$$y_i = \beta_0 + \rho \sum_{h \neq i} w_{ih} y_h + \sum_{j=1}^p \beta_j x_{ij} + \varepsilon_i, \quad i = 1, \dots, n \quad (16.2)$$

که در آن $\varepsilon_i \sim N(0, \sigma^2)$ و جمله خودهمبستگی، $\rho \sum_{h \neq i} w_{ih} y_h$ ، وابستگی ممکن y_i با مقادیر سایر نقاط، y_h ، را نشان می‌دهد. به‌عنوان یک مثال ساده برای این مدل می‌توان از قیمت مسکن در یک شهر نام برد. اگر y_i متوسط قیمت مسکن با ویژگی‌های $(x_{ij} : j = 1, \dots, p)$ در ناحیه i باشد، آن‌گاه این قیمت تحت تاثیر قیمت مسکن‌های نواحی مجاور است. بنابراین خودهمبستگی فضایی در متغیرهای وابسته وجود دارد نه در جمله خطا.

مدل (۱۶.۲) را می‌توان به‌صورت ماتریسی زیر نوشت:

$$\mathbf{y} = \rho \mathbf{W} \mathbf{y} + \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_n). \quad (17.2)$$

معادله (۱۷.۲) را می‌توان به‌صورت زیر خلاصه کرد:

$$\begin{aligned} \mathbf{y} - \rho \mathbf{W} \mathbf{y} &= \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon} \\ (\mathbf{I}_n - \rho \mathbf{W}) \mathbf{y} &= \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon} \\ \mathbf{B}_\rho \mathbf{y} &= \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon} \\ \mathbf{y} &= \mathbf{B}_\rho^{-1} \mathbf{X} \boldsymbol{\beta} + \mathbf{B}_\rho^{-1} \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_n). \end{aligned} \quad (18.2)$$

با تغییر متغیرهای $\mathbf{X}_\rho = \mathbf{B}_\rho^{-1} \mathbf{X}$ و $\mathbf{v} = \mathbf{B}_\rho^{-1} \boldsymbol{\varepsilon}$ ، مدل (۱۸.۲) به شکل زیر نوشته می‌شود:

$$\mathbf{y} = \mathbf{X}_\rho \boldsymbol{\beta} + \mathbf{v}, \quad \mathbf{v} \sim N(\mathbf{0}, \sigma^2 \mathbf{V}_\rho) \quad (19.2)$$

که در آن مشابه مدل SEM، $\mathbf{V}_\rho = (\mathbf{B}_\rho^T \mathbf{B}_\rho)^{-1}$ ، در صورت رد فرضیه صفر در آزمون ضریب لاگرانژ تأخیر فضایی (۱۸.۱)، از این مدل برای مدل بندی ساختار فضایی استفاده می‌کنیم.

برآورد GLS

به منظور برآورد پارامترهای این مدل، معادله (۱۹.۲) را در نظر می‌گیریم. تمامی مراحل عیناً همانند مدل خطای فضایی دنبال می‌شوند و در نهایت برآورد GLS مدل تأخیر فضایی به صورت زیر است:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{B}_\rho \mathbf{y}. \quad (20.2)$$

برآورد ML

مشابه مدل SEM، تابع درست‌نمایی ماکسیمم بردار مشاهدات \mathbf{y} به صورت زیر تعریف می‌شود:

$$L(\theta) = (\gamma \pi)^{-\frac{n}{\gamma}} |\sigma^2 \mathbf{V}_\rho|^{-\frac{1}{\gamma}} \times \exp \left\{ -\frac{1}{\gamma \sigma^2} (\mathbf{y} - \mathbf{X}_\rho \boldsymbol{\beta})^T \mathbf{V}_\rho^{-1} (\mathbf{y} - \mathbf{X}_\rho \boldsymbol{\beta}) \right\}$$

که در این مدل $\theta = [\rho, \sigma^2, \boldsymbol{\beta}^T]$ از (۱۸.۲) نتیجه می‌شود

$$\boldsymbol{\varepsilon} = \mathbf{B}_\rho (\mathbf{y} - \mathbf{X}_\rho \boldsymbol{\beta}) = \mathbf{B}_\rho \mathbf{y} - \mathbf{B}_\rho \mathbf{B}_\rho^{-1} \mathbf{X} \boldsymbol{\beta} = \mathbf{B}_\rho \mathbf{y} - \mathbf{X} \boldsymbol{\beta}. \quad (21.2)$$

بنابراین تابع لگاریتم درست‌نمایی مدل به صورت زیر نتیجه می‌شود:

$$\ln L = -\frac{n}{\gamma} \ln(\gamma \pi) - \frac{n}{\gamma} \ln(\sigma^2) + \ln |\mathbf{B}_\rho| - \frac{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}}{\gamma \sigma^2}. \quad (22.2)$$

مشابه قبل، ماکسیمم کردن تابع (۲۲.۲) نسبت به $\boldsymbol{\beta}$ مستلزم مینیمم کردن جمله آخر آن نسبت به $\boldsymbol{\beta}$ است، یعنی طبق (۲۱.۲) برآوردگر OLS رگرسیون $\mathbf{B}_\rho \mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ، برآوردگر درست‌نمایی ماکسیمم $\boldsymbol{\beta}$ است:

$$\hat{\boldsymbol{\beta}}_{ML} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{B}_\rho \mathbf{y}. \quad (23.2)$$

همچنین برآوردگر درست‌نمایی ماکسیمم σ^2 نیز مشابه مدل SEM با حل معادله (۱۵.۲)، به صورت زیر است:

$$\hat{\sigma}_{ML}^2 = \frac{\mathbf{e}^T \mathbf{e}}{n}$$

با این تفاوت که $\mathbf{e} = \mathbf{B}_\rho \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{ML}$

برای محاسبه پارامترهای λ یا ρ ، برآورد درست‌نمایی ماکسیمم پارامترهای β و σ^2 را در تابع لگاریتم درست‌نمایی قرار داده و تابع لگاریتم درست‌نمایی متمرکز حاصل به شکل تابع غیرخطی از λ یا ρ را با استفاده از روش‌های عددی ماکسیمم می‌کنیم. با اعمال تغییراتی روی دو مدل معرفی شده SEM و SLM، مدل‌های متنوعی تولید شده‌اند. در ادامه دو مدل پرکاربرد که در این پایان‌نامه نیز مورد استفاده قرار گرفته‌اند، را معرفی می‌کنیم.

۳.۳.۲ مدل ترکیبی فضایی

هنگام طراحی مدل تأخیر فضایی این سوال مطرح می‌شود که چرا تمام اثرات ناشناخته، ε_i ، باید مستقل از فضا در نظر گرفته شوند. زیرا ممکن است که خودهمبستگی فضایی در جمله خطا و متغیر وابسته هر دو با هم وجود داشته باشد یا به عبارت دیگر شرایط SEM و SLM هم‌زمان برقرار باشند. بنابراین می‌توان ترکیب این دو مدل را به صورت ماتریسی زیر نوشت:

$$y = \rho W y + X\beta + v, \quad v = \lambda M v + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n) \quad (24.2)$$

که در آن M نیز مشابه W ماتریس وزن‌های فضایی و λ مشابه ρ پارامتر وابستگی فضایی است. این مدل توسط کلجیان و پروچا (۲۰۱۰) معرفی شد و در متون اقتصادسنجی با نام مدل ترکیبی فضایی^۳ (SAC) یا مدل اتورگرسیو فضایی با خطاهای اتورگرسیو^۴ (SARAR) شناخته می‌شود. در این مدل، اولویت اصلی ایجاد یک مدل نسبی از ترکیب دو مدل تأخیر فضایی و خطای فضایی به‌عنوان دو مدل با ساختار مشابه است. از این رو تمرکز بر حالت خاص $W = M$ قرار می‌گیرد و مدل (۲۴.۲) به صورت

$$y = \rho W y + X\beta + v, \quad v = \lambda W v + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n)$$

نوشته می‌شود و به شکل زیر یک‌پارچه می‌شود:

$$y = (B_\rho)^{-1} X\beta + (B_\rho)^{-1} (B_\lambda)^{-1} \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n). \quad (25.2)$$

در حالت خاص اگر $\rho = 0$ ($\lambda = 0$)، مدل (۲۵.۲) معادل SEM (SLM) می‌شود. این مدل در شرایطی که هر دو آزمون ضریب لاگرانژ خطا و آزمون ضریب لاگرانژ تأخیر فضایی معنی‌دار شوند، مورد استفاده قرار می‌گیرد.

استنباط مدل

برای برآورد ضرایب این مدل رگرسیونی از معادله (۲۵.۲) استفاده می‌کنیم و با تغییر متغیر

$$X_\rho = B_\rho^{-1} X$$

³Spatial autoregressive combined model

⁴Spatial autoregressive model with autoregressive disturbances

آن را به صورت

$$y = X_{\rho}\beta + (B_{\lambda}B_{\rho})^{-1}\varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n) \quad (26.2)$$

که یک مدل خطای فضایی با جمله خطا $(B_{\lambda}B_{\rho})^{-1}V_{\lambda}(B_{\rho}^T)^{-1}u \sim N(0, \sigma^2 B_{\rho}^{-1}V_{\lambda}(B_{\rho}^T)^{-1})$ است، بازنویسی می کنیم . بنابراین برآورد ضرایب رگرسیونی آن طبق (۱۱.۲) و (۱۴.۲) به صورت زیر محاسبه می شود:

$$\begin{aligned} \hat{\beta} &= ((B_{\lambda}B_{\rho}X_{\rho})^T(B_{\lambda}B_{\rho}X_{\rho}))^{-1}(B_{\lambda}B_{\rho}X_{\rho})^T(B_{\lambda}B_{\rho}y) \\ &= ((B_{\lambda}X)^T(B_{\lambda}X))^{-1}(B_{\lambda}X)^T(B_{\lambda}B_{\rho}y) \end{aligned} \quad (27.2)$$

همچنین مطابق با مدل SEM برآورد σ^2 به صورت زیر است:

$$\hat{\sigma}^2 = \frac{e^T e}{n}$$

که در آن $e = B_{\lambda}B_{\rho}y - B_{\lambda}X\hat{\beta}$ محاسبه می شود.

۴.۳.۲ مدل دوربین فضایی

مفید و کاربردی ترین مدل اتورگرسیو مدل دوربین فضایی^۵ (SDM) است که توسط انسلین (۱۹۸۸a) معرفی شده است. این مدل از ایجاد تغییر در مدل تاخیر فضایی حاصل می شود و در مقایسه با آن از عملکرد بهتری برخوردار است. برای معرفی این مدل از مثال تعیین قیمت مسکن مورد استفاده در این پایان نامه شروع می کنیم. در این مورد اگر قیمت مسکن در ناحیه y_i ، i تحت تاثیر قیمت مسکن های همسایه باشد، آن گاه غیرمنطقی نیست که تحت تاثیر ویژگی های آن ها نیز باشد. بنابراین، این اثرات فضایی به عنوان شرایط اضافی به مدل افزوده می شود:

$$y_i = \beta_0 + \rho \sum_{h \neq i} w_{ih} y_h + \sum_{j=1}^p \beta_j x_{ij} + \sum_{h=1}^n m_{ih} \left(\sum_{j=1}^p \alpha_j x_{hj} \right) + \varepsilon_i, \quad i = 1, \dots, n. \quad (28.2)$$

برای نوشتن صورت ماتریسی این مدل، لازم است تا ماتریس X و بردار β را به شکل زیر افراز کنیم:

$$X = \begin{bmatrix} 1_n & X_v \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_v \end{bmatrix}. \quad (29.2)$$

بنابراین می توان معادله (۲۸.۲) را به شکل

$$y = \rho W y + \beta_0 1_n + X_v \beta_v + M X_v \alpha + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n)$$

⁵Spatial durbin model

$$y = \rho W y + X \beta + M X_v \alpha + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n). \quad (30.2)$$

نوشت؛ که در آن α یک بردار ستونی از ضرایب رگرسیونی با بعد p است. در این مدل نیز می‌توان وزن‌های فضایی را یکسان در نظر گرفت. یعنی

$$y = \rho W y + X \beta + W X_v \alpha + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n). \quad (31.2)$$

استنباط مدل

برآورد ضرایب این مدل رگرسیونی کاملاً مشابه مدل تأخیر فضایی است و تنها با تغییر متغیر انجام می‌شود. به این منظور، ماتریس متغیرهای تبیینی را به صورت

$$Z = [X \quad W X_v]$$

که ماتریسی با بعد $n \times (2p + 1)$ است و بردار ضرایب رگرسیونی را به شکل

$$\delta = [\beta \quad \alpha]^T$$

که یک بردار با بعد $1 \times (2p + 1)$ است، تعریف می‌کنیم. با این تعاریف معادله (۳۱.۲) به شکل زیر بازنویسی می‌شود:

$$y = \rho W y + Z \delta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n). \quad (32.2)$$

برآورد ضرایب رگرسیونی مدل (۳۲.۲)، که یک مدل تأخیر فضایی است، مطابق (۲۳.۲) و (۲۰.۲) به شکل

$$\hat{\delta} = (Z^T Z)^{-1} Z^T B_{\rho} y \quad (33.2)$$

همچنین برآورد σ^2 ، مشابه SLM به صورت

$$\hat{\sigma}^2 = \frac{e^T e}{n}$$

که در آن $e = B_{\rho} y - Z \delta$ است، محاسبه می‌شوند.

تمامی مدل‌هایی که تا کنون بیان شدند، مدل‌های فراموضعی^۶ فضایی هستند.

تعریف ۱.۳.۲. مدل‌های فراموضعی: مدل‌هایی هستند که ضرایب برآورد شده در مدل برای پیشگویی متغیر پاسخ \hat{y} به ازای هر مقدار x_0 در کل ناحیه تحت مطالعه ثابت است.

^۶Global model

در عمل و به خصوص در مورد داده‌های فضایی ممکن است متغیرهای تبیینی در نواحی مختلف دارای اثرات متفاوت باشند؛ به عبارتی داده‌ها ناهمگن فضایی باشند. در مثال تعیین قیمت مسکن ذکرشده، این حقیقت وجود دارد که با وجود وابستگی قیمت مسکن ناحیه i به نواحی مجاور، اما متغیرهای موثر بر قیمت آن ممکن است دارای اثرات متفاوت از نواحی مجاور باشند. برای مثال، تعداد پارکینگ‌های یک منزل در تعیین قیمت آن موثر هستند. اما این متغیر در نواحی جنوبی یک شهر که افراد با درآمد کمتر در آن سکونت دارند و اغلب فاقد وسیله نقلیه هستند، از اهمیت کمتری نسبت به نواحی دیگر برخوردار است. لذا لازم است که برای عملکرد بهتر، این ناهمگنی در نظر گرفته شود.

تعریف ۲.۳.۲. مدل‌های موضعی^۷: مدل‌هایی هستند که ضرایب برآوردشده در مدل برای پیشگویی متغیر پاسخ به ازای هر مقدار x_0 ثابت نبوده و وابسته به آن هستند. یعنی در مدل‌های موضعی، برای برآورد متغیر پاسخ در نقطه x_0 ابتدا یک همسایگی از آن نقطه انتخاب شده سپس برازش مدل فقط بر روی مشاهدات داخل همسایگی صورت می‌گیرد نه بر روی کل مشاهدات. بنابراین واضح است که با تغییر x_0 همسایگی نیز تغییر نموده و چون مشاهدات داخل همسایگی نیز تغییر می‌کنند، مدل برازش‌شده فعلی با مدل قبلی متفاوت خواهد بود.

یکی از کاربردی‌ترین مدل‌های موضعی، رگرسیون موزون جغرافیایی است که در ادامه به اختصار آن را معرفی می‌کنیم.

۴.۲ رگرسیون موزون جغرافیایی

در بخش‌های پیشین بیان کردیم که مدل‌های رگرسیونی فراموضعی فضایی در مواجهه با ناهمگنی فضایی مناسب نیستند. ناهمگنی فضایی بیانگر این حقیقت است که در هر منطقه رابطه‌ای متفاوت بین متغیر پاسخ و متغیرهای تبیینی وجود دارد. از این رو اگر این ناهمگنی در نظر گرفته نشود و رابطه‌ای یکسان برای تمامی مناطق برآورد شود، تخمینی نادقیق از روابط منطقه‌ای نتیجه خواهد شد. از میان روش‌های ناپارامتری، رگرسیون موزون موضعی^۸ (LWR) (سلولند و دولین، ۱۹۸۸) یا رگرسیون موزون جغرافیایی که توسط فودرینگ‌هام و همکاران (۲۰۰۳) به‌طور کامل معرفی شد؛ این امکان را فراهم می‌آورد که پارامترها را (به جای فراموضعی) ناحیه‌ای برآورد کنیم. ایده اساسی برآورد موضعی این است که توابع خطی ساده ممکن است برای مشاهدات نزدیک به یک نقطه مناسب باشند، اما احتمالاً زمانی که مشاهدات دوردست بیشتر گنجانده شوند، نامناسب خواهند بود. بنابراین محدود کردن برآورد به یک ناحیه از مشاهدات، مشکل واریانس نامتجانس و خودهمبستگی را که در داده‌های فضایی وجود دارند، رفع می‌کند.

⁷Local models

⁸Locally weighted regression

در روش رگرسیون موزون موضعی پارامترهای مدل را می‌توان در هر نقطه‌ای از فضای مورد مطالعه برآورد کرد. یعنی برای هر نقطه در فضا یک معادله رگرسیونی جداگانه تعریف می‌شود. در این روش برای برآورد پارامترهای مدل در هر نقطه از مشاهدات اطراف آن نقطه استفاده می‌شود، اما به مشاهدات نزدیک وزن بیشتر و به مشاهدات دور وزن کمتری داده می‌شود. این امر نیز با ساختار داده‌های فضایی، که در موقعیت‌های نزدیک به هم دارای همبستگی بیشتر هستند، متناسب است. بنابراین، تابع خطی ساده را می‌توان برای برآورد موضعی پارامترها به روش زیر نوشت:

$$\mathbf{y} = (\boldsymbol{\beta} \odot \mathbf{X})\mathbf{1} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim MVN(\mathbf{0}, \sigma_{\varepsilon}^2 \mathbf{I}) \quad (34.2)$$

که در آن \mathbf{y} بردار ستونی با بعد n و $\boldsymbol{\beta}$ یک ماتریس $n \times (p+1)$ با سطر i ام

$$\boldsymbol{\beta}_i = (\beta_{i0}, \beta_{i1}, \dots, \beta_{ip})^T$$

است. عملگر \odot ضرب هادامارد و ماتریس \mathbf{X} دارای بعد یکسان با $\boldsymbol{\beta}$ و سطر i ام

$$\mathbf{x}_i = (x_{i0}, x_{i1}, \dots, x_{ip})$$

است. بردار ستونی $\boldsymbol{\varepsilon}$ ، n بعدی و $\mathbf{1}$ یک بردار ستونی $(p+1)$ بعدی از آن‌ها است. مدل (34.2) را می‌توان برای هر مشاهده، به صورت

$$y_i = \beta_{i0} + \beta_{i1}x_{i1} + \dots + \beta_{ip}x_{ip} + \varepsilon_i \quad \varepsilon_i \sim N(0, \sigma^2) \quad i = 1, 2, \dots, n \quad (35.2)$$

نوشت.

با این که در صورت کلی مدل GWR واریانس جمله خطا ثابت فرض شده است ($\text{Var}(\varepsilon) = \sigma^2$) و ساختار آن عملاً خودهمبستگی و ناهمگنی را نشان نمی‌دهد، اما باید توجه داشت که اکثر این اثرات با برآورد ضرایب مختلف برای هر مشاهده با استفاده از همسایه‌های آن، از بین رفته است.

۱.۴.۲ استنباط مدل

به منظور برآورد یک مجموعه از پارامترهای نقطه i ، می‌توان تابع هدف زیر را به روش کمترین توان‌های دوم، در زیرمجموعه‌ای از نقاط که نزدیک i هستند، کمینه کرد:

$$Q(\boldsymbol{\beta}_i) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_i)^T \mathbf{W}_i (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_i) = \boldsymbol{\varepsilon}^T \mathbf{W}_i \boldsymbol{\varepsilon} \quad i = 1, 2, \dots, n \quad (36.2)$$

که در آن W_i یک ماتریس قطری از بعد n است، که درایه‌های غیرقطری آن صفر و درایه‌های روی قطر اصلی یعنی $(w_{i1}, w_{i2}, \dots, w_{in})$ وزن هریک از n مشاهده برای نقطه i است. یعنی

$$W_i = \begin{bmatrix} w_{i1} & \circ & \dots & \circ \\ \circ & w_{i2} & & \circ \\ \vdots & & \ddots & \vdots \\ \circ & \circ & \dots & w_{in} \end{bmatrix}.$$

این وزن‌ها در مدل GWR فضایی، متناسب با فواصل جغرافیایی بین مشاهدات هستند. به این ترتیب برآوردگر کمترین توان‌های دوم موزون موضعی^۹ (LWLS)، با تکرار WLS برای هر مشاهده i به شکل زیر به دست می‌آید:

$$\hat{\beta}_i^{WLS} = (\mathbf{X}^T \mathbf{W}_i \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_i \mathbf{y}, \quad i = 1, \dots, n. \quad (37.2)$$

بنابراین ما n ماتریس قطری از وزن‌های فضایی و n مجموعه از برآوردگرهای پارامترهای موضعی داریم. می‌توان برآوردگر موزون موضعی را به صورت

$$\hat{\beta}_i^{WLS} = \mathbf{C} \mathbf{y} \quad i = 1, 2, \dots, n$$

نوشت که در آن

$$\mathbf{C} = (\mathbf{X}^T \mathbf{W}_i \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_i.$$

امید و واریانس β_i برابر هستند با

$$\begin{aligned} E(\hat{\beta}_i) &= E(\mathbf{C} \mathbf{y}) = (\mathbf{X}^T \mathbf{W}_i \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_i E(\mathbf{y}) \\ &= (\mathbf{X}^T \mathbf{W}_i \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_i \mathbf{X} \beta_i = \beta_i \end{aligned} \quad (38.2)$$

$$\text{Var}(\hat{\beta}_i) = \text{Var}(\mathbf{C} \mathbf{y}) = \mathbf{C} \text{Var}(\mathbf{y}) \mathbf{C}^T = \sigma^2 \mathbf{C} \mathbf{C}^T \quad (39.2)$$

که در آن σ^2 واریانس جمله خطا در رگرسیون موضعی است. علاوه بر ارزیابی برآوردهای پارامترهای موضعی، بررسی خطاهای استاندارد موضعی نیز حائز اهمیت است. این کار به منظور محاسبه تغییرات داده‌های مورد استفاده برای تخمین برآوردها باید انجام شود. برای مثال، در بعضی موارد ممکن است که برآورد پارامترهای موضعی تابعی از تعداد مشاهدات نسبتاً کم باشد؛ یا برخی مشاهدات به دلیل دور بودن از نقطه رگرسیونی در برآورد پارامترهای موضعی دارای وزن کوچک باشند. بنابراین محاسبه تغییرات خطای استاندارد موضعی لازم است.

⁹Locally weighted least squares

همچنین طبق (۳۷.۲) برای برآورد ضرایب رگرسیونی مشاهده i از ماتریس وزن W_i استفاده می‌شود. با این کار مشاهدات همسایه i که بیشترین وابستگی را با آن دارند، در برآورد ضرایب رگرسیونی آن نیز بیشترین تاثیر را خواهند گذاشت. لذا واضح است که با تغییر نقاط رگرسیونی، $i = 1, \dots, n$ ، ماتریس‌های وزن و تحت تاثیر آن ضرایب رگرسیونی موضعی β_i و باقی‌مانده‌های استاندارد شده رگرسیون موضعی تغییر خواهند کرد. یعنی تغییرات باقی‌مانده‌های استاندارد شده با تغییر نقاط رگرسیونی متفاوت می‌شود. لذا برآورد واریانس جمله خطا در رگرسیون موضعی برای هر مشاهده $i = 1, \dots, n$ به صورت

$$\hat{\sigma}^2 = \frac{\sum_i (y_i - \hat{y}_i)^2}{(n - 2 \operatorname{tr}(\mathbf{S}) + \operatorname{tr}(\mathbf{S}^T \mathbf{S}))} \quad (40.2)$$

محاسبه می‌شود؛ که در آن \mathbf{S} ماتریس طرح $\hat{\mathbf{y}}$ روی \mathbf{y}

$$\hat{\mathbf{y}} = \mathbf{S}\mathbf{y}$$

است و هر سطر آن به صورت

$$\mathbf{s}_i = \mathbf{X}_i (\mathbf{X}^T \mathbf{W}_i \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_i \quad (41.2)$$

تعریف می‌شود. در تمام مطالب گفته شده در مورد رگرسیون موزون جغرافیایی، ابتدا لازم است که ماتریس وزن فضایی W_i مشخص شود.

۲.۴.۲ روش‌های انتساب وزن

توابع وزن دهی هسته باید در عبارت (۳۷.۲) تعریف شوند تا W_i مشخص گردد. این توابع عموماً نمایی، گاوسی، توان دو^{۱۰}، توان سه^{۱۱} هستند که مشابه یکدیگر در نقاط نزدیک وزن بزرگ‌تر و در نقاط دور وزن کوچک‌تر می‌گیرند؛ اما با توجه به نوع پهنای باند^{۱۲} متفاوتند. انتخاب پهنای باند بسیار مهم و حیاتی است. زیرا همسایگی مشاهده i توسط مقدار پهنای باند b مشخص می‌شود. این مقدار، تعداد مشاهداتی را که در برآورد برای نقطه i وزن می‌گیرند و نیز چگونگی کاهش وزن با افزایش فاصله را تعیین می‌کند. پهنای باند می‌تواند یک فاصله ثابت^{۱۳} یا سازوار^{۱۴} در نظر گرفته شود.

رگرسیون موزون فضایی با هسته فضایی ثابت

در روش رگرسیون موزون با هسته فضایی ثابت، دامنه‌ای در اطراف هر نقطه مرجع در نظر گرفته می‌شود و این فاصله به ازای هر نقطه ثابت است.

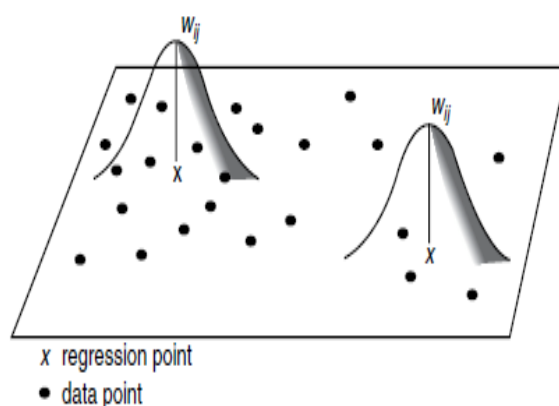
¹⁰Bi-square

¹¹Tri-cube

¹²Bandwidth

¹³Fixed

¹⁴Adaptive



شکل ۱.۲: رگرسیون موزون فضایی با هسته فضایی ثابت

مدل های رگرسیون موزون فضایی با هسته ثابت^{۱۵} (شکل ۱.۲ را ببینید)، زمانی که توزیع مشاهدات حول نقاط مرجع در گستره فضای مورد مطالعه یکسان نباشد، دچار مشکل می شود. در این صورت اگر تابعی یکسان برای موزون کردن مشاهدات اطراف نقاط مرجع به کار گرفته شود، این احتمال وجود دارد که برای نقطه ای از اطلاعاتی کمتر استفاده شود و بنابراین مدلی نادقیق برآورد شود و بالعکس برای نقطه ای دیگر که تجمع مشاهدات حول آن بیشتر است از اطلاعات زیاد و نامرتبط استفاده شود و مجدداً از دقت مدل کاسته شود. برای مقابله با این احتمال، رگرسیون موزون فضایی با هسته فضایی سازوار^{۱۶} معرفی شده است (فودرینگهام و همکاران، ۲۰۰۳).

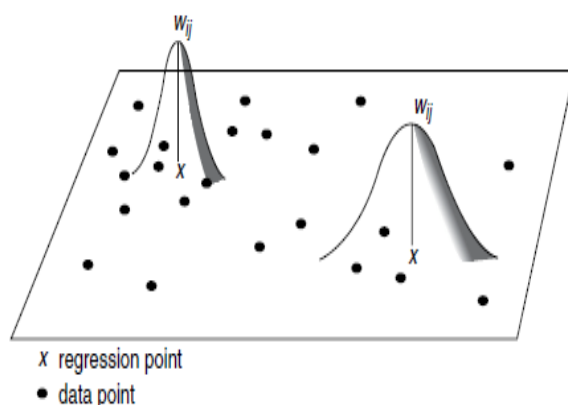
رگرسیون موزون فضایی با هسته فضایی سازوار

در روش رگرسیون موزون فضایی با هسته سازوار، پهنای باند تابع وزن دهی با پراکندگی مشاهدات حول نقطه مرجع تطبیق داده می شود. به طوری که زمانی که مشاهدات پراکنده هستند، پهنای باند بزرگ تری در نظر گرفته می شود و زمانی که مشاهدات حول نقطه مرجع متراکم باشند، پهنای باند کوچک تری انتخاب می شود (شکل ۲.۲ را ببینید). به عبارتی در این روش به جای فاصله ثابت از تعداد ثابت برای تعیین پهنای باند استفاده می شود. این روش تعیین همسایگی k -نزدیک ترین همسایگی^{۱۷} نیز نامیده می شود.

¹⁵Fixed kernel

¹⁶Adaptive kernel

¹⁷k-nearest neighbour



شکل ۲.۲: رگرسیون موزون فضایی با هسته فضایی سازوار

در ادامه نمونه‌هایی از توابع هسته ثابت و سازوار را معرفی می‌کنیم.

توابع هسته

ساده‌ترین روش، استفاده از یک طرح وزن‌دهی دوتایی به صورت زیر است (فودرینگ‌هام و همکاران، ۲۰۰۳):

$$w_{ij} = \begin{cases} 1 & \text{اگر } d_{ij} \leq b \\ 0 & \text{سایر جاها} \end{cases}$$

که در آن d_{ij} فاصله بین مشاهدات i و j و b پهنای باند است. این تابع هسته در مواردی که نقاط کالیبره^{۱۸} (نقاطی که در مرزهای پهنای باند قرار می‌گیرند) در یک منطقه با تراکم کم (تُنک) واقع شده است در مقایسه با یک منطقه متراکم، می‌تواند دست‌آورد بهتری داشته باشد.

یک تابع هسته دیگر به صورت

$$w_{ij} = \begin{cases} 1 & \text{اگر } y_j \in Y_i(N) \\ 0 & \text{سایر جاها} \end{cases}$$

معرفی می‌شود (ویلر و پائز، ۲۰۱۰). در این تابع هسته $Y_i(N)$ مجموعه‌ای از N نزدیکترین مشاهدات همسایه نقطه i است. N مقداری است که باید برآورد شود. در این روش بر خلاف روش قبلی، تابع هسته برای هر نقطه از تعداد یکسانی از مشاهدات استفاده می‌کند، اما این تعداد از مشاهدات ممکن است که برای هر نقطه i در مساحت‌های متفاوتی قرار گرفته باشند. با وجود این سادگی، این نوع تابع هسته به‌طور وسیعی استفاده نشده است.

¹⁸Calibration

در بیشتر موارد در مدل GWR استفاده از توابع هسته پیوسته ترجیح داده می‌شود، به این دلیل که این توابع وزن‌هایی تولید می‌کنند که به‌طور یکنواخت با افزایش فاصله کاهش می‌یابند. به‌عنوان نمونه تابع هسته گاوسی (فودرینگ‌هام و همکاران، ۲۰۰۳) را با ضابطه

$$w_{ij} = \exp \left\{ -\frac{1}{2} \left(\frac{d_{ij}}{b} \right)^2 \right\}$$

می‌توان معرفی کرد که در این تابع وزن مشاهده j در ارتباط با مشاهده i ، براساس فاصله این دو مشاهده از هم، d_{ij} ، تغییر می‌کند و b پهنای باند است که حد فاصله و نزول خودهمبستگی فضایی را کنترل می‌کند. همان‌طور که در بخش‌های قبلی گفته شد، پهنای باند می‌تواند به دو روش ثابت و سازوار مورد استفاده قرار گیرد. در این پایان‌نامه از تابع هسته گاوسی با پهنای باند سازوار استفاده می‌کنیم.

در این حالت تابع هسته گاوسی به شکل زیر معرفی می‌شود (ویلر و پائز، ۲۰۱۰):

(۴۲.۲)

$$w_{ij} = \begin{cases} \exp \left(-\frac{1}{2} \left(\frac{d_{ij}}{d_{iN}} \right)^2 \right) & \text{اگر } j \text{ یکی از } N \text{ نزدیک‌ترین همسایگی‌های } i \text{ باشد} \\ 0 & \text{سایر جاها} \end{cases}$$

که در آن d_{iN} فاصله نقطه‌ی i تا N امین (آخرین) نزدیک‌ترین همسایگی است. مقدار N به منظور محاسبه وزن‌ها برای برآورد ضرایب رگرسیونی موضعی باید تعیین شود. تابع هسته دیگر، تابع نمایی ساده به صورت

$$w_{ij} = \exp \left(-\frac{d_{ij}}{b} \right) \quad \text{اگر } d_{ij} < b$$

است که در آن مقیاس‌گذاری و توان نسبت به مشابه گاوسی خود، حذف شده است (ویلر و پائز، ۲۰۱۰).

همچنین تابع هسته توان‌دو به صورت

$$w_{ij} = \begin{cases} \left(1 - \frac{d_{ij}}{b} \right)^2 & \text{اگر } d_{ij} < b \\ 0 & \text{سایر جاها} \end{cases}$$

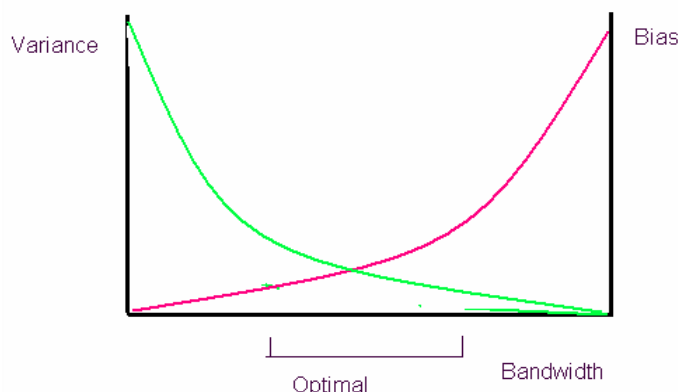
است (فودرینگ‌هام و همکاران، ۲۰۰۳). شکل سازوار این تابع هسته، مشابه تابع هسته گاوسی سازوار تعریف می‌شود (ویلر و پائز، ۲۰۱۰). یعنی

$$w_{ij} = \begin{cases} \left(1 - \left(\frac{d_{ij}}{d_{iN}} \right)^2 \right)^2 & \text{اگر } j \text{ یکی از } N \text{ نزدیک‌ترین همسایگی‌های } i \text{ باشد} \\ 0 & \text{سایر جاها} \end{cases}$$

همان‌طور که در تمامی توابع هسته گفته شد، مقدار پهنای باند در محاسبه وزن‌ها بسیار اثرگذار است. به‌همین منظور تعیین فاصله آن در حالت هسته ثابت و تعیین مقدار N در حالت هسته سازوار بسیار حیاتی است. در بخش بعد به معرفی چند معیار برای تعیین این مقادیر می‌پردازیم.

۳.۴.۲ انتخاب بهترین همسایگی

همان‌طور که گفته شد، انتخاب پهنای باند و در نتیجه آن انتخاب همسایگی مناسب، در برآورد ضرایب بسیار موثر است. بنابراین سوال مهم این است که چگونه می‌توان یک همسایگی مناسب برای i تعریف کرد. در انتخاب b همواره یک مبادله میان واریانس و اریبی رخ می‌دهد (شکل ۳.۲). زیرا در برآورد WLS پارامتر نقطه i ، $\hat{\beta}_i^{WLS}$ ، اگر پهنای باند را کاهش دهیم و فقط به مشاهدات خیلی نزدیک بسنده کنیم، با این کار به دلیل کاهش اندازه نمونه، خطای استاندارد $\hat{\beta}_i^{WLS}$ افزایش می‌یابد. همچنین با افزایش b خطای استاندارد کاهش می‌یابد و اریبی افزایش می‌یابد.



شکل ۳.۲: مبادله واریانس-اریبی

برای حل مشکل انتخاب پهنای باند بهینه، b^{opt} ، اعتبارسنجی متقابل^{۱۹} CV کاربرد گسترده‌ای دارد (سیلورمن، ۱۹۸۴). به عبارتی اعتبارسنجی متقابل به دنبال پهنای باندی برای تابع هسته است که باقی‌مانده‌های برآورد y_i ها توسط یک زیرمجموعه از داده‌ها را مینیمم کند. این رهیافت برای رگرسیون موضعی توسط سلولند (۱۹۷۹) و برای برآورد تابع هسته توسط بومن (۱۹۸۴)، پیشنهاد شد.

برونس‌دان و همکاران (۱۹۹۶) رهیافت CV را به صورت

$$CV = \sum_{i=1}^n (y_i - \hat{y}_{\neq i}(b))^2$$

برای رگرسیون موزون جغرافیایی معرفی کردند که در آن مقدار برآوردشده Y_i پس از حذف مشاهده i از فرآیند است و مجموع توان‌های دوم باقی‌مانده پس از حذف مشاهده i کمینه می‌شود. به‌عنوان مثال، فودرینگ‌هام و همکاران (۲۰۰۳) در تحلیل داده‌های قیمت

^{۱۹}Cross-validation

مسکن شهر لندن برای انتخاب بهترین همسایگی در رگرسیون موزون جغرافیایی با هسته فضایی سازوار، از این معیار استفاده کردند.

یک رهیافت متفاوت و کلی‌تر که بر پایه مبادله بین درجه آزادی و نیکی برازش است، معیار اطلاع آکاییک^{۲۰} (AIC) است (آکاییک، ۱۹۷۶). در این روش مقدار بهینه پهنای باند از طریق کمینه کردن مقدار AIC به دست می‌آید. این معیار در متون مختلف با صورت‌های مختلف ارائه شده است. فودرینگ‌هام و همکاران (۲۰۰۳) یک نسخه تصحیح‌شده از معیار اطلاع آکاییک را برای GWR، پیرو روش معرفی‌شده توسط هورویچ و همکاران (۱۹۹۸)، به شکل زیر معرفی کردند:

$$AIC_c = 2n \ln(\hat{\sigma}_\varepsilon) + n \ln(2\pi) + n \left\{ \frac{n + tr(\mathbf{S})}{n - 2 - tr(\mathbf{S})} \right\}$$

بنابراین مقدار پهنای باند بهینه به صورت

$$b^{opt} = \min_b AIC_c = \min_b \left\{ 2n \ln(\hat{\sigma}_\varepsilon) + n \ln(2\pi) + n \left\{ \frac{n + tr(\mathbf{S})}{n - 2 - tr(\mathbf{S})} \right\} \right\} \quad (43.2)$$

است، که در آن $\hat{\sigma}_\varepsilon$ برآورد انحراف معیار استاندارد خطا در معادله (۳۴.۲) است. توجه داشته باشید که برآورد σ_ε به روش (۴۰.۲) محاسبه نمی‌شود؛ بلکه بر اساس درست‌نمایی ماکسیمم به صورت

$$\hat{\sigma}_\varepsilon^2 = \frac{RSS}{n}$$

برآورد می‌شود؛ که در آن RSS مجموع توان‌های دوم باقی‌مانده در رگرسیون موزون جغرافیایی است. (فودرینگ‌هام و همکاران، ۲۰۰۳).

در مدل GWR درجه آزادی برابر $n - (2tr(\mathbf{S}) - tr(\mathbf{S}^T\mathbf{S}))$ است، که در آن $tr(\mathbf{S})$ مجموع درایه‌های روی قطر اصلی ماتریس طرح (۴۱.۲) است. $tr(\mathbf{S})$ در رگرسیون OLS برابر تعداد پارامترهای مدل p است. اما در مدل GWR تعداد پارامترهای موثر^{۲۱} مدل برابر

$$p \leq 2tr(\mathbf{S}) - tr(\mathbf{S}^T\mathbf{S}) \leq n$$

است که در آن p تعداد پارامترها در مدل فراموضعی متناظر است. چون عموماً $tr(\mathbf{S}^T\mathbf{S})$ و $tr(\mathbf{S})$ مشابه هستند، تعداد پارامترهای موثر مدل تقریباً برابر $tr(\mathbf{S})$ است؛ بنابراین مشخص است که اگر $2tr(\mathbf{S}) - tr(\mathbf{S}^T\mathbf{S}) = p$ آن‌گاه پهنای باند به بی‌نهایت میل می‌کند و برای برآورد هر نقطه مرجع i از تمام مشاهدات نمونه استفاده می‌شود. اما اگر $2tr(\mathbf{S}) - tr(\mathbf{S}^T\mathbf{S}) = n$ آن‌گاه پهنای باند به صفر میل می‌کند.

رگرسیون موزون جغرافیایی، مشکل ناهمگنی فضایی را با برازش مدل خطی به تک تک مشاهدات و تغییرات هموار ضرایب رگرسیونی روی فضای مشاهدات نمایش می‌دهد. در واقع ضرایب GWR وابسته به فضا تغییر می‌کنند، یعنی به ازای هر مشاهده $p+1$ ضریب به مدل اضافه

²⁰Akaike information criterion

²¹Effective number of parameters

می‌شود. این امر بررسی آن‌ها را در مواردی که تعداد مشاهدات و ضرایب مورد بررسی زیاد است، دشوار می‌کند و نمی‌تواند تفاوت معنی‌دار آماری ارائه دهد. موضوع اصلی در تحلیل‌های فضایی ارائه روش‌های جدید موضعی برای تحلیل فضایی یا ارائه مدل‌های پیچیده‌تر برای ساختار داده‌ها نیست، بلکه ارائه ابزارها برای کشف جامعه و مکانیسم رفتاری که الگوهای فضایی نشان داده‌اند، مهم است. به گفته انسلین (۲۰۱۰): در مدل‌های ناهمگنی فضایی، ضرایب مختلف فضایی گواه حضور ناهمگنی است؛ اما آن را توضیح نمی‌دهند. این در حالی است که همیشه هدف ایده‌آل تهیه یک ساختار از ناهمگنی و (یا) خودهمبستگی است. بر این اساس، در فصل سوم مدلی را که همچون GWR خودهمبستگی و ناهمگنی فضایی را توأم در نظر می‌گیرد، اما ایراد وارد بر GWR یعنی عدم توضیح‌پذیری به آن وارد نیست، معرفی می‌کنیم.

فصل ۳

مدل رگرسیون موزون موضعی تکراری

۱.۳ مقدمه

مشکل ناهمگنی فضایی به شکل پارامترهای مختلف فضایی در تحلیل‌های فضایی و به ویژه اقتصادسنجی فضایی ناشناخته مانده است؛ چرا که عموماً هدف آن‌ها کنترل اثرات سرریز و خودهمبستگی فضایی است. به گفته پوستیچ‌لیون و همکاران (۲۰۱۳) مشکل ناهمگنی فضایی اغلب در تحلیل‌های اقتصادی داده‌های فضایی کنار گذاشته و فراموش می‌شود و این عدم توجه می‌تواند در برآورد ضرایب مدل اثرات مخرب آشکاری داشته باشد.

گروهی از محققان نظیر انسلین (۱۹۸۸b)، پد و همکاران (۲۰۱۴) و برخی دیگر قصد داشتند که حضور ناهمگنی فضایی را با ساخت آزمون‌هایی که معمولاً بر پایه آماره ضریب لاگرانژ هستند، آشکار کنند. متأسفانه هنگامی که حضور ناهمگنی فضایی شناسایی شد، هیچ آزمونی قادر به نشان دادن مدل‌بندی درست داده‌های فضایی نیست.

اخیراً ابراگیموف و مولر (۲۰۱۰) ویژگی‌های آماره t را در نمونه‌های کوچک و بزرگ و همچنین در برخورد با داده‌های فضایی، با این فرض که داده‌ها به q گروه منطقی افزاز شده باشند، بررسی کردند. سوال کلیدی این رهیافت در تعداد و ترکیب گروه‌ها است. در این مورد آن‌ها تأکید کردند که نمی‌توان تصمیم‌گیری درباره انتخاب گروه مناسب را به داده‌ها محول کرد و به اطلاعات پیشین در خصوص ساختار همبستگی آن‌ها نیاز است. این در حالی است که در بسیاری از موارد عملی هیچ‌گونه اطلاعی درباره ساختار همبستگی و ناهمگنی داده‌ها

که باعث شود گروهی به دیگری ترجیح داده شود، وجود ندارد. در این فصل یک مدل ایده‌آل برای ساختار داده‌های فضایی معرفی می‌شود که ناهمگنی و وابستگی فضایی را با هم در بر دارد. مدل پیشنهادی ما شامل یک روش قدرتمند تکراری بر پایه ترکیب روش رگرسیون موزون جغرافیایی و الگوریتم هموارسازی وزن‌های سازوار^۱ (AWS) است. با ترکیب این دو روش، نواحی همگن گروه‌بندی و مشخص می‌شوند. در متون اقتصادی به این نواحی همگن، رژیم‌های فضایی می‌گویند. همچنین می‌توان با لحاظ کردن این رژیم‌های فضایی در یک مدل فراموضعی، ضرایب رگرسیونی را به‌طور مجزا برای هر رژیم برآورد و نارسایی مدل GWR در توضیح و نمایش ساختار مدل را رفع کرد. بیل و همکاران (۲۰۱۷) این روش ترکیبی برای تعیین رژیم‌های همگن را رگرسیون موزون موضعی تکراری^۲ (IGWR) و مدل‌های حاصل بر مبنای استفاده از رژیم‌های تعیین‌شده توسط IGWR را مدل‌های با رژیم فضایی درونی^۳ (ESRMs) نامیدند.

در ادامه این فصل رهیافت پیشنهادی خود را در دو مرحله کلی شامل:

مرحله اول: الگوریتم کنترل ناهمگنی فضایی (IGWR)

مرحله دوم: مدل‌های اقتصادسنجی فضایی با رژیم فضایی درونی (ESRMs)

را شرح خواهیم داد.

۲.۳ مرحله اول: کنترل ناهمگنی فضایی

در مرحله اول رهیافت، یعنی IGWR، نواحی همگن توسط یک الگوریتم تکراری مشخص می‌شوند. ایده این روش از الگوریتم AWS که اولین بار توسط پولزهل و اسپوکوینی (۲۰۰۰ و ۲۰۰۶) در متون شبیه‌سازی مطرح شد، سرچشمه می‌گیرد. هدف الگوریتم AWS تعریف بزرگترین ناحیه همسایگی ممکن برای هر نقطه i در فضای مشاهدات است. این همسایگی باید به گونه‌ای باشد که مدل پارامتری هر ناحیه کاملاً منطبق بر داده‌های آن باشد؛ یعنی به وسیله یک مدل با مجموعه‌ای از پارامترهای ثابت به خوبی تقریب زده شوند. طرز کار این الگوریتم مبتنی بر گسترش متوالی ناحیه همسایگی حول نقطه i و تعریف مدل‌های موضعی بر اساس تخصیص وزن، وابسته به نتایج گام قبل، به هر واحد فضایی است. بدین منظور، الگوریتم از یک همسایگی خیلی کوچک برای هر نقطه شروع و در طی فرآیند تکرار تمام همسایگی‌ها با ورود نقاط جدیدی که موجب نقض پارامترهای موضعی نشوند، گسترش می‌یابند.

مقصود رهیافت IGWR بیان این حقیقت است که می‌توان قدرت روش‌های ناپارامتری و سودمندی الگوریتم AWS را ترکیب و به این طریق نواحی همگن را با برازش مدل خاص به

¹Adaptive weights smoothing

²Iterative geographically weighted regression

³Endogenous spatial regime models

هر کدام مشخص کرد. به طور دقیق در این رهیافت، روش رگرسیون موزون جغرافیایی به طور مداوم تکرار و در هر تکرار وزن‌های GWR یعنی درایه‌های قطر اصلی ماتریس W_i به روزرسانی و مجدداً GWR برازش داده می‌شود. این به روزرسانی و اختصاص وزن‌های جدید بر اساس الگوریتم AWS انجام می‌شود؛ یعنی برای هر نقطه i ، وزن‌های w_{ij} در جهت گسترش ناحیه همسایگی همگن آن به روز می‌شوند. بنابراین به منظور کنترل همگنی در هر مرحله، انجام آزمون به روی ضرایب برآوردشده در (۳۷.۲) الزامی است. به این ترتیب الگوریتم IGWR دارای مراحل زیر است:

– اختصاص وزن‌های اولیه و برازش مدل GWR

– تکرار مراحل تا برقراری شرط توقف فرآیند

– انجام آزمون همگنی ضرایب

– الگوریتم به روزرسانی وزن‌ها

در ادامه هر کدام از مراحل و هدف از انجام آن‌ها به ترتیب تشریح می‌شوند.

۱.۲.۳ اختصاص وزن‌های اولیه و برازش مدل GWR

مرحله اول رهیافت IGWR با برازش مدل GWR آغاز می‌شود. برای این کار، باید ماتریس وزن‌های W_i معرفی شوند. لذا آن گونه که در فصل دوم گفته شد، از یک تابع هسته به شکل $w_{ij}^0 = K(d_{ij}; b)$ استفاده می‌کنیم که تابعی از فاصله بین دو مشاهده i, j و مقدار پهنای باند است. همان طور که گفته شد $K(\cdot)$ یک تابع هسته دلخواه و w_{ij}^0 نشان‌دهنده وزن‌های اولیه (تکرار صفر) است. مقدار پهنای باند مناسب به روش AIC گفته شده در فصل دوم، محاسبه و در تابع هسته انتخاب شده قرار می‌گیرد.

پس از آن که وزن‌های اولیه معرفی شدند، ضرایب رگرسیونی و همچنین واریانس جمله خطا، با استفاده از (۳۷.۲) و (۴۰.۲) برآورد می‌شوند و آن‌ها را با نماد $\hat{\beta}_i^0$ و $\hat{\sigma}_{\varepsilon_i}^2$ که نشان‌دهنده برآورد ضرایب GWR و واریانس جمله خطا در تکرار صفر هستند، نشان می‌دهیم.

قابل ذکر است که آندرانو و همکاران (۲۰۱۶) در روش خود پهنای باند بهینه اولیه را اعمال نکردند. در صورتی که تحمیل پهنای باند بهینه به روش AIC برای وزن‌های اولیه، در صورتی که پهنای باند سازوار مورد استفاده باشد، برای همگرا شدن شرط توقف که در بخش بعد به آن می‌پردازیم، لازم است.

۲.۲.۳ شرط توقف الگوریتم تکرار

پس از معرفی وزن‌های اولیه و انجام رگرسیون موزون جغرافیایی، یک ماتریس $n \times (p+1)$ از ضرایب رگرسیونی به شکل

$$\beta = \begin{bmatrix} \beta_{10} & \beta_{11} & \cdots & \beta_{1p} \\ \vdots & \vdots & \cdots & \vdots \\ \beta_{i0} & \beta_{i1} & \cdots & \beta_{ip} \\ \vdots & \vdots & \cdots & \vdots \\ \beta_{n0} & \beta_{n1} & \cdots & \beta_{np} \end{bmatrix}$$

داریم. یعنی هر مشاهده i یک بردار $1 \times (p+1)$ از ضرایب رگرسیونی دارد. برای مشخص شدن نواحی همگن یا به عبارتی مشاهداتی که رفتار و ضرایب مدل در آن‌ها مشابه است، باید این بردار از ضرایب دو به دو با یکدیگر مقایسه شوند. در صورت وجود شباهت با به‌روزرسانی وزن‌ها شرط قرارگیری مشاهدات آن‌ها در یک رژیم فراهم می‌شود. پس از به‌روزرسانی وزن‌ها، برازش مدل GWR با ماتریس‌های W_i از وزن‌های جدید تکرار می‌شود. اما این تکرار و به‌روزرسانی وزن‌ها تا کجا ادامه می‌یابد؟

تا زمانی که شرط (۱.۳)، به ازای مقادیر کوچک ω (برای مثال 0.0001) برقرار شود، فرآیند تکرار ادامه می‌یابد و وزن‌ها در هر تکرار به‌روز می‌شوند:

$$\max |w_{ij}^{(\ell-1)} - w_{ij}^{\ell}| < \omega \quad \forall ij, i \neq j. \quad (1.3)$$

همان‌طور که گفته شد در هر تکرار، برای مثال تکرار ℓ ام، ضرایب مدل GWR برآورد شده برای نقاط مختلف فضایی دو به دو با یکدیگر مقایسه می‌شوند. یعنی β_i^{ℓ} با β_j^{ℓ} به ازای همه $i, j : i \neq j$ مقایسه می‌شوند تا مشخص شود آیا می‌توان مشاهدات آن‌ها را در یک رژیم فضایی قرار داد. مسلماً با افزایش اختلاف β_i^{ℓ} و β_j^{ℓ} احتمال قرارگیری i, j در یک رژیم کاهش می‌یابد. لذا نیازمند یک آماره آزمون بر مبنای این اختلاف یا فاصله ضرایب هستیم.

۳.۲.۳ آماره آزمون همگنی ضرایب

برای مقایسه ضرایب مدل GWR در هر تکرار، نیازمند یک آماره آزمون هستیم. آماره

$$\chi_{ij}^{\ell} = (\hat{\beta}_i^{\ell} - \hat{\beta}_j^{\ell})^T (\Sigma^{\ell})^{-1} (\hat{\beta}_i^{\ell} - \hat{\beta}_j^{\ell}) \quad (2.3)$$

به این منظور معرفی شده است (بیل و همکاران، ۲۰۱۷)؛ که در آن $\hat{\beta}_i^{\ell}$ و $\hat{\beta}_j^{\ell}$ برآورد ضرایب مدل GWR برای دو مشاهده i و j در تکرار ℓ ام هستند. همچنین Σ^{ℓ} ، ماتریس واریانس ادغام‌شده^۴

^۴Pooled variance matrix

نامیده می‌شود و از میانگین دو ماتریس واریانس ضرایب β_i^ℓ و β_j^ℓ در تکرار ℓ ام، به صورت زیر به دست می‌آید:

$$\Sigma^\ell = \frac{tr(\mathbf{W}_i^\ell)\Sigma_i^\ell + tr(\mathbf{W}_j^\ell)\Sigma_j^\ell}{tr(\mathbf{W}_i^\ell) + tr(\mathbf{W}_j^\ell)}$$

که در آن ماتریس واریانس $\hat{\beta}_i^\ell$ طبق (۳۹.۲) برای هر تکرار به صورت

$$\Sigma_i^\ell = [(\mathbf{X}^T \mathbf{W}_i^\ell \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_i^\ell][(\mathbf{X}^T \mathbf{W}_i^\ell \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_i^{\ell T} \hat{\sigma}_{\varepsilon_i}^2]$$

محاسبه می‌شود.

آماره χ_{ij} یک تابع درجه دوم بر حسب فاصله است. مشخص است مقادیر بزرگ آن نشان‌دهنده اختلاف بین مقادیر $\hat{\beta}_i$ و $\hat{\beta}_j$ است. این اختلاف بین ضرایب رگرسیونی دو موقعیت فضایی i و j ، گواه وجود ناهمگنی فضایی است. بنابراین دو موقعیت i و j نمی‌توانند در یک رژیم فضایی قرار بگیرند و بالعکس. یعنی مقادیر کوچک آماره χ_{ij} نشان‌دهنده اختلاف اندک میان ضرایب رگرسیون دو موقعیت فضایی است و بیان‌گر این حقیقت که دو موقعیت i و j می‌توانند در یک رژیم فضایی قرار گرفته و ضرایب یکسان داشته باشند.

با توجه به آن‌چه گفته شد، برای بررسی دقیق برآورد ضرایب و رژیم‌بندی موقعیت‌ها نیازمند آزمون آماری بر مبنای آماره χ_{ij} هستیم، که در بخش بعد معرفی می‌شود.

۴.۲.۳ آزمون همگنی ضرایب

همان‌طور که در فصل دوم گفته شد، هر بردار $\hat{\beta}_i$ یک برآوردگر ضرایب رگرسیونی به روش WLS برای هر مشاهده $n, \dots, 1, i$ است. این برآوردگر با استفاده از یک زیرنمونه از مشاهدات که توسط پهنای باند و تابع هسته انتخابی مشخص شده‌اند، برآورد را محاسبه می‌کند. اگر تعداد مشاهدات در همه زیرنمونه‌های معرفی شده به اندازه کافی بزرگ باشد، آن‌گاه طبق (۳۵.۲) و (۳۷.۲) نتیجه زیر برقرار است (ابراگیموف و مولر، ۲۰۱۰):

$$\hat{\beta}_i^\ell \sim N(\beta_i^\ell, \Sigma_i^\ell).$$

همچنین $\hat{\beta}_i$ و $\hat{\beta}_j$ در هر تکرار مستقل از یکدیگرند (برای اطلاعات بیشتر به ابراگیموف و مولر (۲۰۱۰) مراجعه شود).

بنابراین آماره (۲.۳) برای هر تکرار دارای توزیع مجانبی χ^2 با $p+1$ (برابر با بعد بردار $\hat{\beta}_i$) درجه آزادی است. بدین ترتیب آزمون با فرضیه‌های زیر در هر تکرار انجام می‌گیرد:

$$\begin{aligned} H_0 &: \hat{\beta}_i = \hat{\beta}_j \\ H_1 &: \hat{\beta}_i \neq \hat{\beta}_j \end{aligned} \quad (۳.۳)$$

در این آزمون اگر آماره آزمون مقادیر بزرگ اختیار کند، یعنی $\chi_{ij}^l > \chi_{p+1}^2(1-\alpha)$ ، آن گاه فرضیه صفر در سطح α رد می‌شود. یعنی ضرایب رگرسیونی دو موقعیت i و j یکسان نبوده و این دو موقعیت فضایی به دو رژیم مختلف فضایی تعلق دارند.

آماره χ_{ij} در هر تکرار به عنوان معیار برای مجازات عمل می‌کند. به این ترتیب که اگر مقدار آن بزرگ باشد، لازم است که در تکرار بعد به وزن w_{ij} مقدار کوچک‌تری تخصیص داده شود تا دو مشاهده i و j که ناهمگن هستند در دو رژیم مختلف قرار بگیرند. همچنین اگر مقدار χ_{ij} کوچک باشد، آن گاه باید وزن آن‌ها در تکرار بعد افزایش یابد. البته برای اطمینان از همگرا بودن وزن‌های نهایی تخصیص یافته به مقادیر معقول صفر و یک، باید روش به‌روزرسانی وزن‌ها به آرامی انجام شود. در ادامه بیان می‌کنیم که این به‌روزرسانی (میزان افزایش یا کاهش) وزن‌ها در تکرار بعدی به چه شکل اجرا می‌شود.

۵.۲.۳ الگوریتم به‌روزرسانی وزن‌ها

همان‌طور که در بخش‌های پیش گفته شد، اولین گام در انجام رهیافت IGWR، به‌منظور مشخص کردن ناهمگنی فضایی و در نتیجه آن رژیم‌های فضایی، معرفی وزن‌های اولیه رگرسیون موزون جغرافیایی با استفاده از یک تابع وزن‌دهی و پهنای باند بهینه است. عموماً برای انجام این مهم در مدل GWR (گام اول IGWR)، توابع هسته و معیار AIC به کار گرفته می‌شوند. اما در گام‌های بعدی این رهیافت (یعنی تکرار اول به بعد)، وزن‌ها به این طریق محاسبه نمی‌شوند، بلکه مقدار آن‌ها یا به عبارتی به‌روزرسانی آن‌ها نسبت به مرحله قبل با توجه به نتیجه آماره آزمون (۲.۳) صورت می‌گیرد. ایده این روش وزن‌دهی بر پایه الگوریتم هموارسازی وزن‌های سازوار است. همچنین برای اطمینان از همگرا بودن وزن‌های تخصیصی نهایی به مقادیر منطقی w_{ij}^l و ۱، باید فرآیند همگرایی وزن‌ها به آهستگی انجام شود. در نتیجه تابع وزن به‌روزرسان w_{ij}^l برای هر تکرار l را به صورت ضرب تابع هسته $K(d_{ij}; b)$ در یک هسته که تابعی از آماره آزمون χ_{ij} در تکرار l است، معرفی می‌کنیم.

بنابراین قرار می‌دهیم

$$w_{ij}^l = K(d_{ij}^l; b)K(\chi_{ij}; \tau), \quad l = 1, 2, \dots \quad (4.3)$$

در این تابع $d_{ij}^l = d_{ij}/l$ می‌باشد. این تعریف از فاصله برای هر تکرار به این منظور صورت می‌گیرد تا تضمین کند با جلو بردن فرآیند به‌روزرسانی وزن و تکرارهای مجدد، از دقت برآورد کاسته نشود. در مورد d_{ij}^l همچنین می‌توان گفت که با افزایش تکرارها (l)، وزن فاصله d_{ij} کاهش می‌یابد. افزایش وزن‌ها به χ_{ij}^l وابسته است.

تابع هسته دوم در (۴.۳) به شکل زیر معرفی می‌شود:

$$K(\chi_{ij}^l; \tau) = \exp(-\tau / \chi_{ij}^l)$$

که در آن τ پارامتر توزین کردن آماره آزمون χ_{ij}^ℓ است. مانند پهنای باند مقادیر نسبتاً بزرگ آن می‌تواند به عملکرد ناپایدار روش منجر شود؛ در حالی که مقادیر کوچک آن حساسیت به تغییر ساختار را کاهش می‌دهد.

پولزهل و اسپوکوینی (۲۰۰۶)، با توجه به این نکته که $\chi_{ij}^\ell \sim \chi_{(p+1)}^2$ ، مقدار τ را به صورت زیر معرفی کردند:

$$\tau = \frac{1}{Q_\alpha(\chi_{(p+1)}^2)}$$

که در آن $Q_\alpha(\chi_{(p+1)}^2)$ چندک α ام توزیع کی دو با $(p+1)$ درجه آزادی است. آن‌ها در نهایت وزن‌های معرفی شده در (۴.۳) را یکبار دیگر به روزرسانی کردند. این به روزرسانی با میانگین‌گیری بین وزن‌های حاصل w_{ij} در تکرار ℓ ام و مقادیر میانگین محاسبه شده آن‌ها در تکرار قبلی $\bar{w}_{ij}^{(\ell-1)}$ ، به صورت زیر حاصل می‌شود:

$$\bar{w}_{ij}^\ell = (1 - \eta)\bar{w}_{ij}^{(\ell-1)} + \eta w_{ij}^\ell. \quad (5.3)$$

در این تابع وزن، η پارامتر کنترل است و می‌تواند مقادیر بازه $(0, 1)$ را اختیار کند. برای مثال اگر $\eta = 0.5$ تنظیم شود، آن‌گاه دقیقاً میانگین حسابی وزن‌ها در تکرار ℓ ام و تکرار $(\ell-1)$ ام محاسبه می‌شود. پولزهل و اسپوکوینی η را پارامتر حافظه^۵ نیز نامیده‌اند. با این پارامتر، چگونگی و میزان وزن‌ها در تکرار قبل در تعیین وزن‌ها در تکرار جدید موثر خواهد شد و میزان این تاثیر نیز تحت کنترل خواهد بود. به این ترتیب وزن‌های جدید هموار هستند.

در طی فرآیند همگرایی و مشخص شدن رژیم‌ها، کافی بودن مشاهدات هر گروه بسیار مهم است. اگر مشاهدات یک گروه تعیین شده از تعداد پارامترهای مدل کمتر باشد، آن گروه و مشاهدات آن به عنوان داده دورافتاده از محاسبات خارج می‌شوند. در نهایت تعریف مفید برای تشریح همگرایی، تابع تغییرات وزن یعنی

$$d(w) = \max_{i,j} |\bar{w}_{ij}^{(\ell-1)} - w_{ij}^\ell|, \quad (6.3)$$

است. این تابع، $d(w)$ ، هم‌زمان با تثبیت شدن فرآیند همگرایی و اتمام مرحله اول، به صفر میل می‌کند.

قبل از ورود به مرحله دوم رهیافت و تشریح آن، مطالب گفته شده در این بخش، یعنی مراحل فرآیند همگرایی و تعیین رژیم‌ها به روش IGWR را به‌طور خلاصه بیان می‌کنیم.

– تعریف وزن‌های اولیه W_i توسط یک هسته که تابعی از فاصله هر دو نقطه در فضا و یک مقدار پهنای باند است؛ و استفاده از این وزن‌ها در (۳۷.۲) برای برآورد ضرایب رگرسیون جغرافیایی.

⁵Memory parameter

– پس از برآورد ضرایب اولیه، فرآیند همگرایی با جایگذاری وزن‌های جدید در قطر اصلی ماتریس W_i حاصل می‌شود. وزن‌های جدید به واسطه یک فرآیند به‌روزرسانی وزن‌ها با روش AWS محاسبه می‌شوند.

– فرآیند به‌روزرسانی وزن‌ها، یک فرآیند هموارساز است که به تدریج وزن‌های اولیه را به واسطه ترکیب آن‌ها با یک وزن تعریف‌شده بر اساس آماره آزمون χ_{ij} ، معادله (۴.۳)، و سپس دوباره به‌روزرسانی وزن حاصل با محاسبه میانگین وزنی آن‌ها و وزن‌های نهایی در تکرار قبل، معادله (۵.۳)، تعدیل می‌کند.

– با قرار دادن $\eta = 0/5$ در (۵.۳)، یک ساختار تدریجی برای وزن‌های جدید تضمین می‌شود. به‌روزرسانی وزن‌ها برای استفاده مجدد در (۳۷.۲) تا همگرایی، یعنی زمانی که $d(w) \rightarrow 0$ ادامه می‌یابد.

همان‌طور که می‌بینیم، این روش هیچ‌گونه محدودیتی لحاظ نمی‌کند و کاملاً سازگار است. یعنی به اطلاعات پیشین درباره ساختار فضایی داده‌ها نیاز ندارد.

۳.۳ مرحله دوم: مدل‌های اقتصادسنجی با رژیم فضایی درونی

پس از اتمام مرحله اول و مشخص شدن ناهمگنی فضایی تحت رژیم‌های مختلف فضایی، نوبت به قرار دادن آن‌ها در مدل‌های فراموضعی است. مدل‌های حاصل، مدل‌هایی مشابه قبل هستند با این تفاوت که حضور رژیم‌های شناسایی‌شده توسط روش IGWR با ایجاد افزایش در متغیر پاسخ، متغیرهای تبیینی و سایر اجزای مدل، موجب برآورد ضرایب رگرسیونی مختلف می‌شود. در این پایان‌نامه، مدل‌های اقتصادسنجی مورد استفاده هستند. مدل‌های پرکاربرد اقتصادسنجی در فصل دوم به‌طور کامل معرفی شدند. از آن میان برای تشریح مرحله دوم رهیافت، تنها به مدل‌های SAC و SDM می‌پردازیم، چراکه همان‌طور که گفته شد اثرات فضایی به‌طور کامل‌تری در آن‌ها لحاظ شده است و همچنین این دو مدل حاصل ترکیب دو مدل SLM و SEM (به دو روش مختلف) هستند. یعنی در حالت خاص $\rho = 0$ ، مدل SEM با رژیم فضایی درونی و اگر $\lambda = 0$ ، مدل SLM با رژیم فضایی درونی حاصل می‌شود.

۱.۳.۳ مدل ترکیبی فضایی با رژیم درونی (ESR-SAC)

برای شروع، صورت کلی مدل ترکیبی فضایی را یادآوری می‌کنیم. ساختار سلسله‌مراتبی مدل SAC به صورت زیر است:

$$\begin{aligned} y &= \rho W y + X \beta + v, \\ v &= \lambda M v + \varepsilon, \\ \varepsilon &\sim N(0, \sigma^2 I_n). \end{aligned}$$

در این مدل y ، v و ε بردارهای با بعد n ، W و M ماتریس‌های وزن مربعی از بعد n ، β برداری با بعد $p+1$ و X یک ماتریس $(p+1) \times n$ بعدی است. پس از اتمام فرآیند IGWR و مشخص شدن تعداد رژیم‌ها و عضوهای هر کدام، بردارها و ماتریس‌های ذکر شده به شکل‌های زیر افراز می‌شوند:

$$\begin{aligned} y &= \begin{pmatrix} y_1 \\ \vdots \\ y_c \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_c \end{pmatrix} \quad v = \begin{pmatrix} v_1 \\ \vdots \\ v_c \end{pmatrix} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_c \end{pmatrix} \\ X &= \begin{pmatrix} X_1 & \cdots & \mathbf{o} \\ \vdots & \ddots & \vdots \\ \mathbf{o} & \cdots & X_c \end{pmatrix} \end{aligned} \quad (7.3)$$

به عنوان مثال y_j ، $j = 1, \dots, c$ بردار مقادیر متغیر پاسخ و X_j ماتریس مقادیر متغیرهای تبیینی در رژیم j ام هستند که بعد آن‌ها به ترتیب $1 \times (n_j)$ و $(p+1) \times (n_j)$ است.

ملاحظه ۱.۳.۳. توجه داشته باشید که $\sum_{j=1}^c n_j = n$ و بعد ماتریس X ، $n \times (c(p+1))$ است. بنابراین مدل SAC پس از درونی کردن رژیم‌های فضایی به شکل زیر خواهد بود:

$$\begin{aligned} \begin{pmatrix} y_1 \\ \vdots \\ y_c \end{pmatrix} &= \rho W \begin{pmatrix} y_1 \\ \vdots \\ y_c \end{pmatrix} + \begin{pmatrix} X_1 & \cdots & \mathbf{o} \\ \vdots & \ddots & \vdots \\ \mathbf{o} & \cdots & X_c \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_c \end{pmatrix} + \begin{pmatrix} v_1 \\ \vdots \\ v_c \end{pmatrix} \\ \begin{pmatrix} v_1 \\ \vdots \\ v_c \end{pmatrix} &= \lambda M \begin{pmatrix} v_1 \\ \vdots \\ v_c \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_c \end{pmatrix} \end{aligned} \quad (8.3)$$

که در آن $\varepsilon_j \sim N(0, \sigma_{\varepsilon_j}^2 I_{n_j})$ ، $i = 1, 2, \dots, c$

۲.۳.۳ مدل دوربین فضایی با رژیم درونی (ESR-SDM)

شکل افراز شده جملات معادله (۳۰.۲) برای مدل SDM به صورت

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_c \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_c \end{pmatrix} \quad \alpha = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_c \end{pmatrix} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_c \end{pmatrix}$$

$$X = \begin{pmatrix} X_1 & \cdots & \mathbf{o} \\ \vdots & \ddots & \vdots \\ \mathbf{o} & \cdots & X_c \end{pmatrix} \quad X_v = \begin{pmatrix} X_{v_1} & \cdots & \mathbf{o} \\ \vdots & \ddots & \vdots \\ \mathbf{o} & \cdots & X_{v_c} \end{pmatrix} \quad (9.3)$$

است. مشابه معادله (۲۹.۲)، X_v ماتریسی مشابه X است با این تفاوت که عرض از مبدأ از آن حذف شده است.

نکات گفته شده در مدل ESR-SAC، در این مدل نیز برقرار هستند. به علاوه بردار α از بعد cp و ماتریس X_v دارای بعد $n \times cp$ هستند.

مدل دوربین با رژیم فضایی درونی شده، به صورت زیر نوشته می شود:

$$\begin{pmatrix} y_1 \\ \vdots \\ y_c \end{pmatrix} = \rho W \begin{pmatrix} y_1 \\ \vdots \\ y_c \end{pmatrix} + \begin{pmatrix} X_1 & \cdots & \mathbf{o} \\ \vdots & \ddots & \vdots \\ \mathbf{o} & \cdots & X_c \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_c \end{pmatrix} + M \begin{pmatrix} X_{v_1} & \cdots & \mathbf{o} \\ \vdots & \ddots & \vdots \\ \mathbf{o} & \cdots & X_{v_c} \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_c \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_c \end{pmatrix} \quad (10.3)$$

که در آن $\varepsilon_j \sim N(0, \sigma_{\varepsilon_j}^2 \mathbf{I}_{n_j})$ ، $j = 1, 2, \dots, c$

ماتریس های وزن W و M را می توان در مدل های اقتصادسنجی با رژیم فضایی درونی، همانند مدل های متناظر فراموضعی خود، برابر در نظر گرفت. برآورد پارامترهای این مدل ها نیز مشابه حالت فراموضعی خود محاسبه می شود.

فصل ۴

ارزیابی مدل پیشنهادی

۱.۴ مقدمه

هدف از این فصل آشکار کردن حضور ناهمگنی فضایی در داده‌های اقتصادی است. پس از این نمایش، نوبت به اثبات کارایی رهیافت IGWR، در مدل‌بندی این ناهمگنی به شکل رژیم‌های مختلف فضایی درون مدل‌های فراموضعی اتورگرسیو است. به این منظور از دو مجموعه داده واقعی با نام‌های بالتیمور^۱ و مسکن^۲ استفاده می‌کنیم. هر دوی این مجموعه داده‌ها در بسته‌افزار spdep^۳ (بیوند و همکاران، ۲۰۱۵) در نرم‌افزار R (گروه اصلی R، ۲۰۱۸) موجود و قابل دسترس هستند.

۲.۴ مثال کاربردی اول

مجموعه داده بالتیمور شامل اطلاعات قیمت فروش مسکن و برخی خصوصیات آن‌ها نظیر متراژ منزل، تعداد پارکینگ‌ها، واقع بودن در بخش شهری یا حوالی آن و غیره است. متغیرهای موجود در این مجموعه در جدول ۱.۴ فهرست شده‌اند. این اطلاعات مربوط به ۲۱۱ مسکن

^۱baltimore

^۲house

^۳Spatial dependence

در کلان شهر بالتیمور واقع در ایالت ماری‌لند کشور آمریکا است. شکل ۱.۴ موقعیت کلان شهر بالتیمور را در ایالت ماری‌لند و ایالات متحده آمریکا نمایش می‌دهد. این مجموعه داده اولین بار توسط دوین (۱۹۹۲) مورد بررسی قرار گرفت و برای اولین بار مشکل پارامترهای مختلف در فضا از اهمیت برخوردار شد.



(ب) موقعیت در ایالات متحده آمریکا

(آ) موقعیت در ایالت ماری‌لند

شکل ۱.۴: موقعیت شهر بالتیمور

جدول ۱.۴: متغیرهای مجموعه داده بالتیمور

نام متغیر	توضیحات
price	قیمت خانه - متغیر پاسخ
dwel	متغیر ظاهری - ۱ برای خانه‌های تک خانوار
nbath	تعداد حمام‌ها
patio	متغیر ظاهری - ۱ برای خانه‌های دارای پاسیو
firepl	متغیر ظاهری - ۱ اگر شومینه موجود است
ac	متغیر ظاهری - ۱ اگر تهویه هوا موجود است
bment	تعداد زیرزمین‌ها
gar	تعداد گاراژها
citcou	متغیر ظاهری - ۱ برای خانه‌هایی که در بخش بالتیمور واقع هستند
lotsz	متراژ مسکن به هزار فوت مربع*

* هر فوت مربع برابر ۰/۰۹۲۹۰۳۰۴ مترمربع

۱.۲.۴ انجام آزمون‌های آماری و برازش مدل‌های فراموضعی فضایی

به منظور بررسی وجود خودهمبستگی فضایی در داده‌ها از آزمون‌های معرفی شده در فصل اول استفاده می‌کنیم. توابع این آزمون‌ها در بسته‌افزار spdep موجود هستند.

جدول ۲.۴: نتایج آزمون‌های خودهمبستگی فضایی برای داده‌های بالتیمور

مقدار p -مقدار	مقدار آماره	آزمون‌های خودهمبستگی فضایی
۰/۰۰۰۱	۰/۰۸۷۲	موران I
۰/۰۰۱۱	۱۰/۵۷	نسبت درست‌نمایی
۰/۰۰۰۳	۸/۷۵۶۳	ضریب لاگرانژ خطا
< ۰/۰۰۰۰	۲۰/۷۱۹	ضریب لاگرانژ تأخیر

نتایج آزمون‌های خودهمبستگی نشان‌دهنده وجود خودهمبستگی فضایی در داده‌های بالتیمور است. آزمون ضریب لاگرانژ تأخیر فضایی با قدرت بیشتری، نسبت به سایر آزمون‌ها، فرضیه صفر خود را رد کرده است. یعنی با احتمال بیشتری خودهمبستگی فضایی، در متغیرهای پاسخ داده‌های بالتیمور وجود دارد. با توجه به این نتایج، برآورد به روش OLS برای این مجموعه داده مناسب نبوده و باید از مدل‌های اتورگرسیو فضایی برای برازش داده‌ها استفاده کرد.

برای این داده‌ها از چهار مدل اتورگرسیو خطای فضایی، تأخیر فضایی، دوربین فضایی و ترکیبی فضایی استفاده کردیم. در این مدل‌ها وزن‌ها (W) به روش k -نزدیک‌ترین همسایگی با $k = 10$ محاسبه شدند. همچنین در مدل‌های SAC و SDM، $W = M$ در نظر گرفته شد. برای برازش مدل از توابع موجود در بسته‌افزار spdep در نرم‌افزار R استفاده کردیم. نتایج این مدل‌ها در جدول ۳.۴ گزارش شده‌اند.

جدول ۳.۴: برآورد پارامترهای مدل‌های فراموضعی برای مجموعه داده بالتیمور

	SDM		SAC		SEM		SLM	ضرایب
*	۲۲/۹۲	.	-۶/۱۲۵		۳/۶۷۴	.	-۶/۱۷۳	عرض از مبدأ
***	۸/۴۲۱	***	۷/۵۴۷	***	۹/۲۲۳	***	۷/۳۳۶	dwel
***	۶/۲۲۶	***	۶/۹۶۰	***	۷/۸۵۴	***	۶/۸۳۹	nbath
**	۶/۹۸۵	***	۸/۳۴۳	**	۷/۶۶۱	***	۸/۴۴۳	patio
***	۹/۰۲۴	***	۹/۵۹۲	***	۸/۸۶۸	***	۹/۷۴۰	firepl
***	۸/۷۰۴	**	۷/۰۶۱	***	۷/۹۷۴	**	۶/۹۷۰	ac
***	۳/۰۱۱	**	۳/۴۰۵	***	۳/۴۶۹	***	۳/۳۸۰	bment
***	۵/۰۴۶	***	۵/۴۰۸	**	۴/۸۶۶	***	۵/۴۵۶	gar
***	۱۱/۶۰۳	***	۹/۶۰۹	***	۱۳/۱۳۰	***	۹/۴۲۶	citcou
**	۰/۰۳۸	*	۰/۰۳۴	*	۰/۰۳۷	*	۰/۰۳۴	lotsz
.	-۱۳/۶۴۵۱							M-dwell
*	-۱۰/۵۲۲۴							M-nbath
.	۱۲/۷۶۵۵							M-patio
**	۲۱/۶۵۳۳							M-firepl
	۵/۴۱۵۵							M-ac
	-۴/۵۸۰۳							M-bment
	۶/۰۱۶۱							M-gar
	-۳/۷۷۵۱							M-citcou
.	۰/۱۱۸۳							M-lotsz
	۰/۰۶۲	***	۰/۳۲۴		-	***	۰/۳۳۴	ρ
	-		۰/۰۵۰	**	۰/۴۹۰		-	λ
	۱۶۵۱/۶۰۸		۱۶۶۶/۰۹۲		۱۶۷۲/۵۷		۱۶۶۴/۱۴	AIC

معنی دار بودن: $0/1 < . < 0/05$ * $0/01 < ** < 0/001 < ***$

بنا به نتایج جدول ۳.۴ در همه مدل‌های فراموضعی در سطح معنی‌داری ۰/۰۵، تمامی متغیرهای تبیینی معرفی شده در جدول ۱.۴ دارای اثرات معنی‌دار و مثبت در تعیین قیمت مسکن شهر بالتیمور هستند. در مدل SDM علاوه بر این متغیرها، متغیرهای تبیینی تأخیر فضایی نیز حضور دارند که تعدادی از آن‌ها دارای اثر منفی معنی‌دار هستند. برای مثال، متغیر تعداد حمام‌ها با تأخیر فضایی دارای اثر منفی معنی‌دار در سطح ۰/۰۵ است. یعنی تعداد زیاد حمام‌های مسکن i موجب افزایش قیمت فروش این مسکن اما کاهش قیمت فروش مسکن‌های مجاور آن می‌شود.

ضریب خودهمبستگی ρ در مدل SDM با ورود متغیرهای تبیینی تأخیر فضایی نسبت به مدل SLM به‌طور چشمگیری (از ۰/۳۳۴ به ۰/۰۶۲) کاهش یافته است و این بدان معناست که بی‌توجهی به همبستگی فضایی متغیرهای تبیینی و حذف اثرات آن‌ها در مدل تأخیر

فضایی موجب ایجاد خودهمبستگی ظاهری و کاذب در مدل شده است. به عبارتی ناهمگنی حاصل از حذف این متغیرها (جنبه اول- واریانس نامتجانس) در مدل SLM به اشتباه تحت خودهمبستگی قوی وارد مدل بندی شده است. برای انتخاب بهترین مدل از چهار مدل برازش داده شده، از معیار اطلاع آکاییک استفاده شد. این معیار که به صورت

$$AIC = 2k - 2\ell \quad (1.4)$$

تعریف شده است، اندازه‌ای از میزان نیکویی برازش یک مدل ارائه می‌دهد. در معادله (۱.۴)، k تعداد پارامترهای مدل و ℓ لگاریتم تابع درست‌نمایی ماکسیمم است. برای مجموعه‌ای از مدل‌ها، مدلی که کمترین مقدار AIC را داشته باشد، بهتر است. همان‌طور که از نتایج آزمون‌های خودهمبستگی نتیجه گرفتیم، خودهمبستگی فضایی در متغیر پاسخ وجود دارد. با مقایسه نتایج سطر AIC جدول ۳.۴ خواهیم دید که مدل‌هایی که خودهمبستگی را به این شکل لحاظ کرده‌اند، مدل‌های بهتری هستند (مقایسه SEM با SLM و SDM).

با توجه به نتایج سطر AIC در جدول ۳.۴، مدل دارای کمترین AIC است. اما باید توجه داشت که ضریب خودهمبستگی (ρ) در این مدل معنی‌دار نیست، همان‌طور که گفته شد بی‌توجهی به تاثیر متغیرهای تأخیر تبیینی موجب القای حضور خودهمبستگی فضایی متغیرهای پاسخ شده بود، لذا مدل مناسب برای داده‌های بالتیمور مدل رگرسیون خطی (۲.۱) با ۱۹ متغیر تبیینی ردیف شده در جدول ۳.۴، است. مقدار AIC برای این مدل ۱۶۴۹/۸ محاسبه شد که گواه این ادعا است.

اما همان‌طور که در فصل دوم گفته شد در این مدل‌ها به امکان وجود ناهمگنی فضایی حاصل از پارامترهای مختلف در فضا توجه نشده است و ناهمگنی همچون خودهمبستگی فضایی به‌طور مستقیم وارد مدل نشده است.

۲.۲.۴ ناهمگنی در داده‌ها

برای آزمون ناهمگنی در حضور خودهمبستگی فضایی در مدل‌های فراموضعی فضایی از آزمون بریچ-پیگان (۱۹.۱) استفاده کردیم. نتایج در همه مدل‌ها حاکی از وجود ناهمگنی در مقادیر خطا است.

پس از تایید وجود ناهمگنی فضا، به منظور نمایش ناهمگنی فضایی، رگرسیون موزون جغرافیایی را به داده‌ها برازش دادیم. با این کار ناهمگنی فضایی خود را در قالب اختلاف بارز در ضرایب مشاهدات، نشان می‌دهد. اما همان‌طور که گفته شد نتایج قابل تحلیل و آزمون نیستند.

جدول ۴.۴: نتایج آزمون‌های ناهمگنی فضایی داده‌های بالتیمور

مدل‌های فضایی	آماره بریج-پیگان	درجه آزادی	p-مقدار
SLM	۳۹/۷۲۸	۹	< ۰/۰۰۰۱
SEM	۳۳/۲۵۶	۹	۰/۰۰۰۱
SAC	۳۹/۳۰۶	۹	< ۰/۰۰۰۱
SDM	۴۵/۸۲۷	۱۸	۰/۰۰۰۳

رگرسیون موزون جغرافیایی با تابع `gwr.basic` در بسته‌افزار `GWmodel` (گلینی و همکاران، ۲۰۱۳) در نرم‌افزار R در دسترس است و نتایج برآورد پارامترهای مختلف نقاط جغرافیایی در جدول ۵.۴ نمایش داده شده‌اند. به‌عنوان مثال، ضرایب رگرسیونی متغیر `DWELL` برای ۲۱۱ خانه از مقدار $1/90$ تا $13/90$ تغییر می‌کند که گواه ناهمگنی فضایی به شکل پارامترهای مختلف در فضا (جنبه دوم ناهمگنی) در داده‌های بالتیمور است.

جدول ۵.۴: نتایج برازش مدل GWR برای داده‌های بالتیمور

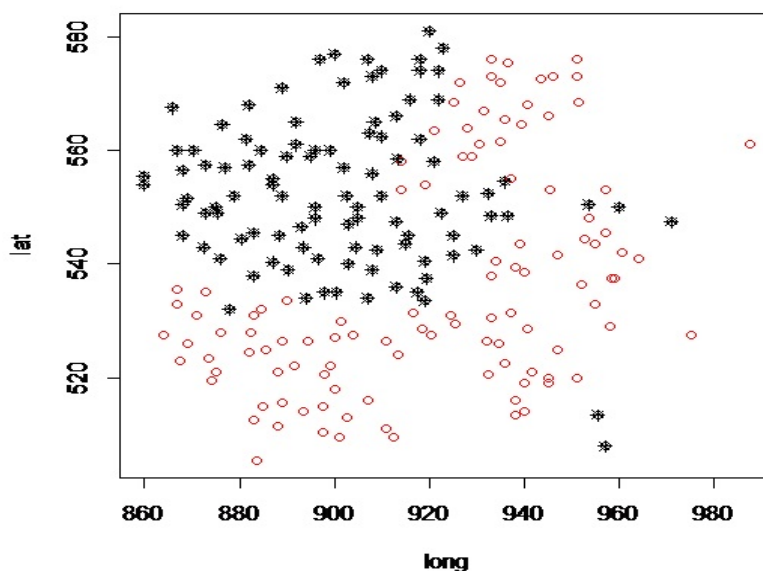
ضرایب	مینیمم	چارک اول	میانه	چارک سوم	ماکسیمم
Intercept	-۴/۸۹۲۸۱۵۳	-۰/۱۵۲۶۲۷۲	۳/۴۵۸۴۲۸۸	۵/۸۵۸۶۱۱۰	۸/۲۵۹۳
DWELL	۱/۹۰۲۶۲۲۳	۳/۸۸۴۴۰۴۶	۸/۰۵۳۱۴۷۲	۱۲/۵۰۲۷۶۰۳	۱۳/۹۰۵۳
NBATH	۴/۶۵۶۸۵۶۲	۵/۷۰۲۰۱۱۵	۶/۸۶۹۵۶۵۳	۷/۷۱۱۵۰۰۲	۱۰/۰۳۸۵
PATIO	۲/۵۲۵۰۷۶۷	۶/۰۶۱۷۷۰۳	۸/۴۲۱۳۵۵۳	۱۰/۶۴۷۸۳۴۰	۱۲/۵۷۴۷
FIREPL	۵/۰۲۷۶۳۳۰	۱۰/۲۶۸۵۰۵۱	۱۲/۶۴۳۴۰۳۷	۱۴/۲۶۹۹۶۱۶	۱۷/۰۲۶۹
AC	۴/۶۵۱۰۴۵۵	۶/۳۶۸۳۵۷۴	۷/۲۸۱۵۷۴۴	۸/۱۳۶۹۳۱۴	۱۱/۱۴۲۱
BMENT	۲/۲۷۵۱۷۴۱	۳/۱۰۸۴۲۲۷	۳/۸۶۸۸۹۰۶	۵/۱۰۰۱۱۸۲	۶/۵۰۰۲
GAR	۱/۴۷۸۱۴۰۲	۲/۳۷۲۷۲۴۴	۳/۷۵۶۳۵۸۲	۵/۶۸۹۰۱۳۵	۱۲/۲۹۵۱
CITCOU	۸/۸۵۱۰۲۵۶	۱۰/۸۶۴۰۰۱۱	۱۳/۸۰۷۱۰۷۰	۱۴/۸۵۲۹۶۸۰	۱۷/۱۳۲۶
LOTSZ	-۰/۰۵۱۵۰۱۴	۰/۰۰۱۱۸۲۶	۰/۰۴۴۷۰۵۳	۰/۰۷۶۴۸۷۷	۰/۰۹۲۳

۳.۲.۴ برازش مدل پیشنهادی

بنابر نتایج گفته‌شده، به دلیل حضور هم‌زمان ناهمگنی و خودهمبستگی فضایی مدل OLS و مدل‌های فراموضعی فضایی برای داده‌های مورد مطالعه مناسب نیستند؛ لذا از رهیافت پیشنهادی خود که ناهمگنی فضایی را در قالب رژیم‌های مختلف فضایی درون مدل‌های خودهمبستگی فضایی لحاظ می‌کند، استفاده می‌کنیم و نتایج را با مدل‌های فراموضعی مقایسه می‌کنیم.

نتایج رهیافت IGWR

در این مثال به منظور تعیین وزن‌های اولیه رگرسیون وزنی جغرافیایی از تابع هسته گاوسی با پهنای باند سازوار بر اساس معیار AIC در رابطه (۴۳.۲) استفاده کردیم. مقدار پهنای باند بهینه برای این مجموعه داده، $b_{Knn}^{AIC} = ۳۳$ ، نتیجه شد. الگوریتم به‌روزرسانی وزن‌ها با $\omega = ۰/۰۰۰۱$ ، پارامتر توزین $\tau = ۰/۰۰۱$ و پارامتر حافظه $\eta = ۰/۵$ ، پس از ۹۶ تکرار به همگرایی رسید و ۲۱۱ مشاهده در $c = ۲$ رژیم همگن با زیرنمونه‌های $n_1 = ۱۰۱$ و $n_2 = ۱۱۰$ دسته‌بندی شدند. این دو رژیم فضایی در شکل ۲.۴ با دو نماد متفاوت * و o نشان داده شده‌اند. همان‌طور که در شکل ۲.۴ مشخص است، الگوریتم IGWR نواحی مرکزی و شمال‌غربی بالتیمور را از نواحی جنوبی و شمال‌شرقی آن جدا کرده است.



شکل ۲.۴: رژیم‌های فضایی مجموعه داده بالتیمور، نقاط * رژیم فضایی اول و نقاط o رژیم فضایی دوم را نشان می‌دهند

برازش مدل‌ها با رژیم فضایی درونی

پس از مشخص شدن رژیم‌های فضایی، داده‌ها را به صورت‌های (۷.۳) و (۹.۳) افراز کرده و مدل‌های (۸.۳) و (۱۰.۳) را به آن‌ها برازش دادیم. مدل (۸.۳) را در حالت‌های خاص $\rho = ۰$ و $\lambda = ۰$ نیز در نظر گرفتیم و به ترتیب مدل‌های خطای فضایی و تأخیر فضایی با رژیم درونی را نیز برازش دادیم. در این مدل‌ها وزن W همانند مدل‌های فراموشی به روش $1 - \lambda$ نزدیکترین همسایگی محاسبه شدند؛ با این تفاوت که رژیم‌های مشخص‌شده در آن لحاظ شدند. یعنی

بردارهای حاوی اطلاعات مختصات داده‌ها نیز افزاز و سپس عملیات وزن‌دهی انجام شد. نتایج در جدول ۶.۴ گزارش شده‌اند. مشابه نتایج قبل، *** برای p- مقدار کمتر از ۰/۰۰۱، ** برای p- مقدار کمتر از ۰/۰۱ و * به ترتیب برای معنی‌دار بودن در سطح ۵٪ و ۱٪ استفاده شده‌اند.

جدول ۶.۴: برآورد پارامترهای مدل‌های با رژیم فضایی درونی برای مجموعه داده بالتیمور

SDM	SAC	SEM	SLM	ضرایب
-۵/۳۶۷ ***	-۱۶/۶۷۰	-۵/۲۷۴ ***	-۱۶/۴۰۹	عرض از مبدأ ۱
۹/۳۰۳	-۰/۶۰۲ *	۱۱/۳۷۴	۰/۲۰۶	عرض از مبدأ ۲
** ۷/۸۱۵ .	۵/۲۶۹ **	۸/۶۳۵ *	۶/۰۰۴	dwell 1
** ۸/۷۸۱ **	۷/۴۴۰ **	۷/۸۶۴ **	۷/۵۷۰	dwell 2
** ۴/۳۲۹	۵/۳۳۰ **	۵/۵۴۶ **	۵/۴۰۶	nbath 1
** ۷/۹۶۵ **	۸/۶۳۵ **	۸/۴۱۳ ***	۸/۵۲۷	nbath 2
** ۷/۸۵۰ ***	۱۰/۶۹۶ **	۸/۹۲۴ ***	۱۰/۳۲۸	patio 1
۰/۷۷۷	۰/۳۹۲	۲/۷۴۴	۱/۰۱۰	patio 2
*** ۱۲/۱۴۱ ***	۱۲/۷۳۳ ***	۱۲/۲۰۳ ***	۱۲/۴۲۰	firepl 1
* ۷/۶۳۲ *	۶/۴۴۱	۴/۹۳۴ .	۶/۰۸۷	firepl 2
*** ۱۰/۰۶۳ **	۷/۵۱۷ **	۷/۹۹۹ **	۷/۶۲۰	ac 1
۳/۵۲۶	۴/۴۷۸ *	۵/۷۹۴ .	۴/۸۰۰	ac 2
*** ۷/۴۰۱ ***	۷/۶۹۹ ***	۷/۸۵۳ ***	۷/۸۸۵	bment 1
. ۱/۸۱۴	۱/۳۹۰	۱/۵۰۳	۱/۴۰۳	bment 2
*** ۷/۰۸۵ ***	۷/۴۰۲ ***	۶/۷۸۲ ***	۷/۳۴۲	gar 1
۰/۲۴۱	۰/۰۹۰	-۰/۴۷۹	-۰/۰۸۵	gar 2
*** ۲۰/۰۴۶ ***	۱۳/۱۴۹ ***	۲۰/۲۵۲ ***	۱۴/۱۳۴	citcou 1
. ۸/۰۵۳ *	۵/۹۵۴ **	۸/۸۱۰ *	۶/۲۱۱	citcou 2
۰/۰۳۰ .	۰/۰۳۴ .	۰/۰۳۳ .	۰/۰۳۳	lotsz 1
۰/۰۳۲ .	۰/۰۳۵ .	۰/۰۴۱ .	۰/۰۳۶	lotsz 2
. -۲۰/۳۲۸				M-dwell 1
۸/۸۵۸				M-dwell 2
۲/۹۱۹				M-nbath 1
-۵/۳۸۷				M-nbath 2
۱۳/۷۵۹				M-patio 1
-۳/۰۹۰				M-patio 2

ادامه جدول

جدول ۶.۴ - ادامه جدول

SDM	SAC	SEM	SLM	ضرایب	
۲۲/۱۰۵				M-firepl 1	
۱۴/۹۰۵				M-firepl 2	
۱۵/۷۳۰				M-ac 1	
-۱/۷۷۷				M-ac 2	
-۴/۲۸۹				M-bment 1	
۳/۷۷۳				M-bment 2	
۷/۵۳۲				M-gar 1	
۹/۱۹۱				M-gar 2	
-۶/۴۵۷				M-citcou 1	
-۴/۲۸۳				M-citcou 2	
۰/۱۱۹				M-lotsz 1	
۰/۰۶۷				M-lotsz 2	
-۰/۱۱۰	***	۰/۳۵۴	***	۰/۳۲۷	ρ
-	-۰/۱۳۹	**	۰/۴۹۴	-	λ
۱۶۴۵/۳	۱۶۴۸/۳۰۰	۱۶۵۵/۱۶۶	۱۶۴۶/۴۸۶		AIC

در خصوص این نتایج می‌توان بیان کرد که

۱. در مدل‌های فراموضعی همه متغیرهای تبیینی دارای اثر معنی‌دار بر قیمت مسکن بودند؛ اما پس از اعمال رژیم‌های فضایی این اثرات تغییر کرده‌اند. اغلب متغیرهای تبیینی فقط در رژیم اول دارای اثر معنی‌دار هستند و در رژیم دوم هیچ تأثیری بر تعیین قیمت مسکن ندارند. این امر می‌توان ناشی از این حقیقت باشد که بخش مرکزی و شمال غربی شهر بالتیمور، منطقه‌ای مرفه‌نشین است و متغیرهای تجملی نظیر تعدد گاراژ و داشتن پاسیو برای افراد ساکن در آن حائز اهمیت است و موجب افزایش قیمت مسکن می‌شود.

۲. مقدار AIC در همه موارد در ESRMs، یعنی پس از اعمال رژیم‌های فضایی و مدل‌بندی ناهمگنی، کوچک‌تر شده است. یعنی مدل‌های اقتصادسنجی فضایی با اعمال رژیم‌های فضایی در مقایسه با مدل‌های مشابه فراموضعی خود، مدل‌های بهتری هستند.

۳. همچنین در این مثال در همه مدل‌ها پس از ورود رژیم‌های فضایی، ضرایب خودهمبستگی فضایی تغییر چندانی نکرده‌اند؛ یعنی اثرات خودهمبستگی واقعی بوده و تحت تاثیر ناهمگنی ضرایب بزرگ‌تر برآورد نشده‌اند.

۴. با مقایسه دو مدل SLM و SDM مشاهده می‌کنیم که ضریب خودهمبستگی فضایی با ورود متغیرهای تبیینی تأخیر فضایی به‌طور چشمگیری کاهش یافته است. یعنی مشابه آنچه در مورد مدل‌های فراموضعی گفته شد، بیشتر اثرات خودهمبستگی فضایی توسط متغیرهای تبیینی حذف شده ناشی از عدم توجه به همبستگی فضایی متغیرهای تبیینی به SLM القا شده است.

۵. همچنین علامت اثر خودهمبستگی نیز تغییر کرده است؛ بدون در نظر گرفتن رژیم‌ها در مدل SDM خودهمبستگی مثبت $0/327$ بین قیمت مسکن‌های شهر باتیمور برآورد شده بود، در حالی که این مقدار پس از اعمال رژیم‌ها به $0/110$ رسید. اما همان‌طور که در جدول ۶.۴ گزارش شد، این اثر معنی‌دار نیست.

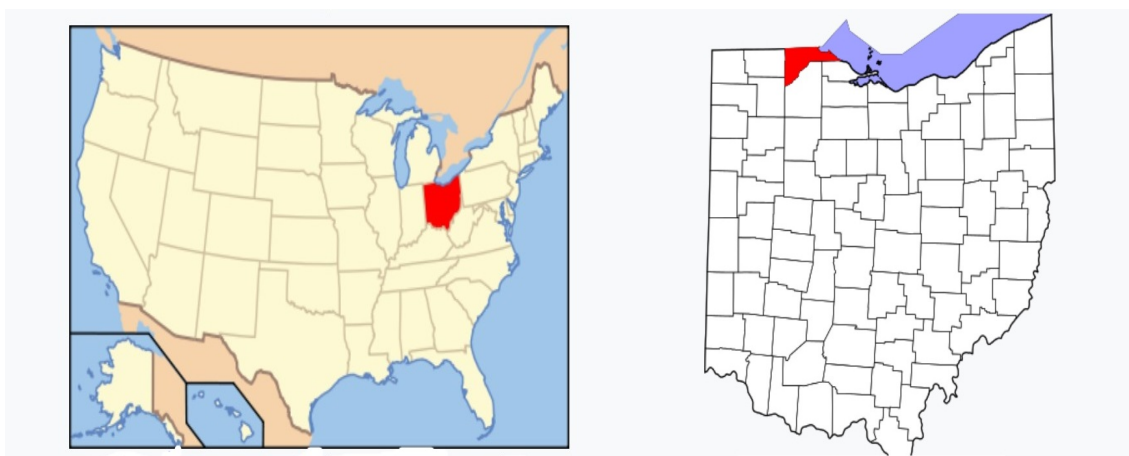
بنابر آنچه گفته شد، بهترین مدل برای داده‌های مسکن شهر بالتیمور مدل رگرسیون خطی با حضور متغیرهای تأخیر فضایی و اعمال رژیم‌های فضایی مشخص شده است؛ یعنی مدلی به‌صورت

$$\begin{pmatrix} y_1 \\ \vdots \\ y_c \end{pmatrix} = \begin{pmatrix} X_1 & \cdots & \mathbf{o} \\ \vdots & \ddots & \vdots \\ \mathbf{o} & \cdots & X_c \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_c \end{pmatrix} + M \begin{pmatrix} X_{v_1} & \cdots & \mathbf{o} \\ \vdots & \ddots & \vdots \\ \mathbf{o} & \cdots & X_{v_c} \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_c \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_c \end{pmatrix}$$

که در آن $\varepsilon_j \sim N(0, \sigma_{\varepsilon_j}^2 I_{n_j})$ و با ۳۸ متغیر تبیینی که در جدول ۶.۴ ردیف شده‌اند. مقدار AIC این مدل ۱۶۴۳/۷ است. برای درک بهتر این مدل را با مدل (۱۰.۳) مقایسه کنید.

۳.۴ مثال کاربردی دوم

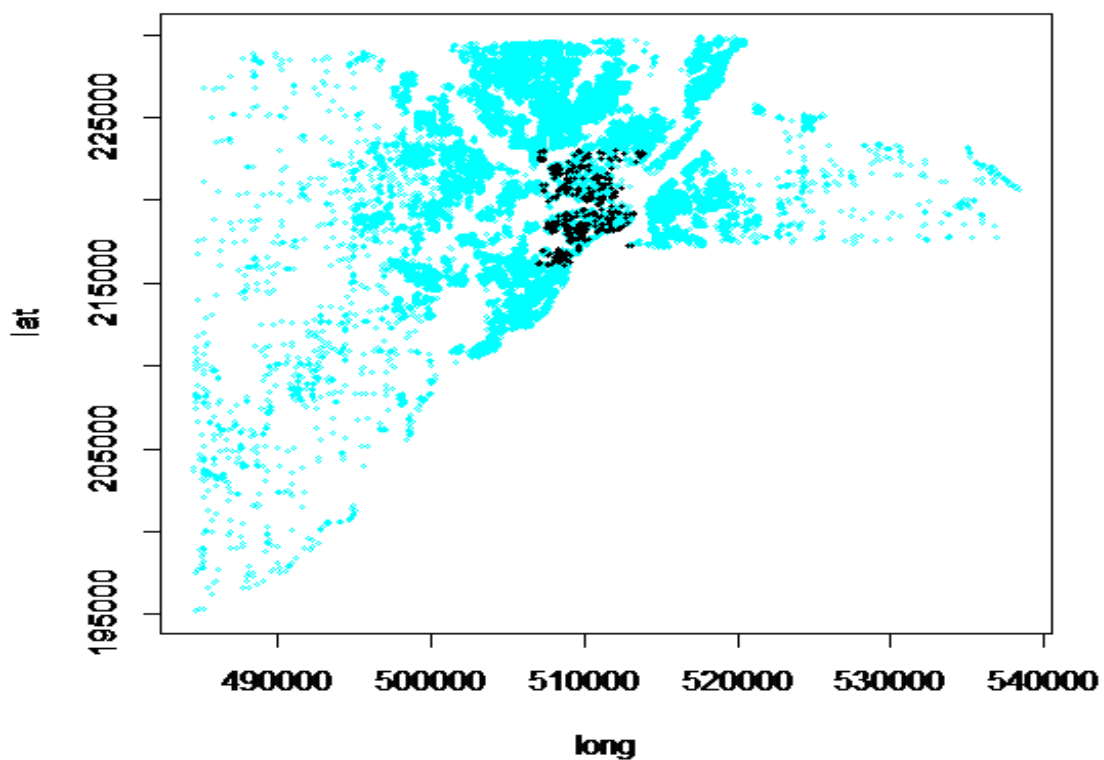
مجموعه داده مسکن نیز شامل اطلاعات قیمت فروش مسکن در بازه سال‌های ۱۹۹۳-۱۹۹۸ برای ۲۵۳۵۷ خانه در شهرستان لوکاس ایالت اوهایو کشور آمریکا است. شکل ۳.۴ موقعیت ایالت اوهایو و شهرستان لوکاس را نشان می‌دهد. بیل و همکاران (۲۰۱۷) به دلیل هزینه محاسبات الگوریتم IGWR تنها از خانه‌های بخش مرکزی لوکاس که در سال ۱۹۹۳ به فروش رفته‌اند (۳۸۲ مورد)، استفاده کردند. در این پایان‌نامه نیز پس از بررسی‌های لازم و مشخص کردن طول و عرض جغرافیایی بخش مرکزی شهرستان لوکاس، از خانه‌های واقع در آن منطقه برای تحلیل استفاده کردیم. شکل ۴.۴، کل مجموعه داده مسکن و داده‌های انتخابی (لوکاس) را نمایش می‌دهد. متغیرهای این مجموعه داده در جدول ۷.۴ معرفی شده‌اند.



(ب) موقعیت ایالت اوهایو در ایالات متحده امریکا

(آ) موقعیت در ایالت اوهایو

شکل ۳.۴: موقعیت شهرستان لوکاس



شکل ۴.۴: داده‌های انتخاب شده از مجموعه داده مسکن

جدول ۷.۴: متغیرهای داده لوکاس

متغیر	توضیحات
price	قیمت خانه - متغیر پاسخ
yrbuilt	سال ساخت
TLA	جابجایی محل زندگی
baths	تعداد حمام‌ها
halfbaths	تعداد حمام‌های کوچک
garagesqft	متراژ گاراژ به فوت مربع
lotsize	متراژ مسکن به فوت مربع*

* هر فوت مربع برابر 0.09290304 مترمربع

۱.۳.۴ انجام آزمون‌های آماری و برازش مدل‌های فراموضعی

برای این مجموعه داده نیز، ابتدا به بررسی وجود خودهمبستگی فضایی پرداختیم. جدول ۸.۴ نتایج را برای آزمون‌های مختلف خودهمبستگی فضایی نشان می‌دهد.

جدول ۸.۴: نتایج آزمون‌های خودهمبستگی فضایی برای داده‌های لوکاس

آزمون‌های خودهمبستگی فضایی	مقدار آماره آزمون	p-مقدار
موران I	۰/۳۶۷۴	< ۰/۰۰۰
نسبت درست‌نمایی	۱۱۲/۳۳	< ۰/۰۰۰
ضریب لاگرانژخطا	۲۸۷/۹۱	< ۰/۰۰۰
ضریب لاگرانژ تأخیر	۲۴۶/۲۵	< ۰/۰۰۰

همان‌طور که در جدول ۸.۴ نشان داده شده است، تمامی آزمون‌های خودهمبستگی معنادار بوده و با قدرت، حضور خودهمبستگی فضایی را در داده‌های لوکاس تایید می‌کنند. لذا مشابه داده‌های بالتیمور از مدل‌های اقتصادسنجی برای برازش استفاده می‌کنیم. برای این داده‌ها از سه مدل تأخیر فضایی، خطای فضایی و ترکیبی فضایی استفاده کردیم. ماتریس‌های وزن در هر سه مدل به روش 10^{-1} - نزدیکترین همسایگی محاسبه شدند. نتایج در جدول ۹.۴ گزارش شده‌اند.

جدول ۹.۴: برآورد پارامترهای مدل‌های فراموضعی برای مجموعه داده مسکن لوکاس

SDM	SEM	SLM	ضرایب
* -۳۰۵۰۵۶/۰۴ ***	-۴۳۸۷۳۷/۹۷ ***	-۵۵۲۵۹۹/۳۳ ***	عرض از مبدأ
*** ۲۰۹/۵۵ ***	۲۲۱/۷۲ ***	۲۷۲/۲۰ ***	yrbuilt
*** ۱۹/۶۴ ***	۲۰/۶۸ ***	۱۶/۴۳ ***	TLA
*** ۱۵۱۸۵/۰۰ ***	۱۴۵۷۸/۸۲ ***	۱۴۲۰۴/۶۴ ***	baths
*** ۵۱۹۳/۲۸ ***	۵۰۰۷/۱ ***	۴۹۷۶/۶۰ ***	halfbaths
** ۱۱/۹۸ ***	۱۱/۹۲ ***	۱۴/۰۳ ***	garagesqft
*** ۰/۱ ***	۱/۰۲ ***	۰/۹۶ ***	lotsize
-۵۵/۲۰			M-yrbuilt
*** -۱۸/۴۸			M-TLA
-۱۰۶۱۴/۳۶			M-baths
۵۸۳۷/۱۸			M-halfbaths
۳/۰۲			M-garagesqft
۰/۹۰			M-lotsize
*** ۰/۶۵۵	-	*** ۰/۵۸۰	ρ
-	*** ۰/۸۳۱	-	λ
۸۲۶۹/۵۴۹	۸۲۸۹/۳۳۳	۸۳۰۲/۷۱۹	AIC

معنی دار بودن: $< ۰/۱$. $< ۰/۰۵$ * $< ۰/۰۱$ ** $< ۰/۰۰۱$ ***

بنابر نتایج جدول ۹.۴، تمامی متغیرهای معرفی شده در جدول ۷.۴ دارای اثرات معنی دار مثبت بر تعیین قیمت مسکن در بخش مرکزی شهرستان لوکاس هستند. به عنوان مثال، ضرایب متغیر yrbuilt که نشان دهنده سال ساخت هر خانه است، در هر سه مدل مثبت و دارای اثر معنی دار است. یعنی افزایش مقدار این متغیر، که به معنی کاهش قدمت خانه و نوساز بودن آن است، موجب افزایش قیمت خانه می شود.

ضرایب خودهمبستگی در هر سه مدل مثبت و معنی دار هستند. از میان مدل‌های برازش داده شده، مدل دوربین فضایی کمترین مقدار AIC را داراست. در نتیجه بهترین مدل برای برازش به داده‌های مسکن لوکاس، SDM است.

همان طور که گفته شد، در این مدل‌ها به جنبه دوم ناهمگنی، یعنی تفاوت اثرات متغیرهای تبیینی وابسته به فضای مشاهدات، توجه نشده است.

۲.۳.۴ ناهمگنی در داده‌ها

مشابه مثال اول، آزمون بریچ-پیگان را برای همه مدل‌های فراموضعی فضایی انجام دادیم تا ناهمگنی فضایی را در حضور انواع مدل‌های خودهمبستگی آزمون کنیم.

جدول ۱۰.۴: نتایج آزمون‌های ناهمگنی فضایی داده‌های مسکن لوکاس

مدل‌های فضایی	آماره بریچ-پیگان	درجه آزادی	p-مقدار
SLM	۸۱/۴۲۵	۶	< ۰/۰۰۰۱
SEM	۷۲/۲۳۱	۶	< ۰/۰۰۰۱
SDM	۹۰/۲۱۵	۱۲	< ۰/۰۰۰۱

با توجه به نتایج جدول ۱۰.۴، در این داده‌ها نیز حضور ناهمگنی با قوت تایید می‌شود. برای درک بهتر حضور ناهمگنی فضایی، مدل رگرسیون موزون جغرافیایی را به داده‌ها برازش دادیم. نتایج رگرسیون موزون جغرافیایی در جدول ۱۱.۴ گزارش شده‌اند. ناهمگنی در داده‌ها و تغییر پارامترها وابسته به فضای مشاهدات، در متغیرهای baths با بیش از ۳۰۰۰۰ واحد تغییر و halfbaths با ۴۰۰۰۰ واحد تغییر، به وضوح قابل مشاهده است.

جدول ۱۱.۴: نتایج برازش مدل GWR برای داده‌های مسکن لوکاس

ضرایب	مینیمم	چارک اول	میانه	چارک سوم	ماکسیمم
Intercept	-۱۳۴۶۴۵۸/۷۶	-۸۸۳۲۴۰/۴۴	-۶۴۷۲۷۲/۴۲	-۳۶۴۵۵۴/۹۲	۶۰۱۴۵/۳۵
yrbuilt	-۱۸/۸۹	۱۸۷/۰۸	۳۳۷/۲۲	۴۴۹/۲۴	۶۹۰/۵۰
TLA	۵/۱۳	۱۱/۶۷	۲۰/۰۵	۲۹/۱۳	۵۳/۸۶
baths	-۴۴۸۴/۷۲	۲۰۱۷/۲۷	۷۳۷۹/۸۳	۱۵۲۷۷/۲۴	۲۷۸۲۴/۵۳
halfbaths	-۲۰۷۵۵/۰۷	-۱۶۶۸/۳۱	۲۲۶۹/۹۷	۱۱۶۹۵/۶۸	۱۹۳۱۲/۰۴
garagesqft	-۶/۷۱	۸/۵۳	۱۲/۴۲	۱۵/۷۵	۳۱/۲۰
lotsize	-۰/۶۴	۰۰/۶۶	۱/۰۷	۱/۶۶	۲/۶۰

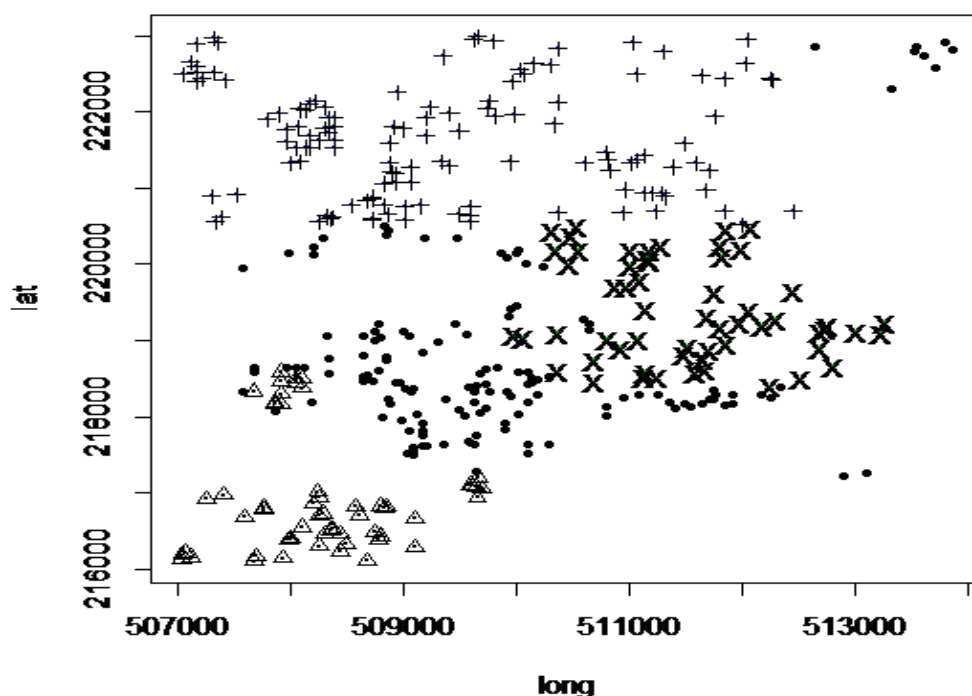
۳.۳.۴ برازش مدل پیشنهادی

بنابر نتایج به‌دست‌آمده، برای این داده‌ها نیز به دلیل حضور هم‌زمان ناهمگنی و خودهمبستگی فضایی، مدل OLS و مدل‌های فراموضعی فضایی مناسب نیستند؛ لذا از رهیافت پیشنهادی خود برای مدل‌بندی داده‌ها استفاده می‌کنیم.

نتایج رهیافت IGWR

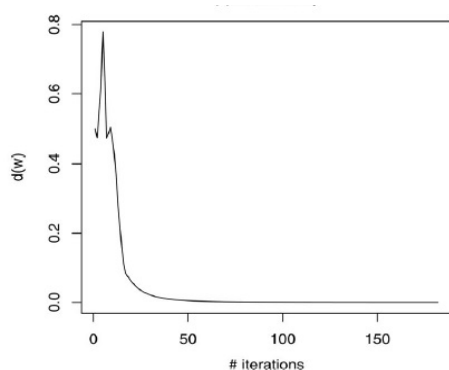
در این مثال نیز مشابه داده‌های بالتیمور، به منظور تعیین وزن‌های اولیه رگرسیون وزنی جغرافیایی از تابع هسته گاوسی با پهنای باند سازوار بر اساس معیار AIC استفاده شد. مقدار پهنای باند اولیه برای مجموعه داده مسکن لوکاس، $b_{Knn}^{AIC} = ۱۹$ محاسبه شد. الگوریتم به‌روزرسانی وزن‌ها در خصوص این داده‌ها پس از ۱۸۲ تکرار به همگرایی رسید و $c = ۴$ رژیم

(شکل ۵.۴) با زیرنمونه‌هایی به حجم‌های $n_1 = 57, n_2 = 161, n_3 = 58, n_4 = 126$ تشخیص داده شدند. در این مثال، پارامترهای الگوریتم به‌روزرسانی وزن‌ها مشابه داده‌های بالتیمور در نظر گرفته شدند ($\omega = 0.0001, \tau = 0.001, \eta = 0.5$). حجم محاسبات در این مثال در مقایسه با داده‌های بالتیمور تقریباً دو برابر بود. این افزایش دو برابری با توجه به حجم داده‌های مسکن لوکاس که تقریباً دو برابر داده‌های بالتیمور است، منطقی است.

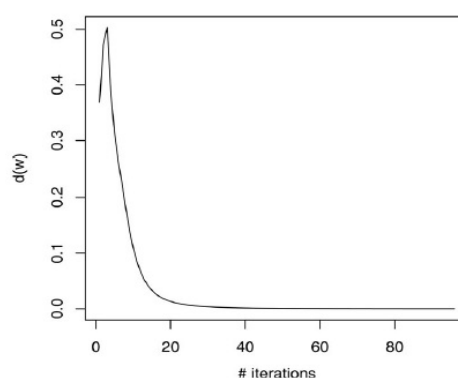


شکل ۵.۴: رژیم‌های فضایی مجموعه داده مسکن لوکاس، نقاط Δ رژیم فضایی اول، نقاط \bullet رژیم فضایی دوم، نقاط \times رژیم فضایی سوم و نقاط $+$ رژیم فضایی چهارم را نشان می‌دهند.

شکل ۶.۴ تابع تغییرات وزن‌ها ($d(w)$) را در تکرارهای مختلف نشان می‌دهد. در این نمودار رفتار فرایند همگرایی مشخص است. پس از یک دوره ناپایداری، اختلاف وزن‌ها با افزایش تکرار به صفر میل می‌کنند. بیشترین کاهش معنادار، در داده‌های بالتیمور در بازه تکرارهای ۱۵ - ۵ و در مجموعه مسکن لوکاس در تکرارهای ۳۰ - ۱۵ اتفاق افتاده است.



(ب) تابع تغییرات وزن مجموعه داده مسکن لوکاس



(آ) تابع تغییرات وزن مجموعه داده بالتیمور

شکل ۶.۴: تابع تغییرات وزن‌های روش IGWR

برازش مدل‌ها با رژیم فضایی درونی

در این بخش نیز مشابه مثال قبل، پس از مشخص شدن رژیم‌های فضایی، داده‌ها را به صورت (۷.۳) و (۹.۳) افراز کرده و مدل‌های تأخیر فضایی، خطای فضایی و دوربین فضایی با رژیم درونی را به داده‌ها برازش دادیم. ماتریس وزن W باز هم به روش ۱۰- نزدیکترین همسایگی محاسبه شد. نتایج در جدول ۱۲.۴ نمایش داده شدند.

جدول ۱۲.۴: برآورد پارامترهای مدل‌های با رژیم فضایی درونی برای مجموعه داده مسکن شهر لوکاس

SDM		SEM		SLM		ضرایب
***	-۱۲۴۵۷۹۷/۶۴	***	-۷۵۲۸۶۳/۴۹	**	-۵۰۷۵۰۹/۸۲	عرض از مبدأ ۱
***	-۱۱۷۸۰۰۵/۷۱	***	-۴۹۳۴۴۱/۷۳	***	-۶۲۶۷۹۶/۷۴	عرض از مبدأ ۲
**	-۸۴۱۰۲۶/۶۸		-۲۰۸۸۲۴/۶۲		-۱۵۱۹۱۲/۸۱	عرض از مبدأ ۳
***	-۱۳۲۱۰۳۷/۸۰	***	-۶۳۹۹۳۴/۵۳	***	-۸۵۹۸۰۲/۴۷	عرض از مبدأ ۴
**	۳۰۵/۶۰	***	۳۹۱/۳۲	*	۲۵۴/۱۱	yrbuilt 1
***	۲۷۷/۹۶	***	۲۶۱/۶۶	***	۳۲۴/۲۵	yrbuilt 2
	۱۰۲/۳۲		۱۱۱/۶۵		۷۶/۷۵	yrbuilt 3
***	۳۳۸/۴۶	***	۳۲۱/۱۸	***	۴۲۸/۹۰	yrbuilt 4
.	۱۳/۸۰	*	۱۲/۹۷	*	۱۲/۲۷	TLA 1
***	۱۲/۴۷	***	۱۳/۱۸	***	۱۱/۷۰	TLA 2
	۷/۶۰	.	۷/۹۲	.	۸/۵۱	TLA 3
***	۲۷/۷۳	***	۲۸/۲۲	***	۲۳/۹۵	TLA 4

ادامه جدول

جدول ۱۲.۴ - ادامه جدول

SDM		SEM		SLM		ضرایب
***	۲۲۷۱۷/۹۹	**	۲۰۱۳۴/۱۱	**	۲۲۸۴۱/۲۳	baths 1
	۳۷۶۰/۲۴		۱۱۵۰/۷۵		۲۹۸۸/۹۱	baths 2
	۵۸۴۰/۱۹		۲۸۲۶/۹۵		۸۳۰/۷۶	baths 3
***	۱۲۹۱۷/۰۴	***	۱۱۱۲۹/۵۹	*	۱۱۲۶۲/۴۹	baths 4
	-۲۴۹۴/۸۵		-۱۵۳۴/۵۲		۱۹۴۱/۸۹	halfbaths 1
	۲۰۶۹/۱۱		۳۸۶/۶۰		۲۶۷۶/۹۱	halfbaths 2
	۱۱۱۵/۱۴		-۴۸۲۸/۶۹		-۱۶۹۷/۷۲	halfbaths 3
***	۷۵۱۸/۸۰	***	۶۱۹۴/۴۳	***	۸۱۹۴/۴۳	halfbaths 4
	۸/۴۴	.	۲۰/۹۳		۱۵/۸۶	garagesqft 1
	۷/۵۳	*	۱۰/۰۴	*	۱۱/۷۳	garagesqft 2
.	۱۳/۶۶		۹/۸۴	.	۱۲/۸۳	garagesqft 3
	۷/۹۰	*	۱۲/۸۲	.	۱۰/۳۸	garagesqft 4
.	۰/۹۴	**	۱/۰۸	*	۱/۰۳	lotsize 1
	۰/۶۸	*	۰/۹۳		۰/۵۵	lotsize 2
	۰/۰۸		۱/۲۴		۰/۱	lotsize 3
**	۱/۰۳	***	۱/۳۵	*	۱/۰۲	lotsize 4
**	۳۱۵/۷۰					M-yrbuilt 1
***	۳۴۳/۱۰					M-yrbuilt 2
**	۳۴۴/۲۴					M-yrbuilt 3
**	۳۴۱/۱۴					M-yrbuilt 4
	۴۶/۲۶					M-TLA 1
	-۷/۹۴					M-TLA 2
	-۹/۵۶					M-TLA 3
	-۱۱/۳۴					M-TLA 4
	۴۶۴۲۴/۶۰					M-baths 1
	۸۵۶۸/۶۴					M-baths 2
	۱۵۱۵۹/۶۸					M-baths 3
	-۳۵۵۶/۳۹					M-baths 4
	۱۹۴۷۲/۴۸					M-halfbaths 1
	۳۰۱۴/۳۲					M-halfbaths 2
	۲۹۴۹۷/۴۲					M-halfbaths 3

ادامه جدول

جدول ۱۲.۴ - ادامه جدول

	SDM	SEM	SLM	ضرایب
***	۲۲۰۵۶/۶۰			M-halfbaths 4
**	-۱۶۶/۲۱			M-garagesqft 1
	-۱۴/۸۰			M-garagesqft 2
*	۶۳/۱۱			M-garagesqft 3
.	-۲۷/۵۶			M-garagesqft 4
.	-۲/۱۷			M-lotsize 1
	-۱/۴۳			M-lotsize 2
	-۸/۴۸			M-lotsize 3
	۰/۲۶			M-lotsize 4
**	۰/۲۵	-	*** ۰/۴۴۴	ρ
		*** ۰/۷۹۱	-	λ
		۸۲۴۶/۱۳۲	۸۲۵۲/۴۱۷	AIC

معنی دار بودن: $< ۰/۱$. $< ۰/۰۵$ * $< ۰/۰۱$ ** $< ۰/۰۰۱$ ***

در مقایسه نتایج مدل‌های فراموضعی و مدل‌ها با اعمال رژیم فضایی، مشاهده می‌کنیم که مقدار AIC در همه موارد در ESRMs کمتر است. در مورد داده‌های مسکن لوکاس نه تنها نتایج ESRMs از نظر مقدار AIC بهتر شده‌اند، بلکه در آن‌ها اثرات سرریز فضایی ρ نیز رقیق شده‌اند. این کاهش در مدل دوربین فضایی از ۰/۶۵۵ به ۰/۲۴۸ کاملاً مشهود است. این نتیجه به این معنی است که اطلاعات اضافه‌شده توسط رژیم‌های فضایی ما را از اثرات سرریز فضایی آشکار به اثرات سرریز فضایی واقعی می‌رساند.

به عنوان نتیجه می‌توان گفت مدل دوربین فضایی با رژیم فضایی درونی نسبت به باقی مدل‌ها از قدرت بیشتری برخوردار است؛ زیرا طوری طراحی شده است که اثرات آشکار فضایی را به ۳ قسمت: اثرات سرریز فضایی واقعی، متغیرهای همبسته فضایی تبیینی و ناهمگنی فضایی، مدل‌بندی می‌کند.

۴.۴ نتیجه و آینده تحقیق

در این پایان‌نامه یک روش دو مرحله‌ای، که در مرحله اول آن رژیم‌های فضایی تعیین و در مرحله دوم ناهمگنی (تحت عنوان رژیم‌های فضایی) و وابستگی فضایی به‌طور هم‌زمان برآورد شدند، معرفی شد. با استفاده از دو مجموعه داده واقعی مربوط به تعیین قیمت مسکن، نشان دادیم که این روش می‌تواند نواحی همگن را شناسایی کند. مدل‌های اقتصادسنجی

فضایی که با اعمال رژیم‌بندی‌ها برآورد شدند نسبت به حالت فراموضعی خود ارجح‌تر هستند. همچنین اکثر نتایج نشان داد که اثرات خودهمبستگی فضایی در داده‌ها وجود دارد اما این اثرات کمتر از مقداری است که در مدل‌های فراموضعی برآورد می‌شود. این اختلاف در برآورد اثر خودهمبستگی، علاوه بر تأثیر متغیرهای حذف‌شده، ناشی از عدم توجه به ناهمگنی فضایی است. بنابراین برای برآورد یک مدل فضایی منعطف، مدل‌بندی هم‌زمان وابستگی و ناهمگنی یک گام اساسی است.

در مرحله اول این رهیافت برای تعیین رژیم‌ها هیچ محدودیتی لحاظ نشده است و به اطلاعات پیشین در مورد ساختار فضایی داده‌ها، که برای اغلب پدیده‌ها در دسترس نیست، نیازی ندارد. لذا به راحتی برای هر مدلی که امکان برآورد موضعی آن توسط کمترین توان‌های دوم موزون وجود دارد، قابل تعمیم است. یک تغییر شدنی و جالب، تعمیم این رهیافت برای مجموعه داده‌های پانلی به منظور لحاظ کردن تغییرات زمانی (علاوه بر ساختار فضایی) است.

پیوست آ

برنامه‌های رایانه‌ای با نرم‌افزار R

۱.آ برنامه‌های داده بالتیمور

دستورات لازم برای تعیین ماتریس وزن W و انجام آزمون‌های خودهمبستگی فضایی

```
library(spdep)
data(baltimore)
View(baltimore)
plot(baltimore$X,baltimore$Y,cex=1,pch=19,xlab=long,ylab=lat)

coords<-cbind(baltimore$X,baltimore$Y)
knn <- knn2nb(knearneigh(coords,k=10))
listw_10nn_dates <- nb2listw(knn)
OLS <- lm(PRICE~DWELL+NBATH+PATIO+FIREPL+AC+BMENT+GAR+CITCOU+LOTSZ,
data=baltimore)
summary(OLS)

balt.moran <- lm.morantest(OLS,listw_10nn_dates)
```



```
balt.moran
```

```
lmtest<-lm.LMtests(OLS,listw_10nn_dates,test=c("LMerr","LMlag"))
```

```
lmtest
```

برازش مدل‌های اقتصادسنجی فضایی

```
SLM <-lagsarlm(PRICE~DWELL+NBATH+PATIO+FIREPL+AC+BMENT+GAR+CITCOU+LOTSZ,  
data=baltimore, listw_10nn_dates, type="lag")
```

```
summary(SLM)
```

```
SDM <- lagsarlm(PRICE~DWELL+NBATH+PATIO+FIREPL+AC+BMENT+GAR+CITCOU+LOTSZ,  
data=baltimore, listw_10nn_dates, type="durbin")
```

```
summary(SDM)
```

```
SEM <- errorsarlm(PRICE~DWELL+NBATH+PATIO+FIREPL+AC+BMENT+GAR+CITCOU+LOTSZ,  
data=baltimore, listw_10nn_dates, method="eigen", quiet=FALSE)
```

```
summary(SEM)
```

```
SAC <- sacsarlm(PRICE~DWELL+NBATH+PATIO+FIREPL+AC+BMENT+GAR+CITCOU+LOTSZ,  
data=baltimore, listw_10nn_dates)
```

آزمون ناهمگنی در حضور خودهمبستگی فضایی (آزمون بریچ-پیگان)

```
bptest.sarlm(SLM)
```

```
bptest.sarlm(SDM)
```

```
bptest.sarlm(SEM)
```

```
bptest.sarlm(SAC)
```

دستورات لازم برای برازش مدل رگرسیون موزون جغرافیایی با پهنای باند بهینه با استفاده
از معیار اطلاع آکاییک

```
library(GWmodel)
```

```

dmat <- gw.dist(coords,focus=0,p=2,theta=0,longlat=F)
bw <- bw.gwr(PRICE~DWELL+NBATH+PATIO+FIREPL+AC+BMENT+GAR+CITCOU+LOTSZ,
data=SpatialPointsDataFrame(coords,baltimore),approach="AIC",kernel="gaussian",
adaptive=T,p=2,theta=0,longlat=F,dMat=dmat)      # 33 nearneighbor
gwr.res <- gwr.basic(PRICE~DWELL+NBATH+PATIO+FIREPL+AC+BMENT+GAR+CITCOU+LOTSZ,
data=SpatialPointsDataFrame(coords,baltimore), bw=bw,
kernel = "gaussian",adaptive=TRUE, dMat=dmat)
gwr.res

```

دستورات رهیافت IGWR

```

library(spdep)
library(GWmodel)
coords<-cbind(baltimore$X,baltimore$Y)
dmat<-gw.dist(coords,focus=0,p=2,theta=0,longlat=F)
bw<-bw.gwr(PRICE~DWELL+NBATH+PATIO+FIREPL+AC+BMENT+GAR+CITCOU+LOTSZ,
data=SpatialPointsDataFrame(coords,baltimore),approach="AIC",kernel="gaussian",
adaptive=T,p=2,theta=0,longlat=F,dMat=dmat)
y<-as.matrix(baltimore$PRICE)
x<-cbind(baltimore$DWELL,baltimore$NBATH,baltimore$PATIO,baltimore$FIREPL,
baltimore$AC,baltimore$BMENT,baltimore$GAR,baltimore$CITCOU,baltimore$LOTSZ)

modparcov <- function (y, x, w)      #GWR/LWR
{
  xpw      <- t(x)%*%w                #w matrice diagonale
  invxpwx  <- solve(xpw)%*%x
  b        <- invxpwx%*%(xpw)%*%y     #kx1 weighted least squares
  resid    <- y - x %*% b             #nx1
  qresid   <- t(resid)%*%w%*%resid
  c        <- invxpwx%*%xpw
  s        <- x%*%c
  v1       <- sum(diag(s))
  v2       <- sum(diag(t(s)%*%s))
  varresid <- qresid/(nrow(x)-2*v1+v2) #sigma2
  ccp      <- c%*%t(c)
}

```

```

res      <- list()
res$par  <- b                               #kx1
res$cov  <- ccp*as.numeric(varresid)        #kxk
res
}

testpar <- function (p1, p2, cov1, cov2, w1, w2) #wald test
{
  n1      <- as.numeric(sum(diag(w1)))
  n2      <- as.numeric(sum(diag(w2)))
  pcov    <- (n1*cov1+n2*cov2)/(n1+n2)       #pooled var-cov beta
  t       <- t(p1-p2)%*%solve(pcov)%*%(p1-p2)
  t
}

awsreg <- function (y,x,coords,bw,tau,niter,conv,eta,numout,sout) #conv=omega
{
  x       <- as.matrix(x)
  y       <- as.matrix(y)
  codic   <- 1:nrow(x)
  clas    <- rep(0,nrow(x))
  woutl   <- rep(0,nrow(x))

  vdist   <- as.matrix(gw.dist(coords,focus=0,p=2,theta=0,longlat=F)) #dmat
  knn     <- knn2nb(knearneigh(coords,k=bw))
  nbdist  <- nbdists(knn,coords)

  #w<-matrix(0,nrow(vdist),ncol(vdist))
  #for (i in 1:nrow(vdist)){
  #for (j in 1:ncol(vdist)){
  #  w[i,j] <- exp(-vdist[i,j]/max(nbdist[[i]]))      #esponenziale
  #}}

  w<-matrix(0,nrow(vdist),ncol(vdist))

```

```
for (i in 1:nrow(vdist)){
  for (j in 1:ncol(vdist)){
    w[i,j] <- exp(-0.5*(vdist[i,j]/max(nbdist[[i]]))^2) #gaussiana
  }}

#w<-matrix(0,nrow(vdist),ncol(vdist)) #bisquare
#for (i in 1:nrow(vdist)){
#for (j in 1:ncol(vdist)){
#  if (vdist[i,j]<max(nbdist[[i]]))
# w[i,j] <- (1-(vdist[i,j]/max(nbdist[[i]]))^2)^2
#  else      w[i,j] <- 0.00001
#}}

#w<-matrix(0,nrow(vdist),ncol(vdist)) #tricube
#for (i in 1:nrow(vdist)){
#for (j in 1:ncol(vdist)){
#  if (vdist[i,j]<max(nbdist[[i]]))
#w[i,j] <- (1-(vdist[i,j]/max(nbdist[[i]]))^3)^3
#  else      w[i,j] <- 0.00001
#}}

wnew <- matrix(0,nrow(x),nrow(x))
parm <- matrix(0,nrow(x),ncol(x)) #matrice parametri nxk
gresid <- rep(0,length=nrow(x))
listcov <- list()

nout <- 0
j <- 0
differ <- 1
while (differ > conv && (j <- j+1)<=niter)
{
  listw<-list()
  for (i in 1:ncol(w))
  {
```

```

listw[[i]]<-diag(w[,i]) #w[i,] pesi x riga
}
for (i in 1:nrow(x))
{
  wls      <- modparcov(y, x, listw[[i]])
  parm[i,] <- wls$par
  listcov[[i]]<- wls$cov
}
for (u in 1:(nrow(x)-1))          #test
{
  wnew[u,u] <- 1
  for (v in (u+1):nrow(x))
  {
    if (w[u,v] < conv)
    {
      wnew[u,v] <- 0
      wnew[v,u] <- wnew[u,v]
    }
    else
    {
      #wnew[u,v] <- exp(-testpar(parm[u,],parm[v,],
      #listcov[[u]],listcov[[v]],listw[[u]],listw[[v]])*tau)          #esponenziale
      wnew[u,v] <- exp(-0.5*(testpar(parm[u,],parm[v,],
      listcov[[u]],listcov[[v]],listw[[u]],listw[[v]])*tau)^2) #gaussiana
      wnew[v,u] <- wnew[u,v]
    }
  }
}
wnew[nrow(x),nrow(x)] <- 1

#wnew2<-matrix(0,nrow(vdist),ncol(vdist))          #esponenziale
#for (i in 1:nrow(vdist)){
#for (h in 1:ncol(vdist)){
#  wnew2[i,h]<-wnew[i,h]*exp(-vdist[i,h]/(max(nbdist[[i]])*j))

```

```

#}}
#wnew<-wnew2

wnew2<-matrix(0,nrow(vdist),ncol(vdist))           #gaussiana
for (i in 1:nrow(vdist)){
  for (h in 1:ncol(vdist)){
    wnew2[i,h]<-wnew[i,h]*exp(-0.5*(vdist[i,h]/(max(nbdist[[i]])*j))^2)
  }}
wnew<-wnew2

differ <- max(abs(w-wnew))
w      <- eta*w+(1-eta)*wnew

for (i in 1:nrow(x))           #controllo outlier
{
  wcontr    <- w[,i]
  woutl[i]  <- length(wcontr[wcontr > sout])
}
selez <- which(woutl[1:nrow(x)] < numout)
if (length(selez) > 0)
{
  nout    <- nout + length(selez)
  x      <- x[-selez,]
  y      <- y[-selez,]
  codic  <- codic[-selez]
  vdist  <- vdist[-selez,-selez]
  w      <- w[-selez,-selez]
  wnew   <- wnew[-selez,-selez]
}
print(c(j,differ,nout))
}
ncl <- 0
for (i in 1:(nrow(x)-1))
{

```

```

if (clas[codic[i]] < 1)
{
  ncl      <- ncl + 1
  clas[codic[i]] <- ncl
  for (j in (i+1):nrow(x))
  {
    if (w[i,j] > 0.9) clas[codic[j]] <- clas[codic[i]]
  }
}
}

if (clas[nrow(x)] < 1) clas[nrow(x)] <- ncl + 1
clas[clas == 0] <- NA
clas
}

#          alfa/2
tau<-qchisq(p=0.100,df=1) #0.015
tau<-qchisq(p=0.050,df=1) #0.004
tau<-qchisq(p=0.025,df=1) #0.001 tau alfa=0.05

system.time(
kk <- awsreg(y,x,coords=coords,bw=bw,tau=0.001,
niter=200,conv=0.0001,eta=0.5,numout=20,sout=1e-05) #eta=memory parameter
)
kk

kk[is.na(kk)] <- max(kk,na.rm = TRUE)+1

# if more than 8 groups (no missing)
kk1<-kk
kk1[kk>8]<-0
plot(baltimore$X,baltimore$Y,cex=1,pch=20,col=kk1)
kk2<-kk
kk2[kk<9]<-0

```

```
points(baltimore$X,baltimore$Y,cex=1,pch=12,col=kk2)
plot(baltimore$X,baltimore$Y,cex=1,pch=1,col=kk,xlab="long",ylab="lat")
```

برنامه‌های لازم برای افراز داده‌ها و رژیم‌بندی آنها

```
group1 = baltimore[kk==1,]
group1x<-cbind(group1$DWELL,group1$NBATH,group1$PATIO,group1$FIREPL,
group1$AC,group1$BMENT,group1$GAR,group1$CITCOU,group1$LOTSZ) # X regim1

group2 = baltimore[kk==2,]
group2x<-cbind(group2$DWELL,group2$NBATH,group2$PATIO,group2$FIREPL,
group2$AC,group2$BMENT,group2$GAR,group2$CITCOU,group2$LOTSZ) # X regim2

newy = cbind(t(baltimore$PRICE[kk==1]),t(baltimore$PRICE[kk==2]))
newy = t(newy)
newx = rbind(cbind(1,group1x,matrix(0,nr=nrow(group1x),nc=ncol(group2x)+1)),
cbind(matrix(0,nr=nrow(group2x),nc=ncol(group1x)+1),1,group2x))

newbalt = cbind(newy,newx) #new baltimore data set

newbalt = data.frame(newbalt)

names(newbalt) = c("PRICE","intercept1","DWELL1","NBATH1","PATIO1",
"FIREPL1","AC1","BMENT1","GAR1","CITCOU1","LOTSZ1","intercept2",
"DWELL2","NBATH2","PATIO2","FIREPL2","AC2",
"BMENT2","GAR2","CITCOU2","LOTSZ2")
```

دستورات لازم برای معرفی ماتریس وزن جدید و برازش مدل‌های اقتصادسنجی با رژیم فضایی درونی

```
coordinat1 = cbind(group1$X,group1$Y)
coordinat2 = cbind(group2$X,group2$Y)
coords1 = rbind(coordinat1,coordinat2) #new coordinat for baltimor data set
knn1 <- knn2nb(knearneigh(coords1,k=10))
```



```
listw_10nn_dates1 <- nb2listw(knn1)
```

```
SLM.ES <- lagsarlm(PRICE~-1+intercept1+intercept2+DWELL1+DWELL2+
NBATH1+NBATH2+PATIO1+PATIO2+FIREPL1+FIREPL2+AC1+AC2+
BMENT1+BMENT2+GAR1+GAR2+CITCOU1+CITCOU2+LOTSZ1+LOTSZ2,
data=newbalt, listw_10nn_dates1, type="lag")
summary(SLM.ES)
```

```
SDM.ES <- lagsarlm(PRICE~-1+intercept1+intercept2+DWELL1+DWELL2+
NBATH1+NBATH2+PATIO1+PATIO2+FIREPL1+FIREPL2+
AC1+AC2+BMENT1+BMENT2+GAR1+GAR2+CITCOU1+CITCOU2+LOTSZ1+LOTSZ2,
data=newbalt, listw_10nn_dates1, Durbin=~DWELL1+DWELL2+NBATH1+NBATH2+
PATIO1+PATIO2+FIREPL1+FIREPL2+AC1+AC2+BMENT1+
BMENT2+GAR1+GAR2+CITCOU1+CITCOU2+LOTSZ1+LOTSZ2)
summary(SDM.ES)
```

```
SEM.ES <- errorsarlm(PRICE~-1+intercept1+intercept2+DWELL1+DWELL2+
NBATH1+NBATH2+PATIO1+PATIO2+FIREPL1+FIREPL2+AC1+AC2+
BMENT1+BMENT2+GAR1+GAR2+CITCOU1+CITCOU2+LOTSZ1+LOTSZ2,
data=newbalt, listw_10nn_dates1, method="eigen", quiet=FALSE)
summary(SEM.ES)
```

```
SAC.ES <- sacsarlm(PRICE~-1+intercept1+intercept2+DWELL1+DWELL2+
NBATH1+NBATH2+PATIO1+PATIO2+FIREPL1+FIREPL2+AC1+AC2+
BMENT1+BMENT2+GAR1+GAR2+CITCOU1+CITCOU2+LOTSZ1+LOTSZ2,
data=newbalt, listw_10nn_dates1)
summary(SAC.ES)
```

۲.آ برنامه‌های داده مسکن

دستورات لازم برای انتخاب داده‌های مسکن بخش مرکزی شهرستان لوکاس که در سال ۱۹۹۳ به فروش رسیده‌اند، تعیین ماتریس وزن W و انجام آزمون‌های خودهمبستگی فضایی

```
library(spdep)
data(house)
plot(coordinates(house),cex=.4,pch=19)

subhouse<-subset(house,house$long>507000 & house$long<514000
& house$lat>216000 & house$lat<223000 & house$s1993==1)

points(coordinates(subhouse),cex=.4,pch=19,col="blue")

coords<- coordinates(subhouse)
knn <- knn2nb(knearneigh(coords,k=10))
listw_10nn_dates <- nb2listw(knn)
OLS <- lm(price~yrbuilt+TLA+baths+halfbaths+garagesqft+lotsize,
data=subhouse)
summary(OLS)

balt.moran <- lm.morantest(OLS,listw_10nn_dates)
balt.moran

lmtest<-lm.LMtests(OLS,listw_10nn_dates,test=c("LMerr","LMlag"))
lmtest

برازش مدل‌های اقتصادسنجی فضایی

SLM <- lagsarlm(price~yrbuilt+TLA+baths+halfbaths+garagesqft+lotsize,
data=subhouse, listw_10nn_dates, type="lag",tol.solve = 1.39885e-19)
summary(SLM)

SDM <- lagsarlm(price~yrbuilt+TLA+baths+halfbaths+garagesqft+lotsize,
data=subhouse, listw_10nn_dates, type="mixed",tol.solve = 1.39885e-19)
```

```
summary(SDM)
```

```
SEM <- errorsarlm(price~yrbuilt+TLA+baths+halfbaths+garagesqft+lotsize,
data=subhouse, listw_10nn_dates, method="eigen",
quiet=FALSE,tol.solve = 1.39885e-19)
summary(SEM)
```

آزمون ناهمگنی در حضور خودهمبستگی فضایی (آزمون بریچ-پیگان)

```
bptest.sarlm(SLM)
```

```
bptest.sarlm(SEM)
```

```
bptest.sarlm(SDM)
```

دستورات لازم برای برازش مدل رگرسیون موزون جغرافیایی با پهنای باند بهینه با استفاده از معیار اطلاع آکاییک

```
library(GWmodel)
dmat<-gw.dist(coords,focus=0,p=2,theta=0,longlat=F)
bw<-bw.gwr(price~yrbuilt+TLA+baths+halfbaths+garagesqft+lotsize,
data=SpatialPointsDataFrame(coords,subhouse@data),
approach="AIC",kernel="gaussian",adaptive=T,p=2,
theta=0,longlat=F,dMat=dmat) #19 nearneighbor

gwr.res1<-gwr.basic(price~yrbuilt+TLA+baths+halfbaths+garagesqft+lotsize,
data=SpatialPointsDataFrame(coords,subhouse@data), bw=bw,
kernel = "gaussian",adaptive=TRUE, dMat=dmat)
gwr.res1
```

دستورات رهیافت IGWR

```
library(spdep)
```

```
library(GWmodel)
```

```
dmat<-gw.dist(coords,focus=0,p=2,theta=0,longlat=F)
```

```

bw<-bw.gwr(price~yrbuilt+TLA+baths+halfbaths+garagesqft+lotsize,
data=SpatialPointsDataFrame(coords,subhouse@data),
approach="AIC",kernel="gaussian",adaptive=T,p=2,
theta=0,longlat=F,dMat=dmat) #19 nearneighbor

y<-as.matrix(subhouse$price)
x<-cbind(subhouse$yrbuilt,subhouse$TLA,subhouse$baths,
subhouse$halfbaths,subhouse$garagesqft,subhouse$lotsize)

modparcov <- function (y, x, w) #GWR/LWR
{
  xpw <- t(x)%*%w #w matrice diagonale
  invxpx <- solve(xpw%*%x)
  b <- invxpx%*%(xpw%*%y) #kx1 weighted least squares
  resid <- y - x %*% b #nx1
  qresid <- t(resid)%*%w%*%resid
  c <- invxpx%*%xpw
  s <- x%*%c
  v1 <- sum(diag(s))
  v2 <- sum(diag(t(s)%*%s))
  varresid <- qresid/(nrow(x)-2*v1+v2) #sigma2
  ccp <- c%*%t(c)
  res <- list()
  res$par <- b #kx1
  res$cov <- ccp*as.numeric(varresid) #kxk
  res
}

testpar <- function (p1, p2, cov1, cov2, w1, w2) #wald test
{
  n1 <- as.numeric(sum(diag(w1)))
  n2 <- as.numeric(sum(diag(w2)))
  pcov <- (n1*cov1+n2*cov2)/(n1+n2) #pooled var-cov beta
}

```

```

t      <- t(p1-p2)%*%solve(pcov)%*%(p1-p2)
t
}

awsreg <- function (y,x,coords,bw,tau,niter,conv,eta,numout,sout) #conv=omega
{
  x      <- as.matrix(x)
  y      <- as.matrix(y)
  codic  <- 1:nrow(x)
  clas   <- rep(0,nrow(x))
  woutl  <- rep(0,nrow(x))

  vdist  <- as.matrix(gw.dist(coords,focus=0,p=2,theta=0,longlat=F)) #dmat
  knn     <- knn2nb(knearneigh(coords,k=bw))
  nbdist  <- nbdists(knn,coords)

  #w<-matrix(0,nrow(vdist),ncol(vdist))
  #for (i in 1:nrow(vdist)){
  #for (j in 1:ncol(vdist)){
  #  w[i,j] <- exp(-vdist[i,j]/max(nbdist[[i]]))      #esponenziale
  #}}

  w<-matrix(0,nrow(vdist),ncol(vdist))
  for (i in 1:nrow(vdist)){
  for (j in 1:ncol(vdist)){
    w[i,j] <- exp(-0.5*(vdist[i,j]/max(nbdist[[i]]))^2) #gaussiana
  }}

  #w<-matrix(0,nrow(vdist),ncol(vdist))      #bisquare
  #for (i in 1:nrow(vdist)){
  #for (j in 1:ncol(vdist)){
  #  if (vdist[i,j]<max(nbdist[[i]]))
  #w[i,j] <- (1-(vdist[i,j]/max(nbdist[[i]]))^2)^2
  #  else          w[i,j] <- 0.00001

```

```
#}}

#w<-matrix(0,nrow(vdist),ncol(vdist))           #tricube
#for (i in 1:nrow(vdist)){
#for (j in 1:ncol(vdist)){
#  if (vdist[i,j]<max(nbdist[[i]]))
#w[i,j] <- (1-(vdist[i,j]/max(nbdist[[i]]))^3)^3
#  else                                     w[i,j] <- 0.00001
#}}

wnew <- matrix(0,nrow(x),nrow(x))
parm <- matrix(0,nrow(x),ncol(x))
qresid <- rep(0,length=nrow(x))
listcov <- list()

nout <- 0
j <- 0
differ <- 1
while (differ > conv && (j <- j+1)<=niter)
{
  listw<-list()
  for (i in 1:ncol(w))
  {
    listw[[i]]<-diag(w[,i])
  }
  for (i in 1:nrow(x))
  {
    wls <- modparcov(y, x, listw[[i]])
    parm[i,] <- wls$par
    listcov[[i]]<- wls$cov
  }
  for (u in 1:(nrow(x)-1))           #test
  {
    wnew[u,u] <- 1
  }
}
```

```

for (v in (u+1):nrow(x))
{
  if (w[u,v] < conv)
  {
    wnew[u,v] <- 0
    wnew[v,u] <- wnew[u,v]
  }
  else
  {
    #wnew[u,v] <- exp(-testpar(parm[u,],parm[v,],listcov[[u]],
    #listcov[[v]],listw[[u]],listw[[v]])*tau)          #esponenziale

    wnew[u,v] <- exp(-0.5*(testpar(parm[u,],parm[v,],listcov[[u]],
    listcov[[v]],listw[[u]],listw[[v]])*tau)^2) #gaussiana
    wnew[v,u] <- wnew[u,v] #per simmetria
  }
}
}
wnew[nrow(x),nrow(x)] <- 1

#wnew2<-matrix(0,nrow(vdist),ncol(vdist))          #esponenziale
#for (i in 1:nrow(vdist)){
#for (h in 1:ncol(vdist)){
#  wnew2[i,h]<-wnew[i,h]*exp(-vdist[i,h]/(max(nbdist[[i]])*j))
#}}
#wnew<-wnew2

wnew2<-matrix(0,nrow(vdist),ncol(vdist))          #gaussiana
for (i in 1:nrow(vdist)){
for (h in 1:ncol(vdist)){
  wnew2[i,h]<-wnew[i,h]*exp(-0.5*(vdist[i,h]/(max(nbdist[[i]])*j))^2)
}}
wnew<-wnew2

```

```

differ <- max(abs(w-wnew))
w      <- eta*w+(1-eta)*wnew

for (i in 1:nrow(x))          #controllo outlier
{
  wcontr    <- w[,i]
  woutl[i]  <- length(wcontr[wcontr > sout])
}
selez <- which(woutl[1:nrow(x)] < numout)
if (length(selez) > 0)
{
  nout    <- nout + length(selez)
  x      <- x[-selez,]
  y      <- y[-selez,]
  codic  <- codic[-selez]
  vdist  <- vdist[-selez,-selez]
  w      <- w[-selez,-selez]
  wnew   <- wnew[-selez,-selez]
}
print(c(j,differ,nout))
}
ncl <- 0
for (i in 1:(nrow(x)-1))
{
  if (clas[codic[i]] < 1)
  {
    ncl    <- ncl + 1
    clas[codic[i]] <- ncl
    for (j in (i+1):nrow(x))
    {
      if (w[i,j] > 0.9) clas[codic[j]] <- clas[codic[i]]
    }
  }
}
}

```



```

if (clas[nrow(x)] < 1) clas[nrow(x)] <- ncl + 1
clas[clas == 0] <- NA
clas
}

#           alfa/2
tau<-qchisq(p=0.100,df=1) #0.015
tau<-qchisq(p=0.050,df=1) #0.004
tau<-qchisq(p=0.025,df=1) #0.001 tau alfa=0.05

system.time(
kk<-awsreg(y,x,coords=coords,bw=bw,tau=0.001,niter=200,
conv=0.0001,eta=0.5,numout=20,sout=1e-05)
)
kk

kk[is.na(kk)] <- max(kk,na.rm = TRUE)+1

# if more than 8 groups (no missing)
kk1<-kk
kk1[kk>8]<-0
plot(coordinates(subhouse),cex=.1,pch=20,col=kk1)

kk2<-kk
kk2[kk<9]<-0
points(coordinates(subhouse),cex=.6,pch=12,col=kk2)

plot(coordinates(subhouse),cex=.6,pch=19,col=kk)

```

برنامه‌های لازم برای افراز داده‌ها و رژیم‌بندی آنها

```

group1 = subhouse[kk==1,]
group2 = subhouse[kk==2,]
group3 = subhouse[kk==3,]

```

```
group4 = subhouse[kk==4,]

newy = cbind(t(subhouse$price[kk==1]),t(subhouse$price[kk==2]),
t(subhouse$price[kk==3]),t(subhouse$price[kk==4]))
newy = t(newy)
dim(newy)

group1x<-cbind(1,group1$yrbuilt,group1$TLA,group1$baths,
group1$halfbaths,group1$garagesqft,group1$lotsize) # X regim1

group2x<-cbind(1,group2$yrbuilt,group2$TLA,group2$baths,
group2$halfbaths,group2$garagesqft,group2$lotsize) # X regim2

group3x<-cbind(1,group3$yrbuilt,group3$TLA,group3$baths,
group3$halfbaths,group3$garagesqft,group3$lotsize) # X regim3

group4x<-cbind(1,group4$yrbuilt,group4$TLA,group4$baths,
group4$halfbaths,group4$garagesqft,group4$lotsize) # X regim4

newx = as.matrix(bdiag(group1x,group2x,group3x,group4x))
dim(newx)

newsubhouse = cbind(newy,newx) #new house data set

newsubhouse = data.frame(newsubhouse)

names(newsubhouse)

names(newsubhouse) = c("price","intercept1","yrbuilt1","TLA1",
"baths1","halfbaths1","garagesqft1","lotsize1","intercept2","yrbuilt2",
"TLA2","baths2","halfbaths2","garagesqft2","lotsize2","intercept3","yrbuilt3",
"TLA3","baths3","halfbaths3","garagesqft3","lotsize3","intercept4",
"yrbuilt4","TLA4","baths4","halfbaths4","garagesqft4","lotsize4"
)
```

دستورات لازم برای معرفی ماتریس وزن جدید و برازش مدل‌های اقتصادسنجی با رژیم فضایی درونی

```

coordinat1 = cbind(group1$long,group1$lat)
coordinat2 = cbind(group2$long,group2$lat)
coordinat3 = cbind(group3$long,group3$lat)
coordinat4 = cbind(group4$long,group4$lat)
coords1 = rbind(coordinat1,coordinat2,coordinat3,
coordinat4)      #new coordinat for baltimor data set
knn1      <- knn2nb(knearneigh(coords1,k=10))
listw_10nn_dates1 <- nb2listw(knn1)

SLM.ES <- lagsarlm(price~-1+intercept1+intercept2+intercept3+intercept4+
yrbuilt1+yrbuilt2+yrbuilt3+yrbuilt4+
TLA1+TLA2+TLA3+TLA4+
baths1+baths2+baths3+baths4+
halfbaths1+halfbaths2+halfbaths3+halfbaths4+
garagesqft1+garagesqft2+garagesqft3+garagesqft4+
lotsize1+lotsize2+lotsize3+lotsize4,
data=newsubhouse, listw_10nn_dates1,
type="lag", tol.solve= 1.58119e-19)
summary(SLM.ES)

options("scipen"=100)

f <- price ~ yrbuilt1+yrbuilt2+yrbuilt3+yrbuilt4+
TLA1+TLA2+TLA3+TLA4+
baths1+baths2+baths3+baths4+
halfbaths1+halfbaths2+halfbaths3+halfbaths4+
garagesqft1+garagesqft2+garagesqft3+garagesqft4+
lotsize1+lotsize2+lotsize3+lotsize4

SDM.ES <- lagsarlm(price~-1+intercept1+intercept2+intercept3+intercept4+

```

```
yrbuilt1+yrbuilt2+yrbuilt3+yrbuilt4+
TLA1+TLA2+TLA3+TLA4+
baths1+baths2+baths3+baths4+
halfbaths1+halfbaths2+halfbaths3+halfbaths4+
garagesqft1+garagesqft2+garagesqft3+garagesqft4+
lotsize1+lotsize2+lotsize3+lotsize4,
data=newsubhouse, listw_10nn_dates1,
Durbin=as.formula(delete.response(terms(f))),
type="mixed", tol.solve= 1.58119e-19)
summary(SDM.ES)
```

```
SEM.ES <- errorsarlm(price~-1+intercept1+intercept2+intercept3+intercept4+
yrbuilt1+yrbuilt2+yrbuilt3+yrbuilt4+
TLA1+TLA2+TLA3+TLA4+
baths1+baths2+baths3+baths4+
halfbaths1+halfbaths2+halfbaths3+halfbaths4+
garagesqft1+garagesqft2+garagesqft3+garagesqft4+
lotsize1+lotsize2+lotsize3+lotsize4,
data=newsubhouse, listw_10nn_dates1,
method="eigen", quiet=FALSE,tol.solve= 1.58119e-19)
summary(SEM.ES)
```


مراجع

- [۱] محمدزاده م.، (۱۳۹۱)، ”آمار فضایی” چاپ اول، مرکز نشر آثار علمی دانشگاه تربیت مدرس، تهران.
- [2] Akaike, H. (1976). An information criterion (AIC), **Mathematical Sciences**, 14, 5-7.
- [3] Andreano, M.S., Benedetti, R. & Postiglione, P. (2016). Spatial regimes in regional European growth: an iterated spatially weighted regression approach, **Quality & Quantity**, 51(6), 2665-2684.
- [4] Anselin, L. (1980). Estimation methods for spatial autoregressive structures, **Regional Science Dissertation and Monograph Series, Cornell University, Ithaca, NY**.
- [5] Anselin, L. (1988a). **Spatial Econometrics: Methods and Models**, Vol. 4, Springer Science & Business Media. London.
- [6] Anselin, L. (1988b). Lagrange multiplier test diagnostics for spatial dependence and spatial heterogeneity, **Geographical Analysis**, 20(1), 1-17.
- [7] Anselin, L., Florax, R. (Eds.) (1995). **New Directions in Spatial Econometrics**, Springer, New York.
- [8] Anselin, L. (1999). **Spatial Econometrics**, Bruton Center, School of Social Sciences, University of Texas at Dallas, Richardson, TX 75083-0688.
- [9] Anselin, L. (2010). Thirty years of spatial econometrics, **Papers in Regional Science**, 89, 3-25.
- [10] Besag, J. (1974). Spatial interaction and the statistical analysis of lattice systems, **Journal of the Royal Statistical Society, Series B (Methodological)**, 36(2), 192-236.
- [11] Bille, A.G., Benedetti, R. & Postiglione, P. (2017). A two-step approach to account for unobserved spatial heterogeneity, **Spatial Economic Analysis**, 4, 81-91.

-
- [12] Bivand, R., Altman, M., Anselin, L., Assunção, R., Berke, O., Bernat, A. & Blanchet, G. (2015). R Package 'spdep', <ftp://garr.tu cows.com/mirrors/CRAN/web/packages/spdep/spdep.pdf>
- [13] Bowman, A.W. (1984). An alternative method of cross-validation for the smoothing of density estimates, **Biometrika**, 71(2), 353-360.
- [14] Box, G.E. & Cox, D.R. (1964). An analysis of transformations, **Journal of the Royal Statistical Society, Series B (Methodological)**, 26(2), 211-252.
- [15] Brunsdon, C., Fotheringham, A.S. & Charlton, M.E. (1996). Geographically weighted regression: a method for exploring spatial nonstationarity, **Geographical Analysis**, 28(4), 281-298.
- [16] Breusch, T.S. & Pagan, A.R. (1979). A simple test for heteroscedasticity and random coefficient variation, **Econometrica: Journal of the Econometric Society**, 47(5), 1287-1294.
- [17] Burridge, P. (1980). On the clifford test for spatial correlation, **Journal of the Royal Statistical Society, Series B (Methodological)**, 42, 107-108.
- [18] Casetti, E. (1972). Generating models by the expansion method: applications to geographical research, **Geographical Analysis**, 4, 81-91.
- [19] Clayton, D. & Kaldor, J. (1987). Empirical Bayes estimates of age-standardized relative risks for use in disease mapping, **Biometrics**, 43(3), 671-681.
- [20] Clayton, D. & Bernardinelli, L. (1992). Bayesian methods for mapping disease risk, **Geographical and Environmental Epidemiology: Methods for Small Area Studies**, 205-220.
- [21] Cleveland, W.S. (1979). Robust locally weighted regression and smoothing scatterplots, **Journal of the American Statistical Association**, 74, 829-836.
- [22] Cleveland, W.S. & Devlin, S.J. (1988). Locally weighted regression: an approach to regression analysis by local fitting, **Journal of the American Statistical Association**, 83, 596-610.
- [23] Cressie, N.A. (1993). **Statistics for Spatial Data**, Wiley Series in Probability and Mathematical Statistics, New York.

- [24] Dubin, R.A. (1992). Spatial autocorrelation and neighborhood quality, **Regional Science and Urban Economics**, 22, 433–452.
- [25] Fotheringham, A.S., Brunson, C. & Charlton, M. (2003). **Geographically Weighted Regression: The Analysis of Spatially Varying Relationship**, John Wiley & Sons, New York.
- [26] Gollini, I., Lu, B., Charlton, M., Brunson, C. & Harris, P. (2013). GWmodel: an R package for exploring spatial heterogeneity using geographically weighted models, arXiv:1306.0413 [stat.AP].
- [27] Haining, R. (1990). **Spatial Data Analysis in the Social and Environmental Sciences**, Cambridge University Press, Cambridge.
- [28] Hurvich, C.M., Simonoff, J.S. & Tsai, C.L. (1998). Smoothing parameter selection in non-parametric regression using an improved Akaike information criterion, **Journal of the Royal Statistical Society, Series B (Statistical Methodology)**, 60(2), 271-293.
- [29] Ibragimov, R. & Muller, U.K. (2010). T-statistic based correlation and heterogeneity robust inference, **Journal of Business and Economic Statistics**, 28(4), 453-468.
- [30] Journel, A.G. & Huijbregts, C.J. (1978). **Mining Geostatistics**, Academic Press, London.
- [31] Kelejian, H.H. & Prucha, I.R. (1999). A generalized moments estimator for the autoregressive parameter in a spatial model, **International Economic Review**, 40(2), 509-533.
- [32] Kelejian, H.H. & Prucha, I.R. (2010). Specification and estimation of spatial autoregressive models with autoregressive and heteroskedastic disturbances, **Journal of Econometrics**, 157(1), 53-67.
- [33] LeSage, J.P. & Pace, R.K. (2009). **Introduction to Spatial Econometrics**, Chapman and Hall/CRC, London.
- [34] Matheron, G. (1971). Random sets theory and its application to stereology, **Journal of Microscopy**, 95, 15-23.
- [35] Million, E. (2007). The hadamard product, **Course Notes**, 3-6.
- [36] Mollie, A. & Richardson, S. (1991). Empirical Bayes estimates of cancer mortality rates using spatial models, **Statistics in Medicine**, 10(1),95-112.
- [37] Montgomery, D.C., Peck, E.A. & Vining, G.G. (2012). **Introduction to Linear Regression Analysis**, Vol. 821, John Wiley & Sons, New York.

-
- [38] Moran, P.A. (1948). Some theorems on time series II: the significance of the serial correlation coefficient, **Biometrika**, 35, 255-260.
- [39] Ord, K. (1975). Estimation methods for models of spatial interaction, **Journal of the American Statistical Association**, 70, 120-126.
- [40] Paelinck, J.H.P. & Klaassen, L.L.H. (1979). **Spatial Econometrics**, Vol. 1, Saxon House.
- [41] Postiglione, P., Andreano, M.S. & Benedetti, R. (2013). Using constrained optimization for the identification of convergence clubs, **Computational Economics**, 42, 151–174.
- [42] Pede, V.O., Florax, R.J.G.M. & Lambert, D.M. (2014). Spatial econometric STAR models: Lagrange multiplier tests, Monte Carlo simulations and an empirical application, **Regional Science and Urban Economics**, 49, 118–128.
- [43] Polzehl, J. & Spokoiny, V. (2000). Adaptive weights smoothing with applications to image restoration, **Journal of the Royal Statistical Society, Series B (Statistical Methodology)**, 62, 335-354.
- [44] Polzehl, J. & Spokoiny, V. (2006). Propagation-separation approach for local likelihood estimation, **Probability Theory and Related Fields**, 135, 335–362.
- [45] R Core Team. (2018). R: a language and environment for statistical computing, **R Foundation for Statistical Computing**, Vienna, Austria.
- [46] Smith, T.E. (2016). **Notebook on Spatial Data Analysis**, Lecture Note, <http://www.seas.upenn.edu/~ese502/>.
- [47] Silverman, B.W. (1984). A fast and efficient cross-validation method for smoothing parameter choice in spline regression, **Journal of the American Statistical Association**, 79, 584-589.
- [48] Wheeler, D.C. & Páez, A. (2010). Geographically weighted regression, **In Handbook of Applied Spatial Analysis**, Springer, Berlin, Heidelberg, 461-486.
- [49] Whittle, P. (1954). On stationary processes in the plane, **Biometrika**, 41, 434-449.

واژه‌نامه فارسی به انگلیسی

الف

Lagrange multiplier lag test	آزمون ضریب لاگرانژ تأخیر فضایی
Lagrange multiplier error test	آزمون ضریب لاگرانژ خطای فضایی
Likelihood ratio test	آزمون نسبت درست‌نمایی
Cross-validation	اعتبارسنجی متقابل
Breusch-Pagan statistic	آماره بریچ-پیگان
Moran I statistic	آماره موران I
Intrinsic stationary	ایستای ذاتی
Strong stationary	ایستای قوی
Second order stationary	ایستای مرتبه دوم

ب

Maximum likelihood estimation	برآورد درست‌نمایی ماکسیمم
Large scale	بزرگ‌مقیاس

پ

Memory parameter	پارامتر حافظه
------------------	---------------

خ

Spatial autocorrelation	خودهمبستگی فضایی
-------------------------	------------------

د

Point pattern data	داده‌های الگونقطه‌ای
Geostatistical data	داده‌های زمین‌آماری
Spatial data	داده‌های فضایی
Lattice data	داده‌های شبکه‌ای

ر

Spatial regimes	رژیم‌های فضایی
Geographical weighted regression	رگرسیون موزون جغرافیایی
Locally weighted regression	رگرسیون موزون موضعی

Iterative geographically weighted regression رگرسیون موزون موضعی تکراری

Trend روند

ض

Hadamard product ضرب هادامارد

ک

Generalized least squares کمترین توان‌های دوم تعمیم‌یافته

Ordinary least square کمترین توان‌های دوم معمولی

Weighted least squares کمترین توان‌های دوم موزون

Locally Weighted Least Squares کمترین توان‌های دوم موضعی

k-nearest neighbour k-نزدیکترین همسایگی

Small scale کوچک‌مقیاس

م

Spatial weight matrix ماتریس وزن فضایی

Vertex contiguity مجاورت رأسی

Radius contiguity مجاورت شعاعی

Edge contiguity مجاورت ضلعی

Conditional autoregressive model مدل اتورگرسیو شرطی

Spatial autoregressive models مدل‌های اتورگرسیو فضایی

Spatial autoregressive model with autoregressive disturbances مدل اتورگرسیو فضایی با خطاهای اتورگرسیو

Simultaneously autoregressive model مدل اتورگرسیو هم‌زمان

Spatial lag model مدل تأخیر فضایی

Spatial autoregressive combined model مدل ترکیبی فضایی

Spatial error model مدل خطای فضایی

Spatial durbin model مدل دوربین فضایی

Spatial econometric models مدل‌های اقتصادسنجی فضایی

Endogenous spatial regime models مدل‌های با رژیم فضایی درونی

Global models مدل‌های فراموضعی

Local models مدل‌های موضعی

Akaike information criterion معیار اطلاع آکاییک

Geometric mean میانگین هندسی

Gaussian random field میدان تصادفی گاوسی

ن

Nonsingular نامنفرد

Heterogeneous ناهمگن
Heterogeneity ناهمگنی
Spatial heterogeneity ناهمگنی فضایی

و

Spatial dependence وابستگی فضایی
Homoscedastic واریانس متجانس

ه

Fixed kernel هسته ثابت
Adaptive kernel هسته سازوار
Homogenous همگن
Adaptive weights smoothing هموارسازی وزن‌های سازوار

واژه‌نامه انگلیسی به فارسی

A

Adaptive kernel هسته سازوار
Adaptive weights smoothing هموارسازی وزن‌های سازوار
Akaike information criterion معیار اطلاع آکایک

B

Breusch-Pagan statistic آماره بریچ-پیگان

C

Conditional autoregressive model مدل اتورگرسیو شرطی
Cross-validation اعتبارسنجی متقابل

E

Edge contiguity مجاورت ضلعی
Endogenous spatial regime models مدل‌های با رژیم فضایی درونی

F

Fixed kernel هسته ثابت

G

Gaussian random field میدان تصادفی گاوسی
 Generalized least squares کمترین توان‌های دوم تعمیم‌یافته
 Geographical weighted regression رگرسیون موزون جغرافیایی
 Geometric mean میانگین هندسی
 Geostatistical data داده‌های زمین‌آماري
 Global models مدل‌های فراموضعی

H

Hadamard product ضرب هادامارد
 Heterogeneity ناهمگنی
 Heterogeneous ناهمگن
 Homogenous همگن
 Homoscedastic واریانس متجانس

I

Intrinsic stationary ایستای ذاتی
 Iterative geographically weighted regression رگرسیون موزون موضعی تکراری

K

k-nearest neighbour k-نزدیکترین همسایگی

L

Lagrange multiplier error test آزمون ضریب لاگرانژ خطای فضایی
 Lagrange multiplier lag test آزمون ضریب لاگرانژ تأخیر فضایی
 Large scale بزرگ‌مقیاس
 Lattice data داده‌های شبکه‌ای
 Likelihood ratio test آزمون نسبت درست‌نمایی
 Local models مدل‌های موضعی
 Locally Weighted Least Squares کمترین توان‌های دوم موزون موضعی
 Locally weighted regression رگرسیون موزون موضعی

M

- Maximum likelihood estimation برآورد درست‌نمایی ماکسیمم
Memory parameter پارامتر حافظه
Moran I statistic آماره موران I

N

- Nonsingular نامنفرد

O

- Ordinary least square کمترین توان‌های دوم معمولی

P

- Point pattern data داده‌های الگونقطه‌ای

R

- Radius contiguity مجاورت شعاعی

S

- Second order stationary ایستای مرتبه دوم
Spatial autoregressive combined model مدل ترکیبی فضایی
Simultaneously autoregressive model مدل اتورگرسیو هم‌زمان
Spatial autoregressive model with autoregressive disturbances
مدل اتورگرسیو فضایی با خطاهای اتورگرسیو
Small scale کوچک‌مقیاس
Spatial autocorrelation خودهمبستگی فضایی
Spatial autoregressive models مدل‌های اتورگرسیو فضایی
Spatial data داده‌های فضایی
Spatial dependence وابستگی فضایی
Spatial durbin model مدل دوربین فضایی

Spatial econometric models مدل‌های اقتصادسنجی فضایی
Spatial error model مدل خطای فضایی
Spatial heterogeneity ناهمگنی فضایی
Spatial lag model مدل تأخیر فضایی
Spatial regimes رژیم‌های فضایی
Spatial weight matrix ماتریس وزن فضایی
Strong stationary ایستای قوی

T

Trend روند

V

Vertex contiguity مجاورت رأسی

W

Weighted least squares کمترین توان‌های دوم موزون

Abstract

Many experimental studies need to use variables that are influenced by their geographical locations. In such cases, to consider the spatial structure of data are essential and ignoring this structure results in to lose some relevant information. To deal with such data, the classical linear regression model is not efficient. Indeed, the spatial structure leads to the violation of homogeneity and uncorrelation assumptions. Hence, we should use some appropriate spatial models. To have sufficient flexibility, we need a model that considers both spatial heterogeneity and dependency, simultaneously. However, the real heterogeneity and dependency structure of data is usually unknown. Furthermore, there is no available prior information about this structure.

In this thesis, we use an iterative locally weighted regression approach to determine homogeneous spatial regions, named as spatial regimes. Next, we apply these regimes in a class of spatial autoregressive models, known as spatial econometric models. We compare the proposed models with the linear regression model and classical spatial econometric models based on the Akaike information criterion. The results show a better performance of econometric models with the endogenous spatial regimes.

Keywords: Spatial heterogeneity, Spatial dependence, Spatial regimes, Spatial econometrics, Endogenous spatial regime econometric models.



Shahrood University of Technology

Faculty of Mathematical Sciences

MSc Thesis in: Statistics

**Spatial Heterogeneity in Production Price
Models by Using an Iterative Locally
Weighted Regression Approach**

By: Roya Payrooliya

Supervisor

Hossein Baghishani

January 2019