

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه پیام نور  
واحد تهران  
دانشکده علوم پایه

## **پایان نامه برای دریافت درجه کارشناسی ارشد**

**گرایش :**

آمار ریاضی

**عنوان :**

برآورد ناپارامتری استوار رگرسیون فضایی

**استاد راهنما:**

دکتر مسعود یار محمدی

**استاد مشاور :**

دکتر پرویز نصیری

**پژوهش گر :**

صفیه غریب آبادی

## فهرست مطالب

صفحه	عنوان
۱	چکیده :
۲	پیشگفتار:
۴	فصل اول: آمار فضایی
۵	۱- مقدمه
۶	۲-۱ ضرورت استفاده از آمار فضایی
۸	۳-۱: انواع داده های فضایی
۱۰	۴-۱: موقعیتهای فضایی
۱۱	۵-۱: تفاوت رگرسیون کلاسیک و فضایی
۱۱	۶-۱: مدل آماری
۱۲	۷-۱: میدان تصادفی
۱۴	فصل دوم: روشهای استوار در تحلیل داده های فضایی
۱۵	۱-۲: تعریف داده فضایی
۱۶	۱-۱-۲: روش های آماری استوار
۱۸	۲-۱-۲: شاخص های استواری
۱۹	۳-۱-۲: روش های استوار فضایی
۱۹	۴-۱-۲: شناسایی داده های دورافتاده
۲۰	۵-۱-۲: داده های دورافتاده فضایی و جستجوی پیشرو
۲۰	۶-۱-۲: داده های دورافتاده فضایی
۲۱	۷-۱-۲: نمای کلی جستجوی پیشرو
۲۳	۸-۱-۲: جستجوی پیشرو برای داده های فضایی
۲۴	۹-۱-۲: کریگیدن معمولی
۲۷	۲-۲: روش های استوار در حضور داده های دور افتاده
۲۷	۱-۲-۲: برآورد استوار تغییرنگار
۲۹	۲-۲-۲: برآوردگرهای استوار مقیاس با نقطه فرو ریزش بالا

۳۰	۳-۲-۲: برآورد استوار پارامترهای روند
۳۱	۳-۲: مدل خطی فضایی
۳۲	۱-۳-۲: تعمیم برآوردگرهای استوار به مدل های خطی فضایی
۳۵	<b>فصل سوم: برآورد ناپارامتری استوار رگرسیون فضایی</b>
۳۶	۱-۳: مقدمه
۳۶	۲-۳: مدل
۳۷	۳-۳: فرضیه ها و نظریه های مورد نیاز:
۳۸	۱-۳-۳: نظریه هایی برای مفاهیم ناپارامتری
۳۸	۲-۳-۳: نظریه در مورد خواص استوار سازی
۳۹	۳-۳-۳: نظریه هایی در مورد شرایط وابستگی فضایی
۴۰	۴-۳-۳: نظریه هایی در مورد هسته $K$
۴۰	۵-۳-۳: نظریه هایی در مورد پهنای باند پارامتری
۴۱	۴-۳: نتایج اصلی تحقیق
۴۱	۱-۴-۳: سازگاری
۵۰	۲-۴-۳: به طور مجانبی نرمال بودن :
۶۰	<b>فصل چهارم: شبیه سازی</b>
۶۱	۱- شبیه سازی
۷۰	۲- نتیجه گیری

## چکیده:

در این پایان نامه، برآورد ناپارامتری استوار برای رگرسیون فضایی مورد بحث و بررسی قرار می‌گیرد. به طور دقیق‌تر برای میدان تصادفی مانای اکید داده شده  $Z_i = (X_i, Y_i)_{i \in N^N, N \geq 1}$ ، یک خانواده از برآوردهای ناپارامتری استوار برای یک تابع رگرسیونی بر اساس روش هسته ای در نظر گرفته می‌شود.

تحت برخی فرضیه های مرکب کلی سازگاری تقریبی و نرمال مجانبی برای این برآوردها بدست آورده و در ادامه یک روش استوار برای انتخاب پارامتر هموار جهت تطبیق که به داده های فضایی مورد بحث و بررسی قرار می‌گیرد.

## واژگان کلیدی:

توزیع متقارن- همگرایی تقریباً کامل- میدان تصادفی- رگرسیون ناپارامتری- برآورد هسته- برآورد استوار

## پیشگفتار:

در بسیاری از مسائل آماری در زمینه های کاربردی نظیر ، زمین شناسی، اقیانوس شناسی ،اقتصاد سنجی، همه گیرشناسی، علوم محیطی، جنگلداری و غیره ، اثربخشی بردار کمکی از متغیرها روی برخی از متغیرهای پاسخ در وابستگی فضایی مورد مطالعه قرار گرفته است. هدف از این پایان نامه مطالعه برآورد استوار یک تابع رگرسیون به عنوان یک ابزار تحلیلی برای این نوع از داده است. به طور کلی برای مدل سازی کردن داده ها از روشهای آماری استفاده می شود، مباحث آماری مربوط به مدل بندی داده های فضایی از دیر باز مورد توجه آماردانان بوده است. به عنوان مثال می توان گیون<sup>۱</sup> (۱۹۹۵) ، انسیلین و فلورکسی<sup>۲</sup> (۱۹۹۵) ، کرسی<sup>۳</sup> (۱۹۹۱) و ریپلای<sup>۴</sup> (۱۹۸۱) اشاره کرد، به طوری که اولین مطالعه توسط ترن<sup>۵</sup> (۱۹۹۰) بدست آمده است. او نرمال متقارن بودن برآورد چگالی هسته ای را پایه گذاری کرده است. لو و چن<sup>۶</sup> (۲۰۰۴) برآورد هسته ای از تابع رگرسیون را مورد مطالعه قرار داده و نرخ همگرایی (همگرایی در احتمال) از این برآورد (برآورد هسته ای) را بدست آورده و بیو و کادر<sup>۷</sup> (۲۰۰۴) سازگاری یکنواخت و نرمال متقارن ، برآورد هسته ای از تابع رگرسیون را ثابت نمود. کاربن<sup>۸</sup> (۲۰۰۷) مدل اتو رگرسیون ناپارامتری برای پیشگویی میدان های تصادفی را مورد مطالعه قرار دادند. اخیراً لی و ترن<sup>۹</sup> (۲۰۰۹) یک روش متقابل از یک برآورد ناپارامتری از رگرسیون فضایی را بدست آورده اند. این روش بر اساس وزن گذاری روی نزدیک ترین همسایگی پایه گذاری شده است، نشان داده اند که برآوردشان ، مخالف نرمال تقریبی است. بحث ابتدایی رگرسیون استوار را هوبر<sup>۱۰</sup> در سال (۱۹۶۴) شروع کرده و توسط لی لاینز تحقیقی در مورد رگرسیون استوار بدست آمده است. برای این که تعمیم پیوستگی رگرسیون استوار را مشاهده نمود، می توان به روبینسون<sup>۱۱</sup> (۱۹۸۴) و کلوب و هاردل<sup>۱۲</sup> (۲۰۰۰) ، و همچنین برای

- 
- 1-Guyon
  - 2-Anselin and Florax
  - 3-Cressi
  - 4-Riply
  - 5-Tran
  - 6-Lu and C
  - 7-Biau and cadre
  - 8-Carbon
  - 9-Li and Tran
  - 10-Huber
  - 11-Robinson
  - 12- Collomb and Hardle

نتایج قبلی ، بونت و فریمن<sup>۱</sup> (۲۰۰۹) برای پیشرفت های اخیر و به عنوان مرجع مراجعه نمود. در اغلب این تحقیقات فرض شده است که مشاهدات وابسته از داده های سری زمانی آمده است. در حالت فضائی، زو و وانگ<sup>۲</sup> (۲۰۰۸) برآورد خطی مکانی از مدل رگرسیون را مورد مطالعه قرار داده اند. این برآورد شبه پارامتری به وسیله مینیم کردن قدر مطلق انحراف آنها بدست می آید. با بررسی های به عمل آمده مسئله برآورد ناپارامتری رگرسیون فضایی بوسیله یک روش استوار موضوع جدیدی بوده و توسط ابدلادرگریباله ، الی لاکساکسی، رچیداروان<sup>۳</sup> (۲۰۱۰) معرفی شده است. در فصل اول این پایان نامه تعاریف پایه ای از آمار فضایی و مروری بر رگرسیون کلاسیک و رگرسیون فضایی ارائه می شود. در فصل دوم روشهای استوار در تحلیل داده های فضایی ارائه می شود. در فصل سوم به تجزیه و تحلیل قضایا و لم هایی که مربوط به برآورد ناپارامتری رگرسیون فضایی است پرداخته و در فصل چهارم روشهای شبیه سازی به تفصیل مورد مطالعه قرار می گیرد.

---

1- Boente and fraiman

2- Xu and Wang

3- Abdelkader Gheriballah, Ali Laksaci, Rachida Rouane

# فصل اول

## آمار فضایی



در آمار، وقتی جامعه آماری مورد نمونه برداری قرار می گیرند، اغلب فرض می شود که مشاهدات حاصل از نمونه از یکدیگر مستقل هستند. با فرض استقلال مشاهدات، استنباط آماری آنها بر پایه قضایای ریاضی و آمار ساده می گردد در عمل با موارد زیادی مواجه می شویم که مشاهدات به نوعی به یکدیگر وابسته هستند به عنوان مثال داده ها در سری های زمانی مثال از داده هایی هستند که در طول زمان همبسته می باشد. داده هایی که در فضای مورد مطالعه وابسته به یکدیگر بوده و این وابستگی نوعاً به محل قرار گرفتن آنها اغلب همبستگی روی زمان را به نمایش می گذارند. به دلیل همبستگی فضایی، روش های کلاسیک آماری برای تجزیه و تحلیل چنین داده هایی قابل استفاده نمی باشند. به همین دلیل شاخه آمار فضایی ایجاد گردیده است. نخستین بار به دنبال روند تکاملی ذخایر معدنی که قبل از سال (۱۹۶۰) آغاز شده بود و به ویژه براساس پژوهشهای افرادی مانند سیشل و کروننگ، ماترون پژوهشگر فرانسوی با انتشار مقاله ای در سال (۱۹۶۲) پایه های زمین آمار را بنا نهاد و شاخه ای جدید علم آمار را تحت همین عنوان بوجود آورد. زمین آمار که توسط ماترون (۱۹۶۲) پایه گذاری شده است، به شاخه ای از علم آمار گفته می شود به تجزیه و تحلیل داده هایی همبسته که همبستگی شان تابعی از موقعیت قرار گرفتن جهت داده ها می باشد، می پردازد. در واقع بین مقادیر مختلف یک متغیر و یک فاصله و جهت قرارگیری آنها ارتباطی را برقرار می کند، که این ارتباط فضایی، ساختار فضایی نامیده می شود.

بررسی و مطالعه این ساختار یکی از مسائل مهم آمار فضایی است. در شکل گیری آمار فضایی افراد مختلفی سهمیم بوده و به طور موازی با سیشل، کریگ و ماترون فعالیت داشتند. از جمله گایدن (۱۹۶۲) از کشور روسیه که مبحثی به نام تحلیل عینی را با مفاهیم بسیار نزدیک به تعاریف ماترون و دیگران در زمینه هواشناسی پیشنهاد کرده است و یا ماترون سوئدی (۱۹۶۲) روشهایی به نام تجزیه و تحلیل داده های فضایی را ابداع نمود. در این راستا دیوبد (۱۹۷۷) به طور مفصل روش های زمین آماری را با اصطلاحات خاص شرح داده است.

ریپلی (۱۹۸۱) در کتاب آمار فضایی دریچه دیگری گشود و با زبان آماری سخن گفت و نظرات مربوط به پیشگویی فرایندهای تصادفی که ارتباط نزدیکی با سری های زمانی دارند را به کارگرفت. کرسی (۱۹۹۳) نیز با انتشار کتاب آمار برای داده فضایی موجب پیشرفت گسترده این شاخه از علوم گردید. در سالهای اخیر به مدد تلاش بسیاری از پژوهشگران، آمار فضایی به عنوان وسیله ای کارآمد برای مطالعه داده های فضایی مطرح گردیده است.

اگرچه عمده زمینه های رشد و توسعه آمار فضایی در مسائل مربوط به ذخائر معدن بوده است ولی در زمینه های مختلف از علوم مانند جغرافیا، همه گیرشناسی، اکتشاف نفت، گاز و کنترل کیفیت هوا تحقیقات کاربردی فراوان انجام گرفته است (خالدی ۱۳۷۸).

## ۲-۱- ضرورت استفاده از آمار فضایی

در آمار کلاسیک مشاهدات مستقل از هم و فارغ از موقعیت مکانی فرض می شوند، در نتیجه نظریه استنباط آماری آنها ساده است. اما در نظر گرفتن موقعیت فضایی و هم چنین وابستگی فضایی بین مشاهدات نمونه، به شناخت بیش تر جامعه مورد بررسی کمک می کند. کرسی (۱۹۹۳) بیان ضرورت در نظر گرفتن همبستگی فضایی در تحلیل داده ها به ذکر مسأله برآورد میانگین جامعه ای با واریانس معلوم  $\sigma^2$  پرداخت. در این مسأله با فرض استقلال مشاهدات و به کمک نمونه تصادفی  $X_1, \dots, X_n$ ,

از برآوردگر نارایب با کم ترین واریانس  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  با واریانس  $\frac{\sigma^2}{n}$  استفاده می شود.

اگر توزیع جامعه ای که از آن نمونه استخراج شده نرمال باشد یا  $n$  به اندازه کافی بزرگ باشد، بنابر قضیه حد مرکزی،  $\bar{X}$  دارای توزیع نرمال با میانگین  $\mu$  و واریانس  $\frac{\sigma^2}{n}$  خواهد بود. بنابراین بازه اطمینان

۹۵٪ برای  $\mu$  به صورت

$$\left( \bar{X} - 1.96 \frac{\sigma}{\sqrt{n}}, \bar{X} + 1.96 \frac{\sigma}{\sqrt{n}} \right) \quad (1-1)$$

اما اگر داده ها به طور مستقل نبوده و همبستگی آنها از طریق تابع کواریانس

$$Cov(X_i, X_j) = \sigma^2 \rho^{|i-j|}, i, j = 1, \dots, n \quad (2-1)$$

قابل بیان باشد، واریانس برآوردگر  $\bar{X}$  برابر

$$\begin{aligned} Var(\bar{X}) &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n Cov(X_i, X_j) \\ &= \frac{\sigma^2}{n} \left[ 1 + 2 \left( \frac{\rho}{1-\rho} \right) \left( 1 - \frac{1}{n} \right) - 2 \left( \frac{\rho}{1-\rho} \right)^2 \frac{1-\rho^{n-1}}{n} \right] \quad (3-1) \end{aligned}$$

خواهد بود. این مقدار با واریانس برآوردگر در حالت استقلال مشاهدات، یعنی  $\frac{\sigma^2}{n}$  تفاوت قابل

ملاحظه ای دارد. علاوه بر این، به ازای  $n=10$  و  $\rho=0.26$ ،  $Var(\bar{X}) = \frac{\sigma^2}{n} [1.608]$  و بازه اطمینان ۹۵٪

برای  $\mu$  به صورت

$$\left( \bar{X} - 2.485 \frac{\sigma}{\sqrt{n}}, \bar{X} + 2.485 \frac{\sigma}{\sqrt{n}} \right)$$

خواهد شد. همان طور که ملاحظه می شود، وجود همبستگی داده ها از نوع (۲-۱) موجب افزایش واریانس  $\bar{X}$  و نیز افزایش طول برآورد بازه ای ۹۵٪ باشد، بلکه احتمال پوشش آن ۷۰٪ است. فرض کنید  $X_1, \dots, X_n$  و  $Y_1, \dots, Y_n$  دو نمونه مستقل از دو جامعه به ترتیب با توزیع های  $N(\mu_x, \sigma^2)$  و  $N(\mu_y, \sigma^2)$  باشند در این صورت  $\bar{X}$  و  $\bar{Y}$  برآوردهای  $\mu_x$  و  $\mu_y$  نیز مستقل و به ترتیب دارای توزیع های  $N(\mu_x, \frac{\sigma^2}{n})$  و  $N(\mu_y, \frac{\sigma^2}{n})$  خواهند بود. آماره آزمون برای فرضیه  $H_0: \mu_x = \mu_y$  در مقابل  $H_1: \mu_x \neq \mu_y$  به صورت زیر است

$$Z = \frac{\bar{X} - \bar{Y}}{\sigma \sqrt{\frac{2}{n}}} \sim N(0,1)$$

اما اگر مشاهدات هر یک از دو نمونه وابسته بوده دارای کواریانس

$$\text{Cov}(X_i, X_j) = \sigma^2 \rho; i \neq j = 1, \dots, n \quad (4-1)$$

باشند، آن گاه

$$\begin{aligned} \text{Var}(\bar{X}) &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \text{Cov}(X_i, X_j) \\ &= \frac{1}{n^2} [n\sigma^2 + n(n-1)\sigma^2\rho] \\ &= \frac{\sigma^2}{n} [1 + (n-1)\rho] \end{aligned}$$

و به طور مشابه  $\text{Var}(\bar{Y}) = \frac{\sigma^2}{n} [1 + (n-1)\rho]$ . در این صورت آماره آزمون برای فرضیه  $H_0: \mu_x = \mu_y$  به صورت زیر می شود.

$$Z^* = \frac{\bar{X} - \bar{Y}}{\sigma \sqrt{2(1 + \frac{(n-1)\rho}{n})}} \sim N(0,1)$$

بدیهی است مقدار مشاهده شده  $Z$  در نمونه ها همواره بزرگ تر از مقدار  $Z^*$  است، یعنی اگر همبستگی داده ها در نظر گرفته نشود، مقدار مشاهده شده  $Z$  خیلی بزرگ تر و به تبع آن  $p$ -مقدار کوچک تر و در نتیجه شواهد رد  $H_0$  در داده ها بیش تر از حالتی خواهد بود که همبستگی داده ها منظور می شوند، یعنی آزمون بیش از مواردی که لازم است رد می شود.

نکته قابل توجه دیگر آن که  $E(\bar{X}) = \mu_x$  و  $\lim_{n \rightarrow \infty} \text{Var}(\bar{X}) = \sigma^2 \rho$ ، یعنی بدون توجه به همبستگی داده ها، همواره  $\bar{X}$  برآوردگری ناریب برای میانگین جامعه است، اما در صورت وابستگی داده ها،  $\bar{X}$  برآوردگری ناسازگار است.

اگر رابطه (۳-۱) به صورت  $\text{Var}(\bar{x}) = \frac{\sigma^2}{n}$  نوشته شود، آن گاه می توان

$$n' = \frac{n}{1 + 2\left(\frac{\rho}{1-\rho}\right)\left(1 - \frac{1}{n}\right) - 2\left(\frac{\rho}{1-\rho}\right)^2 \left(\frac{1-\rho^{n-1}}{n}\right)} \quad (5-1)$$

را به عنوان تعداد مشاهدات مستقل معادل  $n$  مشاهده همبسته تعبیر کرد. به عنوان مثال اگر  $n = 10$  و  $\rho = 0.26$  آن گاه  $n' = 6.2$  یعنی ۶ مشاهده مستقل به طور تقریبی نتایجی با همان دقت

$n = 10$  مشاهده همبسته ایجاد می کنند. برای مقادیر بزرگ  $n$  نیز  $n' \approx \left[ \frac{n}{(1+\rho)(1-\rho)} \right]$  بیانگر آن است که

در نمونه های بزرگ نیز همبستگی داده ها در دقت استنباط آماری تأثیرگذار است. در نتیجه توجه نکردن به همبستگی فضایی داده ها، دقت نتایج تحلیل آماری آن ها را بسیار تحت تأثیر قرار می دهد و لازم است مراحل مختلف تحلیل آماری داده ها با توجه به همبستگی داده ها انجام شود.

### ۳-۱: انواع داده های فضایی

در آمار فضایی متغیر مورد اندازه گیری ممکن است گسسته یا پیوسته باشد. فضای موقعیت یا مکان مشاهدات نیز ممکن است پیوسته یا گسسته باشد. فضای موقعیت یا مکان مشاهدات نیز ممکن است پیوسته یا گسسته، نقطه ای یا ناحیه ای، منظم یا نامنظم باشد. وقتی مقدار متغیر در یک ناحیه ثبت می شود، موقعیت ناحیه ای و اگر مقدار متغیر در یک نقطه ثابت یا طول و عرض جغرافیایی معین اندازه گیری شود، موقعیت نقطه ای است. در موقعیت نقطه ای اگر نقاط در فواصل مساوی از هم قرار داشته باشند، موقعیت های منظم و در غیر اینصورت نامنظم اند. در مورد مکان های ناحیه ای اگر ناحیه ها هم شکل و هم اندازه باشند، ناحیه ها منظم و در غیر این صورت نامنظم اند. با توجه به انواع موقعیت ها مشاهدات فضایی به سه گروه عمده، داده های زمین آماری<sup>۱</sup>، داده های شبکه ای<sup>۲</sup> و الگوهای نقطه ای تقسیم می شوند.

**داده های زمین آماری:** این نوع داده ها در موقعیت های ثابت و مشخص در ناحیه ای پیوسته مشاهده می شوند. متغیر مورد بررسی ممکن است گسسته یا پیوسته باشد. به عنوان مثال برای حالت پیوسته می توان غلظت مواد معدنی در داخل یک معدن، مقدار ریزش باران در ایستگاه هواشناسی و برای حالت گسسته می توان تعداد نوعی جانور دریایی در یک سری از مکان های نمونه گیری شده در طول یک ساحل را نام برد.

**داده های شبکه ای:** این نوع داده ها مربوط به مکان های ناحیه ای هستند، که این مکان ها ممکن است منظم یا نامنظم باشند. تصاویر ماهواره ای از سطح زمین که در آن ها سطح زمین به تعدادی عنصر تصویر<sup>۱</sup> کوچک به صورت شبکه منظم در  $R^2$  تقسیم شده، مثالی برای داده های شبکه ای منظم است. تعداد افراد مبتلا به سرطان در تمام مناطق خدمات درمانی کشور مثالی برای داده های شبکه ای نامنظم است. در این حالت کل افراد مبتلا به سرطان در هر شهرستان در مرکز خدمات درمانی آن شهرستان در نظر گرفته می شوند. تعداد حوادث رانندگی، تعداد اتباع خارجی غیر مجاز ساکن یا نرخ رشد اقتصادی هر استان مثال های دیگری از داده های شبکه ای هستند. به طور معمول هدف از تحلیل داده های شبکه ای مدل بندی احتمالاتی مشاهدات است، در صورتی که در زمین آمار پیش گویی مقدار متغیر در یک موقعیت جدید مدنظر است.

**الگوی نقطه ای:** در این حالت مکان یا موقعیت مشاهده شده خود متغیری تصادفی است. الگوهای نقطه ای شامل تعدادی متناهی از مکان ها در یک ناحیه اند که در آن ها یک صفت خاص اندازه گیری می شود. به طور معمول الگوهای نقطه ای به سه دسته به طور کامل تصادفی فضایی، منظم و خوشه ای تقسیم و به مدل بندی آن ها اقدام می شود. یک مثال برای این نوع مشاهدات موقعیت گونه ای از درختان در یک ناحیه ای جنگلی یا موقعیت مراکز زلزله است (پور طاهری و همکاران ۲۰۰۶)

موقعیت مکانی هر نقطه روی کره زمین به صورت مختصات جغرافیایی بر حسب درجه و دقیقه بیان می شود که بر مبنای هندسه کره استوار است. در این سیستم کره زمین به ۶۰ نوار عمودی یا منطقه هر یک به طول ۶ درجه تقسیم می شود. محدوده شمالی این سیستم عرض جغرافیایی ۸۴ درجه و محدوده جنوبی آن تا درجه ۸۰ می باشد. برای مناطق بالاتر و پایین تر این عرض ها، از سیستم تصویر صفحه ای تحت عنوان سیستم استرایوگرافیک قطبی جهانی<sup>۲</sup> (UPS) استفاده می شود. هر کدام از

نوارهای ۶۰ گانه با یک عدد مشخص می شود. به نوار بین ۱۷۴ تا ۱۸۰ درجه عدد یک را اختصاص می دهند و هر چه به سمت شرق جلو برویم اعداد افزایش پیدا می کند. بنابراین منطقه ۳۱ دارای طول صفر تا ۶ درجه شرقی است. ایران نیز در منطقه های ۳۸، ۳۹، ۴۰، ۴۱ قرار گرفته است.

نحوه دیگر ارائه موقعیت مکانی هر نقطه روی سطح دو بعدی، سیستم مختصات دکارتی است. برای تصویر سطح کره زمین بر روی سطح کاغذ از سیستم تصویر مرکاتور جانبی جهانی<sup>۱</sup> (UTM) استفاده می شود که متریک است و طول و عرض هر نقطه جغرافیایی را در مختصات دکارتی به صورت  $X$  و  $Y$  نمایش می دهد.

با توجه به این که در آمار فضایی موقعیت مکانی هر نقطه به صورت مختصات دکارتی در نظر گرفته می شود، لازم است تا موقعیت جغرافیایی مشاهدات با استفاده از سیستم UTM به مختصات دکارتی تبدیل شود. لازم به ذکر است که مختصات جدید موقعیت داده های واقع در یک ناحیه در سیستم تبدیل UTM، با تعیین فاصله هر موقعیت از وسط آن ناحیه (محل تقاطع خط استوا با نصف النهار میانی ناحیه) محاسبه شود. بنابراین با استفاده از این سیستم برای داده هایی که در نواحی مختلف قرار دارند، صحیح نیست. در این گونه موارد از سیستم لامبرت<sup>۲</sup> استفاده می شود که مبدأ سنجش فاصله موقعیت ها را یک نقطه دلخواه (به طور مثال تقاطع خط استوا با نصف النهار گرینویچ) فرض می کند و فاصله هر موقعیت از آن نقطه را ارائه می دهد.

## ۱-۴: موقعیتهای فضایی

منظور از موقعیت های فضایی<sup>۳</sup> دنباله  $\{s_1, \dots, s_n\}$  است، که این نقاط به طور منظم یا بی قاعده در  $R^d$  قرار گرفته و داده های  $\{Z(s_1), \dots, Z(s_n)\}$  در این موقعیت ها مشاهده می شوند. معمولاً مولفه های  $s \in R^d$  به صورت  $s = (s[1], \dots, s[d])$  نشان داده می شود. در آمار فضایی همبستگی قوی بین داده ها در موقعیت های نزدیک برقرار می گردد. بسیج (۱۹۷۴) چندین مثال در ارتباط با طبقه بندی های مختلف موقعیت ها و متغیرهای مرتبط ارائه می کند. این طبقه بندی بر اساس سه مورد زیر می باشد:

(۱) سیستم و موقعیت ها (منظم یا نامنظم)

(۲) موقعیت های منفرد (نقطه ها یا ناحیه ها)

---

2-Universal polar stereographic

3-Lambert

4- Spatial Location or sits

(۳) متغیرهای تصادفی (مرتبط پیوسته یا گسسته)

بر اساس تقسیم بندی فوق داده های فضایی به سه نوع تقسیم می شوند. اگر داده هادر موقعیت های ثابت مشاهده شوند و متغیر مورد مطالعه پیوسته یا گسسته باشد، داده های زمین آماری نامیده می شوند. در صورتی که داده ها مربوط به مکان های ناحیه ای باشند، داده های شبکه ای و در صورتی که سایت یا موقعیت خود یک متغیر تصادفی باشد، الگوهای نقطه ای نامیده می شوند.

## ۱-۵: تفاوت رگرسیون کلاسیک و فضایی

در رگرسیون کلاسیک مدل آماری  $y_i = g(x_i) + \varepsilon_i$  ;  $i = 1, \dots, n$  با این فرض که خطاها به طور مستقل دارای توزیع نرمال با میانگین صفر و واریانس  $\sigma^2$  هستند، تجزیه و تحلیل می گردد. در عمل پس از انتخاب مدل رگرسیونی بررسی درستی سه فرضیه نرمال بودن خطاها در اولویت قرار می گیرد.

درستی فرض های فوق به اعتبار مدل رگرسیونی منجر می گردد. در رگرسیون فضایی بین متغیرها و در نتیجه خطاها، نوعی همبستگی فضایی وجود دارد. از اینرو در رگرسیون فضایی بررسی فرض خود ناهمبسته نبودن خطاها غیر ضروری است. همچنین در مدل رگرسیون کلاسیک  $\{x_1, \dots, x_n\}$  بردار  $Z = (Z(s_1), \dots, Z(s_n))$  به عنوان نمونه در این موقعیت ها مشاهده می گردد.

## ۱-۶: مدل آماری

برای تحلیل داده های فضایی لازم است یک مدل آماری در نظر گرفته شود. در آمار فضایی به طور معمول یک میدان تصادفی به عنوان مدل های آماری در نظر گرفته می شود. میدان تصادفی مجموعه ای از متغیرهای تصادفی مانند  $\{Z(s); s \in D\}$  است، که در آن مجموعه اندیس گذار یک زیر مجموعه ای از فضای اقلیدسی  $D$  بعدی،  $d \geq 1$ ، از  $R^d$  است. در مورد میدان تصادفی  $Z(\cdot)$  میانگین در موقعیت  $s$  و کواریانس در موقعیت های  $s_1$  و  $s_2$  به ترتیب به صورت

$$E(Z(s)) = \mu(s) ; s \in D$$

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

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

تعریف می شوند. برای  $s = s_1 = s_2$ ، واریانس میدان تصادفی  $Z(\cdot)$  در مکان  $s$  بصورت :

$$\text{Var}(Z(s)) = E[z(s) - \mu(s)]^2$$

حاصل می شود. هر میدان تصادفی  $Z(\cdot)$  را می توان به صورت

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

زیر تجزیه کرد، که در آن  $\mu(\cdot)$  تغییرات بزرگ مقیاس<sup>۱</sup> یا روند<sup>۲</sup> و  $\delta(\cdot)$ ، فرایند خطای تغییرات کوچک مقیاس<sup>۳</sup> میدان تصادفی نامیده می شوند. تغییرات کوچک مقیاس ممکن است ناشی از خطای اندازه گیری یا تغییرپذیری در درون موقعیت مشاهده شده باشد و تغییرات بزرگ مقیاس ممکن است ناشی از تغییرات بین موقعیت های مشاهده شده باشند.

به طور معمول تحلیل داده های فضایی براساس مشاهدات نمونه دشوار است، اما برخی ویژگی های میدان تصادفی موجب ساده سازی مسأله خواهند شد، که در ادامه معرفی می شوند.

## میدان تصادفی

میدان تصادفی<sup>۴</sup> برای مدل سازی داده های فضایی به منظور تجزیه و تحلیل آنها به کار می رود. یک میدان تصادفی مجموعه ای از متغیرهای تصادفی مانند  $\{Z(s); s \in D\}$  است، که در آن  $D$  زیر مجموعه ای از فضای اقلیدسی  $d \geq 1$  بعدی  $R^d$  است. برای حالت خاص  $d=1$ ،  $\{Z(s); s \in D\}$ ، فرایند تصادفی نامیده می شود. مجموعه اندیس گذار  $D$  می تواند پیوسته، مانند  $R^d$  یا زیر مجموعه ای از  $R^d$ ، یا گسسته به صورت شبکه ای<sup>۵</sup> از موقعیت ها، مانند  $Z^d$  یا زیر مجموعه ای از  $Z^d$  باشد. میدان تصادفی پیوسته یا شبکه نامیده می شود هر گاه  $D$  به طور متناظر پیوسته یا گسسته باشد. شبکه می تواند بصورت اشکال هندسی منظم مانند یا نامنظم باشد. برای داده های زمین آماری  $D$  زیر مجموعه ای از مستطیل های  $d$  بعدی با حجم ثابت و پیوسته در  $R^d$  است. برای الگوهای شبکه ای  $D$  زیر مجموعه ای ثابت و شمارشی از  $R^d$  می باشد.

## تعریف ۱-۱

میدان تصادفی  $Z(\cdot)$  مانای ذاتی<sup>۶</sup> نامیده می شود، هرگاه

۱- میانگین میدان تصادفی ثابت یا مستقل از  $s$  باشد، یعنی

- 1-Large scale variation.
- 2-Trend
- 3-Small scale variation
- 4- Random Field
- 5- Lattice

1-Intrinsic stationary



$$E(Z(s)) = \mu$$

۲- واریانس عبارت  $(Z(s_1) - Z(s_2))$  فقط تابع یاز فاصله موقعیت های  $s_1$  و  $s_2$  باشد، یعنی

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

## تعریف ۲-۱

میدان تصادفی  $Z(\cdot)$ ، مانای مرتبه دوم ۱ (یا مانای ضعیف) نامیده می شود، هرگاه:

۱- میانگین میدان تصادفی  $Z(\cdot)$  مستقل از  $s$  و ثابت باشد. یعنی

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

۲- کواریانس  $Z(s_1)$  و  $Z(s_2)$ ، فقط تابعی از فاصله موقعیت های  $s_1$  و  $s_2$  باشد، یعنی

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

واضح است که تحت مانایی مرتبه دوم

$$\text{Var}(Z(s)) = \text{Cov}(Z(s), Z(s)) = C(0) = \sigma^2$$

یعنی واریانس میدان تصادفی به موقعیت فضایی بستگی ندارد و تغییرپذیری میدان در همه جا یکسان است.

## تعریف ۱-۳

میدان تصادفی  $Z(\cdot)$  مانای قوی<sup>۲</sup> نامیده می شود، هر گاه برای همه موقعیت های  $s_1, \dots, s_n$  و تأخیر  $h$ ، توزیع توأم  $Z(s_1), \dots, Z(s_n)$  باشد. یعنی

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

به بیان دیگر، در صورت انتقال موقعیت های  $s_1, \dots, s_n$  در راستای  $h \in R^d$ ، توزیع توأم  $Z(s_1), \dots, Z(s_n)$  تغییر نمی کند.

## تعریف ۴-۱

میدان تصادفی  $Z(\cdot)$ ، گاوسی<sup>۳</sup> نامیده می شود، هرگاه برای هر  $m \geq 1$  توزیع توأم  $(Z(s_1), \dots, Z(s_m))$  نرمال چند متغیره باشد.

---

2-Second order stationary

3-Strong stationary

1-Gaussian

## **فصل دوم**

### **روشهای استوار در تحلیل**

### **داده های فضایی**

## ۲-۱: تعریف داده فضایی

داده های فضایی، مشاهداتی هستند که برحسب موقعیت قرار گرفتشان در فضای مورد مطالعه به یکدیگر وابسته اند. آمار فضایی، شاخه ای از علم آمار است که به تحلیل چنین داده هایی می پردازد. برآورد تغییرنگار، برآورد پارامترهای روند و پیشگویی فضایی از مباحث مهم در آمار فضایی هستند، که معمولاً با پذیرش مفروضاتی صورت می پذیرند. انحراف از این پیش فرضها، ممکن است باعث بروز خطا و کاهش دقت در هر کدام از مباحث مورد اشاره شود. حضور داده های دورافتاده نیز یکی از مهمترین عوامل بروز این انحرافات است. روش های استوار فضایی قادرند تا حد زیادی اثرات، انحراف از پیش فرضها را در تحلیل داده های فضایی کنترل کنند. این فصل از سه بخش تشکیل شده است، در بخش اول تعاریف بررسی شده و در بخش دوم به علت اهمیت داده های دورافتاده شناسایی داده های دورافتاده مورد بررسی قرار گرفته و روش های استوار در حضور داده های دور افتاده در بخش سوم مورد توجه قرار گرفته است.

مشاهدات مبنای استنباط آماری هستند. اما علاوه بر مشاهدات، فرض هایی که درباره موقعیت تحت بررسی صورت می گیرد نقش پایه ای در استنباط آماری دارند. این پیش فرض ها در مورد توزیع، تصادفی بودن، وابستگی پارامترهای نامعلوم و... هستند. انتخاب پیش فرض ها به عنوان مکمل های یک مدل آماری نقش حساس و تعیین کننده ای در استنباط آماری دارند. انتخاب پیش فرض های مناسب، مدل های آماری را آسان و گاهی ممکن می سازند، نیازمند دقت زیاد هستند، ممکن است نتایج نادرستی حاصل شود. به این ترتیب در انتخاب با پیش فرض ها همان قدر که مفید بودن آنها و نقشی که در هموارسازی مسیر تحلیل دارند مورد توجه است باید همان میزان اعتبار آنها به عنوان پایه های استنباط آماری در نظر گرفته شود. هر چقدر دقت در انتخاب پیش فرض ها بیشتر باشد، اطمینان از معتبر بودن نتایج بالاتر خواهد بود. اما نباید فراموش کرد که پیش فرض های در نظر گرفته شده هیچ وقت قادر نخواهند بود طبیعت پدیده مورد بررسی را کامل کنند. اختلاف جزئی بین آنچه پیش فرض ها می گویند و آنچه واقعیت دارد، چیزی طبیعی و قابل اغماض است. ولی اگر تفاوت ها به اندازه ای باشد که صحت پیش فرض ها مورد تردید قرار گیرد نباید از آنها چشم

پوشی کرده در عمل می توان انحرافات زیادی از پیش فرض ها را متصور شد و روش هایی را فراهم آورد که حساسیت کمتری به این انحرافات داشته باشد.

با پیشرفت آمار و توجه به واقعیت های فوق تلاش ها برای یافتن روش هایی که بتوانند حساسیت به تخطی از پیش فرض ها را کنترل کنند، بیشتر شد. روش های ناپارامتری یکی از گزینه هایی بود که برای حل این مشکل مورد توجه قرار گرفت. در این روش ها هیچ پیش فرضی درباره توزیع جامعه ای که داده ها در آن آمده اند، در نظر گرفته نمی شود و نتایج آن هر چند با دقت کمتر اما برای تمامی توزیع جامعه قابل استفاده است. به علت این قابلیت بعضی از این روش ها کاربرد وسیعی در آمار یافتند با این وجود، در مواردی که عملاً استخراج اطلاعات در مورد ساختار جامعه از داده ها مقدور نمی باشد، روش های ناپارامتری ذهن برتری طلب یک آماردان را اغناء نمی کند و همواره با این پرسش مواجه خواهیم بود که آیا واقعاً نمی توان مدل هایی برای جامعه تحت بررسی در نظر گرفت؟ آیا اطلاعات زیادی که در مورد ساختار جامعه در دسترس ماست باید فقط به خاطر ترس از ناکامل بودن مدل، کنار گذاشته شود؟ آیا این یک نوع محافظه کاری افراطی نیست؟ آیا راه هایی برای استفاده از این اطلاعات وجود ندارد؟ این چنین است که نمی توان به سادگی مدل های پارامتری را با توجه ظرفیتشان در میان اطلاعات موجود در داده ها نادیده گرفت. بلکه به دنبال راهکارهایی برای کاهش حساسیت به فاصله پیش فرض ها از واقعیات جامعه گشت. به این صورت که مدل های پارامتری همچنان محور باشند و روش هایی مورد توجه قرار گیرند که وابستگی آنها به پیش فرض های مدل کمتر بحران ساز شود به این نیاز به روش های استوار<sup>۱</sup> بوجود آمد

## ۲-۱-۱: روش های آماری استوار

اصلی ترین تعریف روش های آماری استوار توسط هوبر (۱۹۸۱) ارائه شده است که کیفیت های مطلوب زیر را در نظر گرفت.

۱- باید تحت مدل فرض شده دارای کارایی بهینه یا نزدیک بهینه باشد

۲- عملکرد آنها تحت انحرافات کوچک از فرض های مدل به مقدار کمی دچار آسیب شود.

۳- بعضی از انحرافات بزرگ از مدل سبب فاجعه نشود.