



دانشگاه صنعتی اصفهان  
دانشکده علوم ریاضی

# آزمون‌های استوار در مدل رگرسیون لجستیک

پایان‌نامه کارشناسی ارشد (آمار اقتصادی و اجتماعی)

الهام همایونی

اساتید راهنمای پایان‌نامه

دکتر سروش علیمرادی  
دکتر علی رجالی



دانشگاه صنعتی اصفهان  
دانشکده علوم ریاضی

پایان نامه کارشناسی ارشد (آمار اقتصادی و اجتماعی) خانم الهام همایونی

تحت عنوان

## آزمون‌های استوار در مدل رگرسیون لجستیک

در تاریخ توسط کمیته تخصصی زیر مورد بررسی و تصویب نهائی قرار گرفت.

دکتر سروش علیمزادی

۱- استاد راهنمای پایان نامه

دکتر علی رجالی

۲- استاد راهنمای پایان نامه

۳- استاد داور ۱

()

۴- استاد داور ۲

دکتر اعظم اعتماد

سرپرست تحصیلات تکمیلی دانشکده



کلیه حقوق مادی مترتب بر نتایج مطالعات،  
ابتکارات و نوآوری‌های ناشی از تحقیق موضوع  
این پایان‌نامه متعلق به دانشگاه صنعتی  
اصفهان است.

# فهرست مطالب

۱	فصل اول مقدمه
۴	فصل دوم پیش نیازها
۷	۱-۲ داده پرت
۸	۲-۲ اندازه استواری یک برآورد
۸	۱-۲-۲ نقطه فروریزش
۹	۲-۲-۲ منحنی تأثیر
۱۰	۳-۲ مدل رگرسیونی
۱۱	۴-۲ برآوردگر حداقل مربعات
۱۱	۵-۲ برآوردگر $L_1$
۱۱	۶-۲ بررسی استواری روش $LS$ و $L_1$ نسبت به داده‌های پرت
۱۶	۷-۲ برآوردگرهای $M$
۱۶	۸-۲ برآوردگرهای $GM$
۱۷	۹-۲ مدل‌های خطی تعمیم یافته
۱۹	۱۰-۲ مقدار $p$ -مقدار
۱۹	۱۱-۲ تابع توان
۲۰	۱۲-۲ فاصله ماهالانویس استوار
۲۲	۱۳-۲ هم‌پوشانی بین داده‌ها
۲۵	فصل سوم برآوردگرها و آزمون‌های استوار در مدل لجستیک
۲۵	۱-۳ مقدمه
۲۶	۲-۳ مدل رگرسیون لجستیک

۲۸	.....	برآورد ضرائب مدل رگرسیون لجستیک به روش کلاسیک	۳-۳
۳۰	.....	برآورد ضرائب مدل رگرسیون لجستیک به روش استوار	۳-۴
۳۱	.....	برآوردگر بیانکو و یوهی	۳-۴-۱
۳۵	.....	برآوردگر بیانکو و یوهی وزنی	۳-۴-۲
۳۵	.....	خصوصیات مجانبی برآوردگر استوار بیانکو و یوهی وزنی	۳-۵
۳۹	.....	آماره آزمون والد استوار و رفتار مجانبی آن	۳-۶
۴۴	.....	بررسی پایداری $p$ -مقدار با استفاده از یک مثال	۳-۷
۴۶	.....	اندازه نیکویی برازش	۳-۸

#### فصل چهارم روش‌های بهینه‌سازی عددی

۵۰	.....	بهینه‌سازی نامقید	۴-۱
۵۴	.....	روش‌های درون‌یابی چندجمله‌ای	۴-۱-۱
۵۸	.....	روش تندترین شیب	۴-۱-۲

#### فصل پنجم شبیه‌سازی

۶۳	.....	روش مونت کارلو	۵-۱
۶۴	.....	پایداری سطح	۵-۱-۱
۶۸	.....	پایداری توان	۵-۱-۲

۷۵ ..... پیوست (برنامه‌های کامپیوتری)

۹۸ ..... مراجع

## چکیده:

کاربرد مدل رگرسیون لجستیک در طول دهه‌های اخیر پیشرفت‌های بسیاری داشته است. این مدل نخست در تحقیقات اپیدمیولوژی به کار می‌رفت، اما امروزه در بسیاری از زمینه‌ها از جمله مدیریت، اقتصاد، جرم‌شناسی، مهندسی پزشکی به کار برده می‌شود.

در این پایان‌نامه ابتدا مدل رگرسیون لجستیک معرفی می‌شود و برآوردهای کلاسیک به منظور برآورد ضرائب این مدل بیان می‌شود. با توجه به این‌که برآوردهای کلاسیک به شدت تحت تأثیر مشاهدات پرت قرار می‌گیرند، برآوردهای استوار معرفی می‌شوند. سطح و توان آزمون‌های کلاسیک با ورود مشاهده پرت فروریزش می‌کند، بدین ترتیب آزمون‌های استوار تعریف می‌شود. برای انجام آزمون فرض پارامتری، آماره آزمون والد استوار بر اساس برآوردهای بیانکو و یوهی وزنی بیان می‌شود و توزیع مجانبی آماره آزمون مورد مطالعه قرار می‌گیرد. با ارائه یک مثال عددی ابتدا برآوردهای استوار و کلاسیک و  $p$ -مقدار آزمون‌های استوار و کلاسیک محاسبه می‌شوند و سپس پایداری  $p$ -مقدار هر دو آزمون مورد بررسی قرار می‌گیرند. در پایان با انجام یک مطالعه شبیه‌سازی و استفاده از روش مونت‌کارلو پایداری سطح و توان آزمون مورد مطالعه قرار می‌گیرد.

رده‌بندی موضوعی: ۶۲۰۳۵.

کلمات کلیدی: رگرسیون لجستیک، استواری، آزمون فرض پارامتری، توزیع مجانبی

# فصل ۱

## مقدمه

مشکل دنیای امروز کمبود داده و اطلاعات کافی برای تصمیم‌گیری‌های علمی نیست، بلکه محققان در بیشتر زمینه‌های مطالعاتی با سیلی انبوه از داده‌های خام مواجه هستند که برای ارائه تحلیل‌های مفید و کارآمد نیازمند روشی مناسب برای استخراج اطلاعات از آن داده‌ها هستند. در بسیاری از موارد با داده‌هایی رو به رو می‌شویم که از نوع دودویی هستند. یک مدل متداول برای تحلیل داده‌ها، رگرسیون لجستیک است که در کلاس مدل‌های خطی تعمیم‌یافته<sup>۱</sup> قرار می‌گیرد و مدل‌سازی متغیر پاسخ غیر نرمال را ممکن می‌سازد. این مدل در بسیاری از زمینه‌ها مانند مطالعات اجتماعی و علوم پزشکی کاربرد دارد. از جمله متغیرهای پاسخ دودویی می‌توان به (مرگ و حیات)، (بهبود یا عدم بهبود بیماری)، (موافق یا مخالف بودن با یک موضوع) اشاره کرد.

در مدل رگرسیون لجستیک ضرائب مدل مجهول است. به منظور استفاده از مدل، لازم است ضرائب، برآورد شوند. یک روش متداول، برآوردگر ماکزیمم درست‌نمایی<sup>۲</sup> است، این برآوردگر در صورت نرمال بودن توزیع خطاها مناسب است اما اگر خطاها دارای توزیعی با دم طولانی باشد (داده پرت وجود داشته باشد) چندان کارا نیست. کروکس<sup>۳</sup> و همکاران [۱۷] رفتار فروریزشی  $MLE$  در مدل رگرسیون لجستیک مطالعه کرده‌اند و نشان داده‌اند که اضافه شدن داده پرت باعث فروریزش آن می‌شود. بنابراین آماردانان برآوردگرهای استوار را معرفی کردند که برای تأثیر مشاهدات پرت در تحلیل داده‌ها به کار می‌رود.

<sup>۱</sup> Generalized Linear Model (GLM)

<sup>۲</sup> Maximum Likelihood Estimator (MLE)

<sup>۳</sup> Croux



در چند دهه گذشته مطالعات زیادی برای به دست آوردن برآوردگرهای استوار برای پارامترها در مدل رگرسیون لجستیک و به طور کلی در چارچوب مدل های خطی تعمیم یافته، انجام شده است که می توان به کارهای پریجیبین<sup>۴</sup> [۳۲]، استفانسکی<sup>۵</sup> و همکاران [۴۰]، کانچ<sup>۶</sup> و همکاران [۲۹]، مورجنتلر<sup>۷</sup> [۳۰]، کارل<sup>۸</sup> و پدرسین<sup>۹</sup> [۱۰]، کریسمن<sup>۱۰</sup> [۱۴]، بیانکو<sup>۱۱</sup> و یوهی<sup>۱۲</sup> [۴]، کروکس و هائسبرک<sup>۱۳</sup> [۱۶] و باندل<sup>۱۴</sup> [۷] اشاره کرد.

بیانکو و یوهی [۴] کلاس جدیدی از برآوردگرهای  $M$  برای مدل رگرسیون لجستیک معرفی کردند. این برآوردگر نسبت به مشاهدات پرت در جهت  $y$  (متغیر پاسخ) استوار است. آن ها سازگاری و نرمال مجانبی این برآوردگر را نیز اثبات کردند.

کروکس و هائسبرک [۱۶] حالت وزنی برآوردگر بیانکو و یوهی را به منظور کاهش اثر مشاهدات پرت در جهت  $x$  (متغیر مستقل) ارائه دادند و الگوریتم سریع را برای محاسبه برآوردگر استوار برای مدل رگرسیون لجستیک معرفی کردند و شرایط وجود این برآوردگر را در نمونه های متنهای بیان کردند. همچنین نتایج استفاده از برآوردگر بیانکو و یوهی وزنی نیز مورد بررسی قرار داده اند و با شبیه سازی و مثال، این برآوردگرها را تشریح کردند.

برای آزمون پارامترهای رگرسیون لجستیک، آزمون هایی که بر اساس برآوردگرهای کلاسیک انجام شود تحت تأثیر مشاهدات پرت قرار می گیرد و سطح و توان آزمون پایدار نیست ولی آزمون هایی که بر اساس برآوردگرهای استوار هستند، سطح و توان پایدارتری نسبت به آزمون کلاسیک دارند. در این زمینه هریتیر<sup>۱۵</sup> و رنچتی<sup>۱۶</sup> [۲۵] آزمون های استوار برای مدل پارامتری تعمیم یافته را، شامل مدل رگرسیون لجستیک، معرفی کردند.

کانتونی<sup>۱۷</sup> و رنچتی<sup>۱۸</sup> [۹] انحراف های استوار بر اساس تابع شبه درستنمایی را تعریف کرده اند و یک

<sup>۴</sup> Pregibon

<sup>۵</sup> Stefanski

<sup>۶</sup> kunsch

<sup>۷</sup> Morgenthaler

<sup>۸</sup> Carroll

<sup>۹</sup> Pederson

<sup>۱۰</sup> Christmann

<sup>۱۱</sup> Bianco

<sup>۱۲</sup> Yohai

<sup>۱۳</sup> Haesbroeck

<sup>۱۴</sup> Bondell(2005,2008)

<sup>۱۵</sup> Heritier

<sup>۱۶</sup> Ronchetti

<sup>۱۷</sup> Cantoni

<sup>۱۸</sup> Ronchetti

خانواده از آماره آزمون‌ها را برای انتخاب مدل در مدل‌های خطی تعمیم یافته، معرفی نموده‌اند. در این پایان‌نامه مدل رگرسیون لجستیک معرفی می‌شود و نحوه‌ی برآورد ضرائب آن مدل ( $\beta$ ) به روش کلاسیک و برآوردگر استوار بررسی می‌شود و به دلیل عدم استواری برآوردگر کلاسیک برآوردگرهای استوار معرفی می‌شود. از آنجا که به آزمون‌های استوار توجه کمتری نسبت به برآوردگر استوار شده است، آزمون فرض پارامتری، پارامترهای رگرسیون لجستیک بررسی می‌شود و آماره والد بر اساس نوع وزنی برآوردگر بیانکو و یوهی که توسط کروکس و هائسبرک معرفی شد [۱۶]، بیان می‌شود. هدف بیان شکل استوار آماره آزمون والد کلاسیک است. این آزمون به شکل مربعی و بر اساس برآوردگرهای استوار پارامترهای رگرسیونی و ماتریس کواریانس مجانبی آن است. رفتار مجانبی این آماره آزمون مشابه حالت کلاسیک است و دارای توزیع کای دو است.

مطالب گرد آوری شده در قالب فصل‌های زیر ارائه می‌شود.

در فصل ۲، بعضی از تعاریف و قضایایی که در فصل‌های بعد مورد نیاز است مطرح می‌شود. در فصل ۳، ابتدا برآوردگر ماکزیمم درست‌نمایی معرفی می‌شود و با توجه به آن که این برآوردگر در مقابل مشاهدات پرت استوار نیست، برآوردگر بیانکو و یوهی معرفی خواهد شد. شرایط وجود این برآوردگر با استفاده از یک گزاره مطرح شده و بعد از آن نحوه‌ی محاسبه این برآوردگر توضیح داده می‌شود. این برآوردگر در مقابل مشاهدات پرت در جهت  $y$  استوار است، ولی در مقابل مشاهدات پرت در جهت  $x$  استوار نیست. بدین منظور برای کاهش اثر مشاهدات پرت در جهت  $x$ ، برآوردگر وزنی بیانکو و یوهی مورد استفاده قرار گرفته و خصوصیات مجانبی آن بیان می‌شود. بعد از آن آماره آزمون برای این برآوردگر معرفی و رفتار مجانبی آن بررسی می‌شود. در پایان با ارائه یک مثال عددی برآوردگرهای استوار و کلاسیک و  $p$ -مقدار آزمون کلاسیک و استوار محاسبه می‌شود و پایداری  $p$ -مقدار بررسی می‌شود.

در فصل ۴ روش‌های عددی بهینه‌سازی بیان می‌شود.

در فصل ۵، با استفاده از روش مونت کارلو، پایداری سطح و توان آزمون‌های استوار و کلاسیک مورد مقایسه قرار می‌گیرند. برنامه مربوط به تهیه جدول‌ها و نمودارها در پیوست آورده شده است.

## فصل ۲

### پیش نیازها

در این فصل به منظور یادآوری، برخی از تعاریف، قضایا و اصطلاحات مورد نیاز در این پایان نامه، بیان می شوند.

تعریف ۱.۲ سه تایی  $(\Omega, \mathcal{A}, P)$  یک فضای احتمال است و  $X : \Omega \rightarrow \mathbb{R}$  یک متغیر تصادفی است (یعنی تابع اندازه پذیر روی  $(\Omega, \mathcal{A}, P)$  است).

قضیه ۱.۲ نامساوی مارکوف [۴۱]. اگر  $X$  یک متغیر تصادفی و  $g$  تابعی غیر منفی باشد و  $E[g(X)] < \infty$ ، آن گاه برای هر  $a > 0$ ،

$$P(g(X) \geq a) \leq \frac{E[g(X)]}{a}$$

قضیه ۲.۲ نامساوی چبیشف<sup>۱</sup> [۴۱]. با انتخاب  $g(x) = |x|$ ، در نامساوی مارکوف و با توجه به این که  $P(|X| \geq a) = P(X^2 \geq a^2)$ ، نامساوی زیر برقرار است.

$$P(|X| \geq a) \leq \frac{E[X^2]}{a^2} \quad a > 0$$

---

<sup>۱</sup> Chebyshev Inequality

حالت خاصی از نامساوی چبیشف با انتخاب  $g(X) = |X - \mu|$  به صورت زیر به دست می آید:

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}, \quad k > 0$$

قضیه ۳.۲ (نامساوی کوشی شوارتز)<sup>۲</sup> [۴۲]. اگر  $X$  و  $Y$  دو متغیر تصادفی با تابع توزیع توأم  $F(X, Y)$  باشند بطوری که  $E(|X|^2)$  و  $E(|Y|^2)$  متناهی هستند، آنگاه:

$$|E(XY)| \leq E(|XY|) \leq \sqrt{E(X^2)E(Y^2)}$$

حالت دیگری از نامساوی شوارتز، به فرم زیر است:

$$Cov^2(X, Y) \leq Var(X)Var(Y).$$

قضیه ۴.۲ (نامساوی جنسن)<sup>۳</sup> [۴۱]. اگر  $\varphi$  یک تابع محدب باشد که روی بازه‌ای شامل برد  $X$  تعریف شده است و اگر امید ریاضی  $X$  و  $\varphi(X)$  هر دو موجود باشند، آنگاه

$$\varphi(E(X)) \leq E(\varphi(X))$$

تعریف ۲.۲ [۱۳]. دنباله  $\{X_n\}$  از متغیرهای تصادفی روی فضای احتمال  $(\Omega, \mathcal{A}, P)$  تقریباً همه جا<sup>۴</sup> (a.s) به متغیر تصادفی  $X$  همگراست، هرگاه یک مجموعه  $N$  با احتمال صفر وجود داشته باشد که به ازای هر  $w \in \Omega - N$

$$\lim_{n \rightarrow \infty} X_n(w) = X(w)$$

این همگرایی با  $X_n \xrightarrow{a.s.} X$  نمایش داده می شود.

<sup>۲</sup> Cauchy Schwartz

<sup>۳</sup> Jensen

<sup>۴</sup> almost surely

تعریف ۳.۲ [۱۳]. دنباله‌ی متغیرهای تصادفی  $\{X_n\}$  در احتمال به  $X$  همگراست هر گاه به ازای هر  $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} P[|X_n - X| > \varepsilon] = 0$$

این همگرایی با  $X_n \xrightarrow{p} X$  نمایش داده می‌شود.

تعریف ۴.۲ [۱۳]. هر گاه  $\{F_n\}$  دنباله‌ای از توابع توزیع برای دنباله  $\{X_n\}$ ، و  $F$  تابع توزیع متغیر تصادفی  $X$  باشد و به ازای هر نقطه پیوستگی  $F$ ،  $\lim_{n \rightarrow \infty} F_n(x) = F(x)$ ، گفته می‌شود که  $X_n$  در توزیع به  $X$  همگراست و با  $X_n \xrightarrow{D} X$  نمایش داده می‌شود. (اگر و تنها اگر  $P(X = x) = 0$  به ازای تمام مقادیر که برای آن‌ها  $\lim_n P[X_n \leq x] = P(X \leq x)$ )

قضیه ۵.۲ (قانون ضعیف اعداد بزرگ<sup>۵</sup>) [۶]. اگر  $X_1, X_2, \dots$  متغیرهای تصادفی مستقلی (نه لزوماً با توزیع یکسان) هر یک با میانگین متناهی و واریانس متناهی باشند. علاوه بر آن واریانس‌ها به طور یکنواخت کراندار باشند، در این صورت وقتی برای هر  $n$ ،  $S_n = X_1 + \dots + X_n$ ، آن گاه  $[S_n - E(S_n)]/n$  در احتمال به صفر همگراست. یعنی به ازای هر  $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} P\left\{ \left| \frac{S_n - E(S_n)}{n} \right| \geq \varepsilon \right\} = 0$$

قضیه ۶.۲ (قانون قوی اعداد بزرگ<sup>۶</sup>) [۶]. اگر  $\{X_n\}$ ‌ها یک دنباله‌ای از متغیرهای تصادفی مستقل باشند، به نحوی که برای هر  $n$ ،  $E[X_n] = 0$ ، در این صورت هر گاه  $\sum_n \frac{var(X_n)}{n^2} < \infty$ ، آن گاه

$$n^{-1} \sum_{k=1}^n X_k \xrightarrow{a.s.} 0$$

قانون قوی اعداد بزرگ در حالتی که  $X_i$ ‌ها مستقل و هم توزیع باشند، به صورت زیر است:

قضیه ۷.۲ [۶]. اگر متغیرهای تصادفی  $X_1, X_2, \dots$  مستقل و دارای توزیع یکسانی باشند (*i.i.d*) و

$$E[X_i] = \mu \text{ در این صورت } n^{-1} \sum_{k=1}^n X_k \xrightarrow{a.s.} \mu$$

<sup>۵</sup> Weak Law of Large Numbers

<sup>۶</sup> Strong Law of Large Numbers

قضیه ۸.۲ (قضیه اسلاتسکی<sup>۷</sup>). اگر  $\{X_n, Y_n\}_{n \geq 1}$  دنباله‌ای از جفت متغیرهای تصادفی و  $c$  یک مقدار ثابت باشد،  $Y_n \xrightarrow{p} c$  و  $X_n \xrightarrow{D} X$ ، آن گاه روابط زیر برقرارند: [۲۲]

$$X_n + Y_n \xrightarrow{D} X + c \quad (۱)$$

$$\begin{cases} X_n Y_n \xrightarrow{L} cX, & c \neq 0 \\ X_n Y_n \xrightarrow{p} 0, & c = 0 \end{cases} \quad (۲)$$

$$\frac{X_n}{Y_n} \xrightarrow{D} \frac{X}{c}, c \neq 0 \quad (۳)$$

قضیه ۹.۲ [۱۹]. (قضیه همگرایی غالب<sup>۸</sup>).

اگر  $X_n \xrightarrow{a.s.} X$ ، و برای تمام  $n$  ها، رابطه زیر برقرار باشد:

$$|X_n| \leq |Y|, \quad E(|Y|) < \infty$$

آن گاه  $X$  هم انتگرال پذیر است و  $E(X) = \lim_{n \rightarrow \infty} E(X_n)$

تعریف ۵.۲ دنباله‌ی برآوردگرهای  $\{T_n\}$  برای تابعی از پارامتر  $\theta$  ( $\Gamma(\theta)$ ) سازگار است، هر گاه  $T_n \xrightarrow{p} \Gamma(\theta)$  وقتی  $n$  به سمت بی‌نهایت میل کند. یعنی برای هر  $\varepsilon > 0$  و  $\delta > 0$ ، عدد صحیح مثبتی مانند  $n_0 = n_0(\varepsilon, \delta)$  وجود داشته باشد به طوری که برای هر  $n > n_0$ ،  $P_\theta(|T_n - \Gamma(\theta)| > \varepsilon) < \delta$  یا برای هر  $\varepsilon > 0$ ،  $\lim_{n \rightarrow \infty} P_\theta(|T_n - \Gamma(\theta)| > \varepsilon) = 0$  [۴۱]

## ۱-۲ داده پرت

مشاهده دورافتاده یا پرت<sup>۹</sup> داده‌ای است که از توده اکثریت مشاهدات دور باشد. این گونه مشاهدات پس از شناسایی باید به دقت مورد بررسی قرار گیرند تا معلوم شود که آیا دلیلی برای رفتار غیر عادی آن‌ها وجود

<sup>۷</sup> Slutsky

<sup>۸</sup> Dominated Convergence Theorem

<sup>۹</sup> Outlier

دارد یا خیر. وجود داده‌های پرت گاهی غیر عادی، اما قابل توضیح است. اندازه‌گیری غلط، ثبت و یا انتقال نادرست داده‌ها، از کار افتادن وسیله اندازه‌گیری، مثال‌هایی از این اتفاقات هستند. در چنین مواردی بایستی در صورت امکان، اتفاق مذکور اصلاح و در غیر این صورت آن داده از مجموعه داده‌ها حذف شود. تاکید می‌شود که قبل از خارج کردن داده مذکور بایستی دلایل منطقی کافی برای نامناسب بودن آن داده دور افتاده وجود داشته باشد.

بعضی اوقات داده دور افتاده غیر عادی است، اما آن مشاهده‌ای کاملاً معتبر است. وجود چنین داده‌هایی می‌توانند ناکارآمدی مدل را نشان دهند و اغلب سرنخ‌های با ارزشی جهت ساختن مدل رگرسیونی مناسب‌تر برای تحلیل‌گر فراهم می‌کنند. گاهی نقاط پرت اثر نامتناسبی را بر مدل رگرسیونی اعمال می‌کنند، بدین معنی که برآورد پارامترها تحت تاثیر مجموعه‌ای کوچک از داده‌ها قرار می‌گیرد. این داده‌ها را نقاط تاثیرگذار<sup>۱۰</sup> گویند. با این توضیحات، یک تحلیل‌گر باید داده‌های پرت را شناسایی کند و از تأثیرات آن بر نتایج و تحلیل‌ها آگاه باشد.

تعریف ۶.۲ عدم حساسیت به انحرافات کوچک از فرض‌ها، استواری گفته می‌شود [۳۶].

## ۲-۲ اندازه استواری یک برآورد

دو معیار سنجش استواری به صورت زیر است:

- (۱) نقطه فروریزش بیانگر آن است که چگونه برآوردگر توسط کسری از داده‌ها تحت تاثیر قرار می‌گیرد.
- (۲) منحنی تاثیر که با استفاده از آن، می‌توان اطلاعاتی از تاثیر داده پرت بر برآوردگر را به دست آورد.

### ۱-۲-۲ نقطه فروریزش

ابتدا مفهوم آریبی تعریف می‌شود و سپس با استفاده از آن مفهوم نقطه فروریزش ارائه می‌گردد.

تعریف ۷.۲ اگر  $Z = \{(x_{i1}, x_{i2}, \dots, x_{ip}, y_i), i = 1, \dots, n\}$  یک نمونه  $n$  تایی و  $T$  یک برآوردگر رگرسیونی باشد (یعنی با اعمال  $T$  بر نمونه  $Z$ ، برآورد ضرایب معادله رگرسیون  $T(Z) = \hat{\beta}$  حاصل شود) و تمام نمونه‌های ممکن  $Z'$  که توسط جایگزینی هر  $m$  تا از نقاط داده‌های اصلی توسط مقادیر دلخواه (حتی مشاهدات خیلی پرت) با  $\mathbb{Z}$  نمایش داده می‌شود، آن‌گاه

<sup>۱۰</sup>Influential points

$$\text{bias}(m; T, Z) = \text{Sup}_{Z' \in \mathcal{Z}} \|T(Z') - T(Z)\|$$

که در آن  $\|\cdot\|$  نرم اقلیدسی است، ماکزیمم اریبی نامیده می‌شود [۳۶].

اگر این اریبی نامتناهی باشد، بدین معنی است که  $m$  داده پرت می‌تواند اثر زیادی بر  $T$  داشته باشند. در این صورت می‌توان این مسأله را با جمله «برآوردگر به سمت پایین فرو می‌ریزد» توصیف کرد.

تعریف ۸.۲ [۳۶]. نقطه فروریزش برآوردگر  $T$  براساس نمونه  $n$  تایی  $Z$  به صورت زیر تعریف می‌شود

$$\epsilon_n^*(T, Z) = \min\left\{\frac{m}{n}; \text{bias}(m; T, Z) = \infty\right\}$$

به عبارت دیگر، این معیار کوچکترین کسر عدم خلوص است که باعث می‌شود برآوردگر  $T$  از  $T(Z)$  دور شود. هر چه این کسر بزرگ‌تر باشد، نشان دهنده استوارتر بودن برآوردگر  $T$  است.

## ۲-۲-۲ منحنی تأثیر

اگر  $F$  یک تابع توزیع  $k$  بعدی و  $\theta$  پارامتری از جامعه باشد که  $F$  به آن وابسته است، آن‌گاه می‌توان نوشت  $\theta = T(F)$  که در این صورت  $T$  را تابع آماری نامند. ساده‌ترین مثال از تابع آماری میانگین متغیر تصادفی  $X$ ،  $E_F[X]$  است.

$$E_F[X] = T[F] = \int x dF(x)$$

می‌توان مقادیر کوچک ناخالصی تابع توزیع  $F$  در  $z_0 = (x_0, y_0)$  را با در نظر گرفتن یک توزیع آمیخته به صورت زیر مدل‌سازی کرد:

$$F_t = (1-t)F + t\delta_{z_0}$$

که در آن  $\delta_{z_0}$  تابع توزیع مقادیر ناخالصی  $z_0$  است (صفحه ۸۲ از کتاب سیبر [۳۹]). منحنی تأثیر به صورت زیر تعریف می‌شود.

تعریف ۹.۲ منحنی تأثیر یک تابع آماری  $T$ ، مشتق  $T(F_t)$  نسبت به  $t$  در نقطه  $t = 0$  است که میزان واکنش  $T$  نسبت به مقادیر ناخالصی کوچک  $z_0$  را نشان می‌دهد [۳۹].

$$IC(F, z_0) = \left. \frac{dT(F_t)}{dt} \right|_{t=0}$$



نکته ۱۰.۲ اگر  $IC$  نسبت به  $x(y)$  کران دار باشد نتیجه می شود که نسبت به  $x(y)$  استوار است [۳۹].

تعریف ۱۰.۲ برآوردگر استوار برآوردگری است که با ورود داده پرت به نمونه فروریزش نمی کند [۳۹].

تعریف ۱۱.۲ روش رگرسیون استوار به روشی گفته می شود که نه تنها در صورت نرمال بودن توزیع خطاها خوب عمل می کند، بلکه نسبت به انحراف های کوچک از فرض نرمال بودن نیز حساسیت اندکی دارد [۲۷].

## ۲-۳ مدل رگرسیونی

یکی از روش های کاربردی برای تحلیل داده ها در بین ابزارهای آماری، تحلیل رگرسیونی است. تحلیل رگرسیونی، روشی کارآمد برای بررسی و مدل سازی ارتباط بین متغیرها است که از آن در توصیف داده ها، برآورد پارامترهای مجهول، پیش گویی و کنترل استفاده می شود. به منظور دستیابی به یک مدل رگرسیونی به مشاهداتی از متغیرهای مستقل (توضیحی)  $\mathbf{x} = (x_1, x_2, \dots, x_p)$  و نیز مشاهداتی از متغیر وابسته (پاسخ)  $y$  نیاز است. ساده ترین حالت که در آن تنها یک متغیر مستقل  $x$  وجود داشته باشد، عبارت است از وقتی که رابطه بین این دو متغیر یک رابطه خطی به صورت زیر بین مشاهدات وجود داشته باشد:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad i = 1, 2, \dots, n \quad (1-2)$$

که در آن  $\beta_0$  عرض از مبدا و  $\beta_1$  شیب خط است و مولفه های  $\epsilon_i$ ، به عنوان خطای تصادفی تلقی می شوند. بدین معنی که به عنوان مقادیر مشاهده شده متغیری تصادفی، اندازه ناتوانی مدل در برازش دقیق داده ها را اندازه گیری می کنند. این خطا ممکن است به دلیل عدم حضور برخی از متغیرهای موثر، خطاهای تصادفی مربوط به مشاهدات و اندازه گیری ها و موارد دیگر باشد [۳۶].

روش های مختلفی برای برآورد پارامترهای یک مدل رگرسیونی وجود دارد که در زیر تعدادی از آنها معرفی می شوند.

## ۲-۴ برآوردگر حداقل مربعات

معمول‌ترین روش در برآورد پارامترهای یک مدل خطی استفاده از روش کمترین مربعات است. اگر توزیع خطاها نرمال باشد برآوردگر حداقل مربعات<sup>۱۱</sup>  $\beta$  به گونه‌ای انتخاب می‌شود که مجموع توان دوم انحراف‌ها یعنی

$$\sum_{i=1}^n \varepsilon_i^2 \quad (2-2)$$

را کمینه کند. منحنی تأثیر  $LSE$  عبارت است از  $IC(z_0, F) = \Sigma_F^{-1} x_0 [y_0 - x_0' T(F)]$  که در آن  $T(F) = \{\Sigma_F\}^{-1} \gamma_F$ ،  $\Sigma_F = E_F[XX']$ ،  $\gamma_F = E_F[XY]$ ، وقتی  $F$  تابع چگالی  $X, Y$  یک بردار تصادفی و  $z_0 = (x_0, y_0)$  مقدار توأم مشاهده شده برای  $(X, Y)$  است. منحنی تأثیر نسبت به هر دو مقدار  $X_0$  و  $Y_0$  کران دار نیست و این بدین معنی است که  $LSE$  استوار نیست [۳۶] و [۳۹].

۲-۵ برآوردگر  $L_1$ 

با قرار دادن  $|\varepsilon_i|$  به جای  $\varepsilon_i^2$  در رابطه‌ی (۲-۲) برآوردگر  $L_1$  به دست می‌آید. در واقع ساده‌ترین روش در برآورد پارامترهای یک مدل رگرسیونی، استفاده از کمترین قدرمطلق انحرافات است که در مقایسه با روش حداقل مربعات دارای حساسیت کمتری نسبت به داده‌های پرت است، با این حال نمی‌توان از برآوردگر  $L_1$  به عنوان جایگزین قوی برای برآوردگر حداقل مربعات استفاده کرد، زیرا به شدت توسط تک مشاهده پرت در جهت  $x$  (داده بانفوذ) تحت تأثیر قرار می‌گیرد، [۳۶] که این موضوع در بخش بعد بررسی می‌شود.

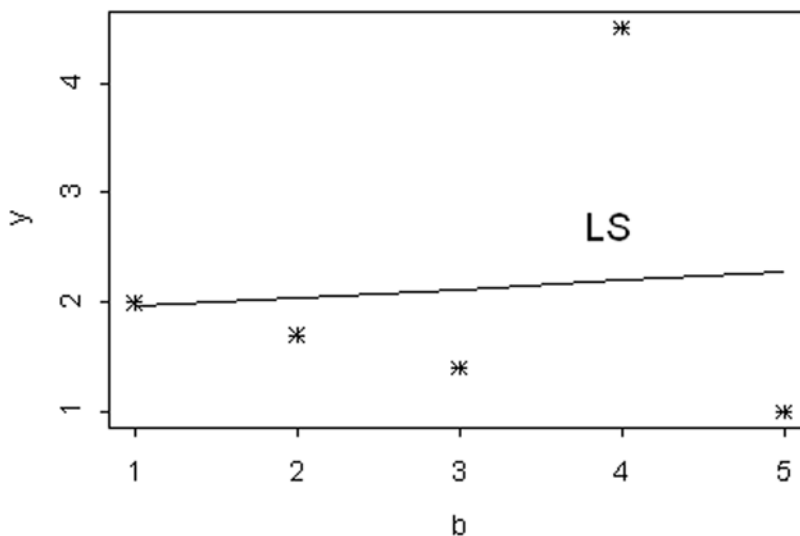
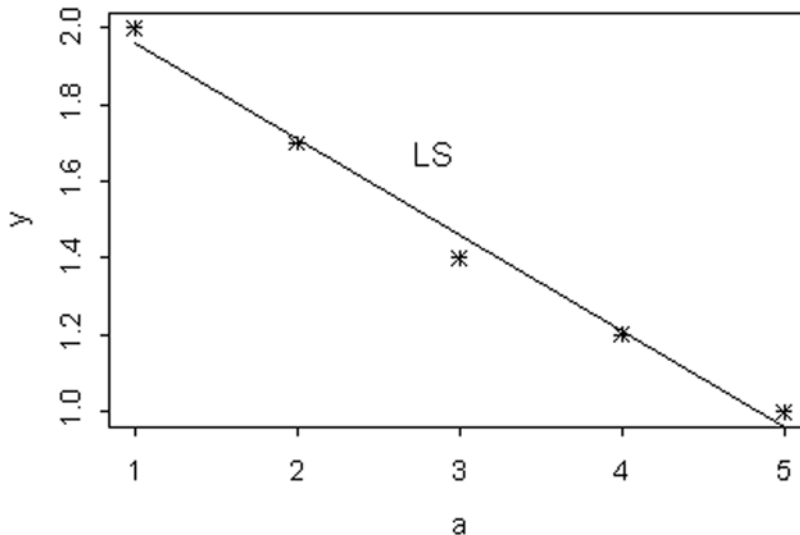
۲-۶ بررسی استواری روش  $LS$  و  $L_1$  نسبت به داده‌های پرت

برای بررسی اثر داده‌های پرت، به مدل رگرسیونی ساده می‌توان توجه کرد. نمودار وضعیت داده‌ها نسبت به یکدیگر را به سادگی می‌توان در یک نمودار پراکنش نمایش داد. شکل ۱.۲a نمودار پراکنش<sup>۱۲</sup> پنج نقطه است که تقریباً روی یک خط مستقیم قرار می‌گیرند و خط رگرسیونی کمترین مربعات ( $LS$ ) به خوبی به آن‌ها برازش می‌یابد. حال اگر مشاهده چهارم اشتباه گزارش داده‌شده باشد و یا نوشته شود. در این صورت نقطه  $(x_4, y_4)$  از بقیه مشاهدات و نیز از خط ایده‌آل دور می‌افتد. شکل ۱.۲b چنین موقعیتی را نشان می‌دهد. این نقطه یک «داده پرت در مسیر  $y$ » نامیده می‌شود و اثر نسبتاً زیادی بر خط  $LS$  دارد، چرا که

<sup>۱۱</sup>Least Squares Estimator(LSE)

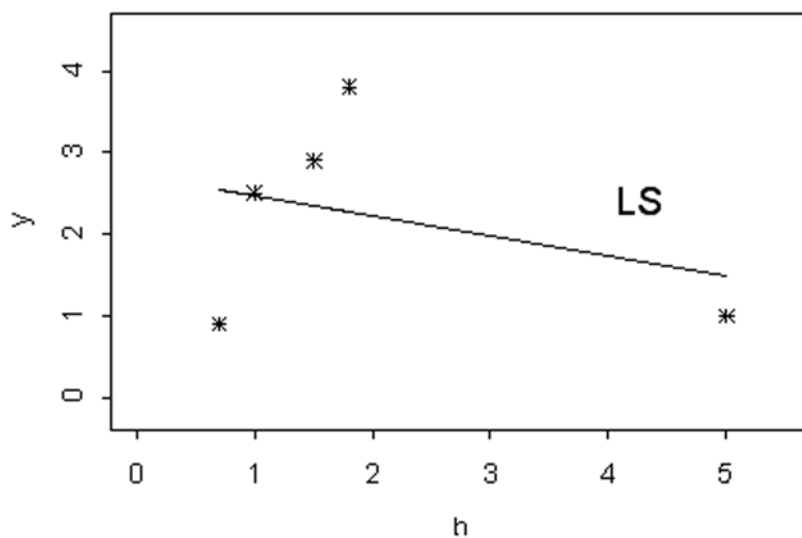
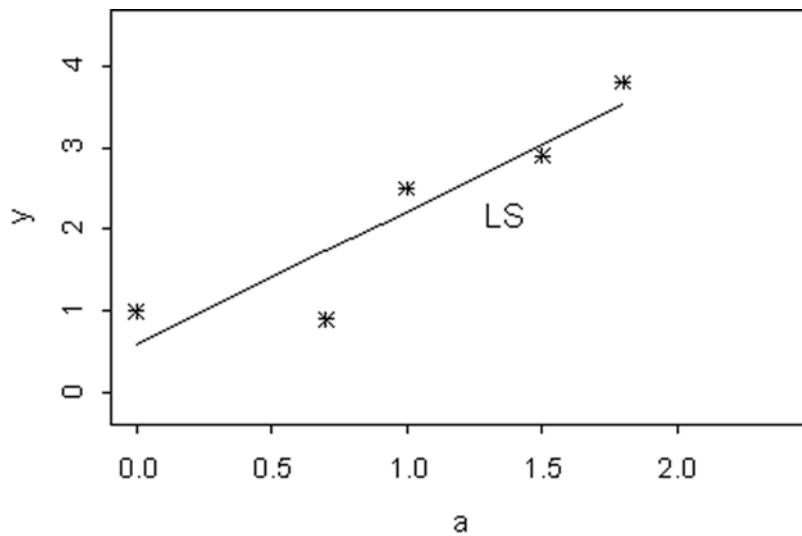
<sup>۱۲</sup>Scatter Plot

خط برازش شده در این حالت با خط ایده آل در شکل ۱.۲a تفاوت دارد [۳۶]. همچنین این داده دارای باقیمانده بزرگی است. در رگرسیون چندگانه (پیش از یک متغیر توضیحی) طریق مقادیر باقیمانده‌ها شناسایی می‌شوند [۴۴].



شکل ۱.۲. تاثیر داده پرت در جهت متغیر پاسخ بر برآورد کمترین مربعات

امکان وجود داده پرت در متغیرهای توضیحی وجود دارد. پس موقعیتی در نظر گرفته می‌شود که متغیر توضیحی شامل یک داده پرت باشد. شکل ۲.۲a حاوی پنج نقطه با برازش خوب  $LS$  است. اکنون فرض می‌شود که در مقدار  $x_1$  تغییری رخ دهد. شکل ۲.۲b نشان دهنده چنین وضعیتی است. نقطه  $(x_1, y_1)$  « داده پرت در مسیر  $x$  » نامیده می‌شود و اثر آن روی خط رگرسیون بسیار زیاد است [۳۶].



شکل ۲.۲. تاثیر داده اهرمی بر برآورد کمترین مربعات

چون وضعیت قرار گرفتن نقاط در فضای  $x$ ، در تعیین خواص مدل اهمیت دارد، به همین دلیل، معمولاً نقاط پرت فضای  $x$  را نقاط بانفوذ<sup>۱۳</sup> می‌نامند. چنین نقاطی معمولاً باقی‌مانده کمی دارند چرا که خط رگرسیونی را به سمت خود کشیده‌اند. بنابراین از طریق بررسی باقی‌مانده‌ها شناسایی نمی‌شوند. تاثیرگذاری نقاط پرت را می‌توان با کنار گذاشتن این نقاط و برازش مجدد معادله رگرسیونی و مقایسه مقادیر به دست آمده جدید و قبلی، مشخص کرد [۴۴].

زمانی که یک نقطه پرت در نزدیکی خط رگرسیونی که توسط اکثریت داده‌ها تعیین می‌شود قرار گیرد،

<sup>۱۳</sup>Leverage Point