

# استفاده از توزیع اسلش بجای توزیع نرمال در مدل‌های رگرسیونی چندگانه

زهرا نیکنام

دانشکده آمار، دانشگاه اصفهان

## چکیده

در بسیاری از تحلیل داده‌های آماری فرض متداول نرمال بودن توزیع مشاهدات است. اما نقض این فرض در تحلیل داده‌های واقعی امکان پذیر است. به همین دلیل کلاس‌های جدیدی از توزیع‌ها که برای داده‌های مختلف دارای انعطاف‌پذیری بیشتری هستند، به عنوان جایگزین توزیع نرمال پیشنهاد شده است. در این رابطه می‌توان به توزیع اسلش که یک توزیع دم-سنگین و متقارن است، اشاره کرد که در دهه‌ی اخیر از سوی پژوهشگران زیادی مورد توجه قرار گرفته است. در این مقاله با استفاده از مجموعه داده‌های واقعی که در متون مربوط به مبحث رگرسیون با توزیع نرمال برازش داده شده است ما با در نظر گرفتن مؤلفه‌های خطا با توزیع اسلش به اهمیت این توزیع در مباحث رگرسیونی می‌پردازیم.

واژه‌های کلیدی: توزیع آمیخته-مقیاس نرمال، توزیع دم سنگین، توزیع اسلش، رگرسیون.

## ۱ مقدمه

توزیع‌های متفاوتی به جای نرمال انجام داده اند که از آن جمله می‌توان به توزیع اسلش<sup>۱</sup> اشاره کرد. توزیع اسلش که علاوه بر داشتن دم‌های سنگین‌تر، نسبت به نرمال از کشیدگی بیشتری نیز برخوردار است، در استنباط‌های استوار<sup>۲</sup> دارای اهمیت است. توزیع اسلش همانند توزیع

بسیاری از استنباط‌های آماری با فرض نرمال بودن داده‌ها صورت می‌گیرد. در صورتی که مشاهدات تأثیرگذاری در تحلیل داده‌ها موجود نباشند که فرض نرمال بودن را خدشه‌دار سازند این کاربردها به نتایج معتبری منجر می‌شوند. محققان تحقیقات زیادی جهت جایگزین کردن

<sup>۱</sup> Slash distribution

<sup>۲</sup> Robust

یک متغیره و چند متغیره پرداخته و ویژگی‌های آن را مورد بررسی قرار داده‌اند. جنس [۶] تعمیمی از توزیع اسلش را از طریق توزیع نرمال-بتا معرفی کرد و ضرورت عملی آن را با مجموعه ای از داده‌های واقعی با در نظر گرفتن روش برآورد حداکثر درستنمایی تشریح کرد. گومز و همکاران [۸] توزیع اسلش اصلاح شده را معرفی و ویژگی‌های آن را مورد بررسی قرار دادند. با توجه به اینکه توزیع اسلش به عنوان آمیخته مقیاسی از توزیع نرمال می‌باشد ابتدا در بخش ۲ به معرفی خانواده توزیع‌های آمیخته-مقیاس نرمال می‌پردازیم و در بخش ۳ به بررسی توزیع اسلش و ویژگی‌های آن می‌پردازیم و در بخش ۴ مدل رگرسیون خطی چندگانه با خطای نرمال مرور می‌کنیم و در بخش ۵ به معرفی مدل رگرسیون خطی چندگانه با خطای اسلش می‌پردازیم و در پایان نیز اهمیت کاربرد این توزیع را با یک مثال تجربی نشان می‌دهیم.

## ۲ خانواده توزیع‌های آمیخته-مقیاس نرمال (SMN)

توزیع‌های آمیخته-مقیاس نرمال<sup>۱۱</sup> که نقش بسیار مهمی در مدل‌سازی آماری دارند داده‌های دور افتاده را پوشش می‌دهند و دم‌هایی سنگین تر از نرمال دارند. مدل‌های آمیخته-مقیاس شامل توزیع‌های تی-استیودنت، اسلش و واریانس-گاما و توزیع پیرسون نوع هفتم است. توزیع‌های آمیخته-مقیاس نرمال توسط آمیختن یک متغیر تصادفی نرمال،  $Z$  با یک متغیر تصادفی نامنفی مقیاسی  $\lambda$  به صورت زیر حاصل می‌شوند:

$$X = \mu + k^{\frac{1}{\lambda}}(\lambda)Z,$$

نرمال استاندارد حول مبدأ متقارن است، اما از دمه‌های سنگینی نسبت به آن برخوردار است که همین ویژگی سبب وجود مقادیری در میان مشاهدات نمونه می‌شود که دور افتاده بودن آنها نسبت به مشاهدات حاصل از توزیع نرمال استاندارد کاملاً واضح است. توزیع نرمال استاندارد برای حالتی که نمونه تصادفی فاقد مقادیر دور افتاده باشد و توزیع اسلش برای حالتی که نمونه تصادفی از مقادیر دور افتاده برخوردار باشد مناسبند. تعریف‌های متعددی از توزیع اسلش یک متغیره و چندمتغیره توسط محققان صورت گرفته است. توزیع اسلش و بررسی ویژگی‌های آن توسط محققان زیادی از جمله روگروز و توکی<sup>۳</sup> [۱۶]، مورجنتالر<sup>۴</sup> [۱۴]، جمشیدیان<sup>۵</sup> [۹]، کاشید و کالکارنی<sup>۶</sup> [۱۲] صورت گرفته است. کفادار<sup>۷</sup> [۱۱] برآورد حداکثر درستنمایی پارامترهای مقیاسی و مکانی را برای توزیع اسلش استاندارد به دست آورد. بسط توزیع اسلش یک متغیره به حالت چند متغیره اولین بار توسط لانگ و سینشایمر<sup>۸</sup> [۱۳] انجام شد. جنس<sup>۹</sup> [۵] توزیع اسلش تعمیم یافته یک متغیره را با استفاده از آمیخته-مقیاس توزیع نمایی-توانی به دست آورد و خصوصیات آن را مورد بررسی قرار داد. گومز و همکاران<sup>۱۰</sup> [۷] خانواده جدیدی از توزیع‌های اسلش یک متغیره و چند متغیره که ساختار آنها بر اساس توزیع‌های بیضوی است را معرفی کردند. اخیراً ارسلان و جنس [۳] به معرفی توزیع اسلش

<sup>۳</sup>Rogers and Tukey

<sup>۴</sup>Morgenthaler

<sup>۵</sup>Jamshidian

<sup>۶</sup>Kashid and Kulkarni

<sup>۷</sup>Kafadar

<sup>۸</sup>Lang and Sinsheimer

<sup>۹</sup>Genç

<sup>۱۰</sup>Gomez and et al.

<sup>۱۱</sup>Scale Mixtures of Normal

میانگین همه ی توزیع های متعلق به این خانواده برابر  $\mu$  است ولی چون پارامتر مقیاس تصادفی است واریانس ها متفاوت است. (آندره و مالوز<sup>۱۲</sup> [۲])

### ۳ توزیع اسلش

توزیع اسلش به عنوان توزیع متقارن و انعطاف پذیر نقش مهمی در مطالعات استوار دارد. توزیع اسلش دم هایی سنگین تر از توزیع نرمال استاندارد دارد اما دم های آن شبیه به توزیع کوشی است که این امر در شکل ۱ به وضوح دیده می شود.

توزیع اسلش و بررسی ویژگی های آن توسط محققان زیادی از جمله راجرز و توکی [۱۶]، مورجنتالر [۱۴]، جمشیدیان [۹]، کاشید و کالکارنی [۱۲] صورت گرفته است. سپس کفادار [۱۱] برآورد حداکثر درستنمایی پارامترهای مقیاسی و مکانی برای توزیع اسلش استاندارد به دست آورد.

توزیع اسلش که به عنوان آمیخته مقیاسی از توزیع نرمال می باشد به صورت زیر تعریف می شود.

**تعریف ۱.** متغیر تصادفی  $X$  دارای توزیع اسلش است و

می نویسیم

$$X \sim SL(\mu, \sigma^2, q)$$

$$X = \mu + \sigma \lambda^{-1} Z,$$

باشد، که در آن  $Z \sim N(0, 1)$  و  $\lambda \sim Beta(q, 1)$  مستقل از یکدیگرند.

تابع چگالی احتمال توزیع اسلش که به عنوان آمیخته-

که در آن  $\mu$  پارامتر مکان و  $\lambda$  متغیر تصادفی مثبت آمیخته با تابع چگالی احتمال  $h(\lambda | v)$  است که مستقل از  $Z \sim N(0, \sigma^2)$  فرض می شود.  $v$  یک اسکالر یا بردار پارامترهای توزیع  $\lambda$  و  $k(\cdot)$  یک تابع وزنی مثبت است. اگر  $k(\lambda) = \frac{1}{\lambda}$ ، آنگاه داریم:

$$X | \lambda \sim N(\mu, \lambda^{-1} \sigma^2),$$

و در این صورت تابع چگالی احتمال  $X$  به صورت زیر است:

$$f(x | \mu, \sigma^2, v) = \int_0^{\infty} N(x | \mu, \lambda^{-1} \sigma^2) h(\lambda | v) d\lambda. \quad (1)$$

با انتخاب مناسب چگالی آمیخته  $h(\cdot | v)$ ، کلاس توانمندی از توزیع های پیوسته تک مدی و متقارن با استفاده از چگالی بیان شده در (۱) به دست می آید که توزیع هایی دم-کلفت تر از نرمال تولید می کنند. اگر با احتمال یک،  $\lambda = 1$  (متغیر تصادفی تباهیده) باشد، آنگاه حاصل یک توزیع نرمال است.

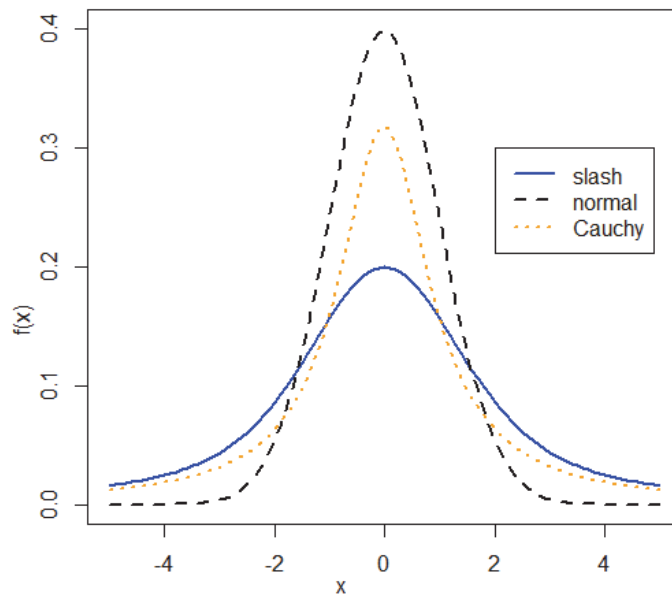
امید ریاضی و واریانس این توزیع ها را با استفاده از قوانین امید ریاضی مکرر و تفکیک واریانس و با توجه به اینکه  $X | \lambda \sim N(\mu, k(\lambda) \sigma^2)$  در صورت وجود به صورت زیر محاسبه می نماییم:

$$\begin{aligned} E(X) &= E(E(X | \lambda)) \\ &= E(\mu) = \mu. \end{aligned} \quad (2)$$

اگر  $E(k(\lambda)) < \infty$  باشد آنگاه واریانس  $X$  به صورت زیر وجود دارد:

$$\begin{aligned} Var(X) &= Var(E(X | \lambda)) + E(Var(X | \lambda)) \\ &= Var(\mu) + E(k(\lambda) \sigma^2) = \sigma^2 E(k(\lambda)). \end{aligned}$$

<sup>۱۲</sup>Andrew and Mallows



شکل ۱: نمودار تابع چگالی توزیع های اسلش و نرمال و کوشی

مقیاسی از توزیع نرمال است، به صورت زیر است: عمل می کنیم:

$$f(x | \mu, \sigma^2, q) = q \int_0^1 \frac{1}{\sqrt{2\pi\frac{\sigma^2}{\lambda}}} \exp\left\{-\frac{\lambda(x-\mu)^2}{2\sigma^2}\right\} \lambda^{q-1} d\lambda$$

$$f(x | \mu, \sigma^2, q) = q \int_0^1 \lambda^{q-1} N(x | \mu, \frac{\sigma^2}{\lambda}) d\lambda,$$

که در آن چگالی  $\lambda$  به فرم زیر می باشد:

$$= \begin{cases} \frac{q}{\sqrt{2\pi\sigma^2}} \frac{\Gamma(q+\frac{1}{2})}{\Gamma(q+\frac{1}{2})} \Gamma(q+\frac{1}{2}, \frac{x-\mu}{\sigma}), & x \neq \mu, \\ \frac{2q}{(2q+1)\sqrt{2\pi\sigma^2}}, & x = \mu, \end{cases} \quad (3)$$

$$h(\lambda | q) = q\lambda^{q-1} I_{(0,1)}.$$

بنابراین، توزیع اسلش معادل با شکل سلسله مراتبی زیر

که در آن

است:

$$z = \frac{(x-\mu)}{\sigma},$$

$$\lambda | q \sim \text{Beta}(q, 1),$$

و

و

$$\Gamma(a, z) = \int_0^z u^{a-1} e^{-u} du,$$

$$X | \mu, \sigma^2, q, \lambda \sim N(\mu, \frac{\sigma^2}{\lambda}),$$

که در آن  $\text{Beta}(\cdot, \cdot)$  نماد توزیع بتا می باشد. برای به دست

آوردن تابع چگالی احتمال توزیع اسلش به صورت زیر

$$f(x | \mu, \sigma^2, q) = \begin{cases} \frac{q}{\sqrt{2\pi\sigma^2}} \frac{\Gamma(q+\frac{1}{2})}{\Gamma(q+\frac{1}{2})} \Gamma(q+\frac{1}{2}, \frac{z}{\sigma}), & z \neq 0, \\ \frac{2q}{(2q+1)\sqrt{2\pi\sigma^2}}, & z = 0. \end{cases}$$

میانگین، نما و واریانس توزیع اسلش که به عنوان

تابع مولد گشتاور توزیع اسلش وجود ندارد اما گشتاورهای آن به صورت زیر محاسبه می‌شود:

$$\begin{aligned} E(X^n) &= E(Z^n) \cdot E(U^{\frac{-n}{q}}) \\ &= E(Z^n) \cdot \frac{q}{q-n}, \quad q > n. \end{aligned} \quad (۶)$$

در نتیجه گشتاورهای اولیه به صورت زیر به دست می‌آیند:

$$\begin{aligned} \mu_1 = E(X) &= 0, & q > 1, \\ \mu_2 = E(X^2) &= \frac{q}{q-2}, & q > 2, \\ \mu_3 = E(X^3) &= 0, & q > 3, \\ \mu_4 = E(X^4) &= \frac{3q}{q-4}, & q > 4. \end{aligned} \quad (۷)$$

بنابراین واریانس آن به صورت زیر است:

$$Var(X) = \frac{q}{q-2}, \quad q > 2. \quad (۸)$$

ضریب چولگی آن عبارتست از:

$$\delta_1(X) = 0, \quad (۹)$$

در توزیع‌های متقارن ضریب چولگی برابر صفر است. اگر گشتاورهای  $X$  تا مرتبه ۴ موجود باشد آن‌گاه ضریب کشیدگی متغیر تصادفی اسلش  $X$  که با  $\delta_2(X)$  نمایش داده می‌شود عبارت خواهد بود از:

$$\begin{aligned} \delta_2(X) &= \frac{E[(X - E(X))^4]}{[Var(X)]^2} - 3 \\ &= 3 \left( \frac{1}{q(q-4)(q-2)^2} - 1 \right), \quad q > 4, \end{aligned}$$

که بزرگ‌تر از صفر است و لذا چگالی اسلش کشیده است.

آمیخته-مقیاسی از توزیع نرمال است، عبارتست از:

$$E(X) = Mode(X) = \mu,$$

$$Var(X) = \frac{q}{q-1} \sigma^2, \quad q > 1.$$

تذکر ۱. به ازای  $\mu = 0$  و  $\sigma = 1$  در رابطه (۳) توزیع اسلش استاندارد به دست می‌آید که به صورت زیر تعریف می‌شود.

تعریف ۲. گویم متغیر تصادفی  $X = \frac{Z}{U^{\frac{1}{q}}}$  دارای توزیع اسلش با پارامتر شکل  $q$  است هرگاه  $Z \sim N(0, 1)$  و  $U \sim U(0, 1)$  مستقل از یکدیگر باشند و آن را با نماد  $X \sim SL(q)$  نمایش می‌دهیم. تابع چگالی آن به صورت زیر است:

$$f(x; q) = q \int_0^1 t^q \phi(xt) dt; \quad -\infty < x < \infty, \quad q > 0 \quad (۴)$$

که در آن  $\phi(\cdot)$  تابع چگالی نرمال استاندارد است.

تذکر ۲. به ازای  $\mu = 0$  و  $\sigma = 1$ ،  $q = 1$  در رابطه (۳) توزیع اسلش کانونی به دست می‌آید (جانسون ۱۳ و همکاران [۱۰]) که تابع چگالی آن به صورت زیر است:

$$f(x, 1) = \begin{cases} \frac{1}{\sqrt{4}\pi} (1 - e^{-\frac{x^2}{4}}), & x \neq 0, \\ \frac{1}{2\sqrt{2}\pi}, & x = 0. \end{cases}$$

تابع توزیع تجمعی اسلش به صورت زیر است:

$$F(x; q) = q \int_0^1 t^{q-1} \Phi(xt) dt, \quad (۵)$$

که در آن  $\Phi(\cdot)$  تابع توزیع نرمال استاندارد است.

توزیع اسلش دم کلفت تر از نرمال است و در صورتی که  $q \rightarrow \infty$  به نرمال میل می‌کند. نمودار تابع چگالی توزیع اسلش به ازای مقادیر مختلف  $q$  در شکل ۲ نشان داده شده است.

$$\hat{\beta} = (\mathbf{X}\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}, \quad (11)$$

که حاصل آن بردار  $p$  بعدی است. حال توزیع توأم و بردار باقی‌مانده  $(\hat{\epsilon})$ ،  $\hat{\epsilon} = \mathbf{Y} - \mathbf{X}\hat{\beta}$  را در نظر می‌گیریم. با توجه به روابط (۱۰) و (۱۱) داریم:

$$\begin{bmatrix} \hat{\beta} \\ \hat{\epsilon} \end{bmatrix} = \begin{bmatrix} (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \\ \mathbf{I} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \end{bmatrix} \mathbf{Y} \begin{bmatrix} \hat{\beta} \\ \mathbf{0} \end{bmatrix} + \begin{bmatrix} (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \\ \mathbf{I} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \end{bmatrix} \hat{\epsilon},$$

## ۴ مدل رگرسیون چندگانه با خطای نرمال

در این بخش مدل رگرسیون با خطای توزیع نرمال را مرور می‌کنیم. مدل رگرسیون

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon, \quad (10)$$

و با فرض‌های

$$E(\epsilon) = \mathbf{0} \text{ یا } E(\mathbf{Y}) = \mathbf{X}\beta \text{ (الف)}$$

$$\text{cov}(\epsilon) = \sigma^2\mathbf{I} \text{ یا } \text{cov}(\mathbf{Y}) = \sigma^2\mathbf{I} \text{ (ب)}$$

را در نظر بگیرید. در این مدل  $\mathbf{Y}$  بردار تصادفی  $n \times 1$  از مشاهدات،  $\mathbf{X}$  یک ماتریس غیرتصادفی  $n \times p$  با رتبه  $p$ ، که در آن  $n$  تعداد مشاهدات و  $p = k + 1$  تعداد ضرایب رگرسیونی است. بردار پارامترهای مجهول رگرسیون  $1 \times p$  است.  $\epsilon$  خطای مدل بنا بر فرض معمول توزیع آن نرمال است. توزیع این برآوردگر نرمال با میانگین  $\beta$  و واریانس  $(\mathbf{X}\mathbf{X})^{-1}\sigma^2$  است.  $\beta = (\beta_1, \beta_2, \dots, \beta_k)'$  ضرایب رگرسیونی و قابل برآورد به روش حداکثر درستنمایی و حداقل مربعات در صورت مشخص نبودن توزیع خطا هستند. تابع درستنمایی به صورت

$$L(\beta) = \frac{1}{(\sqrt{\pi}\sigma)^n} \exp\left\{-\frac{1}{2\sigma^2}(\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta)\right\},$$

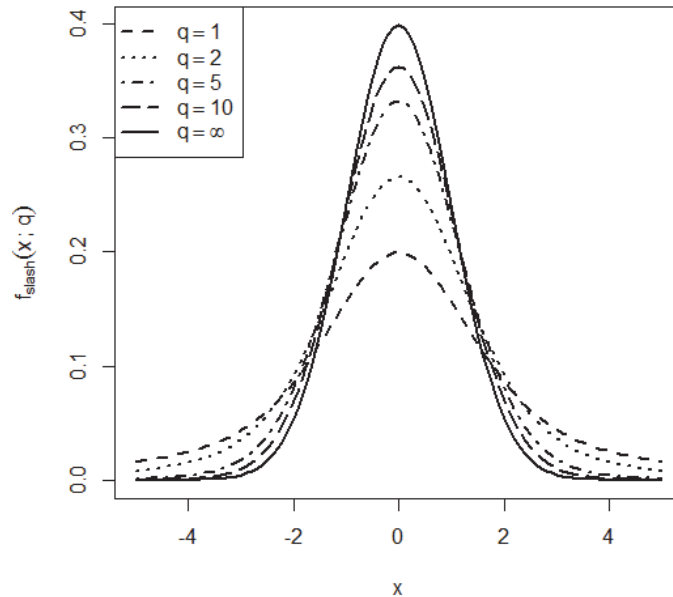
که در آن  $\epsilon \sim N_p(\mathbf{0}, \sigma^2\mathbf{I})$ . حال از آنجا که یک تابع خطی از  $\epsilon$  است و داریم:

$$\begin{bmatrix} \hat{\beta} \\ \hat{\epsilon} \end{bmatrix} - \begin{bmatrix} \hat{\beta} \\ \mathbf{0} \end{bmatrix} \hat{\epsilon} \sim N_{n+p}(\mathbf{0}, \Sigma),$$

که در آن

$$\Sigma = \begin{bmatrix} (\mathbf{X}'\mathbf{X})^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \end{bmatrix}.$$

لذا به کمک روش کمترین مربعات نمی‌توانیم تابعی از  $y$  ها و  $x$  های نمونه برای برآورد مینیمم  $\sigma^2$  بسازیم. با وجود این، می‌توانیم یک برآوردگر نارایب برای  $\sigma^2$  بر مبنای برآوردگر کمترین مربعات  $\hat{\beta}$  بسازیم. با فرض برقراری شرط (ب)،  $\sigma^2$  برای هر  $y_i$ ،  $i = 1, 2, \dots, n$  ثابت است.



شکل ۲: نمودار تابع چگالی توزیع اسلش به ازای مقادیر مختلف پارامتر  $q$

$\sigma^2$  به صورت  $\sigma^2 = E[y_i - E(y_i)]^2$  تعریف می‌شود و با توجه به فرض (الف) داریم:

که در آن  $n$  حجم نمونه و  $k$  تعداد  $x$  هاست. رابطه (۱۱) را می‌توان به صورت زیر نوشت:

$E(y_i) = \beta_0 x_{i1} + \beta_1 x_{i2} + \dots + \beta_k x_{ik}$  که در آن  $x'_i$  سطر  $i$  ام ماتریس  $X$  است. پس  $\sigma^2$  به صورت زیر نوشته می‌شود:

$$s^2 = \frac{1}{n-k-1} (y - \mathbf{X}\hat{\beta})'(y - \mathbf{X}\hat{\beta})$$

$$= \frac{\mathbf{y}'\mathbf{y} - \hat{\beta}'\mathbf{X}'\mathbf{y}}{n-k-1} = \frac{\hat{\epsilon}'\hat{\epsilon}}{n-k-1} = \frac{SSE}{n-k-1}$$

$$\sigma^2 = E[y_i - \mathbf{x}'_i\beta]^2$$

بنابراین می‌توانیم مقدار  $\sigma^2$  را به صورت یک مقدار متوسط از نمونه به فرم زیر برآورد می‌کنیم.

برآوردگر  $s^2$  برای  $\sigma^2$  نارایب است. (رنچر و اسکالچ<sup>۱۴</sup> [۱۵])  
برای اولین بار فیشر<sup>۱۵</sup> [۴] به نتایج نادرست استفاده از توزیع نرمال به عنوان توزیع خطاهای تصادفی در مدل‌های خطی اشاره کرد.

$$s^2 = \frac{1}{n-k-1} \sum_{i=1}^n (y_i - \mathbf{x}'_i\hat{\beta})^2, \quad (12)$$

<sup>۱۴</sup>Rencher and Schaalje

<sup>۱۵</sup>Fisher

## ۵ مدل رگرسیون چندگانه با خطای اسلش

فرض کنید  $\epsilon_1, \dots, \epsilon_n$  متغیرهای تصادفی مستقل باشند که در آن

$\epsilon_i \sim SL(0, \sigma^2)$  خطای مدل رگرسیون خطی با خطای توزیع اسلش با پارامتر  $q$  است. برای برآورد بردار پارامتر  $\theta = (\beta', q)'$  از روش حداکثر درستنمایی استفاده می‌کنیم. تابع لگاریتم درستنمایی توزیع اسلش به صورت زیر است:

$$\ell(\theta) = n \log q - n \log(2\pi\sigma^2) + (q + \frac{1}{q}) \log 2 + n \log 2q - n \log(2q + 1) + \sum_i a(z_i),$$

که در آن  $a(z_i) = \log \Gamma(q + \frac{1}{q}, \frac{z_i^2}{q}) - (2q + 1) \log z_i$  است. معادلات درستنمایی توزیع اسلش فرم بسته ساده‌ای ندارد، لذا این معادلات جواب صریحی برای پارامترهای آن ارائه نمی‌دهند. بنابراین باید برآوردهای حداکثر درستنمایی آن را به روشهای عددی شبیه روش نیوتن-رافسون محاسبه کرد.

## ۶ مثال کاربردی: تحلیل داده‌های جرم و جنایت

هفت متغیر زیر مربوط به داده‌های جرم و جنایت ۴۴ محله دنور ۱۶ می‌باشد.

$X_1$ : کل جمعیت (بر حسب هزار)

$X_2$ : درصد تغییر در جمعیت طی چند سال گذشته

$X_3$ : درصد کودکان (زیر ۱۸ سال) در جامعه

$X_4$ : شرکت در صرف ناهار مدارس آزاد

$X_5$ : تغییر در درآمد خانوار طی چند سال گذشته

$X_6$ : نرخ جرم و جنایت در هر ۱۰۰۰ نفر جمعیت

$Y$ : تغییر در نرخ جرم و جنایت در طی چند سال گذشته  
این داده‌ها تنها با استفاده از مدل رگرسیونی نرمال در متون مختلف بررسی شده است.

در این مقاله با در نظر گرفتن مؤلفه‌های خطا با توزیع اسلش و برآورد حداکثر درستنمایی پارامترهای مدل و با استفاده از معیارهای انتخاب، مدل مناسب برای برازش به این مدل رگرسیونی را انتخاب می‌کنیم. بدین منظور ابتدا با فرض نرمال بودن متغیر پاسخ ( $Y$ ) به شرط متغیرهای توضیحی ( $X_i$ )، مدل رگرسیونی

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \beta_4 X_{i4} + \beta_5 X_{i5} + \beta_6 X_{i6} + \epsilon_i, \quad (13)$$

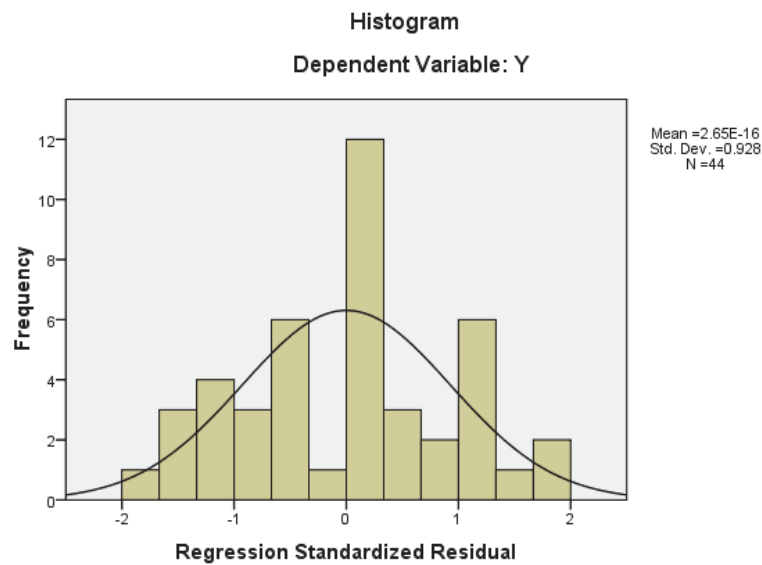
را به روش کلاسیک به داده‌ها برازش می‌دهیم که نتایج در جدول ۱ مشاهده می‌شود. محاسبات ارائه شده در جدول ۱ این نتیجه را می‌دهد که با افزایش افراد جامعه میزان جنایت‌های رخ داده شده کمتر می‌شود در صورتی که این امر غیر معقول به نظر می‌رسد. بنابراین علت چنین نتیجه‌گیری برازش مدل اشتباه به داده‌های این مثال است.

برای بررسی نرمال بودن و همگنی واریانس مانده‌های مدل، نمودار بافت‌نگار باقی مانده‌ها و نمودار مانده‌ها در برابر مقادیر برازش شده در شکل‌های ۳ و ۴ ارائه شده است.



جدول ۱: تحلیل  $OLS$  پارامترهای مدل نرمال برای مدل رگرسیونی (۱۳)

پارامترها	برآورد پارامترها	انحراف معیار	$p - value$
$\beta_0$	-۱۶/۸۵۶	۱۲/۶۰	۰/۱۸۹
$\beta_1$	-۰/۲۰۳	۰/۲۱۲	۰/۳۴۳
$\beta_2$	۱/۴۳۳	۰/۳۶۷	۰/۰۰۰
$\beta_3$	-۰/۴۶۸	۰/۱۳۳	۰/۰۰۱
$\beta_4$	-۰/۳۱۸	۰/۳۹۰	۰/۴۲۰
$\beta_5$	۰/۰۲۹	۰/۰۳۵	۰/۴۰۴
$\beta_6$	۱/۰۵۶	۰/۶۷۱	۰/۱۲۴
$\sigma^2$	۲۰۳/۶۳۲	۱/۷۶۶	۰/۰۰۰۵

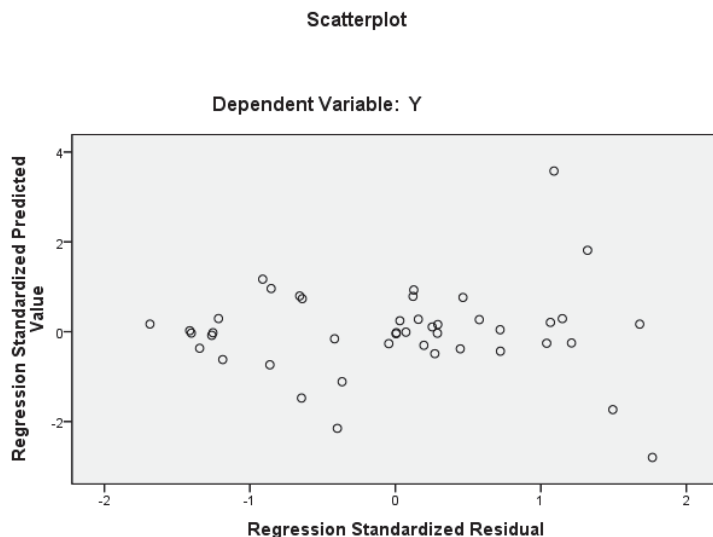


شکل ۳: نمودار بافت‌نگار باقی مانده‌های مدل (۱۳)

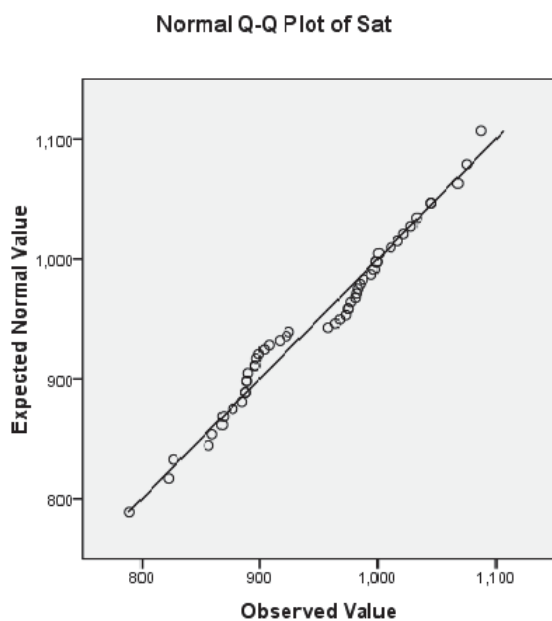
باتوجه به شکل‌های فوق می‌توان نتیجه گرفت که مدل نرمال برای خطاها در این مثال مناسب نیست. از ۰/۰۵ باشد نتیجه می‌گیریم که مانده‌ها از توزیع نرمال تبعیت می‌کنند.

در ادامه آزمون نرمال بودن را برای مانده‌های مدل (۱۳) انجام دادیم. نتایج حاصل در جدول ۲ نشان داده شده اند که نیز گواه بر صحت ادعای فوق می‌باشد. زیرا با توجه به میزان  $p - value$  در جدول ۲، فرض نرمال بودن مانده‌های مدل در سطح ۰/۰۵ رد می‌شود. همچنین شکل ۵ نمودار چندک-چندک نرمال برای باقیمانده‌ها را نشان می‌دهد.

نکته: آزمون شاپیرو-ویلک برای داده‌هایی جهت تست نرمالیتی با حجم کم، به عنوان مثال ۵۰ و یا کمتر، به کار می‌رود و آزمون کلموگروف-اسمیرنوف برای داده‌هایی با حجم نمونه بالا استفاده می‌شود. برای بررسی نرمال بودن مانده‌ها میزان  $p - value$  متناظر با آن‌ها را بررسی می‌کنیم. چنانچه میزان  $p - value$  این آزمون‌ها بیشتر



شکل ۴: نمودار مانده ها در برابر مقادیر برازش شده



شکل ۵: نمودار چندک-چندک نرمال برای مانده‌های مدل (۱۳)

جدول ۲: آزمون شاپیرو - ویلک برای نرمال بودن مانده‌ها

آماره	درجه آزادی	$p - value$
۰/۹۴	۴۴	۰/۰۲۳

حال فرض می‌کنیم مؤلفه‌های خطا دارای توزیع اسلش هستند و برآورد حداکثر درست‌نمایی پارامترهای مدل را به دست می‌آوریم و مشاهده می‌کنیم که مدل اسلش برازش بهتری برای این مدل رگرسیونی است. در مسئله برآوردیابی حداکثر درست‌نمایی پارامترهای آن از نرم افزار R استفاده شده است که نتایج در جدول ۳ ارائه شده است.

نمودار بافت نگار باقی مانده‌های مدل با فرض برقراری توزیع اسلش در شکل ۶ رسم شده است.

جدول ۳: برآورد پارامترهای مدل (۱۳) با فرض برقراری توزیع اسلش برای خطاهای مدل

پارامترها	برآورد پارامترها	انحراف معیار	loglike
$\beta_0$	-۶/۱۳۳	۱۰/۱۶۶	-۱۷۵/۵۷۷
$\beta_1$	۰/۳۲۹	۰/۲۰۵	
$\beta_2$	۰/۹۸۵	۰/۳۰۶	
$\beta_3$	-۰/۲۷۹	۰/۱۱۰	
$\beta_4$	۰/۶۵۷	۰/۳۳۱	
$\beta_5$	۰/۰۱۱	۰/۰۲۸	
$\beta_6$	-۰/۴۲۳	۰/۵۲۲	
$\beta_7$	۱۰۸/۷۸۴	۱/۸۴۶	
q	۳/۰۱۰	۰/۹۳۳	

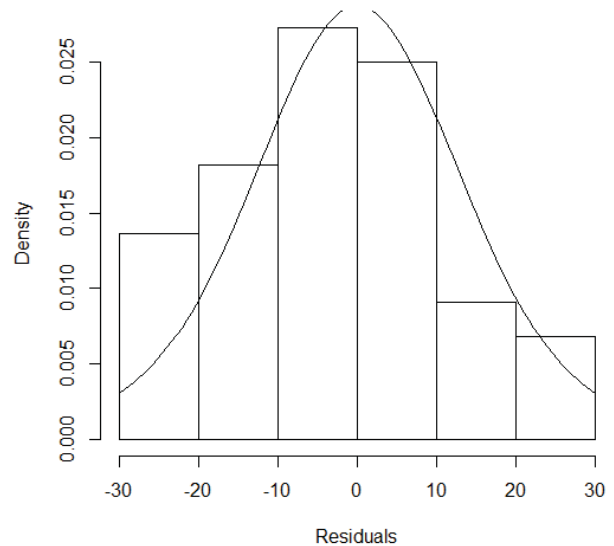
تعداد زیادی از معیارهای انتخاب مدل بر اساس نظریه ی اطلاع پدید آمده‌اند. یکی از معروفترین معیارها لگاریتم درستنمایی است، هنگامی که در جستجوی بهترین برازش برای داده‌ها در بین سایر توزیع‌ها هستیم معمولاً توزیع با بیشترین مقدار لگاریتم درستنمایی به عنوان برازش بهتر انتخاب می‌شود. زمانی که تعداد پارامترهای برآورد شده یکسان نباشد، این معیار چندان مناسب نیست. به طور معمول برای مقایسه مدل‌هایی با تعداد پارامتر متفاوت، از معیار اطلاع آکائیک ( $AIC$ ) و معیار اطلاع بیزی ( $BIC$ ) مناسب هستند که به ترتیب توسط آکائیک<sup>۱۷</sup> [۱] و شوارتز<sup>۱۸</sup> [۱۷] برای انتخاب بهترین مدل آماری معرفی شدند. مقدار معیار اطلاع آکائیک ( $AIC$ ) با استفاده از فرمول  $AIC = -2\ell(\hat{\theta}) + 2p^*$  که در آن  $p^*$  تعداد پارامترهای مدل است، معین می‌شود. معیار اطلاع بیزی ( $BIC$ ) به صورت زیر تعریف می‌شود:

$BIC = -2\ell(\hat{\theta}) + p^* \ln n$  که در آن  $n$  حجم نمونه است. هنگامی که در جستجوی برازش بهتر برای داده‌ها در بین سایر توزیع‌ها هستیم توزیع با کوچکترین مقدار  $AIC$  و  $BIC$  به عنوان برازش بهتر انتخاب می‌شود. معیارهای

<sup>۱۷</sup>Akaike Hirotugu

<sup>۱۸</sup>Schwarz

Histogram of Residuals



شکل ۶: نمودار بافت نگار باقی مانده‌های مدل (۱۳) با فرض برقراری توزیع اسلش برای خطاها

برآورد واریانس ( $\sigma^2$ ) در مدل نرمال از مدل اسلش بیشتر است در واقع ما انتظار چنین مسئله ای را داشتیم چون پراکندگی زیاد تا حدودی جای خود را به دم سنگینی می‌دهد و به خاطر اینکه توزیع اسلش دم سنگین تر از توزیع نرمال است، این امر کاملاً معقول به نظر می‌رسد. اکنون با معیارهای کمی به مقایسه دو مدل می‌پردازیم.

## مراجع

[1] Akaike, H. (1974). A new look at the statistical model identification. IEEE Transactions on Automatic Control, 19, 716-723.

[2] Andrews, D.F. and Mallows, C.L. (1974). Scale Mixtures of Normal Distributions. Journal of the Royal Statistical Society, B36, 99-102.

[3] Arslan, O. and Genç, A.I. (2009). A generalization of the multivariate slash distribution. Journal of Statistical Planning and Inference, 139, 1164-1170.

[4] Fisher, R.A. (1956). Statistical Methods in Scientific Inference. Oliver and Boyd, London.

[5] Genç, A.I. (2007). A generalization of the univariate slash by a scale-mixed exponential power distribution. Journal of Communications in Statistics-Simulation and Computation, 36, 937-947.

[6] Genç, A.I. (2012). A skew extension of the slash distribution via beta-normal distribution. Statist. Papers, 54(2), 427-442.

لگاریتم درست‌نمایی، اطلاع آکائیک و اطلاع بیزی برای دو حالت مدل با خطای توزیع نرمال و با خطای توزیع اسلش در جدول ۴ ارائه شده است.

جدول ۴: مقادیر معیارهای لگاریتم درست‌نمایی، اطلاع آکائیک و اطلاع بیزی برای دو مدل با خطاهای توزیع اسلش و توزیع نرمال

معیار سنجش مدل	مدل نرمال	مدل اسلش
<i>AIC</i>	۳۷۸/۶۱۹	۳۶۷/۱۵۴
<i>BIC</i>	۳۹۰/۴۸۶	۳۸۱/۴۲۷۵
لگاریتم درست‌نمایی	-۱۸۰/۳۰۹	-۱۷۵/۵۷۷

از بین دو مدل اسلش و نرمال، مدل با خطای توزیع اسلش دارای کمترین مقدار معیار اطلاع آکائیک (*AIC*) و معیار اطلاع بیزی (*BIC*) و بیشترین مقدار لگاریتم درست‌نمایی است، پس نسبت به توزیع نرمال برازش بهتری برای خطاهای مدل رگرسیونی در رابطه (۱۳) می‌باشد.

## ۷ نتیجه‌گیری

در این مقاله توزیع اسلش و برخی ویژگی‌های آن را مرور نمودیم و از آن در برازش مدل‌های رگرسیونی چندگانه معمولی استفاده کردیم و با استفاده از داده‌های واقعی نشان دادیم که توزیع اسلش جایگزین مناسبی برای توزیع نرمال در خطای مدل‌های رگرسیونی است. برای مقایسه دو مدل از معیارهای مختلف استفاده شد و مشاهده کردیم که توزیع اسلش نسبت به توزیع نرمال برای مؤلفه‌های خطای مدل رگرسیونی (۱۳) مناسب‌تر است.

- Computation and Simulation, 73,791-805.
- [13] Lange, K. and Sinsheimer, J.S. (1993). Normal Independent distributions and their applications in robust regression. *Journal of Computational and Graphical Statistics*, 2, 175-198.
- [14] Morgenthaler, S. (1986). Robust confidence intervals for a location parameter: the configural approach. *Journal of the American Statistical Association*, 81, 518-525.
- [15] Rencher, A.C. and Schaalje, G. B. (2008). *Linear models in statistics*. 2nd ed., Wiley, New York.
- [16] Rogers, W. H. and Tukey, J. W. (1972). Understanding some long-tailed symmetrical distributions. *Statistica Neerlandica*, 26, 211-226.
- [17] Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, 6(2), 461-464.
- [7] Gomez, H.W., Quintana, A. and Torres, J. (2007). A new family of slash-distributions with elliptical contours. *Statistics and Probability Letters*, 77, 717-725.
- [8] Gómez, W., Reyesa, H. and Bolfarineb, H. (2013). Modified slash distribution. *Statistics: A Journal of Theoretical and Applied Statistics*, 47(5), 929-941.
- [9] Jamshidian, M. (2001). A note on parameter and standard error estimation in adaptive robust regression. *Journal of Statistical Computation and Simulation*, 71, 11-27.
- [10] Johnson, N.L., Kotz, S. and Balakrishnan, N. (1995). *Continuous univariate distribution*, Vol. 1, Second ed., Wiley, New York.
- [11] Kafadar, K. (1982). A biweight approach to the one-sample problem. *Journal of the American Statistical Association*, 77, 416-424.
- [12] Kashid, D.N and Kulkarni, S.R. (2003). Subset selection in multiple linear regression with heavy tailed error distribution. *Journal of Statistical*