

تبدیل داده‌ها برای مقابله با نقاط پرت

محمد مهدی شریف بیگی

۱ مهر ۱۴۰۴

۱ مقدمه

تبدیل داده‌ها (Data Transformation) یکی از روش‌های مؤثر برای کاهش تأثیر نقاط پرت (Outliers) است. این روش با اعمال تبدیل ریاضی مناسب، مقیاس داده‌ها را تغییر داده و توزیع آن‌ها را به شکل طبیعی‌تر درآورده و تأثیر مقادیر افراطی را کاهش می‌دهد.

۲. تبدیل لگاریتمی (Logarithmic Transformation)

۱.۲ فرمول:

$$y = \log(x) \quad \text{یا} \quad y = \log_{10}(x) \quad \text{یا} \quad y = \ln(x)$$

۲.۲ توضیح مفصل:

تبدیل لگاریتمی یکی از پرکاربردترین روش‌های تبدیل داده است که به خصوص برای داده‌های دارای چولگی راست (right-skewed) بسیار مؤثر است.

۱.۲.۲ خصوصیات ریاضی:

برای $x > 0$ ، تابع لگاریتم دارای خصوصیات زیر است:

$$\frac{d}{dx} \ln(x) = \frac{1}{x} \quad (۱)$$

$$\lim_{x \rightarrow \infty} \frac{\ln(x)}{x} = 0 \quad (۲)$$

$$\ln(ab) = \ln(a) + \ln(b) \quad (۳)$$

این خصوصیات نشان می‌دهد که:

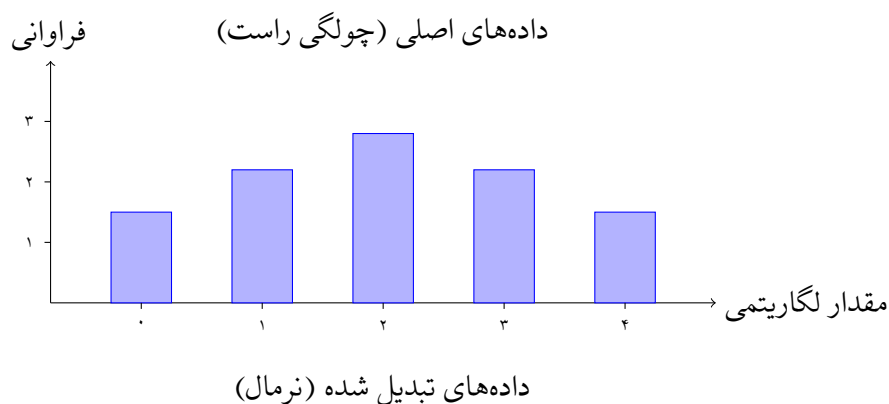
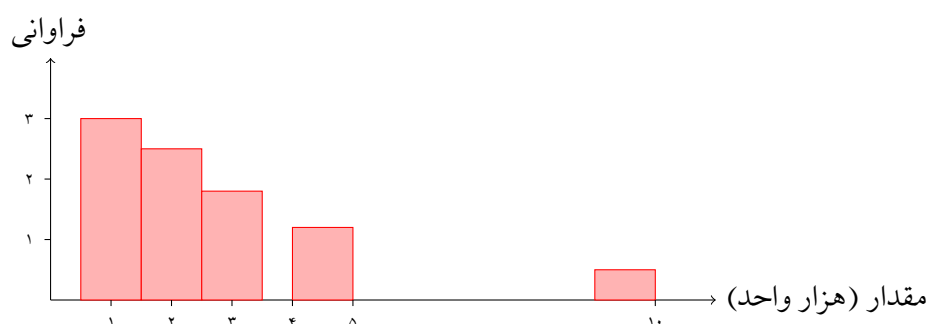
- شیب تابع با افزایش x کاهش می‌یابد
- رشد لگاریتم نسبت به x کندتر می‌شود
- ضرب در جمع تبدیل می‌شود

۲.۲.۲ مثال عملی با محاسبات کامل:

فرض کنید داده‌های اصلی درآمد (میلیون ریال): $\{1, 10, 100, 1000, 10000\}$ مقایسه آمارهای توصیفی

آماره	داده اصلی	داده تبدیل شده
میانگین	$\frac{1+10+100+1000+10000}{5} = 2222/2$	$\frac{0+1+2+3+4}{5} = 2$
میانه	۱۰۰	۲
انحراف معیار	$4499/4$	$1/58$
محدوده	$10000 - 1 = 9999$	$4 - 0 = 4$
ضریب تغییرات	$\frac{4499/4}{2222/2} = 2/025$	$\frac{1/58}{2} = 0/79$

۳.۲ نمودار مقایسه‌ای بهبود یافته:



۳. تبدیل رادیکال (Square Root Transformation)

۱.۳ فرمول:

$$y = \sqrt{x}$$

۲.۳ توضیح مفصل:

تبدیل رادیکال نسبت به تبدیل لگاریتمی ملایم‌تر است و برای داده‌هایی که دارای چولگی متوسط هستند مناسب است.

۱.۲.۳ خصوصیات ریاضی:

برای $x \geq 0$:

$$\frac{d}{dx} \sqrt{x} = \frac{1}{2\sqrt{x}} \quad (4)$$

$$\lim_{x \rightarrow 0^+} \frac{1}{2\sqrt{x}} = +\infty \quad (5)$$

$$\lim_{x \rightarrow \infty} \frac{1}{2\sqrt{x}} = 0 \quad (6)$$

۲.۲.۳ مثال عملی با محاسبات کامل:

فرض کنید داده‌های اصلی مساحت (متر مربع): $\{4, 16, 64, 256, 1024\}$ محاسبه دقیق آمارهای توصیفی برای داده‌های اصلی: $X = \{4, 16, 64, 256, 1024\}$

$$\bar{X} = \frac{4 + 16 + 64 + 256 + 1024}{5} = \frac{1364}{5} = 272.8$$

$$s_X^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1} = \frac{(4 - 272.8)^2 + \dots + (1024 - 272.8)^2}{4}$$

$$s_X^2 = \frac{72267.04 + 66022.24 + 43436.84 + 379.24 + 56484.04}{4} = \frac{74694.4}{4} = 18673.6$$

$$s_X = \sqrt{18673.6} = 136.7$$

برای داده‌های تبدیل شده: $Y = \{2, 4, 8, 16, 32\}$

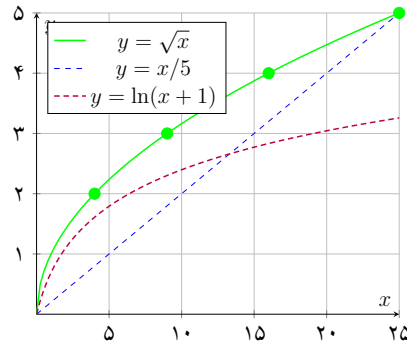
$$\bar{Y} = \frac{2 + 4 + 8 + 16 + 32}{5} = \frac{62}{5} = 12.4$$

$$s_Y^2 = \frac{(2 - 12.4)^2 + (4 - 12.4)^2 + (8 - 12.4)^2 + (16 - 12.4)^2 + (32 - 12.4)^2}{4}$$

$$s_Y^2 = \frac{108.16 + 70.56 + 19.36 + 12.96 + 384.16}{4} = \frac{595.2}{4} = 148.8$$

$$s_Y = \sqrt{148.8} = 12.2$$

۳.۳ نمودار تابع رادیکال:



۴. تبدیل باکس-کاکس (Box-Cox Transformation)

۱.۴ فرمول:

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

۲.۴ تفسیر ریاضی مقادیر مختلف λ :

۱.۲.۴ $\lambda = 2$ (تبدیل درجه دوم):

$$y = \frac{x^2 - 1}{2}$$

این تبدیل برای داده‌هایی که چولگی چپ دارند (left-skewed) مناسب است.

۲.۲.۴ $\lambda = 1$ (تبدیل خطی):

$$y = x - 1$$

این تبدیل فقط یک انتقال ساده است و شکل توزیع را تغییر نمی‌دهد.

۳.۲.۴ $\lambda = 0.5$ (تبدیل رادیکال):

$$y = \frac{\sqrt{x} - 1}{0.5} = 2(\sqrt{x} - 1)$$

۴.۲.۴ $\lambda = 0$ (تبدیل لگاریتمی):

$$y = \ln(x)$$

۵.۲.۴ $\lambda = -۰/۵$ (تبدیل معکوس رادیکال):

$$y = \frac{x^{-۰/۵} - ۱}{-۰/۵} = ۲ \left(۱ - \frac{۱}{\sqrt{x}} \right)$$

۶.۲.۴ $\lambda = -۱$ (تبدیل معکوس):

$$y = ۱ - \frac{۱}{x}$$

۳.۴ روش یافتن λ بهینه:

پارامتر λ از طریق بیشینه‌سازی تابع درستنمایی محاسبه می‌شود:

$$L(\lambda) = -\frac{n}{۲} \ln \left(\frac{\text{RSS}(\lambda)}{n} \right) + (\lambda - ۱) \sum_{i=1}^n \ln(x_i)$$

که در آن:

$$\text{RSS}(\lambda) = \sum_{i=1}^n [y_i(\lambda) - \bar{y}(\lambda)]^2 \quad (۷)$$

$$y_i(\lambda) = \begin{cases} \frac{x_i^\lambda - 1}{\lambda} & \lambda \neq ۰ \\ \ln(x_i) & \lambda = ۰ \end{cases} \quad (۸)$$

الگوریتم محاسبه:

۱. مقادیر مختلف λ را در بازه $[-۲, ۲]$ با گام $۰/۱$ آزمایش کنید

۲. برای هر λ ، داده‌ها را تبدیل کنید

۳. $L(\lambda)$ را محاسبه کنید

۴. λ که بیشترین $L(\lambda)$ را دارد، انتخاب کنید

۴.۴ مثال عملی کامل:

داده‌های اصلی: $\{۹, ۲۵, ۱۰۰, ۴۰۰\}$ ، فرض کنید $\lambda = ۰/۵$ بهینه محاسبه شده است.
گام ۱: تبدیل داده‌ها

$$y_1 = \frac{۹^{۰/۵} - ۱}{۰/۵} = \frac{۳ - ۱}{۰/۵} = \frac{۲}{۰/۵} = ۴ \quad (۹)$$

$$y_2 = \frac{۲۵^{۰/۵} - ۱}{۰/۵} = \frac{۵ - ۱}{۰/۵} = \frac{۴}{۰/۵} = ۸ \quad (۱۰)$$

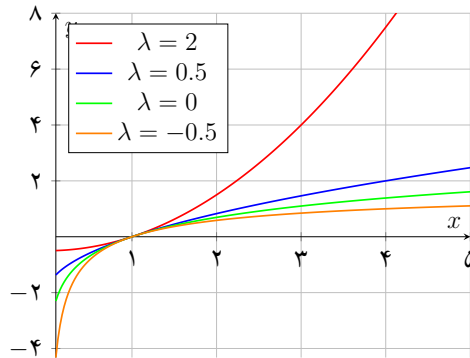
$$y_3 = \frac{۱۰۰^{۰/۵} - ۱}{۰/۵} = \frac{۱۰ - ۱}{۰/۵} = \frac{۹}{۰/۵} = ۱۸ \quad (۱۱)$$

$$y_4 = \frac{۴۰۰^{۰/۵} - ۱}{۰/۵} = \frac{۲۰ - ۱}{۰/۵} = \frac{۱۹}{۰/۵} = ۳۸ \quad (۱۲)$$

گام ۲: محاسبه آمارهای توصیفی

داده‌های اصلی: میانگین $= 133/5 = \frac{9+25+100+400}{4}$ ، انحراف معیار $= 181/9$
 داده‌های تبدیل شده: میانگین $= 17 = \frac{4+8+18+38}{4}$ ، انحراف معیار $= 15/6$
 ضریب تغییرات کاهش یافته: از $1/36 = \frac{181/9}{133/5}$ به $0/92 = \frac{15/6}{17}$

۵.۴ نمودار اثر مقادیر مختلف λ :



۵ مقایسه کمی سه روش

۱.۵ مثال کاربردی: داده‌های درآمد

داده‌های درآمد (میلیون تومان): $\{2, 5, 8, 12, 15, 25, 45, 120\}$

۱.۱.۵ محاسبات دقیق:

۱. آمارهای اصلی:

$$\bar{x} = \frac{2 + 5 + 8 + 12 + 15 + 25 + 45 + 120}{8} = \frac{232}{8} = 29 \quad (13)$$

$$s^2 = \frac{\sum (x_i - 29)^2}{7} = \frac{11914}{7} = 1702 \quad (14)$$

$$s = \sqrt{1702} = 41/26 \quad (15)$$

$$\text{چولگی} = \frac{\sum (x_i - 29)^3 / 8}{s^3} = 1/89 \quad (16)$$

۲. تبدیل لگاریتمی:

$$y_{\log} = \{\ln(2), \ln(5), \ln(8), \ln(12), \ln(15), \ln(25), \ln(45), \ln(120)\}$$

$$= \{0/69, 1/61, 2/08, 2/48, 2/71, 3/22, 3/81, 4/79\}$$

میانگین: $\bar{y}_{\log} = 2/67$ ، انحراف معیار: $s_{\log} = 1/35$ ، چولگی: $0/23$

۳. تبدیل رادیکال:

$$y_{\text{sqrt}} = \{\sqrt{2}, \sqrt{5}, \sqrt{8}, \sqrt{12}, \sqrt{15}, \sqrt{25}, \sqrt{45}, \sqrt{120}\}$$

$$= \{1/41, 2/24, 2/83, 3/46, 3/87, 5/00, 6/71, 10/95\}$$

میانگین: $\bar{y}_{\text{sqr}} = 4/56$ ، انحراف معیار: $s_{\text{sqr}} = 3/24$ ، چولگی: $0/67$

آماره	داده اصلی	لگاریتمی	رادیكال
میانگین	۰.۲۹	۶۷.۲	۵۶.۴
انحراف معیار	۲۶.۴۱	۳۵.۱	۲۴.۳
ضریب تغییرات	۴۲.۱	۵۱.۰	۷۱.۰
چولگی	۸۹.۱	۲۳.۰	۶۷.۰

۶ راهنمای عملی انتخاب روش

۱.۶ الگوریتم تصمیم‌گیری علمی:

۱. محاسبه ضریب چولگی (γ_1):

$$\gamma_1 = \frac{E[(X - \mu)^3]}{\sigma^3} = \frac{\sum_{i=1}^n (x_i - \bar{x})^3 / n}{s^3}$$

۲. قانون تصمیم‌گیری:

- اگر $|\gamma_1| < 0.5$: نیازی به تبدیل نیست
- اگر $0.5 \leq |\gamma_1| < 1$: تبدیل رادیكال
- اگر $|\gamma_1| \geq 1$: تبدیل لگاریتمی یا باکس-کاکس

$$W = \frac{(\sum_{i=1}^n a_i x_{(i)})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

۳. آزمون نرمالیتی شاپیرو-ویلک: $p\text{-value} > 0.05$ اگر داده‌ها نرمال هستند.

۷ نتیجه‌گیری

تبدیل داده‌ها ابزار قدرتمندی برای مقابله با نقاط پرت و بهبود نرمالیتی است. انتخاب روش مناسب باید بر اساس:

- تحلیل کمی چولگی داده‌ها
- آزمون‌های آماری نرمالیتی
- ماهیت علمی متغیرها
- هدف نهایی تحلیل

همیشه قبل و بعد از تبدیل، آزمون‌های آماری مناسب را انجام دهید تا از مؤثر بودن تبدیل اطمینان حاصل کنید.