

آزمون خوبی برازندگی توزیع با استفاده از نوعی نامساوی چبیشف

حسین نگارستانی^۱، زهره شیشه بر^۲، مینا توحیدی^۳

چکیده:

تشخیص توزیع داده‌ها یکی از مسائل مهم و کاربردی در انجام استنباط آماری می‌باشد. آزمون‌های زیادی جهت بررسی توزیع داده‌ها وجود دارد. در این مقاله یک آزمون جدید بر پایه نوعی از نامساوی چبیشف ارائه می‌گردد. نامساوی ذکر شده، بهبود یافته‌ی نامساوی چبیشف می‌باشد. مقایسه توان این آزمون با دیگر آزمون‌ها بیانگر برتری این آزمون می‌باشد.

واژه‌های کلیدی: آزمون خوبی برازندگی توزیع، نامساوی چبیشف.

۱ مقدمه

صورت کلاسیک نامساوی چبیشف به صورت زیر است:

$$P(|X - E(X)| \geq t) \leq \frac{Var(X)}{t^2}$$

آمیخته مقیاسی ارائه کردند. ما در اینجا به بیان این

قضیه می‌پردازیم:

قضیه ۱: فرض کنید F یک تابع توزیع متقارن با تابع

چگالی f باشد، به طوری که روی یک بازه متناهی یا

نامتناهی مثبت و پیوسته و گشتاور مرتبه $a - a$ آن متناهی

باشد. اگر X یک متغیر تصادفی با توزیع آمیخته مقیاسی

از تابع توزیع F باشد ℓ به صورت زیر تعریف شود:

$$\ell = \left(\frac{E|X|^a}{M_a} \right)^{\frac{1}{a}}$$

که در آن:

$$M_a = \int_{-\infty}^{\infty} |z|^a dF(z)$$

آن‌گاه:

$$\max P(|X| \geq t\ell) = \begin{cases} 2\bar{F}(z_a)z_a^a & t \geq z_a \\ 2\bar{F}(t) & \text{در سایر نقاط} \end{cases}$$

این نامساوی برای هر متغیر تصادفی X با واریانس

متناهی برقرار است. در سال ۲۰۰۶ گارود و همکارانش

[۲] بر اساس نامساوی فوق، آزمونی را برای خوبی

برازندگی توزیع ارائه کردند. به منظور یافتن آزمونی

پرتوان‌تر، نیاز به کران دقیق‌تری برای نامساوی چبیشف

داریم. یعنی نیاز به کوچکترین C_α داریم که در رابطه زیر

صدق کند:

$$P(|X| \geq t) \leq C_\alpha \cdot \frac{E|X|^\alpha}{t^\alpha}$$

که در آن α و t اعداد حقیقی مثبت هستند. C_α ثابت

چبیشف نامیده می‌شود. سیسزار و همکارانش [۱] قضیه

ای را برای به دست آوردن ثابت چبیشف در توزیع‌های

^۱گروه آمار - دانشگاه شیراز

^۲گروه آمار - دانشگاه شیراز

^۳گروه آمار - دانشگاه شیراز

کنیم. طول پیشنهادی بازه ها توسط اسکات به صورت زیر ارائه شد [۲]:

$$h = \frac{Y_{max} - Y_{min}}{1 + \log_2^2}$$

پس از این تقسیم بندی، فرض صفر را در هر بازه مورد بررسی قرار می دهیم. منطقی است که اگر حداقل یک بازه اختلافی با توزیع فرض صفر وجود داشت، آن گاه فرض صفر را رد کنیم. یعنی داده ها از توزیع فرض صفر پیروی نمی کنند. در نتیجه رد فرض صفر در حداقل یک بازه معادل با رد فرض صفر کلی می باشد. احتمال این که یک مشاهده در بازه i - ام ($1 \leq i \leq m$) قرار گیرد را با p_i نمایش می دهیم.

فرض کنید n تعداد کل مشاهدات نمونه باشد. متغیر تصادفی X_i را تعداد مشاهدات در بازه i - ام در نظر می گیریم. همچنین آماره آزمون را در بازه i - ام به صورت زیر تعریف می شود [۲]:

$$\tau_i = \frac{X_i - E(X_i)}{\sqrt{Var(X_i)}}$$

فرض کنید که مشاهدات مستقل از یکدیگر باشند. در نتیجه متغیر تصادفی X_i دارای توزیع دو جمله ای با پارامترهای n و p_i است. از این رو $E(X_i) = np_i$ و $Var(X_i) = np_i(1 - p_i)$ می باشد. احتمال این است که یک مشاهده تحت فرض صفر در بازه i - ام قرار گیرد و به فرض صفر آزمون بستگی دارد.

حال با استفاده از قضیه حد مرکزی وقتی که $n \rightarrow \infty$ می توان توزیع متغیر تصادفی X_i که دو جمله ای می باشد را با توزیع نرمال تقریب زد (در عمل این تقریب برای نمونه هایی با حجم نمونه بیش از ۳۰ قابل استفاده

وقتی که z_a را کوچکترین ریشه مثبت معادله زیر در نظر گرفته شود:

$$\frac{zf(z)}{F(z)} = a$$

که در آن $\bar{F}(z) = 1 - F(z)$. در نتیجه ثابت چبیشف برابر است با:

$$C_a = \frac{\bar{F}(z_a)z_a^a}{M_a}$$

پیامد ۱: برای یک متغیر تصادفی X با تابع توزیع F ، با در نظر گرفتن $\sigma = 1$ در قضیه بالا ثابت چبیشف برابر است با:

$$D_a = \frac{\bar{F}(z_a)z_a^a}{E|X|^a}$$

در نتیجه یک کران بالا برای $P(|X| \geq t)$ به صورت زیر خواهد بود:

$$\frac{D_a}{t^a} = \frac{\bar{F}(z_a)z_a^a}{t^a}$$

۲ آزمون خوبی برازندگی چبیشف

هدف از انجام این آزمون این است که آیا یک مجموعه از داده ها را می توان با یک تابع توزیع مشخص (F_0) توصیف نمود. فرض کنید Y_1, Y_2, \dots, Y_n یک نمونه تصادفی n تایی از جامعه ای با توزیع F باشند. می خواهیم آزمون زیر را انجام دهیم:

$$\begin{cases} H_0 : F = F_0 \\ H_1 : F \neq F_0 \end{cases}$$

برای انجام این آزمون ابتدا بایستی داده های بین نقاط حداقل و حداکثر را به m بازه مساوی به طول h تقسیم

برای این که آزمون در سطح α باشد، باید رابطه زیر برقرار باشد:

$$P_{H_0}(\text{reject } H_0) \leq \alpha.$$

از طرفی با استفاده از نامساوی بول داریم:

$$\begin{aligned} P_{H_0}(\text{Reject } H_0) &= P_{H_0}(|\tau| \geq k) \\ &= P_{H_0}\left\{\bigcup_{i=1}^m (|\tau_i| \geq k)\right\} \\ &\leq \sum_{i=1}^m P_{H_0}(|\tau_i| \geq k) \leq \frac{mD_2}{k^2}. \end{aligned}$$

این آزمون در سطح α می باشد اگر مقدار $\frac{mD_2}{k^2}$ برابر α در نظر گرفته شود. در نهایت استراتژی آزمون جدید به صورت زیرارائه می شود:

در سطح معنی داری α فرض H_0 رد می شود اگر و تنها اگر:

$$|\tau| \geq \sqrt{\frac{0.33143m}{\alpha}}$$

این آزمون، آزمون چبیشف نامیده می شود.

۳ شبیه سازی

برای بررسی برتری آزمون جدید نسبت به دیگر آزمون ها توان این آزمون ها را با هم مقایسه می کنیم. برای مقایسه توان آزمون ها نیاز به مشخص کردن توزیع تحت فرض مقابل داریم. در فرض مقابل از توزیع های زیر استفاده می کنیم:

می باشد). در نتیجه با استفاده از نامساوی بهبود یافته چبیشف و قضیه (۱) برای توزیع نرمال در حالت $a = 2$ خواهیم داشت:

$$\begin{aligned} P(|\tau_i| \geq k) &= P\left(\frac{|X_i - E(X_i)|}{\sqrt{\text{Var}(X_i)}} \geq k\right) \\ \dots &= P(|Z| \geq k) \leq \frac{2\bar{\Phi}(z_2)z_2^2}{k^2} \end{aligned}$$

که در آن Z دارای توزیع نرمال استاندارد و $\Phi(\cdot)$ تابع توزیع مربوط به Z است.

در توزیع نرمال، برای حالت $a = 2$ به وسیله محاسبات عددی مقدار z_2 برابر $1/19$ به دست آمده و طرف راست نامساوی فوق برابر است با:

$$\frac{2\bar{\Phi}(z_2)z_2^2}{k^2} = \frac{\phi(z_2)z_2^2}{k^2} = \frac{0.33143}{k^2}$$

وقتی که $\phi(\cdot)$ تابع چگالی مربوط به توزیع نرمال استاندارد باشد. همان طور که اشاره شد در فرض صفر در حداقل یک بازه $(|\tau_i| \geq k)$ معادل با رد فرض صفر کلی ($H_0: F = F_0$) می باشد. بنابراین می توان برای سهولت در انجام آزمون، آماره آزمون را به صورت زیر در نظر گرفت:

$$|\tau| = \max |\tau_i|$$

حال ناحیه پذیرش کلی^۴ را در نظر بگیرید:

$$\bigcap_{i=1}^m (|\tau_i| < k) = (\max |\tau_i| < k)$$

از آنجا که ناحیه رد متمم ناحیه پذیرش می باشد، پس ناحیه رد^۵ به صورت زیر است:

$$\bigcup_{i=1}^m (|\tau_i| \geq k) = (\max |\tau_i| \geq k) = (|\tau| \geq k)$$

Acceptance region^۴
Rejection region^۵

توزیع یکنواخت $(-2, 2)$

توزیع t - استیودنت با سه درجه آزادی

توزیع لاپلاس

توزیع نمایی با میانگین $\frac{1}{3}$

توزیع کوشی استاندارد

توزیع لگ نرمال استاندارد

شبیه سازی به روش مونت کارلو در سطح 0.05 و با

استفاده از نرم افزار *Minitab* انجام شده است.

آزمون توزیع نرمال^۶:

برای مقایسه توان آزمون جدید و آزمون های معروفی

همچون آندرسن - دارلینگ (AD)، کولموگروف -

اسمیرنوف (KS) و خی - د^۷ و برای آزمون کردن توزیع

نرمال، 1000 نمونه 60 تایی از توزیع های فرض مقابل

تولید کرده و توان آزمون های ذکر شده را محاسبه نمودیم.

نتایج در جدول شماره (۱) آمده است.

جدول ۱. توان آزمون های توزیع نرمال برای فرض های مقابل مختلف در سطح معنی داری 0.05 .

<i>Chebyshev</i>	<i>Chi - Square</i>	<i>KS</i>	<i>AD</i>	فرض مقابل
۸۹	۷۲	۳۳	۴۴	$u(-2, 2)$
۶۰	۵۵	۱۱	۳۴	$t(3)$
۴۳	۳۳	۶	۳۶	<i>laplace</i>
۱۰۰	۹۸	۱۰۰	۱۰۰	$exp(3)$
۹۸	۹۷	۵۶	۴۶	<i>cauchy</i>
۱۰۰	۹۹	۱۰۰	۱۰۰	<i>lognormal</i>

بررسی نرمال بودن داده ها در حالتی که حجم نمونه بیشتر

از 50 است، توصیه می شود.

همان طور که در جدول (۱) آمده است توان آزمون جدید

از دیگر آزمون ها بیشتر می باشد. پس این آزمون جهت

مراجع

[1] Csiszar, V., Mori, T. F. and Szekely, G. J. (2005), Chebyshev-type inequalities for scale mixtures, *Statist. Probab. Lett.*, 71, 323-335.

[2] Garrod, N., Pirkovic, S. R. and Valentincic A. (2006), Testing for discontinuity or type of distribution, *Math. Comput. Simul.*, 71, 9-15.

[3] Scott, D. (1992), *Multivariate Density Estimation*, John Wiley, New York.