

## روشی گرافیکی برای انتخاب بهینه کسر نمونه

اکرم فخاری اسفریزی<sup>۱</sup>، ناهید سنجری فارسی پور<sup>۲</sup>

چکیده:

ضریب دمی - وایبل،  $\theta$ ، به عنوان ضریب تغییرات منظم معکوس تابع انباشتگی نرخ شکست تعریف می شود. برآورد این پارامتر دم توزیع بر اساس آماره های ترتیبی بزرگ به دست می آید. بیرلنت و همکاران [۱] یک برآورد بر اساس  $k_n$  تالی بالای آماره های ترتیبی معرفی کردند. بنابراین مساله مهم در برآورد کردن این است که چه تعدادی از نمونه در برآورد استفاده شود. ما با استفاده از روش گرافیکی روشی را برای انتخاب کسر نمونه پیشنهاد می کنیم. **واژه های کلیدی:** ضریب دمی - وایبل، آماره های ترتیبی، نرمال بودن مجانبی.

### ۱ مقدمه

شاخص  $\theta$  در بی نهایت گفته می شود و می نویسیم  $H^{\leftarrow} \in R_{\theta}$  (در این جا علامت  $\leftarrow$  نشان دهنده معکوس تابع است). در (۱) تابع  $l$  می تواند یک تابع ثابت باشد که در این حالت پارامتر  $\theta$  متناسب با پارامتر شکل توزیع وایبل می شود [۳]، توزیع هایی هم با  $l$  غیر ثابت وجود دارند مثل نرمال، گاما، وایبل توسعه یافته<sup>۴</sup> که در (۱) صدق می کنند. چنین توزیع هایی در بیمه و مطالعات قابلیت اعتماد کاربرد دارند.

توزیع نرمال  $N(\mu, \sigma^2)$

فرض کنید  $X_1, X_2, \dots, X_n$  یک دنباله از متغیرهای تصادفی مستقل و هم توزیع با تابع توزیع  $F$  باشند. ما مساله برآورد ضریب دمی - وایبل<sup>۳</sup> را بررسی می کنیم وقتی که دم توزیع در روابط زیر صدق کند:

$$1 - F(x) = e^{-H(x)} \quad x \geq x_0 \geq 0 \quad (1)$$

$$H^{\leftarrow}(t) = \inf\{x : H(x) \geq t\} = t^{\theta} l(t)$$

( $x_0$  نقطه آستانه توزیع است) به طوری که  $l$  یک تابع با آهنگ تدریجی باشد یعنی اگر برای  $\lambda > 0$  وقتی  $x \rightarrow \infty$  داشته باشیم:

$$\frac{l(\lambda, x)}{l(x)} \rightarrow 1$$

توزیعی که در روابط بالا صادق است جز خانواده توزیع دمی - وایبل می باشد و پارامتر  $\theta$  را ضریب دمی - وایبل گویند. تابع  $H^{\leftarrow}$  یک تابع با آهنگ تغییر منظم با

$$H^{\leftarrow}(x) = x^{\frac{1}{\theta}} l(x)$$

$$l(x) = \sqrt{\frac{\sigma}{\Gamma}} \log x - \frac{\sigma}{\Gamma} \frac{\log x}{x} + O(1/x)$$

$$\theta = \frac{1}{\Gamma}$$

<sup>۱</sup> گروه آمار، دانشگاه شیراز  
<sup>۲</sup> گروه آمار، دانشگاه شیراز  
<sup>۳</sup> Weibull Tail-Coefficient  
<sup>۴</sup> Extended Weibull

رکورد پیشنهاد کرد. روش دیگر برآورد از  $k_n$  تای بالایی آماره‌های ترتیبی استفاده می‌کند یعنی برآورد براساس  $x_{n-k_n+1,n} \leq \dots \leq X_{n,n}$  که  $k_n$  یک دنباله از اعداد صحیح است به طوری که  $1 \leq k_n \leq n$  باشد. به این ترتیب برای برآورد  $\theta$ ، کسری از نمونه مورد استفاده قرار می‌گیرد. انتخاب بهینه برای این کسر از نمونه ۵ مساله مورد بررسی در این مقاله است.

برآوردگر گارد براساس همین روش به صورت زیر می‌باشد [۴]:

$$\hat{\theta}_n^G = \frac{\sum_{i=1}^{k_n-1} (\log(X_{n-i+1,n}) - \log(X_{n-k_n+1,n}))}{\sum_{i=1}^{k_n-1} (\log_2(\frac{n}{i}) - \log_2(\frac{n}{k_n}))}$$

$$\log_2(t) = \log(\log(t)) ; t > 1 \text{ که}$$

در [۱] برای به دست آوردن برآوردگر  $\theta$  از تابع  $e(x) = E(X - x | X > x)$  یعنی تابع میانگین باقی مانده عمر استفاده شده است [۱]. تحت (۱) داریم:

$$(\log(x)) \frac{e(H^{\leftarrow}(\log(x)))}{H^{\leftarrow}(\log(x))} \rightarrow \theta \text{ as } x \rightarrow \infty$$

و برآوردگر برای  $\theta$  به فرم زیر خواهد بود:

$$\hat{\theta}_n^{BTV} = \frac{\log(\frac{n}{k_n})}{X_{n-k_n+1,n}} \frac{1}{k_n - 1} \times \sum_{i=1}^{k_n-1} (X_{n-i+1,n} - X_{n-k_n+1,n})$$

برآوردگرهای دیگری برای  $\theta$  توسط برانیاتوسکی [۳]، کلپلبرگ و ویلاسور [۶] ارائه شده‌اند.

## ۲ روش گرافیکی انتخاب $k_n$

کلید اصلی حل مساله برآورد  $\theta$  در انتخاب صحیح  $k_n$  است، چرا که اگر  $k_n$  خیلی بزرگ انتخاب شود ارزیابی

## توزیع گاما $\Gamma(\beta, \alpha)$

$$\begin{aligned} f(x) &= \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} \\ H^{\leftarrow}(x) &= xl(x) \\ l(x) &= \frac{1}{\beta} + \frac{\alpha-1}{\beta} \frac{\log x}{x} + O(1/x) \\ \theta &= 1 \end{aligned}$$

## توزیع وایبل $W(\alpha, \lambda)$

$$\begin{aligned} f(x) &= \frac{\alpha}{\lambda^\alpha} x^{\alpha-1} e^{-\frac{x^\alpha}{\lambda^\alpha}} \\ H^{\leftarrow}(x) &= \lambda x^{\frac{1}{\alpha}} \\ l(x) &= \lambda \\ \theta &= \frac{1}{\alpha} \end{aligned}$$

در اینجا اگر برآورد پارامتر  $\theta$  را بخواهیم در واقع باید برآورد پارامتر شکل توزیع وایبل را پیدا کنیم.

## توزیع وایبل توسیع یافته

$$\begin{aligned} EW(\alpha, \beta) &; \alpha \in (0, 1), \beta \in R \\ F(x) &= 1 - r(x) \exp(-x^\alpha) \quad r \in R_\beta \\ H^{\leftarrow}(x) &= x^{\frac{1}{\alpha}} l(x) \\ l(x) &= 1 + \frac{\beta}{\alpha} \frac{\log x}{x} + O(1/x) \\ \theta &= \frac{1}{\alpha} \end{aligned}$$

چون پارامتر  $\theta$ ، پارامتر دم توزیع است پس برای برآورد آن از بخشی از نمونه استفاده می‌شود. یکی از مهمترین پژوهش‌ها برای برآورد  $\theta$  توسط برد [۲] صورت گرفته است که برآوردگری براساس مقادیر

این نمودار به صورت کاملاً خطی است دلیل این است که چون  $l$  در این توزیع یک تابع ثابت است پس (۱.۲) به صورت تقریبی نخواهد بود و یک رابطه دقیق است.

در این روش گرافیکی، از دیدگاهی دیگر، می توان یک برآورد برای  $\theta$  به دست آورد. با توجه به این که می دانیم حداقل برای  $i$  های بزرگ نمودار خطی خواهد شد، پس می توان شبیهی را که برای  $i$  های بزرگ وجود دارد را به عنوان برآوردی برای  $\theta$  به حساب آورد (شکل ۱).

لازم به ذکر است در نمودارها  $|N|$  منظور قدرمطلق مقادیر شبیه سازی شده از توزیع نرمال است.

روش گرافیکی دیگر با استفاده از نمودار چندک، بر اساس برآوردگر بیرلنت و همکاران [۱] می باشد که به صورت زیر است:

$$\hat{\theta}_n^{BBTV} = \frac{\log\left(\frac{n}{k_n}\right)}{X_{n-k_n+1,n} - k_n - 1} \times \sum_{i=1}^{k_n-1} (X_{n-i+1,n} - X_{n-k_n+1,n})$$

با تعریف کردن  $v_{m,n}$

$$v_{m,n} = \frac{X_{n-m}}{\frac{1}{m} \sum_{i=1}^m (X_{n-i+1,n} - X_{n-m,n})}$$

خواهیم داشت:

$$\hat{\theta}_n^{BBTV} = \left( \frac{1}{v_{k_n-1,n}} \right) / \left( \frac{1}{\log\left(\frac{n}{k_n}\right)} \right)$$

اگر نمودار نقاط  $\left( \frac{1}{\log\left(\frac{n}{k_n}\right)}, \frac{1}{v_{k_n-1,n}} \right)$  رسم شود، باز هم جایی که نمودار شروع به خطی شدن می کند را انتخاب مناسبی برای  $k_n$  در نظر می گیریم.

از دیدگاهی دیگر این نمودار دارای شیب  $\theta$  می باشد. با در نظر گرفتن این مطلب می توان برآوردی برای  $\theta$  به دست آورد.

$\hat{\theta}_n$  بزرگ می شود و اگر  $k_n$  خیلی کوچک انتخاب شود، واریانس  $\hat{\theta}_n$  خیلی بزرگ می شود.

یکی از روش ها برای انتخاب  $k_n$  بهینه می تواند با استفاده از نمودار چندک، بر اساس برآوردگر گارد باشد. این روش بر اساس تئوری زیر می باشد:

$$q(t) = (1 - F)^{\leftarrow}(t) = H^{\leftarrow}\left(\log\left(\frac{1}{t}\right)\right)$$

تحت (۱)

$$\begin{aligned} \log(q(t)) - \log(q(s)) &= \theta \left( \log_2\left(\frac{1}{s}\right) \right. \\ &\quad \left. - \log_2\left(\frac{1}{t}\right) + \log\left(\frac{l(\log(\frac{1}{t}))}{l(\log(\frac{1}{s}))}\right) \right) \end{aligned}$$

چون  $l$  یک تابع با آهنگ تغییر تدریجی است پس

$$\log(q(t)) - \log(q(s)) \approx \theta \left( \log_2\left(\frac{1}{t}\right) - \log_2\left(\frac{1}{s}\right) \right)$$

این دقیقاً مثل این است که معادله خطی با شیب  $\theta$  نوشته ایم که این رابطه حداقل برای  $t$  های کوچک خطی خواهد بود.

در این جا اگر  $t = \frac{1}{n}$  را در نظر بگیریم با توجه به این که  $\hat{q}\left(\frac{1}{n}\right) = X_{n-i+1,n}$  را به عنوان برآورد ناپارامتری از چندک ها در نظر می گیریم. پس باید  $(\log_2\left(\frac{n}{t}\right), X_{n-i+1,n})$  یک رابطه خطی تقریبی با شیب  $\theta$  داشته باشند که در [۴] آمده است.

در این روش گرافیکی، روشی که برای انتخاب  $k_n$  بهینه انتخاب می شود بدین صورت است که  $k_n$  را طوری انتخاب می کنیم تا نمودار خطی شود. البته اشکال این روش این است که یک انتخاب صریح به ما نمی دهد و انتخاب شهودی است. در نمودار مربوط به توزیع وایبل

که  $AMSE(\hat{\theta}_n^{MI}(k_n))$  را حداقل کند.

$$\widehat{AMSE}(\hat{\theta}_n^{MI}(k_n)) = \frac{(\hat{\theta}_n^{LS}(k_n))^2}{k_n} + \left( \hat{b} \left( \log\left(\frac{n}{k_n}\right) \frac{1}{k_n} \sum_{j=1}^{k_n} \left( \frac{\log\left(\frac{n}{j}\right)}{\log\left(\frac{n}{k_n}\right)} \right)^{-1} \right) \right)^2$$

که این رابطه برآوردی از رابطه (۲) می باشد.

پس برآوردی که برای  $k_n$  به دست می آید به صورت زیر می باشد:

$$\hat{k}_n = \arg \min_{k_n} \widehat{AMSE}(\hat{\theta}_n^{MI}(k_n))$$

استفاده از روش فوق در مطالعات شبیه سازی به شرح زیر است:

از توزیع های  $W(4, 4), |N|(0, 1), \Gamma(4, 1), \Gamma(0.25, 1)$  و  $W(0.25, 0.25)$  به تعداد ۱۰۰ بار نمونه ۵۰۰ تایی شبیه سازی شده است. پس  $N = 100$  و  $n = 500$

$$(x_{n,i})_{i=1, \dots, N}$$

روی هر نمونه  $x_{n,i}$  برآورد  $\hat{\theta}_n^{MI}(k_n)$  را به عنوان تابعی از  $k_n$  به دست می آید که این برآوردها به ازای  $k_n = 2, \dots, 350$  می باشد.

به منظور یافتن انتخاب مناسب برای  $k_n$  مقدار زیر را محاسبه کرده

$$\hat{k}_n = \arg \min_{k_n \in [2, 350]} \widehat{AMSE}(\hat{\theta}_{n,i}^{MI}(k_n)) ; i = 1, \dots, N$$

سپس

$$\mu(\hat{k}_n) = \frac{1}{N} \sum_{i=1}^N \hat{k}_{n,i}$$

$$\sigma(\hat{k}_n) = \sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{k}_{n,i} - \mu(\hat{k}_n))^2}$$

$$K_n^{opt} = \arg \min_{k_n \in [2, 350]} AMSE(\hat{\theta}_n^{MI}(k_n))$$

## ۳ روش تطبیقی برای انتخاب $k_n$

در یک روش دیگر برای به دست آوردن  $k_n$  بهینه، دیابلت و همکاران [۴] میانگین مربعات خطای مجانبی مربوط به  $\hat{\theta}_n^{ML}(k_n)$  را به صورت

$$AMSE(\hat{\theta}_n^{ML}(k_n)) = \frac{\theta^2}{k_n} + \left( b \left( \log\left(\frac{n}{k_n}\right) \frac{1}{k_n} \sum_{i=1}^{k_n} \left( \frac{\log\left(\frac{n}{i}\right)}{\log\left(\frac{n}{k_n}\right)} \right)^n \right) \right)^2 \quad (2)$$

ارایه کردند که برآورد حداکثر درست‌نمایی  $\theta$  است [۳].

در [۴] به منظور به دست آوردن توزیع حدی برآوردگر معرفی شده، شرط دومی برای توزیع های دمی - وایبل مطرح کرده است: یک  $\eta \geq 0$  و یک تابع  $b(x)$  وجود داشته باشند، به طوری که وقتی  $x \rightarrow \infty$  آن گاه  $b(x) \rightarrow 0$ . همچنین برای هر  $1 < A < \infty$  به طوری که برای هر  $\lambda \in [1, A]$

$$\log\left(\frac{l(\lambda x)}{l(x)}\right) \sim b(x)k_\eta(\lambda) \quad (3)$$

$$k_\eta(\lambda) = \int_1^\lambda u^{\eta-1} du = \frac{1}{\eta}(\lambda^\eta - 1)$$

(منظور از  $\sim$  هم‌ارزی می باشد) پارامتر  $\eta \geq 0$  نرخ همگرایی  $\frac{l(\lambda x)}{l(x)}$  به ۱ را نشان می دهد. هر چه  $\eta$  به صفر

نزدیکتر باشد نرخ همگرایی آهسته تر است و نیز داریم که

$$|b| \in R_\eta$$

برآورد  $\eta$  بسیار پیچیده است به همین علت مقدار کانونی ۱- برای آن در نظر می گیریم. مقدار  $\eta$  را در رابطه (۲)

برابر ۱- قرار دهیم این کار باعث سادگی محاسبات می شود. همچنین در رابطه (۲) برآوردهای حداقل

مربعات  $\theta$  و  $B$  را قرار دهیم و به دنبال  $k_n$  ای می گردیم

همان‌طور که پیداست نتایج بهتری برای برآورد  $\hat{\theta}_n$  به دست می‌آید. برای  $\Gamma(1, 25)$  مقدار  $k_n^{opt} = 186$  شده که نشان می‌دهد برای برآورد  $\theta$  از آماره ترتیبی ۱۸۶ تا ۵۰۰ استفاده شده است.

مقدار ایده آل برای  $R_n$  برابر ۱ می‌باشد. در جدول  $R_n$  معمولاً نزدیک ۱ است بجز برای توزیع وایبل، که در مورد این توزیع مقادیر بزرگ  $R_n$  و همچنین مقادیر بزرگ  $\mu(\hat{k}_n)$  نشان می‌دهد که  $k_n$  بهینه بزرگتر از ۳۵۰ انتخاب می‌شود.

برآوردهای  $\theta$  را بر اساس  $\hat{K}_{n,i}$  به دست می‌آوریم:

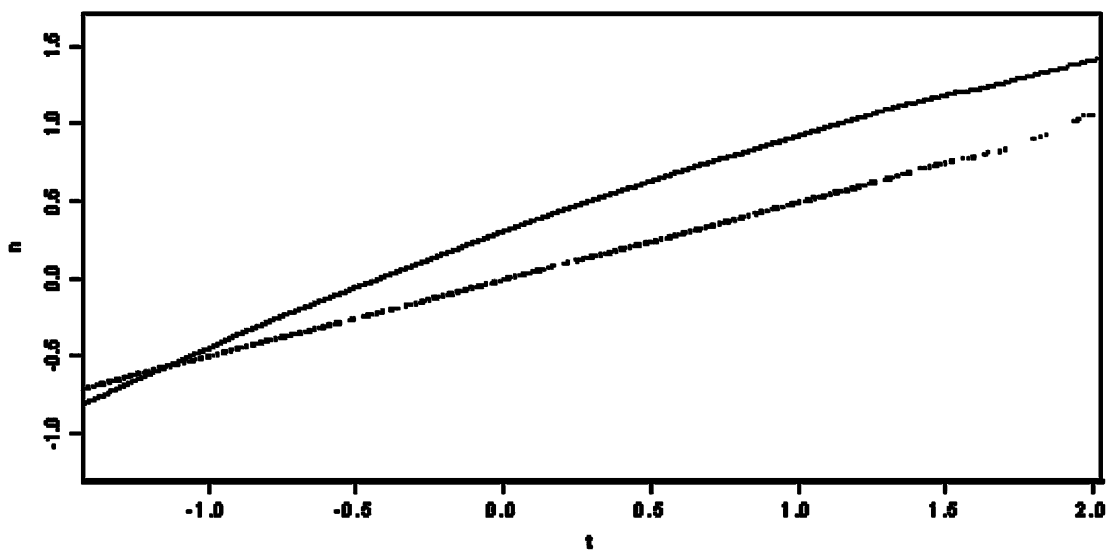
$$\mu(\hat{\theta}_n^{ML}) = \frac{1}{N} \sum_{i=1}^N \hat{\theta}_{n,i}^{ML}(\hat{k}_{n,i})$$

$$\sigma(\hat{\theta}_n^{ML}) = \sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{\theta}_{n,i}^{ML}(\hat{k}_{n,i}) - \mu(\hat{\theta}_n^{ML}))^2}$$

همچنین کمیت  $R_n^\gamma$  یعنی نسبت میانگین توان دوم خطای تجربی به حداقل میانگین توان دوم می‌تواند به عنوان یک معیار به منظور انتخاب بهینه  $k_n$  استفاده گردد.

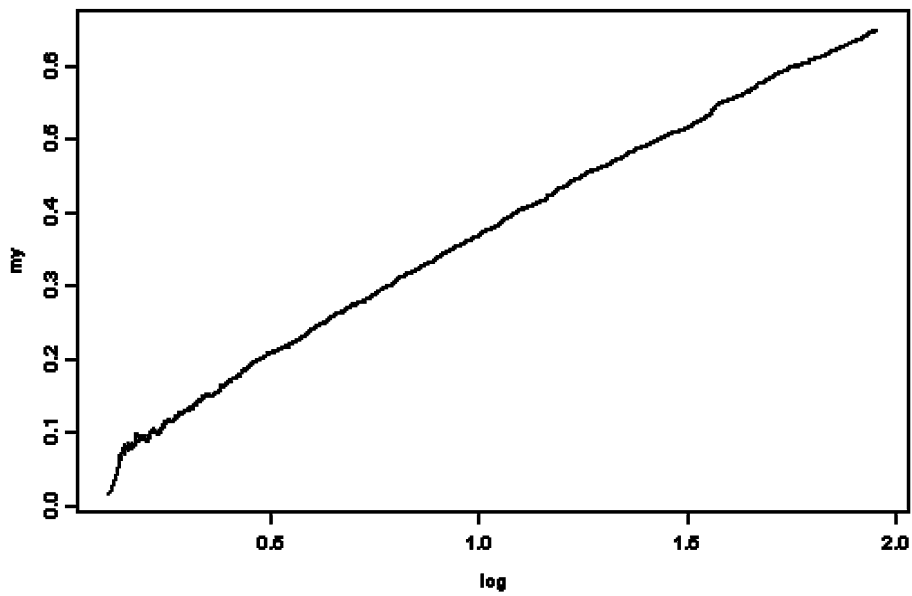
$$R_n^\gamma = \frac{\sum_{i=1}^N (\hat{\theta}_{n,i}^{ML}(\hat{k}_{n,i}) - \theta)^\gamma}{\min_{k_n \in [2, 350]} \sum_{i=1}^N (\hat{\theta}_n^{ML}(k_n) - \theta)^\gamma}$$

با استفاده از شبیه‌سازی، روش فوق در به دست آوردن  $k_n$  استفاده می‌شود که نتایج را در جدول (۱) قرار داده‌ایم.

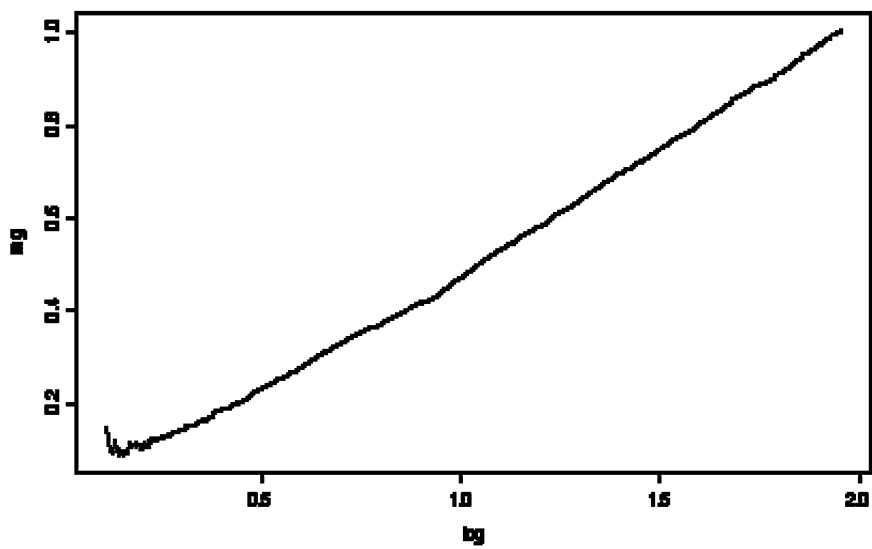


شکل ۱. نمودار چندک مربوط به توزیع  $W(2, 1)$  و توزیع  $|N|(1, 1)$

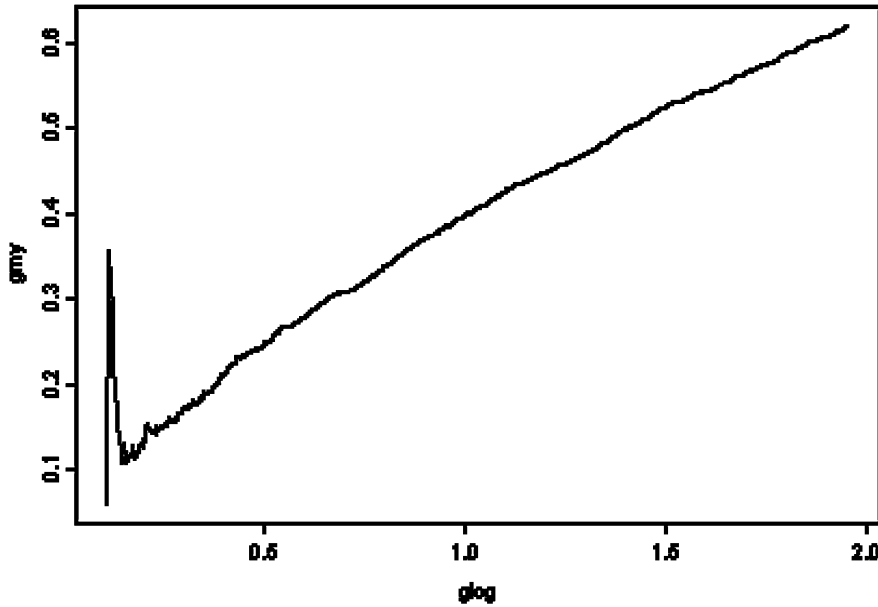
شبیه‌سازی به تعداد ۱۰۰۰۰ نمونه تصادفی



شکل ۲. نمودار مربوط به توزیع  $W(2, 1)$  با ۱۰۰۰۰ نمونه



شکل ۳. نمودار مربوط به توزیع  $|N|(2, 1)$  با ۱۰۰۰۰ نمونه



شکل ۴. نمودار مربوط به توزیع  $\Gamma(4, 1)$  با ۱۰۰۰۰۰ نمونه

جدول ۱. نتایج شبیه‌سازی برای به دست آوردن انتخاب بهینه  $k_n$ .

توزیع	$\theta$	$\eta$	$\mu(\hat{k}_n)$	$\sigma(\hat{k}_n)$	$\mu(\hat{\theta}_n^{ML})$	$\sigma(\hat{\theta}_n^{ML})$	$R_n$	$k_n^{opt}$
$\Gamma(0/25, 1)$	۱	-۱	۱۰۵/۵	۶۲/۲	۱/۶۶۷	۰/۲۹۴	۱/۲۶	۱۸۶
$\Gamma(4, 1)$	۱	-۱	۲۲۲/۷	۸۲/۱	۰/۰۵۴۸	۰/۰۵۱	۱/۱۳	۱۸۴
$ N (0, 1)$	۰/۵	-۱	۲۴۶/۶	۸۱/۱	۰/۶۷۹	۰/۰۱۰۹	۱/۲۱	۱۸۹
$W(25, 25)$	۴	$-\infty$	۳۰۵/۸	۵۹/۰	۴/۰۱۶	۰/۲۶۵	۱/۶۲	۳۵۰
$W(4, 4)$	۰/۲۵	$-\infty$	۳۱۰/۴	۵۰/۹	۰/۲۴۹	۰/۰۱۳	۱/۴۳	۳۵۰

## مراجع

- [1] Beirlant, J., Broniatowski, M., Teugels, J. L. and Vynckier, P. (1995), The mean residual
- [2] Berred, m. (1991), Record values and the estimation of the weibull tail-coefficient, *comptes-rendus de l'Academie des sciences. Serie I*, 312, 943-946.
- [3] Broniatowski, M. (1993). On The Estimation of Weibull Tail Coefficient, *J.Statist. Plan. Inference*, 35, 349-366.

- 
- [4] Diebolt, J., Gardes, S. L., Girard, S. and Guillo, A. (2006), Bias- Reduced Estimators of the Weibull Tail-Coefficient, Test-girad.pdf. 10 October 2006.
- [5] Girard, S. (2004), A Hill Type Estimate of the Weibull Tail-Coefficient, *Commun. Statist. Theor. Meth*, 33(2), 205-234.
- [6] Kluppelberg, C. and Villasenor, J. A. (1993), Estimation of distribution tails- a semiparametric approach, *deutschen gesellschaft fur versicherungs mathematik XXI(2)*, 213-235.