

СИБИРСКИЕ ЭЛЕКТРОННЫЕ МАТЕМАТИЧЕСКИЕ ИЗВЕСТИЯ

Siberian Electronic Mathematical Reports

<http://semr.math.nsc.ru>

Том 20, стр. 144–144 (2023)
DOI 10.33048/semi.2023.20.xxx

УДК 519.233
MSC 62F03

LIMIT THEOREMS FOR FORWARD AND BACKWARD PROCESSES OF NUMBERS OF NON-EMPTY URNS IN INFINITE URN SCHEMES

M.G. CHEBUNIN, A.P. KOVALEVSKII

ABSTRACT.

We study the joint asymptotics of forward and backward processes of numbers of non-empty urns in an infinite urn scheme. The probabilities of balls hitting the urns are assumed to satisfy the conditions of regular decrease. We prove weak convergence to a two-dimensional Gaussian process. Its covariance function depends only on exponent of regular decrease of probabilities. We obtain parameter estimates that have a normal asymptotics for its joint distribution together with forward and backward processes. We use these estimates to construct statistical tests for the homogeneity of the urn scheme on the number of thrown balls.

Keywords: Zipf's law, weak convergence, Gaussian process, statistical test.

1. INTRODUCTION

Let X_1, \dots, X_n be independent and identically distributed positive integer-valued random variables,

$$(1) \quad p_i = \mathbf{P}(X_1 = i), \quad \sum_{i=1}^{\infty} p_i = 1.$$

CHEBUNIN, M.G., KOVALEVSKII, A.P., LIMIT THEOREMS FOR FORWARD AND BACKWARD PROCESSES OF NUMBERS OF NON-EMPTY URNS IN INFINITE URN SCHEMES.

© 2023 CHEBUNIN M.G., KOVALEVSKII A.P..

The work is supported by Mathematical Center in Akademgorodok under agreement No. 075-15-2019-1675 with the Ministry of Science and Higher Education of the Russian Federation.

Received November, 1, 2022, published May, 1, 2023.

The number of different elements among first k ones ($2 \leq k \leq n$) is

$$(2) \quad R_k = 1 + \sum_{i=2}^k \mathbf{1}(X_i \notin \{X_1, \dots, X_{i-1}\}).$$

Similarly, the number of different elements among *last* k ones ($2 \leq k \leq n$) is

$$(3) \quad R'_k = 1 + \sum_{i=n-k+1}^{n-1} \mathbf{1}(X_i \notin \{X_{n-k+2}, \dots, X_n\}).$$

In other words,

$$(4) \quad R'_k = 1 + \sum_{i=2}^k \mathbf{1}(X'_i \notin \{X'_1, \dots, X'_{i-1}\}),$$

with $X'_i = X_{n-i+1}$, the random variables in the backward order, $1 \leq i \leq n$.

We put by definition

$$(5) \quad R_0 = R'_0 = 0, \quad R_1 = R'_1 = 1.$$

Distributions of R_n and R'_n are identical, their limiting properties are known. We study their limiting joint distribution under the appropriate centering and normalizing.

If there is an infinite number of positive probabilities in (1) then this probability model is the infinite urn scheme. Karlin (1967) established the SLLN for R_n in the infinite urn scheme (Bahadur (1960) proved the weak LLN),

$$(6) \quad R_n / \mathbf{E}R_n \rightarrow 1 \quad \text{a.s.}$$

Now we need the regularity condition. Let $p_1 \geq p_2 \geq \dots > 0$ and

$$(7) \quad \alpha(x) := \max\{k > 0 : p_k \geq 1/x\} = x^\theta L(x) \quad \text{as } x \rightarrow \infty, \quad 0 < \theta < 1,$$

$L(\cdot)$ is the slowly varying function of the real argument: $L(tx)/L(x) \rightarrow 1$ as $x \rightarrow +\infty$ for any real $t > 0$.

From Karamata's characterization theorem, $\alpha(x)$ is a regular varying function with index θ . The model (7) is the elementary probability model that corresponds to the Zipf's Law (Zipf, 1936) of power decreasing of word probabilities.

Karlin (1967) proved the CLT: if (7) holds then $(R_n - \mathbf{E}R_n) / \sqrt{\mathbf{Var}R_n}$ converges weakly to the standard normal distribution,

$$(8) \quad \mathbf{E}R_n \sim \Gamma(1 - \theta)\alpha(n), \quad \mathbf{Var}R_n / \mathbf{E}R_n \rightarrow 2^\theta - 1,$$

$\Gamma(\cdot)$ is the Euler gamma.

From the Karlin's CLT and (8), $(R_n - \mathbf{E}R_n) / \sqrt{\mathbf{E}R_n}$ converges weakly to the centered normal distribution with variance $2^\theta - 1$. The CLT holds for $\theta = 1$ too but with another normalization.

Chebunin and Kovalevskii (2016) proved the Functional CLT: if (7) holds then the process

$$(9) \quad Z_n = \{Z_n(t), 0 \leq t \leq 1\} = \{(R_{[nt]} - \mathbf{E}R_{[nt]}) / \sqrt{\mathbf{E}R_n}, 0 \leq t \leq 1\}$$

converges weakly in $D(0, 1)$ with uniform metric to a centered Gaussian process Z_θ with continuous a.s. sample paths and covariance function

$$(10) \quad K(s, t) = (s + t)^\theta - \max(s^\theta, t^\theta).$$

The Karlin's CLT is a particular case of the FCLT for $Z_n(1)$. The same FCLT is true for

$$(11) \quad Z'_n = \{Z'_n(t), 0 \leq t \leq 1\} = \{(R'_{[nt]} - \mathbf{E}R_{[nt]})/\sqrt{\mathbf{E}R_n}, 0 \leq t \leq 1\}.$$

We prove the theorem about the joint limiting distribution of (Z_n, Z'_n) .

All the papers on properties of R_n and similar statistics in the infinite urn scheme can be divided into 4 types:

1. Results under the regularity condition (7): the papers above, Durieu & Wang (2016), Chebunin (2017), Durieu, Samorodnitsky & Wang (2020), Chebunin & Zuyev (2020).

2. Results under (7) with $\theta = 0$ instead of $0 < \theta \leq 1$, that is, for the slowly varying function $\alpha(x)$: Dutko (1989), Barbour (2009), Barbour & Gneden (2009).

3. Results for the model without assuming the regularity condition (7): Key (1992, 1996), Hwang & Janson (2008), Muratov & Zuyev (2016), Ben-Hamou, Boucheron & Ohannessian (2017), Decrouez, Grabchak & Paris (2018).

4. Statistical applications — we postpone the survey of these results to Section 2.

Gneden, Hansen & Pitman (2007) made a detailed survey of the results of types 1–3 existed at the time.

2. MAIN RESULTS

Theorem 1. *If (7) holds then the process $(Z_n, Z'_n) = \{(Z_n(t), Z'_n(t)), 0 \leq t \leq 1\}$ converges weakly in the uniform metric in $D(0, 1)^2$ to 2-dimensional Gaussian process (Z, Z') with zero expectation and covariance function*

$$\mathbf{E}Z(s)Z(t) = \mathbf{E}Z'(s)Z'(t) = K(s, t), \quad \mathbf{E}Z(s)Z'(t) = K'(s, t),$$

where $K(s, t)$ is given by (10), and

$$(12) \quad K'(s, t) = ((s + t)^\theta - 1)\mathbf{1}(s + t > 1).$$

From Theorem 1 we have that the limiting process $\{(Z(t) - Z'(t))/\sqrt{2}, 0 \leq t \leq 1/2\}$ is the stochastically self-similar process which coincide in distribution with the limiting process of Durieu and Wang (2016). So Theorem 1 gives an alternative way to simulate these processes without additional randomization.

We need some estimate of the unknown parameter θ to use the theorem in applications. Various classes of such estimates have been obtained and analysed by Hill (1975), Nicholls (1978), Zakrevskaya and Kovalevskii (2001, 2019), Guillou and Hall (2002), Ohannessian and Dahleh (2012), Chebunin (2014), Chebunin and Kovalevskii (2019a, 2019b), Chakrabarty et al. (2020).

But we need an estimate that is symmetric to the forward and backward processes. Moreover, we want to have the limiting joint distribution of the estimate and the two-dimensional process. We introduce the estimate and study its properties in the next section.

3. PARAMETER'S ESTIMATION

From (6) and (8), we have $\log R_n \sim \theta \log n$ a.s. Therefore, we may propose the following estimators for parameter θ :

$$\theta_n = \int_0^1 \log^+ R_{[nt]} dA(t), \quad \theta'_n = \int_0^1 \log^+ R'_{[nt]} dA(t),$$

here $\log^+ x = \max(\log x, 0)$. Function $A(\cdot)$ has bounded variation and

$$(13) \quad A(0) = A(1) = 0, \quad \lim_{x \downarrow 0} \log x \int_0^x |dA(t)| = 0, \quad \int_0^1 \log t dA(t) = 1.$$

Let

$$\widehat{\theta} = (\theta_n + \theta'_n)/2.$$

Theorem 2. *Let $p_i = i^{-1/\theta} l(i, \theta)$, $\theta \in [0, 1]$, and $l(x, \theta)$ is a slowly varying function as $x \rightarrow \infty$. Then the estimator $\widehat{\theta}$ is strongly consistent.*

Proof. The proof follows from the definition of $\widehat{\theta}$ and Theorem 1 from Chebunin and Kovalevskii (2019). \square

We need extra conditions to obtain the asymptotic normality of $\widehat{\theta}$.

Theorem 3. *Let $p_i = ci^{-1/\theta}(1 + o(i^{-1/2}))$, $\theta \in (0, 1)$, and $A(t) = 0$, $t \in [0, \delta]$ for some $\delta \in (0, 1)$. Then*

$$\sqrt{\mathbf{E}R_n}(\widehat{\theta} - \theta) - \frac{1}{2} \int_0^1 t^{-\theta} (Z_n(t) + Z'_n(t)) dA(t) \rightarrow_p 0.$$

Proof. The proof follows from the definition of $\widehat{\theta}$ and Theorem 2 from Chebunin and Kovalevskii (2019). \square

From Theorem 3, it follows that $\widehat{\theta}$ converges to θ at rate $(\mathbf{E}R_n)^{-1/2}$, and $\sqrt{\mathbf{E}R_n}(\widehat{\theta} - \theta)$ converges weakly to the normal random variable $\frac{1}{2} \int_0^1 t^{-\theta} (Z_\theta(t) + Z'_\theta(t)) dA(t)$ with variance $\frac{1}{2} \int_0^1 \int_0^1 (st)^{-\theta} (K(s, t) + k(s, t)) dA(s) dA(t)$.

Example 1 Take

$$A(t) = \begin{cases} 0, & 0 \leq t \leq 1/2; \\ -(\log 2)^{-1}, & 1/2 < t < 1; \\ 0, & t = 1. \end{cases}$$

Then

$$\widehat{\theta} = \log_2 \left(R_n / \sqrt{R_{\lfloor n/2 \rfloor} R'_{\lfloor n/2 \rfloor}} \right), \quad n \geq 2.$$

4. TEST FOR A KNOWN RATE

Let $0 < \theta < 1$ be known. We introduce *empirical bridges* $\overset{\circ}{Z}_n, \overset{\circ}{Z}'_n$ (Kovalevskii and Shatalin, 2015, 2016) as follows.

$$\overset{\circ}{Z}_n(k/n) = (R_k - (k/n)^\theta R_n) / \sqrt{R_n}, \quad \overset{\circ}{Z}'_n(k/n) = (R'_k - (k/n)^\theta R_n) / \sqrt{R_n},$$

$0 \leq k \leq n$, where $R_0 = 0$. We construct a piecewise linear approximation: for any $0 \leq u < 1/n$ and $0 \leq k \leq n-1$,

$$\overset{\circ}{Z}_n \left(\frac{k}{n} + u \right) = \overset{\circ}{Z}_n(k/n) + nu \left(\overset{\circ}{Z}_n((k+1)/n) - \overset{\circ}{Z}_n(k/n) \right),$$

$$\overset{\circ}{Z}'_n \left(\frac{k}{n} + u \right) = \overset{\circ}{Z}'_n(k/n) + nu \left(\overset{\circ}{Z}'_n((k+1)/n) - \overset{\circ}{Z}'_n(k/n) \right).$$

Theorem 4. *Under the assumptions of Theorem 2,*

$$\sup_{0 \leq t \leq 1} |\overset{\circ}{Z}_n(t) - (Z_n(t) - t^\theta Z_n(1))| \rightarrow 0 \text{ a.s.}$$

$$\sup_{0 \leq t \leq 1} |\overset{\circ}{Z}'_n(t) - (Z'_n(t) - t^\theta Z'_n(1))| \rightarrow 0 \text{ a.s.}$$

Proof. The first statement is Theorem 3 from Chebunin and Kovalevskii (2019). For the second statement, let $t \in [0, 1]$, and $k = [nt]$, then $t = k/n + u$, $0 \leq k \leq n-1$, $u \in [0, 1/n)$. Let $f_\theta(x) = (1+x)^\theta - x^\theta$. So $0 \leq f_\theta(x) \leq f_\theta(0) = 1$ for $x \geq 0$.

By the definition of $\overset{\circ}{Z}'_n(t)$,

$$\frac{R'_k - \left(\frac{k+1}{n}\right)^\theta R_n}{\sqrt{R_n}} \leq \overset{\circ}{Z}'_n(t) \leq \frac{R'_{k+1} - \left(\frac{k}{n}\right)^\theta R_n}{\sqrt{R_n}},$$

so

$$\begin{aligned} \left| \overset{\circ}{Z}'_n(t) - \frac{R_{[nt]} - t^\theta R_n}{\sqrt{R_n}} \right| &\leq \frac{R'_{k+1} - R'_k + \frac{1}{n^\theta} f_\theta(k) R_n}{\sqrt{R_n}} \\ &\leq \frac{1}{\sqrt{R_n}} + \frac{\sqrt{R_n}}{n^\theta} \rightarrow 0 \end{aligned}$$

a.s. uniformly on $t \in [0, 1]$. □

Let $C(0,1)$ be the set of all continuous functions on $[0, 1]$ with the uniform metric $\rho(x, y) = \max_{t \in [0, 1]} |x(t) - y(t)|$. By Theorem 1, we have

Corollary 1. *Under the assumptions of Theorem 4, $(\overset{\circ}{Z}_n, \overset{\circ}{Z}'_n)$ converges weakly in $C(0, 1)$ to 2-dimensional Gaussian process $(\overset{\circ}{Z}, \overset{\circ}{Z}')$ that can be represented as $(\overset{\circ}{Z}(t), \overset{\circ}{Z}'(t)) = (Z_\theta(t) - t^\theta Z_\theta(1), Z'_\theta(t) - t^\theta Z'_\theta(1))$, $0 \leq t \leq 1$. Its correlation function is given by the covariance function*

$$c_{R,R}(s, t) = c_{R',R'}(s, t) = \overset{\circ}{K}(s, t), \quad c_{R,R'}(s, t) = \overset{\circ}{K}'(s, t),$$

where

$$\begin{aligned} \overset{\circ}{K}(s, t) &= K(s, t) - s^\theta K(1, t) - t^\theta K(s, 1) + s^\theta t^\theta K(1, 1), \\ \overset{\circ}{K}'(s, t) &= K'(s, t) - s^\theta K'(1, t) - t^\theta K'(s, 1) + s^\theta t^\theta K'(1, 1). \end{aligned}$$

Now we show how to implement the goodness-of-fit test in this case.

Let $W_n^2 = \int_0^1 \left(\overset{\circ}{Z}_n(t) \right)^2 + \left(\overset{\circ}{Z}'_n(t) \right)^2 dt$. It is equal to

$$(14) \quad \begin{aligned} W_n^2 &= \frac{1}{3n} \sum_{k=1}^{n-1} \overset{\circ}{Z}_n \left(\frac{k}{n} \right) \left(2\overset{\circ}{Z}_n \left(\frac{k}{n} \right) + \overset{\circ}{Z}_n \left(\frac{k+1}{n} \right) \right) \\ &\quad + \frac{1}{3n} \sum_{k=1}^{n-1} \overset{\circ}{Z}'_n \left(\frac{k}{n} \right) \left(2\overset{\circ}{Z}'_n \left(\frac{k}{n} \right) + \overset{\circ}{Z}'_n \left(\frac{k+1}{n} \right) \right). \end{aligned}$$

Then W_n^2 converges weakly to $W_\theta^2 = \int_0^1 \left(\overset{\circ}{Z}_\theta(t) \right)^2 + \left(\overset{\circ}{Z}'_\theta(t) \right)^2 dt$.

So the test rejects the null hypothesis if $W_n^2 \geq C$. The p-value of the test is $1 - F_\theta(W_{n,obs}^2)$. Here F_θ is the cumulative distribution function of W_θ^2 and $W_{n,obs}^2$ is the value of W_n^2 for observations under consideration.

One can estimate F_θ by simulations or find it explicitly using the Smirnov formula (Smirnov, 1937): if $W_\theta^2 = \sum_{k=1}^{\infty} \frac{\eta_k^2}{\lambda_k}$, η_1, η_2, \dots are independent and have standard normal distribution, $0 < \lambda_1 < \lambda_2 < \dots$, then

$$(15) \quad F_\theta(x) = 1 + \frac{1}{\pi} \sum_{k=1}^{\infty} (-1)^k \int_{\lambda_{2k-1}}^{\lambda_{2k}} \frac{e^{-\lambda x/2}}{\sqrt{-D(\lambda)}} \cdot \frac{d\lambda}{\lambda}, \quad x > 0,$$

$$D(\lambda) = \prod_{k=1}^{\infty} \left(1 - \frac{\lambda}{\lambda_k}\right).$$

The integrals in the RHS of (15) must tend to 0 monotonically as $k \rightarrow \infty$, and λ_k^{-1} are the eigenvalues of the kernel (see Martynov (1973), Chapter 3).

5. TEST FOR AN UNKNOWN RATE

Let us introduce the process $(\widehat{Z}_n, \widehat{Z}'_n)$:

$$\widehat{Z}_n(k/n) = \left(R_k - (k/n)^\theta R_n\right) / \sqrt{R_n}, \quad \widehat{Z}'_n(k/n) = \left(R'_k - (k/n)^\theta R_n\right) / \sqrt{R_n},$$

$0 \leq k \leq n$. As for $\overset{\circ}{Z}_n$, let for $0 \leq u < 1/n$ and $0 \leq k \leq n-1$

$$\widehat{Z}_n\left(\frac{k}{n} + u\right) = \widehat{Z}_n(k/n) + nu \left(\widehat{Z}_n((k+1)/n) - \widehat{Z}_n(k/n)\right),$$

$$\widehat{Z}'_n\left(\frac{k}{n} + u\right) = \widehat{Z}'_n(k/n) + nu \left(\widehat{Z}'_n((k+1)/n) - \widehat{Z}'_n(k/n)\right).$$

Theorem 5. *Under assumptions of Theorem 3, $(\widehat{Z}_n, \widehat{Z}'_n)$ converges weakly as $n \rightarrow \infty$ to 2-dimensional Gaussian process $(\widehat{Z}_\theta, \widehat{Z}'_\theta)$ that can be represented as $(\widehat{Z}_\theta(t), \widehat{Z}'_\theta(t))$, $0 \leq t \leq 1$, where*

$$\widehat{Z}_\theta(t) = \overset{\circ}{Z}_\theta(t) - \frac{t^\theta \log t}{2} \int_0^1 u^{-\theta} (Z_\theta(u) + Z'_\theta(u)) dA(u),$$

$$\widehat{Z}'_\theta(t) = \overset{\circ}{Z}'_\theta(t) - \frac{t^\theta \log t}{2} \int_0^1 u^{-\theta} (Z_\theta(u) + Z'_\theta(u)) dA(u).$$

Proof. Similarly to the proof of Theorem 4 from Chebunin and Kovalevskii (2019), we can show that

$$\sup_{t \in [0,1]} \left| \widehat{Z}_n(t) - \overset{\circ}{Z}_n(t) - \sqrt{R_n}(\widehat{\theta} - \theta)t^\theta \log t \right| \rightarrow_p 0,$$

$$\sup_{t \in [0,1]} \left| \widehat{Z}'_n(t) - \overset{\circ}{Z}'_n(t) - \sqrt{R_n}(\widehat{\theta} - \theta)t^\theta \log t \right| \rightarrow_p 0.$$

Let us demonstrate it for $\widehat{Z}'_n(t)$. Let $t \in [0, 1]$, $k = [nt]$, $u = t - k/n$, $f_\theta(x) = (1+x)^\theta - x^\theta$ as in the proof of Theorem 4. By the definition,

$$\widehat{Z}'_n(k/n) = \overset{\circ}{Z}'_n(k/n) + \sqrt{R_n} \left((k/n)^\theta - (k/n)^{\widehat{\theta}} \right),$$

$$\widehat{Z}'_n(t) = \overset{\circ}{Z}'_n(t) + \sqrt{R_n} \left((k/n)^\theta - (k/n)^{\widehat{\theta}} \right)$$

$$+nu\sqrt{R_n} \left(\left(\frac{k+1}{n} \right)^\theta - \left(\frac{k+1}{n} \right)^{\hat{\theta}} - \left(\frac{k}{n} \right)^\theta + \left(\frac{k}{n} \right)^{\hat{\theta}} \right).$$

We have

$$\left(\frac{k+1}{n} \right)^\theta - \left(\frac{k}{n} \right)^\theta = f_\theta(k)/n^\theta, \quad \left(\frac{k+1}{n} \right)^{\hat{\theta}} - \left(\frac{k}{n} \right)^{\hat{\theta}} = f_{\hat{\theta}}(k)/n^{\hat{\theta}},$$

so

$$\begin{aligned} & \left| \widehat{Z}'_n(t) - \overset{\circ}{Z}'_n(t) + \sqrt{R_n} (t^{\hat{\theta}} - t^\theta) \right| \\ &= \left| \widehat{Z}'_n(t) - \overset{\circ}{Z}'_n(t) + \sqrt{R_n} \left(\left(\frac{k}{n} + u \right)^{\hat{\theta}} - \left(\frac{k}{n} + u \right)^\theta \right) \right| \\ &\leq 2\sqrt{R_n} (f_\theta(k)/n^\theta + f_{\hat{\theta}}(k)/n^{\hat{\theta}}) \leq 2\sqrt{R_n} (1/n^\theta + 1/n^{\hat{\theta}}) \rightarrow 0 \end{aligned}$$

a.s. uniformly in $t \in [0, 1]$.

Note that one can change $t^{\hat{\theta}} - t^\theta$ by $(\hat{\theta} - \theta)t^\theta \log t$. Indeed,

$$\begin{aligned} t^{\hat{\theta}} - t^\theta &= t^\theta (e^{(\hat{\theta}-\theta)\log t} - 1) \\ &= (\hat{\theta} - \theta)t^\theta \log t + t^\theta \sum_{k \geq 2} \frac{((\hat{\theta} - \theta)\log t)^k}{k!} \\ &= (\hat{\theta} - \theta)t^\theta \log t + t^\theta (\hat{\theta} - \theta)^2 (1 + o(1)) \sum_{k \geq 2} \frac{\log^k t}{k!} \\ &= (\hat{\theta} - \theta)t^\theta \log t \left(1 + (\hat{\theta} - \theta)(1 + o(1)) \frac{e^{\log t} - 1 - \log t}{\log t} \right) \\ &= (\hat{\theta} - \theta)t^\theta \log t (1 + o(1)) \end{aligned}$$

a.s. uniformly in $t \in [0, 1]$. Hence from Theorems 3 and 4, we have a joint weak convergence of

$$(\overset{\circ}{Z}_n, \overset{\circ}{Z}'_n, \sqrt{R_n}(\hat{\theta} - \theta))$$

to

$$\left(\overset{\circ}{Z}_\theta, \overset{\circ}{Z}'_\theta, \frac{1}{2} \int_0^1 u^{-\theta} (Z_\theta(u) + Z'_\theta(u)) dA(u) \right).$$

So, $(\widehat{Z}_n, \widehat{Z}'_n)$ converges weakly to $(\widehat{Z}_\theta, \widehat{Z}'_\theta)$. □

Corollary 2. *Assume the conditions of Theorem 2 to hold. Let $\widehat{W}_n^2 = \int_0^1 (\widehat{Z}_n(t))^2 + (\widehat{Z}'_n(t))^2 dt$. Then \widehat{W}_n^2 converges weakly to $\widehat{W}_\theta^2 = \int_0^1 (\widehat{Z}_\theta(t))^2 + (\widehat{Z}'_\theta(t))^2 dt$.*

Similarly to (14), \widehat{W}_n^2 has the following representation

$$\begin{aligned} \widehat{W}_n^2 &= \frac{1}{3n} \sum_{k=1}^{n-1} \widehat{Z}_n \left(\frac{k}{n} \right) \left(2\widehat{Z}_n \left(\frac{k}{n} \right) + \widehat{Z}_n \left(\frac{k+1}{n} \right) \right) \\ &\quad + \frac{1}{3n} \sum_{k=1}^{n-1} \widehat{Z}'_n \left(\frac{k}{n} \right) \left(2\widehat{Z}'_n \left(\frac{k}{n} \right) + \widehat{Z}'_n \left(\frac{k+1}{n} \right) \right) \end{aligned}$$

The p-value of the goodness-of fit test is $1 - \widehat{F}_\theta(\widehat{W}_{n,obs}^2)$. Here \widehat{F}_θ is the cumulative distribution function of \widehat{W}_θ^2 , and $\widehat{W}_{n,obs}^2$ is the observed value of \widehat{W}_n^2 . Further, the function \widehat{F}_θ can be found using the approach described in Section 3 with λ_k replaced by the eigenvalues $\widehat{\lambda}_k$ of the kernel in the Smirnov formula.

6. PROOF OF THEOREM 1

We denote by $\mathbb{X}_i(n)$ the number of balls in urn i . Let $\Pi = \{\Pi(t), t \geq 0\}$ be a Poisson process with parameter 1. The Poissonized version of Karlin model assumes the total number $\Pi(n)$ of balls. According to the well-known thinning property of Poisson flows, stochastic processes $\{\mathbb{X}_i(\Pi(t)) \stackrel{\text{def}}{=} \Pi_i(t), t \geq 0\}$ are Poisson processes (with intensities p_i) which are mutually independent for different i 's. The definition implies that for any fixed $n \geq 1$, $\tau, t \in [0, 1]$

$$R_{\Pi(tn)} = \sum_{k=1}^{\infty} \mathbf{I}(\Pi_k(tn) > 0) = \sum_{k=1}^{\infty} \mathbf{I}_k(tn),$$

$$R'_{\Pi(\tau n)} = \sum_{k=1}^{\infty} \mathbf{I}(\Pi_k(n) - \Pi_k((1-\tau)n) > 0) = \sum_{k=1}^{\infty} \mathbf{I}'_k(\tau n).$$

Step 1 (covariances) Let $\tau, t \in [0, 1]$

$$\begin{aligned} \text{cov} \left(R_{\Pi(tn)}, R'_{\Pi(\tau n)} \right) &= \sum_{k=1}^{\infty} \text{cov}(\mathbf{I}_k(tn), \mathbf{I}'_k(\tau n)) \\ &= \sum_{k=1}^{\infty} \left(\mathbf{P}(\Pi_k(tn) > 0, \Pi_k(n) - \Pi_k((1-\tau)n) > 0) - (1 - e^{-p_k tn})(1 - e^{-p_k \tau n}) \right) \end{aligned}$$

Note that if $t + \tau > 1$, then

$$\begin{aligned} \mathbf{P}(\Pi_k(tn) > 0, \Pi_k(n) - \Pi_k((1-\tau)n) > 0) &= \mathbf{P}(\Pi_k(tn) - \Pi_k((1-\tau)n) > 0) \\ &+ \mathbf{P}(\Pi_k(tn) - \Pi_k((1-\tau)n) = 0, \Pi_k((1-\tau)n) > 0, \Pi_k(n) - \Pi_k(tn) > 0) \\ &= 1 - e^{-p_k(t+\tau-1)n} + e^{-p_k(t+\tau-1)n} (1 - e^{-p_k(1-\tau)n}) (1 - e^{-p_k(1-t)n}) \\ &= 1 - e^{-p_k tn} - e^{-p_k \tau n} + e^{-p_k n}. \end{aligned}$$

Hence

$$\begin{aligned} \text{cov} \left(R_{\Pi(tn)}, R'_{\Pi(\tau n)} \right) &= \mathbf{I}(t + \tau > 1) \sum_{k=1}^{\infty} \left(e^{-p_k(t+\tau)n} - e^{-p_k n} \right) \\ &= \mathbf{I}(t + \tau > 1) (\mathbf{E}R_{\Pi((t+\tau)n)} - \mathbf{E}R_{\Pi(n)}). \end{aligned}$$

Since

$$\mathbf{E}R_{\Pi(tn)} / \mathbf{E}R_n \sim t^\theta,$$

then

$$\text{cov} \left(R_{\Pi(tn)}, R'_{\Pi(\tau n)} \right) / \mathbf{E}R_n \sim K'(t, \tau).$$

Step 2 (convergence of finite-dimensional distributions) Analogously to the proof of Theorem 1 in Dutko (1989), we have that, for any fixed $m \geq 1, 0 < t_1 < t_2 < \dots < t_m \leq 1$, the triangle array of $2m$ -dimensional random vectors $\{((\mathbf{I}_k(nt_i) - \mathbf{E}\mathbf{I}_k(nt_i)) / \sqrt{\mathbf{E}R_n}, (\mathbf{I}'_k(nt_i) - \mathbf{E}\mathbf{I}'_k(nt_i)) / \sqrt{\mathbf{E}R_n}, i \leq m), k \leq n\}_{n \geq 1}$ satisfies the Lindeberg condition (see Borovkov (2013), Theorem 8.6.2, p. 215).

Step 3 (relative compactness) Since $R_{\Pi(nt)} \stackrel{d}{=} R'_{\Pi(nt)}$, then the relative compactness follows from Chebunin and Kovalevskii (2016).

Step 4 (approximation of the original process)

It follows from the corresponding step of the proof in Chebunin and Kovalevskii (2016) and the previous step.

Theorem 1 is proved.

Acknowledgements

The work is supported by Mathematical Center in Akademgorodok under agreement No. 075-15-2019-1675 with the Ministry of Science and Higher Education of the Russian Federation. The authors thank the anonymous referee for very helpful comments that improved the quality of the manuscript.

REFERENCES

- [1] R.R. Bahadur, *On the number of distinct values in a large sample from an infinite discrete distribution*, Proceedings of the National Institute of Sciences of India, **26A**, Supp. II (1960), 67–75. MR0137256
- [2] A.D. Barbour, *Univariate approximations in the infinite occupancy scheme*, Alea **6** (2009), 415–433. MR2576025
- [3] A.D. Barbour, A.V. Gnedin, *Small counts in the infinite occupancy scheme*, Electronic Journal of Probability, Vol. 14, Paper no. 13 (2009), 365–384. MR2480545
- [4] A. Ben-Hamou, S. Boucheron, M. I. Ohannessian, *Concentration inequalities in the infinite urn scheme for occupancy counts and the missing mass, with applications*, Bernoulli **23**, Number 1 (2017), 249–287. MR3556773
- [5] A. Chakrabarty, M. Chebunin, A. Kovalevskii, I. Pupyshv, N. Zakrevskaya, Q. Zhou, *A statistical test for correspondence of texts to the Zipf - Mandelbrot law*, Siberian Electronic Mathematical Reports **17** (2020), 1959–1974. DOI 10.33048/semi.2020.17.132
- [6] M.G. Chebunin, *Estimation of parameters of probabilistic models which is based on the number of different elements in a sample*, Sib. Zh. Ind. Mat., **17**:3 (2014), 135–147 (in Russian). MR3364413
- [7] M.G. Chebunin, *Functional central limit theorem in an infinite urn scheme for distributions with superheavy tails*, Sib. Elektron. Mat. Izv. **14** (2017), 1289–1298. MR3744074
- [8] M. Chebunin, A. Kovalevskii, *Functional central limit theorems for certain statistics in an infinite urn scheme*, Statistics and Probability Letters, **119** (2016), 344–348. MR3555307
- [9] M. Chebunin, A. Kovalevskii, *A statistical test for the Zipf's law by deviations from the Heaps' law*, Siberian Electronic Mathematical Reports **16** (2019), 1822–1832. DOI 10.33048/semi.2019.16.129
- [10] M. Chebunin, A. Kovalevskii, *Asymptotically Normal Estimators for Zipf's Law*, Sankhya A **81** (2019), 482–492. DOI 10.1007/s13171-018-0135-9
- [11] M. Chebunin, S. Zuyev, *Functional Central Limit Theorems for Occupancies and Missing Mass Process in Infinite Urn Models*, J. Theor. Probab. (2020). DOI 10.1007/s10959-020-01053-6
- [12] G. Decrouez, M. Grabchak, Q. Paris, *Finite sample properties of the mean occupancy counts and probabilities*, Bernoulli **24** (2018), no. 3, 1910–1941 MR3757518
- [13] O. Durieu, Y. Wang, *From infinite urn schemes to decompositions of self-similar Gaussian processes*, Electron. J. Probab. **21** (2016), paper no. 43, 23 pp. MR3530320
- [14] O. Durieu, G. Samorodnitsky, Y. Wang, *From infinite urn schemes to self-similar stable processes*, Stochastic Processes and their Applications **130**:4 (2020), 2471–2487.
- [15] M. Dutko, *Central limit theorems for infinite urn models*, Ann. Probab. **17** (1989), 1255–1263. MR1009456
- [16] A. Guillou, P. Hall, *A diagnostic for selecting the threshold in extreme value analysis*, Journal of the Royal Statistical Society: Series B **63**:2 (2002), 293–305. MR1841416
- [17] A. Gnedin, B. Hansen, J. Pitman, *Notes on the occupancy problem with infinitely many boxes: general asymptotics and power laws*, Probability Surveys **4** (2007), 146–171. MR2318403

- [18] B. M. Hill, *A Simple General Approach to Inference About the Tail of a Distribution*, Ann. Statist. **3**:5 (1975), 1163–1174. DOI 10.1214/aos/1176343247
- [19] H.-K. Hwang, S. Janson, *Local Limit Theorems for Finite and Infinite Urn Models*, The Annals of Probability **36**, No. 3 (2008), 992–1022. MR1620350
- [20] S. Karlin, *Central Limit Theorems for Certain Infinite Urn Schemes*, Journal of Mathematics and Mechanics, **17**, No. 4 (1967), 373–401. MR0216548
- [21] E. S. Key, *Rare Numbers*, Journal of Theoretical Probability **5**, No. 2 (1992), 375–389. MR1157991
- [22] E. S. Key, *Divergence rates for the number of rare numbers*, Journal of Theoretical Probability **9**, No. 2 (1996), 413–428. MR1385405
- [23] A. Muratov, S. Zuyev, *Bit flipping and time to recover*, J. Appl. Probab. **53** (2016), no. 3, 650–666. MR3570086
- [24] P.T. Nicholls, *Estimation of Zipf parameters*, J. Am. Soc. Inf. Sci., **38** (1987), 443–445.
- [25] M.I. Ohannessian, M.A. Dahleh, *Rare probability estimation under regularly varying heavy tails*, Proceedings of the 25th Annual Conference on Learning Theory, PMLR 23:21.1–21.24 (2012).
- [26] N.S. Zakrevskaya, A.P. Kovalevskii, *One-parameter probabilistic models of text statistics*, Sib. Zh. Ind. Mat., **4**:2 (2001), 142–153 (in Russian). MR1965927
- [27] N. Zakrevskaya, A. Kovalevskii, *An omega-square statistics for analysis of correspondence of small texts to the Zipf–Mandelbrot law*, Applied methods of statistical analysis. Statistical computation and simulation — AMSA’2019, 18–20 September 2019, Novosibirsk: Proceedings of the International Workshop, Novosibirsk: NSTU (2019), 488–494.
- [28] G.K. Zipf, *The Psycho-Biology of Language*, Routledge, London, 1936.

MIKHAIL GEORGIEVICH CHEBUNIN
KARLSRUHE INSTITUTE OF TECHNOLOGY,
INSTITUTE OF STOCHASTICS,
76131, KARLSRUHE, GERMANY;
NOVOSIBIRSK STATE UNIVERSITY,
PIROGOVA STR., 1,
630090, NOVOSIBIRSK, RUSSIA
Email address: chebuninmikhail@gmail.com

ARTYOM PAVLOVICH KOVALEVSKII
SOBOLEV INSTITUTE OF MATHEMATICS,
KOPTYUGA PR., 4,
NOVOSIBIRSK STATE UNIVERSITY,
PIROGOVA STR., 1,
630090, NOVOSIBIRSK, RUSSIA;
NOVOSIBIRSK STATE TECHNICAL UNIVERSITY,
K. MARKSA AVE., 20,
630073, NOVOSIBIRSK, RUSSIA
Email address: artyom.kovalevskii@gmail.com