S@MR

ISSN 1813-3304

# СИБИРСКИЕ ЭЛЕКТРОННЫЕ МАТЕМАТИЧЕСКИЕ ИЗВЕСТИЯ

Siberian Electronic Mathematical Reports http://semr.math.nsc.ru

Том 17, стр. 1959–1974 (2020) DOI 10.33048/semi.2020.17.132 УДК 519.233 MSC 62F03

## A STATISTICAL TEST FOR CORRESPONDENCE OF TEXTS TO THE ZIPF-MANDELBROT LAW

A. CHAKRABARTY, M.G. CHEBUNIN, A.P. KOVALEVSKII, I.M. PUPYSHEV, N.S. ZAKREVSKAYA, Q. ZHOU

ABSTRACT. We analyse correspondence of texts to a simple probabilistic model. The model assumes that the words are selected independently from an infinite dictionary, and the probability distribution of words corresponds to the Zipf-Mandelbrot law. We count the numbers of different words in the text sequentially and get the process of the numbers of different words. Then we estimate the Zipf-Mandelbrot law's parameters using the same sequence and construct an estimate of the expectation of the number of different words in the text. After that we subtract the corresponding values of the estimate from the sequence and normalize along the coordinate axes, obtaining a random process on a segment from 0 to 1. We prove that this process (the empirical text bridge) converges weakly in the uniform metric on C(0,1) to a centered Gaussian process with continuous a.s. paths. We develop and implement an algorithm for calculating the probability distribution of the integral of the square of this process. We present several examples of application of the algorithm for analysis of the homogeneity of texts in English, French, Russian, and Chinese.

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

Chakrabarty, A., Chebunin, M.G., Kovalevskii, A.P., Pupyshev, I.M., Zakrevskaya, N.S., Zhou, Q., A statistical test for correspondence of texts to the Zipf—Mandelbrot law.

 $<sup>\</sup>textcircled{O}$ 2020 Chakrabarty A., Chebunin M.G., Kovalevskii A.P., Pupyshev I.M., Zakrevskaya N.S., Zhou Q.

The reported study was funded by RFBR and NSFC according to the research project No. 19-51-53010.

Received September, 28, 2020, published November, 27, 2020.

#### A. CHAKRABARTY ET AL.

#### 1. INTRODUCTION

Our analysis is based on the fact that a text in any natural language can be divided into words. The source material for our analysis is a text with separated words and excluded punctuation. In addition, all capital letters (if any) are replaced by lowercase.

We test the hypothesis H that a text matches a simple probabilistic model. The model satisfies the following three assumptions:

1) the dictionary contains countably many words that are enumerated  $i = 1, 2, \ldots$ ;

2) words are sampled from the dictionary in the i.i.d. fashion according to discrete distribution  $F = \{p_i\}_{i>1}$  where  $p_i$  is the probability of word i;

3) the distribution F has a power tail:

(1) 
$$p_i = c(i+q)^{-\theta^{-1}}, i \ge 1, 0 < \theta < 1, q > -1,$$

where q and  $\theta$  are the distribution parameters and c is the normalising constant.

These assumptions have a long history. Power decay of probabilities together with unboundedness of the dictionary were proposed by Zipf (1936). Mandelbrot (1965) noted that shift q is required for the model to better match real texts. Modern large-scale studies (see Petersen et al., 2012) show that the process of the emergence of new words never stops, but very long sequences of texts show a slightly lower frequency of new words than it is predicted by the formula (1). So we apply our test for examples of not very long texts only.

The second assumption of the independence of the choice of consecutive words is obviously false for any meaningful text. We can easily reject it statistically. If we calculate the relative frequency of the word 'and ' in an English text, and then the relative frequency of the sequence 'and and ' (that is, two and in succession), then, according to H, the second should be approximately equal to the square of the first. In practice, the second is much smaller; often, it is simply zero.

But we are not ready to abandon the assumption of independence. This assumption is the base of our analysis. Therefore, we choose a characteristic that changes little when rearranging neighboring words. This characteristic is the number of different words of the text.

Let  $R_n$  be the number of different words in the text of n words. Under the assumption of independence, Bahadur proved the law of large numbers  $R_n/\mathbf{E}R_n \rightarrow 1$  in probability, and Karlin proved strong law of large numbers  $R_n/\mathbf{E}R_n \rightarrow 1$  a.s. and the central limit theorem:  $(R_n - \mathbf{E}R_n)/\sqrt{\mathbf{Var}R_n}$  converges weakly to the standard normal law.

Thus, we consider a text as a random sequence, we construct the sequence  $R_1, \ldots, R_n$  and study it using the methods of the theory of random processes: we invent parameter estimates and construct the empirical text bridge, that is, a random process built on the parameter estimates and the sequence of numbers of different words. We find the limit Gaussian process in the sense of weak convergence in C(0, 1). Then we calculate the distribution function G of the integral of the squared limit process using the Smirnov (1937) formula.

Eigenvalues of the covariance function that are necessary for applying the Smirnov formula are calculated approximately as the eigenvalues of the matrix  $Q = (q_{ij})_{i,j=1}^{L}$ 

composed of coefficients

(2) 
$$q_{ij} = \int_0^1 \int_0^1 \widehat{K}(s,t) \sin \pi i s \sin \pi j t \, ds dt.$$

We calculate  $q_{ij}$  by a fast algorithm that reduces double integrals to definite integrals.

Asymptotics of  $R_n$  and similar statistics (in particular, the number of unique words, that is, words with exactly one occurrence in the text) have been studied by a number of authors. The Gaussian approximation under assumptions (1) was studied by Karlin, and beyond these assumptions by Dutko (1989), Gnedin, Hansen and Pitman (2007), Hwang and Janson (2008), Barbour and Gnedin (2009). Barbour (2009) proposed translated Poisson approximation for the number of unique words. New papers by Ben-Hamou, Boucheron and Ohannessian (2017) and Decrouez, Grabchak and Paris (2018) proved new general facts about these statistics.

The main result on which our study is based is the functional central limit theorem for the sequence  $R_1, \ldots, R_n$ . It was proven by Chebunin and Kovalevskii (2016) in preparation for this study. Note that the Gaussian process which is the limit for this sequence can be found also in Durieu and Wang (2016) as a limit for another prelimit process. Its generalization is in Durieu, Samorodnitsky and Wang (2019).

Zipf parameter estimates were proposed by Nicholls (1987), Chebunin and Kovalevskii (2019b), but we need a special estimate for which we can calculate joint limit distribution of it and of the sequence  $R_1, \ldots, R_n$ . This estimate is proposed in Chebunin and Kovalevskii (2019a). Zakrevskaya and Kovalevskii (2019) used the estimate in analysis of Shakespeare's sonnets.

Bahadur proved that under H the mathematical expectation of the number of different words grows according to an asymptotically power law. This fact is known to experts in natural language processing as Herdan's law (Herdan, 1960) or Heaps' law (Heaps, 1978). Van Leijenhorst and van der Weide (2005), Eliazar (2011) analyzed the relationship between the Zipf's law and Heaps' law based on probabilistic models that were different from H.

Gerlach and Altmann (2013) noted specifically that there is no mathematically correct statistical test for correspondence of a text to the Zipf's law. We proposed such a test in Chebunin and Kovalevskii (2019a). We are developing the algorithm and applying it to texts in different languages in the present paper.

We propose an estimate for the Mandelbrot parameter q, prove its consistence, then we construct an approximation of the process of the number of different words using this estimate, and prove the weak convergence of the normalized difference between the process and the approximation to a centered Gaussian process. We calculate the covariance function of the limiting process and the distribution of the integral of the square of this process. We use these results for an algorithm for calculating the p-value of the hypothesis of text homogeneity and apply the algorithm to the analysis of the homogeneity of texts in different languages.

The rest of the paper is organised as follows. The necessary theoretical results are in Section 2, the algorithm is in Section 3, examples of text analysis are in Section 4, and a discussion is in Section 5.

## 2. Theoretical results

The Hurvitz zeta function is

$$\zeta(a, x) = \sum_{i=0}^{\infty} (i+x)^{-a}.$$

Note that constant c in (1) is

$$c = c(\theta, q) = (\zeta(\theta^{-1}, q+1))^{-1}.$$

Let n be a number of words in a text.

Let  $R_k$  be the number of different words among first k words of the text.

Let  $R_0 = 0$ . We have  $R_1 = 1$ ,  $R_0 < R_1 \le R_2 \le ... \le R_n$ . Bahadur (1960) proved that

(3) 
$$\mathbf{E}R_j \sim c^{\theta} \Gamma(1-\theta) j^{\theta}$$

where  $\Gamma(x) = \int_0^\infty y^{x-1} e^{-y} dy$  is the Euler gamma function. Bahadur also proved convergence in probability  $R_i / \mathbf{E} R_i \xrightarrow{p} 1$ .

Karlin (1967) proved that  $R_i / \mathbf{E} R_i \xrightarrow{a.s.} 1$ , which is equivalent to

(4) 
$$R_j \sim C_1 j^\theta \text{ a.s.},$$

thanks to (3). Here  $C_1 = C_1(\theta, q) = c^{\theta} \Gamma(1-\theta)$ .

Chebunin and Kovalevskii (2016) proved the Functional Central Limit Theorem, that is, the weak convergence of the process  $\{(R_{[nt]} - \mathbf{E}R_{[nt]})/\sqrt{\mathbf{E}R_n}, 0 \le t \le 1\}$  to a centered Gaussian process  $Z_{\theta}$  with continuous a.s. sample paths and covariance function

$$K(s,t) = (s+t)^{\theta} - \max(s^{\theta}, t^{\theta})$$

Denote  $\log^+ x = \max(\log x, 0)$ .

We propose the following estimator for parameter  $\theta$  (Chebunin and Kovalevskii, 2019a)

$$\widehat{\theta} = \int_0^1 \log^+ R_{[nt]} \, dA(t)$$

with a function  $A(\cdot)$  that is the sum of a step function and a piecewise continuously differentiable function on [0, 1] and

(5) 
$$\int_0^1 \log t \, dA(t) = 1, \quad A(t) = 0, \ t \in [0, \ \delta], \text{ for some } \delta \in (0, \ 1), \ A(1) = 0.$$

The next theorem follows from Theorems 2.1, 2.2 by Chebunin and Kovalevskii (2019a).

**Theorem 2.1** If (5) holds then the estimator  $\hat{\theta}$  is strongly consistent, and

$$\sqrt{\mathbf{E}R_n}(\widehat{\theta}-\theta) - \int_0^1 t^{-\theta} Z_n(t) \, dA(t) \to_p 0.$$

From Theorem 2.1, it follows that  $\hat{\theta}$  converges to  $\theta$  with rate  $(\mathbf{E}R_n)^{-1/2}$ , and normal random variable  $\int_0^1 t^{-\theta} Z_{\theta}(t) \, dA(t)$  has variance  $\int_0^1 \int_0^1 (st)^{-\theta} K(s,t) \, dA(s) \, dA(t)$ .

Example 2.1 Take

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

Then

(6)

$$\widehat{\theta} = \log_2(R_n/R_{[n/2]}), \quad n \ge 2.$$

Note that, in this example, for any function g on [0, 1],

$$\int_0^1 g(t) \, dA(t) = \frac{g(1) - g(1/2)}{\log 2}$$

Let us introduce the process  $\widehat{Z}_n$ :

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

 $0 \le k \le n$ . Let for  $0 \le t \le 1/n$  and  $0 \le k \le n-1$ 

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

Theorem 2.2 (Theorem 4.1 by Chebunin and Kovalevskii, 2019a) If (5) holds then  $\widehat{Z}_n$  converges weakly to  $\widehat{Z}_{\theta}$  as  $n \to \infty$ , where

$$\widehat{Z}_{\theta}(t) = Z_{\theta}^{0}(t) - t^{\theta} \log t \int_{0}^{1} u^{-\theta} Z_{\theta}(u) \, dA(u),$$

$$\begin{split} Z^0_\theta(t) &= Z_\theta(t) - t^\theta Z_\theta(1), \ 0 \leq t \leq 1. \\ \text{The correlation function of } Z^0_\theta \text{ is given by} \end{split}$$

$$K^{0}(s,t) = \mathbf{E}Z^{0}_{\theta}(s)Z^{0}_{\theta}(t) = K(s,t) - s^{\theta}K(1,t) - t^{\theta}K(s,1) + s^{\theta}t^{\theta}K(1,1).$$

**Corollary 2.1** Let  $\widehat{W}_n^2 = \int_0^1 \left(\widehat{Z}_n(t)\right)^2 dt$  and (5) holds. Then  $\widehat{W}_n^2$  converges weakly to  $\widehat{W}_{\theta}^2 = \int_{0}^{1} \left(\widehat{Z}_{\theta}(t)\right)^2 dt.$ 

 $\widehat{W}_n^2$  has the following representation

$$\widehat{W}_n^2 = \frac{1}{3n} \sum_{k=1}^{n-1} \widehat{Z}_n\left(\frac{k}{n}\right) \left(2\widehat{Z}_n\left(\frac{k}{n}\right) + \widehat{Z}_n\left(\frac{k+1}{n}\right)\right)$$

The p-value of the goodness-of fit test is  $1 - \widehat{F}_{\theta}(\widehat{W}_{n,obs}^2)$ . Here  $\widehat{F}_{\theta}$  is the cumulative distribution function of  $\widehat{W}_{\theta}^2$ , and  $\widehat{W}_{n,obs}^2$  is the observed value of  $\widehat{W}_n^2$ . One can estimate  $F_{\theta}$  by simulations or find it explicitly using the Smirnov's

formula (Smirnov, 1937): if  $\widehat{W}_{\theta}^2 = \sum_{k=1}^{\infty} \frac{\eta_k^2}{\lambda_k}$ ,  $\eta_1, \eta_2, \ldots$  are independent and have standard normal distribution,  $0 < \lambda_1 < \lambda_2 < \ldots$ , then

(7) 
$$F_{\theta}(x) = 1 + \frac{1}{\pi} \sum_{k=1}^{\infty} (-1)^k \int_{\lambda_{2k-1}}^{\lambda_{2k}} \frac{e^{-\lambda x/2}}{\sqrt{-D(\lambda)}} \cdot \frac{d\lambda}{\lambda}, \ x > 0,$$
$$D(\lambda) = \prod_{k=1}^{\infty} \left(1 - \frac{\lambda}{\lambda_k}\right).$$

The integrals in the RHS of (7) must tend to 0 monotonically as  $k \to \infty$ , and  $\lambda_k^{-1}$  are the eigenvalues of kernel  $\hat{K}(s,t) = \mathbf{E}\hat{Z}_{\theta}(s)\hat{Z}_{\theta}(t)$ , see Smirnov (1937).

We are exploring the two-parameter model in this paper. Therefore, we now develop a new approximation of the process of the number of different words using the estimate of the Mandelbrot parameter q.

Let us introduce the empirical bridge of a text  $\widetilde{Z}_n$  by substituting

(8) 
$$r_k = r_k(\widehat{\theta}, \widehat{q}) = \sum_{i=1}^{\infty} \left( 1 - (1 - \widehat{p}_i)^k \right)$$

instead of  $(k/n)^{\widehat{\theta}}R_n$  in  $\widehat{Z}_n$ , that is,

(9) 
$$\widetilde{Z}_n(k/n) = (R_k - r_k) / \sqrt{R_n},$$

 $\begin{array}{l} 0 \leq k \leq n. \\ \text{Here} \end{array}$ 

(10) 
$$\widehat{p}_i = c(\widehat{\theta}, \widehat{q})(i + \widehat{q})^{-1/\widehat{\theta}}, \quad i \ge 1$$

(11) 
$$\widehat{q} = \min\{\widetilde{q} > -1: r_n(\widehat{\theta}, \widetilde{q}) = R_n\}.$$

**Theorem 2.3** If H is true then  $\hat{q} \to q$  in probability as  $n \to \infty$ , and  $\hat{Z}_n - \tilde{Z}_n \to 0$  in probability uniformly on q in any segment in  $(-1, \infty)$ .

Proof

From Lemma 1 in Gnedin et al. (2007) we have

$$|\mathbf{E}R_n - \mathbf{E}R_{\Pi(n)}| < \frac{2}{n}\mathbf{E}R_{n,2} \le \frac{2}{n} \times \frac{n}{2} = 1.$$

Karlin (1967) proposed representation

$$\mathbf{E}R_{\Pi(t)} = \int_0^\infty \alpha(ty) \frac{1}{y^2} e^{-1/y} dy$$

where

$$\alpha(x) = \max\{j|p_j \ge \frac{1}{x}\}.$$

In our case

$$\alpha(x) = [(cx)^{\theta} - q].$$

We represent

$$\Gamma(k-\gamma) = \int_0^\infty y^{\gamma-k-1} e^{-1/y} dy$$

So

$$(cn)^{\theta} \Gamma(1-\theta) - q - 1 \leq \mathbf{E} R_{\Pi(n)} \leq (cn)^{\theta} \Gamma(1-\theta) - q,$$
$$\mathbf{E} R_n = (cn)^{\theta} \Gamma(1-\theta) - q + T_n, \quad |T_n| < 2,$$

then

$$|r_n - (\widehat{c}n)^{\widehat{\theta}} \Gamma(1 - \widehat{\theta}) + \widehat{q}| < 2 \text{ a.s.}$$

Notice, that for any x > 0

$$\zeta(\theta^{-1}, x) = \sum_{i=0}^{\infty} (i+x)^{-\theta^{-1}} > \sum_{i=0}^{\infty} \int_{i}^{i+1} (y+x)^{-\theta^{-1}} dy = \int_{0}^{\infty} (y+x)^{-\theta^{-1}} dy = c_1 x^{1-\theta^{-1}},$$

where  $c_1 = c_1(\theta) = (\theta^{-1} - 1)^{-1}$ , and

$$\zeta(\theta^{-1}, x) = \sum_{i=0}^{\infty} (i+x)^{-\theta^{-1}} < x^{-\theta^{-1}} + \sum_{i=1}^{\infty} \int_{i-1}^{i} (y+x)^{-\theta^{-1}} dy = x^{-\theta^{-1}} + c_1 x^{1-\theta^{-1}}.$$
 It's not hard to see that

It's not hard to see that

$$c^{\theta} = (\zeta(\theta^{-1}, q+1))^{-\theta} \sim c_1 q^{1-\theta} \text{ as } q \to \infty.$$

From the strong law of large numbers for  $R_n$  and  $\hat{\theta} \to \theta$  a.s., we can choose a constant  $K = K(\theta, q) = 2c^{\theta}\Gamma(1-\theta)$  such that

$$\mathbf{P}(R_n < Kn^{\widehat{\theta}}) \to 1 \text{ as } n \to \infty.$$

Therefore, we can choose  $q^*$  such, that

$$(c(\widehat{\theta}, q^*))^{\widehat{\theta}} \Gamma(1 - \widehat{\theta}) \underset{a.s.}{\sim} (c(\theta, q^*))^{\theta} \Gamma(1 - \theta) > K.$$

Hence

$$\mathbf{P}(\widehat{q}_n < q^*) \to 1 \text{ as } n \to \infty.$$

On the other hand, from (3) we have

$$R_n \underset{a.s.}{\sim} (c(\theta, q)n)^{\theta} \Gamma(1-\theta),$$
$$r_n \underset{a.s.}{\sim} (c(\widehat{\theta}, \widehat{q})n)^{\widehat{\theta}} \Gamma(1-\widehat{\theta}),$$

and  $r_n = R_n$ ,  $\hat{\theta} \to \theta$  a.s., so

$$\frac{(c(\widehat{\theta},\widehat{q})n)^{\widehat{\theta}}}{(c(\theta,q)n)^{\theta}} \xrightarrow{a.s.} 1.$$

From Theorem 2.1,  $n^{\widehat{\theta}-\theta} \to 1$  in probability, so  $c(\widehat{\theta}, \widehat{q}) \to c(\theta, q)$  in probability, and  $c(\theta, q)$  is a continious and strictly monotone function of q for any  $\theta \in (0, 1)$ . Thus  $\widehat{q} \to q$  in probability.

Now we estimate

$$\max_{1 \le k \le n} \frac{|(k/n)^{\theta} R_n - r_k|}{\sqrt{R_n}}.$$

We have

$$|(k/n)^{\hat{\theta}}R_n - r_k| = |(k/n)^{\hat{\theta}}r_n - r_k|$$
  
$$< |(k/n)^{\hat{\theta}}((\widehat{c}n)^{\hat{\theta}}\Gamma(1-\widehat{\theta}) - \widehat{q}) - (\widehat{c}k)^{\hat{\theta}}\Gamma(1-\widehat{\theta}) + \widehat{q}| + 4$$
  
$$= |\widehat{q}(1 - (k/n)^{\hat{\theta}})| + 4 \le \widehat{q} + 4.$$

As  $R_n \to \infty$  a.s. and  $\widehat{q} \to q$  in probability, we have

$$\max_{1 \le k \le n} \frac{|(k/n)^{\theta} R_n - r_k|}{\sqrt{R_n}} \to 0$$

in probability.

The proof is complete.

3. Algorithm

Note that

$$r_{k} = \sum_{i=1}^{\infty} \left( 1 - \sum_{j=0}^{k} C_{k}^{j} (-1)^{j} (\widehat{p}_{i})^{j} \right)$$
$$= \sum_{j=1}^{k} C_{k}^{j} (-1)^{j+1} \sum_{i=1}^{\infty} \widehat{c}^{j} (i+\widehat{q})^{-j/\widehat{\theta}}$$
$$= \sum_{j=1}^{k} C_{k}^{j} (-1)^{j+1} \zeta(j/\widehat{\theta}, 1+\widehat{q}) \left( \zeta(1/\widehat{\theta}, 1+\widehat{q}) \right)^{-j}.$$

So we have the finite formula to calculate  $r_1, \ldots, r_n$ . This formula is not good for large k due to high complexity of calculations of binomial coefficients. So we use an appoximation by substituting an intergal instead of series residual

$$r_k \approx \sum_{i=1}^{M} (1 - (1 - \hat{p}_i)^k) + \int_{M+0.5 + \hat{q}}^{\infty} (1 - \exp(-k\hat{c}y^{-\hat{\alpha}})) dy$$

(12) 
$$=\sum_{i=1}^{M} (1 - (1 - \hat{p}_i)^k) + (k\hat{c})^{\hat{\theta}} \int_{0}^{k\hat{c}N^{-\hat{\alpha}}} z^{-\hat{\theta}} e^{-z} dz - N(1 - \exp(-k\hat{c}N^{-\hat{\alpha}})),$$

 $\widehat{\alpha}=1/\widehat{\theta},\,N=M+0.5+\widehat{q}.$ 

We calculate the integral using incomplete Gamma function. We find  $\hat{q}$  by dichotomy method for  $r_n = R_n$ .

Elementary calculations give

$$\begin{split} \widehat{K}(s,t) &= \mathbf{E}\widehat{Z}_{\theta}(s)\widehat{Z}_{\theta}(t) \\ &= K^{0}(s,t) - t^{\theta}\log t \frac{K(s,1) - 2^{\theta}K(s,1/2)}{\log 2} - s^{\theta}\log s \frac{K(t,1) - 2^{\theta}K(t,1/2)}{\log 2} \\ &+ s^{\theta}t^{\theta}(\log s + \log t) \frac{K(1,1) - 2^{\theta}K(1,1/2)}{\log 2} \\ &+ s^{\theta}t^{\theta}\log s\log t \frac{K(1,1) - 2^{\theta+1}K(1,1/2) + 2^{2\theta}K(1/2,1/2)}{\log^{2} 2}. \end{split}$$

Now we represent the kernel  $\widehat{K}$  by the matrix  $Q = (q_{ij})_{i,j=1}^{L}$  by calculation of

(13) 
$$q_{ij} = \int_0^1 \int_0^1 \widehat{K}(s,t) \sin \pi i s \sin \pi j t \, ds dt.$$

We know  $K(s,t) = (s+t)^{\theta} - \max(s^{\theta}, t^{\theta}),$ 

$$\begin{split} \widehat{K}(s,t) &= K(s,t) - s^{\theta}K(1,t) - t^{\theta}K(s,1) + s^{\theta}t^{\theta}K(1,1) \\ -t^{\theta}\log t \frac{K(s,1) - 2^{\theta}K(s,1/2)}{\log 2} - s^{\theta}\log s \frac{K(t,1) - 2^{\theta}K(t,1/2)}{\log 2} \\ &+ s^{\theta}t^{\theta}(\log s + \log t) \frac{K(1,1) - 2^{\theta}K(1,1/2)}{\log 2} \\ + s^{\theta}t^{\theta}\log s\log t \frac{K(1,1) - 2^{\theta+1}K(1,1/2) + 2^{2\theta}K(1/2,1/2)}{\log^2 2}. \end{split}$$

Let us denote

$$J_{ij} = \int_0^1 \int_0^1 (s+t)^\theta \sin \pi i s \sin \pi j t \, ds dt,$$
$$a(\theta, k, x) = \int_0^x t^\theta \sin \pi k t \, dt, \quad b(\theta, k, x) = \int_0^x t^\theta \cos \pi k t \, dt,$$
$$A_i = \int_0^1 t^\theta \sin \pi i t \, dt = a(\theta, i, 1), \quad B_i = \int_1^2 t^\theta \sin \pi i t \, dt = a(\theta, i, 2) - a(\theta, i, 1),$$
$$C_i = \int_0^1 t^{\theta+1} \cos \pi i t \, dt = b(\theta+1, i, 1),$$

$$D_{i} = \int_{1}^{2} t^{\theta}(2-t) \cos \pi it \, dt = 2b(\theta, i, 2) - 2b(\theta, i, 1) - b(\theta + 1, i, 2) + b(\theta + 1, i, 1),$$
$$E_{ij} = \int_{0}^{1} t^{\theta} \sin \pi it \cos \pi jt \, dt, \quad F_{i} = \int_{0}^{1} K(t, 1) \sin \pi it \, dt,$$
$$G_{i} = \int_{0}^{1} t^{\theta} \log t \sin \pi it \, dt, \quad H_{i} = \int_{0}^{1} K(t, 1/2) \sin \pi it \, dt.$$

Then we have

$$\begin{split} q_{ij} &= J_{ij} - \frac{1}{\pi j} (A_i - E_{ij}) - \frac{1}{\pi i} (A_j - E_{ji}) \\ &- A_i F_j - A_j F_i + K(1, 1) A_i A_j \\ &- G_j \frac{F_i - 2^{\theta} H_i}{\log 2} - G_i \frac{F_j - 2^{\theta} H_j}{\log 2} \\ &+ (A_i G_j + A_j G_i) \frac{K(1, 1) - 2^{\theta} K(1, 1/2)}{\log 2} \\ &+ G_i G_j \frac{K(1, 1) - 2^{\theta + 1} K(1, 1/2) + 2^{2\theta} K(1/2, 1/2)}{\log^2 2}. \end{split}$$

We calculate  $J_{ij}$  substituting t = u - s,  $s \le u \le s + 1$ . So

$$J_{ij} = \int_0^1 u^\theta \, du \int_0^\infty \sin \pi i s \sin \pi j (u-s) \, ds + \int_1^\infty u^\theta \, du \int_{u-1}^1 \sin \pi i s \sin \pi j (u-s) \, ds.$$
  
If  $i \neq j$  then we have

$$\begin{aligned} J_{ij} &= \frac{1}{2} \int_0^1 u^\theta \, du \int_0^u (\cos \pi (is - ju + js) - \cos \pi (is + ju - js)) \, ds \\ &+ \frac{1}{2} \int_1^2 u^\theta \, du \int_{u-1}^1 (\cos \pi (is - ju + js) - \cos \pi (is + ju - js)) \, ds \\ &= \frac{1}{2} \int_0^1 u^\theta \left( \frac{\sin \pi iu + \sin \pi ju}{\pi (i + j)} - \frac{\sin \pi iu - \sin \pi ju}{\pi (i - j)} \right) \, du \\ &+ \frac{1}{2} \int_1^2 u^\theta \left( \frac{\sin (\pi (i + j) - \pi ju) + \sin (\pi (i + j) - \pi iu))}{\pi (i + j)} - \frac{\sin (\pi (i - j) + \pi ju) + \sin (\pi (i - j) - \pi iu))}{\pi (i - j)} \right) \, du \\ &= \frac{iA_j - jA_i - (-1)^{i+j} (iB_j - jB_i)}{\pi (i^2 - j^2)}. \end{aligned}$$

If i = j then we have

$$J_{ii} = \frac{1}{2} \int_0^1 u^\theta \, du \int_0^u (\cos \pi (2is - iu) - \cos \pi iu) \, ds$$
$$+ \frac{1}{2} \int_1^2 u^\theta \, du \int_{u-1}^1 (\cos \pi (2is - iu) - \cos \pi iu) \, ds$$
$$= \frac{1}{2} \int_0^1 \left( u^\theta \frac{\sin \pi iu}{\pi i} - u^{\theta+1} \cos \pi iu \right) \, du$$
$$+ \frac{1}{2} \int_1^2 u^\theta \frac{\sin(2\pi i - \pi iu) - \sin(2\pi iu - 2\pi i - \pi iu)}{2\pi i} \, du - \frac{1}{2} \int_1^2 u^\theta (2 - u) \cos \pi iu \, du$$

$$=\frac{A_i-B_i}{2\pi i}-\frac{C_i+D_i}{2}.$$

We know

$$\begin{split} a(\theta,k,x) &= \pi k x^{\theta+2} {}_{1} F_{2} \left( \theta/2 + 1, 3/2, \theta/2 + 2, -k^{2} \pi^{2} x^{2}/4 \right) / (\theta+2), \\ b(\theta,k,x) &= x^{\theta+1} {}_{1} F_{2} \left( \theta/2 + 1/2, 1/2, \theta/2 + 3/2, -k^{2} \pi^{2} x^{2}/4 \right) / (\theta+1), \\ E_{ij} &= \frac{1}{2} A_{i+j} + \frac{1}{2} A_{i-j}, \\ F_{i} &= (-1)^{i} B_{i} + \frac{(-1)^{i} - 1}{\pi i}, \\ G_{i} &= \int_{0}^{1} t^{\theta} \log t \sum_{k=0}^{\infty} \frac{(-1)^{k} (\pi i t)^{2k+1}}{(2k+1)!} dt = -\sum_{k=0}^{\infty} \frac{(-1)^{k} (\pi i j)^{2k+1}}{(2k+1)!(2k+\theta+2)^{2}} \\ &= -\pi i {}_{2} F_{3} \left( \theta/2 + 1, \theta/2 + 1, 3/2, \theta/2 + 2, \theta/2 + 2, -i^{2} \pi^{2}/4 \right) / (\theta+2)^{2}, \\ H_{i} &= \int_{1/2}^{3/2} t^{\theta} \sin(\pi i t - \pi i/2) dt + (\cos(\pi i/2) - 1)2^{-\theta} / (\pi i) - \int_{1/2}^{1} t^{\theta} \sin \pi i t dt \\ &= \pi i \cos(\pi i/2) \left( (3/2)^{\theta+2} {}_{1} F_{2} \left( \theta/2 + 1, 3/2, \theta/2 + 2, -9i^{2} \pi^{2}/16 \right) \right) \\ -2^{-\theta-2} {}_{1} F_{2} \left( \theta/2 + 1, 3/2, \theta/2 + 2, -i^{2} \pi^{2}/16 \right) \right) / (\theta+2) \\ &- \sin(\pi i/2) \left( (3/2)^{\theta+1} {}_{1} F_{2} \left( \theta/2 + 1/2, 1/2, \theta/2 + 3/2, -9i^{2} \pi^{2}/16 \right) \right) \\ -2^{-\theta-1} {}_{1} F_{2} \left( \theta/2 + 1/2, 1/2, \theta/2 + 3/2, -i^{2} \pi^{2}/16 \right) \right) / (\theta+1) \end{split}$$

+ $(\cos(\pi i/2)-1)2^{-\theta}/(\pi i)-A_i+\pi i2^{-\theta-2}{}_1F_2(\theta/2+1,3/2,\theta/2+2,-i^2\pi^2/16)/(\theta+2).$ {}\_1F\_2 and {}\_2F\_3 are generalized hypergeometric functions, hyp1f2 and hyp2f3 in

Python.

We calculate eigenvalues  $\lambda_i$ ,  $1 \leq i \leq L$ , of the matrix Q. Then we use the Smirnov formula to calculate the p-value.

Let 
$$W_n^2 = \int_0^1 \left(Z_n^0(t)\right)^2 dt$$
. It is equal to

(14) 
$$W_n^2 = \frac{1}{3n} \sum_{k=1}^{n-1} Z_n^0\left(\frac{k}{n}\right) \left(2Z_n^0\left(\frac{k}{n}\right) + Z_n^0\left(\frac{k+1}{n}\right)\right).$$

Then  $W_n^2$  converges weakly to  $W_{\theta}^2 = \int_0^1 \left( Z_{\theta}^0(t) \right)^2 dt.$ 

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

So, we developed the following text homogeneity analysis algorithm.

## Algorithm

- 1. Remove all punctuation marks.
- 2. Calculate the process of numbers of different words.
- 3. Calculate the estimate  $\hat{\theta}$  of the parameter  $\theta$  by (6).
- 4. Calculate the estimate  $\hat{q}$  of the parameter q by (11) with (12) for k = n and (10).

- 5. Calculate the normalised difference of the process of numbers of different words and its estimate (the empirical bridge  $\widetilde{Z}_n$ ) by (9) with (12).
- 6. Calculate the integral  $W_{n,obs}^2 = W_n^2$  of the square of the empirical bridge by (14).
- 7. Represent the kernel  $\hat{K}$  by the matrix  $Q = (q_{ij})_{i,j=1}^{L}$  by (13). The integrals are calculated by methods of Section 3.
- 8. Find eigenvalues  $\lambda_i$  of the matrix.
- 9. Calculate the p-value  $1 F_{\theta}(W_{n,obs}^2)$  by (7).

L and M are the parameters of the algorithm, L is the matrix dimension, M is the constant in (12).

#### 4. Examples of text analysis

We use the algorithm with L = 100, M = 1000.

We find  $\hat{q}$  by dichotomy method for  $r_n = R_n$  with 20 iterations on segment [-0.9, 40].

Here are given an example of French poetry, *Les Regrets* by Joachim du Bellay (1558), sonnet 1 (Pic. 1) and Shakespeare's sonnet 1 (Pic. 2).

These are examples of homogeneous texts.



РИС. 1. Les Regrets by Joachim du Bellay (1558), sonnet 1

This approach can work with texts on any language.

See the first stanza of the first chapter of *Eugene Onegin* by Pushkin (Pic. 3) for example of Russian homogeneous text.

We use hieroglyphs instead of words when analyzing Chinese texts. We substitute hieroglyphs with their HTML codes. We analyse *Danqing Painting* by Du Fu (Pic. 4). This text is homogeneous, too.

An example of a nonhomogeneous text is Shakespeare's sonnet 1 with 3 repeated lines (Pic. 5).

A. CHAKRABARTY ET AL.



Рис. 2. Shakespeare's sonnet 1



PUC. 3. The first stanza of the first chapter of Eugene Onegin by Pushkin

Another example of nonhomogeneity is a text from different languages. The example is Du Bellay's sonnet 1 +Shakespeare's sonnet 1 (Pic. 6).

To find nonhomogeneity for a text in one language we need more longer texts. Here the first three stanzas of *Childe Harold's Pilgrimage* by Byron and the first three Shakespeare's sonnets (Pic. 7).



РИС. 4. Danging Painting by Du Fu



Рис. 5. Shakespeare's sonnet 1 with 3 repeated lines

## 5. DISCUSSION

There are some open questions in the application of this approach. The first open question is the Poisson approximation. If the number of identical words is small, then the Gaussian approximation is inaccurate. Barbour (2009) proposed an approximation by a translated Poisson distribution. We need a functional version of his theorem.



Рис. 6. Du Bellay's sonnet 1 + Shakespeare's sonnet 1



РИС. 7. The first three stanzas of *Childe Harold's Pilgrimage* by Byron and the first three Shakespeare's sonnets

Another open question is the implementation of the Simon model. Let  $R_{n,1}$  be the number of words that occur exactly once. Under the assumptions made, there should be convergence  $R_{n,1}/R_n \to \theta$  a.s. (Karlin, 1967). But real texts behave differently. Typically, the number of words that occur once is significantly less than  $\hat{\theta}R_n$ .

Simon (1955) proposed the next stochastic model: the (n + 1)-th word in the text is new with probability p; it coincides with each of the previous words with probability (1-p)/n. The drawback of Simon's model is that the number of different words grows linearly. We need some kind of hybrid of an infinite urn model and Simon's model.

#### Acknowledgements

The research was supported by RFBR grant 19-51-53010. The authors would like to thank Sergey Foss and an anonimous referee for helpful and constructive comments and suggestions.

#### References

- R.R. Bahadur, On the number of distinct values in a large sample from an infinite discrete distribution, Proc. Natl Inst. Sci. India, 26A, Supp. II (1960), 67-75. Zbl 0151.23803
- [2] A.D. Barbour, Univariate approximations in the infinite occupancy scheme, Alea, 6 (2009), 415-433. MR2576025
- [3] A.D. Barbour, A.V. Gnedin, 2009. Small counts in the infinite occupancy scheme, Electron. J. Probab., 14 (2009), 365-384. Zbl 1189.60048
- [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, 23:1 (2017), 249-287. Zbl 1366.60016
- M. Chebunin, A. Kovalevskii, 2016. Functional central limit theorems for certain statistics in an infinite urn scheme, Stat. Probab. Lett., 119 (2016), 344-348. Zbl 1398.60051
- [6] M. Chebunin, A. Kovalevskii, A statistical test for the Zipf's law by deviations from the Heaps' law, Sib. Electron. Mat. Izv., 16 (2019), 1822–1832. Zbl 1433.62060
- M. Chebunin, A. Kovalevskii, Asymptotically normal estimators for Zipf's law, Sankhya, Ser. A, 81:2 (2019), 482-492. Zbl 1437.62097
- [8] G. Decrouez, M. Grabchak, Q. Paris, Finite sample properties of the mean occupancy counts and probabilities, Bernoulli, 24:3 (2018), 1910–1941. Zbl 1429.60016
- [9] O. Durieu, Y. Wang, From infinite urn schemes to decompositions of self-similar Gaussian processes, Electron. J. Probab., 21 (2016), Paper No. 43. Zbl 1346.60039
- [10] O. Durieu, G. Samorodnitsky, Y. Wang, From infinite urn schemes to self-similar stable processes, Stochastic Processes Appl., 130:4 (2020), 2471-2487. Zbl 1434.60105
- [11] M. Dutko, Central limit theorems for infinite urn models, Ann. Probab., 17:3 (1989), 1255– 1263. Zbl 0685.60023
- [12] I. Eliazar, The Growth Statistics of Zipfian Ensembles: Beyond Heaps' Law, Physica (Amsterdam), 390 (2011), 3189.
- [13] M. Gerlach, E.G. Altmann, Stochastic Model for the Vocabulary Growth in Natural Languages, Physical Review X 3 (2013) 021006.
- [14] A. Gnedin, B. Hansen, J. Pitman, Notes on the occupancy problem with infinitely many boxes: general asymptotics and power laws, Probab. Surv., 4 (2007), 146–171. Zbl 1189.60050
- [15] H.S. Heaps, Information Retrieval: Computational and Theoretical Aspects, Academic Press, New York etc., 1978. Zbl 0471.68075
- [16] G. Herdan, Type-token mathematics. A textbook of mathematical linguistics, Mouton and Co., 's-Gravenhage, 1960. Zbl 0163.40904
- [17] H.-K. Hwang, S. Janson, Local limit theorems for finite and infinite urn models, Ann. Probab., 36:3 (2008), 992–1022. Zbl 1138.60027
- [18] D.C. van Leijenhorst, Th.P. van der Weide, A Formal Derivation of Heaps' Law, Inf. Sci., 170:2-4 (2005), 263-272. Zbl 1070.60009

#### A. CHAKRABARTY ET AL.

- [19] B. Mandelbrot, Information Theory and Psycholinguistics, In: B.B. Wolman and E. Nagel, Scientific psychology, Basic Books. 1965
- [20] P.T. Nicholls, Estimation of Zipf parameters, J. Am. Soc. Inf. Sci., 38:8 (1987), 443-445.
- [21] A.M. Petersen, J.N. Tenenbaum, S. Havlin, H.E. Stanley, M. Perc, Languages cool as they expand: Allometric scaling and the decreasing need for new words, Scientific Reports 2 (2012), Article No. 943. https://doi.org/10.1038/srep00943
- [22] N.V. Smirnov, On the  $\omega^2$  distribution, Mat. Sb. n. Ser., **2** (1937), 973–993. Zbl 0018.41202
- [23] N. Zakrevskaya, A. Kovalevskii, An omega-square statistics for analysis of correspondence of small texts to the Zipf-Mandelbrot law, In: Applied methods of statistical analysis. Statistical computation and simulation, Proceedings of the International Workshop, NSTU, Novosibirsk, 2019, 488-494.
- [24] G.K. Zipf, The Psycho-Biology of Language, Houghton Mifflin, Boston, 1935.

ANIK CHAKRABARTY Novosibirsk State University, 1, Pirogova str., Novosibirsk, 630090, Russia Email address: a.chakrabarty@g.nsu.ru

Mikhail Georgievich Chebunin Sobolev Institute of Mathematics, 4, Koptyuga ave., Novosibirsk State University, 1, Pirogova str., Novosibirsk, 630090, Russia *Email address*: chebuninmikhail@gmail.com

ARTYOM PAVLOVICH KOVALEVSKII NOVOSIBIRSK STATE TECHNICAL UNIVERSITY, 20, K. MARKSA AVE., NOVOSIBIRSK, 630073, RUSSIA NOVOSIBIRSK STATE UNIVERSITY, 1, PIROGOVA STR., NOVOSIBIRSK, 630090, RUSSIA Email address: artyom.kovalevskii@gmail.com

ILYA MIKHAILOVICH PUPYSHEV Novosibirsk State Technical University, 20, K. Marksa ave., Novosibirsk, 630073, Russia Novosibirsk State University, 1, Pirogova str., Novosibirsk, 630090, Russia Email address: iluxal@ngs.ru

NATALIA STANISLAVOVNA ZAKREVSKAYA Novosibirsk State Technical University, 20, K. Marksa ave., Novosibirsk, 630073, Russia Email address: natali.erlagol@gmail.com

Qianqian Zhou School of Mathematical Sciences, Nankai University, Tianjin, 300071, China Email address: qianqzhou@yeah.net