

Bioinformatics

On the Estimators for some Frequency Distributions Arising in Bioinformatics

Davood Farbod

Department of Mathematics, Quchan University of Advanced Technology, Quchan, Iran

(Presented by Academy Member Mindia Salukvadze)

ABSTRACT. We consider two well-known frequency distributions which can be used for modeling phenomena arising in bioinformatics. These are: the two-parameter Waring frequency distribution and the two-parameter Pareto-like frequency distribution, both introduced by V. A. Kuznetsov. We propose the system of finding the *maximum likelihood estimators* (MLE) for the unknown parameters of such frequency distributions. The proposed MLE are coincided with some moment estimators. Moreover, a method of approximate computation of the MLE for the models parameters is obtained. Simulation studies are given to support the theoretical results. © 2015 Bull. Georg. Natl. Acad. Sci.

Key words: waring distribution, Pareto-like distribution, Markov Chain Monte Carlo (MCMC), MLE.

1. Introduction and Preliminaries

Several parametric frequency distributions have been proposed for the needs of biomolecular applications. Among them, there are two well-known classes of frequency distributions introduced by V. A. Kuznetsov [1, 2]: Waring and Pareto-like frequency distributions. These distributions may be described for modeling phenomena in biomolecular systems such as for the number of expressed genes in the transcriptome, the number of protein domain occurrences in the proteomes. For details about the applications and properties of such distributions see [1-4].

The Waring frequency distribution is a two-parameter distribution with the following probability function (see [1, 3]):

$$\begin{cases} f_x(\theta) = (f_0(\theta))^{-1} \cdot \prod_{k=0}^{x-1} \frac{p+k}{q+k}, & x = 1, 2, \dots, \\ f_0(\theta) = \sum_{y=1}^{\infty} \prod_{k=0}^{y-1} \frac{p+k}{q+k} \end{cases} \quad (1)$$

where $\theta = (p, q)$ and $0 < p < q < \infty$. Note that p and q are the model's parameters.

Another important family is the two-parameter Pareto-like frequency distribution with the following prob-

ability function (see [2], [3]):

$$\begin{cases} f_x(\alpha) = (f_0(\alpha))^{-1} \cdot (x+b)^{-\rho}, & x = 1, 2, \dots, \\ f_0(\alpha) = \sum_{y=1}^{\infty} (y+b)^{-\rho} \end{cases} \quad (2)$$

where $\alpha = (\rho, b)$, $1 < \rho < \infty$ is the shape parameter and $-1 < b < \infty$ is the location parameter.

It is of interest to investigate the statistical properties of the parameters for the models (1) and (2). But, the lack of simple closed formulas for the probability mass and cumulative distribution functions were major drawbacks to the use of such frequency distributions. In this paper we are going to study some statistical inferences for the models (1) and (2).

Assuming $X^n = (X_1, \dots, X_n)$, with observation $x^n = (x_1, \dots, x_n)$, is sample corresponding to a random variable ξ with the probability distribution given by (1) or (2). For $x \in \mathbb{N}$, we fix the following notations to be used in this paper:

$$\begin{aligned} h(x; \theta) &= \sum_{k=0}^{x-1} \frac{1}{p+k}, & l(x; \theta) &= \sum_{k=0}^{x-1} \frac{-1}{q+k}, & k(x; \theta) &= \sum_{k=0}^{x-1} \frac{-1}{(p+k)^2}, \\ t(x; \theta) &= \sum_{k=0}^{x-1} \frac{1}{(q+k)^2}, & m(x; \alpha) &= -\ln(x+b), & w(x; \alpha) &= \frac{-\rho}{(x+b)}, \\ \overline{f^n(\theta)} &= \frac{1}{n} \sum_{i=1}^n f(x_i; \theta), & \overline{f^n(\alpha)} &= \frac{1}{n} \sum_{i=1}^n f(x_i; \alpha) \end{aligned}$$

Also, the symbols $E[\cdot]$, $Var(\cdot)$ and $Cov(\cdot, \cdot)$ refer to the mathematical expectation, variance and covariance, respectively.

2. MLE

Supposing, as above, ξ is a discrete random variable with the probability distribution given by formula (1) or (2). We first establish a lemma, which will be used in the proof of the main theorem. Here, we consider the model (1).

Lemma 1. *For the model (1), we have*

$$E_{\theta} [h(\xi; \theta)] < \infty, \quad E_{\theta} [l(\xi; \theta)] < \infty.$$

Proof. By using definition of mathematical expectation the proof is met obviously.

Now, from Lemma 1 the following theorem is presented.

Theorem 1. *The MLE of the parameter θ is obtained from the following moment equations:*

$$\begin{cases} E_{\theta} [h(\xi; \theta)] = \overline{h^n(\theta)}, \\ E_{\theta} [l(\xi; \theta)] = \overline{l^n(\theta)}. \end{cases} \quad (3)$$

Proof. We derive the logarithm of likelihood function as follows

$$l(X^n; \theta) = \ln \prod_{i=1}^n \left(f_0(\theta) \cdot \prod_{k=0}^{x_i-1} \frac{p+k}{q+k} \right) = -n \ln f_0(\theta) + \sum_{i=1}^n \sum_{k=0}^{x_i-1} \ln \left(\frac{p+k}{q+k} \right). \quad (4)$$

In order to existence of the MLE for the parameter θ , the necessary conditions are expressed as

$$\frac{\partial l(X^n; \theta)}{\partial \theta_i} = 0, \quad i = 1, 2,$$

where $\theta_1 = p$, $\theta_2 = q$.

Let us establish the derivatives with respect to the parameters p and q . We obtain:

$$\frac{\partial l(X^n; \theta)}{\partial p} = -n \frac{1}{f_0(\theta)} \cdot \frac{\partial f_0(\theta)}{\partial p} + \sum_{i=1}^n \sum_{k=0}^{x_i-1} \frac{1}{p+k},$$

where $\frac{1}{f_0(\theta)} \cdot \frac{\partial f_0(\theta)}{\partial p} = E_\theta [h(\xi; \theta)]$. From $\frac{\partial l(X^n; \theta)}{\partial p} = 0$ it follows that

$$E_\theta [h(\xi; \theta)] = \overline{h^n(\theta)}.$$

Meanwhile, we get

$$\frac{\partial l(X^n; \theta)}{\partial q} = -n \frac{1}{f_0(\theta)} \cdot \frac{\partial f_0(\theta)}{\partial q} + \sum_{i=1}^n \sum_{k=0}^{x_i-1} \frac{-1}{q+k}$$

Taking into account that $\frac{1}{f_0(\theta)} \cdot \frac{\partial f_0(\theta)}{\partial q} = E_\theta [l(\xi; \theta)]$, then from the condition $\frac{\partial l(X^n; \theta)}{\partial q} = 0$ we conclude

$$E_\theta [l(\xi; \theta)] = \overline{l^n(\theta)}.$$

The proof of Theorem 1 is complete.

Now, we are going to show that the solution $\hat{\theta} = \hat{\theta}_i^n = \left(\hat{\theta}_i^n \right)_{i=1}^2$ of the system (3) (if it exists) is the MLE of the parameter θ . It suffices to establish that the matrix

$\widehat{M}_n = \left(\widehat{M}_{i,j}^n \right)_{i,j=1}^2$ with $\widehat{M}_{i,j}^n = \widehat{M}_{i,j}^n(\hat{\theta})$, $\widehat{M}_{i,j}^n(\hat{\theta}) = \left. \frac{\partial l(X^n, \theta)}{\partial \theta_i \partial \theta_j} \right|_{\theta=\hat{\theta}}$ is *negative definite*. We record the following assertion.

Lemma 2. *Assuming the solution $\hat{\theta}$ of the system (3) (if it exists) holds in the following conditions*

$$\begin{cases} E_\theta [k(\xi; \theta)] = \overline{k^n(\theta)} \\ E_\theta [l(\xi; \theta)] = \overline{l^n(\theta)} \end{cases} \quad (5)$$

then the elements of the matrix \widehat{M}_n are

$$\begin{aligned} \widehat{M}_{11}^n &= -n \text{Var}_{\hat{\theta}}(h(\xi; \theta)), \\ \widehat{M}_{12}^n &= \widehat{M}_{21}^n = -n \text{Cov}_{\hat{\theta}}(h(\xi; \theta), l(\xi; \theta)), \\ \widehat{M}_{22}^n &= -n \text{Var}_{\hat{\theta}}(l(\xi; \theta)). \end{aligned}$$

Proof. It is easy to see that

$$M_{11}^n = \frac{\partial^2 l(X^n; \theta)}{\partial p^2} = -n \left(\frac{1}{f_0(\theta)} \cdot \frac{\partial^2 f_0(\theta)}{\partial p^2} - \left(\frac{1}{f_0(\theta)} \cdot \frac{\partial f_0(\theta)}{\partial p} \right)^2 \right) + \overline{nk^n(\theta)},$$

$$M_{12}^n = M_{21}^n = \frac{\partial^2 l(X^n; \theta)}{\partial p \partial q} = \frac{\partial^2 l(X^n; \theta)}{\partial q \partial p} = -n \left(\frac{1}{f_0(\theta)} \cdot \frac{\partial^2 f_0(\theta)}{\partial p \partial q} - \frac{1}{(f_0(\theta))^2} \cdot \frac{\partial f_0(\theta)}{\partial p} \cdot \frac{\partial f_0(\theta)}{\partial q} \right),$$

$$M_{22}^n = \frac{\partial^2 l(X^n; \theta)}{\partial q^2} = -n \left(\frac{1}{f_0(\theta)} \cdot \frac{\partial^2 f_0(\theta)}{\partial q^2} - \left(\frac{1}{f_0(\theta)} \cdot \frac{\partial f_0(\theta)}{\partial q} \right)^2 \right) + \overline{nt^n(\theta)}$$

After some calculations and simplifications we obtain

$$M_{11}^n = -n \text{Var}_\theta (h(\xi; \theta)) - \left(n E_\theta [k(\xi; \theta)] - \overline{nk^n(\theta)} \right)$$

$$M_{12}^n = M_{21}^n = -n \text{Cov}_\theta (h(\xi; \theta), l(\xi; \theta)),$$

$$M_{22}^n = -n \text{Var}_\theta (l(\xi; \theta)) - \left(n E_\theta [t(\xi; \theta)] - \overline{nt^n(\theta)} \right).$$

From the condition (5) the proof is completed.

Lemma 3. Under satisfying the conditions (5), the matrix \widehat{M}_n is negative definite.

Proof. In order to due that it is enough to demonstrate $\widehat{M}_{11}^n < 0$ and $\det(\widehat{M}_n) > 0$.

From Lemma 2 it is obvious $\widehat{M}_{11}^n < 0$. To prove that $\det(\widehat{M}_n) > 0$ we have

$$\det(\widehat{M}_n) = \widehat{M}_{11}^n \widehat{M}_{22}^n - \left(\widehat{M}_{12}^n \right)^2.$$

According to the value of \widehat{M}_{11}^n , \widehat{M}_{12}^n , \widehat{M}_{22}^n and based on Cauchy-Bunyakovski-Schwartz inequality (see, for example, [5]) the proof is finished.

As an immediate consequence of Lemmas 2 and 3 the following result is received.

Corollary 1. If the solution of the system (3) satisfies the conditions (5), then it coincides with the MLE of the parameter θ .

Remark 1. We notice that the above mentioned results can also be proved for the model (2).

Thus, compare to Theorem 1, let us propose the following assertion for the model (2).

Theorem 2. The MLE of the parameter α is given from the following moment equations:

$$\begin{cases} E_\alpha [m(\xi; \alpha)] = \overline{m^n(\alpha)}, \\ E_\alpha [w(\xi; \alpha)] = \overline{w^n(\alpha)}. \end{cases} \tag{6}$$

Proof. Similar to the proof of Theorem 1.

Corollary 2. Similarly, the Lemmas 2 and 3 may be considered for the two-parameter Pareto-like frequency distribution given by formula (2). Because of similarity, we leave the results and proofs for the readers.

3. Approximate Computation of the MLE

Based on the systems (3) and (6), it is not simple to derive closed forms for the solutions. So, we need to propose a method of the approximate computation of the MLE for the models parameters. In order to do that (compare to [6], p. 228), the Fisher’s accumulation method is suggested. For more details about the Fisher’s accumulation method we refer the readers to, for example, [7], p. 88.

3.1 Waring Distribution

Supposing $\theta_0 = (p_0, q_0)$ is an initial value for the parameter $\theta = (p, q)$. Recurrently, we find out the $(s + 1)_{th}$ approximation by using the formula (compare to [6], p. 228):

$$\theta_j(s+1) = \theta_j(s) + \frac{\Lambda_j(\theta(s))}{n \cdot \det I(\theta(s))}, \quad j = 1, 2, \quad (7)$$

where $I(\cdot)$ is the Fisher's information measure. In details, the formula (7) is considered as

$$\begin{cases} p(s+1) = p(s) + \frac{\Lambda_p(\theta(s))}{n \cdot \det I(\theta(s))}, \\ q(s+1) = q(s) + \frac{\Lambda_q(\theta(s))}{n \cdot \det I(\theta(s))}, \end{cases} \quad (8)$$

where

$$\Lambda_p(\theta(s)) = \begin{pmatrix} U_1(\theta) & I_{12}(\theta) \\ U_2(\theta) & I_{22}(\theta) \end{pmatrix}, \quad \Lambda_q(\theta(s)) = \begin{pmatrix} I_{11}(\theta) & U_1(\theta) \\ I_{21}(\theta) & U_2(\theta) \end{pmatrix},$$

and $U_1(\theta) = \frac{\partial l(X^n; \theta)}{\partial p}$ and $U_2(\theta) = \frac{\partial l(X^n; \theta)}{\partial q}$ are the *contribution functions*. We obtain

$$U_1(\theta) = -n E_\theta [h(\xi; \theta)] + n \overline{h^n(\theta)},$$

$$U_2(\theta) = -n E_\theta [l(\xi; \theta)] + n \overline{l^n(\theta)}.$$

Comparing to Farbod and Gasparian ([6], p. 229), let us give the following algorithm for the model (1).

Algorithm I.

1. Generate random numbers from the model (1) using MCMC method;
2. Use (7) (or (8)) in order to calculate $\theta_j(s)$, $j = 1, 2$, $s = 0, 1, 2, \dots$;
3. If $|\theta_j(s+1) - \theta_j(s)| < \varepsilon$ (ε is some small positive constant), then $\theta_j(s+1) = \hat{\theta}$ is the MLE, otherwise go to the step 2.

3.2 Pareto-Like Distribution

In a similar way, for the model (2) we establish

$$\alpha_j(s+1) = \alpha_j(s) + \frac{\Lambda_j(\alpha(s))}{n \cdot \det I(\alpha(s))}, \quad j = 1, 2. \quad (9)$$

In details, we get

$$\begin{cases} \rho(s+1) = \rho(s) + \frac{\Lambda_\rho(\alpha(s))}{n \cdot \det I(\alpha(s))} \\ b(s+1) = b(s) + \frac{\Lambda_b(\alpha(s))}{n \cdot \det I(\alpha(s))} \end{cases} \quad (10)$$

where

$$\Lambda_\rho(\alpha(s)) = \begin{pmatrix} U_1(\alpha) & I_{12}(\alpha) \\ U_2(\alpha) & I_{22}(\alpha) \end{pmatrix}, \quad \Lambda_b(\alpha(s)) = \begin{pmatrix} I_{11}(\alpha) & U_1(\alpha) \\ I_{21}(\alpha) & U_2(\alpha) \end{pmatrix},$$

and $U_1(\alpha) = \frac{\partial l(X^n; \alpha)}{\partial \rho}$ and $U_2(\alpha) = \frac{\partial l(X^n; \alpha)}{\partial b}$. We have

$$U_1(\alpha) = -n E_\alpha [m(\xi; \alpha)] + n \overline{m^n(\theta)}, \quad U_2(\alpha) = -n E_\alpha [w(\xi; \alpha)] + n \overline{w^n(\theta)}.$$

Again, by comparing Farbod and Gasparian ([6], p. 229), we record the following algorithm for the model (2).

Algorithm II.

1. Generate random numbers from the model (2) using MCMC method;
2. Use (9) (or (10)) in order to calculate $\alpha_j(s)$, $j = 1, 2$; $s = 0, 1, 2, \dots$;
3. If $|\alpha_j(s+1) - \alpha_j(s)| < \varepsilon$ (ε is some small positive constant), then $\alpha_j(s+1) = \hat{\alpha}$ is the MLE, other wise go to the step 2.

3.3 Simulation Studies

As we refereed in Introduction, to the author’s knowledge, by now there have not been presented any closed expressions for the *cumulative distribution functions* of the models (1) and (2). Therefore, it is not possible to generate sample numbers at random based on the *cumulative distribution functions*. To overcome this difficulty, the MCMC method (see [8]) is used. Because of numerical calculations, let us consider the random variable ξ as *truncated*. Here, random variable ξ is restricted to 100. From Algorithms I and II we propose the Tables 1 and 2 for the models (1) and (2), respectively.

Note 1. The value $\theta = (p = 0.5, q = 1.9)$ is given as true value of the parameters of the model (1). Simulation studies are done for 1000 times to illustrate the behavior of the MLE. Namely, we do simulation for $M = 1000$ (M is the number of iteration) and $N = 20, 50, 100, 200$ (N is the sample size). Consider $\varepsilon = 0.0005$.

Table 1. The mean of estimations and the mean square errors (MSE) for the model (1)

	N=20		N=50	
Case	Mean	MSE	Mean	MSE
p	2.7226527	4.940185	2.7222685	4.938477
q	3.5981751	3.029295	3.4578110	2.962381
Iteration	274	---	208	---
	N=100		N=200	
Case	Mean	MSE	Mean	MSE
p	2.1383215	4.620950	1.8123654	3.818908
q	3.1511743	2.254472	2.5301244	1.743261
Iteration	139	---	187	---

Table 2. The mean of estimations and the MSE for the model (2)

	N=20		N=50	
Case	Mean	MSE	Mean	MSE
ρ	4.865513	4.266343	4.105983	1.705592
b	5.559738	18.145364	3.777491	6.137962
Iteration	1451	---	715	---
	N=100		N=200	
Case	Mean	MSE	Mean	MSE
ρ	2.081051	0.5168874	2.766523	0.0011207
b	0.477143	0.6770937	1.165229	0.0181633
Iteration	208	---	367	---

The mean of estimations and the MSE are given. As we see from Table 1, with increasing sample size the MSE decreases.

Note 2. The value $\alpha = (\rho = 2.8, b = 1.3)$ is considered as true value of the parameters of the model (2).

We present the mean of estimations and the MSE. As we expected, from Table 2 it is readily seen that with increasing sample size the MSE decreases.

Note 3. All computations and simulations have been performed by using “R” statistical software.

Acknowledgement. This work was supported by a grant from Quchan University of Advanced Technology, Quchan, Iran.

ბიონფორმატიკა

ბიონფორმატიკაში წარმოშობილი სიხშირის განაწილების შეფასების ფორმულები

დ. ფარბოდი

ქუჩანის მოწინავე ტექნოლოგიების უნივერსიტეტი, მათემატიკის განყოფილება, ქუჩანი, ირანი

(წარმოდგენილია აკადემიის წევრის მ. სალუქჰადის მიერ)

განვიხილავთ სიხშირის განაწილების ორ ცნობილ სახეს, რომელიც გამოიყენება ინფორმატიკაში წარმოშობილი ფენომენების მოდელირებისათვის. ესენია, ვარინგის ორპარამეტრიანი სიხშირის განაწილება და პარეტოს ტიპის ორპარამეტრიანი სიხშირის განაწილება, რომელიც ვ. ა. კუზნეცოვმა გაგვაცნო. გთავაზობთ მაქსიმალური აღბათობის შეფასების ფორმულებს (მაშფ) ასეთი სიხშირის განაწილების უცნობი პარამეტრებისათვის. შემოთავაზებული მაშფ-ი მყისიერი შეფასების ზოგიერთ ფორმულას ემთხვევა. ამასთან ერთად, მიღებულია მაშფ-ის მიახლოებითი გამოთვლის მეთოდი მოდელირების პარამეტრებისთვის. თეორიული შედეგების დასაბუთების მიზნით მოცემულია სიმულაციური კვლევები.

REFERENCES

1. Kuznetsov V. A. (2003) Signal Processing, **33**(4): 889-910.
2. Kuznetsov V. A. (2001) EURASIP J. Appl. Signal Processing, **4**: 285-296.
3. Astola J. and Danielian E. (2007) Frequency Distributions in Biomolecular Systems and Growing Networks, TICSP Series no. 31, Tampere, Finland
4. Danielian E. and Astola J. (2007) Proceedings of International TICSP Workshop on Spectral Methods and Multirate Signal Processing, Moscow, Russia, 235-237.
5. Shiryaev A. N. (1995) Probability (Graduate Texts in Mathematics), **95**, Springer. Translated from Russian.
6. Farbod D. and Gasparian K. (2013) J. Iranian Statist. Soc. (JIRSS), **12**(2): 211-234.
7. Ivchenko G. and Medvedev Yu. (1990) Mathematical Statistics. Translated from Russian.
8. Rizzo M. L. (2008) Statistical Computing with R, Chapman and Hall/CRC.

Received February, 2015