

Research Article

The Discrete Exponentiated Pareto Distribution: Its Properties and Application

Yadapa Chotedelok¹ and Winai Bodhisuwan^{1*}

¹Department of Statistics, Faculty of Science, Kasetsart University, Chatuchak, Bangkok 10900, Thailand

*E-mail: fsciwnb@ku.ac.th

Received: 15/05/2020; Revised: 09/08/2020; Accepted: 23/08/2020

Abstract

The discrete exponentiated Pareto (DEP) distribution is developed by using the discretization method based on the survival function. It is discretized from the exponentiated Pareto distribution. In this paper, a probability mass function of the DEP distribution is derived. Some mathematical properties and model parameters estimation are discussed. In addition, we applied the DEP distribution to two real datasets. The results of model fitting of these datasets based on the DEP distribution are reasonably constructive. The proposed distribution performs well with a goodness of fit test and some criterions. The distribution can be used as an alternative model for discrete data analytics.

Keywords: exponentiated Pareto distribution, discretization method, survival function, maximum likelihood estimation

Introduction

Lifetime data are usually described by continuous distribution such as exponential distribution, Pareto distribution, and Weibull distribution. Because of the precision of the measuring instrument for collecting the data and the nature of the data in the long term, sometimes the data are obviously presented in a discrete sense. For instance, in Figure 1, (a) the temperature is presented in the mobile application and (b) battery charging which is shown as a percentage in the mobile phone. These examples are actually continuous data in nature but they are always presented in discrete integer value for some specified purposes. So, it can be implied that the variables are continuous in nature but they are presenting in a discrete sense.

In such a situation, there is a method to generate the discrete distribution from the continuous distribution, it is so-called the discretization method (Chakraborty, 2015). The most popular method is the discretization that utilized the survival function to construct the discrete distribution. For example, Nakagawa and Osaki (1975) proposed the discrete Weibull distribution and Roy (1993) studied between the exponential distribution and the geometric distribution that both were related, these distributions were presented by applying the discretization with survival function. In addition, several researchers have proposed the new discrete distributions such as discrete normal distribution (Roy, 2003), discrete Maxwell distribution (Krishna & Pundir, 2007), discrete Burr distribution (Krishna & Pundir, 2009), discrete Pareto distribution (Krishna & Pundir, 2009), discrete Lindley distribution (Gómez-Déniz & Calderín-Ojeda, 2011), discrete gamma

distribution (Chakraborty & Chakravarty, 2012), discrete Gumbel distribution (Chakraborty & Chakravarty, 2014), discrete inverse Rayleigh distribution (Hussain & Ahmad, 2014) and discrete asymmetric Laplace distribution (Sangpoom & Bodhisuwan, 2016).

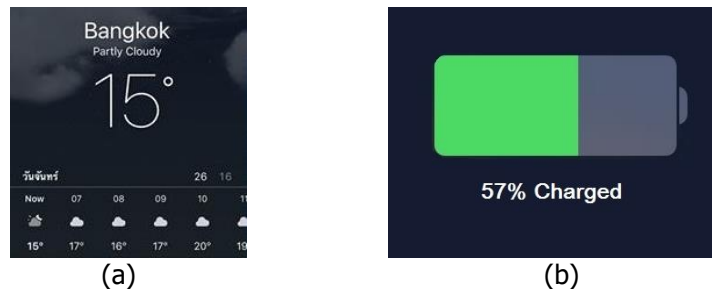


Figure 1 Some examples of the continuous data which they were presented in a discrete sense.

In this paper, we present a discrete analogue of the exponentiated Pareto distribution by discretizing the continuous exponentiated Pareto distribution utilized the method of Roy. The exponentiated Pareto distribution is modified from the Pareto distribution and proposed for wider applicability (Nadarajah, 2005). In addition, the exponentiated Pareto distribution is studied by many researchers such as in 2015, Fatima and Roohi introduced the transmuted exponentiated Pareto-I distribution via the transmutation technique (Fatima & Roohi, 2015). Jabbari Nooghabi studied the parameter estimation of the exponentiated Pareto distribution in the presence of outliers in 2017 (Jabbari Nooghabi, 2017). Moreover, Bhatti and Ali presented the characterizations of transmuted exponentiated Pareto-I distribution in 2019 (Bhatti & Ali, 2019).

The rest of the paper is as follows. In the part of methods, the proposed distribution is introduced. Its essential mathematical properties and parameter estimation are discussed in results and discussion. The applications of the DEP distribution are demonstrated which we applied the proposed distribution to some real datasets. Furthermore, the conclusion is presented in the last section.

Methods

There are several ways to derive discrete analogues of continuous distribution such as the method based on probability mass function (pmf), cumulative distribution function (cdf), survival function, hazard rate function, etc. (Chakraborty, 2015). For this paper, we determine the survival function to discretize the continuous to discrete distribution which is the well-known method of discretization.

The discretization method based on survival function is using the difference values between survival values $S_X(x)$ and $S_X(x+1)$. If the underlying continuous random variable X has the survival function $S_X(x)$, and the random variable $Y = \lfloor X \rfloor$ where $\lfloor X \rfloor$ is the largest integer but not greater than X and $y = 0, 1, 2, \dots$, then the pmf of Y is

$$\begin{aligned} f(y) &= P(Y = y) \\ &= P(\lfloor X \rfloor = y) \\ &= P(y \leq X < y+1) \\ &= F_X(y+1) - F_X(y) \\ &= S_X(y) - S_X(y+1) \end{aligned} \quad (1)$$

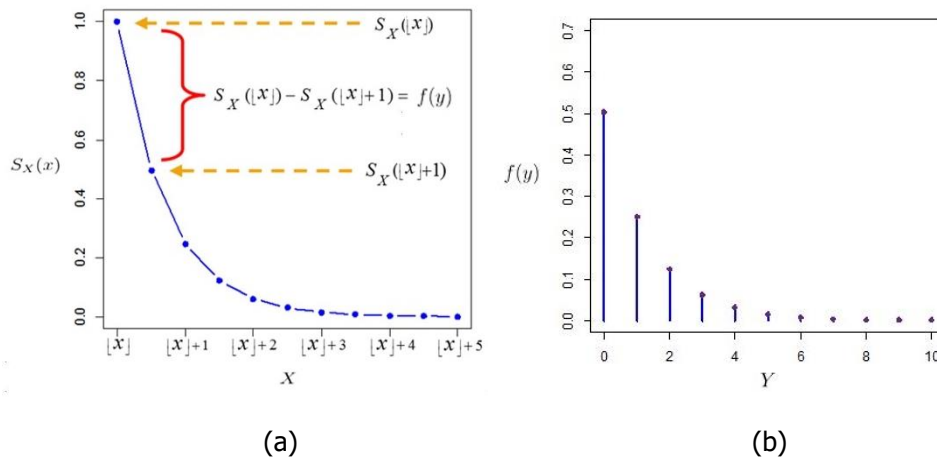


Figure 2 (a) The survival plot of a continuous random variable X and (b) the pmf plot of discretizing X .

The exponentiated Pareto distribution was developed by Nadarajah (2005) with distribution function

$$F(x) = 1 - k^a e^{(-ax)},$$

for $x > \log k$, $k > 0$ is a scale parameter and $a > 0$ is a shape parameter.

Immediately, the survival function of the exponentiated Pareto distribution can be obtained as

$$S(x) = 1 - F(x) = k^a e^{(-ax)}. \quad (2)$$

Some probability density function (pdf) and survival plots of the exponentiated Pareto distribution are illustrated in Figure 3. All of the pdf and survival plots are a decreasing function. The scale of distribution is in keeping with parameter k and it is increased when k is increasing.

Results and Discussion

Discrete Exponentiated Pareto Distribution

We use the discretization method based on the survival function to develop the DEP distribution. The verification of the pmf and the survival function are presented in the following in theorems.

Theorem 1. Let Y be a random variable of the DEP distribution, $Y \sim \text{DEP}(k, a)$. The pmf of the DEP distribution is

$$f(y) = k^a e^{-ay} (1 - e^{-a}), \quad (3)$$

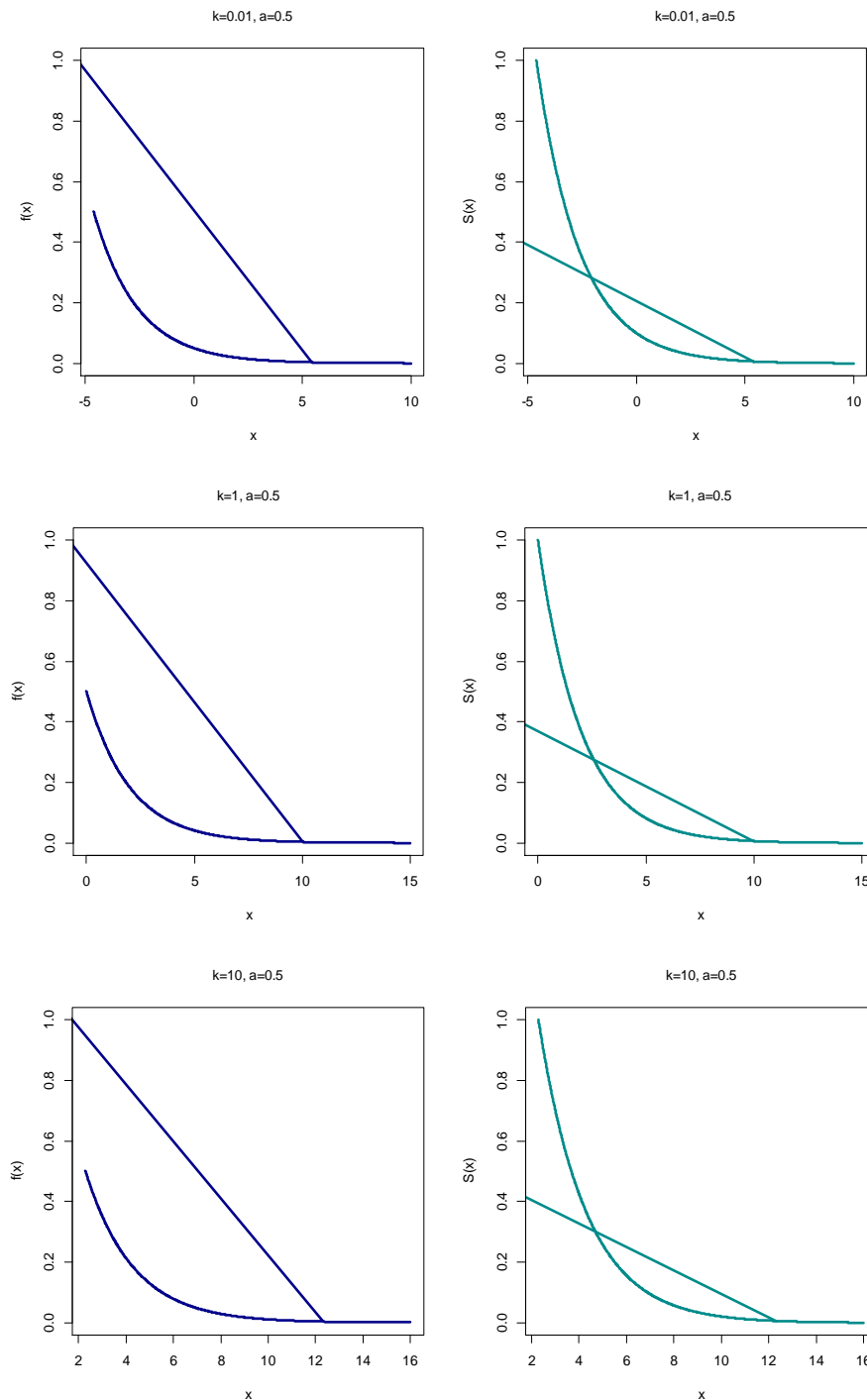


Figure 3 Some pdf (left) and survival (right) plots of the exponentiated Pareto distribution according to different values of k and a .

where $y = \lfloor \log k \rfloor, \lfloor (\log k) + 1 \rfloor, \lfloor (\log k) + 2 \rfloor, \dots$, $k > 0$ is a scale parameter and $a > 0$ is a shape parameter.

Proof: Underlying the survival function of the exponentiated Pareto distribution and the random variable $Y = \lfloor X \rfloor$, the pmf is

$$\begin{aligned} f(y) &= P(Y = y) \\ &= S_{\lfloor X \rfloor}(y) - S_{\lfloor X \rfloor}(y + 1) \\ &= k^a e^{-ay} - k^a e^{-a(y+1)} \\ &= k^a e^{-ay} (1 - e^{-a}). \end{aligned}$$

Theorem 2. Let Y is a random variable of the DEP distribution, $Y \sim \text{DEP}(k, a)$. The cdf of Y is

$$F(y) = 1 - k^a e^{-a(y+1)} \quad (4)$$

where $y = \lfloor \log k \rfloor, \lfloor (\log k) + 1 \rfloor, \lfloor (\log k) + 2 \rfloor, \dots$ and parameters $k, a > 0$.

Proof: If Y be a random variable of the DEP distribution with the pmf in Equation (3), then the cdf of Y can be obtained from

$$\begin{aligned} F(y) &= \sum_{y=\log k}^y f(y) \\ &= \sum_{y=\log k}^y k^a e^{-ay} (1 - e^{-a}) \\ &= k^a (1 - e^{-a}) \sum_{y=\log k}^y e^{-ay}, \end{aligned}$$

since $\sum_{y=\log k}^y e^{-ay} = \frac{k^{-a} - e^{-ay} e^{-a}}{1 - e^{-a}}$ is a geometric series, then the cdf of Y is

$$F(y) = 1 - k^a e^{-a(y+1)}.$$

Theorem 3. If Y is a random variable of the DEP distribution, denoted by $Y \sim \text{DEP}(k, a)$, then its survival function is

$$S(y) = k^a e^{-a(y+1)}, \quad (5)$$

where $y = \lfloor \log k \rfloor, \lfloor (\log k) + 1 \rfloor, \lfloor (\log k) + 2 \rfloor, \dots$ and parameters $k, a > 0$.

Proof: Since $Y \sim \text{DEP}(k, a)$ with the cdf in Equation (4) and $S(y)$ is defined as the survival function of DEP distribution. Thus,

$$\begin{aligned} S(y) &= 1 - F(y) \\ &= 1 - (1 - k^a e^{-a(y+1)}) \\ &= k^a e^{-a(y+1)}. \end{aligned}$$

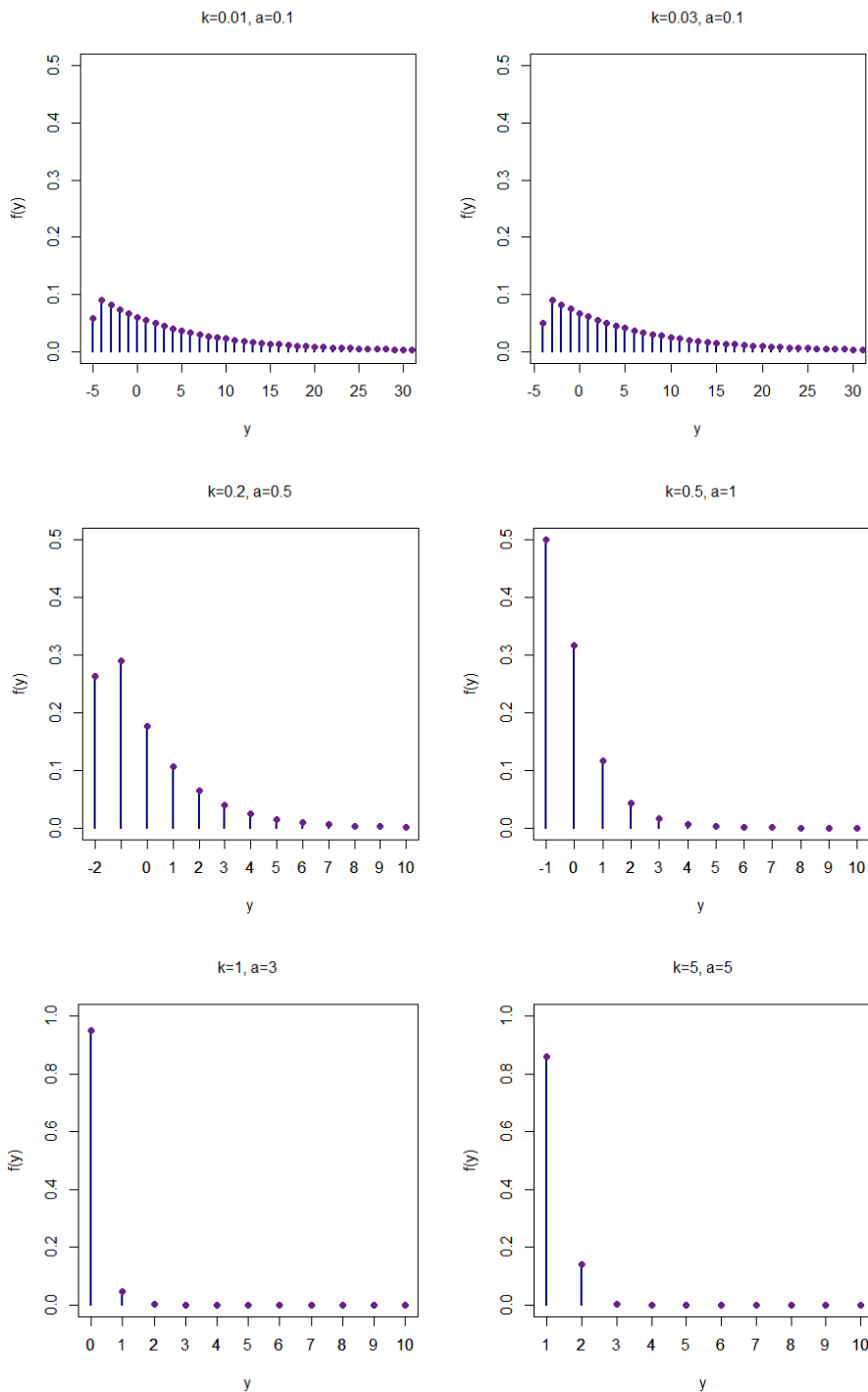


Figure 4 Some pmf plots of the DEP distribution with various values of k and a .

Figure 4 shows some pmf plots of the DEP distribution. Evidently, the scale of distribution changes according to the parameter k . The random variable Y can be a negative integer when k less than 1. The shape of distribution respect to parameter a . It appears that clearly, the pmf decreases faster as parameter a increases and it is a unimodal curve when a less than 1. So, the DEP distribution has the right skew and unimodal curve.

Mathematical Properties

Some mathematical properties of the DEP distribution, especially the moment generating function (mgf) and quantile function are provided in this section.

Moment Generating Function

Theorem 4. Let Y be a random variable of the DEP distribution, $Y \sim \text{DEP}(k, a)$. The mgf of Y , denoted by $M_Y(t)$, is

$$M_Y(t) = \frac{(e^a - 1) k^t}{e^a - e^t} \quad (6)$$

where $k > 0$, $a > 0$ and $t < a$.

Proof: The mgf of the DEP distribution can be obtained from

$$\begin{aligned} M_Y(t) &= E(e^{tY}) \\ &= \sum_{\forall y} e^{ty} \cdot f(y) \\ &= \sum_{y=\lfloor \log k \rfloor}^{\infty} e^{ty} \cdot k^a e^{-ay} (1 - e^{-a}) \\ &= k^a (1 - e^{-a}) \sum_{y=\lfloor \log k \rfloor}^{\infty} e^{ty-ay}, \end{aligned}$$

since $\sum_{y=\lfloor \log k \rfloor}^{\infty} e^{ty-ay} = \frac{k^{t-a}}{1 - e^{t-a}}$ is the geometric series, then the mgf will be

$$\begin{aligned} M_Y(t) &= k^a (1 - e^{-a}) \left(\frac{k^{t-a}}{1 - e^{t-a}} \right) \\ &= k^a (1 - e^{-a}) \left(\frac{k^t k^{-a}}{1 - e^{t-a}} \right) \\ &= \frac{(e^a - 1) k^t}{e^a - e^t}. \end{aligned}$$

By using the first four raw moments, we can find the mean, variance, skewness, and kurtosis of Y by successively differentiating $M_Y(t)$ and then evaluating the result at $t = 0$,

$E(Y^r) = \frac{d^r}{dt^r} M_Y(t)|_{t=0}$, for $r = 1, 2, \dots$. Consequently, the first four raw moments of Y are

$$E(Y) = \frac{(e^a - 1) \log k + 1}{e^a - 1},$$

$$E(Y^2) = \frac{(e^a - 1)^2 \log^2 k + 2(e^a - 1) \log k + (e^a - 1) + 2}{(e^a - 1)^2},$$

$$E(Y^3) = (e^a - 1) \left[\frac{(e^a - 1) + (e^a - 1)^2 \log^3 k + 3(e^a - 1) \log k + 3}{(e^a - 1)^3} + \frac{3((e^a - 1) + (e^a - 1)^2 \log^2 k + 2(e^a - 1) \log k + 2)}{(e^a - 1)^4} \right]$$

and

$$E(Y^4) = (e^a - 1) \left[\left(\frac{1}{(e^a - 1)^3} \right) ((e^a - 1) + 4(e^a - 1) \log k + 5 + ((e^a - 1) \log^3 k ((e^a - 1) \log k - 2)) + 3(e^a - 1) \log^2 k) + \left(\frac{1}{(e^a - 1)^4} \right) (9(e^a - 1) + 6(e^a - 1)^2 \log^3 k + 3(e^a - 1)^2 \log^2 k + 24(e^a - 1) \log k + 24) + \left(\frac{1}{(e^a - 1)^5} \right) (12(e^a - 1) + 12(e^a - 1)^2 \log^2 k + 24(e^a - 1) \log k + 24) \right].$$

Hence, the mean, variance, skewness, and kurtosis of $Y \sim \text{DEP}(k, a)$ according to its first four raw moments, respectively, are

$$E(Y) = \frac{(e^a - 1) \log k + 1}{e^a - 1},$$

$$\text{Var}(Y) = \frac{e^a}{(e^a - 1)^2},$$

$$\text{Skewness}(Y) = \frac{E(Y^3) - 3E(Y)E(Y^2) + 2(E(Y))^3}{(\text{Var}(Y))^{3/2}}$$

and

$$\text{Kurtosis}(Y) = \frac{E(Y^4) - 4E(Y)E(Y^3) + 6E(Y^2)(E(Y))^2 - 3(E(Y))^4}{(\text{Var}(Y))^2}.$$

Quantile Function

Let Y be a random variable of the DEP distribution with cdf, $F_Y(y)$. The quantile function is the generalized inverse of $F_Y(y)$. Let U be distributed as the uniform on $(0, 1)$, $u \in (0, 1)$ and the quantile function denoted by $Q(u)$. According to the quantile function, if $F_Y(y) = u$ when $u = 1 - k^a e^{-a(Q_Y(u)+1)}$, then $Q_Y(u) = y = F_Y^{-1}(u)$. Thus, the quantile function of the DEP distribution is

$$y = F_Y^{-1}(u) = - \left(1 + \frac{\log \left(\frac{1-u}{k^a} \right)}{a} \right) \quad (7)$$

where $k > 0$ and $a > 0$.

There are many methods to generate random variates from a probability distribution. The method that simplicity and generality is the inverse of cdf. Therefore, the quantile function in Equation (7) is very useful for generating a random variable Y of the DEP distribution.

Parameter Estimation

The maximum likelihood estimation (MLE) is the widely used method for model parameter estimation. In this section, the MLE of the DEP distribution will be discussed.

Let Y_1, Y_2, \dots, Y_n be an independent and identically distributed (iid) random variables of the DEP distribution with the pmf of Equation (3). The likelihood function of the DEP distribution is given by

$$\begin{aligned} L(k, a|y_i) &= \prod_{i=1}^n f(y_i, k, a) \\ &= \prod_{i=1}^n k^a e^{-ay_i} (1 - e^{-a}). \end{aligned}$$

The log-likelihood function of n observations of Y can be written as

$$\begin{aligned} l = \log L(k, a|y_i) &= \log \prod_{i=1}^n k^a e^{-ay_i} (1 - e^{-a}) \\ &= \sum_{i=1}^n \log (k^a e^{-ay_i} (1 - e^{-a})). \end{aligned} \quad (8)$$

For the parameter estimation of the DEP distribution, two parameters are estimated. First, we estimate the parameter k , since $y = \log k$, the likelihood function is maximized with

$$\hat{k} = \min \{ e^{y_1}, e^{y_2}, \dots, e^{y_n} \}.$$

(See Rytgaard, 1990 and Mukhopadhyay & Ekwo, 1987 about estimation problems for k).

Next, \hat{k} is plugged into the Equation (8) and we take the derivative of the log-likelihood function with respect to a is given as

$$\begin{aligned}\frac{dl}{da} &= \frac{d}{da} \left[\sum_{i=1}^n \log \left(\hat{k}^a e^{-ay_i} (1 - e^{-a}) \right) \right] \\ &= n \log \hat{k} + \frac{n}{e^a - 1} - \sum_{i=1}^n y_i.\end{aligned}\quad (9)$$

The parameter a of the DEP distribution is estimated by setting this differential equation in the Equation (9) to zero, then solving this equation. The estimator of parameter a is

$$\hat{a} = \log \left(1 + \frac{n}{\sum_{i=1}^n y_i - n \log \hat{k}} \right).$$

However, for convenience and less complicated of the parameter estimation, the maximum likelihood estimators can be obtained by a numerical method. The `bbmle` package (Bolker & Team, 2019) of the R programming language (R Core Team, 2019) is the package for fitting maximum likelihood models, extended and modified from the `mle` function in `stat4` package. In this work, the `bbmle` package is employed.

Applications

We consider two real datasets to fit with the DEP distribution and the discrete Pareto (DP) distribution (Krishna & Pundir, 2009). The first dataset is the infant and child mortality in Sri Lanka from the Sri Lanka fertility survey in 1975 (Meegama, 1980). And the dataset is the electron-microscopic studies of the density of dystrophin in the fibers of the human quadriceps muscle, the number of attached particles from Immunogold data (Mathews & Appleton, 1993). In this work, the `bbmle` package of the R programming language is used to estimate parameters. Tables 1 and 2 show the results of fitting between the DEP and DP distributions to these real datasets. The appropriate distribution for fitting data is verified with the Anderson-Darling (AD) goodness of fit test for discrete data (Choulakian et al., 1994). The discrete AD test is obtained by using the `dgoF` package (Arnold & Emerson, 2011) in the R language. Other criteria for model selection that used to show the performance of the model are the minus log-likelihood (-LL), the Akaike information criterion (AIC), and the Bayesian information criterion (BIC). Furthermore, the comparison between real datasets and expected values of fitted distributions can be illustrated in Figure 5.

The fitted distributions for the number of infant and child deaths and the number of attached particles from Immunogold data are shown in Tables 1 and 2, respectively. The p -value based on the discrete AD test of the DEP distribution is greater than the DP distribution. Moreover, the DEP distribution gives the values of -LL, AIC, and BIC smaller than the DP distribution. Thus, the DEP distribution is more appropriate than the DP distribution.

Figure 5 displays the plots of the fitted frequency of the DEP and the DP distributions with the number of infant and child deaths and the number of attached particles from Immunogold data. It illustrates that the DEP distribution can be fitted more closely to these real

datasets than DP distribution. Therefore, the DEP distribution is more appropriate than the DP distribution related to the results in Table 1 and Table 2.

Table 1 Results of parameter estimation for the infant and child deaths data.

Number of infant and child deaths	Observed number of mothers	Expected frequency	
		DP	DEP
1	176	186.5613	172.0648
2	44	32.7931	50.7275
3	16	11.1243	14.9553
4	6	5.0337	4.4091
5	2	2.6860	1.2999
Estimated parameters		$\hat{k} = 1$ $\hat{a} = 2.0868$	$\hat{k} = 2.7183$ $\hat{a} = 1.2214$
-LL		217.2600	209.8100
AD-statistics (p -value)		1.3036 (0.1199)	0.2136 (0.6980)
AIC		436.5173	421.6231
BIC		440.0145	425.1203

Table 2 Goodness of fit test for the DP and DEP distributions for the number of attached particles from Immunogold data.

Number of attached particles	Observed frequency	Expected frequency	
		DP	DEP
1	122	139.6299	125.6493
2	50	29.8019	45.9132
3	18	11.3608	16.7770
4	4	5.5945	6.1305
5	4	3.1911	2.2401
Estimated parameters		$\hat{k} = 1$ $\hat{a} = 1.7622$	$\hat{k} = 2.7183$ $\hat{a} = 1.0067$
-LL		219.5200	204.8100
AD-statistics (p -value)		3.6074 (0.0079)	0.1457 (0.8235)
AIC		441.0385	411.6270
BIC		444.3268	414.9153

Conclusions

A discrete version of the continuous exponentiated Pareto distribution is proposed which called the DEP distribution. It is developed based on the discretization method of the survival function. We derived some essential mathematical properties, for instance, pmf, mgf, mean, variance, and quantile function. In addition, the parameter estimation by the maximum likelihood estimation is discussed. Furthermore, the proposed distribution is applied to two real datasets. The results for the comparison of -LL, AIC, and BIC and according to the p -value of the discrete AD test indicated that the DEP distribution is a better fit than the DP distribution for these real datasets. In conclusion, the DEP distribution may be a useful alternative to other distributions for discrete data analytics.

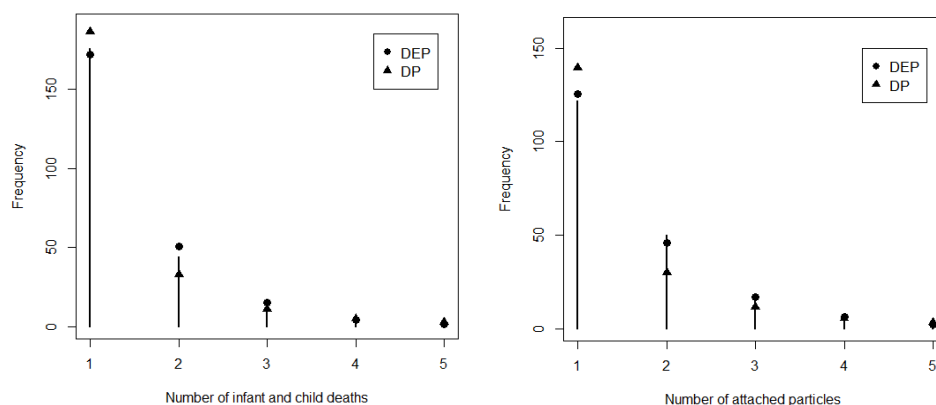


Figure 5 The fitted frequency of the DEP and the DP distributions to real datasets.

Acknowledgement

The authors would like to thank the Department of Statistics, Faculty of Science, Kasetsart University and the anonymous reviewers for their valuable suggestions.

References

- Arnold, T. B., & Emerson, J. W. (2011). Nonparametric goodness-of-fit tests for discrete null distributions. *R Journal*, 3(2), 34-39.
- Bolker, B. & Team, R. D. C. (2019). bbmle: Tools for General Maximum Likelihood Estimation. R package version 1.0.20.
- Bhatti, F. A. & Ali, A. (2019). Characterizations of transmuted exponentiated Pareto-I (TEP-I) distribution. *International Journal of Modern Mathematical Sciences*, 17(1), 1-20.
- Chakraborty, S. (2015). Generating discrete analogues of continuous probability distributions-A survey of methods and constructions. *Journal of Statistical Distributions and Applications*, 2(1), 6. <https://doi.org/10.1186/s40488-015-0028-6>
- Chakraborty, S. & Chakravarty, D. (2012). Discrete gamma distributions: properties and parameter estimations. *Communications in Statistics-Theory and Methods*, 41(18), 3301-3324. <https://doi.org/10.1080/03610926.2011.563014>
- Chakraborty, S. & Chakravarty, D. (2014). A discrete Gumbel distribution. arXiv:1410.7568
- Choulakian, V., Lockhart, R. A., & Stephens, M. A. (1994). Cramér-von Mises statistics for discrete distributions. *The Canadian Journal of Statistics*, 22(1), 125-137. doi: 10.2307/3315828
- Fatima, A. & Roohi, A. (2015). Transmuted exponentiated Pareto-I distribution. *Pakistan Journal of Statistics*, 32(1), 63-80.
- Gómez-Déniz, E. & Calderín-Ojeda, E. (2011). The discrete Lindley distribution: properties and applications. *Journal of Statistical Computation and Simulation*, 81(11), 1405-1416. doi.org/10.1080/00949655.2010.487825
- Hussain, T. & Ahmad, M. U. (2014). Discrete inverse Rayleigh distribution. *Pakistan Journal of Statistics*, 30(2), 203-222.

- Jabbari Nooghabi, M. (2017). On estimation in the exponentiated Pareto distribution in the presence of outliers. *Applied Mathematics and Information Sciences*, 11(4), 1129-1137. doi:10.18576/amis/110420
- Krishna, H. & Pundir, P. S. (2007). Discrete Maxwell distribution. *InterStat*, 3.
- Krishna, H. & Pundir, P. S. (2009). Discrete Burr and discrete Pareto distributions. *Statistical Methodology*, 6(2), 177-188. <https://doi.org/10.1016/j.stamet.2008.07.001>
- Matthews, J. & Appleton, D. (1993). An Application of the Truncated Poisson Distribution to Immunogold Assay. *Biometrics*, 49(2), 617-621. doi: 10.2307/2532574
- Meegama, S. A. (1980). Socio-economic determinants of infant and child mortality in Sri Lanka: An analysis of post-war experience. *Scientific Report (8)*. London, England: World Fertility Survey.
- Mukhopadhyay, N., & Ekwo, M. E. (1987). Sequential estimation problems for the scale parameter of a Pareto distribution. *Scandinavian Actuarial Journal*, 1987(1-2), 83-103. <https://doi.org/10.1080/03461238.1987.10413820>
- Nadarajah, S. (2005). Exponentiated Pareto distributions. *Statistics*, 39(3), 255-260. <https://doi.org/10.1080/02331880500065488>
- Nakagawa, T. & Osaki, S. (1975). The discrete Weibull distribution. *IEEE Transactions on Reliability*, 24(5), 300-301. doi: 10.1109/TR.1975.5214915
- R Core Team. (2019). R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing. Vienna, Austria.
- Roy, D. (1993). Reliability measures in the discrete bivariate set up and related characterization results for a bivariate geometric distribution. *Journal of Multivariate Analysis*, 46(2), 362-373. <https://doi.org/10.1006/jmva.1993.1065>
- Roy, D. (2003). The discrete normal distribution. *Communications in Statistics Theory and Methods*, 32(10), 1871-1883. <https://doi.org/10.1081/STA-120023256>
- Rytgaard, M. (1990). Estimation in the Pareto distribution. *ASTIN Bulletin: The Journal of the IAA*, 20(2), 201-216. doi: <https://doi.org/10.2143/AST.20.2.2005443>
- Sangpoom, S. & Bodhisuwan, W. (2016). The discrete asymmetric Laplace distribution. *Journal of Statistical Theory and Practice*, 10(1), 73-86. <https://doi.org/10.1080/15598608.2015.1067659>