

*Journal of Korean  
Data & Information Science Society  
2002, Vol. 13, No.1 pp. 113 119*

## **Goodness-of-Fit Test for the Pareto Distribution Based on the Transformed Sample Lorenz curve<sup>1)</sup>**

**Suk-Bok Kang<sup>2)</sup> and Young-Suk Cho<sup>3)</sup>**

### **Abstract**

A powerful and easily computed goodness-of-fit test for Pareto distribution which does not depend on the unknown location and scale parameters is proposed based on the transformed sample Lorenz curve. We compare the power of the proposed test statistic with the other goodness-of-fit tests for Pareto distribution against various alternatives through Monte Carlo methods.

### **1. Introduction**

A continuous random variable  $X$  has the Pareto distribution with the location parameter  $a$ , the scale parameter  $b$ , and the shape parameter  $c$  if it has a cumulative distribution function (cdf) of the form

$$F(x) = 1 - [1 + (x - a)/b]^{-c}, \quad x \geq a, \quad b, c > 0. \quad (1.1)$$

The Pareto distribution is an important distribution in statistical analysis, income, wealth, and service time queueing system. Fisk (1961) cited several examples of economic data which follow the Pareto distribution. Berger and Mandelbrot (1963) used the Pareto in studies of error clusters in communication circuits. Harris (1968) found the Pareto to be useful in modeling service times and in queueing systems. Kaminsky and Nelson (1975) showed how the Pareto can be used in life testing, reliability, and replacement policy. Davis and Feldstein (1979) used the Pareto to model survival data based on times-to-failure of an observed

- 
1. This research was supported by the Yeungnam University research grants in 2001
  2. Professor, Department of Statistics, Yeungnam University, Kyongsan, 712-749, Korea
  3. Adjunct Assistant Professor, Department of Statistics, Yeungnam University, Kyongsan, 712-749, Korea

sample. Lawless (1983) stated that the 3-parameter Pareto distribution could have a decreasing or increasing hazard function depending on the parameters.

Likes (1969) derived the uniformly minimum variance unbiased estimators (UMVUE) of the parameters in the Pareto distribution. Malik (1970) derived distributions of the maximum likelihood estimators (MLE) of the parameters in the Pareto distribution. Kulldorff and Vannman (1973) studied estimation of the location and scale parameters of the Pareto distribution. Woo and Kang (1990) considered a more general class of UMVUE for the function of two parameters in the Pareto distribution. Kang and Cho (1996) obtained the jackknife estimator and the generalized jackknife estimator, the minimum risk estimator (MRE) of two parameters in the Pareto distribution.

Moothathu (1985) derived the MLEs of the Lorenz curve and the Gini index of a Pareto distribution, their exact and asymptotic distributions and moments. Moothathu (1990) also obtained the UMVUE and a strongly consistent asymptotically normal unbiased estimator (SCANUE) of the Lorenz curve, the Gini index and Theil entropy index of a Pareto distribution. Kang and Cho (1999) proposed the several estimators of the Lorenz curve in the Pareto distribution.

Use of the Pareto distribution for practical applications can be enhanced by an accurate method of determining whether a set of data comes from a population governed by the Pareto distribution. One class of goodness-of-fit tests that can be used for this purpose consists of tests based on the distance between the empirical distribution function (edf) and the hypothesized cdf. Three of the better known tests in this class Kolmogorov-Smirnov (K-S), Anderson-Darling (A-D), Cramer-von Mises (C-vM) are valid when there are no unknown parameters in the hypothesized distribution.

Lilliefors (1967, 1969) used Monte Carlo methods to construct tables for the modified K-S test when the parameters of a normal or an exponential distribution are estimated. Green and Hegazy (1976) constructed modified K-S, A-D, and C-vM critical value tables for the uniform, Laplace, Cauchy, and other distributions. Porter III, Coleman, and Moore (1992) modified K-S, A-D, and C-vM critical value tables for the Pareto distribution.

## 2. Goodness - of - fit Tests

Let  $X_{(j)}$  ( $j = 1, 2, \dots, n$ ) be the  $j$ -th order statistic based on a random sample  $X_1, X_2, \dots, X_n$  from the Pareto distribution with cdf (1.1). Consider now the case when shape parameter  $c$  is known and both the location parameter  $a$  and the scale parameter  $b$  are unknown. The best linear unbiased estimates (BLUEs)  $\hat{a}$  and  $\hat{b}$  were proposed by Kulldorff and Vannman (1973).

For  $c > 2$ , the BLUEs are

$$\hat{a} = x_{(1)} - \frac{Y}{(nc - 1)(c - 2) - ncD} \quad (2.1)$$

and

$$\hat{b} = (x_{(1)} - \hat{a})(nc - 1) \quad (2.2)$$

where

$$B_i = \left(1 - \frac{2}{c(n - i + 1)}\right) B_{i-1}, \text{ for } i = 1, 2, \dots, n,$$

$$B_0 \equiv 1,$$

$$D = (c + 1) \sum_{i=1}^{n-1} B_i + (c - 1)B_n,$$

$$Y = (c + 1) \sum_{i=1}^{n-1} B_i x_{(i)} + (c - 1)B_n x_{(n)} - D x_{(1)}.$$

For  $2/n < c \leq 2$  and  $2/c$  is an integer, the BLUEs are

$$\hat{a} = x_{(1)} - \frac{\hat{b}}{nc - 1} \quad (2.3)$$

and

$$\hat{b} = \frac{(c + 1)(c + 2)(nc - 1)}{(nc - 2)(nc - c - 2)} \left[ \sum_{i=1}^{n-2/c} B_i x_{(i)} - \frac{(nc - 2)}{(c + 2)} x_{(1)} \right] \quad (2.4)$$

The BLUEs were used to find the hypothesized cumulative distribution function  $P_i = F(x_{(i)}, \hat{a}, \hat{b}, c)$ , for  $i = 1, 2, \dots, n$ . Then the values of the three modified test statistics were calculated.

The K-S statistic was computed from :

$$D = \max \{D^+, D^-\},$$

$$D^+ = \sup_{1 \leq i \leq n} \left[ P_i - \frac{i-1}{n} \right],$$

$$D^- = \sup_{1 \leq i \leq n} \left[ \frac{i}{n} - P_i \right].$$

The A-D statistic was computed from

$$A^2 = -n - \frac{1}{n} \sum_{i=1}^n (2i-1)(\log P_i + \log(1 - P_{n+1-i})).$$

The C-vM was computed from

$$W^2 = \frac{1}{12n} + \sum_{i=1}^n \left( P_i - \frac{2i-1}{2n} \right)^2.$$

This procedure was repeated 5,000 times for each sample size  $n$ , each shape parameter  $c$  with  $a = b = 1$ , and for all three tests. Critical values are contained in table by Porter III, Coleman, and Moore (1992). The null hypothesis that a set of sample data follows a Pareto distribution with specified shape parameter  $c$  is rejected at the desired significance level if the calculated value of the test statistic exceeds the table value.

The Lorenz curve is extensively used in the study of inequality distribution and used to be a powerful tool for the analysis of a variety of scientific problems. The Lorenz curve is given by

$$L(y) = \int_0^y x dF(x) / E(Y)$$

where  $Y$  is a nonnegative income variable for which the mathematical expectation  $\mu = E(Y)$  exists.

Assume that  $X_1, X_2, \dots, X_n$  are positive random variables with order statistics  $X_{(1)} < \dots < X_{(n)}$ . Let  $r = [np]$  denote the greatest integer less than or equal to  $np$ . Then the sample Lorenz curve (Gail and Gastwirth (1978)) is defined by

$$L_n(p) = \frac{\sum_{i=1}^{r=[np]} X_{(i)}}{\sum_{i=1}^n X_{(i)}}.$$

Cho *et al.* (1999) proposed the transformed Lorenz curve that can be used in the study of symmetric distribution. The transformed Lorenz curve is defined by

$$TL(p) \equiv L(p) - p + 1.$$

To test  $H_0: X \sim F(x)$ , Kang and Cho (2001) proposed Normalized Sample Lorenz Curve (NSLC). The NSLC is defined by

$$NSLC(p) = \frac{TSL(p)}{TSL_F(p)}, \quad p = i/n, \quad i = 1, 2, \dots, n$$

where

$$TSL(p) = \frac{\sum_{j=1}^i (X_{j:n} - X_{1:n})}{\sum_{j=1}^n (X_{j:n} - X_{1:n})} - p + 1,$$

$$TSL_F(p) = \frac{\sum_{j=1}^i (F^{-1}(j/(n+1)) - F^{-1}(1/(n+1)))}{\sum_{j=1}^n (F^{-1}(j/(n+1)) - F^{-1}(1/(n+1)))} - p + 1.$$

We propose test statistic based on *NSLC* for the pareto distribution as follows.

$$TS = NSL C_{par}(0.5)$$

where

$$NSL C_{par}(p) = \frac{TSL(p)}{TSL_{par}(p)}, \quad p = i/n, \quad i = 1, 2, \dots, n$$

$$TSL(p) = \frac{\sum_{j=1}^i (X_{j:n} - X_{1:n})}{\sum_{j=1}^n (X_{j:n} - X_{1:n})} - p + 1,$$

$$TSL_{par}(p) = \frac{\sum_{j=1}^i ((1 - j/(n+1))^{1/c} - (1 - 1/(n+1))^{1/c})}{\sum_{j=1}^n ((1 - j/(n+1))^{1/c} - (1 - 1/(n+1))^{1/c})} - p + 1.$$

### 3. The Simulated Results

The exact distribution of the test statistic *TS* is hard to calculate. So, the critical values of the test statistic *TS* are obtained by repeating 5,000 times for sample size 25 and each shape parameter *c*. The *TS* critical value is 0.2051726(1.213659) for sample size 25, significance level  $\alpha = 0.05$ , and  $c = 1.0(3.5)$ . A power comparison was made among the K-S, A-D, C-vM, *TS* goodness-of-fit tests for three parameter Pareto distribution with only shape parameter specified. The power values were obtained by generating 5,000 random samples of size 25 for each alternative distributions for each tests. Table contains powers for hypothesized Pareto distribution shape parameters  $c = 1.0$  and 3.5. The proposed test statistic usually has greater power than the other test statistics. The test statistic *TS* provides a powerful and easily computed goodness-of-fit test for Pareto distribution which does not depend on the unknown location and scale parameters.

**Table 1.** Monte Carlo power estimates based on 5,000 samples of size  $n = 25$  using significance level  $\alpha = 0.05$  with  $c = 1.0$

	Exp(0,1)	U(0,1)	N(0,1)	Beta(2,2)	Wei(shape=3.5)
K-S	0.139	0.958	0.984	0.927	0.985
A-D	0.154	0.943	0.994	0.964	0.991
C-vM	0.165	0.920	0.992	0.956	0.989
<i>TS</i>	0.404	0.999	1.000	0.999	1.000

**Table 2.** Monte Carlo power estimates based on 5,000 samples of size  $n = 25$  using significance level  $\alpha = 0.05$  with  $c = 3.5$

	Exp(0,1)	U(0,1)	N(0,1)	Beta(2,2)	Wei(shape=3.5)
K - S	0.085	0.760	0.936	0.745	0.924
A - D	0.089	0.881	0.985	0.917	0.982
C - vM	0.095	0.856	0.983	0.920	0.978
TS	0.198	0.996	0.999	0.995	0.999

### References

- Berger, J. M. and Mandelbrot, B. (1963). A new model for error clustering in telephone circuits, *IBM J. Research & Development*, Vol. 7, 224-236.
- Cho, Y. S., Lee, J. Y., and Kang, S. B. (1999). A study on distribution based on the Transformed Lorez Curve. *The Korean Journal of Applied Statistics*, Vol. 12(1), 153-163.
- Davis, H. T. and Feldstein, M. L. (1979). The generalized Pareto law as a model for progressively censored survival data. *Biometrika*, Vol. 66, 299-306.
- Fisk, P. R. (1961). The graduation of income distributions, *Econometrica*, Vol. 29, 171-185.
- Gail, M. H. and Gastwirth, J. L. (1978). A Scale-Free Goodness-of-Fit test for the exponential distribution based on Lorenz curve. *Journal of American Statistical Association*, Vol. 73, 787-793.
- Green, J. and Hegazy. (1976). Powerful modified EDF goodness-of-fit tests. *Journal of American Statistical Association*, Vol. 71, 204-209.
- Harris, C. M. (1968). The Pareto distribution as a queue service discipline. *Operations Research*, Vol. 16 307-313.
- Kaminsky, K. S. and Nelson, P. I. (1975). Best liner unbiased prediction of order statistics in location and scale families, *Journal of American Statistical Association*, Vol. 70, 145-150.
- Kang, S. B. and Cho, Y. S. (1996). Estimation of the Parameters in a Pareto Distribution by Jackknife and Bootstrap Method. *Journal of Information & Optimization Sciences*, Vol., 18(2), 289-300.
- Kang, S. B. and Cho, Y. S. (2001). A study on Distribution Based on the Normalized Sample Lorenz Curve. *The Korean Communications in Statistics*, Vol. 8(1), 185-192.
- Kulldorff, G. and Vannman, K. (1973). Estimation of the Location and Scale Parameters of a Pareto Distribution by Linear Functions of Order

- Statistics, *Journal of American Statistical Association*, Vol. 68, 218-227.
12. Malik, H. J. (1970). Estimation of the Parameters of the Pareto Distribution, *Metrika*, Vol. 16, 126-132.
  13. Moothathu, T. S. K. (1985) Sampling Distribution of Lorenz Curve and Gini Index of the Pareto Distribution. *Sankhya*, Vol. 47(B), 247-278
  14. Moothathu, T. S. K. (1990) The Best Estimator of Lorenz Curve, Gini Index and Theil Entropy Index of Pareto Distribution. *Sankhya*, Vol. 52(B), 125-127
  15. Lawless, J. F. (1983). *Statistical Model for Life Data*, John Wiley & Sons.
  16. Likes, J. (1969). Minimum Variance Unbiased Estimation of the Parameters of power-function and Pareto's Distribution, *Statistische Hefte*, 10, 104-110.
  17. Lilliefors, H. (1967). On the Kolmogorov-Smirnov test for normality with mean and variance unknown, *Journal of American Statistical Association*, Vol. 62, 399-402.
  18. Lilliefors, H. (1969). On the Kolmogorov-Smirnov test for the exponential distribution with mean unknown, *Journal of American Statistical Association*, Vol. 64, 387-399.
  19. Porter III, J. E., Coleman, J. W. and Moore, A. H. (1992). Modified KS, AD, C-vM tests for the Pareto distribution with unknown location & scale parameters, *IEEE Transactions on Reliability*, Vol. 41(1), 112-117.
  20. Woo, J. and Kang, S. B. (1990). Estimation for Functions of Two Parameters in the Pareto Distribution, *Yongnam Statistical Letters*, Vol. 1, 67-76.