

Minimum sample size determination for generalized extreme value distribution

Yuzhi Cai, Dominic Hames

▶ To cite this version:

Yuzhi Cai, Dominic Hames. Minimum sample size determination for generalized extreme value distribution. Communications in Statistics - Simulation and Computation, 2010, 40 (01), pp.87-98. 10.1080/03610918.2010.530368 . hal-00651154

HAL Id: hal-00651154 https://hal.science/hal-00651154

Submitted on 13 Dec 2011

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Minimum sample size determination for generalized extreme value distribution

Journal:	Communications in Statistics - Simulation and Computation			
Manuscript ID:	LSSP-2010-0250			
Manuscript Type:	Original Paper			
Date Submitted by the Author:	07-Aug-2010			
Complete List of Authors:	Cai, Yuzhi; University of Plymouth, School of Computing and Mathematics Hames, Dominic; HR Wallingford Limited			
Keywords:	Bootstraping, Generalized extreme value distribution, Return level, Sample size			
Abstract:	Sample size determination is an important issue in statistical analysis. Obviously, the larger the sample size is, the better the statistical results we have. However, in many areas such as coastal engineering and environmental sciences, it can be very expensive or even impossible to collect large samples. In this paper, we propose a general method for determining the minimum sample size required by estimating the return levels from a generalized extreme value distribution. Both simulation studies and the applications to real data sets show that the method is easy to implement and the results obtained are very good.			
Note: The following files were submitted by the author for peer review, but cannot be converted				

Page 1 of 17



Minimum sample size determination for generalized extreme value distribution

^aYuzhi Cai*and ^bDominic Hames ^aUniversity of Plymouth, Plymouth, UK ^bHR Wallingford, Oxfordshire, UK

Abstract

Sample size determination is an important issue in statistical analysis. Obviously, the larger the sample size is, the better the statistical results we have. However, in many areas such as coastal engineering and environmental sciences, it can be very expensive or even impossible to collect large samples. In this paper, we propose a general method for determining the minimum sample size required by estimating the return levels from a generalized extreme value distribution. Both simulation studies and the applications to real data sets show that the method is easy to implement and the results obtained are very good.

Key words: Bootstraping, Generalized extreme value distribution, Return level, Sample size.

Introduction

Sample size is very important in statistical analysis. However, the sample size required for a study usually depends on different factors such as the probability models used and the types of statistical tests applied. Various methods have been proposed in the literature on sample size determination. For example, Shieh (2000) used the likelihood ratio test, Self and Mauritsen (1988) and Lubin and Gail (1990) applied the score test, and Bickel and Doksum (2001) and Demidenko (2007) considered the Wald test. Ashour and Shalaby (1983) derived some Bayesian and non-Bayesian estimators for sample size when the underlying distribution is a Weibull idstribution, and Ashour et al. (1996) also derived these estimators in the case of Burr type XII failure model. Abd-Elfattah and Bakoban (2003) obtained the estimators for sample size in case of a generalized gamma distribution, and

^{*}Address for correspondence: Dr Yuzhi Cai, School of Computing and Mathematics, University of Plymouth, Plymouth PL4 8AA, United Kingdom. Email: ycai@plymouth.ac.uk

 Abd-Elfattah et al. (2007) considered the Marcus and Blumenthal (1975) approach to estimate the sample size in the case when the distribution is the Lomax distribution, i.e. the Pareto distribution of the second kind.

The generalized extreme value (GEV) distribution has been used widely in many areas because the extreme value theory guarantees that the maximum group values will follow a GEV distribution if the group size tends to infinity. For example, in coastal engineering, we could assume that the maximum annual sea-levels will approximately follow the GEV distribution, and hence a GEV model is often fitted to the maximum annual data. In this paper, we study the sample size determination in the case when the underlying distribution is the GEV distribution, and we assume that the group size is large enough for us to use the GEV model. However, in reality limited data sets are available , or it can be very expensive or even impossible to collect a large number of such data. Therefore, it is very important to develop a method for estimating the accuracy of any study or the minimum sample size required for such a study.

Let x_i (i = 1, ..., n) be the maximum value of group *i*. It is well known that under certain regular conditions the maximum likelihood estimators (MLE) of the parameters are normally distributed as *n* approaches infinity. However, for any finite values of *n*, the distributions of the estimators are unknown, and hence the distributions of any functions of the estimators, such as return levels, are also unknown.

The main purpose of the paper is to develop a method for determining the minimum sample size based on the asymptotic distribution of MLEs. We will focus on the return level as it is a very important quantity in coastal engineering and environmental sciences and it provides crucial information in, for example, the design of flood defences. The developed methodology can be easily extended to other models and other quantities of interest. The arrangement of the paper is as follows. In Section 2, we develop the methodology for determining the sample size. Simulation studies are given in Section 3 and applications in Section 4. Some discussions and conclusions are given in Section 5.

2 The method

The GEV distribution function is defined by

$$F(x;\mu,\sigma,\xi) = e^{-\left[1+\xi\left(\frac{x-\mu}{\sigma}\right)\right]^{-1/\xi}} \tag{1}$$

for $1 + \xi \left(\frac{x-\mu}{\sigma}\right) > 0$, where $\mu \in R$ is the location parameter, $\sigma > 0$ the scale parameter and $\xi \in R$ the shape parameter. Let $\beta = (\mu, \sigma, \xi)$ be the parameter vector. Suppose we are interested in estimating the τ th quantile of the GEV, which is given by

$$x = \mu - \frac{\sigma}{\xi} [1 - (-\ln\tau)^{-\xi}],$$
(2)

where $0 < \tau < 1$. So, if $\tau = 0.99$, the 0.99th quantile may correspond (in an analysis of annual maxima) to the 100-year return level. We need to decide the minimum sample size so that the statistical inferences on the return level are proper. Note that in the following

we illustrate our method by using the 100-year return level, any other quantities of interest can be dealt with similarly.

The basic idea of the proposed method is given below. If the sample size n is large enough, then the maximum likelihood estimator (MLE) of β should be normally distributed, hence any functions of the MLE, such as the return level. As we only have one observed data set, it would be difficult to check the normality of the MLE of the return level. Therefore, we propose to apply the bootstrap method with replacement to the observed data set. Suppose we generated M bootstrap samples each of size n. Then we should expect that the MLEs of the return level based on the bootstrap samples should also be normally distributed. Standard bootstrap methodology guarantees that the MLEs form a random sample of a normal distribution for the return level. In this paper, we use the Shapiro-Wilk test to check the normality of the MLEs of the return level. If we have to reject the null hypothesis, i.e. the MLEs are normally distributed, then we may conclude that the sample size n is not large enough at a given significant level. So we need to collect more samples and to apply the above procedure to the enlarged sample. This process should be continued until we failed to reject the null hypothesis. The details of the method are given below.

Let $\mathbf{y}_1 = (y_{11}, \ldots, y_{1n})$ be the observed data.

Step 1. Obtain M - 1 bootstrap samples (with replacement) from y_1 , denoted by $y_m = (y_{m1}, \ldots, y_{mn})$, where $m = 2, \ldots, M$.

Step 2. For m = 1, ..., M, let $\mathbf{x}_{mj}^{(1)}$ be the first j values of \mathbf{y}_m . That is $\mathbf{x}_{mj}^{(1)} = (y_{m1}, ..., y_{mj})$, where $j = n_0, ..., n$ and n_0 is large enough for parameter estimation.

Step 3. For fixed j, obtain K - 1 bootstrap samples (with replacement) from $\mathbf{x}_{mj}^{(1)}$, denoted by $\mathbf{x}_{mj}^{(k)}$, where $k = 2, \dots, K$.

Step 4. Fit a GEV model to each data set $\mathbf{x}_{mj}^{(k)}$, and calculate the return level $\ell_{mj}^{(k)}$, $m = 1, \ldots, M, j = n_0, \ldots, n, k = 1, \ldots, K$.

Step 5. Carry out the Shapiro-Wilk normality test on $\ell_{mj}^{(k)}$ (k = 1, ..., K) and record the p-value, denoted by p_{mj} , of the test statistic, where m = 1, ..., M and $j = n_0, ..., n$.

Step 6. For $j = n_0, \ldots, n$ calculate

$$\bar{p}_j = \frac{1}{M} \sum_{m=1}^M p_{mj}, \quad s_j^2 = \frac{1}{M-1} \sum_{m=1}^M (p_{mj} - \bar{p}_j)^2.$$

and construct a $100(1-\alpha)\%$ confidence interval

 $\bar{p}_j \pm t^* s_j,$

where t^* is the critical value obtained from a t-distribution with M-1 degrees of freedom.

Step 7. Let N be the number such that for any j > N we have $\bar{p}_j - t^* s_j > \alpha$. As $\bar{p}_j - t^* s_j$ is the lower band of the confidence interval, we are $100(1 - \alpha)\%$ confident that the minimum number of samples required for obtaining a properly estimated return level is N.

 It is worth mentioning that the standard theory of maximum likelihood was developed for the case when the support of the distribution does not depend on unknown parameters, see, for example, Rao (1993). For GEV distribution, this is not the case because the support of the distribution must satisfy $1 + \xi(x - \mu)/\sigma > 0$. However, Smith (1985) studied this problem thoroughly and showed that (a) if $\xi > -0.5$, then the MLEs of the model parameters are asymptotically normally distributed; (b) if $-1 < \xi < -0.5$, then the MLEs can be obtained but are not asymptotically normally distributed; and (c) if $\xi < -1$, then it is usually not possible to obtain the MLEs. Note that when $\xi \leq -0.5$, the GEV distribution has a very short bounded right tail which, as Cole (2001) pointed out, is rarely encountered in applications of extreme value analysis. Therefore the theoretical limitations of the maximum likelihood method have little effect in practice. So in this paper we focus on the case $\xi > -0.5$. Indeed, all the real data sets considered in this paper have a reasonable long right tail, and hence the developed method can be safely applied.

To investigate the performance of the proposed method, we carried out extensive simulation studies which are given in the next section.

Simulation studies

Simulation study 1

Consider the GEV model given by

$$F(x;\mu,\sigma,\xi) = e^{-[1-0.11(x-100)/13]^{1/0.11}}$$
(3)

for 1 - 0.11(x - 100)/13 > 0. So, the true parameter values are: $\mu = 100$, $\sigma = 13$ and $\xi = -0.11$. Furthermore, the 100-year return level is given by

$$x = 100 + 13[1 - (-\ln 0.99)^{0.11}]/0.11 = 146.93,$$
(4)

Note that the true parameter values were chosen arbitrarily.

The statistical software R was used to generate a random sample of size n = 200from the above GEV model. In this simulation study, we take M = K = 50, $n_0 = 25$ and $\alpha = 0.05$. By applying the developed method to this data set, we obtained the results shown in Figure 1, where the lighter curve shows the average p-values \bar{p}_j against the sample size $j = 25, \ldots, 200$, the darker curves are the corresponding lower and upper limits of the 95% confidence interval, i.e. $\bar{p}_j \pm t^* s_j$ ($j = 25, \ldots, 200$). The horizontal line is at $\alpha = 0.05$, and the vertical line is at N = 40. It is seen that for $j \ge 40$ all curves are above the horizontal line, indicating that the minimum sample size we could take is N = 40 in this case. So if the data represent the maximum annual values of sea levels at a particular site, then 40 years' data would give a for good estimate of the return levels in terms of the normality of the MLE.

Now we compare the two fitted GEV models: one was fitted to the first 40 values of the simulated data, the other to all the data. The true and the estimated values are shown in Table 1. As expected, the standard errors of the MLEs based on the whole data sets are



Figure 1: Plot of the average p-values of the Shapiro-Wilk Normality test statistics together with a 95% confidence interval in the Simulation Study 1.

	True value	MLE (s.e.)	MLE (s.e.)			
		(n = 40)	(n = 200)			
μ	100.00	99.16 (2.13)	98.75 (0.95)			
σ	13.00	12.38 (1.44)	12.10 (0.67)			
ξ	-0.11	-0.10 (0.08)	-0.08(0.05)			
Return level	146.93	144.95 (8.02)	145.57 (4.44)			

Table 1: True and estimated parameter values together with their standard errors (in brackets).

 smaller than those based on the first 40 values. However, the estimated values are all very similar to the true parameter values. Note that the standard error of the estimated return level was obtained by using the delta method, and that a 95% confidence interval for each parameter can be constructed easily by using the estimated standard errors. The diagnostic plots in Figure 2 also shows that the two fitted models are very similar.



Figure 2: The probability plot and the return level plot of the GEV models fitted to the first N = 40 simulated data (left column), and to the whole data set (right column) respectively in the Simulation Study 1. On the return level plots, the top and the bottom curves give a 95% confidence bands for the return levels.

Simulation study 2.

In simulation study 1 we only dealt with one simulated data. In this simulation study, we repeat the above simulation study on 70 independently simulated data sets, each of size 200. So the average performance of the method can be assessed.

For each simulated data and for different sample sizes, we recorded the sample average p-values of the test statistic. Then the final average p-values are obtained by averaging over the 70 simulated data sets. Figure 3 shows the average p-values against the sample size and the corresponding 95% confidence interval. The lower limit of the confidence interval at a sample size N = 40 is 0.047, and for any sample size n > 40, all the curves are above the horizontal line. This suggests that on average, the performance of the method is very stable over different simulated data sets.



Figure 3: Average p-values over 70 different simulated data sets together with a 95% confidence interval in the Simulation Study 2.

4 Applications

In the applications of the developed method, we consider four real data sets obtained from Reiss and Thomas (2001). The first data set is the maximum annual sea-levels (in meters) in Venice from 1931 to 1981, containing 51 values. The mean sea level at Venice is 0.52m. In this application we consider the sea levels relative to the mean sea level. The second data set is the Iceland maximum annual wind speed (in meters per second) data in the years 1912 to 1992 with measurements for two years missing. We ignore the missing data in this application, hence the data set contains 79 values. The third data set is the maximum annual discharges (in cubic meters per second) of the Harricana River at Amos (Quebec, Canada) from 1915 to 1983, containing 69 values. The fourth data set is the maximum annual De Bilt temperatures (in Celsius) from 1849 to 1981, containing 133 values. The scatter plots of the four data sets are given in Figure 4.

All the four data sets are not long. We would like to estimate the minimum sample size required for making reasonable statistical inferences on the 100-year return level, and hence to see what confidence we can place in these estimates and whether further data are required.

Because the data are annual maximums, it is reasonable to fit a GEV model to each data set, hence the developed method can be used. For each of them, the minimum sample size was taken to be 25, because for smaller sample sizes we often have difficulties in fitting a GEV model to them. The largest sample size was taken to be the total number of observations in each data set. Furthermore, the number of bootstrap samples obtained from each data set and from its subsets with different sample sizes is 50. That is, we let M = K = 50. For larger values we have very similar results.



Figure 4: Scatter plots of (a) the maximum annual sea-levels in Venice, (b) the maximum annual wind speed in Iceland, (c) the maximum annual discharges of the Harricana River and (d) the maximum annual De Bilt temperature data.

Sample size	Sea-levels	Wind-speeds	Discharges	Temperatures
Return level 1	1.78 (0.113)	79.53 (2.79)	327.52 (31.100)	36.00 (0.889)
Return level 2	1.78 (0.110)	84.42 (3.13)	326.78 (23.454)	35.99 (0.358)

Table 2: Estimated return levels, where Return level 1 corresponds to the model fitted to data of minimum sample size, while the Return level 2 corresponds to the model fitted to all the data. Numbers in brackets are the corresponding standard errors.

Figure 5 shows the average p-values of the Shapiro-Wilk test statistic together with a 95% confidence interval for each data set, where the horizontal lines are at $\alpha = 0.05$, and the vertical lines are at 50, 72, 44 and 44 respectively, indicating the minimum sample size required in each case. Note that the total sample size of Venice sea-levels is only 51. Figure 5(a) shows it is difficult to make a sensible conclusion about the minimum sample size required in this case. So more data need to be collected. Figure 5(b) shows that all three curves are above the horizontal line for sample sizes greater than 71. This suggests that the minimum sample size can be taken as N = 72. Both Figure 5(c) and (d) suggest that the minimum sample size can be taken as N = 44. Therefore, except for the Venice sea-levels data set, all other data sets are large enough for making good statistical inferences on the return level.

To compare the model fitted to the data of minimum sample size with the model fitted to the whole data set, we produced Table 2 which shows that no significant differences at a 95% level between the two estimated return levels for each data set: one is based on the minimum sample sizes and the other is based on the whole data sets. This is because the corresponding 95% confidence intervals overlap. Note that the larger standard error value in the second row of Table 2 for the Iceland maximum annual wind speed data is caused by the approximation features of the delta method in calculating return levels for this data set.

Figure 6 shows the return level plots for the sea-level and wind-speed data sets, while Figure 7 is for the discharge and temperature data sets. In both figures, the left column corresponds to the minimum sample size, and the right column corresponds to the whole data set. These figures also show no significant differences in estimating return levels caused by using the minimum sample sizes.

A further statistical inferences about the sample size effects on the accuracy of the return levels can be carried out. We should expect that, generally, the standard errors of the MLEs of the return level decreases as the sample size increases. These standard errors tell us how much gain we will have if we increase the sample size. In other words, how much it would be worth paying to obtain a greater length of data. Figure 8 shows the average standard errors of the MLEs of the return level when the sample size increases from the minimum sample size up to another 300 data. Those 300 extra data were simulated from the fitted model by using the observed samples of minimum length. The average standard errors are based on 50 independently simulated samples each of size 300. From these plots, say, from Figure 8(b), we see that if we require the standard error of the MLE is less than 2.0, then we need at least 120 data instead of 72. That means that extra 48 data need to be collected in order to decrease the standard error of the MLE of the return level from 2.79



Figure 5: Average p-values and the corresponding 95% confidence interval for (a) the maximum annual sea-levels in Venice, (b) the maximum annual wind speed in Iceland, (c) the maximum annual discharges of the Harricana River and (d) the maximum annual De Bilt temperature data.



Figure 6: The return level plots of the GEV models fitted to the data of minimum sample sizes (left column), and the return level plots of the GEV models fitted to the whole data sets (right column). First row is for the Venice sea-level data, and the second row is for the Iceland wind-speed data.





Figure 7: The return level plots of the GEV models fitted to the data of minimum sample sizes (left column), and the return level plots of the GEV models fitted to the whole data sets (right column). First row is for the discharges of the Harricana River data, and the second row is for the De Bilt temperature data.

to 2.0. This would provide useful information for a company to decide whether new data are worth collecting.

5 Comments and conclusion

In this paper, we proposed a simple method to determine the minimum sample size required for making statistical inferences on a statistical quantity of interest. The developed methodology is based on the standard theory of maximum likelihood and we illustrated the method through the estimation of the 100-year return level by using a GEV model. It is worth mentioning that for other probability models and other quantities of interest, the method proposed in this paper can also be used, provided that the MLE is normally distributed in large samples. However, we would expect that the minimum sample sizes required by making statistical inferences on different quantities may be different. Simulation studies and application results show that the method can be easily implemented and the fitted models based on the minimum sample size are very similar to the fitted models based on all the available data.

In this paper we have also demonstrated different applications of the developed method. In summary, this method can be used (a) to check whether an existing data set is long enough for making good statistical inference, (b) to decide how many more data need to be collected for making good statistical inference and (c) to make sure how much gain we will have in terms of confidence if we increase the sample size. Therefore, we expect that the developed method can be very useful in practice.

Acknowledgement

We would like to express our sincere thanks to the referees for their very constructive comments and suggestions which have greatly enhanced the quality and the presentation of the paper.

References

- [1] Abd-Elfattah, AM., Alaboud, FM. and Alharby, AH. (2007). On sample size estimation for Lomax distribution. *Australian Journal of Basic and Applied Sciences* 1: 373-378.
- [2] Abd-Elfattah, AM. and Bakoban, RA. (2003). A study on the estimation of sample size for generalized Gamma distribution. *The Scientific Journal for Economic and Commerce, Faculty of Commerce, Ain Shames University* 2: 17-28.
- [3] Ashour, SK., Abd-Elfattah, AM. and Mahmoud, MR. (1996). Estimation of sample size for the Burr failure model. *Tha Annual Conference in Statisitcs, Computer Science and Operations Reserach, ISSR, Cairo University* 31: 18-26.
- [4] Ashoud, SK. and Shalaby, OA. (1983). Estimating sample size with Weibull failure. *Math. Operationsforsch. U. Statist. Ser. Statist., Berlin* 2: 263-268.

Page 15 of 17



Figure 8: The plots of the average standard errors of the MLEs of the return level against sample sizes for (a) the maximum annual sea-levels in Venice, (b) the maximum annual wind speed in Iceland, (c) the maximum annual discharges of the Harricana River and (d) the maximum annual De Bilt temperature data.

- [5] Bickel PJ. and Doksum, KA. (2001). *Mathematical Statistics* (2nd edn). Prentice-Hall: Upper Saddle River, NJ.
- [6] Coles, S. (2001). An introduction to statistical modelling of extreme values. Springer Series in Statistics.
- [7] Demidenko, E. (2007). Sample size determination for logistic regression revisited. *Statistics in Medicine* 26: 3385-3397.
- [8] Lubin JH. and Gail, MH. (1990). On power and sample size for studying features of the relative odds disease. *American Journal of Epedemiology* 131: 552-566.
- [9] Marcus, R. and Blumenthal, S. (1975). Estimating population size with exponential failure. *Journal of the American Statistical Association* 70: 913-922.
- [10] Rao, RC. (2002), *Linear Statistical Inference and Its Applications*, 2nd edition, New York: Wiley.
- [11] Reiss, RD. and Thomas M. (2001). *Statistical analysis of extreme values*. Birkhäuser Verlag, Basel-Boston-Berlin.
- [12] Smith, RL. (1985). Maximum likelihood estimation in a class of non-regular cases. *Biometrika*, 72, 67-90.
- [13] Shieh, G. (2000). On power and sample size calculations for likelihood ratio tests in generalized linear models. *Biometrics* 56: 1192-1196.
- [14] Self, SG. and Mauritsen, RH. (1988). Power/sample size calculations for generalized linear models. *Biometrics* 44: 79-86.

Responses to the referees' comments on the paper "Minimum sample size determination for generalized extreme value distribution"

Yuzhi Cai, University of Plymouth, UK. Dominic Hames, HR Wallingford, UK

We sincerely thank the referees for their very helpful comments and suggestions, according to which the paper has been modified. A separate paragraph has been added at the end of Section 2 to explain why the developed method can be used to GEV distributions. We also re-emphasized the asymptotic normality condition at the last section of the paper. An acknowledgement has also been added into the paper to thank the referees' comments which has improved the quality and the presentation of the paper significantly. The detailed responses are given below.

Referee's comments:

The standard theory of maximum likelihood is developed for the case when the distribution domain does not depend on unknown parameters, e.g. C.R.Rao (1993), "Linear Statistical Inference and Its Applications," New York: Wiley. Unfortunately, GEM, as defined in equation (1), does, namely, 1+ksi*(xmu)/sigma>0. The authors should provide theoretical justification why MLE is normally distributed in large sample.

Responses:

This is an excellent comment, many thanks! A separate paragraph has been added into the paper to discuss this problem. Generally speaking, Smith (1985) studied this problem thoroughly and showed that (a) kxi>-0.5, then the MLEs of the model

parameters are asymptotically normally distributed; (b) if -1<kxi<-0.5, then the MLEs can be obtained but are not asymptotically normally distributed; and (c) if kxi<-1, then it is usually not possible to obtain the MLEs. Note that when kxi

<= -0.5, the GEV distribution has a very short bounded right tail which, as Cole (2001) pointed out, is rarely encountered in applications of extreme value analysis. Therefore the theoretical limitations of the maximum likelihood method have little effect in practice. So in this paper we focus on the case kxi > - 0.5. Indeed, all the real data sets considered in this paper have a reasonable long right tail, and hence the developed method can be safely applied.

The last section of the paper has also been modified to re-emphasize the importance of the asymptotic normality condition.

The reference C.R.Rao (1993) has also been added into the paper.