# Comparison of methods for the estimation of statistical parameters of censored data

S. Kuttatharmmakul [a], D.L. Massart [a], D. Coomans [b], J. Smeyers-Verbeke [a],*

[a] *ChemoAC, Pharmaceutical Institute, Vrije Universiteit Brussel, Laarbeeklaan 103, B-1090 Brussel, Belgium*
[b] *School of Computer Science, Mathematics and Physics, James Cook University, Townsville Q4811, Australia*

## Abstract

Approaches based on the maximum likelihood (ML) method and on the order statistics are described and evaluated for the estimation of the mean and standard deviation of a normal population from a left-singly censored sample, i.e. a sample for which some measurement results fall below the reporting limit of the analytical method. The performance of the methods is evaluated by means of data simulations. The sample size considered is small to moderate: $N = 6$–$18$. Simulation data show that the ML method performs better than the method based on order statistics, especially in difficult situations, e.g. large expected censored proportion $h_{ex}$ ($h_{ex} \geq 50\%$) and for small sample size ($N = 6$). The reliability of the estimates depends on the censored proportion. The larger the censored proportion, the poorer the quality of the estimates. When the expected censored proportion does not exceed 50%, i.e. when the true mean $\mu$ of the measurement results is above the reporting limit, the performance of the ML method in the estimation of the mean of a censored sample is very acceptable, i.e. it is comparable to that using classical moment calculation on a complete (non-censored) sample. When the expected censored proportion is very high (e.g. 83%) the estimates are, as expected, largely biased. The performance of the ML method in the estimation of the standard deviation of censored data is not as good as in the estimation of the mean. A formula is given for the approximate sample size required to have a specified confidence level that a ML estimated mean for the censored sample will not differ from the true mean by a certain magnitude. © 2001 Elsevier Science B.V. All rights reserved.

*Keywords:* Maximum likelihood; Maximum likelihood estimation (MLE); Order statistics; Mean; Standard deviation; Censored data; Censored samples; Reporting limit; Sample size

## 1. Introduction

Most of the current literature on analysis of measurement data concerns unrestricted (i.e. complete) sample data. However, in some situations, restricted (incomplete) sample data, i.e. truncated data, censored data, are encountered [1,2]. Truncated data are data of which some (unknown amount of) data are missing. It is known that the missing data are those values below the lower limit and/or above the upper limit. Censored data are data of which some have no numerical values and are only known to be below a lower limit and/or above an upper limit. If the non-numerical data are the values below a lower limit, they are called left-censored data. If the non-numerical data are the values above an upper limit, they are right-censored. For analytical purposes, the left-censored data are

* Corresponding author. Tel.: +32-2477-4737;
fax: +32-2477-4735.
*E-mail address:* asmeyers@vub.vub.ac.be (J. Smeyers-Verbeke).

**Nomenclature**

| | |
|---|---|
| CI | confidence interval |
| $h$ | observed proportion of test results which are not known numerically, but only as '$<y_L$' |
| $h_{ex}$ | expected proportion of test results which are not known numerically, but only as '$<y_L$' |
| $n$ | observed number of test results that are numerically known, i.e. number of test results (in a dataset) that are above the reporting limit |
| $n_1$ | observed number of test results (in a dataset) for which no result except a '$<y_L$' is obtained |
| $n_{1ex}$ | expected number of test results (in a dataset) for which no result except a '$<y_L$' is obtained |
| $N$ | total number of test results in a dataset |
| $s^2$ | variance of the $n$ numerically known test results |
| $v$ | variance |
| $V$ | covariance matrix of normal order statistics |
| $\bar{y}$ | arithmetic mean of the $n$ numerically known test results |
| $y_i$ | individual test results |
| $y_{(i)}$ | sample order statistics |
| $y_L$ | reporting limit or lower limit |

*Greek letters*

| | |
|---|---|
| $\gamma$ | parameter used to find the associated constant $\lambda$ in the maximum likelihood estimation |
| $\eta$ | standardized reporting limit |
| $\lambda$ | constant derived from Tables 1 and 2 |
| $\mu$ | population mean of test results |
| $\mu_{11}$ | variance factor of $\hat{\mu}$ |
| $\sigma$ | population standard deviation |
| $\xi(i/N)$ | expected value of the $i$th (standard) normal order statistic for the sample size $N$ |

*Superscripts*

| | |
|---|---|
| $\wedge$ | superscript put on the parameters to represent the parameter estimates |

most important and the further discussion will focus on this type of data. Indeed, a trace analysis of which the measurements are carried out around the reporting limit of the analytical method and of which some measurement results fall below the reporting limit, yields a left-censored dataset. This situation is often encountered in, e.g. certification of reference materials [3], food and environmental analyses [4–7]. The problem is how to estimate the mean and standard deviation of a left-censored dataset. We can calculate the mean $\bar{y}$ and the standard deviation $\sigma$ of the numerically known data. However, $\bar{y}$ will overestimate the true mean $\mu$ and $\sigma$ will underestimate the true standard deviation $\sigma$ since the censored data with low values ($<$reporting limit) are not taken into account in the calculation. Several attempts have been made to consider these non-numerical data too in order to obtain more realistic estimates of the mean and standard deviation. Cohen [1,2] applied the maximum likelihood (ML) method to estimate the parameters of normal populations from truncated and censored samples. Gupta [8], Sarhan and Greenberg [9] employed the method based on order statistics to derive the coefficients of the best linear unbiased estimator (BLUE) for $\mu$ and $\sigma$ of normal populations from censored samples. Using also order statistics, Travis et al. [4] proposed a probability plot which estimates $\mu$ and $\sigma$ from the intercept and slope, respectively, of a least squares line fit to the data above the reporting limit. For non-normally distributed censored samples, robust methods [5] or data transformation combined with the iterative ML method via expectation-maximization algorithm [7,10] can be used. Statistical softwares were made available for the computation of parameter estimates based on censored samples [5,10,11].

In this paper, two methods for the analysis of a left-singly censored sample from a normal population have been extensively evaluated and compared by means of simulations, namely the ML method [1,2,12] and the method based on order statistics [4,6,8,9,13] (the methods can also be applied to log-normal samples by analyzing the logarithms of the data and then re-transform these back, however, for a lognormal distribution, the mean $\mu$ determines the median, i.e. antilog($\mu$) = median [4,5,11]). The sample size considered is small to moderate ($N = 6$–18). Various proportions of censoring, ranging between 17 and 83%, are studied in order to evaluate the degree

of censoring which still permits a reliable estimation of the mean and standard deviation. Though somewhat similar studies have been previously addressed [5–7], a different approach for the evaluation of the method performance is used in this work. Instead of measuring the method performance in terms of the estimated coverages of 90% confidence interval (CI) [7], the root mean squared error (RMSE), the bias or variance of the estimates [5,6], the distribution of estimates around the true value is considered. This gives us a good idea of how close an estimate is to its true value and how reliable an estimate is. In addition, by comparing the distributions of estimates obtained from censored samples with those using classical moment calculation on complete (non-censored) samples, the degree of censoring which still permits a reliable estimation of the statistical parameters could be evaluated (see Section 3).

## 2. Methods

To evaluate the reliability of the estimates for the mean and the standard deviation of a left-censored dataset, complete (non-censored) datasets and censored datasets, with a known true mean $\mu$ (of low value around the reporting limit) and a known true standard deviation $\sigma$, are generated (see Section 2.1). The estimates of the mean and the standard deviation are then calculated from those datasets (see Section 2.2) and their distributions are considered (see Section 2.3). The efficiency of the methods is determined by comparing the distributions of the estimates obtained for censored datasets with the distributions of the estimates obtained in the classical way with the method of moments (see Eqs. (2) and (3) in Section 2.2) for the complete (non-censored) datasets (see Section 3).

### 2.1. Data simulations

#### 2.1.1. Parameters applied in the data simulations

The situation where an analysis is carried out around the reporting limit is taken as a case study. Other situations where censored data arise can be found in [1]. The parameters for the simulations are as follows:

Sample size:

$N = 6, 12, 18$

For each $N$, the following means and standard deviations are considered:

Population mean:

$$\mu = \begin{cases} 1.19, 1.09, 1.00, 0.91, 0.87, 0.81 & (\text{for } \sigma = 0.20) \\ 1.34, 1.15, 1.00, 0.85, 0.76, 0.66 & (\text{for } \sigma = 0.35) \end{cases}$$

Population standard deviation:

$\sigma = 0.20, 0.35$

Reporting limit:

$y_L = 1$ for all censored cases

the units for $\mu$, $\sigma$, and $y_L$ are arbitrary.

The different means $\mu$ (calculated from Eq. (1)) considered in the simulations are determined in such a way that the probabilities to obtain a test result below a reporting limit of $y_L = 1$ equal the expected censored proportions $h_{ex}$. The values of $h_{ex}$ considered are 17, 33, 50, 67, 75 and 83%. Thus, $n_{1ex} = 0.01 h_{ex} N$, is the expected number of test results obtained from a sample of size $N$ that are reported as below the reporting limit, e.g. with $N = 12$ and $h_{ex} = 17, 33, 50, 67, 75$ and 83%, $n_{1ex}$ are 2, 4, 6, 8, 9 and 10, respectively.

Use is made of the relationship:

$$\mu = y_L + z_{(100-h_{ex})}\sigma \tag{1}$$

where $z_{(100-h_{ex})}$ is the $(100 - h_{ex})$th percentile of the standard normal distribution, i.e. the value below which a standard normal random variable falls with a probability $(100 - h_{ex})$% with $h_{ex} = 17, 33, 50, 67, 75$ and 83%.

#### 2.1.2. Complete (non-censored) datasets

For each situation considered, 10,000 datasets, each consisting of $N$ random test results normally distributed around the true mean $\mu$ with the standard deviation $\sigma$, are generated by the RANDN function in Matlab 4.0 [14]. Note that a complete (non-censored) dataset may contain values below the reporting limit.

#### 2.1.3. Singly censored datasets

In what follows, $N$ is the sample size (i.e. the total number of test results), $n_1$ is the number of test results which are censored since they are below the reporting limit $y_L$ and $n$ is the number of remaining test results after data censoring ($N = n_1 + n$).
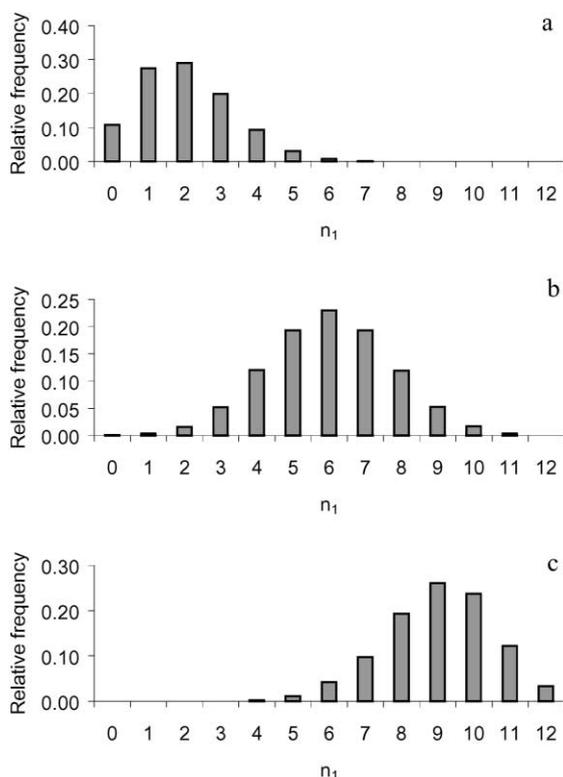
Fig. 1. Distributions of the observed number of censored test results, $n_1$, for different expected censor proportions $h_{ex}$ (sample size $N = 12$). Relative frequency is obtained as the proportion of datasets with the observed $n_1$ ($h_{ex}$: expected censored proportion, $n_{1ex}$: expected number of censored test results). (a) $h_{ex} = 17\%$ ($n_{1ex} = 2$). (b) $h_{ex} = 50\%$ ($n_{1ex} = 6$). (c) $h_{ex} = 75\%$ ($n_{1ex} = 9$).

The censored datasets were derived from the complete (non-censored) datasets (generated as described in Section 2.1.2) by leaving out those test results which are smaller than the reporting limit. The number of test results that were left out from each dataset, $n_1$, is not fixed but random. As an example, distributions of $n_1$, associated with (>10,000) complete (non-censored) datasets of $N = 12$ from which 10,000 censored datasets were originated, are represented in Fig. 1. The datasets with $n_1 = 0$, $N - 1$ and $N$ are discarded since (i) the data are not censored when $n_1 = 0$, and (ii) the analysis of censored data required at least two numerically known data, $n_1 \leq N - 2$ (see Section 2.2). Because some datasets have to be discarded, more than 10,000 complete (non-censored) datasets were simulated in order to obtain 10,000 censored datasets, i.e.

the simulations of complete datasets were iterated till the number of censored datasets (with $1 \leq n_1 \leq N - 2$) reached 10,000.

### 2.2. Calculation of the estimates $\hat{\mu}$ and $\hat{\sigma}$

For the complete (non-censored) datasets, the estimates $\hat{\mu}$ and $\hat{\sigma}$ are calculated in the classical way (applying the method of moments) from

$$\hat{\mu} = \frac{\sum_{i=1}^{N} y_i}{N} \tag{2}$$

$$\hat{\sigma} = \sqrt{\frac{\sum_{i=1}^{N} (y_i - \hat{\mu})^2}{N - 1}} \tag{3}$$

For the censored datasets, the estimates $\hat{\mu}$ and $\hat{\sigma}$ are calculated using the ML method and the method based on order statistics.

#### 2.2.1. Maximum likelihood method
For each censored dataset, the ML estimates of the mean and the standard deviation are obtained as follows [12].

1. Calculate the mean $\bar{y}$ and the variance $s^2$ of the $n$ numerically known data as

$$\bar{y} = \frac{\sum_{i=1}^{n} y_i}{n} \tag{4}$$

$$s^2 = \frac{\sum_{i=1}^{n} (y_i - \bar{y})^2}{n - 1} \tag{5}$$

2. Calculate the proportion of censored data $h$ and the parameter $\hat{\gamma}$ as

$$h = \frac{n_1}{N} \tag{6}$$

$$\hat{\gamma} = \frac{s^2}{(\bar{y} - y_L)^2} \tag{7}$$

3. Use the $h$ and $\hat{\gamma}$ (obtained from Eqs. (6) and (7)) to find the associated constant $\lambda(h, \hat{\gamma})$ from Tables 1 and 2, apply linear interpolation if the values of $h$ and $\hat{\gamma}$ are not the same as those in Tables 1 and 2.

4. Calculate the ML estimates $\hat{\mu}$ and $\hat{\sigma}$ as

$$\hat{\mu} = \bar{y} - \lambda(\bar{y} - y_L) \tag{8}$$

$$\hat{\sigma} = \sqrt{s^2 + \lambda(\bar{y} - y_L)^2} \tag{9}$$

Table 1
$\lambda(h, \gamma)$ for singly censored samples from the normal distribution[a]

| γ | h (%) | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 12 | 14 | 16 | 18 | 20 |
| 0.00 | 0.010103 | 0.020400 | 0.030902 | 0.041606 | 0.052515 | 0.063630 | 0.074953 | 0.086488 | 0.098237 | 0.110204 | 0.13480 | 0.16031 | 0.18677 | 0.21421 | 0.24268 |
| 0.04 | 0.010467 | 0.021125 | 0.031975 | 0.043019 | 0.054260 | 0.065701 | 0.077345 | 0.089194 | 0.101253 | 0.113524 | 0.13872 | 0.16481 | 0.19184 | 0.21983 | 0.24885 |
| 0.08 | 0.010796 | 0.021777 | 0.032944 | 0.044300 | 0.055848 | 0.067592 | 0.079534 | 0.091678 | 0.104027 | 0.116586 | 0.14235 | 0.16899 | 0.19657 | 0.22510 | 0.25464 |
| 0.12 | 0.011098 | 0.022376 | 0.033836 | 0.045482 | 0.057317 | 0.069344 | 0.081567 | 0.093989 | 0.106613 | 0.119445 | 0.14574 | 0.17292 | 0.20102 | 0.23007 | 0.26012 |
| 0.16 | 0.011378 | 0.022934 | 0.034669 | 0.046588 | 0.058692 | 0.070987 | 0.083475 | 0.096160 | 0.109047 | 0.122138 | 0.14895 | 0.17664 | 0.20524 | 0.23478 | 0.26533 |
| 0.20 | 0.011642 | 0.023459 | 0.035453 | 0.047629 | 0.059990 | 0.072539 | 0.085280 | 0.098216 | 0.111352 | 0.124692 | 0.15200 | 0.18018 | 0.20926 | 0.23929 | 0.27031 |
| 0.24 | 0.011891 | 0.023956 | 0.036197 | 0.048618 | 0.061223 | 0.074014 | 0.086996 | 0.100173 | 0.113549 | 0.127127 | 0.15491 | 0.18356 | 0.21311 | 0.24360 | 0.27509 |
| 0.28 | 0.012129 | 0.024429 | 0.036906 | 0.049561 | 0.062399 | 0.075423 | 0.088638 | 0.102046 | 0.115651 | 0.129459 | 0.15770 | 0.18680 | 0.21681 | 0.24775 | 0.27969 |
| 0.32 | 0.012356 | 0.024883 | 0.037585 | 0.050465 | 0.063527 | 0.076775 | 0.090213 | 0.103843 | 0.117671 | 0.131701 | 0.16038 | 0.18993 | 0.22037 | 0.25176 | 0.28414 |
| 0.36 | 0.012574 | 0.025318 | 0.038237 | 0.051334 | 0.064613 | 0.078076 | 0.091729 | 0.105575 | 0.119618 | 0.133862 | 0.16297 | 0.19294 | 0.22382 | 0.25563 | 0.28844 |
| 0.40 | 0.012784 | 0.025738 | 0.038866 | 0.052173 | 0.065660 | 0.079332 | 0.093193 | 0.107247 | 0.121498 | 0.135951 | 0.16548 | 0.19586 | 0.22715 | 0.25938 | 0.29260 |
| 0.44 | 0.012987 | 0.026144 | 0.039475 | 0.052983 | 0.066673 | 0.080547 | 0.094610 | 0.108867 | 0.123320 | 0.137974 | 0.16790 | 0.19870 | 0.23039 | 0.26303 | 0.29665 |
| 0.48 | 0.013183 | 0.026536 | 0.040064 | 0.053769 | 0.067654 | 0.081725 | 0.095985 | 0.110437 | 0.125087 | 0.139938 | 0.17026 | 0.20145 | 0.23354 | 0.26657 | 0.30059 |
| 0.50 | 0.013279 | 0.026728 | 0.040352 | 0.054153 | 0.068135 | 0.082301 | 0.096657 | 0.111206 | 0.125952 | 0.140899 | 0.17142 | 0.20280 | 0.23508 | 0.26831 | 0.30253 |
| 0.60 | 0.013739 | 0.027649 | 0.041733 | 0.055995 | 0.070439 | 0.085068 | 0.099887 | 0.114899 | 0.130109 | 0.145521 | 0.17697 | 0.20928 | 0.24250 | 0.27667 | 0.31184 |
| 0.70 | 0.014171 | 0.028513 | 0.043030 | 0.057726 | 0.072605 | 0.087670 | 0.102925 | 0.118374 | 0.134022 | 0.149873 | 0.18220 | 0.21540 | 0.24951 | 0.28458 | 0.32065 |
| 0.80 | 0.014579 | 0.029330 | 0.044258 | 0.059364 | 0.074655 | 0.090133 | 0.105802 | 0.121666 | 0.137731 | 0.154000 | 0.18717 | 0.22121 | 0.25617 | 0.29210 | 0.32903 |
| 0.90 | 0.014967 | 0.030107 | 0.045425 | 0.060923 | 0.076606 | 0.092477 | 0.108541 | 0.124802 | 0.141264 | 0.157932 | 0.19190 | 0.22676 | 0.26253 | 0.29927 | 0.33703 |
| 1.00 | 0.015338 | 0.030850 | 0.046540 | 0.062413 | 0.078471 | 0.094720 | 0.111162 | 0.127803 | 0.144646 | 0.161696 | 0.19644 | 0.23206 | 0.26862 | 0.30615 | 0.34471 |
| 1.10 | 0.015693 | 0.031562 | 0.047611 | 0.063843 | 0.080262 | 0.096872 | 0.113678 | 0.130684 | 0.147893 | 0.165312 | 0.20079 | 0.23717 | 0.27448 | 0.31277 | 0.35209 |
| 1.20 | 0.016035 | 0.032247 | 0.048641 | 0.065219 | 0.081986 | 0.098945 | 0.116101 | 0.133459 | 0.151022 | 0.168796 | 0.20499 | 0.24209 | 0.28013 | 0.31916 | 0.35922 |
| 1.30 | 0.016365 | 0.032909 | 0.049635 | 0.066547 | 0.083650 | 0.100946 | 0.118441 | 0.136139 | 0.154044 | 0.172161 | 0.20905 | 0.24684 | 0.28559 | 0.32533 | 0.36612 |
| 1.40 | 0.016684 | 0.033549 | 0.050597 | 0.067833 | 0.085260 | 0.102883 | 0.120706 | 0.138734 | 0.156970 | 0.175419 | 0.21298 | 0.25145 | 0.29088 | 0.33131 | 0.37280 |
| 1.50 | 0.016994 | 0.034169 | 0.051529 | 0.069079 | 0.086822 | 0.104761 | 0.122903 | 0.141250 | 0.159808 | 0.178581 | 0.21679 | 0.25592 | 0.29602 | 0.33712 | 0.37929 |
| 1.60 | 0.017294 | 0.034771 | 0.052435 | 0.070289 | 0.088338 | 0.106586 | 0.125037 | 0.143695 | 0.162566 | 0.181653 | 0.22050 | 0.26027 | 0.30101 | 0.34277 | 0.38561 |
| 1.70 | 0.017586 | 0.035357 | 0.053316 | 0.071467 | 0.089814 | 0.108362 | 0.127114 | 0.146075 | 0.165250 | 0.184643 | 0.22410 | 0.26450 | 0.30588 | 0.34828 | 0.39176 |
| 1.80 | 0.017871 | 0.035928 | 0.054175 | 0.072615 | 0.091253 | 0.110092 | 0.129138 | 0.148395 | 0.167866 | 0.187558 | 0.22762 | 0.26863 | 0.31062 | 0.35365 | 0.39776 |
| 1.90 | 0.018149 | 0.036485 | 0.055012 | 0.073734 | 0.092656 | 0.111781 | 0.131113 | 0.150658 | 0.170419 | 0.190403 | 0.23105 | 0.27266 | 0.31525 | 0.35889 | 0.40362 |
| 2.00 | 0.018420 | 0.037029 | 0.055830 | 0.074828 | 0.094026 | 0.113430 | 0.133042 | 0.152869 | 0.172914 | 0.193182 | 0.23441 | 0.27659 | 0.31978 | 0.36401 | 0.40934 |
| 2.20 | 0.018945 | 0.038081 | 0.057413 | 0.076944 | 0.096678 | 0.116621 | 0.136776 | 0.157148 | 0.177742 | 0.198563 | 0.24090 | 0.28421 | 0.32854 | 0.37393 | 0.42044 |
| 2.40 | 0.019448 | 0.039090 | 0.058931 | 0.078974 | 0.099224 | 0.119684 | 0.140360 | 0.161256 | 0.182377 | 0.203728 | 0.24714 | 0.29154 | 0.33696 | 0.38347 | 0.43110 |
| 2.60 | 0.019932 | 0.040062 | 0.060393 | 0.080928 | 0.101674 | 0.122633 | 0.143811 | 0.165212 | 0.186841 | 0.208703 | 0.25315 | 0.29859 | 0.34507 | 0.39265 | 0.44138 |
| 2.80 | 0.020400 | 0.041000 | 0.061803 | 0.082815 | 0.104039 | 0.125480 | 0.147142 | 0.169031 | 0.191151 | 0.213507 | 0.25895 | 0.30540 | 0.35291 | 0.40153 | 0.45131 |
| 3.00 | 0.020852 | 0.041907 | 0.063168 | 0.084640 | 0.106328 | 0.128235 | 0.150366 | 0.172727 | 0.195322 | 0.218156 | 0.26456 | 0.31199 | 0.36050 | 0.41012 | 0.46092 |

[a] h-range:1–20% [1,2,12].

Table 2
$\lambda(h, \gamma)$ for singly censored samples from the normal distribution[a]

| $\gamma$ | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| h (%) | 22 | 24 | 26 | 28 | 30 | 32 | 34 | 36 | 38 | 40 | 45 | 50 | 55 | 60 | 65 | 70 | 80 | 90 |
| 0.00 | 0.27221 | 0.30286 | 0.33469 | 0.36774 | 0.40210 | 0.43782 | 0.47499 | 0.51369 | 0.55401 | 0.59607 | 0.70957 | 0.83684 | 0.9808 | 1.1454 | 1.3362 | 1.5615 | 2.1759 | 3.2826 |
| 0.04 | 0.27892 | 0.31010 | 0.34245 | 0.37602 | 0.41089 | 0.44712 | 0.48479 | 0.52399 | 0.56481 | 0.60736 | 0.72210 | 0.85060 | 0.9958 | 1.1616 | 1.3537 | 1.5803 | 2.1977 | 3.3079 |
| 0.08 | 0.28523 | 0.31693 | 0.34979 | 0.38387 | 0.41924 | 0.45596 | 0.49413 | 0.53382 | 0.57514 | 0.61818 | 0.73414 | 0.86386 | 1.0103 | 1.1774 | 1.3708 | 1.5987 | 2.2190 | 3.3328 |
| 0.12 | 0.29122 | 0.32342 | 0.35677 | 0.39135 | 0.42720 | 0.46442 | 0.50307 | 0.54325 | 0.58505 | 0.62858 | 0.74575 | 0.87669 | 1.0243 | 1.1927 | 1.3874 | 1.6167 | 2.2399 | 3.3574 |
| 0.16 | 0.29692 | 0.32961 | 0.36345 | 0.39850 | 0.43484 | 0.47254 | 0.51167 | 0.55232 | 0.59460 | 0.63860 | 0.75697 | 0.88912 | 1.0380 | 1.2076 | 1.4036 | 1.6342 | 2.2604 | 3.3815 |
| 0.20 | 0.30238 | 0.33554 | 0.36985 | 0.40538 | 0.44219 | 0.48035 | 0.51995 | 0.56108 | 0.60382 | 0.64829 | 0.76784 | 0.90118 | 1.0513 | 1.2221 | 1.4194 | 1.6514 | 2.2805 | 3.4054 |
| 0.24 | 0.30762 | 0.34124 | 0.37602 | 0.41200 | 0.44927 | 0.48789 | 0.52795 | 0.56954 | 0.61275 | 0.65768 | 0.77839 | 0.91292 | 1.0642 | 1.2363 | 1.4349 | 1.6682 | 2.3003 | 3.4289 |
| 0.28 | 0.31267 | 0.34674 | 0.38197 | 0.41840 | 0.45612 | 0.49519 | 0.53570 | 0.57774 | 0.62140 | 0.66679 | 0.78865 | 0.92434 | 1.0768 | 1.2502 | 1.4500 | 1.6847 | 2.3197 | 3.4521 |
| 0.32 | 0.31755 | 0.35206 | 0.38772 | 0.42460 | 0.46276 | 0.50227 | 0.54322 | 0.58570 | 0.62981 | 0.67565 | 0.79864 | 0.93549 | 1.0892 | 1.2637 | 1.4648 | 1.7008 | 2.3388 | 3.4750 |
| 0.36 | 0.32228 | 0.35722 | 0.39331 | 0.43061 | 0.46920 | 0.50915 | 0.55053 | 0.59345 | 0.63799 | 0.68427 | 0.80838 | 0.94637 | 1.1012 | 1.2770 | 1.4793 | 1.7167 | 2.3576 | 3.4976 |
| 0.40 | 0.32687 | 0.36222 | 0.39873 | 0.43646 | 0.47547 | 0.51584 | 0.55765 | 0.60099 | 0.64597 | 0.69268 | 0.81789 | 0.95700 | 1.1130 | 1.2900 | 1.4936 | 1.7322 | 2.3762 | 3.5199 |
| 0.44 | 0.33132 | 0.36709 | 0.40401 | 0.44214 | 0.48157 | 0.52236 | 0.56459 | 0.60835 | 0.65375 | 0.70089 | 0.82719 | 0.96741 | 1.1246 | 1.3027 | 1.5075 | 1.7475 | 2.3944 | 3.5420 |
| 0.48 | 0.33566 | 0.37183 | 0.40915 | 0.44769 | 0.48752 | 0.52872 | 0.57136 | 0.61554 | 0.66136 | 0.70892 | 0.83628 | 0.97760 | 1.1359 | 1.3152 | 1.5213 | 1.7626 | 2.4124 | 3.5638 |
| 0.50 | 0.33779 | 0.37415 | 0.41167 | 0.45041 | 0.49044 | 0.53184 | 0.57469 | 0.61907 | 0.66509 | 0.71286 | 0.84075 | 0.98262 | 1.1415 | 1.3214 | 1.5281 | 1.7700 | 2.4213 | 3.5745 |
| 0.60 | 0.34806 | 0.38538 | 0.42386 | 0.46357 | 0.50457 | 0.54695 | 0.59079 | 0.63617 | 0.68320 | 0.73199 | 0.86247 | 1.00700 | 1.1686 | 1.3514 | 1.5611 | 1.8063 | 2.4649 | 3.6276 |
| 0.70 | 0.35777 | 0.39600 | 0.43541 | 0.47604 | 0.51798 | 0.56131 | 0.60609 | 0.65244 | 0.70044 | 0.75021 | 0.88320 | 1.03032 | 1.1946 | 1.3802 | 1.5929 | 1.8413 | 2.5070 | 3.6793 |
| 0.80 | 0.36702 | 0.40612 | 0.44641 | 0.48794 | 0.53078 | 0.57501 | 0.62071 | 0.66799 | 0.71693 | 0.76764 | 0.90306 | 1.05270 | 1.2196 | 1.4080 | 1.6236 | 1.8751 | 2.5480 | 3.7296 |
| 0.90 | 0.37585 | 0.41580 | 0.45694 | 0.49932 | 0.54303 | 0.58814 | 0.63473 | 0.68290 | 0.73275 | 0.78439 | 0.92216 | 1.07425 | 1.2437 | 1.4348 | 1.6533 | 1.9078 | 2.5877 | 3.7787 |
| 1.00 | 0.38434 | 0.42510 | 0.46705 | 0.51027 | 0.55481 | 0.60077 | 0.64822 | 0.69725 | 0.74798 | 0.80051 | 0.94058 | 1.09506 | 1.2670 | 1.4607 | 1.6820 | 1.9395 | 2.6264 | 3.8267 |
| 1.10 | 0.39250 | 0.43405 | 0.47679 | 0.52081 | 0.56617 | 0.61294 | 0.66123 | 0.71111 | 0.76269 | 0.81609 | 0.95838 | 1.11519 | 1.2896 | 1.4859 | 1.7099 | 1.9704 | 2.6640 | 3.8736 |
| 1.20 | 0.40038 | 0.44269 | 0.48620 | 0.53100 | 0.57714 | 0.62472 | 0.67381 | 0.72451 | 0.77693 | 0.83117 | 0.97563 | 1.13471 | 1.3115 | 1.5103 | 1.7370 | 2.0004 | 2.7008 | 3.9196 |
| 1.30 | 0.40800 | 0.45105 | 0.49531 | 0.54086 | 0.58777 | 0.63612 | 0.68600 | 0.73750 | 0.79073 | 0.84579 | 0.99237 | 1.15367 | 1.3328 | 1.5341 | 1.7634 | 2.0296 | 2.7367 | 3.9646 |
| 1.40 | 0.41540 | 0.45916 | 0.50415 | 0.55043 | 0.59809 | 0.64720 | 0.69784 | 0.75012 | 0.80413 | 0.86000 | 1.00865 | 1.17211 | 1.3535 | 1.5572 | 1.7892 | 2.0582 | 2.7718 | 4.0087 |
| 1.50 | 0.42258 | 0.46704 | 0.51273 | 0.55974 | 0.60812 | 0.65796 | 0.70935 | 0.76239 | 0.81718 | 0.87383 | 1.02451 | 1.19009 | 1.3737 | 1.5798 | 1.8143 | 2.0861 | 2.8062 | 4.0520 |
| 1.60 | 0.42956 | 0.47470 | 0.52109 | 0.56879 | 0.61788 | 0.66845 | 0.72057 | 0.77435 | 0.82989 | 0.88731 | 1.03996 | 1.20762 | 1.3935 | 1.6019 | 1.8389 | 2.1134 | 2.8399 | 4.0945 |
| 1.70 | 0.43637 | 0.48217 | 0.52923 | 0.57762 | 0.62740 | 0.67867 | 0.73150 | 0.78601 | 0.84229 | 0.90046 | 1.05505 | 1.22475 | 1.4128 | 1.6234 | 1.8629 | 2.1401 | 2.8729 | 4.1363 |
| 1.80 | 0.44301 | 0.48946 | 0.53718 | 0.58623 | 0.63669 | 0.68865 | 0.74218 | 0.79740 | 0.85440 | 0.91331 | 1.06980 | 1.24149 | 1.4316 | 1.6446 | 1.8864 | 2.1663 | 2.9053 | 4.1773 |
| 1.90 | 0.44950 | 0.49658 | 0.54494 | 0.59465 | 0.64578 | 0.69840 | 0.75262 | 0.80854 | 0.86625 | 0.92588 | 1.08423 | 1.25788 | 1.4501 | 1.6652 | 1.9095 | 2.1919 | 2.9371 | 4.2177 |
| 2.00 | 0.45584 | 0.50355 | 0.55254 | 0.60288 | 0.65466 | 0.70795 | 0.76284 | 0.81943 | 0.87784 | 0.93818 | 1.09836 | 1.27394 | 1.4682 | 1.6855 | 1.9321 | 2.2171 | 2.9683 | 4.2575 |
| 2.20 | 0.46812 | 0.51704 | 0.56725 | 0.61884 | 0.67188 | 0.72645 | 0.78265 | 0.84057 | 0.90033 | 0.96204 | 1.12579 | 1.30512 | 1.5034 | 1.7249 | 1.9762 | 2.2662 | 3.0292 | 4.3352 |
| 2.40 | 0.47993 | 0.53001 | 0.58140 | 0.63419 | 0.68845 | 0.74426 | 0.80171 | 0.86092 | 0.92198 | 0.98502 | 1.15221 | 1.33518 | 1.5373 | 1.7630 | 2.0187 | 2.3136 | 3.0883 | 4.4107 |
| 2.60 | 0.49131 | 0.54251 | 0.59505 | 0.64899 | 0.70442 | 0.76143 | 0.82011 | 0.88055 | 0.94287 | 1.00721 | 1.17773 | 1.36423 | 1.5701 | 1.7998 | 2.0598 | 2.3595 | 3.1455 | 4.4841 |
| 2.80 | 0.50231 | 0.55460 | 0.60824 | 0.66330 | 0.71988 | 0.77804 | 0.83790 | 0.89954 | 0.96309 | 1.02867 | 1.20244 | 1.39237 | 1.6019 | 1.8355 | 2.0998 | 2.4041 | 3.2012 | 4.5557 |
| 3.00 | 0.51296 | 0.56630 | 0.62101 | 0.67717 | 0.73485 | 0.79414 | 0.85514 | 0.91795 | 0.98270 | 1.04949 | 1.22641 | 1.41967 | 1.6327 | 1.8702 | 2.1386 | 2.4474 | 3.2554 | 4.6254 |

[a] h-range: 22–90% [1,2,12].

We have observed from data simulations that the estimated gamma ($\hat{\gamma}$) obtained from Eq. (7) can be larger than 1. This can be the case when $N$ is small and/or the censored proportion is large ($h > 50\%$). However, $\hat{\gamma}$ larger than 3 is rarely observed. For example, from data simulations of 10,000 censored samples with $h_{ex} = 50\%$ and $N = 6$, 17 and 0.03% of censored samples are found to have $\hat{\gamma}$ larger than 1 and 3, respectively. Therefore, the tables of constant $\lambda$ in [12], which are limited to entries for $\gamma \leq 1$, have been enlarged to include the values of $\lambda(h, \gamma)$ for $\gamma$ ranging between 0 and 3 (Tables 1 and 2). The original tables given by Cohen [1,2] are also limited to $\gamma \leq 1$.

### 2.2.2. Method based on order statistics

The method makes use of the expected values of order statistics, also called rankits [15,16]. The approach is more commonly known as a technique used to check the normality of data [13,15]. It can, however, as is further explained, also be adapted to estimate the mean and the standard deviation of censored data.

Let $y_{(1)} \geq y_{(2)} \geq \cdots y_{(i)} \cdots y_{(n)}$ be the $n$ numerically known data, arranged in descending order, observed in a censored sample of size $N$ from a normal distribution with mean $\mu$ and standard deviation $\sigma$. These ordered observations $y(i)$ are called *sample order statistics*. It can then be written as

$$E(y_{(i)}) = \mu + \sigma \xi(i/N) \tag{10}$$

where $E(y_{(i)})$ is the expected value of the $i$th sample order statistic and $\xi(i/N)$ is the expected value of the $i$th (standard) normal order statistic for the sample size $N$.

The practical meaning of $\xi(i/N)$ is as follows: if samples of size $N$ are repeatedly taken from a normal population with mean $\mu$ and standard deviation $\sigma$, the average of $y_{(i)}$, the $i$th largest value, obtained from all samples would be on $\xi(i/N)$ standard deviations from the mean. For example, $\xi(i/N)$ for sample size $N = 6$ and $i = 1$–6 is 1.267, 0.642, 0.202, $-0.202$, $-0.642$ and $-1.267$, respectively. Thus, the first and the second largest values of a normally distributed sample with $N = 6$ are on the average, respectively, 1.267 and 0.642 standard deviations above the mean. Values of $\xi(i/N)$ for the sample size $N = 2$–50 can be found in literature [13,16].

To represent all observations of the sample, Eq. (10) can be rewritten in a more generalized form in matrix notation as

$$E(y) = X\beta \tag{11}$$

where

$$y^T = [y_{(1)}, y_{(2)}, \ldots, y_{(i)}, \ldots, y_{(n)}]$$
$$\beta^T = [\mu, \sigma] \quad X = [\mathbf{1}_n, x],$$

with $\mathbf{1}_n$ an $n \times 1$ column vector of ones and

$$x^T = [\xi(1/N), \xi(2/N), \ldots, \xi(i/N), \ldots, \xi(n/N)]$$

It follows from Eq. (11) that the estimates of the $y$-intercept ($\hat{\beta}_0$) and the slope ($\hat{\beta}_1$) of the straight line regressed between the sample order statistics $y_{(i)}$ as ordinates and the expected values $\xi(i/N)$ as abscissae correspond to the estimates of the mean and the standard deviation of the sampled population, respectively. The least squares estimates of the intercept and the slope, $\hat{\beta}_0$ and $\hat{\beta}_1$, respectively, are calculated as follows [17]

$$\hat{\beta} = (X^T X)^{-1} X^T y \tag{12}$$

where

$$\hat{\beta}^T = [\hat{\beta}_0, \hat{\beta}_1] = [\hat{\mu}, \hat{\sigma}]$$

the matrix $X$ and the vector $y$ are the same as in Eq. (11).

To obtain a better estimate of $\beta$, the variance–covariance between the expected values of the (standard) normal order statistics should be taken into account in the least square fitting of a straight line through the data pairs ($\xi(i/N)$, $y_{(i)}$). Therefore, Eq. (12) becomes [8]

$$\hat{\beta} = (X^T V^{-1} X)^{-1} X^T V^{-1} y \tag{13}$$

where $V$ is the covariance matrix of the normal order statistics which can be obtained from [9] for the sample size $N = 2$–20.

To illustrate the application of the method for the analysis of censored samples, an example is given below.

**Example.** The following censored dataset ($N = 6$, $n_1 = 1$) is obtained [1.5, 1.4, 1.2, 1.1, 1.0 and <1.0]. Notice that the data are arranged in descending

order. To estimate the mean and the standard deviation, the matrix $X$ and the vector $y$ (in Eq. (13)) are defined as

$$X = \begin{bmatrix} 1 & \xi(1/6) \\ 1 & \xi(2/6) \\ \vdots & \vdots \\ 1 & \xi(5/6) \end{bmatrix} = \begin{bmatrix} 1 & 1.267 \\ 1 & 0.642 \\ 1 & 0.202 \\ 1 & -0.202 \\ 1 & -0.642 \end{bmatrix},$$

$$y = \begin{bmatrix} 1.5 \\ 1.4 \\ 1.2 \\ 1.1 \\ 1.0 \end{bmatrix}$$

The corresponding covariance matrix $V$ of $\xi(i/N)$; $i = 1, 2, \ldots, 5$ and $N = 6$ (source [9]) is

$$V = \begin{bmatrix} 0.4159 & 0.2085 & 0.1394 & 0.1024 & 0.0774 \\ 0.2085 & 0.2796 & 0.1890 & 0.1397 & 0.1059 \\ 0.1394 & 0.1890 & 0.2462 & 0.1833 & 0.1397 \\ 0.1024 & 0.1397 & 0.1833 & 0.2462 & 0.1890 \\ 0.0774 & 0.1059 & 0.1397 & 0.1890 & 0.2796 \end{bmatrix}$$

$\hat{\boldsymbol{\beta}}$ is estimated from Eq. (13) as

$$\begin{bmatrix} 1.17 \\ 0.27 \end{bmatrix}.$$

Therefore, the estimated mean and standard deviation correspond to $\hat{\beta}_0 = 1.17$ and $\hat{\beta}_1 = 0.27$, respectively.

### 2.3. Building the distributions of the estimates

#### 2.3.1. Observed distribution of the estimates

The performance of the methods for the estimation of $\mu$ and $\sigma$ is evaluated by considering the distribution of estimates around the true value. These are derived by categorizing the estimates ($\hat{\mu}$ or $\hat{\sigma}$) of 10,000 datasets into appropriate intervals. The number of estimates which falls in each interval is then counted and divided by the total number of estimates (10,000). This is then expressed as the percent of $\hat{\mu}$ (or $\hat{\sigma}$) included in the interval. Since the relative frequency is based on a large number of estimates, i.e. 10,000, it can also be interpreted as the probability that an estimate will be

observed in the interval. To describe and summarize the distributions of the estimated means $\hat{\mu}$, the following intervals are considered for each distribution:

$$\mu \pm d_\mu \frac{\sigma}{\sqrt{N}}$$

with $d_\mu = 0.67$, 1 and 1.96, these intervals describe for a given distribution the range within which the estimated mean is expected to lie with 50, 68 and 95% probability.

The distributions of the estimated standard deviations $\hat{\sigma}$ are described and summarized by considering the following intervals for each distribution:

$$\sigma \pm d_\sigma \sigma$$

with $d_\sigma = 0.25$, 0.50, 0.75 and 1.0. These intervals of course are not CIs but they describe the probability that an estimate $\hat{\sigma}$ would not differ from its true value by more than 25, 50, 75 and 100%, respectively.

#### 2.3.2. Expected distribution of the ML estimates $\hat{\mu}$

Since the results from the simulations demonstrate that the performance of the method based on order statistics are inferior to that of the ML method (see Section 3.1), the distributions of the estimates $\hat{\mu}$ are not considered here for the order statistical method. Moreover, only the expected distribution of the estimates $\hat{\mu}$ is presented here since it provides a basis for the approximation of the sample size required to obtain a specified confidence level that the ML estimate $\hat{\mu}$ of a censored sample will not deviate from the true value $\mu$ by more than a certain magnitude (see Section 3.4). The performance of the ML method in the estimation of the standard deviation being not as good as in the estimation of the mean we do not consider here the expected distribution of the estimates $\hat{\sigma}$.

The expected distribution of the ML estimate $\hat{\mu}$ is based on the variance of $\hat{\mu}$. The asymptotic variance of $\hat{\mu}$ is calculated as [2]

$$v(\hat{\mu}) = \frac{\sigma^2}{N} \mu_{11} \tag{14}$$

where $\mu_{11}$ is the variance factor of $\hat{\mu}$ which is a function of the censored proportion $h$, or the reporting limit expressed in terms of the deviation from the mean value in standard deviation units: $\eta = (y_L - \mu)/\sigma$. Values of $\mu_{11}$ associated with $h$ or $\eta$ can be found

Table 3
Variance factor $\mu_{11}$ for singly censored samples from the normal distribution [1,2]

| $\eta$ | $\mu_{11}$ | $h$ (%) | $\eta$ | $\mu_{11}$ | $h$ (%) | $\eta$ | $\mu_{11}$ | $h$ (%) | $\eta$ | $\mu_{11}$ | $h$ (%) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| −4.00 | 1.000002 | 0.00 | −0.54 | 1.132963 | 29.46 | 0.48 | 2.795042 | 68.44 | 1.50 | 33.338560 | 93.32 |
| −3.50 | 1.000013 | 0.02 | −0.52 | 1.139788 | 30.15 | 0.50 | 2.892934 | 69.15 | 1.52 | 35.376900 | 93.57 |
| −3.00 | 1.000096 | 0.13 | −0.50 | 1.146963 | 30.85 | 0.52 | 2.996396 | 69.85 | 1.54 | 37.550910 | 93.82 |
| −2.50 | 1.000559 | 0.62 | −0.48 | 1.154506 | 31.56 | 0.54 | 3.105761 | 70.54 | 1.56 | 39.870140 | 94.06 |
| −2.40 | 1.000777 | 0.82 | −0.46 | 1.162437 | 32.28 | 0.56 | 3.221383 | 71.23 | 1.58 | 42.344780 | 94.29 |
| −2.30 | 1.001072 | 1.07 | −0.44 | 1.170775 | 33.00 | 0.58 | 3.343639 | 71.90 | 1.60 | 44.985830 | 94.52 |
| −2.20 | 1.001470 | 1.39 | −0.42 | 1.179542 | 33.72 | 0.60 | 3.472929 | 72.57 | 1.62 | 47.805090 | 94.74 |
| −2.10 | 1.002000 | 1.79 | −0.40 | 1.188761 | 34.46 | 0.62 | 3.609677 | 73.24 | 1.64 | 50.815210 | 94.95 |
| −2.00 | 1.002704 | 2.28 | −0.38 | 1.198456 | 35.20 | 0.64 | 3.754337 | 73.89 | 1.66 | 54.029830 | 95.15 |
| −1.95 | 1.003137 | 2.56 | −0.36 | 1.208652 | 35.94 | 0.66 | 3.907390 | 74.54 | 1.68 | 57.463590 | 95.35 |
| −1.90 | 1.003634 | 2.87 | −0.34 | 1.219376 | 36.69 | 0.68 | 4.069347 | 75.17 | 1.70 | 61.132230 | 95.54 |
| −1.85 | 1.004203 | 3.22 | −0.32 | 1.230656 | 37.45 | 0.70 | 4.240754 | 75.80 | 1.72 | 65.052670 | 95.73 |
| −1.80 | 1.004854 | 3.59 | −0.30 | 1.242521 | 38.21 | 0.72 | 4.422191 | 76.42 | 1.74 | 69.243150 | 95.91 |
| −1.75 | 1.005598 | 4.01 | −0.28 | 1.255004 | 38.97 | 0.74 | 4.614275 | 77.04 | 1.76 | 73.723240 | 96.08 |
| −1.70 | 1.006447 | 4.46 | −0.26 | 1.268138 | 39.74 | 0.76 | 4.817664 | 77.64 | 1.78 | 78.514040 | 96.25 |
| −1.65 | 1.007416 | 4.95 | −0.24 | 1.281956 | 40.52 | 0.78 | 5.033058 | 78.23 | 1.80 | 83.638270 | 96.41 |
| −1.60 | 1.008520 | 5.48 | −0.22 | 1.296498 | 41.29 | 0.80 | 5.261203 | 78.81 | 1.82 | 89.120360 | 96.56 |
| −1.55 | 1.009776 | 6.06 | −0.20 | 1.311800 | 42.07 | 0.82 | 5.502893 | 79.39 | 1.84 | 94.986640 | 96.71 |
| −1.50 | 1.011205 | 6.68 | −0.18 | 1.327906 | 42.86 | 0.84 | 5.758976 | 79.95 | 1.86 | 101.265500 | 96.86 |
| −1.45 | 1.012829 | 7.35 | −0.16 | 1.344859 | 43.64 | 0.86 | 6.030354 | 80.51 | 1.88 | 107.987400 | 96.99 |
| −1.40 | 1.014673 | 8.08 | −0.14 | 1.362704 | 44.43 | 0.88 | 6.317988 | 81.06 | 1.90 | 115.185400 | 97.13 |
| −1.35 | 1.016765 | 8.85 | −0.12 | 1.381491 | 45.22 | 0.90 | 6.622904 | 81.59 | 1.92 | 122.895000 | 97.26 |
| −1.30 | 1.019139 | 9.68 | −0.10 | 1.401271 | 46.02 | 0.92 | 6.946196 | 82.12 | 1.94 | 131.154400 | 97.38 |
| −1.25 | 1.021829 | 10.56 | −0.08 | 1.422099 | 46.81 | 0.94 | 7.289031 | 82.64 | 1.96 | 140.005000 | 97.50 |
| −1.20 | 1.024878 | 11.51 | −0.06 | 1.444032 | 47.61 | 0.96 | 7.652652 | 83.15 | 1.98 | 149.491200 | 97.61 |
| −1.15 | 1.028331 | 12.51 | −0.04 | 1.467132 | 48.40 | 0.98 | 8.038388 | 83.65 | 2.00 | 159.661200 | 97.72 |
| −1.10 | 1.032241 | 13.57 | −0.02 | 1.491463 | 49.20 | 1.00 | 8.447655 | 84.13 | 2.02 | 170.566900 | 97.83 |
| −1.05 | 1.036667 | 14.69 | 0.00 | 1.517094 | 50.00 | 1.02 | 8.881966 | 84.61 | 2.04 | 182.264200 | 97.93 |
| −1.00 | 1.041677 | 15.87 | 0.02 | 1.544097 | 50.80 | 1.04 | 9.342936 | 85.08 | 2.06 | 194.813600 | 98.03 |
| −0.98 | 1.043860 | 16.35 | 0.04 | 1.572548 | 51.60 | 1.06 | 9.832289 | 85.54 | 2.08 | 208.280500 | 98.12 |
| −0.96 | 1.046155 | 16.85 | 0.06 | 1.602529 | 52.39 | 1.08 | 10.351870 | 85.99 | 2.10 | 222.735500 | 98.21 |
| −0.94 | 1.048565 | 17.36 | 0.08 | 1.634124 | 53.19 | 1.10 | 10.903640 | 86.43 | 2.12 | 238.254800 | 98.30 |
| −0.92 | 1.051098 | 17.88 | 0.10 | 1.667427 | 53.98 | 1.12 | 11.489700 | 86.86 | 2.14 | 254.920800 | 98.38 |
| −0.90 | 1.053759 | 18.41 | 0.12 | 1.702531 | 54.78 | 1.14 | 12.112310 | 87.29 | 2.16 | 272.822800 | 98.46 |
| −0.88 | 1.056555 | 18.94 | 0.14 | 1.739540 | 55.57 | 1.16 | 12.773860 | 87.70 | 2.18 | 292.057200 | 98.54 |
| −0.86 | 1.059493 | 19.49 | 0.16 | 1.778561 | 56.36 | 1.18 | 13.476920 | 88.10 | 2.20 | 312.728200 | 98.61 |
| −0.84 | 1.062579 | 20.05 | 0.18 | 1.819708 | 57.14 | 1.20 | 14.224230 | 88.49 | 2.22 | 334.948600 | 98.68 |
| −0.82 | 1.065821 | 20.61 | 0.20 | 1.863103 | 57.93 | 1.22 | 15.018740 | 88.88 | 2.24 | 358.840700 | 98.75 |
| −0.80 | 1.069228 | 21.19 | 0.22 | 1.908875 | 58.71 | 1.24 | 15.863580 | 89.25 | 2.26 | 384.536600 | 98.81 |
| −0.78 | 1.072808 | 21.77 | 0.24 | 1.957159 | 59.48 | 1.26 | 16.762110 | 89.62 | 2.28 | 412.179500 | 98.87 |
| −0.76 | 1.076569 | 22.36 | 0.26 | 2.008100 | 60.26 | 1.28 | 17.717940 | 89.97 | 2.30 | 441.924500 | 98.93 |
| −0.74 | 1.080520 | 22.96 | 0.28 | 2.061851 | 61.03 | 1.30 | 18.734920 | 90.32 | 2.32 | 473.939600 | 98.98 |
| −0.72 | 1.084673 | 23.58 | 0.30 | 2.118574 | 61.79 | 1.32 | 19.817160 | 90.66 | 2.34 | 508.406800 | 99.04 |
| −0.70 | 1.089036 | 24.20 | 0.32 | 2.178441 | 62.55 | 1.34 | 20.969090 | 90.99 | 2.36 | 545.523600 | 99.09 |
| −0.68 | 1.093621 | 24.83 | 0.34 | 2.241636 | 63.31 | 1.36 | 22.195430 | 91.31 | 2.38 | 585.503900 | 99.13 |
| −0.66 | 1.098439 | 25.46 | 0.36 | 2.308352 | 64.06 | 1.38 | 23.501260 | 91.62 | 2.40 | 628.579700 | 99.18 |
| −0.64 | 1.103502 | 26.11 | 0.38 | 2.378794 | 64.80 | 1.40 | 24.892000 | 91.92 | 2.42 | 675.002900 | 99.22 |
| −0.62 | 1.108823 | 26.76 | 0.40 | 2.453182 | 65.54 | 1.42 | 26.373480 | 92.22 | 2.44 | 725.046500 | 99.27 |
| −0.60 | 1.114415 | 27.43 | 0.42 | 2.531747 | 66.28 | 1.44 | 27.951930 | 92.51 | 2.46 | 779.007200 | 99.31 |
| −0.58 | 1.120292 | 28.10 | 0.44 | 2.614735 | 67.00 | 1.46 | 29.634060 | 92.79 | 2.48 | 837.206900 | 99.34 |
| −0.56 | 1.126470 | 28.77 | 0.46 | 2.702408 | 67.72 | 1.48 | 31.427030 | 93.06 | 2.50 | 899.995000 | 99.38 |

in [2]. For the convenience of the readers, Table 3 with more extensive entries of $\mu_{11}$ for the singly censored sample has been added here. They were completely recomputed. Values of $\mu_{11}$ are provided with more decimal places, as well as for more elaborated intervals of the argument $h$ or $\eta$ than those given in [2]. For the values of $h$ or $\eta$ that do not correspond with those in Table 3, linear interpolation can be applied.

It is known that for large sample size, the ML estimates are unbiased and approximately normally distributed [18]. By assuming a normal distribution of the ML estimates $\hat{\mu}$, their expected distributions for the chosen intervals $\mu \pm d_\mu(\sigma/\sqrt{N})$ can then be obtained as the proportion of the whole area under the standard normal distribution that lies between $z = -d_\mu/\sqrt{\mu_{11}}$

and $d_\mu/\sqrt{\mu_{11}}$. Despite the small sample sizes considered in this article, simulation results (see Section 3.4) demonstrate that the approximation to the normal distribution is satisfactory for (i) predicting the distribution of the estimated means, and (ii) the estimation of the optimum sample size required for the censored samples.

## 3. Results and discussion

The distributions of the estimates, $\hat{\mu}$ and $\hat{\sigma}$, obtained from 10,000 censored datasets, are summarized in Tables 4 and 5, respectively. For illustration, histograms of the distributions are given in Fig. 2.

Table 4
Distributions of the estimated means $\hat{\mu}$

| $h_{ex}$[a] (%) | Probability$_{n<2}$ (%)[b] | | | Percent of $\hat{\mu}$ included in the interval[c] | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | $\mu \pm 0.67(\sigma/\sqrt{N})$ | | | $\mu \pm (\sigma/\sqrt{N})$ | | | $\mu \pm 1.96(\sigma/\sqrt{N})$ | | |
| | $N=6$ | $N=12$ | $N=18$ | $N=6$ | $N=12$ | $N=18$ | $N=6$ | $N=12$ | $N=18$ | $N=6$ | $N=12$ | $N=18$ |
| 17 | 0 | 0 | 0 | 49 | 50 | 50 | 68 | 68 | 68 | 92 | 94 | 94 |
| | | | | *51* | *51* | *50* | *69* | *69* | *68* | *94* | *95* | *95* |
| | | | | ***44*** | ***34*** | ***24*** | ***61*** | ***50*** | ***39*** | ***90*** | ***86*** | ***79*** |
| | | | | **49** | | | **67** | | | **94** | | |
| 33 | 2 | 0 | 0 | 49 | 47 | 47 | 70 | 64 | 64 | 93 | 91 | 92 |
| | | | | *51* | *47* | *47* | *67* | *64* | *64* | *92* | *92* | *92* |
| | | | | ***23*** | ***7*** | ***2*** | ***38*** | ***14*** | ***5*** | ***78*** | ***55*** | ***35*** |
| | | | | **47** | | | **64** | | | **93** | | |
| 50 | 11 | 0 | 0 | 52 | 44 | 43 | 68 | 61 | 60 | 95 | 87 | 88 |
| | | | | *42* | *41* | *41* | *59* | *57* | *58* | *88* | *87* | *88* |
| | | | | ***5*** | ***0*** | ***0*** | ***13*** | ***1*** | ***0*** | ***55*** | ***18*** | ***4*** |
| | | | | **42** | | | **58** | | | **89** | | |
| 67 | 35 | 6 | 1 | 40 | 35 | 34 | 58 | 52 | 49 | 95 | 83 | 80 |
| | | | | *21* | *25* | *28* | *36* | *38* | *41* | *78* | *74* | *75* |
| | | | | ***0*** | ***0*** | ***0*** | ***0*** | ***0*** | ***0*** | ***20*** | ***1*** | ***0*** |
| | | | | **33** | | | **47** | | | **78** | | |
| 75 | 53 | 16 | 4 | 22 | 26 | 26 | 38 | 40 | 39 | 92 | 79 | 72 |
| | | | | *14* | *17* | *19* | *23* | *26* | *30* | *64* | *57* | *59* |
| | | | | ***0*** | ***0*** | ***0*** | ***0*** | ***0*** | ***0*** | ***3*** | ***0*** | ***0*** |
| | | | | **26** | | | **38** | | | **67** | | |
| Non-censored[d] | | | | 50 | | | 68 | | | 95 | | |

[a] $h_{ex}$: proportion of test results which are expected to be censored.

[b] Probability that a set of $N$ test results would contain only one or even no test result above the reporting limit and thus it is not possible to calculate the estimates based on either the ML method or the order statistics (see Section 3.2).

[c] Values appearing in normal, italic, and bolded italic fonts are related to the distributions of the estimates obtained by the ML method, the method based on order statistics, and the method of moments (using only numerically known data), respectively. Bolded values are the expected distributions of the ML estimates (see Section 2.3.2 for details), they are the same for all $N$, $\mu$ and $\sigma$ applied in the simulations.

[d] The observed distributions of the estimated mean obtained with the method of moments from the complete (non-censored) datasets coincide with the theoretical expectation, i.e. for all $N$, they are 50, 68 and 95% for the intervals $\mu \pm 0.67(\sigma/\sqrt{N})$, $\mu \pm (\sigma/\sqrt{N})$ and $\mu \pm 1.96(\sigma/\sqrt{N})$, respectively.

Table 5
Distributions of the estimated standard deviations $\hat{\sigma}$

| $h_{ex}$[a] (%) | Percent of $\hat{\sigma}$ included in the interval[b] | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\sigma \pm 25\%$ | | | $\sigma \pm 50\%$ | | | $\sigma \pm 75\%$ | | | $0 \rightarrow 2\sigma$ | | |
| | $N = 6$ | $N = 12$ | $N = 18$ | $N = 6$ | $N = 12$ | $N = 18$ | $N = 6$ | $N = 12$ | $N = 18$ | $N = 6$ | $N = 12$ | $N = 18$ |
| 17 | 55 | 72 | 81 | 86 | 96 | 99 | 95 | 99 | 100 | 98 | 100 | 100 |
| | *46* | *69* | *80* | *78* | *96* | *99* | *93* | *99* | *100* | *97* | *100* | *100* |
| | ***39*** | ***50*** | ***54*** | ***76*** | ***93*** | ***97*** | ***96*** | ***100*** | ***100*** | ***100*** | ***100*** | ***100*** |
| 33 | 47 | 64 | 74 | 81 | 93 | 97 | 94 | 98 | 100 | 97 | 99 | 100 |
| | *39* | *59* | *71* | *71* | *91* | *97* | *89* | *98* | *100* | *96* | *99* | *100* |
| | ***26*** | ***29*** | ***28*** | ***61*** | ***78*** | ***86*** | ***90*** | ***99*** | ***100*** | ***100*** | ***100*** | ***100*** |
| 50 | 40 | 54 | 65 | 74 | 88 | 94 | 92 | 96 | 98 | 97 | 99 | 100 |
| | *33* | *49* | *61* | *62* | *83* | *92* | *84* | *95* | *98* | *95* | *98* | *99* |
| | ***18*** | ***18*** | ***15*** | ***47*** | ***58*** | ***65*** | ***80*** | ***94*** | ***98*** | ***100*** | ***100*** | ***100*** |
| 67 | 35 | 44 | 53 | 68 | 77 | 86 | 90 | 93 | 96 | 98 | 97 | 98 |
| | *27* | *38* | *48* | *53* | *69* | *81* | *78* | *88* | *94* | *94* | *96* | *98* |
| | ***13*** | ***12*** | ***10*** | ***34*** | ***40*** | ***44*** | ***68*** | ***81*** | ***90*** | ***100*** | ***100*** | ***100*** |
| 75 | 32 | 40 | 46 | 64 | 73 | 79 | 89 | 91 | 94 | 98 | 97 | 97 |
| | *24* | *32* | *39* | *49* | *61* | *72* | *74* | *82* | *89* | *94* | *94* | *96* |
| | ***11*** | ***10*** | ***8*** | ***29*** | ***32*** | ***34*** | ***62*** | ***71*** | ***79*** | ***100*** | ***100*** | ***100*** |
| Non-censored[c] | 56 | 76 | 86 | 90 | 98 | 100 | 99 | 100 | 100 | 100 | 100 | 100 |

[a] $h_{ex}$: proportion of test results which are expected to be censored.

[b] Values appearing in normal, italic, and bolded italic fonts are related to the distributions of the estimates obtained by the ML method, the method based on order statistics, and the method of moments (using only numerically known data), respectively.

[c] The observed distributions of the estimated standard deviations obtained with the method of moments from the complete (non-censored) datasets. They are the same for all $\mu$ and $\sigma$ applied in the simulations.
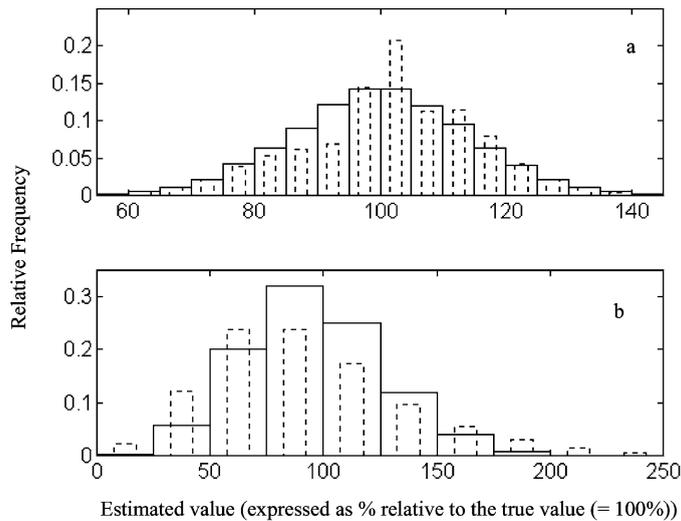


Fig. 2. Histograms obtained with $N = 6$, $\mu = 1$, $\sigma = 0.35$, $y_L = 1$. The full line corresponds to distributions of the estimates obtained from 10,000 complete (non-censored) datasets by the method of moments. The broken line corresponds to distributions of the estimates obtained from 10,000 censored datasets by the ML method. Adjacent broken line bars should touch each other as the bars with full lines do, they are separated for better visualisation. Relative frequency is obtained as the proportion of estimates that fall in the observed intervals. (a) Distributions of the estimated means $\hat{\mu}$. (b) Distributions of the estimated standard deviations $\hat{\sigma}$.

Of course, the observed distributions for $\hat{\mu}$ and $\hat{\sigma}$ obtained by the method of moments from the complete (non-censored) datasets, summarized by considering the intervals $\mu \pm d_\mu(\sigma/\sqrt{N})$ and $\sigma \pm d_\sigma \sigma$, respectively, are independent of $\mu$ and $\sigma$.

The distributions of $\hat{\mu}$ and $\hat{\sigma}$ observed for censored data are independent on the value of $\sigma$ but are dependent on the value of $\mu$ applied in the simulation. In fact, the difference is not really due to the value of $\mu$ but to the expected censored proportion $h_{ex}$. Since the limit where the test results are censored is fixed at $y_L = 1$, different $\mu$ result in different $h_{ex}$-values. Consequently, different combinations of the two parameters, $\mu$ and $\sigma$, applied in the simulations (see Section 2.1.1) are represented in Tables 4 and 5 by a single parameter $h_{ex}$ (see Eq. (1) for the relationship of $\mu$, $\sigma$ and $h_{ex}$).

In order to demonstrate a typical problem of censored data, distributions of the estimates $\hat{\mu}$ and $\hat{\sigma}$, calculated by the simple method of moments on numerically known data, are also given (see Tables 4 and 5). A comparison of the results obtained from the three methods (ML method, order statistical method, method of moments) clearly shows relatively large errors for the simple moment calculation. For example, for a sample size $N = 6$ and an expected censored proportion $h_{ex} = 50\%$, the probability that an estimate $\hat{\mu}$ would lie within $\mu \pm 0.67(\sigma/\sqrt{N})$ is only 5% for the simple moment calculation versus 52% for the ML method.

### 3.1. Comparison between the ML method and the method based on order statistics

It follows from Tables 4 and 5 that the performance of the method based on order statistics is generally inferior to that of the ML method. Both have comparable performance in the estimation of the mean only when the expected censored proportion is not large, i.e. $h_{ex} < 33\%$. The larger the expected censored proportion (e.g. $h_{ex} \geq 50\%$) and the smaller the sample size (e.g. $N = 6$), the better the performance of the ML method is relative to the method based on order statistics. Since the overall efficiency of the ML method is much better than the method based on order statistics, only the ML method is discussed further in what follows.

### 3.2. Performance of the ML method in the estimation of the mean of censored datasets

As follows from Table 4, in general, the larger the censored proportion, the smaller the probability (expressed as percent of the estimates $\hat{\mu}$) that an estimate $\hat{\mu}$ lies close to the true value $\mu$ and, thus, the worse the performance of the ML method is. When the expected censored proportion ($h_{ex}$) does not exceed 50%, i.e. when the true mean of the test results, $\mu$, is not below the limit where the test results are censored (here the reporting limit), the performance of the ML method in the estimation of the mean of a censored dataset is very comparable to the performance of the method of moments in the estimation of the mean of the complete (non-censored) dataset. When $h_{ex} \leq 50\%$, results from data simulations (not shown in the table) also demonstrate that the 95% CIs of the means (calculated using Eq. (15)) obtained with the ML method really contain the true mean about 95% of the times, while less confidence is obtained when $h_{ex} > 50\%$. It is remarkable that with $N = 6$, the probability that an estimate $\hat{\mu}$ falls in the intervals $\mu \pm d_\mu(\sigma/\sqrt{N})$ is higher than when $N = 12$ or $N = 18$. That is, the performance of the ML method in the estimation of the mean from the censored sample with $N = 6$ is better than when $N = 12$ or $N = 18$ (notice that in the evaluation of the performance, standardized distributions are considered; the distribution of the means of course becomes wider with smaller $N$ due to an increase of the variance). This might be due to the fact that the probability, that an observed censored proportion $h = 100n_1/N$ (or an observed number of censored test results $n_1$) corresponds with the expected censored proportion $h_{ex}$ (or $n_{1ex}$), is larger with smaller $N$. For example, from data simulations of 10,000 censored samples with $h_{ex} = 50\%$ and $N = 6$, 12 and 18, the probability that the $n_1$ observed for censored samples corresponds with the $n_{1ex}$ are 36, 23 and 18%, respectively. The accuracy of the estimates depends on the degree of correspondence between the observed censored proportion ($h = 100n_1/N$) and the expected censored proportion $h_{ex}$. The higher the degree of correspondence, the smaller the probability that an estimate is biased and, thus, the better performance of the ML method.

When the true value of the test results, $\mu$, is below the limit at which the data are censored, i.e. when

more than half of the test results are expected to be censored, it should first be noticed that with $N = 6$, it is quite probable (e.g. probability = 35% when $\mu$ is below $y_L$ for about half a standard deviation) that the ML estimates cannot be calculated due to the fact that almost all test results of the sample are below the limit. In case that it is possible to calculate the ML estimates, the more the true $\mu$ is below the (reporting) limit (i.e. the larger the expected censored proportion is), the poorer the ML estimates of $\mu$ are. When the expected censored proportion is as high as 83%, the quality of estimates becomes so poor that it is not recommended to do any computation.

The problem, however, is that in practice the true $\mu$ is unknown. We then have no idea of the expected censored proportion nor whether $\mu$ is above or below the (reporting) limit. When $\mu$ is above the limit (i.e. $h_{ex} < 50\%$), the performance of the ML method applied to censored datasets is quite acceptable. On the other hand, when $\mu$ is below the (reporting) limit (i.e. $h_{ex} > 50\%$), the quality of the ML estimates deteriorates. Fortunately, the probability that the ML estimates cannot be calculated in that case (probability$_{n < 2}$ in Table 4) also increases, especially when the sample size is small. Probability$_{n < 2}$ is obtained by determining the percentage of datasets (based on >10,000 datasets) that are discarded because the number $n$ of test results that remains after leaving out those test results that are below the reporting limit is either 1 or 0. A total number of >10,000 datasets is simulated in order to obtain 10,000 censored datasets (see Section 2.1.3). From Table 4, it can be seen for example that with the expected censored proportion $h_{ex} = 75\%$, the probability that the ML estimates cannot be calculated (probability$_{n < 2}$) is as high as 53% for $N = 6$, and therefore, it is less probable that such poor ML estimates as displayed in the table would be obtained.

In a previous study [5,6], simulations were performed and the method performance was evaluated using the average bias and variance for each value of the observed number of data censor $n_1$ (not the true value of $n_{1ex}$ or $h_{ex}$). It was then concluded that the ML method and the linear method (which is analogous to the method based on order statistics) produced large bias and large variance for small $N$ ($N = 5, 10, 15$). This conclusion is contradictory to the above finding that for small censored samples

($N = 6, 12, 18$), the ML method and the method based on order statistics perform acceptably if the censored proportion $h_{ex}$ does not exceed 50 and 33%, respectively. However different approaches were used in the evaluation of the method performance. The current approach seems more justified for the following reasons.

1. Using the observed $n_1$ as in the previous study [6] makes less sense than using the true value $n_{1ex}$ or analogously, $h_{ex}$ (as in this study) since an observed value of $n_1$ is only a randomly sampled value of the true value $n_{1ex}$, observed with a certain probability (see Fig. 1). The grouping of censored datasets based on a value of $n_1$ may represent as little as 5% of the total censored samples when the true censored proportion is $h_{ex}$ [6]. Though all possible values of $n_1$ were considered in [6], the evaluation of method performance was performed separately for different values of $n_1$ (despite they are observed values of the same true value $n_{1ex}$). Consequently, the observed $n_1$ might not well represent the real situation of censored samples as the true value of $n_{1ex}$ or $h_{ex}$ does.

2. In [5,6], the performance of the methods is judged only by comparing the averages for bias and variance or mean square error among different methods applied to censored samples. Here the method performance is evaluated by comparing the distribution of the estimates around their true value not only for different methods applied to censored samples but the comparison is also made with the classical method of moments applied on complete (non-censored) samples. The performance of the former are judged reliable when they are comparable to that of the non-censored samples.

3. The performance of the methods should not be judged only by considering the averages for bias and variance or mean square error [5,6], but also by considering the distribution of the estimates around their true value.

## 3.3. Performance of the ML method in the estimation of the standard deviation of censored datasets

It follows from Table 5 that in general, the probability that an estimated standard deviation of a censored dataset lies close to the true value is smaller than that

observed from complete (non-censored) dataset. The distributions of the ML estimates of $\sigma$ are broader than the distributions of $\hat{\sigma}$ obtained from the complete (non-censored) datasets by the method of moments. The performance of the ML method in the estimation of the standard deviation of a censored sample is not as good as in the estimation of the mean. The more the data are censored, the larger the variance of $\hat{\sigma}$ is. The efficiency in the estimation of $\sigma$ depends of course not on the true value $\sigma$ but on the sample size $N$ and on the expected censored proportion $h_{ex}$. The larger the sample size and the smaller the expected censored proportion, the better the ML method performs in the estimation of $\sigma$.

### 3.4. Calculation of the optimum sample size required to obtain a specified confidence level on the ML estimate $\hat{\mu}$

It can be seen from Table 4 that the observed distributions of the ML estimates of the mean $\mu$ correspond well with the expected distributions determined with the method described in Section 2.3.2. Generally, when they do not correspond well, the observed distributions are narrower than the expected ones. Simulations demonstrate that an observed distribution converges to the expected one when the sample size grows large, e.g. $N \geq 60$ (results not shown). In addition to the small sample size, the deviation of the observed censored proportion $h$ from the expected $h_{ex}$ also plays a role in the difference between the expected and the observed distributions. However, the data simulations from Table 4 show that the ML method for the estimation of the mean of a censored sample is efficient with small $N$ ($N = 6$–$18$). Moreover, as explained below, the sample size required to have a specified confidence level that an estimate will not differ from the true mean by a certain magnitude can be determined before carrying out an analysis.

Suppose that we want to determine the sample size $N$ required to have 95% confidence that a $\hat{\mu}$ estimated (by the ML method) from the censored sample will not differ from the true value $\mu$ by more than $d_{spec}$. The question can be restated as finding the sample size $N$ such that the resulting 95% CI for $\mu$, $\hat{\mu} \pm d_{obs}$, has its half width $d_{obs}$ not larger than $d_{spec}$.

By assuming a normal distribution of $\hat{\mu}$, the half width $d_{obs}$ of the 95% CI for $\mu$ can be obtained as

$$d_{obs} = 1.96\sqrt{V(\hat{\mu})} \leq d_{spec} \tag{15}$$

where 1.96 is the two-sided tabulated $z$-value of the standard normal distribution at the significance level $\alpha = 0.05$.

By substituting Eq. (14) into Eq. (15), we obtain

$$1.96\sqrt{\frac{\sigma^2}{N}\mu_{11}} \leq d_{spec} \tag{16}$$

From this, the approximate sample size $N$ required to have 95% confidence that an estimate $\hat{\mu}$ will not differ from the true mean $\mu$ by more than $d_{spec}$ can be calculated as

$$N \geq \left(\frac{1.96\sigma}{d_{spec}}\right)^2 \mu_{11} \tag{17}$$

**Example.** If the censored proportion $h_{ex}$ and the standard deviation $\sigma$ are expected to be at most 50% and 0.35, respectively, the approximate sample size $N$ required to have 95% confidence that the ML estimate $\hat{\mu}$ of a censored sample will not differ from the true value $\mu$ by more than $d_{spec} = 0.20$ can be determined as

$$N \geq \left(\frac{1.96 \times 0.35}{0.20}\right)^2 \mu_{11}$$

From Table 3, $\mu_{11}$ associated with $h = 50\%$ is 1.517094. Thus,

$$N \geq \left(\frac{1.96 \times 0.35}{0.20}\right)^2 1.517094, \quad \text{or } N \geq 17.8 \geq 18$$

Therefore, if the expected censored proportion $h_{ex}$ and standard deviation $\sigma$ do not exceed 50% and 0.35, respectively, 18 measurements are required to have a confidence of about 95% that the ML estimate $\hat{\mu}$ will not differ from its true value $\mu$ by more than 0.2. If a confidence level of only 90% is required, the $z$-value of 1.96 in Eq. (17) can be replaced by 1.645. Thus, the required $N$ would be 13, instead of 18. Results from data simulations (based on 10,000 censored samples of $\mu = 1$ and $y_L = 1$ (i.e. $h_{ex} = 50\%$), $\sigma = 0.35$, $N = 13$ and 18) verify that the sample size approximation satisfies the requirement:

$\hat{\mu}$ falls in the interval, $\mu \pm d_{\text{spec}} = \mu \pm 0.2$, 89 and 93% of the times for $N = 13$ and 18, respectively.

## 4. Conclusion

The performance of the ML method and the method based on order statistics in the estimation of the mean and the standard deviation of a normal population from a censored sample, i.e. a sample of which some test results fall below the reporting limit of the analytical method, has been investigated by means of simulations. The study, focussed on small to moderate sample size ($N = 6$–18), demonstrates that reliable estimates for censored samples can be obtained. The approaches are computationally simple and provide the possibility to utilize all data for obtaining more realistic estimates of the mean and standard deviation. The reliability of the estimated mean and the estimated standard deviation for the censored samples depends on the proportion of data expected to be censored $h_{\text{ex}}$ ($h_{\text{ex}}$ is the proportion of test results that are not known numerically since they are below the reporting limit). The larger the expected censored proportion, the less reliable the methods are. When the expected censored proportion $h_{\text{ex}}$ is not large, e.g. $h_{\text{ex}} = 33\%$, either method proposed in this paper (Section 2.2) can be applied to obtain good estimates $\hat{\mu}$ and $\hat{\sigma}$. However, in practice it is likely that the value of $h_{\text{ex}}$ is unknown. Moreover, the observed censored proportion $h$ does not necessarily correspond to the true value $h_{\text{ex}}$. Since the simulations demonstrate that the ML method is more robust to large censored proportions, as well as to small sample sizes, the ML method is preferred to the method based on order statistics when pre-knowledge about the censored sample is not available. The reliability of the ML estimates is rather poor for very large censored proportions (e.g. $h_{\text{ex}} \geq 75\%$) and therefore, the method is not recommended when the observed censored

proportion $h$ is larger than 50% because then $h_{\text{ex}}$ might be very large. If pre-knowledge about the precision of the analytical method and the expected censored proportion is available, the approximate sample size required to obtain a specified confidence level on the ML estimate of the mean can be determined before the experiment is carried out.

## Acknowledgements

## References

[1] A.C. Cohen Jr., Technometrics 1 (1959) 217.
[2] A.C. Cohen Jr., Technometrics 3 (1961) 535.
[3] S. Noack, Personal communication.
[4] C.C. Travis, M.L. Land, Environ. Sci. Technol. 24 (1990) 961.
[5] D.R. Helsel, Environ. Sci. Technol. 24 (1990) 1767.
[6] A. Gleit, Environ. Sci. Technol. 19 (1985) 1201.
[7] R.H. Shumway, A.S. Azari, P. Johnson, Technometrics 31 (1989) 347.
[8] A.K. Gupta, Biometrika 39 (1952) 260.
[9] A.E. Sarhan, B.G. Greenberg, Ann. Math. Statist. 27 (1956) 427.
[10] M.S. Wolynetz, Appl. Statist. 28 (1979) 185.
[11] J. Schmee, D. Gladstein, W. Nelson, Technometrics 27 (1985) 119.
[12] S. Kuttatharmmakul, J. Smeyers-Verbeke, D.L. Massart, D. Coomans, S. Noack, Trends Anal. Chem. 19 (2000) 215.
[13] International Standard, ISO 2854–1976 (E), Geneva.
[14] Matlab Reference Guide, The MathWorks, Natick, MA, USA, 1992, pp. 402–403.
[15] R.R. Sokal, F.J. Rohlf, Biometry, 2nd Edition, W.H. Freeman & Company, New York, 1981, pp. 117–124.
[16] F.J. Rohlf, R.R. Sokal, Statistical Tables, W.H. Freeman & Company, San Francisco, 1969, pp. 171–173.
[17] D.L. Massart, B.G.M. Vandeginste, L.M.C. Buydens, S. De Jong, P.J. Lewi, J. Smeyers-Verbeke, Handbook of Chemometrics and Qualimetrics: Part A, Elsevier, Amsterdam, 1997.
[18] H.J. Larson, Introduction to Probability Theory and Statistical Inference, Wiley, new York, USA, 1969, pp. 223–234.