

# Estimating the number of usability problems affecting medical devices: modelling the discovery matrix

**Vincent Vandewalle**

Inria Centre de recherche Lille Nord Europe

**Alexandre Caron** (✉ [alexandre.caron2@univ-lille.fr](mailto:alexandre.caron2@univ-lille.fr))

Univ. Lille, CHU Lille, EA2694, F-59000 Lille, France <https://orcid.org/0000-0002-9872-0633>

**Coralie Delettrez**

CHU LILLE

**Renaud Périchon**

Universite de Lille

**Sylvia Pelayo**

Centre d'Investigation Clinique - Innovation Technologique Lille

**Alain Duhamel**

Universite de Lille

**Benoit Dervaux**

Universite de Lille

---

## Technical advance

**Keywords:** usability testing, medical device, missing data, Bayesian statistics, maximum likelihood

**Posted Date:** August 5th, 2020

**DOI:** <https://doi.org/10.21203/rs.2.16958/v4>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published on September 18th, 2020. See the published version at <https://doi.org/10.1186/s12874-020-01091-y>.

1 **Title Page**

2 **Title of the manuscript:** Estimating the number of usability problems affecting medical  
3 devices: modelling the discovery matrix

4 **Author listing:**

5 **Vincent Vandewalle**<sup>a,b\*</sup>, PhD, **Alexandre Caron**<sup>a\*</sup>, MD, MSc, **Coralie Delettrez**<sup>c</sup>, MSc,  
6 **Renaud Périchon**<sup>a</sup>, MSc, **Sylvia Pelayo**<sup>a,d</sup>, PhD, **Alain Duhamel**<sup>a,c</sup>, PhD, **Benoit**  
7 **Dervaux**<sup>a,c</sup>, PhD.

8 **\*Equal contributions**

9 **Affiliations:**

- 10 - <sup>a</sup>Univ. Lille, CHU Lille, ULR 2694 Evaluations des technologies de santé et des  
11 pratiques médicales, F-59000 Lille, France,  
12 - <sup>b</sup>Inria, F-59000 Lille, France,  
13 - <sup>c</sup>CHU Lille, Direction de la Recherche et de l'Innovation, F-59000 Lille, France  
14 - <sup>d</sup>CIC-IT/Evalab 1403, CHU Lille, F-59000 Lille, France.

15 **Corresponding Author:** Alexandre Caron ([alexandre.caron2@univ-lille.fr](mailto:alexandre.caron2@univ-lille.fr))

16 **Word Count:** 7126

## 17 **Abstract**

18 **Background.** Usability testing of medical devices are mandatory for market access. The  
19 testings' goal is to identify usability problems that could cause harm to the user or limit  
20 the device's effectiveness. In practice, human factor engineers study participants under  
21 actual conditions of use and list the problems encountered. This results in a binary  
22 discovery matrix in which each row corresponds to a participant, and each column  
23 corresponds to a usability problem. One of the main challenges in usability testing is  
24 estimating the total number of problems, in order to assess the completeness of the  
25 discovery process. Today's margin-based methods fit the column sums to a binomial  
26 model of problem detection. However, the discovery matrix actually observed is  
27 truncated because of undiscovered problems, which corresponds to fitting the marginal  
28 sums without the zeros. Margin-based methods fail to overcome the bias related to  
29 truncation of the matrix. The objective of the present study was to develop and test a  
30 matrix-based method for estimating the total number of usability problems.

31 **Methods.** The matrix-based model was based on the full discovery matrix (including  
32 unobserved columns) and not solely on a summary of the data (e.g. the margins). This  
33 model also circumvents a drawback of margin-based methods by simultaneously  
34 estimating the model's parameters and the total number of problems. Furthermore, the  
35 matrix-based method takes account of a heterogeneous probability of detection, which  
36 reflects a real-life setting. As suggested in the usability literature, we assumed that the  
37 probability of detection had a logit-normal distribution.

38 **Results.** We assessed the matrix-based method's performance in a range of settings  
39 reflecting real-life usability testing and with heterogeneous probabilities of problem  
40 detection. In our simulations, the matrix-based method improved the estimation of the

41 number of problems (in terms of bias, consistency, and coverage probability) in a wide  
42 range of settings. We also applied our method to five real datasets from usability testing.

43 **Conclusions.** Estimation models (and particularly matrix-based models) are of value in  
44 estimating and monitoring the detection process during usability testing. Matrix-based  
45 models have a solid mathematical grounding and, with a view to facilitating the decision-  
46 making process for both regulators and device manufacturers, should be incorporated  
47 into current standards.

48 **Keywords:** usability testing, medical device, missing data, Bayesian statistics, maximum  
49 likelihood

50

## 51 Main manuscript text

### 52 I. Background

#### 53 A. Introduction

54 The usability testing is a cornerstone of medical device development, and proof of  
55 usability is mandatory for market access in both the European Union and the United  
56 States [1]. The overall objective of a usability assessment is to ensure that a medical device  
57 is designed and optimized for use by the intended users in the environment in which the  
58 device is likely to be used [2]. The goal is to identify problems (called “use errors”) that  
59 could cause harm to the user or impair medical treatment (e.g. an inappropriate number  
60 of inhalations, finger injection with an adrenaline pen, etc.) [3]. The detection of usability  
61 problems must be as comprehensive as possible because medical devices are safety-  
62 critical systems [4]. However, the total number of usability problems is never known in  
63 advance. The main challenge during the usability testing is thus to estimate this number,  
64 in order to assess the completeness of the problem discovery process [5].

65 In practice, participants are placed under actual conditions of use (real or simulated), and  
66 usability problems are observed and listed by human factor engineers. The experimental  
67 conditions are defined in a risk analysis that gathers together possible usability problems.  
68 Throughout the usability testing, problems are discovered and added to a discovery  
69 matrix - a binary matrix with the participants as the rows and the problems as the  
70 columns. The current approach involves estimating the total number of problems as the  
71 usability testing progresses, starting from the first sessions. The number is estimated  
72 iteratively as the sample size increases, until the objective of completeness has been  
73 achieved [6].

74 From a statistical perspective, the current estimation procedure is based on a model of  
75 how the usability problems are detected; this is considered to be a binomial process. The  
76 literature suggests that the total number of usability problems can be estimated from the  
77 discovery matrix's problem margin (the sum of the columns) [7-11]. However, this  
78 estimation is complicated by (i) the small sample size usually encountered in usability  
79 testing of medical devices [12] and (ii) as-yet unobserved problems that truncate the  
80 margin and bias estimates [13-15].

81 The objective of the present study was to develop a matrix-based estimation of the  
82 number of usability problems affecting a medical device. This new method is based on the  
83 likelihood of the discovery matrix (rather than the matrix's margins alone), so as to avoid  
84 a reduction in the level of information prior to modeling. The method's main targets are  
85 (i) regulatory agencies and notified bodies involved in the pre-market evaluation of  
86 medical devices, and (ii) medical device manufacturers (more specifically, the human  
87 factors engineers in charge of ensuring that the devices are usable).

## 88 B. Data collected during the usability testing: the discovery matrix

89 The human factor engineer collects the results of the usability testing in a problem-  
90 discovery matrix  $\mathcal{d}$ . Each row corresponds to a participant, and each column corresponds  
91 to a usability problem. The result is 1 if the participant discovered the problem and 0 if  
92 not. Considering that after the inclusion of  $n$  participants,  $j$  problems have been  
93 discovered, a  $n \times j$  matrix is built. By way of an example, the discovery matrix obtained  
94 after  $n = 8$  participants (in rows) might be the one presented below:

95

$$\mathfrak{d} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{pmatrix}$$

96 In this example,  $j = 10$  different problems (in columns) have been detected so far. The  
 97 first participant discovered only one problem (column 1), whereas the second discovered  
 98 two new problems (columns 2 and 3), etc.

99 At this stage, some problems might not have been detected, and the total number of  
 100 usability problems ( $m$ ) is unknown. It should be noted that by definition,  $m \geq j$  and  $m -$   
 101  $j$  problems remain undetected. Indeed,  $\mathfrak{d}$  comes from a complete but unobserved matrix  
 102 of dimensions  $n \times m$ . This matrix is denoted as  $\mathfrak{x}$ . Thus, the “observed” matrix  $\mathfrak{d}$  is a  
 103 truncated version of the “complete” matrix  $\mathfrak{x}$ ; it lacks the columns corresponding to the  
 104 as-yet undetected problems. Hereafter, we use the following notation:  $\mathfrak{x} = (x_{il})_{1 \leq i \leq n, 1 \leq l \leq m}$   
 105 where  $x_{il} = 1$  if the participant  $i$  experiences the problem  $l$ , and  $x_{il} = 0$  otherwise.

106

$$\mathfrak{x} = \begin{pmatrix} x_{11} & \cdots & x_{1l} & \cdots & x_{1m} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{i1} & \cdots & x_{il} & \cdots & x_{im} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nl} & \cdots & x_{nm} \end{pmatrix}$$

107 The human factor engineer’s goal is to estimate the total number of problems  $m$  from the  
 108 discovery matrix  $\mathfrak{d}$  and thus deduce  $m - j$  - the number of problems that have not been  
 109 detected. The new method presented below addresses this goal.

### 110 C. Conventional estimation of $m$ using a margin-based probabilistic 111 model

112 In this section, we describe the margin-based methods currently employed to estimate  
113 the number of usability problems. As mentioned above,  $m$  is currently estimated by fitting  
114 a probabilistic (binomial) model to the discovery matrix's problems margin. More  
115 specifically, the probability with which a given usability problem is discovered by a  
116 participant is modelled by a Bernoulli trial with a probability of success (i.e. detection)  $p$ .  
117 For a given problem, the Bernoulli trial is considered to apply independently to each of  
118 the  $n$  participants in the usability testing. Thus, the problem margin sums can be  
119 considered as an independent, identically distributed sequence of Bernoulli trials, in  
120 which the number of times a given usability problem (a random variable  $X$ ) has been  
121 observed after  $n$  participants follows a binomial distribution,  $X \sim \text{Bin}(n, p)$ . Considering  
122 the binomial distribution of the margin sums, the proportion of problems that has been  
123 discovered at least once after  $n$  participants is given by the cumulative function of the  
124 shifted geometric distribution [6, 16, 17]:

$$125 \quad P(X > 0) = 1 - (1 - p)^n \quad (1)$$

126 The total number of problems  $m$  is then deduced from the following relationship:

$$127 \quad j = (1 - (1 - p)^n) \times m \quad (2)$$

128 The discovery progress is thus assessed in two steps: the probability of detection  $p$  is first  
129 estimated and then plugged into Equation (2) to estimate the number of problems  $m$ . A  
130 wide range of literature methods are available for estimating the probability of problem  
131 detection. The simplest way involves computing the naive estimate (denoted as  $\hat{p}$ ) using

132 the observed discovery matrix  $\mathbb{d}$ , considering that only  $j$  problems have been detected so  
133 far:

$$134 \quad \hat{p} = \frac{\sum_{i=1}^n \sum_{l=1}^j x_{il}}{n * j} \quad (3)$$

135 As mentioned above, the naïve estimate is systematically biased - especially for small  
136 samples. Indeed, unobserved problems result in zero columns that shrink the probability  
137 space and lead to overestimation of  $p$ , particularly at the beginning of the process when  
138  $j \ll m$ . Consequently,  $m$  is systematically underestimated, which generates safety  
139 concerns in the medical device field. In response, several strategies have been employed  
140 to overcome the truncated matrix problem.

141 In 2001, Hertzum and Jacobsen suggested normalizing the value of  $\hat{p}$  [9]. This procedure  
142 considers that the lower boundary of the probability of detection estimated with  $n$   
143 participants is  $1/n$ . For example, in a sample of 5 participants,  $\hat{p} \in [0.2 ; 1]$ . Conversely,  
144 the normalized estimator  $\hat{p}_{Norm} \in [0; 1]$ , and is computed as follows:

$$145 \quad \hat{p}_{Norm} = \frac{\hat{p} - \frac{1}{n}}{1 - \frac{1}{n}} \quad (4)$$

146 However, the normalized approach suffers from a major limitation when estimating the  
147 total number of problems with Equation (4). In fact, if each participant has discovered  
148 only one problem and if each problem was discovered only once,  $\hat{p} = \frac{1}{n}$ ,  $\hat{p}_{Norm} = 0$ , and  
149 the estimated number of problems  $\hat{m}$  is infinite. We will not discuss this estimation  
150 method further.

151 Turing and Good developed a discounting method for estimating the probability of unseen  
152 species on the basis of observed data [18]. Lewis suggested that the Good-Turing (GT)

153 adjustment could be used to reduce the magnitude of the overestimation of  $p$  by  
 154 increasing the probability space and thus accounting for unobserved usability problems  
 155 [8]. The GT adjustment is computed as the proportion of singletons relative to the total  
 156 number of events (i.e. the proportion of problems discovered only once,  $x_{il} = 1$ ), and is  
 157 incorporated in the estimation as follows:

$$158 \quad \hat{p}_{GT} = \frac{\hat{p}}{1 + GT} \quad (5)$$

159 However, Lewis observed that use of the GT estimator overestimated  $p$ . He empirically  
 160 assessed the best adjustment for a small sample size by carrying out Monte Carlo  
 161 simulations on a range of usability testing databases involving web or software user  
 162 interfaces with known true values. Based on these simulations, Lewis concluded that the  
 163 best method was to average the GT adjustment and a “double-deflation” term:

$$164 \quad \hat{p}_{\text{double-deflation}} = \frac{1}{2} \left[ \frac{\hat{p}}{1 + GT_{adj}} \right] + \frac{1}{2} \left[ \left( \hat{p} - \frac{1}{n} \right) \times \left( 1 - \frac{1}{n} \right) \right] \quad (6)$$

165 Nevertheless, the degree of adjustment of the probability space for unobserved problems  
 166 is essentially empirical. The residual bias is not known to trend towards over- or  
 167 underestimation.

168 In 2009, Schmettow considered the problem margin sums in a zero-truncation framework  
 169 [19]. Indeed, the distribution of the problems so far observed follows a binomial  
 170 distribution with only a positive integer as support (i.e. a positive or conditional  
 171 distribution). The distribution is zero-truncated because problems only appear in the  
 172 discovery matrix once they have been discovered. The probability is then estimated using  
 173 standard mathematical techniques, such as the maximum likelihood or moment estimator  
 174 [20-22]. The probability mass function is:

175 
$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k} \quad (7)$$

176 and zero truncation is achieved as follows:

177 
$$P(X = k)_{zt} = \begin{cases} 0 & \text{if } k = 0 \\ \frac{P(X = k)}{1 - P(X = 0)} & \text{if } k > 0 \end{cases} \quad (8)$$

178 The probability of problem discovery is then estimated by using maximum likelihood  
 179 techniques to fit the marginal sums to the zero-truncated binomial distribution. It should  
 180 be noted that the expected probability of unobserved problems,  $\Pr(X = 0)$ , is deduced  
 181 from the non-truncated function [19].

182 **D. Methods taking account of a heterogeneous problem detection**  
 183 **probability**

184 All the methods presented above assume that the probability of detection is the same for  
 185 all usability problems (i.e., the same  $p$ ). However, this assumption is unrealistic and does  
 186 not hold true in real-life usability testing. Schmettow showed that overdispersion was  
 187 frequent in the problem margin sums, reflecting heterogeneity in the probability of  
 188 detection [23]. Furthermore, erroneously ignoring the presence of heterogeneity by using  
 189 a single, average value of  $p$  leads to overestimation of the completeness of the discovery  
 190 process (Jensen's inequality) [24]. Schmettow tackled this problem by developing a model  
 191 that incorporated heterogeneity. The probability of detection was considered to be a  
 192 random variable, which enabled each problem to have its own probability of detection.  
 193 Schmettow used the logit-normal distribution as a plugin distribution for the probability  
 194 of detection. Formally, the logit of the probability of detection follows a normal  
 195 distribution  $\mathcal{N}(\mu, \sigma)$ . In this model, the problem margin sums follows a logit-normal  
 196 binomial distribution and the probability mass function is:

197 
$$P(X = k) = \binom{n}{k} \frac{1}{\sqrt{2\pi}\sigma} \int_0^1 (1-p)^{n-k-1} p^{k-1} \exp\left(-\frac{(\text{logit}(p) - \mu)^2}{2\sigma^2}\right) dp \quad (9)$$

198 Using the zero truncation technique presented in equation (8), Schmettow developed the  
 199 logit-normal binomial zero truncated (LNBzt) model and applied it to the usability of  
 200 medical infusion pumps [25]. To the best of our knowledge, this model is the only one that  
 201 accounts for both heterogeneity and unobserved problems.

## 202 E. Statistical limitations of margin-based methods

203 The primary limitation of the margin-based methods presented above is that they  
 204 estimate the probability of detection only. The number of problems  $m$  is deduced but not  
 205 estimated *per se*. It would be possible to estimate both  $m$  and  $p$  by summarizing the  
 206 discovery matrix on the basis of the participants' margin. In such a case, each sum follows  
 207 a binomial  $\text{Bin}(m, p)$ , thus enabling estimation of both the number of attempts and the  
 208 probability of success in a binomial setting. However, DasGupta and Rubin established  
 209 that there were no unbiased estimates for essentially any functions of either the number  
 210 of attempts or the probability of success [26]. This problem was initially considered by  
 211 Fisher and Haldane for estimating species abundance [27, 28]. It has also been considered  
 212 by Olkin, Petkau, and Zidek, who developed both a moment and a maximum likelihood  
 213 estimator, and by Carroll and Lombard, who proposed an estimator in a Bayesian setting  
 214 (leading to a beta-binomial distribution) [29, 30]. Hall also considered this problem in an  
 215 asymptotic framework [31].

216 The second limitation of margin-based methods is information loss, relative to the initially  
 217 available data. For example,  $j$  and the number of singletons were the only data used in the  
 218 GT estimates. In the same way, the zero-truncated method considered only the column  
 219 sums for the problems and omitted the pattern of detection (i.e., the users).

220 Here, we tackle these problems by directly modelling the full discovery matrix (including  
221 unobserved columns) and not only a summary of the data (e.g. the margins). In the  
222 Methods section, we describe the statistical basis of the matrix-based method and detail  
223 a Bayesian approach for estimating the number of problems. In the Results section, we  
224 compare the matrix-based method’s statistical properties with those of existing models  
225 in a simulation study and then in actual usability studies. Lastly, we discuss the  
226 implications of our results with regard to estimation of the number of problems in  
227 usability testing.

## 228 II. Methods

229 We first specify the statistical basis underpinning the matrix-based method, and the  
230 principle of column permutation in particular. Next, we present our estimation of the  
231 number of problems in a Bayesian setting. The last part is dedicated to the methods used  
232 to assess the matrix-based model’s performance.

### 233 A. The matrix-based method

234 We first present the matrix-based method. For the sake of clarity, we simplified the  
235 problem by considering that the probability of problem detection was homogeneous. The  
236 concept of heterogeneous probability will be introduced in the second part of this section,  
237 along with the Bayesian estimation.

#### 238 1. Presentation of the method

239 Consider the complete discovery matrix  $\mathbf{x}$ . The probability of  $\mathbf{x}$  can be written as follows:

$$240 \quad P(\mathbf{x}|p, m) = p^{\mathbf{x}_{..}}(1 - p)^{nm - \mathbf{x}_{..}} \quad (10)$$

241 where  $\mathbf{x}_{..} = \sum_{i=1}^n \sum_{l=1}^m x_{il}$  is the total number of problems observed by  $n$  participants.

242 An example of a possible matrix  $\mathbb{x}$  obtained from two participants during a usability  
 243 testing of a medical device with  $m = 3$  problems is given below (with users in rows and  
 244 problems in columns):

$$245 \quad \mathbb{x} = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix} \quad (11)$$

246 As seen above, the complete discovery matrix  $\mathbb{x}$  is never observed, and the discovery  
 247 matrix  $\mathbb{d}$  is the only one available. It is similar to the matrix  $\mathbb{x}$ , except that unobserved  
 248 problems are missing. Considering the above example, neither of the users observed the  
 249 second problem, and the resulting observed discovery matrix  $\mathbb{d}$  would be:

$$250 \quad \mathbb{d} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad (12)$$

251 It should be noted that if the total number of problems  $m$  is known, then the complete  
 252 matrix  $\mathbb{x}$  could be reconstituted (with permutation), based on the matrix  $\mathbb{d}$ . For instance,  
 253 if we take the matrix  $\mathbb{x}$  and consider (wrongly, in this case) that the number of problems  
 254  $m = 5$ , then the reconstituted complete matrix denoted by  $\hat{\mathbb{x}}^m$  would be obtained by  
 255 padding the matrix  $\mathbb{d}$  with columns of zeros (corresponding to as-yet unobserved  
 256 problems):

$$257 \quad \hat{\mathbb{x}}^{m=5} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{pmatrix} \quad (13)$$

258 Thus, noting that  $\mathbb{x}_{..} = \mathbb{d}_{..}$ , it is possible to compute the likelihood of the complete matrix  
 259  $\hat{\mathbb{x}}^m$  on the basis of the discovery matrix  $\mathbb{d}$ . This likelihood is given by the following  
 260 equation:

$$261 \quad P(\hat{\mathbb{x}}^m | p, m) = p^{\mathbb{x}_{..}} (1 - p)^{nm - \mathbb{x}_{..}} \quad (14)$$

262 Note that the definition of  $\hat{\mathbb{x}}^m$  depends on the value  $m$ , which is unknown. Thus, any  
 263 inference based on  $\hat{\mathbb{x}}^m$  will induce some bias. For instance, a maximum likelihood  
 264 estimation of  $(p, m)$  based on  $\hat{\mathbb{x}}^m$  (consisting in maximizing  $p(\hat{\mathbb{x}}^m|p, m)$  with respect to  $m$   
 265 and  $p$ ) leads to  $\hat{m} = j$  (where  $j$  is the number of problems observed so far) and  $p = \frac{\mathbb{x}_{\bullet j}}{nj}$ ,  
 266 which are known to be biased. We tackled this issue by modeling the distribution of the  
 267 observed discovery matrix  $p(\mathbb{d}|p, m)$ .

268 It should be noted that the matrix  $\mathbb{d}$  is defined in a lexicographic order, which simply  
 269 means that the problems are ordered in the order of detection. For instance, the six  
 270 possible complete matrices  $\mathbb{x}$  leading to the previous matrix  $\mathbb{d}$  if  $m = 3$  are presented in  
 271 Table 1.

272 *Table 1: Six possible complete matrices  $\hat{\mathbb{x}}^{m=3}$  leading to the observed discovery matrix  $\mathbb{d} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$*

Possibility 1	Possibility 2	Possibility 3
$\hat{\mathbb{x}}_1^{m=3} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$	$\hat{\mathbb{x}}_2^{m=3} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$	$\hat{\mathbb{x}}_3^{m=3} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}$
Possibility 4	Possibility 5	Possibility 6
$\hat{\mathbb{x}}_4^{m=3} = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}$	$\hat{\mathbb{x}}_5^{m=3} = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$	$\hat{\mathbb{x}}_6^{m=3} = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}$

273  
 274 In fact, if we could consider the label (the name of the usability problem) associated with  
 275 each column, only one matrix  $\mathbb{x}$  could lead to the matrix  $\mathbb{d}$ . However, since we have no  
 276 means of finding the names of the columns in the initial matrix  $\mathbb{x}$ , we will consider that  
 277 the matrix  $\mathbb{d}$  has unnamed columns. Removing these column names allows us to consider  
 278 the matrix  $\mathbb{d}$  for the observed data (for which the definition does not vary as a function of  
 279 the model's definition of the model – in contrast to  $\hat{\mathbb{x}}^m$ ). Thus:

280 
$$P(\mathbb{d}|m = 3, p) = \sum_{h=1}^6 P(\hat{\mathbb{x}}_h^{m=3}|m = 3, p) \quad (15)$$

281 and more generally

282 
$$P(\mathbb{d}|m, p) = \sum_{h=1}^{H(\mathbb{d}, m)} P(\hat{\mathbb{x}}_h^m|m, p) \quad (16)$$

283 where  $H(\mathbb{d}, m)$  is the number of different matrices  $\hat{\mathbb{x}}_h^m$  with  $m$  columns leading to the  
 284 same discovery matrix  $\mathbb{d}$ .

285 In the simple example presented above (Table 1),  $H(\mathbb{d}, m) = 6$  and each matrix  $\hat{\mathbb{x}}_h^m$  has  
 286 the same probability, i.e.  $p^2(1 - p)^4$ . It follows that:

287 
$$P(\mathbb{d}|m = 3, p) = H(\mathbb{d}, m = 3) \times P(\hat{\mathbb{x}}_h^{m=3}|m = 3, p) =$$
  

$$6 \times p^2(1 - p)^4 = A_3^2 \times p^2(1 - p)^4 \quad (17)$$

288 More generally, the number of matrices  $\mathbb{x}$  with  $m$  columns associated with an observed  
 289 discovery matrix  $\mathbb{d}$  is:

290 
$$H(\mathbb{d}, m) = \frac{m!}{(m - j)! j_1! \dots j_r!} = \frac{1}{j_1! \dots j_r!} \times A_m^j \quad (18)$$

291 where  $r$  is the number of different columns of  $\mathbb{d}$ , and  $j_h$  ( $1 \leq h \leq r$ ) is the number of  
 292 repetitions of the column of type  $h$ . Of course,  $j = j_1 + \dots + j_r$ . Here, we recognize a  
 293 familiar equation: that associated with the number of anagrams of a word in which each  
 294 type of column corresponds to a different letter, including the null column (repeated  $m -$   
 295  $j$  times).

296 Lastly, since each matrix  $\hat{\mathbb{x}}_h^m$  has the same probability, we obtain the likelihood of  $\mathbb{d}$  as  
 297 follows:

298 
$$P(\mathbb{d}|p, m) = \frac{1}{j_1! \dots j_r!} \times A_m^j \times P(\hat{\mathbb{x}}_h^m | m, p) \quad (19)$$

299 In practice, the computation of  $\frac{1}{j_1! \dots j_r!}$  has no impact on the estimation, since it is the same  
 300 for all values of  $m$  and  $p$ . This result is not limited to the homogenous setting and would  
 301 remain valid for any probability of  $\mathbb{x}$  with a column-wise exchangeability property.

302 In the particular case of the homogeneous setting, we obtain:

303 
$$P(\mathbb{d}|p, m) = \frac{1}{j_1! \dots j_r!} \times A_m^j \times p^{\mathbb{x}\bullet\bullet} (1-p)^{nm-\mathbb{x}\bullet\bullet}. \quad (20)$$

304 In the homogeneous setting, our matrix-based approach could be extended to perform  
 305 maximum likelihood inference or Bayesian inference on the parameters. However, as  
 306 explained above, this setting is unrealistic in practice and so a heterogeneous probability  
 307 of detection should be considered in the following section.

## 308 2. Heterogeneity and Bayesian estimation

309 We considered a heterogeneous probability of detection; i.e. each problem  $l$  has its own  
 310 probability of detection  $p_l$ . In line with Schmettow's method, we assume that the  
 311 probabilities of detection are independent and follow a logit-normal distribution, i.e.  
 312  $\text{logit}(p_l) \sim \mathcal{N}(\mu, \sigma)$ . The model's parameters are  $m$ ,  $\mu$  and  $\sigma$ . Note that  $p_1, \dots, p_m$  are  
 313 considered as latent random variables - like random effects in the mixed model.

314 Given these parameters, the likelihood of the discovery matrix  $\mathbb{d}$  can be written as

315 
$$P(\mathbb{d}|\mu, \sigma, m) = \int_0^1 \dots \int_0^1 P(\mathbb{d}|p_1, \dots, p_m, m) f(p_1, \dots, p_m | \mu, \sigma) dp_1 \dots dp_m \quad (21)$$

316 where  $f(p_1, p_2, \dots, p_m | \mu, \sigma)$  is the probability density function of  $p_1, p_2, \dots, p_m$ . Given that  
 317 the columns are exchangeable, we can also write

318 
$$P(\mathbb{d}|\mu, \sigma, m) = \frac{1}{j_1! \dots j_r!} \times A_m^j \times P(\hat{\mathbb{x}}_h^m | \mu, \sigma, m) \quad (22)$$

319 which will be useful for subsequent computations.

320 We now consider a Bayesian framework [32] for estimation of the parameters. This  
 321 framework has good theoretical properties and can include prior knowledge about the  
 322 problem's parameters. Indeed, the distribution of the parameters  $P(\mu, \sigma, m)$  must first be  
 323 defined. Moreover, assuming the prior independence of  $\mu$ ,  $\sigma$  and  $m$ ,  $P(\mu, \sigma, m) =$   
 324  $P(\mu)P(\sigma)P(m)$ . We assume a prior uniform distribution for  $m$ :

325 
$$P(m) = \frac{1}{M} \forall m \in \{1, \dots, M\} \quad (23)$$

326 The value of  $M$  is the pre-determined upper boundary for  $m$ , and should be chosen by the  
 327 human factor engineer according to the expected maximum possible number of problems.  
 328 To prevent underestimation, a high value should be used. However, if  $M$  is unnecessarily  
 329 high, it will lead to an increase in the computing time.

330 Since our goal here is to estimate the number of problems, our main interest is  $P(m|\mathbb{d})$ ,  
 331 which is obtained using Bayes' theorem:

332 
$$P(m|\mathbb{d}) = \frac{P(m) \times P(\mathbb{d}|m)}{\sum_{m'=1}^M P(m') \times P(\mathbb{d}|m')} \quad (24)$$

333 Thus, we need to compute  $P(\mathbb{d}|m)$  for each possible value of  $m$  in  $\{1, \dots, M\}$ . This  
 334 computation requires computation of the integrated likelihood  $P(\mathbb{d}|m)$ , as follows

335 
$$P(\mathbb{d}|m) = \int_0^{+\infty} \int_{-\infty}^{+\infty} P(\mathbb{d}|\mu, \sigma, m)P(\mu)P(\sigma)d\mu d\sigma \quad (25)$$

336 The choice of prior distributions for  $P(\mu)$  and  $P(\sigma)$  is discussed below.  $P(\mathbb{d}|m)$  can be  
337 computed by approximating this integral with Markov chain Monte Carlo (MCMC)  
338 techniques.

339 Even though  $P(m|\mathbb{d})$  is the main quantity of interest,  $P(\mu|\mathbb{d})$  and  $P(\sigma|\mathbb{d})$  are also of  
340 interest because they can be used as prior distributions for future studies; this will  
341 decrease the sample size and improve early estimates as part of an early control strategy.

### 342 3. Computational aspects

343 From a computational perspective, and since  $P(\mathbb{d}|\mu, \sigma, m) = \frac{1}{j_1! \dots j_r!} \times A_m^j \times$   
344  $P(\mathbb{X}_h^m|\mu, \sigma, m)$ , we will first focus on the computation based on  $\mathbb{X}_h^m$  and will then deduce  
345 the results for  $\mathbb{d}$ .

346 Let now consider the choice of a prior distribution for  $\mu$  and  $\sigma$ . Since  $\mu$  and  $\sigma$  are Gaussian  
347 distribution parameters and in the absence of additional information (e.g. from previous  
348 usability studies), we chose the following flat priors:

- 349 -  $\mu \sim \mathcal{N}(0; \mathcal{A})$ : a Gaussian distribution with a high variance  $\mathcal{A}$ , (e.g.  $\mathcal{A} = 10^8$ ),  
350 mimicking a uniform distribution on  $\mathbb{R}$ ,
- 351 -  $\sigma^2 \sim \text{inv} - \chi_\nu^2$ : an inverse chi-squared distribution with  $\nu$  degrees of freedom  
352 (typically  $\nu = 1$ ).

353 When the data has a Gaussian distribution, choosing the above priors leads to a  
354 conjugated posterior distribution. However, a logistic-normal distribution of the  
355 probabilities of detection means that conjugacy cannot be obtained. Thus, estimation of  
356 the posterior distribution required the use of MCMC methods. This consisted in drawing  
357  $\mu$  and  $\sigma$  for each possible value of  $m$ ,  $m \in 1, \dots, M$  according to their posterior distribution

358  $P(\mu, \sigma | m, \mathbb{d})$ , and deducing a numerical approximation of  $P(\mathbb{d} | m)$  from the Monte-Carlo  
 359 sample. Lastly,  $P(m | \mathbb{d})$  was computed using Bayes' theorem.

360 For a fixed value of  $m$ , we consider sampling from  $P(\mu, \sigma | \hat{\mathbb{X}}_h^m, m)$ , computing the  
 361 integrated likelihood  $P(\hat{\mathbb{X}}_h^m | m)$  with bridge sampling [33], and deducing  $P(\mathbb{d} | m)$ .

362 The parameters  $\mu$  and  $\sigma$  (given  $\hat{\mathbb{X}}_h^m$  and  $m$ ) are sampled using the parameter space  
 363 augmented by  $p_1, \dots, p_m$ , i.e. the discovery probabilities associated with each column of  
 364  $\hat{\mathbb{X}}_h^m$ . Thus, we will now sample from  $\mu, \sigma, p_1, \dots, p_m | \hat{\mathbb{X}}_h^m$ , using stan software (adaptive  
 365 Hamiltonian Monte Carlo algorithm).

## 366 B. Assessment of the performance of the matrix-based method

367 We compared the performance of five methods (naïve, GT, double-deflation, LNBzt, and  
 368 matrix-based methods) first in a simulation study and then using literature data from  
 369 actual usability studies.

### 370 1. Simulation study

371 Each simulation consisted in generating an observed discovery matrix  $\mathbb{d}$  from the  
 372 usability testing of a hypothetical medical device with a known total number of usability  
 373 problems  $m$  and a sample size  $n$ . The probability of detection was normally distributed  
 374 ( $\mathcal{N}(\mu, \sigma)$ ) on a logit scale. The combinations of parameters used in the simulations are  
 375 specified in Table 2. The values were chosen to reflect a wide range of parameters  
 376 encountered in usability testing of medical devices.

377 *Table 2: Combinations of parameters for the simulation testing with homogeneous and heterogeneous*  
 378 *probabilities of detection.*

Parameter	Values
<b>Total number of usability problems</b>	$m = 20, 50, 100$
<b>Sample size</b>	$n = 15, 20, 30, 40, 50$
<b>Probability of problem detection</b>	$\mu = \text{logit}(0.1), \text{logit}(0.2)$ $\sigma = 0.5, 1, 2$
<b>Number of combinations tested</b>	<b>90</b>

379

380 In each setting (i.e. for each combination of  $m, \mu, \sigma$  and  $n$ ), we simulated  $S = 2 \times 10^4$   
 381 complete discovery matrices,  $\mathbb{X}_{m,\mu,\sigma,n,i}, i \in \{1,2, \dots, S\}$ . The matrices  $\mathbb{d}$  were obtained by  
 382 truncation of the zero columns (problems not yet discovered). We averaged the estimates  
 383 of  $m$  over the  $S$  simulations and computed the 95% fluctuation interval (0.025 and 0.975  
 384 quantiles). We also calculated the prediction's root mean square error (RMSE) as the  
 385 square root of the mean square difference between the predicted and true values of  $m$ :

$$386 \quad RMSE(m) = \sqrt{\frac{1}{S} \sum_{i=1}^S (m - \hat{m}_i)^2} \quad (26)$$

387 When the sample is small, little information is available; a tight credible interval might  
 388 reflect overconfidence rather than a good estimation. Thus, to gauge the level of  
 389 confidence that human factor engineers can place in each method, we computed the  
 390 coverage probability. In each setting, this is the proportion of 95% confidence intervals  
 391 for the simulated  $\hat{m}_i$  that include the true value of  $m$ . The confidence intervals for  $\hat{m}_i$  were  
 392 computed using 1000 parametric bootstrap repetitions with the parameters  
 393  $(\hat{m}_i, \hat{\mu}_i, \hat{\sigma}_i, n)$ . For the matrix-based method, we were able to directly compute the 95%  
 394 confidence interval of the posterior distribution of each simulation, which saved  
 395 substantial computation time.

## 396 2. Application to actual usability studies

397 We applied the above-described methods to the discovery matrices of five published  
 398 usability studies. Four did not involve a medical device: the EDU3D dataset encompassed  
 399 119 problems discovered by 20 participants during the evaluation of virtual  
 400 environments [34], the MACERR dataset encompassed 145 problems discovered by 15  
 401 participants during a scenario-driven usability testing of an integrated office system [35],

402 the MANTEL dataset encompassed 30 problems submitted by 76 expert participants  
403 evaluating the specifications of a computer program, and the SAVINGS dataset  
404 encompassed 48 usability problems discovered by 34 participants on voice response  
405 systems MANTEL and SAVINGS comes from the same experiment on heuristic evaluations  
406 [36]. These four studies were included because they have been used in important  
407 publications in this field [8] and they enabled us to address heterogeneity in the  
408 probability of discovery, in particular [23]. The fifth usability testing involved a medical  
409 device: INFPUMP encompassed 107 usability problems discovered by 34 participants  
410 (intensive care unit nurses and anesthesiologist) evaluating a prototype medical infusion  
411 pump [25].

412 For each of the five datasets, we computed the estimates and the 95% confidence intervals  
413 for the final data. When a sufficient number of participants had been included (i.e. for  
414 MANTEL, SAVINGS, and INFPUMP), we addressed the change in the estimates as a  
415 function of the sample size.

416 All the analyses were carried out running R software (version 3.6.1) on several servers  
417 equipped with 12-core Intel® Xeon® E5-2650 v4 processors ([http://hpc.univ-](http://hpc.univ-lille.fr/cluster-hpc-htc)  
418 [lille.fr/cluster-hpc-htc](http://hpc.univ-lille.fr/cluster-hpc-htc)). The MCMC was performed using the *Stan* library ([http://mc-](http://mc-stan.org)  
419 [stan.org](http://mc-stan.org)) via the *rstan* package [37]. The integrated likelihood was obtained using the  
420 *bridge\_sampler* function of the *bridgesampling* package [38]. In order to facilitate the  
421 matrix-based method's application in practice, a short step-by-step tutorial [see  
422 Additional file 1] and the code [see Additional file 2] is provided as supplementary  
423 material. A reproducible R code with the data and the simulation study performed in this  
424 manuscript is available on GitHub ([https://github.com/alexandre-caron/matrix based-](https://github.com/alexandre-caron/matrix-based-)

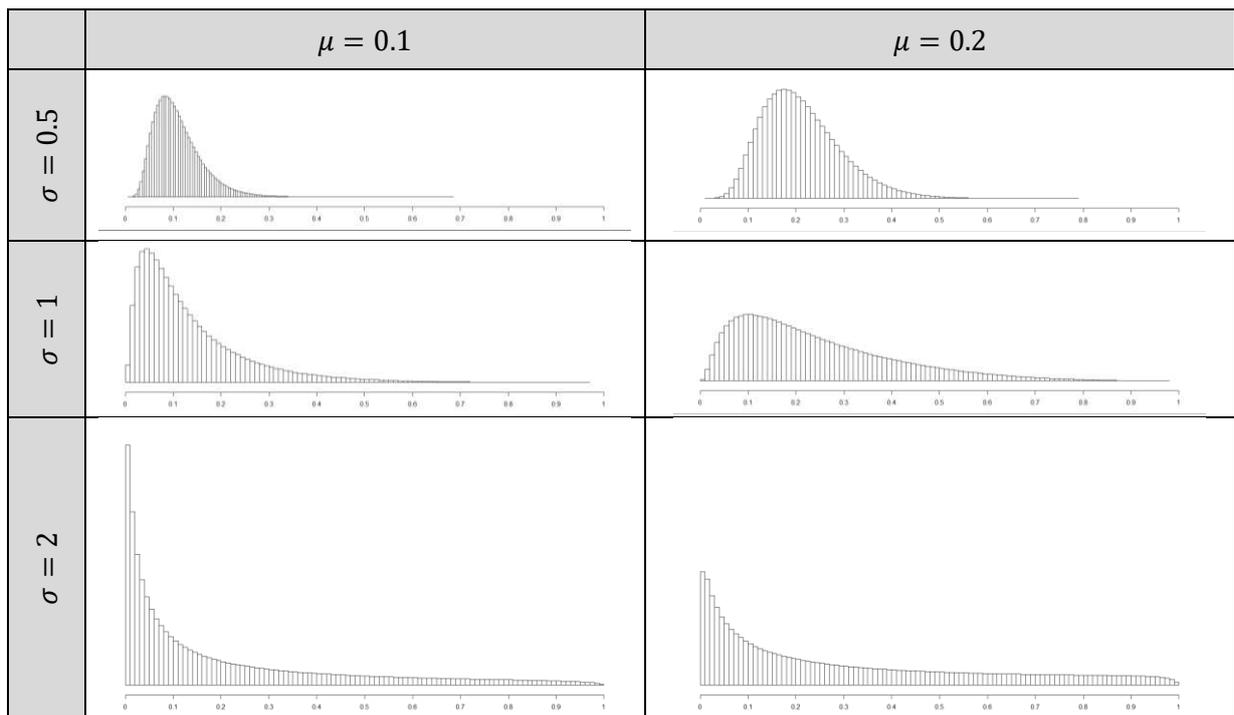
425 [usability](#)). The link to the archived version referenced in this manuscript is available in  
 426 the “Availability of data and materials” section.

### 427 III. Results

#### 428 A. The simulation study

429 The distributions of the probability of detection for each setting are summarized in Table  
 430 3. The distribution shifted to a highest average probability of detection when  $\mu$  increased.  
 431 It is noteworthy that a higher dispersion ( $\sigma$ ) not only flattened the distribution but also  
 432 led to an increase in probability of very rare problems.

433 *Table 3: Distribution of the probability of detection as a function of  $\mu$  and  $\sigma$ . The probability of detection*  
 434 *followed a logit-normal distribution:  $\text{logit}(p_1) \sim \mathcal{N}(\mu, \sigma)$ .*



435

436 The results of the simulation are presented for the five methods (naïve, GT, double-  
 437 deflation, LNBzt, and matrix-based). The prediction error of  $m$  as a function of the sample  
 438 size  $n$  are presented in Figure 1. The RMSE is presented in Figure 2. A tabulated version  
 439 of these data is also provided as supplementary material [see S-Table 5 and S-Table 6 in

440 Appendix 3]. As mentioned by Schmettow, extreme estimates of  $m$  can be obtained with  
441 the LNBzt method when the number of singletons is high. We decided to discard any  
442 results with  $\hat{m}_{LNBzt} > 500$ , to avoid penalizing the method with estimates that would not  
443 be realistic in real life [19].

444 As expected, the accuracy of the estimation of the number of problems increased with the  
445 sample size for all estimates, with less bias and greater consistency (i.e. the RMSE tended  
446 towards zero as the sample size increased). Likewise, the estimates were better as the  
447 number of problems to discover  $m$  increased. For all methods, the bias was higher as the  
448 number of “rare” problems increased (i.e. for a higher  $\sigma$ ).

#### 449 **Methods accounting for heterogeneity: the matrix-based and LNBzt estimates**

450 The matrix-based method showed less bias overall; the bias ranged from -8.5% to +14.7%  
451 for the 90 simulated combinations. This range was narrower (from -5.1 to +1.2%) when  
452 the participant sample size was 30 or more. In contrast, the LNBzt method displayed  
453 systematic upward bias; although the lower boundary was -0.1%, the upper boundary  
454 was 54.7%. This bias was still observed for 30 participants, with an upper boundary of  
455 23.8%.

456 When  $\sigma = 2$ , the matrix-based method underestimated the number of problems.  
457 However, this underestimation was less than -5.1% for  $n \geq 30$ . For lower values of  $\sigma$ , the  
458 matrix-based method’s bias ranged from -2.6% to +1.2% for  $n \geq 30$ . The bias associated  
459 with the LNBzt method was high for  $\sigma = 2$ . Although the bias decreased with  $n$ , it was still  
460 +11.8% for  $n = 50$ . For a lower value of  $\sigma$ , the bias associated with the LNBzt method  
461 ranged from -2.6% to +1.2% for  $n \geq 30$ .

462 The matrix-based method gave the lowest RMSE in all settings. This was particularly true  
463 when the number of “rare” problems was high ( $\sigma > 0.5$ ). The LNBzt gave the highest  
464 average RMSE. As mentioned in the Methods, this bias resulted from a few very high  
465 estimates of  $m$ , which increased the average RMSE dramatically. This was true for the  
466 lowest average probability of detection (i.e.  $\mu = \text{logit}(0.1)$ ) and the highest variance (i.e.  
467  $\sigma = 2$ ).

468 **Methods not accounting for heterogeneity: the naïve, GT, and double-deflation**  
469 **estimates**

470 The estimates that did not take account of heterogeneity showed the strongest bias. The  
471 naïve estimate was the worst; it systematically underestimated the true value of  $m$   
472 (range: -33.2% to -0.2%). This underestimation was slightly lower for the GT estimate,  
473 especially when  $\sigma$  was low. However, the range was still broad: from -32.2% to -0.2%. The  
474 double-deflation method compensated even more for underestimation but sometimes led  
475 to overestimation (range: -32.0% to +8.6%).

476 When  $\sigma$  was lower (i.e. 0.5 or 1), the trend towards underestimation was less pronounced  
477 for the double-deflation and the GT methods (with lower boundaries of -14.1% and -17.2,  
478 respectively) than for the naïve method (lower boundary: -22.8%). The bias persisted for  
479 larger sample sizes: it was still as high as -6.4% for the three methods for  $n = 50$ .

480 The naïve RMSE estimate was again the worst of the methods that did not take account of  
481 heterogeneity. Although the GT and the double-deflation methods gave acceptable RMSEs,  
482 this feature must be interpreted with caution. In fact, the acceptable RMSEs resulted  
483 essentially from systematic underestimation, which in turn limited the range of possible  
484  $\hat{m}$  (which can never be lower than  $j$ ). Hence, the interpretation of the RMSE was limited  
485 for these methods.

486 **Coverage probability.**

487 As explained in the Methods, human factor engineers do not know the variables for the  
488 usability testing they are carrying out. The coverage probability enables them to study the  
489 reliability of the estimate (and its 95% confidence interval). A tabulated version of the  
490 data is provided as supplementary material [see S-Table 7 in Appendix 3].

491 For the matrix-based method, the coverage probability was always over 80% (except for  
492  $m = 100, n = 15, \mu = \text{logit}(0.1)$ , and  $\sigma = 0.5$ , where the probability of coverage dropped  
493 to 72%) with an average of 94% over the range of settings tested in the simulations study.  
494 The probability was at least 81% for  $n \geq 20$  and at least 88% for  $n \geq 30$ . The LNBzt  
495 method's coverage probability was always over 80%, with an average of 92%. The LNBzt  
496 performed particularly well for small sample sizes, with a minimum coverage of 89% for  
497  $n = 15$ , of 86% for  $n = 20$ , and of 82% for  $n = 30$ . Indeed, the LNBzt method provided  
498 the broadest confidence intervals of the five methods studied here. It is noteworthy that  
499 the LNBzt method was the only one that sometimes failed to fit the data (in 33% of cases).  
500 However, it was impossible to adjust the method's parameter for each individual  
501 simulation. In practice, changing the optimization function's starting values would avoid  
502 most of the fitting failures.

503 The methods not taking account of heterogeneity provided a low, erratic coverage  
504 probability in most settings. On average, the coverage probabilities were 17.9%, 31.5%  
505 and 33.7% for the naïve, GT, and double-deflation methods, respectively. Furthermore,  
506 the three methods frequently yielded excessively high estimated levels of confidence -  
507 especially for high values of  $m$ .

508 **Lessons learned from the simulation study**

509 From the human factor engineer’s point of view, the matrix-based and LNBzt methods are  
 510 the only reliable ones; they gave a good coverage probability in almost any setting and for  
 511 almost any sample size. Conversely, the methods not taking account of heterogeneity  
 512 were unreliable and so could not be trusted.

513 **B. Application to real data from published usability studies**

514 The estimated number of problems computed from the discovery matrices of five  
 515 published usability studies are presented in Table 4. Although the real number of  
 516 problems is not known, we can compare the matrix-based method’s predictions with  
 517 those of the other methods (and especially the LNBzt method).

518 *Table 4: The estimated number of problems for five real datasets from published usability studies.*

	$n^*$	$j^{**}$	naïve	Good-Turing	double deflation	LNBzt	matrix-based
<b>EDU3D</b>	<b>20</b>	<b>119</b>					
$\hat{m}$			120	121	122	155	152
95%CI			117 – 121	118 – 125	120 – 129	132 – 195	135 – 167
<b>MACERR</b>	<b>15</b>	<b>145</b>					
$\hat{m}$			156	178	184	449	382
95%CI			146 – 160	171 – 207	192 – 245	256 – 1301	346 – 440
<b>MANTEL</b>	<b>76</b>	<b>30</b>					
$\hat{m}$			30	30	30	31	30
95%CI			30 – 30	30 – 30	30 – 30	31 – 35	30 – 37
<b>SAVINGS</b>	<b>34</b>	<b>44</b>					
$\hat{m}$			44	44	44	46	45
95%CI			44 – 45	44 – 45	44 – 45	42 – 50	44 – 51
<b>INFPUMP</b>	<b>34</b>	<b>107</b>					
$\hat{m}$			107	107	107	122	120

$n^*$	$j^{**}$	naïve	Good-Turing	double deflation	LNBzt	matrix-based
95%CI		107 – 108	106 – 108	106 – 108	110 – 136	112 – 143

519 \*  $n$  is the number of participants in the study

520 \*\*  $j$  is the number of problems discovered after analyses by  $n$  participants

521

522 In these five datasets, the number of participants ranged from 15 to 76. Previous studies  
523 of these datasets [8, 19, 23, 25] demonstrated that the probability of problem detection  
524 was heterogeneous. As suggested by the results of the simulation study, the methods not  
525 taking account of heterogeneity considered that the discovery process was complete or  
526 very close to being complete for all datasets (except MACERR: see below). Thus, we  
527 compared the results of the methods that do account for heterogeneity. It is noteworthy  
528 that the estimates of  $\mu$  and  $\sigma^2$  by both the LNBzt and the matrix-based methods fell within  
529 the range observed in our simulation study for all datasets other than MACERR.

530 All five methods considered that the SAVINGS and MANTEL datasets were complete after  
531 34 and 76 participants had been included, respectively. However, the confidence intervals  
532 produced by the matrix-based and the LNBzt methods suggest that few problems had yet  
533 to be discovered.

534 The matrix-based and the LNBzt methods estimated similar number of problems for  
535 EDU3D (  $\hat{m}_{\text{matrix-based}} = 152$  and  $\hat{m}_{\text{LNBzt}} = 155$ ). The 95% confidence interval was  
536 broader for the LNBzt method (132 to 195) than for the matrix-based method (135 to  
537 167).

538 The infusion pumps in the INFPUMP study were in early-stage development, and an  
539 additional re-design phase (for fixing the usability problems discovered) was planned;  
540 this explains why  $n = 107$  unique problems were detected by the 34 participants in the  
541 usability testing. The LNBzt and matrix-based methods gave similar estimates and

542 confidence intervals:  $\hat{m}_{\text{LNBzt}} = 122$  (i.e. 15 undiscovered problems), with a 95%  
543 confidence interval from 115 to 131, whereas  $\hat{m}_{\text{matrix-based}} = 120$ , with a 95% confidence  
544 interval from 112 to 143. The parameters computed by the matrix-based method  
545 predicted an average probability of detection  $\hat{\mu}_{\text{matrix-based}} = \text{logit}(0.136)$  and a dispersion  
546 of  $\hat{\sigma}_{\text{matrix-based}} = 1.52$ . For the LNBzt method, the probability  $\hat{\mu}_{\text{LNBzt}} = \text{logit}(0.136)$  was  
547 the same, and the dispersion was slightly higher ( $\hat{\sigma}_{\text{LNBzt}} = 1.50$ ). The confidence interval  
548 (from 110 to 136) was narrower. The true number of problems with the pump was not  
549 known because it was redesigned after 34 participants had tested the device. However, if  
550 we accept the parameters  $\hat{\mu}$  and  $\sigma$  as true and apply the results of our simulation study,  
551 the INFPUMP data suggest that the LNBzt and matrix-based methods are both reliable.  
552 Nevertheless, the breadth of the respective confidence intervals emphasizes the  
553 remaining uncertainty for these two methods.

554 Using the MACERR data, the LNBzt predicted a very low average probability of detection  
555 ( $\hat{\mu}_{\text{LNBzt}} = \text{logit}(0.014)$ ) and a high level of heterogeneity ( $\hat{\sigma}_{\text{LNBzt}} = 1.90$ ). These values  
556 were out of the range of the settings tested in the simulation study, and suggested that the  
557 number of “rare” problems was high. This might explain the high number of problems  
558 predicted by the LNBzt method ( $\hat{m}_{\text{LNBzt}} = 449$ ), and the very large 95% confidence  
559 interval (from 256 to 1301). The matrix-based method’s estimate was lower  
560 ( $\hat{m}_{\text{Matrix-based}} = 382$ ), and the 95% confidence interval was narrower (346 to 440).  
561 However, the number of participants included in MACERR was low ( $n=15$ ); a larger  
562 number of participants would have been necessary to discover new problems and  
563 improve the estimates.

564 On average, computation of the estimate and its confidence interval took less than ten  
565 minutes for the matrix-based method, less than one minute for the LNBzt method, and  
566 only a few seconds for the three other methods.

## 567 IV. Discussion

568 We decided to model the full discovery matrix (including unobserved columns) and not  
569 just a summary of the data (e.g. the margins). The estimation problem was considered  
570 simultaneously in terms of the (heterogeneous) probability of problem detection and the  
571 number of problems. Although the experimental conditions in real-life usability testing  
572 are unknown, the matrix-based method outperformed the other methods and appeared  
573 to be the most reliable in a broad range of settings.

574 Most of the currently available methods assume that the probability of detection is the  
575 same for all problems. This assumption is likely to be wrong, since real data show that the  
576 probability of detection varies [19, 23]. Furthermore, ignoring heterogeneity is known to  
577 strongly bias the results [24, 39]. We therefore developed a method that accounted for  
578 heterogeneity in the probability of problem discovery  $p$ ; we used a logit-normal  
579 distribution as a plugin to model this uncertainty. The choice of this distribution was  
580 convenient in that it allowed us to compare our method with the only published model  
581 that accounts for heterogeneity. However, there are no data for confirming the validity of  
582 this choice. Nevertheless, this limitation could be easily overcome by replacing the logit-  
583 normal by another distribution (such as beta or gamma) if it proves to be more  
584 appropriate. This choice could be made using model choice criteria (e.g. the Akaike  
585 information criterion or the Bayesian information criterion). However, it should be borne  
586 in mind that for a small sample size, fitting for both incompleteness and heterogeneity is  
587 complex and inevitably leads to a high degree of uncertainty.

588 Here, we sampled  $\mu$  and  $\sigma$  for fixed values of  $m$ . This turned out to be a rather time-  
589 consuming strategy because we had to run as many chains as there were values of  $m$ . We  
590 chose not to sample directly from the joint distribution  $P(\mu, \sigma, m | \mathbb{d})$  because the  
591 dimension of the latent parameters  $p_1, p_2, \dots, p_m$  varied as a function of  $m$  - making it  
592 impossible to use a standard MCMC algorithm. In this particular situation, use of the  
593 reversible jump algorithm [40] might be a solution but would considerably complicate  
594 our algorithm.

595 There are two key moments in medical device development for assessing the best method.  
596 Early in the development cycle, the device is not mature; usability testing is referred to as  
597 “formative” because many usability problems are being discovered and corrected in an  
598 iterative design improvement process. Just before market access, usability testing is  
599 referred to as “validation” testing; they are performed on the final version of the device to  
600 ensure that no critical usability problems remain [1, 2].

601 The number of participants in the validation testing is an important parameter for both  
602 the regulatory authorities and the device manufacturer. Indeed, a sufficient sample size  
603 will (i) guarantee the medical device’s compliance with the safety standards required for  
604 market authorization, and (ii) avoid a “black swan” effect that would strongly affect the  
605 manufacturer’s credibility and profitability [41]. The validation testing focuses on the  
606 detection of infrequent usability problems. The US Food and Drug Administration  
607 requires a minimum of 15 participants [1]. This minimum is based on a naïve estimate,  
608 which has been proven to dramatically underestimate the true number of usability  
609 problems for this number of participants [12]. Indeed, the average coverage probability  
610 observed in our simulation study for  $n = 20$  was as low as 12% and did not exceed 51%.  
611 Furthermore, this threshold does not consider heterogeneity in the probability of

612 problem detection. Our findings suggest that to produce a relevant estimate with the  
613 matrix-based method, at least 20 participants are required in the validation step. In fact,  
614 the matrix-based method displayed good statistical properties with as few as 20  
615 participants.

616 Since the validation testing only concerned problems that are probably less frequent, one  
617 could question the need to use methods that account for a heterogeneous probability of  
618 problem detection. In fact, problems are expected to be “homogeneously rare”. To the best  
619 of our knowledge, however, the assumption of homogeneity for rare problems has no  
620 theoretical or experimental basis. Furthermore, human factor engineers will define the  
621 usability testing’s experimental conditions according to the risk analysis, in order to  
622 facilitate the detection of problems previously described in the literature. If an engineer  
623 suspects the existence of problem removing the cap from an adrenaline pen, he/she might  
624 choose to evaluate the device in a more realistic test environment (e.g. with an actor  
625 pretending to go into anaphylactic shock); the problem is more likely to occur there than  
626 in a quiet, low-fidelity environment. By making some problems more detectable, the  
627 human factor engineer might introduce a degree of heterogeneity into the discovery  
628 process.

629 The choice of method was even more obvious for “formative” testing. In our simulations,  
630 the “formative” testing corresponds to a setting in which usability problems are frequent  
631 and numerous. Schmettow’s usability testing of a medical infusion pump is also an  
632 example of a formative assessment because it was followed by a redesign. Here, we  
633 proved that matrix-based methods are more reliable and have low bias and high  
634 consistency. As in the case of the infusion pump, a reliable estimate from a small number  
635 of participants is an economic advantage for the manufacturer, who can shorten redesign

636 cycles, accelerate device development, and hasten market access. The matrix-based  
637 method met this requirement because it required the fewest participants to guarantee  
638 good statistical properties. Another strength of the matrix-based method is its ability to  
639 embed previous knowledge through the prior parameters. Indeed, we used weakly  
640 informative priors for  $\mu$  and  $\sigma$  to avoid introducing information that we did not have  
641 about the medical device in question. However, one could take advantage of prior  
642 knowledge from earlier stages in device development or from a formative usability  
643 assessment to increase the accuracy of the estimate, especially when the sample size is  
644 small (i.e. an early control strategy). This approach is actually encouraged by regulatory  
645 bodies for medical device clinical trials [42] and helps to reduce the overall sample size.

646 Although we have suggested a threshold of 20 participants as the minimum sample size  
647 for obtaining a reliable estimate with the matrix-based method, we do not consider this  
648 to be the final threshold or a “magic number”. Indeed, as suggested by various researchers,  
649 the estimation models should be run iteratively as the sample size increases [4]. Thus,  
650 estimation models constitute a means of controlling and ensuring quality in formative  
651 testing and should not solely be considered as a checkpoint for validation testing.

652 Although the matrix-based method was more reliable, the LNBzt method could be used to  
653 double check the estimates - especially when high dispersion and/or the presence of very  
654 rare problems is suspected. Indeed, the LNBzt method’s coverage probability is high, and  
655 the overestimation bias makes it a conservative method that could usefully prevent the  
656 usability testing from being stopped too early.

## 657 V. Conclusions

658 Estimation models (and particularly matrix-based models) are of value in estimating and  
659 monitoring the detection process during usability testing. Matrix-based models have a

660 solid mathematical grounding and, with a view to facilitating the decision-making process  
661 for both regulators and device manufacturers, should be incorporated into current  
662 standards. To this end, the step-by-step tutorial provided here should facilitate the  
663 practical use of the matrix-based method in the evaluation of medical devices.

## 664 VI. Abbreviations

665 LNBzt: logit normal binomial zero truncated

666 GT: Good-Turing

## 667 VII. Declarations

668 **Ethics approval and consent to participate:** Not applicable.

669 **Consent for publication:** Not applicable.

670 **Availability of data and materials:** The data and code supporting the conclusions of this  
671 article are available in the GitHub repository,  
672 <https://zenodo.org/badge/latestdoi/279117812>.

673 **Competing interests:** The authors declare that they have no competing interests.

674 **Funding:** This research was funded by the Swiss National Science Foundation (grant  
675 number: SNSF-164279) and the French Agence Nationale de la Recherche (grant number:  
676 ANR-15-CE36-0007). The funding bodies had no role in the design of the study and  
677 collections, analysis, and interpretation of data, and in writing the manuscript.

678 **Authors' contributions:** AC, BD and VV conceptualized and designed the study. AC, CD,  
679 RP and VV carried out the analysis. AC, AD, BD, CD, RP, SP and VV contributed to the  
680 interpretation of the results. AC and VV drafted the initial manuscript. AD, BD, CD, RP and  
681 SP critically reviewed and revised the manuscript. All authors approved the final

682 manuscript as submitted and have agreed both to be personally accountable for the  
683 author's own contributions and to ensure that questions related to the accuracy or  
684 integrity of any part of the work, even ones in which the author was not personally  
685 involved, are appropriately investigated, resolved, and the resolution documented in the  
686 literature.

687 **Acknowledgements:** We thank Cedric Bach, James Lewis, and Jakob Nielsen for making  
688 their datasets available. We would like to express our particular gratitude to Martin  
689 Schmettow for providing us with the full dataset from the infusion pump usability study.  
690 We also thank the HPC Computing Mésocentre of the University of Lille which provided  
691 us with the computing grid. Lastly, we thank Simone Borsci, James Lewis, and Martin  
692 Schmettow for their valuable comments on an early version of this manuscript, which  
693 helped us to significantly improve the content.

## 694 VIII. References

- 695 1. US-FDA: **Applying human factors and usability engineering to medical**  
696 **devices: Guidance for industry and Food and Drug Administration staff.**  
697 *Washington, DC: FDA* 2016.
- 698 2. UK-MHRA: **Human Factors and Usability Engineering – Guidance for Medical**  
699 **Devices Including Drug-Device Combination Products.** In. Edited by Agency  
700 MHpR; 2017.
- 701 3. US-FDA: **Medical device recall report FY2003 to FY2012.** *Center for Devices and*  
702 *Radiological Health* 2012.
- 703 4. Borsci S, Macredie RD, Barnett J, Martin J, Kuljis J, Young T: **Reviewing and**  
704 **extending the five-user assumption: a grounded procedure for interaction**  
705 **evaluation.** *ACM Transactions on Computer-Human Interaction (TOCHI)* 2013,  
706 **20(5):1-23.**
- 707 5. Borsci S, Macredie RD, Martin JL, Young T: **How many testers are needed to**  
708 **assure the usability of medical devices?** *Expert Rev Med Devices* 2014,  
709 **11(5):513-525.**
- 710 6. Lewis JR: **Sample sizes for usability studies: Additional considerations.** *Human*  
711 *factors* 1994, **36(2):368-378.**
- 712 7. Kanis H: **Estimating the number of usability problems.** *Appl Ergon* 2011,  
713 **42(2):337-347.**

- 714 8. Lewis JR: **Evaluation of Procedures for Adjusting Problem-Discovery Rates**  
715 **Estimated From Small Samples.** *International Journal of Human-Computer*  
716 *Interaction* 2001, **13**(4):445-479.
- 717 9. Hertzum M, Jacobsen NE: **The Evaluator Effect: A Chilling Fact About Usability**  
718 **Evaluation Methods.** *International Journal of Human-Computer Interaction* 2003,  
719 **15**(1):183-204.
- 720 10. Schmettow M: **Sample size in usability studies.** *Communications of the ACM* 2012,  
721 **55**(4):64-70.
- 722 11. Borsci S, Londei A, Federici S: **The Bootstrap Discovery Behaviour (BDB): a new**  
723 **outlook on usability evaluation.** *Cogn Process* 2011, **12**(1):23-31.
- 724 12. Faulkner L: **Beyond the five-user assumption: Benefits of increased sample**  
725 **sizes in usability testing.** *Behavior Research Methods, Instruments, & Computers*  
726 2003, **35**(3):379-383.
- 727 13. Lewis JR: **Using discounting methods to reduce overestimation of p in**  
728 **problem discovery usability studies.** In.: Citeseer; 2000.
- 729 14. Sauro J, Lewis JR: **Quantifying the user experience: Practical statistics for user**  
730 **research:** Morgan Kaufmann; 2016.
- 731 15. Thomas DG, Gart JJ: **Small sample performance of some estimators of the**  
732 **truncated binomial distribution.** *Journal of the American Statistical Association*  
733 1971, **66**(333):169-177.
- 734 16. Virzi RA: **Refining the test phase of usability evaluation: How many subjects**  
735 **is enough?** *Human factors* 1992, **34**(4):457-468.
- 736 17. Nielsen J, Landauer TK: **A mathematical model of the finding of usability**  
737 **problems.** In: *Proceedings of the INTERACT'93 and CHI'93 conference on Human*  
738 *factors in computing systems: 1993*; 1993: 206-213.
- 739 18. Good IJ: **The population frequencies of species and the estimation of**  
740 **population parameters.** *Biometrika* 1953, **40**(3-4):237-264.
- 741 19. Schmettow M: **Controlling the usability evaluation process under varying**  
742 **defect visibility.** In: *Proceedings of the 23rd British HCI Group Annual Conference*  
743 *on People and Computers: Celebrating People and Technology: 2009*; British  
744 Computer Society; 2009: 188-197.
- 745 20. Finney D: **The truncated binomial distribution.** *Annals of Eugenics* 1947,  
746 **14**(1):319-328.
- 747 21. Rider PR: **Truncated binomial and negative binomial distributions.** *Journal of*  
748 *the American Statistical Association* 1955, **50**(271):877-883.
- 749 22. Shah S: **The asymptotic variances of method of moments estimates of the**  
750 **parameters of the truncated binomial and negative binomial distributions.**  
751 *Journal of the American Statistical Association* 1961, **56**(296):990-994.
- 752 23. Schmettow M: **Heterogeneity in the usability evaluation process.** *People and*  
753 *Computers XXII Culture, Creativity, Interaction* 22 2008:89-98.
- 754 24. Caulton DA: **Relaxing the homogeneity assumption in usability testing.**  
755 *Behaviour & Information Technology* 2001, **20**(1):1-7.
- 756 25. Schmettow M, Vos W, Schraagen JM: **With how many users should you test a**  
757 **medical infusion pump? Sampling strategies for usability tests on high-risk**  
758 **systems.** *J Biomed Inform* 2013, **46**(4):626-641.
- 759 26. DasGupta A, Rubin H: **Estimation of binomial parameters when both n, p are**  
760 **unknown.** *Journal of Statistical Planning and Inference* 2005, **130**(1-2):391-404.
- 761 27. Fisher RA: **The negative binomial distribution.** *Annals of Eugenics* 1941,  
762 **11**(1):182-187.

- 763 28. Haldane JB: **The fitting of binomial distributions**. *Annals of Eugenics* 1941,  
764 **11**(1):179-181.
- 765 29. Carroll RJ, Lombard F: **A note on N estimators for the binomial distribution**.  
766 *Journal of the American Statistical Association* 1985, **80**(390):423-426.
- 767 30. Olkin I, Petkau AJ, Zidek JV: **A comparison of n estimators for the binomial**  
768 **distribution**. *Journal of the American Statistical Association* 1981, **76**(375):637-  
769 642.
- 770 31. Hall P: **On the erratic behavior of estimators of N in the binomial N, p**  
771 **distribution**. *Journal of the American Statistical Association* 1994, **89**(425):344-  
772 352.
- 773 32. Robert C: **The Bayesian choice: from decision-theoretic foundations to**  
774 **computational implementation**: Springer Science & Business Media; 2007.
- 775 33. Meng X-L, Wong WH: **Simulating ratios of normalizing constants via a simple**  
776 **identity: a theoretical exploration**. *Statistica Sinica* 1996:831-860.
- 777 34. Bach C, Scapin DL: **Comparing inspections and user testing for the evaluation**  
778 **of virtual environments**. *Intl Journal of human-computer interaction* 2010,  
779 **26**(8):786-824.
- 780 35. Lewis JR, Henry SC, Mack RL: **Integrated office software benchmarks: A case**  
781 **study**. In: *Interact: 1990*; 1990: 337-343.
- 782 36. Nielsen J, Molich R: **Heuristic evaluation of user interfaces**. In: *Proceedings of the*  
783 *SIGCHI conference on Human factors in computing systems: 1990*; 1990: 249-256.
- 784 37. Team SD: **RStan: the R Interface to Stan. R package version 2.17.3**. In.; 2018.
- 785 38. Gronau QF, Singmann H, Wagenmakers E-J: **Bridgesampling: An R package for**  
786 **estimating normalizing constants**. *arXiv preprint arXiv:171008162* 2017.
- 787 39. Woolrych A, Cockton G: **Why and when five test users aren't enough**. In:  
788 *Proceedings of IHM-HCI 2001 conference: 2001*: Eds)(Cépaduès Editions, Toulouse,  
789 FR, 2001); 2001: 105-108.
- 790 40. Green PJ: **Reversible jump Markov chain Monte Carlo computation and**  
791 **Bayesian model determination**. *Biometrika* 1995, **82**(4):711-732.
- 792 41. Bias RG, Mayhew DJ: **Cost-justifying usability: An update for the Internet age**:  
793 Elsevier; 2005.
- 794 42. US-FDA: **Guidance for the use of Bayesian statistics in medical device clinical**  
795 **trials**. *Maryland: US Food and Drug Administration* 2010.

## 796 IX. Figures

### 797 A. Figure 1

798 **Title:** Bias in the prediction of  $m$ : the mean error and 95% fluctuation interval (as a  
799 percentage of the true  $m$ ) as a function of the sample size ( $n$ ).

800 **Legend:** The results are presented for various probabilities of problem detection ( $(\mu, \sigma)$ ,  
801 columns) and various numbers of usability problems ( $m$ , rows). The dashed line  
802 represents the true  $m$ .

803           B.     Figure 2

804     **Title:** Consistency in the prediction of  $m$ : the RMSE for the prediction of  $m$  (as a  
805     percentage of the true  $m$ ) as a function of the sample size ( $n$ ).

806     **Legend:** The results are presented for various probabilities of problem detection ( $(\mu, \sigma)$ ,  
807     columns) and various numbers of usability problems ( $m$ , rows). The LNBzt results are not  
808     represented for  $m < 100$  and  $\mu = \text{logit}(0.1)$ , due to a high RMSE

809     X.     Supplementary Material

810           A.     Additional file 1

811     File name: tutorial.pdf

812     File format: .pdf

813     Title of data: Step by step instructions for the matrix-based method presented in this  
814     manuscript.

815     Description of data: Open the file "*tutorial.pdf*" and follow the instructions.

816           B.     Additional file 2

817     File name: Supplementary Tables.xlsx

818     File format: .xlsx

819     Title of data: Tabulated version of the results of the simulation study presented in the  
820     main manuscript

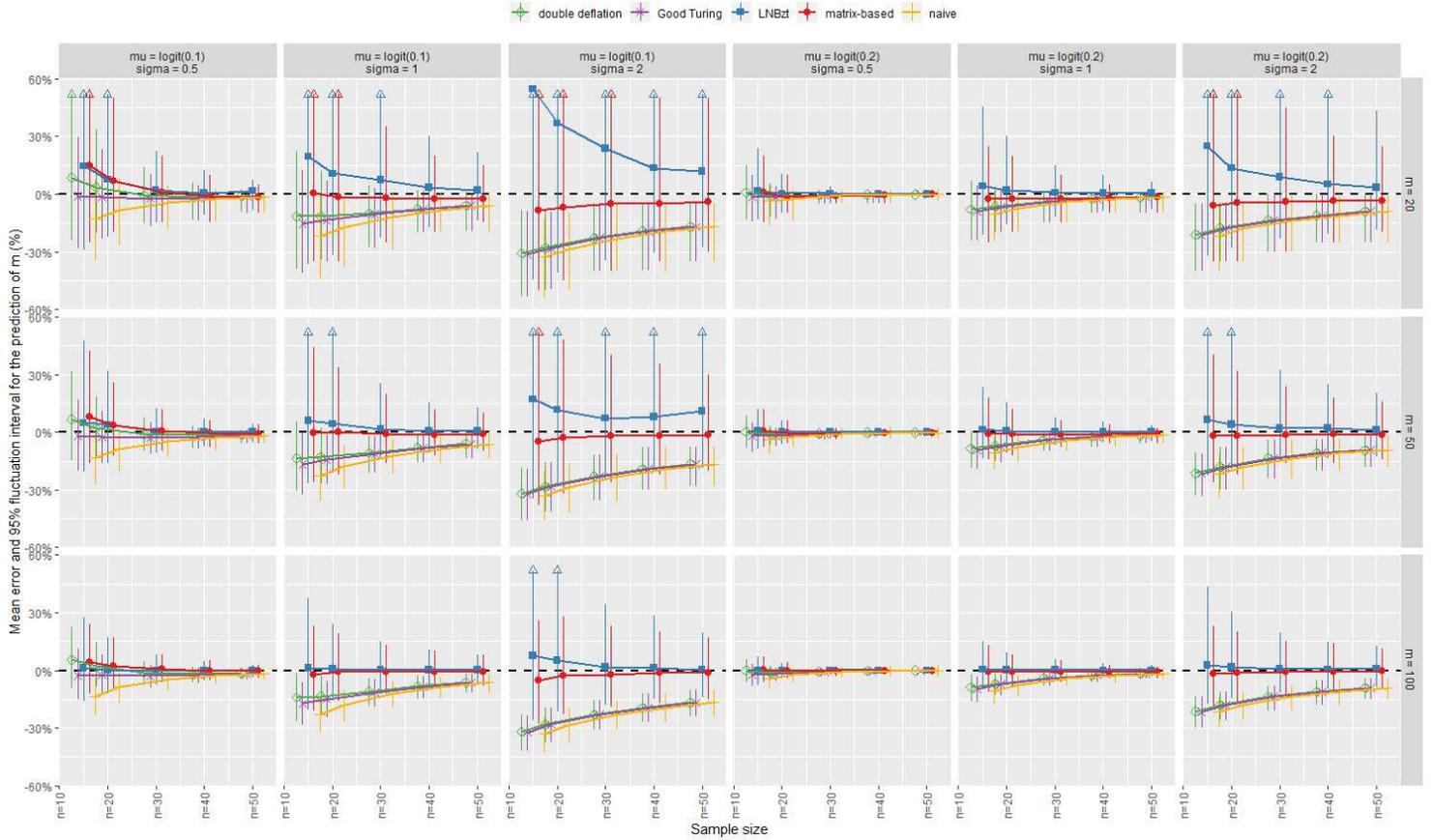
821     Description of data:

- 822           • S-Table 5: Bias in the prediction of  $m$ : the mean error (as a percentage of the true  
823            $m$ ) as a function of the sample size ( $n$ ).
- 824           • S-Table 6: Consistency in the prediction of  $m$ : the RMSE for the prediction of  $m$  as  
825           a function of the sample size ( $n$ ).

826  
827

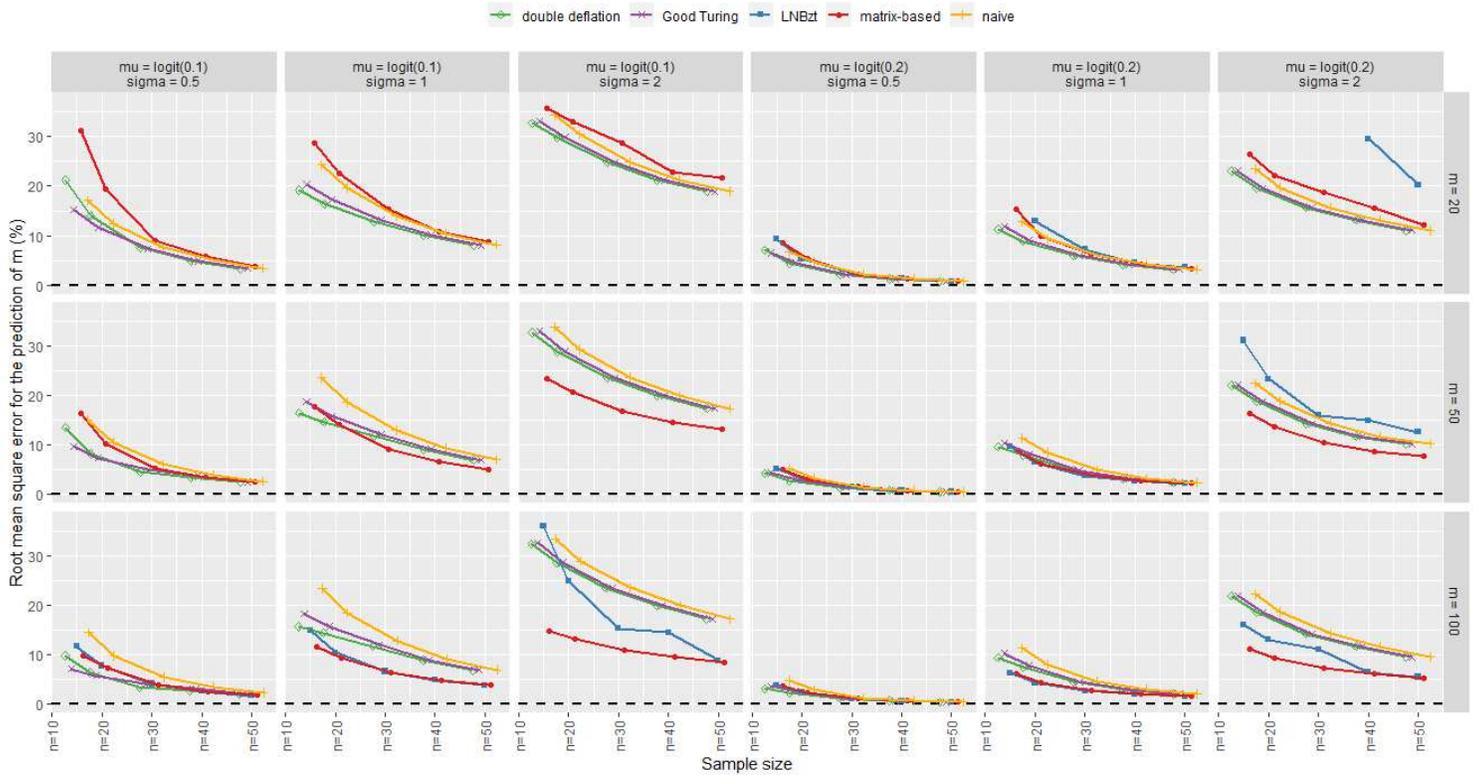
- S-Table 7: Coverage probability (in % of the 95% confidence interval) of  $\hat{m}$  with each combination  $(m, \mu, \sigma, n)$ .

# Figures



**Figure 1**

Bias in the prediction of  $m$ : the mean error and 95% fluctuation interval (as a percentage of the true  $m$ ) as a function of the sample size ( $n$ ). The results are presented for various probabilities of problem detection ( $(\mu, \sigma)$ , columns) and various numbers of usability problems ( $m$ , rows). The dashed line represents the "true"  $m$ .



**Figure 2**

Consistency in the prediction of  $m$ : the RMSE for the prediction of  $m$  (as a percentage of the true  $m$ ) as a function of the sample size ( $n$ ). The results are presented for various probabilities of problem detection ( $(\mu, \sigma)$ , columns) and various numbers of usability problems ( $m$ , rows). The LNBzt results are not represented for  $m < 100$  and  $\mu = \text{logit}(0.1)$ , due to a high RMSE.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [AdditionalFile3.rar](#)
- [AdditionalFile2.rar](#)
- [SupplementaryTables.xlsx](#)