

1 **Relaxation of the parameter independence assumption in the**
2 **bootComb R package**

3 Marc Y. R. Henrion^{1,2}

4 ¹ Malawi Liverpool Wellcome Programme, Blantyre, Malawi

5 ² Liverpool School of Tropical Medicine, Liverpool, UK

6

7

8

9

10

11

12

13 **Key words**

14 Biostatistics, R, confidence intervals, bootstrap, estimation

15 **Word count**

16 **Abstract:** 142 words

17 **Main text** (excluding abstract, references): 1,738 words

This peer-reviewed article has been accepted for publication but not yet copyedited or typeset, and so may be subject to change during the production process. The article is considered published and may be cited using its DOI.

This is an Open Access article, distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivatives licence (<http://creativecommons.org/licenses/by-nc-nd/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is unaltered and is properly cited. The written permission of Cambridge University Press must be obtained for commercial re-use or in order to create a derivative work.

18 Abstract

19 **Background** The bootComb R package allows researchers to derive confidence intervals
20 with correct target coverage for arbitrary combinations of arbitrary numbers of
21 independently estimated parameters. Previous versions (< 1.1.0) of bootComb used
22 independent bootstrap sampling and required that the parameters themselves are
23 independent - an unrealistic assumption in some real-world applications.

24 **Findings** Using Gaussian copulas to define the dependence between parameters, the
25 bootComb package has been extended to allow for dependent parameters.

26 **Implications** The updated bootComb package can now handle cases of dependent
27 parameters, with users specifying a correlation matrix defining the dependence structure.
28 While in practice it may be difficult to know the exact dependence structure between
29 parameters, bootComb allows running sensitivity analyses to assess the impact of parameter
30 dependence on the resulting confidence interval for the combined parameter.

31 **Availability** bootComb is available from the Comprehensive R Archive Network
32 (<https://CRAN.R-project.org/package=bootComb>).

33 Introduction

34 The bootcomb R package Henrion (2021) was recently published. This package for the
35 statistical computation environment R (R Core Team, 2021) allows researchers to derive
36 confidence intervals (CIs) with correct coverage for combinations of independently
37 estimated parameters. Important applications include adjusting a prevalence for estimated

38 test sensitivity and specificity (e.g. Mandolo et al. (2021)) or combining conditional
39 prevalence estimates (e.g. Stockdale et al. (2020)).

40 Briefly, for each of the input parameters, `bootComb` finds a best-fit parametric distribution
41 based on the confidence interval for that parameter estimate. `bootComb` then uses the
42 parametric bootstrap to sample many sets of parameter estimates from these best-fit
43 distributions and computes the corresponding combined parameter estimate for each set.
44 This builds up an empirical distribution of parameter estimates for the combined parameter.
45 Finally, `bootComb` uses either the percentile or the highest density interval method to derive
46 a confidence interval for the combined parameter estimate. Full details of the algorithm are
47 given in Henrion (2021).

48 A key point of the algorithm is that the best-fit distributions for the different parameters are
49 sampled from independently. This requires the parameters to be independent. This may not
50 be a realistic assumption in some real-world applications.

51 While for most practical applications the input parameters are typically estimated from
52 independent experiments (otherwise the combined parameter could be directly estimated),
53 the parameters themselves may not be independent. This is for instance the case when
54 adjusting a prevalence for the diagnostic test's sensitivity and specificity. The latter two
55 parameters are not independent: higher sensitivity can be achieved by lowering specificity
56 and vice versa.

57 If the experiments estimating these parameters are sufficiently large, then the violation of
58 the assumption of parameter independence may only have negligible impact on the resulting
59 confidence interval for the combined parameter. However, for the sake of general

60 applicability and to allow running sensitivity analyses, the author felt it was beneficial to
 61 extend bootComb to handle dependent parameters.

62 **Methods**

63 Copulas are multivariate distribution functions where the marginal probability distribution
 64 of each variable is the uniform distribution on the interval $[0,1]$. Copulas allow to specify the
 65 intercorrelation between random variables. An important probability theory result, Sklar's
 66 Theorem (Sklar, 1959), states that any multivariate probability distribution can be expressed
 67 in terms of its univariate marginal distributions and a copula defining the dependence
 68 between the variables.

69 Mathematically, let X_1, X_2, \dots, X_d be d random variables and define $U_i = F_i(X_i), i = 1, \dots, d$.
 70 Then the copula C of (X_1, \dots, X_d) is defined as the joint cumulative distribution function of
 71 (U_1, \dots, U_d) :

$$72 \quad C(u_1, \dots, u_d) = Pr(U_1 \leq u_1, \dots, U_d \leq u_d)$$

73 Assume that the marginal distributions, $F_i(x) = Pr[X_i \leq x], i = 1, \dots, d$ are continuous.
 74 Then, via the probability integral transform (Angus, 1994), the random vector (U_1, U_2, \dots, U_d)
 75 has marginals that are uniformly distributed on $[0,1]$.

76 bootComb makes use of the fact that the above can be reversed: given a sample (u_1, \dots, u_d) ,
 77 a sample for (X_1, \dots, X_d) can be obtained by $(x_1, \dots, x_d) = (F_1^{-1}(u_1), \dots, F_d^{-1}(u_d))$. The inverse
 78 functions $F_i^{-1}(u)$ will be defined if the marginals $F_i(x)$ are continuous. For the use of

79 bootComb, where users input confidence intervals for an estimated numeric parameter, this
80 will always be the case.

81 bootComb will proceed as follows to generate samples from a multivariate distribution of d
82 dependent variables:

- 83 • Estimate best-fit distributions F_1, \dots, F_d for each of the d parameters X_1, \dots, X_d given
84 the lower and upper limits of the estimated confidence intervals for each parameter.
- 85 • Sample (z_1, \dots, z_d) from a multivariate normal distribution $\mathcal{N}(\mathbf{0}, \Sigma)$ where the
86 variances in Σ are all 1.
- 87 • Since the marginals of this normal distribution are all $\mathcal{N}(0,1)$, compute $u_i = \Phi(z_i)$
88 where Φ is the cumulative distribution function of the standard normal.
- 89 • Finally, for each $i = 1, \dots, d$, compute $x_i = F_i^{-1}(u_i)$ where F_i is the best-fit marginal
90 distribution of parameter i .

91 The resulting vector (x_1, \dots, x_d) will be a sample from the multivariate distribution of
92 (X_1, \dots, X_d) . Note that the dependence structure was completely specified through the
93 covariance matrix Σ (since the variances are assumed to be 1, this really is a correlation
94 matrix) and marginal distributions for each parameter were specified by $F_i, i = 1, \dots, d$.

95 Results

96 I repeat the 2 examples from Henrion (2021) here, but look at the effect of specifying a
97 dependence between the input parameters.

98 All examples below use the highest density interval (HDI) method (input argument
 99 `method="hdi"`) to derive the final confidence interval. Whether this or the percentile method
 100 is used is a user choice. The HDI derived interval will be the narrowest interval with the
 101 desired coverage and the probability density will always be higher within that interval than
 102 outside it. To note however that the HDI may not be a single interval but a set of intervals if
 103 the density is multimodal. In this case, the single interval returned by `bootComb` will be too
 104 wide. For this reason, users should always inspect the histogram of the sampled combined
 105 parameter when using the HDI method.

106 1. HDV prevalence in the general population

107 With an application to hepatitis D and B viruses (HDV and HBV respectively) from Stockdale
 108 et al. (2020), Henrion (2021) showed how to use `bootComb` to obtain a valid confidence
 109 interval for \hat{p}_{aHDV} , the prevalence of HDV specific immunoglobulin G antibodies (anti-HDV)
 110 in the general population.

111 HBV is a pre-condition for HDV and hence to derive \hat{p}_{aHDV} Stockdale et al. (2020), obtained
 112 estimates of the prevalence of surface antigen of the hepatitis B virus (HBsAg), $\hat{p}_{HBsAg} =$
 113 3.5%, and the conditional prevalence of anti-HDV given the presence of HBsAg,
 114 $\hat{p}_{aHDV|HBsAg} = 4.5\%$:

- 115 • $\hat{p}_{HBsAg} = 3.5\%$ with 95% CI (2.7%, 5.0%).
- 116 • $\hat{p}_{aHDV|HBsAg} = 4.5\%$ with 95% CI (3.6%, 5.7%).

117 Assuming these 2 parameters to be independent, Henrion (2021) derived a 95% confidence
 118 interval for the estimate $\hat{p}_{aHDV} = \hat{p}_{aHDV|HBsAg} \cdot \hat{p}_{HBsAg}$ using `bootComb`, (0.11%, 0.25%).

119 If, however, the 2 input prevalences are not independent, e.g. if anti-HDV is more common
 120 among people with presence of HBsAg the higher the population prevalence of HBsAg is,
 121 then that assumption of independence would not hold. We can investigate how strong an
 122 effect dependence of the parameters can have on the resulting confidence estimate. For
 123 example, let us run the same example using `bootComb` with specifying the following
 124 covariance matrix for the bivariate normal copula:

$$125 \quad \Sigma = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$$

```
126 library(bootComb)
127
128 combFunEx<-function(pars){pars[[1]]*pars[[2]]}
129 bootComb(distributions=c("beta", "beta"),
130         qLowVect=c(0.027,0.036),
131         qUppVect=c(0.050,0.057),
132         combFun=combFunEx,
133         Sigma=matrix(byrow=TRUE,ncol=2,c(1,0.5,0.5,1)),
134         doPlot=TRUE,
135         method="hdi",
136         N=1e6,
137         seed=123)
```

138 This yields the 95% confidence interval (0.10%, 0.26%), a slightly wider interval – which
 139 makes sense, as the positive correlation means it is more likely for pairs of bootstrapped
 140 input parameters to be both near the upper (respectively lower) end of their confidence
 141 intervals.

142 For this particular application, a dependence between both prevalence parameters, \hat{p}_{HBsAg}
 143 and $\hat{p}_{aHDV|HBsAg}$, is unlikely and I have therefore not considered this example any further.

144 2. SARS-CoV-2 seroprevalence adjusted for test sensitivity and specificity

145 Henrion (2021) gave an example of adjusting an estimated SARS-CoV-2 seroprevalence for
 146 the estimated sensitivity and specificity of the test assay. Specifically:

- 147 • 84 out of 500 study participants tested positive for SARS-CoV-2 antibodies, yielding
 148 a seroprevalence estimate $\hat{\pi}_{raw} = 16.8\%$ with exact binomial 95% CI
 149 (13.6%, 20.4%).
- 150 • Estimated assay sensitivity: 238 out of 270 known positive samples tested positive
 151 $\hat{p}_{sen} = 88.1\%$, 95% CI (83.7%, 91.8%).
- 152 • Estimated assay specificity: 82 out of 88 known negative samples tested negative
 153 $\hat{p}_{spec} = 93.2\%$, 95% CI (85.7%, 97.5%).

154 Assuming the sensitivity and specificity to be independent, Henrion (2021) reported an
 155 adjusted seroprevalence estimate $\hat{\pi} = 12.3\%$ with 95% CI (3.9%, 19.0%).

156 However in this case, the assumption of independence is not fully realistic: there is a trade-
 157 off between sensitivity and specificity of the test assay, and as such one would expect a
 158 negative dependence between the two parameters: sensitivity can be increased at the cost
 159 of decreased specificity and vice versa.

160 Assuming that the sensitivity and specificity are negatively correlated with the copula
 161 correlation parameter $\rho = -0.5$ between these two parameters, using the extension of
 162 `bootComb` we can now account for the dependence of the parameters:

```
163 adjPrevSensSpecCI(  

  164   prevCI=c(0.136,0.204),  

  165   sensCI=c(0.837,0.918),
```

```

166 specCI=c(0.857,0.975),
167 Sigma=matrix(byrow=TRUE,ncol=3,c(1,0,0,0,1,-0.5,0,-0.5,1)),
168 doPlot=TRUE,
169 prev=84/500,
170 sens=238/270,
171 spec=82/88,
172 seed=123)

```

173 The reported confidence interval is now (3.8%, 19.4%) - marginally wider than when the
 174 dependence was ignored.

175 If we additionally specify `returnBootVals=TRUE` in the function call, we can extract and plot
 176 the sampled pairs of sensitivity and specificity values to check the dependence structure.
 177 This is shown on Figure 1: as the correlation parameter ρ in the copula between the
 178 sensitivity and specificity is decreased from 0 to -1, the dependence between both
 179 parameters becomes more and more pronounced as one would expect.

180 This shows that a simple correlation matrix specified for the Gaussian copula results in this
 181 case in a non-trivial dependence structure between two beta-distributed variables,
 182 respecting the specified marginal distributions.

183 We can also visualise the effect on the estimated confidence interval, as shown on Figure 2.
 184 We can see that in this case, with a negative correlation, the width of the CI increases at the
 185 correlation becomes stronger. However, looking at the scale of the y-axis we see that this is
 186 just a marginal effect.

187 A more substantial effect of parameter dependence is obtained when we also allow the
 188 measured prevalence π_{raw} to be correlated with sensitivity (p_{sens} ; positive correlation) and
 189 specificity (p_{spec} ; negative correlation). Specifically, we can specify the following correlation
 190 matrix for the parameters $(\pi_{raw}, p_{sens}, p_{spec})$:

191

$$\Sigma = \begin{pmatrix} 1 & 0.3 & -0.3 \\ 0.3 & 1 & -0.5 \\ -0.3 & -0.5 & 1 \end{pmatrix}$$

192

193

194

195

196

197

198

199

200

201

```
adjPrevSensSpecCI(
  prevCI=c(0.136,0.204),
  sensCI=c(0.837,0.918),
  specCI=c(0.857,0.975),
  Sigma=matrix(byrow=TRUE,ncol=3,c(1,0.3,-0.3,0.3,1,-0.5,-0.3,-0.5,1)),
  doPlot=TRUE,
  prev=84/500,
  sens=238/270,
  spec=82/88,
  seed=123)
```

202

In this case, the reported confidence interval is (4.7%, 18.2%). This CI is 11% narrower than

203

when the dependence structure was ignored – a substantial effect for practical purposes.

204

Conclusions

205

The R package `bootComb` has been extended and, using Gaussian copulas, it can now handle

206

the case of dependent input parameters. For many applications, the effect of dependence

207

between the parameters will be marginal or even negligible, but this is not always the case.

208

The package now allows users to do sensitivity analyses to assess the effects of a miss-

209

specified dependence structure between the parameters that are being combined.

210

At the time of publication, the most recent version of `bootComb` was 1.1.2.

211

Figure captions

212

Figure 1: Scatterplots showing the bootstrapped values of sensitivity and specificity for

213

different strengths of dependence (from independence to perfect correlation) between

214 *sensitivity and specificity. The empirical kernel density estimate for the bivariate distribution in*
215 *each case is shown as orange contour lines.*

216

217 **Figure 2:** *Width of the estimated confidence interval as a function of increased strength of the*
218 *negative correlation between sensitivity and specificity.*

219 **Funding Information (see funding information section for more**
220 **information)**

221 This research was funded in whole, or in part, by the Wellcome Trust [grant:
222 206545/Z/17/Z]. For the purpose of open access, the author has applied a CC BY public
223 copyright licence to any Author Accepted Manuscript version arising from this submission.

224 **Data Availability Statement**

225 All data to support this work are contained within the article. The software package itself is
226 available from <https://cran.r-project.org/package=bootComb>.

227 **Conflicts of interest**

228 Author Marc Y. R. Henrion declares none.

229 **References**

- 230 Angus, J. E. (1994). The Probability Integral Transform and Related Results. *SIAM Review*,
231 36(4), 652–654. <http://www.jstor.org/stable/2132726>
- 232 Henrion, M. Y. (2022). *bootComb: Combine Parameter Estimates via Parametric Bootstrap* (R
233 package version 1.1.2) [Computer software]. <https://cran.r-project.org/package=bootComb>
- 234 Henrion, M. Y. (2021). bootComb—an R package to derive confidence intervals for
235 combinations of independent parameter estimates. *International Journal of Epidemiology*,
236 50(4), 1071–1076. <https://doi.org/10.1093/ije/dyab049>
- 237 Mandolo, J. J., Henrion, M. Y. R., Mhango, C., Chinyama, E., Wachepa, R., Kanjerwa, O.,
238 Malamba-Banda, C., Shawa, I. T., Hungerford, D., Kamng'ona, A. W., Iturriza-Gomara, M.,
239 Cunliffe, N. A., & Jere, K. C. (2021). Reduction in Severity of All-Cause Gastroenteritis
240 Requiring Hospitalisation in Children Vaccinated against Rotavirus in Malawi. *Viruses*,
241 13(12), 2491. <https://doi.org/10.3390/v13122491>
- 242 R Core Team. (2021). *R: A language and environment for statistical computing*. R Foundation
243 for Statistical Computing. <https://www.R-project.org/>
- 244 Sklar, M. (1959). *Fonctions de Répartition À N Dimensions Et Leurs Marges* (Issue 8, pp. 229–
245 231). Publications de l'Institut Statistique de l'Université de Paris.
- 246 Stockdale, A. J., Kreuels, B., Henrion, M. Y. R., Giorgi, E., Kyomuhangi, I., de Martel, C., Hutin,
247 Y., & Geretti, A. M. (2020). The global prevalence of hepatitis D virus infection: Systematic
248 review and meta-analysis. *Journal of Hepatology*, S0168827820302208.
249 <https://doi.org/10.1016/j.jhep.2020.04.008>