

Comparison of Three Alternative Methods for Analysis of Equine Faecal Egg Count Reduction Test Data

Denwood, M.J.^{*,a}, Reid, S.W.J.^a, Love, S.^b, Nielsen, M.K.^c, Matthews, L.^a, McKendrick, I.J.^d, Innocent, G.T.^d

^a*Boyd Orr Centre for Population and Ecosystem Health, Institute of Comparative Medicine, Faculty of Veterinary Medicine, University of Glasgow, Bearsden Road, Glasgow, G61 1QH, UK.*

^b*Veterinary Companion Animal Sciences, Institute of Comparative Medicine, Faculty of Veterinary Medicine, University of Glasgow, Bearsden Road, Glasgow, G61 1QH, UK.*

^c*Department of Large Animal Sciences, Faculty of Life Sciences, University of Copenhagen, Højbakkegaard Alle 5, DK-2630 Taastrup, Denmark.*

^d*Biomathematics & Statistics Scotland (BioSS), The King's Buildings, Edinburgh, EH9 3JZ, UK.*

Abstract

The Faecal Egg Count Reduction Test (FECRT) is the most widely used method of assessing the efficacy of anthelmintics, and is the only *in vivo* technique currently approved for use with horses. Equine Faecal Egg Count (FEC) data are frequently characterised by a low mean, high variability, small sample sizes and frequent zero observations. Accurate analysis of the data therefore depends on the use of an appropriate statistical technique. Analyses of simulated FECRT data by methods based on calculation of the empirical mean and variance, non-parametric bootstrapping, and Markov chain Monte Carlo (MCMC) were compared. The MCMC method consistently outperformed the other methods, independently of the sample size and distribution from which the data were generated. Bootstrapping produced notional 95% confidence intervals containing the true parameter as little as 40% of the time with sample sizes of less than 50. Analysis of equine FECRT data yielded inconclusive results in 53 of 63 (84%) datasets, suggesting that the routine use of prior sample size calculations should be adopted to ensure sufficient data are collected. The authors conclude that computationally intensive parametric methods such as MCMC should be used for analysis of FECRT data with sample sizes of less than 50, in order to avoid making erroneous inference about the true efficacy of anthelmintics in the field. Software to perform all three types of analyses documented here is freely available in the form of an add-on package to the R statistical programming language from <http://cran.r-project.org/web/packages/bayescount/index.html>.

Key words: MCMC, bootstrap, WAAVP, FECRT, anthelmintic efficacy, equine

*Corresponding author

Email address: m.denwood@vet.gla.ac.uk (Denwood, M.J.)

1. Introduction

The Faecal Egg Count Reduction Test (FECRT) is the most widely used method of assessing the *in vivo* efficacy of anthelmintics against parasitic nematodes of horses, sheep and cattle (Coles et al., 2006; Kaplan, 2002), and is an essential tool in the process of monitoring the increasing prevalence of anthelmintic resistance. The test is known to have several limitations, including the variability of Faecal Egg Count (FEC) data (Uhlinger, 1993), leading to a relatively variable FECRT result (Miller et al., 2006). This is especially true in equine FEC data, where effects such as differing age related immunity (Klei and Chapman, 1999) and differences in grazing management (Dopfer et al., 2004) impact on the observed FEC. Combined with the small group sizes and frequent zero FEC observations (Kaplan, 2002; Nielsen et al., 2006) often encountered with horses, this high variability between animals and low mean FEC introduce difficulties in analysis of equine FECRT data which do not arise to the same extent in analysis of FECRT data obtained from cattle or sheep.

The method currently advocated by the World Association for the Advancement of Veterinary Parasitology (WAAVP) involves calculating the empirical mean and variance before and after treatment, and calculating the empirical mean reduction and estimates of the 95% confidence interval for the true reduction using these figures (Coles et al., 1992). This method takes no account of the difference between uncertainty regarding the true mean of a sample, introduced by the Poisson variability of the counting process, and variability in the true mean of different samples. Calculation of 95% confidence intervals in this manner also assumes that the distribution of error for the mean is symmetrical on the log scale, although parameter likelihoods (and therefore errors) have been reported to be skewed for FEC data (Denwood et al., 2008), potentially justifying this assumption.

A non-parametric bootstrapping approach has recently been suggested as an appropriate method to generate confidence limits from equine FECRT data (Vidyashankar et al., 2007). The technique involves re-sampling and summarising the observed data, and makes no assumptions about the underlying distribution or processes generating the data (Mooney and Duval, 1993), or the parameter error structure. Non-parametric bootstrapping approaches are therefore widely used and extremely useful when the underlying distribution of data is unknown. A fundamental assumption underlying this approach is that the data obtained are fully representative of the population, an assumption which risks being violated when dealing with small sample sizes, giving misleading results. A non-parametric bootstrapping approach is more complex and time consuming than the currently advocated WAAVP method, however the use of facilities like Excel spreadsheet macros and basic computer programs potentially allow different data to be analysed relatively quickly.

Alternative options for analysis of FECRT data include computationally intensive parametric methods. These include parametric bootstrapping, or the likelihood profiling method proposed by Torgerson et al.

35 (2005), but here Markov chain Monte Carlo (MCMC) (Gilks et al., 1998) is used as an example. Each
36 of these methods requires the use of a parametric distribution in order to describe the FEC data. The
37 negative binomial is the most frequently used parametric distribution for FEC data, and is equivalent to the
38 gamma-Poisson compound distribution implemented here (for the derivation see Vose (2004)). Conceptually,
39 this represents a population of Poisson distributions with gamma distributed means, where the Poisson
40 distributions account for counting variability in observed FEC within a sample, and the gamma distribution
41 describes the variability between samples. The latter could arise as a combination of several factors, including
42 the aggregated distribution of eggs in faeces, variations in worm fecundity over time, variations in faecal
43 consistency, and variations in the numbers of worms present, which are impossible to separate using only
44 a single faecal sample per individual. For the MCMC model, pre-treatment data are assumed to follow
45 a single gamma-Poisson (negative binomial) distribution, while post treatment data are distributed as a
46 different gamma-Poisson distribution, with a mean value which has been scaled relative to the pre-treatment
47 mean, and a value for variability which has separately been scaled relative to the pre-treatment variability.
48 This allows inference on the true change in mean egg shedding, with an additional parameter reflecting
49 the true change in variability between egg counts. From this model, estimates of the mean anthelmintic
50 efficacy and the variability in anthelmintic efficacy between animals can be obtained. The advantage of an
51 MCMC based approach is that the different sources of variability can be taken into consideration, leading
52 to more accurate estimates of the uncertainty of true parameter estimates. Disadvantages of this approach
53 include the comparatively high computational effort required to implement the method, and the need to
54 make distributional assumptions about the processes generating the data. FEC between animals is well
55 described by a gamma-Poisson (negative binomial) distribution, however alternatives include zero-inflated
56 distributions (Denwood et al., 2008; Nødtvedt et al., 2002), and the use of a lognormal distribution to
57 describe the variability in means (Morrison, 2004).

58 The bootstrapping and MCMC procedures also have the advantage of attempting to define the full
59 distribution of likely values for the true FEC reduction. This allows the results to be presented in a more
60 intuitive way, such as a probability that the true egg count reduction is less than a given percentage, typically
61 the published efficacy of the drug used. This probability, \hat{p} , is relatively easy to estimate using numerical
62 integration for both the bootstrapping and MCMC methods, and can allow the group to be classified as
63 ‘Susceptible’ if $0 < \hat{p} < 2.5\%$, ‘Possible resistant’ if $2.5 < \hat{p} < 50\%$, ‘Probable resistant’ if $50 < \hat{p} < 97.5\%$,
64 or ‘Confirmed resistant’ if $97.5 < \hat{p} < 100\%$. These definitions allow a distinction to be made between
65 confirmed resistance and the lack of evidence of susceptibility, which is lacking in the current interpretation
66 of lower 95% confidence interval and empirical mean reduction statistics described by Coles et al. (1992).

67 Given the worldwide importance of anthelmintic resistance, there is an urgent need to improve and
68 standardise the statistical method used to analyse such data (Coles et al., 2006; Kaplan, 2002). Compared
69 to the case with ruminants, the relatively small sample sizes, high variability between counts, and relatively

low pre-treatment mean FEC frequently encountered with equine FECRT data may provide a challenge to the use of a non-parametric bootstrapping procedure, since there are relatively few data points from which to sample. There is also often insufficient data to be able to analyse the underlying distribution, which prevents validation of the choice of distribution used by the MCMC analysis. The aim of this study was to assess the usefulness of 95% confidence intervals generated using these three methods using simulated data, and then to assess the impact of the assumptions being made for each method.

2. Materials and Methods

2.1. Statistical Analysis

The analysis currently recommended by the World Association for the Advancement of Veterinary Parasitology was performed as described by Coles et al. (1992). Bootstrapping was conducted using a function written by the author in the R statistical programming language (R Development Core Team, 2008). New pre-and post-treatment pseudo-datasets were sampled from each dataset, and the mean reduction calculated 10,000 times. The mean estimate and 95% confidence intervals for each dataset were then calculated and recorded from these 10,000 iterations.

Bayesian MCMC analysis was performed using a bespoke model, implemented using JAGS (Plummer, 2008) for the MCMC simulation. The model fits a gamma-Poisson distribution to the pre and post-treatment data, with parameters for pre- and post-treatment means and shape parameters. The pre-treatment mean and shape parameters are given minimally informative prior distributions spanning all values that are seen in real FECRT data for each parameter. Post-treatment mean and shape parameters are calculated by multiplying the pre-treatment mean and shape parameters by a “change in mean” and “change in shape” parameter, respectively. The “change in mean” is given an uninformative $Beta(1, 1)$ prior, and the “change in shape” a diffuse lognormal prior with a mean of one. The true % FEC reduction is derived from $(1 - change\ in\ mean) * 100$. Calling JAGS to run each simulation and summarising of MCMC chains was automated using the runjags package (Denwood, 2008) for R, with two chains. Convergence was assessed using the Gelman-Rubin statistic (Gelman and Rubin, 1992), and necessary sample size using Raftery and Lewis’s diagnostic (Raftery and Lewis, 1995). The median estimate and 95% credible intervals for the true egg count reduction were calculated in R using the MCMC output.

For all three methods, credible intervals for the proportion of datasets with the true reduction parameter contained within the nominal 95% confidence intervals were calculated using a Bayesian approach with an uninformative $Beta(1, 1)$ prior. The mean relative size of these confidence intervals was calculated using equation (1).

$$confidence\ interval\ size = \frac{\sum \frac{U-L}{T}}{N} \quad (1)$$

101 Where L denotes the lower confidence interval, U the upper confidence interval, T the true parameter value,
102 and N the number of datasets

103 To assess the accuracy of the median estimates, the relative root-mean-square-error (RMSE) was calcu-
104 lated using the simulated (true) value for each parameter. The RMSE can also be thought of as the standard
105 deviation of the ratio between each median estimate and the simulated values; however it should be noted
106 that this is not equivalent to the accepted meaning of the term “standard deviation”. The term relative
107 RMSE will be used to avoid confusion.

108 2.2. Comparisons of methods for analysis of FECRT data

109 A total of 1000 parameters for a simulated FECRT were generated in the R statistical programming
110 language. The true proportional FEC reduction was simulated from a *Uniform*(0.75, 1) distribution, so
111 that true egg count reductions varied from reduced efficacy to efficacious reductions. The pre-treatment
112 mean number of eggs counted (equal to FEC if the egg counting technique had an egg detection threshold
113 of 1 EPG), and sample size (number of animals) were chosen to reflect the values seen in real equine
114 FECRT data obtained from 63 typical Danish equine datasets. The 2.5% and 97.5% quantiles for observed
115 pre-treatment mean and sample size were used as the lower and upper bounds of the distributions used
116 to generate the parameters. Pre-treatment mean was taken from a *Uniform*(1.45, 53.1) distribution, and
117 sample size per group was sampled randomly from integers between 4 and 16 inclusive with each integer
118 having an equal probability of selection. The coefficient of variation (cv) between samples before treatment
119 was sampled from a *Uniform*(1, 1.41) distribution (corresponding to a pre-treatment shape parameter of
120 the gamma distribution, k , of between 1 and 0.5), and the proportional increase in cv after treatment was
121 sampled from the same distribution (corresponding to a post-treatment shape parameter of between $1 * 1 = 1$
122 and $0.5 * 0.5 = 0.25$). These values were also chosen to reflect the values most likely to be encountered in
123 real FECRT data; published values of k are usually less than one (Shaw et al., 1998), and differing efficacy
124 of anthelmintic between animals would be expected to result in an increase in variability post-treatment.

125 In order to test the implications of the distributional assumptions made by the MCMC and WAAVP
126 methods, simulated datasets were generated using the following three different distributions of underlying
127 sample means; gamma-Poisson (negative binomial), multi-modal lognormal-Poisson, and uniform-Poisson.
128 For each dataset, the meta-population mean and variance was the same for all distributions. The number
129 of modes for each multi-modal lognormal-Poisson distribution was sampled as between two and ten for
130 each dataset, and a separate lognormal distribution used to describe the distribution of modes within the
131 group. These modes conceptually represent sub-groups within the population, with the population variance
132 split equally between the two compound lognormal distributions for each animal. If the simulated parameter
133 mean and variance required negative parameter value for the lower limit of the uniform-Poisson distribution,
134 then a log-uniform distribution was used instead (that is, a distribution which is uniform on the log scale).

135 Pre- and post-treatment egg count data were generated using each of these three distributions with the 1000
136 parameter values, to simulate a FECRT for a total of 3000 datasets. These datasets were then analysed
137 using each of the three methods. More details regarding the generation of these data are available from the
138 corresponding author.

139 2.3. Equine FECRT data

140 The MCMC and bootstrap methods were applied to equine FECRT data obtained from 63 typical
141 Danish equine establishments, with a median (range) of 9 (6-22) animals per dataset. For these data, a
142 modified MCMC method using zero-inflated gamma-Poisson distributions in place of the gamma-Poisson
143 distributions was used, in addition to the MCMC method described previously. For each dataset, the
144 probability, \hat{p} , that the observed FEC reduction was less than the “desired” FEC reduction was calculated
145 by numerical integration of the posterior estimates for true FEC reduction. From this, the dataset was
146 classified as ‘Susceptible’ if $0 < \hat{p} < 2.5\%$, ‘Possible resistant’ if $2.5 < \hat{p} < 50\%$, ‘Probable resistant’ if
147 $50 < \hat{p} < 97.5\%$, or ‘Confirmed resistant’ if $97.5 < \hat{p} < 100\%$. In this study, the “desired” FEC reduction
148 was set 95%, corresponding to the best estimate of the efficacy of the drug used in a naïve population. This
149 figure represents the minimum population mean FEC reduction we would expect from a fully susceptible
150 group of animals if we were able to observe the true mean FEC before and after treatment, and could be
151 adjusted for both methods with other datasets depending on the drug used and desired tolerance in true
152 efficacy.

153 2.4. Bootstrap analysis

154 A more complex analysis of the performance of the bootstrapping method was performed using gamma-
155 Poisson data. Sample size was drawn from the set $\{5, 10, 20, 30, 40, 50, 60, 70, 80, 90 \text{ \& } 100\}$, and
156 pre-treatment mean number of eggs counted from the set $\{1, 5, 10, 20, 30, 40, 50, 75 \text{ \& } 100\}$. Each of these
157 99 combinations was used to generate 1000 datasets using two gamma-Poisson distributions and a true
158 FEC reduction randomly generated from a *Uniform*(0.75, 1) distribution. For each dataset, the parameter
159 value used for pre-treatment cv was either 1 or 1.41, and post-treatment change in cv either 1 or 1.41. Each
160 dataset was analysed using the bootstrap method to provide a median estimate and 95% confidence intervals
161 as before.

162 3. Results

163 3.1. Comparisons of methods for analysis of FECRT data

164 Of the 3000 datasets, 35 of the gamma-Poisson datasets, 32 of the multi-modal lognormal-Poisson
165 datasets, and 33 of the (log) Uniform-Poisson datasets gave an empirical reduction of 100%. The me-
166 dian (95% confidence interval) simulated true reduction for these empirical 100% reduction datasets was

167 99.13% (82.23% - 99.97%). As the post-treatment variance for these datasets was 0, the WAAVP method
168 of calculating 95% confidence intervals could not be applied. In practice, these datasets would be assumed
169 to represent a 100% reduction, so 95% confidence limits of 100% to 100% were assigned to these datasets.
170 The non-parametric bootstrapping approach generated the same confidence limits for these datasets, since
171 all possible combinations of datapoints give a 100% reduction.

172 In [Figure 1](#), the proportion of true reductions that were contained within the notional 95% confidence
173 intervals for each method with all datasets are shown (95% credible intervals calculated using a Bayesian
174 method with an uninformative prior). There is no evidence that the MCMC method did not estimate true
175 95% confidence intervals for both the gamma-Poisson and (log) Uniform data, but the confidence was lower
176 for the multi-modal data. Non-parametric bootstrapping and the WAAVP method both returned notional
177 95% confidence intervals that contained the true value between 85% and 90% of the time for all data types.
178 Discounting the datasets with an empirical reduction of 100% improved the apparent performance of the
179 bootstrapping and WAAVP methods, although both methods still generated lower estimates of confidence
180 than the MCMC method for all data types (data not shown).

181 In [Table 1](#), the mean relative size of the notional 95% confidence intervals for each method and dataset
182 are shown. The relative RMSE for each combination is shown in [Table 2](#). The MCMC method returned
183 on average slightly larger 95% confidence limits than the other methods for each dataset, although when
184 datasets with 100% apparent reductions were excluded, the three methods produce similarly sized 95% con-
185 fidence intervals (data not shown). The 95% confidence intervals were largest for the (log) Uniform-Poisson
186 data, and most narrow for the multi-modal data. The MCMC median estimates produced a lower relative
187 RMSE than the bootstrapping median and WAAVP mean estimates in every case. The bootstrapping me-
188 dian and WAAVP mean estimates generally had a similar relative RMSE, although those produced by the
189 bootstrapping method were lower. As for the relative size of 95% confidence intervals, the relative RMSE
190 was smallest for each method for the multi-modal data and largest for the (log) Uniform-Poisson data.

191 *3.2. Analysis of equine FECRT data*

192 The probabilities of resistance returned by the bootstrapping and modified MCMC method relative to the
193 first MCMC method are shown in [Figure 2](#). The probabilities were greater for MCMC than bootstrapping
194 in all but one case, indicating that bootstrapping consistently estimated the true efficacy to be higher than
195 the estimates produced by MCMC. Estimates produced by the modified MCMC method using the zero-
196 inflated gamma-Poisson distribution were very similar to those produced by the first MCMC method using
197 the uni-modal gamma-Poisson distribution.

198 In [Table 3](#), the classifications made for each dataset using each method are shown. None of the datasets
199 were classified as confirmed susceptible using the MCMC method, and of 14 (22%) classified as confirmed
200 susceptible with the bootstrap method, four (6%) were classified as ‘probable resistant’ using the MCMC

201 method. In addition, seven (11%) of the datasets were classified as ‘confirmed resistant’ using MCMC and
202 only ‘probable resistant’ using the bootstrap method. There was insufficient information in the data to
203 determine either confirmed resistance or susceptibility for 53 datasets (84%) using MCMC and 46 datasets
204 (73%) using the bootstrap method.

205 *3.3. Bootstrap analysis*

206 The effect of increasing pre-treatment mean FEC and sample size on the ability of the bootstrapping
207 method to accurately predict the true FEC reduction is shown in [Figure 3](#) and [Figure 4](#). As pre-treatment
208 mean FEC increased, the 95% confidence intervals were more reliable, although this affect appeared to be
209 less pronounced with an increase in mean above ten counted eggs at sample sizes 20 and greater. With
210 sample sizes of five and ten, the notional 95% confidence intervals contained the true parameter no more
211 than 90% of the time, and as little as 40% of the time with a very low mean FEC. At sample sizes 20
212 to 40, the 95% confidence intervals contained the true parameter between 90% and 95% of the time for
213 pre-treatment mean FEC of over ten counted eggs. This improved to between around 93% and 95% for
214 sample sizes of 50 and above with pre-treatment mean FEC of ten counted eggs and above. Even with a
215 sample size of 100, the notional 95% confidence intervals contained the true parameter between only 89%
216 and 93% of the time with a pre-treatment mean FEC of one egg counted, and between 92% and 95% with
217 a pre-treatment mean FEC of five eggs counted. Conversely, the confidence of the estimates produced by
218 the MCMC method were not decreased by a reduced mean and sample size, with notional 95% confidence
219 intervals containing the true value 99% of the time with a mean of 1 and sample size of 5, 97% of the time
220 with a mean of 100 and sample size of 5, 97% of the time with a mean of 1 and sample size of 100, and 96%
221 of the time with a mean of 100 and sample size of 100 (95% credible intervals not shown).

222 **4. Discussion**

223 For all datasets, simulated from each of the distributions tested, the MCMC method provided confidence
224 intervals with the best defined properties, as well as the most precise median estimates for the true FEC
225 reduction. The size of the 95% confidence intervals produced was slightly greater for the MCMC method,
226 but not when datasets with empirical reductions of 100% were removed. This indicates that the MCMC
227 methods were producing more appropriate 95% confidence intervals, rather than merely larger 95% confi-
228 dence intervals. This was the case not only for data simulated from a gamma-Poisson distribution, where
229 the MCMC method using the same distribution would be expected to perform well, but also using data
230 simulated from different distributions. The performance of the MCMC method was less optimal using the
231 multi-modal data, but even here it out-performed the other two methods. In addition, the modified MCMC
232 method (based on a zero-inflated gamma-Poisson distribution) produced much more similar results to the

233 first MCMC method (based on a uni-modal gamma-Poisson distribution) than the bootstrap procedure for
234 the analysis of equine FECRT data. This implies that the distributional assumptions made by the MCMC
235 method has less practical impact on the analysis of these types of FECRT data than the assumption that
236 bootstrapping a limited number of data points can capture all the variability of an inherently very variable
237 system. [Vidyashankar et al. \(2007\)](#) propose dealing with this effect by taking into account the inter-farm
238 variability. The intention of this paper was to assess the performance of each method when analysing in-
239 dividual datasets in the absence of any other comparable datasets, so that taking into account inter-farm
240 variability would not have been possible. The MCMC method is also capable of analysing data from multiple
241 sites, for example by defining a distribution of efficacy that describes the mean FEC reduction at each site
242 and using this extra information to reduce uncertainty in the estimate for the true mean efficacy. However,
243 by directly describing the variability structure in FEC data, parametric techniques eliminate the necessity
244 for data from additional sites (where none is available), and allow efficacy to be analysed at an individual
245 farm level.

246 Several of the datasets generated with parameters similar to observed equine FECRT data gave an
247 empirical reduction of 100%, even where the true mean reductions were close to 75%. These datasets present
248 difficulties when using both the WAAVP and bootstrap methods, which were unable to generate appropriate
249 95% confidence limits. Nineteen (19%) of these datasets were simulated using empirical reductions of less
250 than 95%, and so represent a consistent source of false negatives for these methods. The MCMC method
251 was the only method examined in this paper which is capable of analysing datasets with 100% empirical
252 reductions in an appropriate fashion.

253 It is also apparent from the analysis presented here that analysis of a single equine FECRT dataset will
254 often prove inconclusive. Using the MCMC method, only 10 of the datasets from equine field studies were
255 classified as ‘confirmed resistant’ and 0 as ‘confirmed susceptible’, with the remaining 53 (84%) datasets
256 containing insufficient information to be sure if the true drug efficacy was reduced or not. This is consistent
257 with the conclusions made by [Miller et al. \(2006\)](#), that the results of a FECRT based on an arithmetic
258 mean reduction can be inconsistent. The utility of the method could be increased by performing a suitable
259 sample size calculation prior to performing the FECRT, and increasing the number of samples taken and/or
260 reducing the egg detection threshold accordingly. This is probably not practical for routine clinical tests,
261 however, due to the added cost and time associated with taking more samples and counting more eggs. A
262 more useful solution might be to use a process control approach for routine surveillance, combined with
263 the use of a more detailed FECRT with prior sample size calculations to calculate the required number of
264 samples to take when the process control indicated a possible problem. This may represent both a more
265 efficient use of resources, and a greater overall diagnostic test sensitivity and specificity, than the current
266 use of repeated reduction tests viewed in isolation and without the necessary sample size calculations.

267 In this paper, the efficacy of reductions were classified according to the probability that the true reduction

268 was below a given threshold, which is not consistent with the method currently advocated by the WAAVP.
269 This departure was made to allow a distinction to be drawn between cases where there is clear evidence
270 for resistance and cases where there is insufficient evidence to demonstrate acceptable efficacy. Using the
271 classification scheme currently used by the WAAVP, which involves consideration of the mean estimate and
272 lower 95% confidence interval only (Coles et al., 1992), it is not possible to make this distinction in the
273 absence of suitable power calculations. This limitation may lead to confusion over the clinical interpretation
274 of FECRT analysis results.

275 The more flexible and intuitive output produced by the MCMC and bootstrap methods, including the
276 ability to produce a single probability that the true reduction is less than a given value, make them both
277 more attractive methods than the current WAAVP recommendation. It is evident that the MCMC method
278 outperformed the bootstrap method in this study, however this may not be true when the data has a larger
279 sample size or mean. Since the true distribution of data is unknown, the most conservative estimate would
280 be to use the data at which the MCMC method performed worst. This produced notional 95% confidence
281 intervals with a true estimated confidence of 93%. The bootstrapping procedure returned notional 95%
282 confidence intervals with a true confidence greater than or equal to 93% only when the sample size was
283 at least 40 with a pre-treatment mean FEC of 40 counted eggs or more, or with a sample size of at least
284 50 with a pre-treatment mean FEC of ten counted eggs or more. This suggests that the MCMC method
285 should be used in preference to the bootstrap method with a sample size of less than 40 with a pre-
286 treatment mean FEC of 40 counted eggs (equal to, for example, 1000EPG with an egg detection threshold
287 of 25 EPG), or with a sample size of less than 50 with smaller pre-treatment mean FEC. The authors
288 expect that similar results could be obtained using any computationally intensive parametric method such
289 as parametric bootstrapping, likelihood profiling, or MCMC sampling from the likelihood without the use
290 of prior information. With larger datasets, the data distribution independence and reduced computational
291 effort associated with the non-parametric bootstrap procedure make this method more attractive.

292 **5. Conclusions**

293 Using data simulated with similar values of mean and sample size to those observed in equine FECRT
294 data, both the method currently advocated by the WAAVP and a non-parametric bootstrap method failed
295 to provide true 95% confidence intervals for the FEC reduction. In order to avoid making erroneous inference
296 regarding the true efficacy of anthelmintics in the field, computationally intensive parametric methods such as
297 MCMC should therefore be used with sample sizes of less than 50. The large proportion of inconclusive results
298 returned from analysis of equine FECRT data suggests that the routine use of prior sample size calculations
299 should be adopted to ensure sufficient data is collected. Software to perform all three types of analyses
300 documented here is freely available in the form of an add-on package to the R statistical programming

301 language from <http://cran.r-project.org/web/packages/bayescount/index.html>.

302 **6. Acknowledgements**

303 This research was produced as part of the DEFRA-funded VTRI project 0101. BioSS is partly funded
304 by the Scottish Government. The authors are grateful to Stig Petersen (Equilab Laboratory, Joerlunde
305 Overdrev 7, DK-3500 Slangerup, Denmark) for providing the equine FECRT data discussed.

306 **References**

- 307 Coles, G. C., Bauer, C., Borgsteede, F. H., Geerts, S., Klei, T. R., Taylor, M. A., Waller, P. J., 1992. World Association for the
308 Advancement of Veterinary Parasitology (W.A.A.V.P.) methods for the detection of anthelmintic resistance in nematodes
309 of veterinary importance. *Vet. Parasitol.* 44 (1-2), 35–44.
- 310 Coles, G. C., Jackson, F., Pomroy, W. E., Prichard, R. K., von Samson-Himmelstjerna, G., Silvestre, A., Taylor, M. A.,
311 Vercruyse, J., 2006. The detection of anthelmintic resistance in nematodes of veterinary importance. *Vet. Parasitol.* 136 (3-
312 4), 167–185.
- 313 Denwood, M., 2008. runjags: Run Bayesian MCMC Models in the BUGS syntax from Within R. R package version 0.9.2.
314 URL <http://cran.r-project.org/web/packages/runjags/>
- 315 Denwood, M. J., Stear, M. J., Matthews, L., Reid, S. W. J., Toft, N., Innocent, G. T., 2008. The distribution of the pathogenic
316 nematode *Nematodirus battus* in lambs is zero-inflated. *Parasitology* 135 (10), 1225–1235.
- 317 Dopfer, D., Kerssens, C. M., Meijer, Y. G. M., Boersema, J. H., Eysker, M., 2004. Shedding consistency of strongyle-type eggs
318 in Dutch boarding horses. *Vet. Parasitol.* 124 (3-4), 249–258.
- 319 Gelman, A., Rubin, D., 1992. Inference from Iterative Simulation using Multiple Sequences. *Statistical Science* 7, 457–511.
- 320 Gilks, W. R., Richardson, S., Spiegelhalter, D. J., 1998. Markov chain Monte Carlo in practice. Chapman and Hall, Boca
321 Raton, Fla.
322 URL <http://www.loc.gov/catdir/enhancements/fy0646/98033429-d.html>
- 323 Kaplan, R. M., 2002. Anthelmintic resistance in nematodes of horses. *Vet. Res.* 33 (5), 491–507.
- 324 Klei, T. R., Chapman, M. R., 1999. Immunity in equine cyathostome infections. *Vet. Parasitol.* 85 (2-3), 123–133.
- 325 Miller, C. M., Waghorn, T. S., Leathwick, D. M., Gilmour, M. L., 2006. How repeatable is a faecal egg count reduction test?
326 *N. Z. Vet. J.* 54 (6), 323–328.
- 327 Mooney, C. Z., Duval, R. D., 1993. Bootstrapping: a nonparametric approach to statistical inference. Sage Publications, 2455
328 Teller Road, Thousand Oaks, CA 91320.
329 URL <http://www.sagepub.com/booksProdDesc.nav?prodId=Book3980&>
- 330 Morrison, D. A., 2004. Technical variability and required sample size of helminth egg isolation procedures: revisited. *Parasitol.*
331 *Res.* 94 (5), 361–366.
- 332 Nielsen, M. K., Haaning, N., Olsen, S. N., 2006. Strongyle egg shedding consistency in horses on farms using selective therapy
333 in Denmark. *Vet. Parasitol.* 135 (3-4), 333–335.
- 334 Nødtvedt, A., Dohoo, I., Sanchez, J., Conboy, G., DesCôteaux, L., Keefe, G., Leslie, K., Campbell, J., 2002. The use of negative
335 binomial modelling in a longitudinal study of gastrointestinal parasite burdens in Canadian dairy cows. *Can. J. Vet. Res.*
336 66 (4), 249–257.
- 337 Plummer, M., 2008. Just Another Gibbs Sampler (JAGS).
338 URL <http://calvin.iarc.fr/~martyn/software/jags/>
- 339 R Development Core Team, 2008. R: A Language and Environment for Statistical Computing. R Foundation for Statistical
340 Computing, Vienna, Austria, ISBN 3-900051-07-0.
341 URL <http://www.R-project.org>
- 342 Raftery, A. E., Lewis, S. M., 1995. The number of iterations, convergence diagnostics and generic Metropolis algorithms. In
343 *Practical Markov Chain Monte Carlo* (W.R. Gilks, D.J. Spiegelhalter and S. Richardson, eds.). London, U.K.: Chapman
344 and Hall.
- 345 Shaw, D. J., Grenfell, B. T., Dobson, A. P., 1998. Patterns of macroparasite aggregation in wildlife host populations. *Para-*
346 *sitology* 117 (Pt 6), 597–610.
- 347 Torgerson, P. R., Schnyder, M., Hertzberg, H., 2005. Detection of anthelmintic resistance: a comparison of mathematical
348 techniques. *Vet. Parasitol.* 128 (3-4), 291–298.

- 349 Uhlinger, C. A., 1993. Uses of fecal egg count data in equine practice. *Compendium On Continuing Education For the Practicing*
350 *Veterinarian* 15 (5), 742–749.
- 351 Vidyashankar, A. N., Kaplan, R. M., Chan, S., 2007. Statistical approach to measure the efficacy of anthelmintic treatment on
352 horse farms. *Parasitology* 134 (Pt.14), 2027–2039.
- 353 Vose, D., 2004. ModelAssist for @Risk. Risk Thinking, Ltd.
354 URL <http://www.vosesoftware.com/modelassist.htm>

	Gamma-Poisson	Multi-modal	Uniform-Poisson
Bootstrapping	0.702	0.459	0.803
WAAVP	0.673	0.532	0.786
MCMC	0.746	0.555	0.803

Table 1: Mean relative size of 95% confidence intervals for the true mean FEC reduction produced by each method from the analysis of 1000 datasets simulated using each distribution

	Gamma-Poisson	Multi-modal	Uniform-Poisson
Bootstrapping	1.69	1.49	2.6
WAAVP	1.7	1.53	2.66
MCMC	1.61	1.46	1.77

Table 2: Relative root-mean-square-error for median or mean estimate for the true mean FEC reduction produced by each method from the analysis of 1000 datasets simulated using each distribution

		Bootstrap			
		Sus.	Poss. res.	Prob. res.	Res.
MCMC	Susceptible	0	0	0	0
	Possible resistant	10	8	0	0
	Probable resistant	4	9	22	0
	Resistant	0	0	7	3

Table 3: The number of datasets assigned to each category of estimated efficacy status by MCMC and bootstrap analysis for 63 individual Danish equine datasets (median (range) of 9 (6-22) animals per dataset)

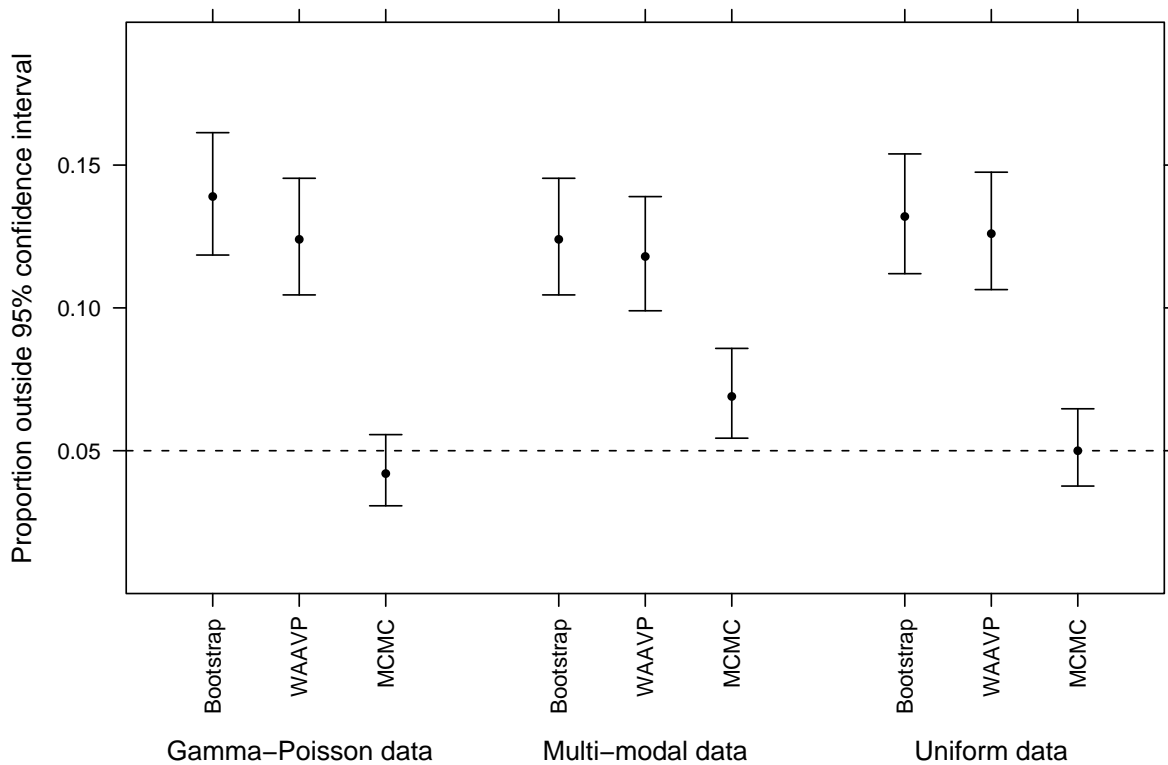


Figure 1: The proportion of 95% confidence intervals not containing the simulated true mean FEC reduction parameter for each method from the analysis of 1000 datasets simulated using each distribution

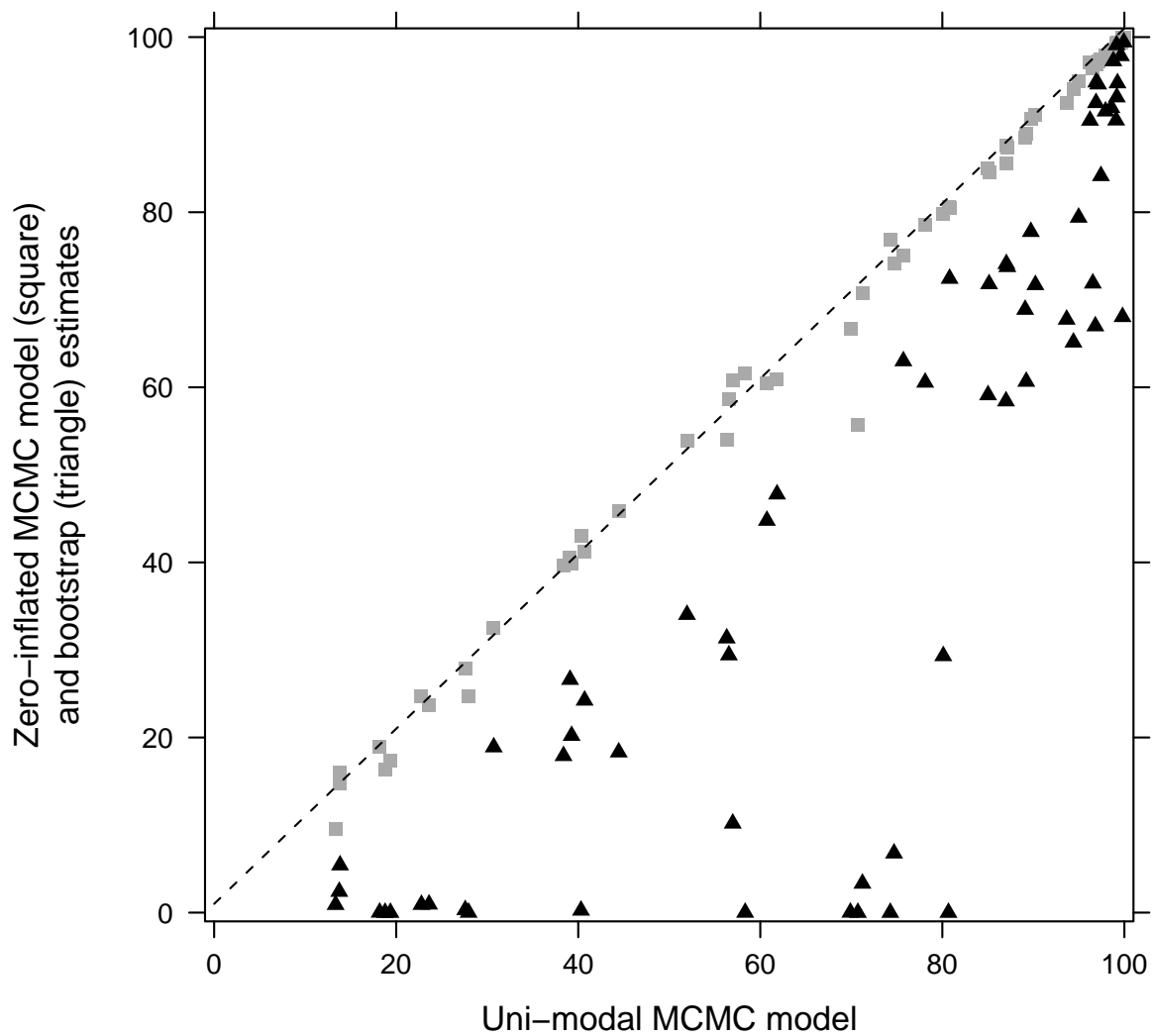


Figure 2: Comparison of the estimated probability of efficacy $< 95\%$ returned for 63 individual Danish equine datasets by bootstrapping and an MCMC method based on a zero-inflated gamma-Poisson distribution, relative to an MCMC method based on a uni-modal gamma-Poisson distribution (median (range) of 9 (6-22) animals per dataset)

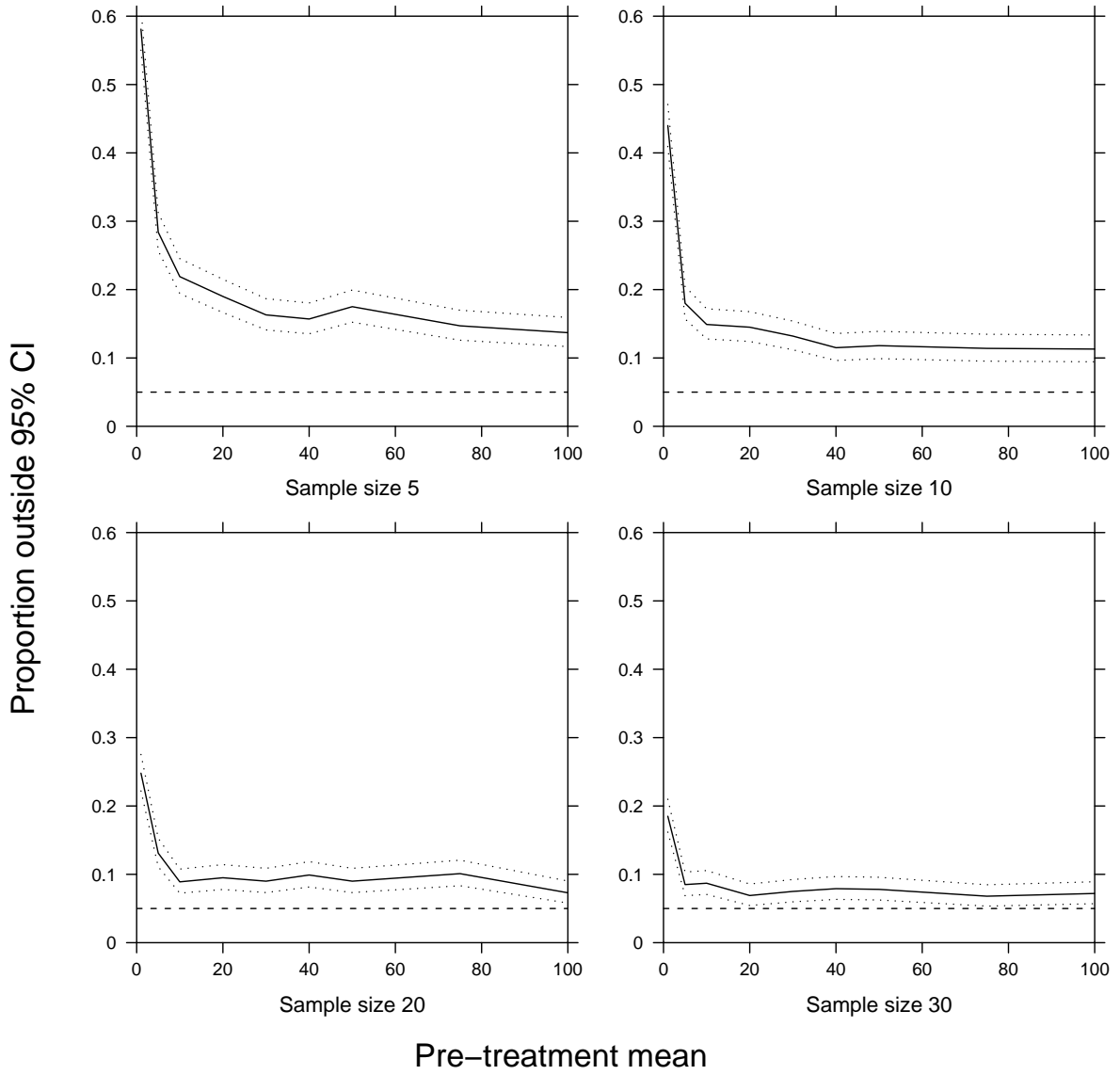


Figure 3: Proportion of 95% confidence intervals produced using the bootstrap method that did not contain the true parameter from 1000 simulated datasets at each pre-treatment mean number of eggs counted (95% credible intervals in dotted lines). Sample sizes 5, 10, 20 and 30 shown.

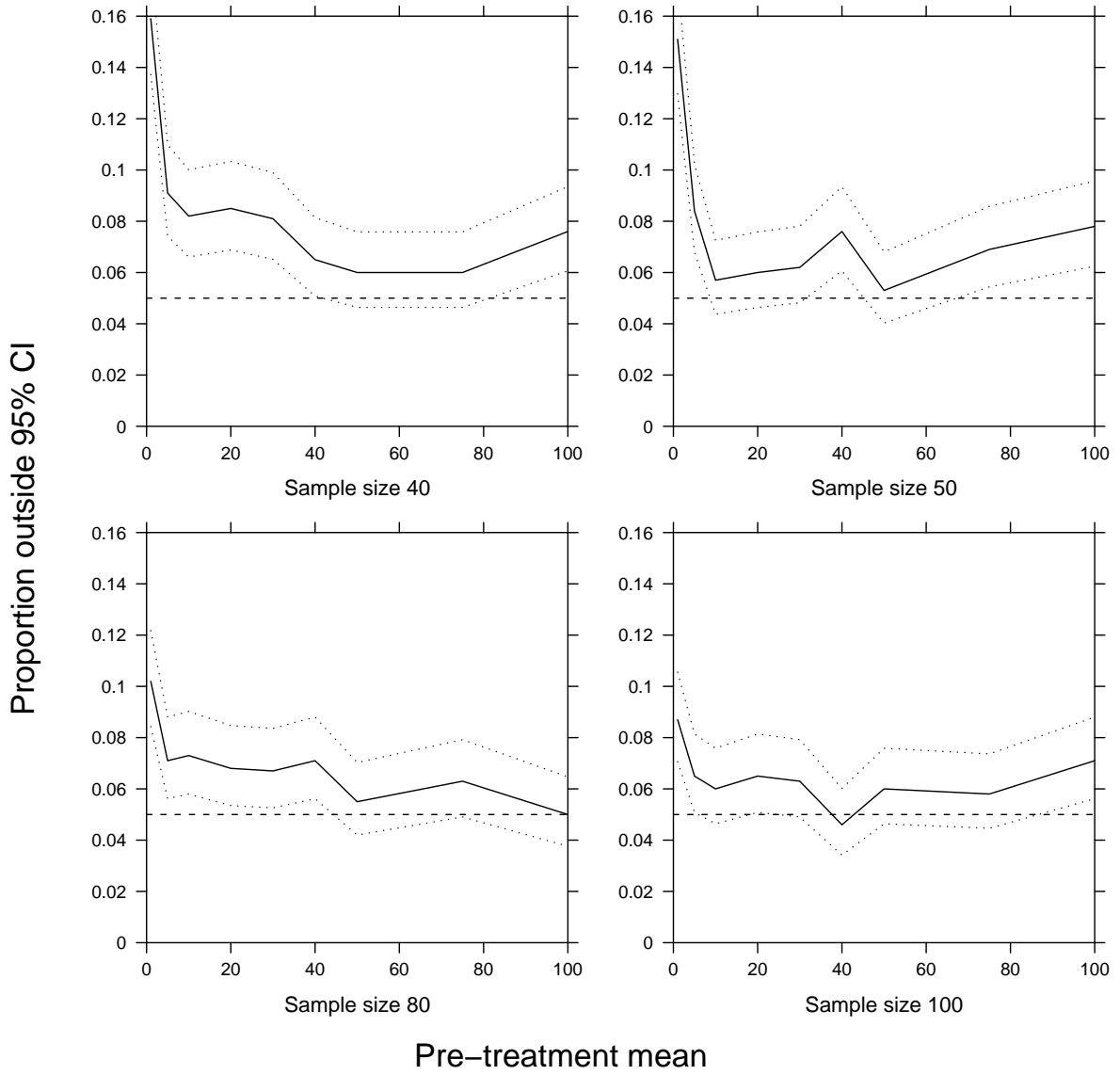


Figure 4: Proportion of 95% confidence intervals produced using the bootstrap method that did not contain the true parameter from 1000 simulated datasets at each pre-treatment mean number of eggs counted (95% credible intervals in dotted lines). Sample sizes 40, 50, 80 and 100 shown.