

1 GLM for partially pooled categorical predictors with a
2 case study in biosecurity

3 Christopher M. Baker^{1,2,3} Howard Bondell^{1,2}
4 Nathaniel Bloomfield^{1,2,3} Elena Tartaglia⁴
5 Thomas K. Waring^{1,2,3} Andrew P. Robinson^{1,3}

6 August 24, 2023

- 7 1. School of Mathematics and Statistics, The University of Melbourne
8 2. Melbourne Centre for Data Science, The University of Melbourne
9 3. Centre of Excellence for Biosecurity Risk Analysis, The University of Melbourne
10 4. Data61, Commonwealth Scientific and Industrial Research Organisation

11 cbaker1@unimelb.edu.au
12 howard.bondell@unimelb.edu.au
13 nathaniel.bloomfield@unimelb.edu.au
14 elena.tartaglia@data61.csiro.au
15 tom.waring@unimelb.edu.au
16 apro@unimelb.edu.au

17 **Acknowledgements**

Abstract

National governments use border information to efficiently manage the biosecurity risk presented by travel and commerce. In the Australian border biosecurity system, data about cargo entries are collected from records of directions: that is, the records of actions taken by the biosecurity regulator. An entry is a collection of import lines where each line is a single type of item or commodity. Analysis is simple when the data are recorded in line mode: the directions are recorded individually for each line. The challenge comes when data are recorded in container mode, because the same direction is recorded against each line in the entry, meaning that we don't know which line(s) within the entry are non-compliant. We develop a statistical model to use container mode data to help inform biosecurity risk of items. We use asymptotic analysis to estimate the value of container mode data compared to line mode data, do a simulation study to verify that we can accurately estimate parameters in a large dataset, and we apply our methods to a real dataset, for which important information about the risk of non-compliance is recovered using the new model.

Keywords: Bayesian inference; risk analysis; cargo

1 Introduction

Invasive species pose a multifaceted threat to society, leading to reductions in agricultural productivity, as well as damages to the environment, human health, and the economy [Kumar Rai and Singh, 2020]. Considerable effort is devoted to managing invasive species [Jardine and Sanchirico, 2018], in either eradicating them [Baker and Bode, 2020, Holmes et al., 2019, Wenger et al., 2017, Helmstedt et al., 2016], or suppressing their numbers to reduce damages [Binny et al., 2021, Brook et al., 2012, Sharov et al., 2002]. The costs associated with managing invasive species provides governments with an incentive to manage biosecurity risks at national borders to prevent the establishment of the new species.

Given the massive scale of global trade, biosecurity regulators need to be able to allocate their resources efficiently, but to do so they must understand the risks associated with various entries. At the border, one of the most effective ways of achieving this is by using the outcome of previously conducted inspections as intelligence that informs future operations. This allows regulators to identify high risk commodities and importers, and modify their inspection targets and policies accordingly.

However, gaining insight from border inspections is a massive logistical challenge, as data must be consistently recorded in a format that makes this analysis possible. Putting infrastructure into place to collect this data, and collecting it accurately can be expensive and challenging. Often, shortcuts may be taken that render the data less valuable in analysing patterns of risk.

In Australia systems have been put in place to capture biosecurity data since the early 90's; these data are extracted from the directions applied to the cargo in question. Directions are used to control the movement and direct the assessment and management

of goods subject to biosecurity control, and the *modes* discussed in this paper correspond to how those directions are carried out: at the line level, or at the entry level. As such, in line mode, the details of which items within an entry are inspected and found to be compliant or non-compliant are fully recorded. However, in container mode, the results of an inspection are applied to all items within an entry. Container mode was introduced so that containers could be released piecemeal in order to minimise the bottlenecks created in ports when lines that comprise many containers are held until all of the line has been cleared. Container mode makes it much quicker for border staff to manage entries with a large number of items, but means that when the data are analysed, entries in container mode are censored — in these cases, it is unknown which items were inspected, and which of those were found to be compliant or non-compliant. This makes analysis of the data to identify trends in biosecurity risk challenging.

Data that have been collected in container mode are closely related to the data collected under *pooled testing*, which is often used for disease surveillance. Within the pooled testing literature there are two main branches: one aims to identify positives within a pool, while the other seeks to use pooled data to estimate quantities about the population. It is the latter – estimating quantities – that we are interested in. The fundamental problem is estimating a prevalence, p , within a population, when only pooled data are available [Thompson, 1962]. More recent work has focused on improving estimates by reducing bias, either through altering the sampling strategy [Schaarschmidt, 2007, Hepworth and Watson, 2009] or by incorporating bias correction into models [Hepworth and Biggerstaff, 2017, 2021]. There have also been extensions of the problem where p is not a constant, but it is estimated using linear regression using only the pooled data [Delaigle and Hall, 2015, Chatterjee and Bandyopadhyay, 2020, McMahan et al., 2017, Liu et al., 2020]. These papers have made significant progress in fitting increasingly complex models, but do not focus on the impacts of different types of pooling on the precision of model estimates.

In practice, however, pooled testing is manifestly different from the biosecurity scenario. Pooled testing exists by design: as a way to gather information about a population while reducing testing. In biosecurity, inspection is applied to every individual line, and the results are only pooled at the point of data capture — as a side effect of the mode selection of the entry. Hence, we are interested in how much information we are losing due to aggregating results as container mode. Our analysis offers regulators the opportunity to assess the risks of the continued use of container mode, and to weigh them against its operational advantages.

In this paper we investigate the effect of container mode data collection upon our ability to estimate the biosecurity risk of items. We start with an asymptotic analysis, where we calculate the precision of estimates and determine the implications of mixing different item types in container mode. We then develop a simulation experiment that allows us to understand how larger entries and more item types affect the precision of our estimates. Finally, we analyse biosecurity data provided by the Australian Department of Agriculture, Fisheries and Forestry (DAFF) to identify the real-world differences

101 between using the line-only data and including the container mode data.

102 2 Model overview

103 2.1 Data

104 To make our language about the data more precise, we will explicitly define what we
105 mean by entries, lines and directions. Entries are a collection of lines, and a line is
106 a group of the same type of item or commodity being imported. When cargo enters
107 the country, each line is given *directions*. These directions detail all of the activities
108 undertaken by the biosecurity regulator to manage the biosecurity risk of each line, and
109 also the outcome of those activities. To isolate the uncertainty that arises from use of
110 container mode, we focus on one aspect of the directions recorded against import lines:
111 whether they were deemed compliant (within biosecurity regulations) or not.

112 Line and container mode are the two ways that records are kept of actions made by
113 the biosecurity regulator. In line mode, the directions assigned to each line are recorded
114 along with the outcome for that line. In container mode, directions are only recorded per
115 entry. This means that in container mode, if any line in that entry has an inspection, and
116 non-compliance is found, then every line in that entry is recorded to be non-compliant.
117 If all of the lines in the entry are compliant, then they are all marked compliant, so in
118 this case line and container mode are equivalent.

119 2.2 Modelling

120 Throughout this paper, we focus on estimating the probability that a line is non-
121 compliant using information including the type of item, country of origin and whether
122 it has complete documentation. The full model for the probability that a line is non-
123 compliant, $p_{ijk\ell}$, is:

$$\text{logit}(p_{ijk\ell}) = \alpha_i + \beta \mathbb{I}_j + \delta_k + \gamma_\ell, \quad (1)$$

124 where the fixed effects are α_i , $\beta \mathbb{I}_j$ and δ_k : α_i represents the item type, the indicator
125 variable \mathbb{I}_j denotes whether there is correct documentation, the coefficient β is the given
126 to the instance that there is correct documentation, and δ_k represents the country of
127 provenance. The random effect γ_ℓ represents the entry effect, which we include because
128 there may be correlations between lines within an entry. We anticipate some correlation
129 because lines in the same entry originate from a common context, so they may be non-
130 compliant for related reasons. The indices can take values

$$i = 1, \dots, a, \quad a \in \mathbb{N}, \quad a = \# \text{ items} \quad (2)$$

$$j = 1, 2, \quad \text{without and with documentation} \quad (3)$$

$$k = 1, \dots, d, \quad d \in \mathbb{N}, \quad d = \# \text{ countries} \quad (4)$$

$$\ell = 1, \dots, g, \quad g \in \mathbb{N}, \quad g = \# \text{ entries.} \quad (5)$$

131 The values of the indicator variable are

$$\mathbb{I}_1 = 0, \quad \text{without documentation} \quad (6)$$

$$\mathbb{I}_2 = 1, \quad \text{with documentation.} \quad (7)$$

132 The random effect γ_ℓ has distribution

$$\gamma_\ell | \sigma \sim \text{Normal}(0, \sigma), \quad \ell = 1, \dots, g, \quad g \in \mathbb{N}. \quad (8)$$

133 If all data were in line mode then the above model would be a fairly standard mixed
 134 effects logistic regression with categorical variables. However, because of the use of
 135 container mode to capture the data, we don't observe outcomes for each line, as every
 136 line in the entry is marked as non-compliant if any line in the entry is found to be
 137 non-compliant. Therefore, the outcome is whether the entry is compliant and we need
 138 to calculate the probability that the entry is non-compliant, which is one minus the
 139 probability that every line in the entry is compliant:

$$\mathbb{P}(\text{Entry } l \text{ non-compliant}) = q_l = 1 - \prod_{ijk \text{ for lines in } l} [1 - p_{ijkl}], \quad (9)$$

140 where p_{ijkl} is the probability that the line with indices $ijkl$ is non-compliant, calculated
 141 from Eq. (1). Hence, for entries in container mode, we treat the entry as a Bernoulli
 142 random variable with probability defined by Eq. (9), while for entries in line mode, we
 143 treat each line as a Bernoulli random variable with probability as defined in Eq. (1).

144 This paper includes three analyses: an asymptotic analysis, a simulation study, and a
 145 case study of Australian biosecurity data. For the asymptotic analysis we only consider
 146 the item type, ignoring effects due to the country of origin, documentation and entry
 147 effect. As such, rather than using Eq. (1), we just consider the probability that a line
 148 of item type i is non-compliant, p_i . The simulation study and the case study both use
 149 the full model, as defined above.

150 3 Asymptotic analysis

151 We use asymptotic analysis to investigate how the precision of estimates depends on
 152 entry size, the number of entries, the probability of non-compliance and whether item

types are mixed. This analysis comprises two parts. The first assumes that all items are a single type, which allows us to quantify how the amount of data, probability of non-compliance and entry size affect precision. The second part assumes that there are two different item types, and it explores how changing the proportion of entries with both item types mixed affects precision.

Throughout this section we make two simplifications. Firstly, we do not separate line mode and container mode because container mode data with an entry size of one is mathematically equivalent to line mode data. Hence, throughout these analysis, an entry size of one means line mode and entry size greater than one implies container mode. Secondly, we assume that each item has a fixed probability of non-compliance. As such, any uncertainty in the inference of this value arises as a result of the difference between container and line mode. With this in hand, we consider each line a Bernoulli trial, which only depends on the item type. When the entry size is greater than one, the relevant probability is whether at least one line was non-compliant.

We estimate precision using an asymptotic estimate of the standard error. We calculate the precision from the square roots of the diagonal elements in the Fisher information matrix, \mathcal{I} , which is the expected value of the negative of the Hessian matrix of the log-likelihood evaluated at the value of the parameter.

3.1 Single item type

For the single item type case, we set the probability of non-compliance to be p , and define N as the total number of entries, I as the number of non-compliant entries and S as the size (i.e. the number of lines) in each entry. The likelihood is a binomial distribution, where the outcome is the discovery of a non-compliant entry. The probability that an entry is compliant is

$$\mathbb{P}(\text{entry compliant}) = (1 - p)^S, \quad (10)$$

meaning that the binomial likelihood for a set entry size S proportional to

$$\mathcal{L}_S = (1 - (1 - p)^S)^I (1 - p)^S)^{N-I}. \quad (11)$$

Therefore, the log-likelihood is

$$\log \mathcal{L}_S = I \log (1 - (1 - p)^S) + (N - I) \log ((1 - p)^S). \quad (12)$$

As there is only one parameter, we calculate its second derivative (rather than needing a Hessian matrix):

$$\left[\frac{\partial^2 \log \mathcal{L}_S}{\partial p^2} \right] = \frac{S \left(N + \frac{I((1+S)(1-p)^S - 1)}{((1-p)^S - 1)^2} \right)}{(1 - p)^2}. \quad (13)$$

181 To calculate the Fisher information, we take the expectation of the number of entries
 182 with non-compliance, which depends on the size of the entry:

$$\mathbb{E}[I] = N(1 - (1 - p)^S). \quad (14)$$

183 Hence the Fisher information is

$$\mathcal{I} = -\mathbb{E}\left[\frac{\partial^2 \log \mathcal{L}_S}{\partial p^2}\right] = -\frac{NS^2(1 - p)^{S-2}}{(1 - p)^S - 1}, \quad (15)$$

184 and the standard error estimate is

$$SE = \left(-\frac{NS^2(1 - p)^{S-2}}{(1 - p)^S - 1}\right)^{-1/2}. \quad (16)$$

185 Using Eq. (16) we can understand how the probability of non-compliance, the entry
 186 size and the number of lines affect the standard error, and we plot these relationships
 187 in Figure 1. The left plot shows that the standard error depends on the probability
 188 of non-compliance and that the relationship depends on the entry size. For all entry
 189 sizes, the standard error is small when the probability of non-compliance is small (below
 190 ~ 0.3). However, for larger values of the probability of non-compliance, the standard
 191 error increase significantly if the entry size is three or greater. The large p behaviour is
 192 driven by the $(1 - p)^{S-2}$ term in Eq. (16), which means SE goes to 0 if $S = 1$, while it
 193 diverges if $S \geq 3$. Figure 1 also shows how the standard error decreases as the number
 194 of lines of data increases. The lower the entry size is, the lower the standard error, and,
 195 as $SE \sim \sqrt{1/N_{E,S}}$, container mode data with larger entry sizes require a large amount
 196 of data to reach the same standard error.

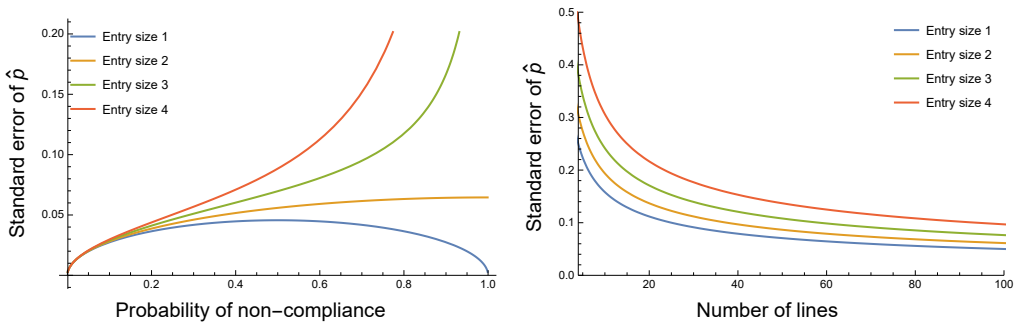


Figure 1: The standard error (Eq. (16)) as the probability of non-compliance, p , is varied (left) and as the number of lines are varied (right). For the left plot the number of lines is held constant at 120. For the right plot, the probability of non-compliance is held constant at 0.5

3.2 Two item types

In this section we consider a situation where there are two items with probabilities of non-compliance of p_1 and p_2 , and we examine how these different items probabilities interact. We focus on a scenario where every entry is of size two, meaning there are three types of entries: only type 1; only type 2; or mixed, with one line of type 1 and one of type 2. We denote the number of lines within a single entry of type 1 and 2 as S_1 and S_2 respectively, and $I(S_1, S_2)$ and $N(S_1, S_2)$ are the number of entries with non-compliance and total number of entries with S_1 type 1 lines and S_2 type 2 lines. For our case, we can have $S_1 = 2, S_2 = 0$; $S_1 = 1, S_2 = 1$; or $S_1 = 0, S_2 = 2$. Rewriting the log-likelihood from Eq. (12), we get

$$\log \mathcal{L} = \sum_{S_1, S_2} I(S_1, S_2) \log \left(1 - (1 - p_1)^{S_1} (1 - p_2)^{S_2} \right) + (N(S_1, S_2) - I(S_1, S_2)) \log \left((1 - p_1)^{S_1} (1 - p_2)^{S_2} \right). \quad (17)$$

We compute the Fisher information matrix and the standard error using Mathematica. The standard error for p_1 is

$$\frac{1}{2} \sqrt{-\frac{p_1 (p_1^2 - 3p_1 + 2) (N_{1,2}(p_1 - 1)(p_2 - 2)p_2 + 4N_{2,2}(p_2 - 1)(p_1(p_2 - 1) - p_2))}{N_{1,1}(p_1 - 1)(N_{1,2}(p_1 - 1)(p_2 - 2)p_2 + 4N_{2,2}(p_2 - 1)(p_1(p_2 - 1) - p_2)) + N_{1,2}N_{2,2}(p_1 - 2)p_1(p_2 - 1)^2}}, \quad (18)$$

where $N_{1,1}$ and $N_{2,2}$ are the number of entries with two type 1 lines and two type 2 lines respectively, while $N_{1,2}$ are the number of entries with both type 1 and type 2 lines. The standard error for p_2 is the same, with $N_{1,1}$ and $N_{2,2}$ switched and p_1 and p_2 switched.

By examining Eq. (18) we can see that the behaviour of the standard errors is more complex than when we considered only one type of item. Notably, the number of entries of only type 2, $N_{2,2}$, is in the equation for the type 1 standard error, along with the probability of non-compliance of type 2, p_2 .

Figure 2 shows how the standard error varies with the proportion of mixed entries, for different non-compliance probabilities. Here we fix $N_{\text{total}} = 50$, and keep $N_{1,1} = N_{2,2}$ while the proportion $N_{1,2}/N_{\text{total}}$ is varied. In each case, we hold p_1 at 0.1, and we vary p_2 for 0.05 up to 0.9 across the four plots. The way that the standard error changes as a function of the proportion of mixed entries changes markedly, depending on the value of p_2 . In particular, when p_2 is 0.7 and 0.9, the standard error for p_1 increases as the proportion of mixed entries increases, while the standard error for p_2 actually decreases, while the proportion of mixed entries is below $\sim 90\%$.

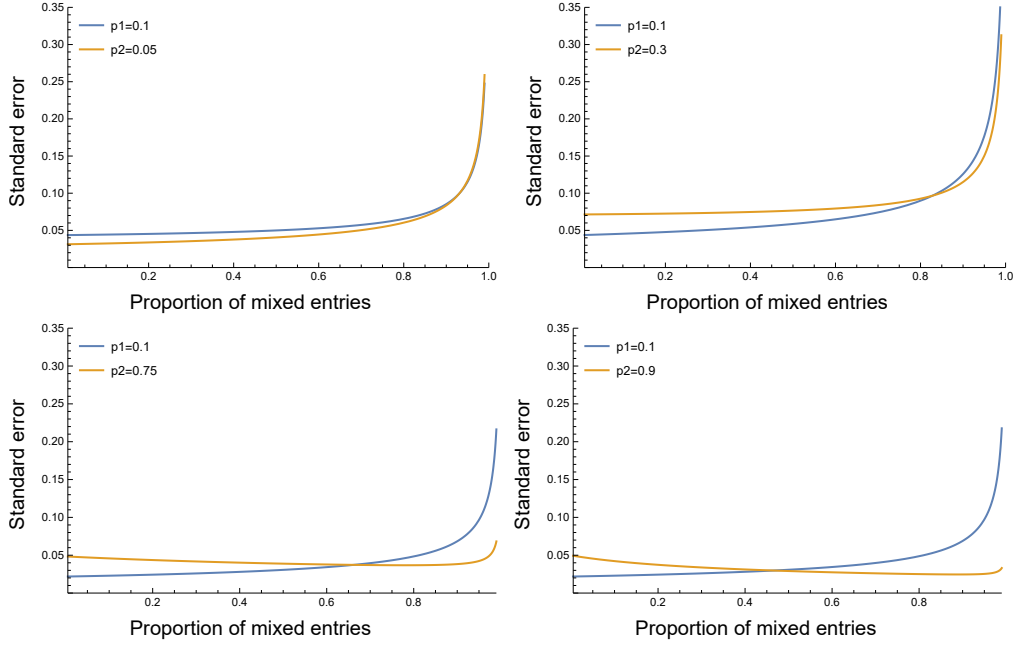


Figure 2: The standard error when for entry size 2, when the proportion of mixed entries is varied. In each plot there are 50 total entries and p_1 is held constant at 0.1, while p_2 goes from 0.05 up to 0.9.

3.3 Asymptotic analysis conclusions

From the asymptotic analysis we make two observations:

1. Container mode data are most useful when the probability of non-compliance, p , is small. In the extreme case where $p \rightarrow 0$, container mode data are equivalent to line mode data because whenever an entry is found to be compliant, we know that every line within that entry is compliant, regardless of whether it is container mode or line mode.
2. By allowing entries to contain lines of different item types, container mode can increase or decrease the precision of estimates, depending on the true probability of non-compliance.

4 Simulation study

As the first step in developing a method to analyse real-world data sets, we simulate a data set that contains key complexities including multiple item types, fixed effects and random effects. Once we have simulated data, we can fit a model to estimate parameters. The advantage of our approach is that we can (1) verify that our model behaves correctly and (2) explore how model precision varies in a more complex setting.

240 4.1 Data simulation

241 We simulate data from the full model (as described in Section 2.2). For our simulations,
 242 we chose parameters that we expect to be similar to real-world values. We set $\beta = -1$,
 243 choosing an negative effect of having correct documentation, because we expect lines
 244 with correct documentation have a higher chance of being compliant. We include $a = 5$
 245 item types with α_i taking values of -6.91, -4.60, -3.89, -2.94, -1.39 (corresponding to
 246 non-compliance probabilities of 0.001, 0.01, 0.02, 0.05 and 0.2, if all other effects were
 247 0). We chose negative values for the item effect α_i , since the probability of detection
 248 is expected to be low for any item in real-world data. We used $d = 3$ countries and
 249 set their weights to be 0.5, -1 and 0.25. For within-entry correlation, we set $\sigma = 0.25$.
 250 For each entry we draw whether it is in line mode or not with probability 0.25, and for
 251 each line we set the probability of having correct documentation to be 0.2. We show an
 252 example of the format of the data in Table 1.

Table 1: Example simulated data. Lines 1-3 are all marked as non-compliant because entry 1 is in container mode, even though only 1 or 2 of the lines were actually non-compliant.

Line	Entry	Type	Documentation	Mode	Non-compliant
1	1	3	1	Container	1
2	1	2	0	Container	1
3	1	1	0	Container	1
4	2	5	0	Line	0
5	3	4	1	Line	1
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

253 4.2 STAN model

254 We model the system in a Bayesian framework, using the RSTAN package in R. We
 255 choose a Bayesian package due to the ease of specifying the model and could be updated,
 256 in principle, for any specification of the model for the probability of non-compliance.
 257 The basic structure for fitting the simulated data is to define the probability of non-
 258 compliance for each line (following Eq. (1))

$$\text{logit}(p_{ijk\ell}) = \alpha_i + \beta \mathbb{I}_j + \delta_k + \gamma_\ell, \quad (19)$$

259 where

$$\gamma_\ell | \sigma \sim \text{Normal}(0, \sigma), \quad \ell = 1, \dots, g, \quad g \in \mathbb{N}. \quad (20)$$

260 As we fit our model in a Bayesian framework, we set priors on each parameter:

$$\alpha_i \sim \text{Normal}(-4, 4), \quad i = 1, \dots, a, \quad a \in \mathbb{N}, \quad (21)$$

$$\beta \sim \text{Normal}(0, 0.5), \quad (22)$$

$$\delta_k \sim \text{Normal}(0, 0.5), \quad k = 1, \dots, d, \quad d \in \mathbb{N} \quad (23)$$

$$\sigma \sim \text{Normal}(0, 0.5) \quad (24)$$

261 The probability p_{ijkl} is the probability of non-compliance of a line that has charac-
 262 teristics (i, j, k, ℓ) . We denote the probability that line n is non-compliant as $p_{(n)}$, where
 263 $p_{(n)} = p_{ijkl}$ if that line has characteristics (i, j, k, ℓ) .

264 We fit the STAN model to the simulated data for two reasons. Firstly we confirm
 265 that, with sufficient data, the STAN model can accurately estimate model parameters.
 266 Secondly we explore how changing the amount of data and the proportion of container
 267 mode data affects precision, over a range of entry sizes. To measure precision, we use
 268 the standard deviation of the posterior samples, because it is the Bayesian equivalent of
 269 the standard error calculations in Section 3.

270 4.3 Simulation results

271 Figure 3 shows the that the STAN model gives good parameter estimates, if there is
 272 sufficient data. Overall, the STAN estimates are close to the true values when there is a
 273 large amount of data. The two smallest values of α show the worst performance, which
 274 is not surprising, given our asymptotic analysis results (Section 3).

275 The precision of model estimates depends on the entry size, the probability of non-
 276 compliance and the amount of container mode data (Figure 4). The best-case scenario
 277 is when all data are in line mode (ratio of container mode = 0), and the distance from
 278 that line shows the impact of using a mix, or only, container mode data. We find that
 279 the difference between the precision estimates depends strongly on the combination of
 280 factors. There are some combinations where all analyses return similar precision (e.g.
 281 $\alpha = -1.4$ and entry size 10), while other combinations have a large gap where the line
 282 only data far outperforms the others (e.g. $\alpha = -4.6$ and entry size = 5).

283 NEW FIGURE DISCUSSION

- 284 • Figure 4 to be removed: some counter-intuitive results probably due to stochastic
 285 variation (in the container generation) rather than actual results. See eg entry size
 286 20 and $\alpha = -3.9$, estimates get for more lines and 100% container mode.
- 287 • New version of Figure 3. Previously, the model was fit to some number of lines, all
 288 using one batch of data (list of containers), the spread was the standard error in
 289 the MCMC samples. Complaint with this was it didn't actually show the model
 290 produces reliable estimates, as it was only tested once.
- 291 • In the new version (see below), the model was fit to 20 separate batches of data,
 292 in the same way (ie for increasing number of lines). Now the spread is the 90% CI

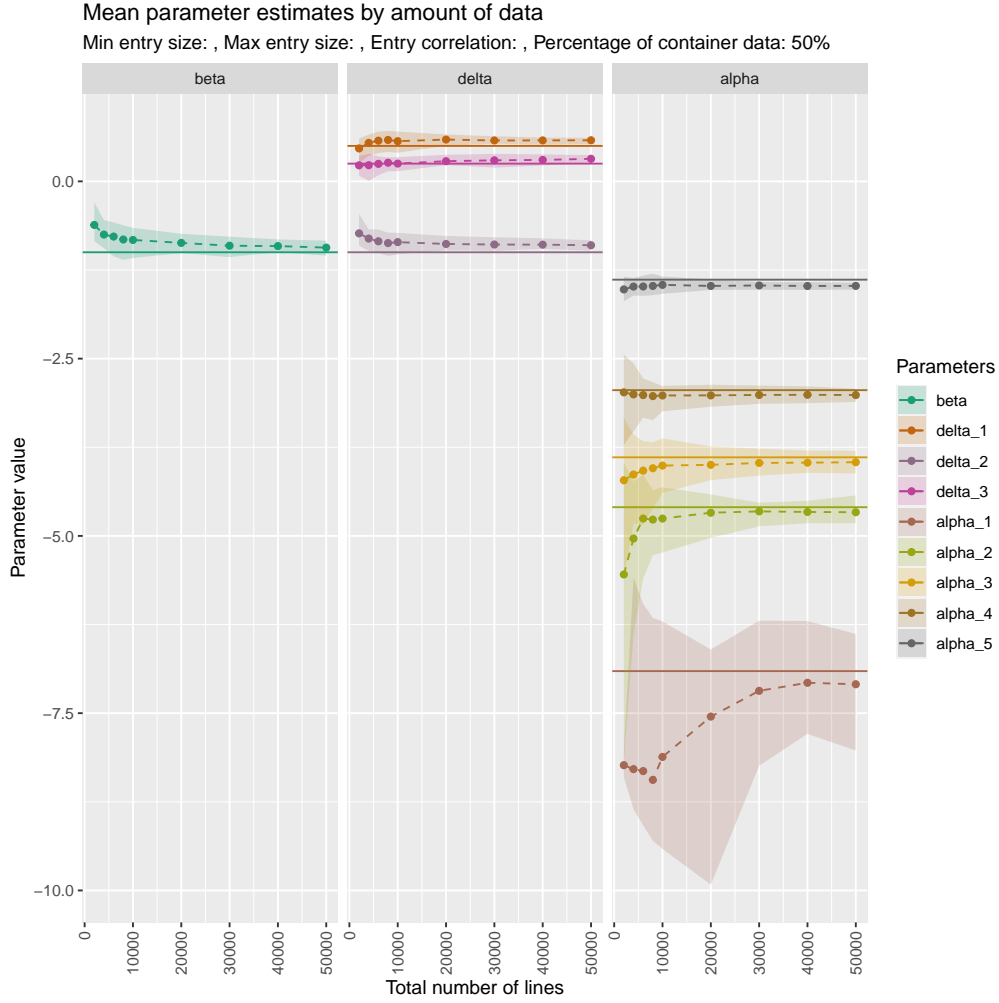


Figure 3: Results from the STAN model, fitted to simulated data. 50% of the data are in container mode and 50% are in line mode, and each entry is a random size between 1 and 20. Dots and dashed lines in each plot show the STAN estimates, while the solid lines show the true value of each parameter.

293 in the estimates for those parameters (ie we have 20 means from the 20 runs, this
294 shows how much the estimates vary, for different batches of data).

- 295 • Observe: the estimates get closer together, as well as closer to the true value,
296 as the number of lines increases. This gives better evidence for the highlighted
297 sentence at the beginning of this section. As I read it that's kind of the point of
298 this part of the paper.

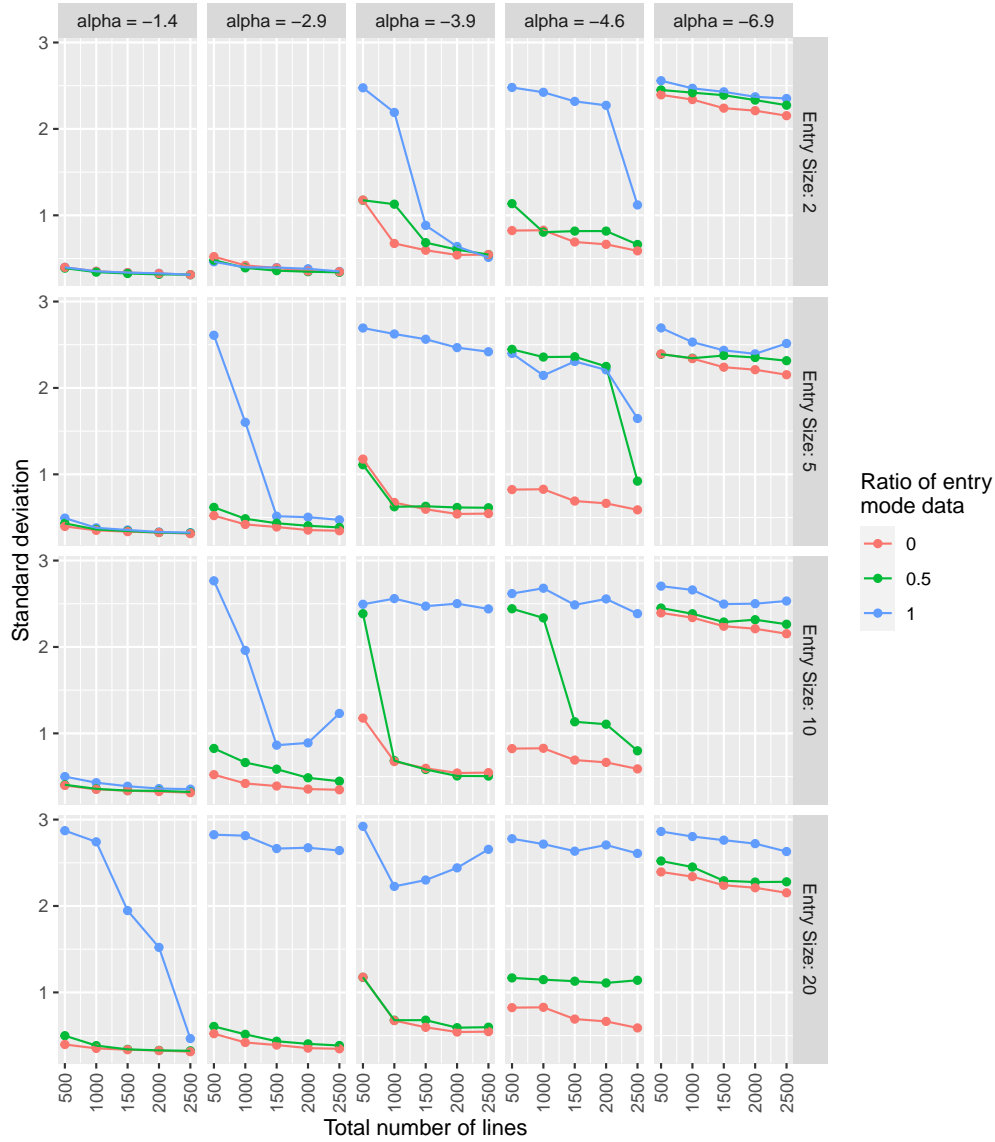


Figure 4: The standard deviation of the alpha parameter estimates for simulations with varying entry size and total number of lines. The ratio of container mode data is the proportion of data in container mode, meaning that 0 corresponds to line mode data-sets, 0.5 if a mixed data-set and 1 is purely container mode.

299 5 Case study

300 We apply our model to some real biosecurity data, to see how our understanding of
 301 the system changes by fitting the full data, compared to only using the container mode
 302 data. The data set is of furniture imports in 2020, and we break the data set into 52
 303 weekly data sets for this analysis. Because container mode data with no non-compliance
 304 is equivalent to line mode data, the more container-mode entries with non-compliance,

the larger the differences could be between analysing all of the data compared to only the line mode data. There is a natural break in the dataset, where three weeks have four container mode entries with non-compliance, with the remaining weeks having fewer. Hence, we choose to analyse these three weeks as a demonstration of potential real-world differences between analysing all data and only line mode data. For confidentiality, for each week we re-name the item type, the country and the week to be integers (week 1, 2 and 3 do not correspond to the first 3 weeks of the year, and country 1 and item 1 are not the same in weeks 1 and 2). Relabeling the data does not cause issues because we are not trying to compare estimates between weeks or draw inferences between countries or item types in this analysis. A summary of the data is given in Table 2.

Table 2: The number of lines in container mode and line mode for each week. The ratio is the fraction of the week’s data that is in container mode.

Week	Entry	Line	Ratio
1	516	6741	0.07
2	1667	7213	0.19
3	679	6832	0.09

We fit our STAN model to the data from each week, using all the data and the line-only subset of the data. Because the line-only analysis is a subset of the data, not every country or item is present in the line-only data set. For the parameters that do match, we generate a scatter plot of the mean estimates to see how frequently we get different estimates (Figure 5). For many of the parameters, we get very similar results irrespective of which data we use. This is not surprising, given that most of the data are in line mode and that line mode data gives more information than container mode data. However, it is interesting that including the container mode data results in quite large changes to some of the parameter estimates.

Analysing the full data set, rather than the line-only data, gives us information about more items and countries. Table 3 shows all of the parameter estimates that we only get when including the container mode data in the analysis.

6 Discussion

In this paper we sought to understand how container mode data affects our ability to estimate risks in the biosecurity setting. Due to the relative ease of analysing line data, in our experience, container mode data are often excluded from data analysis. From our analysis of real biosecurity data, we find that including the container mode data in analysis can markedly change some results and that by only using line data the analysis is somewhat limited – there are parameter estimates that simply cannot be made.

Through our simulation study and our asymptotic analysis we gain an understanding of how mixing container mode data into line mode data impacts our ability to make precise inferences. We see that mixing lines that have different probabilities of being non-

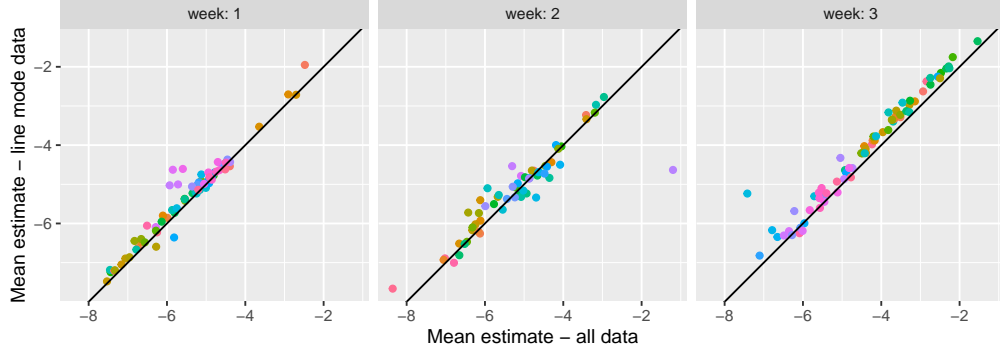


Figure 5: Scatter plot of parameter estimates using the line only data and all data. Dots that fall on the solid black line indicate when we get the same parameter estimate from each data set.

compliant has divergent affects on model precision. In the asymptotic analysis we see that mixing a high-probability item with a lower probability item makes it easier to precisely estimate the high-probability item's parameter, while mixing two low-probability items makes it challenging to make inferences about either. When one item has a low probability of non-compliance, a non-compliant mixed entry is likely to correspond to a non-compliant entry of the of the item with a higher probability of compliance. The extra data from this entry can therefore be used to refine the parameter estimates, whereas if the respective probabilities are not well resolved this is not possible. These interactions between items are exacerbated in the simulation study, where we have larger entries and more item types, and we see differing patterns of item types, entry sizes and amounts of data where the presence of container mode data degrades model precision, compared to a line-only analysis.

While we are able to model the full dataset, the continuing presence of container mode data will be a barrier for future data analysis. Most statistical and machine learning algorithms are designed such that each row of data will have explanatory variables and an outcome. With container mode data, we have the explanatory variables, but only partial information about the outcome. Hence, we need customised algorithms (such as described in this paper) to analyse it appropriately. Even so, many simplifying assumptions have been made: for example, that a given item's probability of non-compliance does not vary over time. These simplifications do not change our conclusions, as the analysis still serves to highlight the differences between container and line mode. How-

Table 3: Parameter estimates that are only possible when using all data, as opposed to the line-only fits. The value `summary_rhat_all` measures convergence of MCMC chains [Brooks et al., 2011].

	param_name	week	summary_mean_all	summary_sd_all	summary_rhat_all
1	p_intercept[72]	1.00	-1.70	1.68	1.00
2	p_intercept[59]	2.00	-1.34	1.71	1.01
3	p_intercept[60]	2.00	-2.07	1.63	1.04
4	p_intercept[61]	2.00	-1.18	1.92	1.00
5	country_effect[62]	2.00	-0.01	0.49	1.00
6	country_effect[63]	2.00	-0.01	0.47	1.00
7	country_effect[64]	2.00	-0.01	0.49	1.01
8	country_effect[65]	2.00	0.02	0.48	1.02
9	country_effect[66]	2.00	0.01	0.48	1.00
10	country_effect[67]	2.00	-0.01	0.48	1.00
11	p_intercept[72]	3.00	-1.43	1.78	1.00

ever, without proper care, the specific results of our case-study analysis should not be applied to decision making.

While in our case the pooling of outcomes in container mode is an artifact of the system, pooled testing are often designed strategically in public health to improve efficiency. The strain on the PCR testing system has prompted regimes to identify cases while minimising the rounds of testing [Mutesa et al., 2021]. While it is clear that the way pooling is conducted affects how quickly infections can be identified, in some circumstances, the outcomes of pooled testing can also be used to estimate parameters, and our work has implications for these situations [Delaigle and Hall, 2015, Chatterjee and Bandyopadhyay, 2020, McMahan et al., 2017, Liu et al., 2020]. From our asymptotic analysis, it is clear that even when only estimating a prevalence, for a fixed number of tests, the standard error of the estimate depends on pool size and therefore there would be an optimal pool size, which would depend on the prevalence [Keeling and Rohani, 2008]. When there is a model with parameters being estimated, the added complexity could compound the issue, making pool sizing a more important aspect of study design. Furthermore, as we demonstrate here, it is not only pool size that matters, it is how samples with different characteristics are grouped together. Hence, for a given study, there may not be one optimal pooling strategy, and there would likely be trade-offs between different parameters when it comes to maximising the precision of estimates.

In this paper we have investigated how pooled data can reduce our ability to make statistical inferences about the population. Within the Australian biosecurity context, the pooling is due to how data are recorded and was implemented for operational reasons. Due to this operational decision, analysing full datasets is harder than the line-only mode data, meaning that either analysis is restricted to a subset of possible methods or that container mode data is ignored. While it does not follow that container mode should

383 be removed due to operational efficiency, if the line-level data could be captured then it
384 would improve our understanding of biosecurity risk. However, pooled data will continue
385 to be collected in various fields, for example due to the cost-savings of testing multiple
386 samples at once for diseases. While there is extensive work in pooled testing protocols
387 for case identification, there is less work on identifying pool sizes when aiming to make
388 inferences about aspects of the population. Our work shows that pooling tests with
389 variable underlying prevalence affects precision differentially. Hence, careful considering
390 should be given to designing pools when the data will be used to make inferences about
391 the population.

392 7 Code availability

393 All code is available at <https://github.com/cmbaker00/container-line-analysis-public>.

394 8 Data availability

395 The data that support the findings of this simulation study are openly available at
396 <https://github.com/cmbaker00/container-line-analysis-public>. The case study
397 data are available from the Australian Department of Agriculture, Fisheries and Forestry.
398 Restrictions apply to the availability of these data, which were used under license for
399 this study.

400 9 Conflict of interest

401 The authors declare no conflict of interest.

402 References

- 403 C. M. Baker and M. Bode. Recent advances of quantitative modeling to sup-
404 port invasive species eradication on islands. *Conservation Science and Prac-*
405 *tice*, n/a(n/a):e246, 2020. ISSN 2578-4854. doi: 10.1111/csp2.246. URL
406 <https://conbio.onlinelibrary.wiley.com/doi/abs/10.1111/csp2.246>. _eprint:
407 <https://conbio.onlinelibrary.wiley.com/doi/pdf/10.1111/csp2.246>.
- 408 R. N. Binny, J. Innes, N. Fitzgerald, R. Pech, A. James, R. Price, C. Gillies,
409 and A. E. Byrom. Long-term biodiversity trajectories for pest-managed ecolog-
410 ical restorations: eradication vs. suppression. *Ecological Monographs*, n/a(n/a),
411 2021. ISSN 1557-7015. doi: <https://doi.org/10.1002/ecm.1439>. URL <https://esajournals.onlinelibrary.wiley.com/doi/abs/10.1002/ecm.1439>. _eprint:
412 <https://esajournals.onlinelibrary.wiley.com/doi/pdf/10.1002/ecm.1439>.

- 414 L. A. Brook, C. N. Johnson, and E. G. Ritchie. Effects of predator control on be-
 415 haviour of an apex predator and indirect consequences for mesopredator suppres-
 416 sion. *Journal of Applied Ecology*, 49(6):1278–1286, Dec. 2012. ISSN 1365-2664.
 417 doi: 10.1111/j.1365-2664.2012.02207.x. URL <http://onlinelibrary.wiley.com.ezp.lib.unimelb.edu.au/doi/10.1111/j.1365-2664.2012.02207.x/abstract>.
- 419 S. Brooks, A. Gelman, G. Jones, and X.-L. Meng. *Handbook of Markov Chain Monte*
 420 *Carlo*. CRC press, 2011.
- 421 A. Chatterjee and T. Bandyopadhyay. Regression models for group testing: Identi-
 422 fiability and asymptotics. *Journal of Statistical Planning and Inference*, 204:141–
 423 152, Jan. 2020. ISSN 0378-3758. doi: 10.1016/j.jspi.2019.05.003. URL <https://www.sciencedirect.com/science/article/pii/S0378375818301423>.
- 425 A. Delaigle and P. Hall. Nonparametric methods for group testing data, taking dilution
 426 into account. *Biometrika*, 102(4):871–887, Dec. 2015. ISSN 0006-3444. doi: 10.1093/
 427 biomet/asv049. URL <https://doi.org/10.1093/biomet/asv049>.
- 428 K. J. Helmstedt, J. D. Shaw, M. Bode, A. Terauds, K. Springer, S. A. Robinson, and
 429 H. P. Possingham. Prioritizing eradication actions on islands: it’s not all or nothing.
 430 *Journal of Applied Ecology*, 53:733–741, Jan. 2016. ISSN 1365-2664. doi: 10.1111/
 431 1365-2664.12599. URL <http://onlinelibrary.wiley.com.ezp.lib.unimelb.edu.au/doi/10.1111/1365-2664.12599/abstract>.
- 433 G. Hepworth and B. J. Biggerstaff. Bias Correction in Estimating Proportions by Pooled
 434 Testing. *Journal of Agricultural, Biological and Environmental Statistics*, 22(4):602–
 435 614, Dec. 2017. ISSN 1537-2693. doi: 10.1007/s13253-017-0297-2. URL <https://doi.org/10.1007/s13253-017-0297-2>.
- 437 G. Hepworth and B. J. Biggerstaff. Bias Correction in Estimating Proportions by Imper-
 438 fect Pooled Testing. *Journal of Agricultural, Biological and Environmental Statistics*,
 439 26(1):90–104, Mar. 2021. ISSN 1537-2693. doi: 10.1007/s13253-020-00411-5. URL
 440 <https://doi.org/10.1007/s13253-020-00411-5>.
- 441 G. Hepworth and R. Watson. Debaised estimation of proportions in group test-
 442 ing. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*,
 443 58(1):105–121, 2009. ISSN 1467-9876. doi: 10.1111/j.1467-9876.2008.00639.
 444 x. URL <https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-9876.2008.00639.x>.
 445 eprint: <https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-9876.2008.00639.x>.
- 447 N. D. Holmes, D. R. Spatz, S. Oppel, B. Tershy, D. A. Croll, B. Keitt, P. Genovesi,
 448 I. J. Burfield, D. J. Will, A. L. Bond, A. Wegmann, A. Aguirre-Muñoz, A. F. Raine,
 449 C. R. Knapp, C.-H. Hung, D. Wingate, E. Hagen, F. Méndez-Sánchez, G. Rocamora,
 450 H.-W. Yuan, J. Fric, J. Millett, J. Russell, J. Liske-Clark, E. Vidal, H. Jourdan,
 451 K. Campbell, K. Springer, K. Swinnerton, L. Gibbons-Decherong, O. Langrand,

452 M. d. L. Brooke, M. McMinn, N. Bunbury, N. Oliveira, P. Sposimo, P. Geraldles,
453 P. McClelland, P. Hodum, P. G. Ryan, R. Borroto-Páez, R. Pierce, R. Griffiths,
454 R. N. Fisher, R. Wanless, S. A. Pasachnik, S. Cranwell, T. Micol, and S. H. M.
455 Butchart. Globally important islands where eradicating invasive mammals will ben-
456 efit highly threatened vertebrates. *PLOS ONE*, 14(3):e0212128, Mar. 2019. ISSN
457 1932-6203. doi: 10.1371/journal.pone.0212128. URL [https://journals.plos.org/
458 plosone/article?id=10.1371/journal.pone.0212128](https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0212128).

459 S. L. Jardine and J. N. Sanchirico. Estimating the cost of invasive species control.
460 *Journal of Environmental Economics and Management*, 87:242–257, Jan. 2018. ISSN
461 0095-0696. doi: 10.1016/j.jeem.2017.07.004. URL [http://www.sciencedirect.com/
462 science/article/pii/S0095069616301322](http://www.sciencedirect.com/science/article/pii/S0095069616301322).

463 M. J. Keeling and P. Rohani. *Modeling Infectious Diseases in Humans and Animals*.
464 Princeton University Press, 2008. ISBN 9780691116174. URL [http://www.jstor.
465 org/stable/j.ctvc4gk0](http://www.jstor.org/stable/j.ctvc4gk0).

466 P. Kumar Rai and J. S. Singh. Invasive alien plant species: Their impact on envi-
467 ronment, ecosystem services and human health. *Ecological Indicators*, 111:106020,
468 Apr. 2020. ISSN 1470-160X. doi: 10.1016/j.ecolind.2019.106020. URL [https:
469 //www.sciencedirect.com/science/article/pii/S1470160X19310167](https://www.sciencedirect.com/science/article/pii/S1470160X19310167).

470 Y. Liu, C. S. McMahan, J. M. Tebbs, C. M. Gallagher, and C. R. Bilder. Generalized
471 additive regression for group testing data. *Biostatistics*, (kxaa003), Feb. 2020. ISSN
472 1465-4644. doi: 10.1093/biostatistics/kxaa003. URL [https://doi.org/10.1093/
473 biostatistics/kxaa003](https://doi.org/10.1093/biostatistics/kxaa003).

474 C. S. McMahan, J. M. Tebbs, T. E. Hanson, and C. R. Bilder. Bayesian regression
475 for group testing data. *Biometrics*, 73(4):1443–1452, 2017. ISSN 1541-0420. doi:
476 10.1111/biom.12704. URL [https://onlinelibrary.wiley.com/doi/abs/10.1111/
477 biom.12704](https://onlinelibrary.wiley.com/doi/abs/10.1111/biom.12704). eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/biom.12704>.

478 L. Mutesa, P. Ndishimye, Y. Butera, J. Souopgui, A. Uwineza, R. Rutayisire, E. L.
479 Ndoricimpaye, E. Musoni, N. Rujeni, T. Nyatanyi, E. Ntagwabira, M. Semakula,
480 C. Musanabaganwa, D. Nyamwasa, M. Ndashimye, E. Ujeneza, I. E. Mwikarago,
481 C. M. Muvunyi, J. B. Mazarati, S. Nsanzimana, N. Turok, and W. Ndifon. A pooled
482 testing strategy for identifying SARS-CoV-2 at low prevalence. *Nature*, 589(7841):276–
483 280, Jan. 2021. ISSN 1476-4687. doi: 10.1038/s41586-020-2885-5. URL [https://
484 www.nature.com/articles/s41586-020-2885-5](https://www.nature.com/articles/s41586-020-2885-5). Number: 7841 Publisher: Nature
485 Publishing Group.

486 F. Schaarschmidt. Experimental design for one-sided confidence intervals or hypothesis
487 tests in binomial group testing. *Communications in Biometry and Crop Science*, 2007.
488 ISSN 1896-0782. URL [http://agrobiol.sggw.waw.pl/~cbcs/pobierz.php?plik=
489 CBCS_2_1_5.pdf](http://agrobiol.sggw.waw.pl/~cbcs/pobierz.php?plik=CBCS_2_1_5.pdf). Publisher: Faculty of Agriculture and Biology, Warsaw Agricultural
490 University, Poland.

- 491 A. Sharov, D. Leonard, A. Liebhold, E. Roberts, and W. Dickerson. “Slow The Spread”:
492 A National Program to Contain the Gypsy Moth. *Journal of Forestry*, 100(5):30–36,
493 July 2002.
- 494 K. H. Thompson. Estimation of the Proportion of Vectors in a Natural Population of
495 Insects. *Biometrics*, 18(4):568–578, 1962. ISSN 0006-341X. doi: 10.2307/2527902.
496 URL <https://www.jstor.org/stable/2527902>. Publisher: [Wiley, International
497 Biometric Society].
- 498 A. S. Wenger, V. M. Adams, G. D. Iacona, C. Lohr, R. L. Pressey, K. Morris,
499 and I. D. Craigie. Estimating realistic costs for strategic management planning
500 of invasive species eradications on islands. *Biological Invasions*, pages 1–19, Nov.
501 2017. ISSN 1387-3547, 1573-1464. doi: 10.1007/s10530-017-1627-6. URL [https:](https://link.springer.com/article/10.1007/s10530-017-1627-6)
502 [//link.springer.com/article/10.1007/s10530-017-1627-6](https://link.springer.com/article/10.1007/s10530-017-1627-6).