RESEARCH ARTICLE



Check for updates

Correcting prevalence estimation for biased sampling with testing errors

Lili Zhou¹ | Daniel Andrés Díaz-Pachón¹ | Chen Zhao¹ | J. Sunil Rao² | Ola Hössjer³

Correspondence

Daniel Andrés Díaz-Pachón, Division of Biostatistics, University of Miami, 1120 NW 14th St, Room 1057, Miami, FL, 33136, USA.

Email: Ddiaz3@miami.edu

Funding information

Copeland Foundation Award, Department of Public Health Sciences, University of Miami, 2022, Grant/Award Number: 2022 Sampling for prevalence estimation of infection is subject to bias by both over-sampling of symptomatic individuals and error-prone tests. This results in naïve estimators of prevalence (ie, proportion of observed infected individuals in the sample) that can be very far from the true proportion of infected. In this work, we present a method of prevalence estimation that reduces both the effect of bias due to testing errors and oversampling of symptomatic individuals, eliminating it altogether in some scenarios. Moreover, this procedure considers stratified errors in which tests have different error rate profiles for symptomatic and asymptomatic individuals. This results in easily implementable algorithms, for which code is provided, that produce better prevalence estimates than other methods (in terms of reducing and/or removing bias), as demonstrated by formal results, simulations, and on COVID-19 data from the Israeli Ministry of Health.

KEYWORDS

active information, bias correction, COVID-19, maximum entropy, prevalence, sampling bias, testing errors

1 | INTRODUCTION

Estimation of disease prevalence is challenging. First, except for the hypothetical case of random errors, imperfect testing almost always distorts actual proportions. Second, it is not uncommon to have to derive estimates from samples that under-represent or fail to capture subpopulations that are at greatest risk or of interest. An example is estimating the general population prevalence of chronic hepatitis C (HCV) because of the challenges of sampling from subpopulations of former and current injecting drug users, the homeless or incarcerated.¹ Other examples include the over-representation of symptomatic individuals in a sample since these individuals are more likely to get tested than asymptomatic ones, with which the final estimates of prevalence inflates, since symptomatic individuals are also more likely to be truly infected than asymptomatic ones.²

This situation became clear during the recent COVID-19 pandemic: besides usual discussions of the error rates of PCR and rapid tests, surveillance mechanisms have usually relied on convenience sampling or contact tracing. Therefore sampling bias was also present. In the case of convenience sampling, because it passively waits for symptomatic individuals to get tested, whereas asymptomatic individuals have few reasons to do so. As for contact tracing, because it actively pursues infected individuals, ignoring the noninfected almost altogether. Besides this, contact tracing has also raised questions on privacy and individual liberties. Though this example corresponds to a non-probability COVID-19 sampling setting, the problem is of course more general. It applies not only to every form of prevalence estimation performed through testing—probabilistic or not—and even more general forms of selection bias. 6-9

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

 $\ @$ 2023 The Authors. Statistics in Medicine published by John Wiley & Sons Ltd.

¹Division of Biostatistics, University of Miami, Miami, Florida, USA

²Division of Biostatistics, University of Minnesota, Minneapolis, Minnesota, USA

³Department of Mathematics, Stockholm University, Stockholm, Sweden

Recently, Díaz-Pachón and Rao introduced a correction for oversampling of the symptomatic group. ¹⁰ It was a three-step procedure based on the assumption that all symptomatic individuals in the population were sampled and infected but it did not address the issue of imperfect testing (ie, the presence of false positives and false negatives). This implies that the symptomatic and infected individuals in the sample corresponded to the total number of symptomatic individuals in the population. Thus the asymptomatic group in the population was the complement of the total of symptomatic individuals in the sample. The prevalence among the asymptomatic group was then obtained as a uniform random variable among the asymptomatic individuals in the population, with no resource to the sample.

In this article, a method that is stronger in all aspects is presented. First, it does not assume that *all* symptomatic individuals are sampled, only that symptomatic individuals are overrepresented in the sample. Second, sample values among the asymptomatic are used to produce an estimator of prevalence that is informed by evidence. Third, testing errors are considered. And fourth, the proposed correction is extended to stratified errors by symptom status.

The article can be summarized as follows:

- 1. The researcher observes, from a sample of size N_T , the proportion of individuals whose test was positive. This constitutes the naïve estimator of prevalence $\tilde{\mathbf{p}}_T^{*,1}$.
- 2. The naïve estimator $\tilde{\mathbf{p}}_T^{*,1}$ is biased in two ways. First, it is subject to testing errors; and second, there is sample bias because symptomatic individuals are more likely to be tested than asymptomatic ones.
- 3. Under the assumption that, from a different and independent sample, testing error rates are estimated for group s as $\hat{\alpha}_s$ and $\hat{\beta}_s$ for false positives and false negatives, respectively, it is possible to obtain another estimator $\overline{p}_T^{*,1}$ that corrects the effect of errors.
- 4. From $\bar{p}_T^{*,1}$ it is possible to obtain another estimator $\hat{p}^{(1)}$ that reduces (and possibly removes) the sampling bias. This is achieved through applying the principle of maximum entropy to the fraction of symptomatic individuals with the disease, using the knowledge that symptomatic individuals are more likely to get tested.

With this summary and the information of Table 1, without having to go through the details that led to their derivation, the reader can obtain corrections following the steps of Algorithm 1 when no testing errors are present and all symptomatic individuals are sampled, Algorithm 2 when no testing errors are present and not all the symptomatic individuals are sampled (provided that the testing errors are estimated unbiasedly), and Algorithm 4 when there are testing errors and not all symptomatic individuals are sampled. Section 5 presents an example with real data. Proofs of all the results are consigned to the Appendix, as well as a set of simulations that assess the behavior of the four algorithms.

TABLE 1 A population of known size N is divided into symptomatic individuals (s = 1) and asymptomatic ones (s = 0), and noninfected individuals (i = 0) and infected ones (i = 1).

Quantity	Symptoms s and infection $iI_s^{(i)}$	Symptoms $s I_s$	Infection $i I^{(i)}$
Population totals	$N_{\scriptscriptstyle S}^{(i)}$	N_s	$N^{(i)}$
Population proportion	$p_s^{(i)}$	p_s	$p^{(i)}$
Sampling probability	$q\Big(I_s^{(i)}\Big)$	$q(I_s)$	$qig(I^{(i)}ig)$
Sampling probability approximation	$q^*\left(I_s^{(i)} ight)$	$q^*(I_s)$	$q^*ig(I^{(i)}ig)$
Sampling totals	$N_T^{s,i}$	$N_T^{{\scriptscriptstyle S},*}$	$N_T^{*,i}$
Naïve estimator (with errors and sampling bias)	$\tilde{\boldsymbol{p}}_{T}^{s,i}$	$\tilde{\boldsymbol{p}}_{T}^{s,*}$	$\tilde{\mathbf{p}}_T^{*,i}$
Naïve estimator (only with sampling bias)	$\mathbf{\hat{p}}_{T}^{s,i}$	$\mathbf{\hat{p}}_{T}^{s,*}$	$\mathbf{\hat{p}}_T^{*,i}$
Corrected estimator for errors	$\overline{p}_T^{s,i}$	$\overline{p}_T^{s,*}$	$\overline{p}_T^{*,i}$
Corrected estimator for errors and sampling bias	$\hat{p}_s^{(i)}$	\hat{P}_{s}	$\hat{p}^{(i)}$

Note: The second column gives the notation for symptoms *s* and infection *i*, while the third column marginalizes symptoms, and the fourth one marginalizes infection. To facilitate reading, the notation is arranged as follows: (a) The letter *q* is only used for sampling probabilities; (b) all estimators are capped by tildes, bars, or hats; (c) bold caps refer to naïve estimators; (d) estimators not in bold are corrections of naïve estimators; (e) population proportions do not have any cap; (f) prevalence values are obtained replacing *i* in the last column by 1.

2 | SETTING

Consider a population \mathcal{P} of size N that is divided into four categories: asymptomatic and noninfected individuals, $I_0^{(0)}$, with size $N_0^{(0)}$; asymptomatic and infected individuals, $I_0^{(1)}$, with size $N_0^{(1)}$; symptomatic and infected individuals, $I_1^{(1)}$, with size $N_1^{(1)}$; and symptomatic and noninfected individuals, $I_1^{(0)}$, with size $N_1^{(0)}$. The population total N is known, whereas $N_1^{(1)}$, $N_0^{(1)}$, $N_0^{(1)}$, and $N_0^{(0)}$ are unknown, though their sum is N.

The group of individuals with symptoms s in the population will be denoted by $I_s = I_s^{(0)} \cup I_s^{(1)}$, and its total by $N_s = N_s^{(0)} + N_s^{(1)}$, for s = 0, 1. Analogously, the group of individuals with infection status i in the population will be denoted by $I^{(i)} = I_0^{(i)} \cup I_1^{(i)}$, and its total by $N^{(i)} = N_0^{(i)} + N_1^{(i)}$, for i = 0, 1.

Now, $p_s^{(i)} = N_s^{(i)}/N$ will be the proportion of individuals in the population with symptoms s and infection status i. More formally, define a random element S^* taking values in the set $\mathbf{I} = \left\{I_0^{(0)}, I_0^{(1)}, I_1^{(0)}, I_1^{(1)}\right\}$, with density given by

$$f_{S^*}\left(I_s^{(i)}\right) = p_s^{(i)},\tag{1}$$

and $p_0^{(0)} + p_0^{(1)} + p_1^{(0)} + p_1^{(1)} = 1$. The proportion of individuals in the group I_s is then given by $p_s = p_s^{(0)} + p_s^{(1)}$, for s = 0, 1. And the proportion of individuals in the group $I^{(i)}$ is given by $p^{(i)} = p_0^{(i)} + p_1^{(i)}$, for i = 0, 1. The proportion to be estimated is $p^{(1)} = p_0^{(1)} + p_1^{(1)}$, corresponding to the infected individuals, and the naïve estimator is biased because the proportion of symptomatic individuals p_1 is overestimated.

2.1 | Sampling probabilities

For the *j*th individual in the population $(0 < j \le N)$, define a Bernoulli random variable as follows:

$$T_{j}|j \in I_{s}^{(i)} = \begin{cases} 1 & \text{with probability } q(I_{s}^{(i)}), \\ 0 & \text{with probability } 1 - q(I_{s}^{(i)}). \end{cases}$$
 (2)

That is, an individual in the category $I_s^{(i)}$ will be tested with probability $q(I_s^{(i)})$, for s, i = 0, 1.

The sampling probability $q\left(I_s^{(i)}\right)$ of individuals with symptoms s and infection i is approximated by

$$q^* \left(I_s^{(i)} \right) = \frac{N_T^{s,i}}{N_s^{(i)}},\tag{3}$$

where $N_T^{s,i}$ is the number of tested individuals from group $I_s^{(i)}$. Analogously to (3), $q(I_s)$, the sampling probability among individuals with symptoms s, is approximated by

$$q^*(I_s) = \frac{N_T^{s,*}}{N_s},\tag{4}$$

where $N_T^{s,*} = N_T^{s,0} + N_T^{s,1}$. And $q(I^{(i)})$, the sampling probability among individuals with infection status i, is approximated by

$$q^*(I^{(i)}) = \frac{N_T^{*,i}}{N^{(i)}},\tag{5}$$

where $N_T^{*,i} = N_T^{0,i} + N_T^{1,i}$. Notice that, except when all symptomatic individuals are sampled (in whose case N_1 is known), the approximations $q^*(\cdot)$ are not estimators of the sampling probabilities because the population values in their

0970258, 2023, 26, Downloaded from https://onlinelibrary.wiely.com/doi/10.1002/sim.9885, Wiley Online Library on [09/07/2024]. See the Terms and Conditions (https://onlinelibrary.wiely.com/rems-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Licenses

denominators are in general unknown. However, when $N \to \infty$,

$$q^*\left(I_s^{(i)}\right) \to q\left(I_s^{(i)}\right). \tag{6}$$

The total of sampled individuals $\sum_{j=1}^{N} T_j$ is defined as

$$N_T = \sum_{s,i} N_T^{s,i}. \tag{7}$$

Finally, we define the expected fraction of sampled individuals:

$$q = p_0 q(I_0) + p_1 q(I_1). (8)$$

3 | NO TESTING ERRORS

In case there is no error in testing, the naïve estimator of $p_s^{(i)}$ can be naturally defined as

$$\hat{\mathbf{p}}_{T}^{s,i} = f_{S^*|T} \Big(I_s^{(i)} | T = 1 \Big), \tag{9}$$

the conditional probability of an individual belonging to the group $I_s^{(i)}$, given that they were sampled. The main goal of this article is to provide a correction for $\hat{\mathbf{p}}_T^{s,i}$, under the assumption that symptomatic individuals are oversampled. A Bayesian approach, inspired from ideas in publication bias, ⁶ leads to

Proposition 1.

$$f_{S^*|T}\Big(I_s^{(i)}|T=1\Big) = \frac{N_T^{s,i}}{N_T}. (10)$$

Then $N_s^{(i)}$, the population size of $I_s^{(i)}$, disappears from the sample estimator, and (A1) in the Appendix shows that all information in the sample about the group $I_s^{(i)}$ comes from the sampling mechanism $q\left(I_s^{(i)}\right)$. In fact, $p_s^{(i)}$ can be seen as the message sent, $\hat{\mathbf{p}}_T^{s,i}$ as the message received, and $q\left(I_s^{(i)}\right) = P(T=1|I_s^{(i)})$ as the channel between them distorting the message. This interpretation, taken from Shannon's information diagram, is particularly important to analyze bias as a modification of the information inherent to the prevalence parameter in Appendix D. The property of the information inherent to the prevalence parameter in Appendix D.

Analogously to (9) with Proposition 1, the naive estimator of individuals with symptoms s, $\hat{\mathbf{p}}_{T}^{s,*}$, and the naive estimator of individuals with infection status i, $\hat{\mathbf{p}}_{T}^{*,i}$, are defined as

$$\hat{\mathbf{p}}_{T}^{s,*} = f_{S^*|T}(I_s|T=1) = \frac{N_T^{s,*}}{N_T},\tag{11}$$

$$\hat{\mathbf{p}}_{T}^{*,i} = f_{S^*|T} \left(I^{(i)}|T=1 \right) = \frac{N_T^{*,i}}{N_T}. \tag{12}$$

3.1 | Correction of sampling bias

According to (A1) in the appendix, some information about the sampling mechanism is needed if any meaningful conclusion is going to be obtained. For the scenario considered in this article, this corresponds to symptomatic individuals being more prone to get tested than asymptomatic ones:

$$q(I_0) < q(I_1). (13)$$

Also corresponding to the intuition that infected and noninfected individuals inside each category are randomly sampled,

$$q\left(I_s^{(0)}\right) = q\left(I_s^{(1)}\right),\tag{14}$$

for s = 0, 1.

We consider two scenarios. First, when all symptomatic individuals are sampled. Second, building on the previous case, when not all symptomatic are tested, but they are overrepresented in the sample.

3.1.1 | All the symptomatic group is sampled

Díaz-Pachón and Rao studied the situation in which, for COVID-19, all symptomatic individuals are tested. This scenario corresponds well to some subpopulations like those of universities or industries, in which all symptomatic individuals are required to get tested. In this case, (13) becomes

$$q(I_0) < q(I_1) = 1$$
,

so the proportion p_1 of symptomatic in the population can be fully recovered from the sample as

$$p_1 = \hat{\mathbf{p}}_T^{1,*} \frac{N_T}{N}. \tag{15}$$

Since, by (14) the sample is assumed to be random among symptomatic individuals,

$$\hat{p}_{1}^{(1)} := p_{1} \frac{N_{T}^{1,1}}{N_{T}^{1,*}} = p_{1} \frac{\hat{\mathbf{p}}_{T}^{1,1}}{\hat{\mathbf{p}}_{T}^{1,*}}.$$
(16)

Using now (14) on the asymptomatic group, the prevalence among the asymptomatic is obtained as

$$\hat{p}_0^{(1)} := p_0 \frac{N_T^{0,1}}{N_T^{0,*}} = (1 - p_1) \frac{\hat{\mathbf{p}}_T^{0,1}}{\hat{\mathbf{p}}_T^{0,*}}.$$
(17)

Using (16) and (17), the final sampling-bias corrected prevalence is then taken to be

$$\hat{p}^{(1)} := \hat{p}_1^{(1)} + \hat{p}_0^{(1)},\tag{18}$$

and the random sampling inside each group gives that $E\hat{p}^{(1)} = p^{(1)}$, making the estimator unbiased.

Algorithm 1 summarizes the steps from observations to corrected estimate when there is no testing error and all symptomatic individuals are tested.

Remark 1. Although the scenario of all sampled individuals is also considered by Díaz-Pachón and Rao in the context of COVID-19, 10 the correction obtained by Algorithm 1 is stronger than theirs in at least two aspects. First, Díaz-Pachón and Rao considered $m \geq 2$ categories of symptoms, so that, if m is large, an individual with all symptoms is highly likely to be infected; however, with two categories of symptoms, this would correspond to $p_1 = p_1^{(1)}$, which seems a very strong assumption. Second, Algorithm 1 takes into account the information from the sample to obtain $\hat{p}_0^{(1)}$ (see step 3 or (17)); instead, Díaz-Pachón and Rao proposed to take U uniformly distributed in the interval [0,1] and make $\hat{p}_0^{(1)} = Up_0$.

The following result shows that $\hat{\mathbf{p}}_T^{*,1}$ and $\hat{p}^{(1)}$ are asymptotically normal. More specifically, it is proven that $\hat{\mathbf{p}}_T^{*,1}$ is not a consistent estimator of the true prevalence $p^{(1)}$, but $\hat{p}^{(1)}$ is.

0970258, 2023, 26, Downloaded from https://onlinelibrary.wiely.com/doi/10.1002/sim.9885, Wiley Online Library on [09/07/2024]. See the Terms and Conditions (https://onlinelibrary.wiely.com/rems-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Licenses

Algorithm 1. Corrected estimator of prevalence without errors and all symptomatic individuals sampled

1. For $\hat{\mathbf{p}}_{T}^{1,*} = \hat{\mathbf{p}}_{T}^{1,0} + \hat{\mathbf{p}}_{T}^{1,1}$, take

$$p_1 = \hat{\mathbf{p}}_T^{1,*} \frac{N_T}{N}.$$

2. Make

$$\hat{p}_1^{(1)} = p_1 \frac{\hat{\mathbf{p}}_T^{1,1}}{\hat{\mathbf{p}}_T^{1,*}}.$$

- 3. Take $\hat{p}_0^{(1)} = \frac{\hat{\mathbf{p}}_T^{0,1}}{\hat{\mathbf{p}}_U^{0,*}} (1 p_1)$, where $\hat{\mathbf{p}}_T^{0,*} = \hat{\mathbf{p}}_T^{0,0} + \hat{\mathbf{p}}_T^{0,1}$.
- 4. The estimated total prevalence is: $\hat{p}^{(1)} = \hat{p}_0^{(1)} + \hat{p}_1^{(1)}$.

Theorem 1. Suppose $N \to \infty$ in such a way that $p_1 = N_1/N$ is kept fixed. Then

$$N^{1/2} \left(\hat{\mathbf{p}}_T^{*,1} - \frac{p_0^{(1)} q(I_0) + p_1^{(1)}}{p_0 q(I_0) + p_1} \right) \xrightarrow{\mathcal{L}} Z_0, \tag{19}$$

$$N^{1/2}\left(\hat{p}^{(1)} - p^{(1)}\right) \xrightarrow{\mathcal{L}} Z_1,\tag{20}$$

as $N \to \infty$, where $Z_0 \sim \mathcal{N}(0, V_{01} + V_{02})$, $Z_1 \sim \mathcal{N}(0, V_{03})$ are normally distributed, and V_{01} , V_{02} , and V_{03} are defined in the Appendix.

3.1.2 | Not all the symptomatic group is sampled

The main difference between this section and the previous one is that, since now the proportion of symptomatic individuals in the population is not known, it has to be estimated from the sample. Drawing inspiration from cosmological fine-tuning, ^{14,15} the approach will be to use the information in (13) to generate a maximum entropy distribution, which is "the least biased estimate possible on the given information." Next we will use (14) to obtain estimators of prevalence inside each class of symptoms.

Theorem 2. For $\hat{\mathbf{p}}_T^{1,*}, p_1 \in (0,1), q^*(I_0) < q^*(I_1)$ if and only if $\hat{\mathbf{p}}_T^{1,*} > p_1$.

Theorem 2 shows that, given the basic assumption (13), with high probability p_1 is bounded above by $\hat{\mathbf{p}}_T^{1,*}$. On the other hand, $\hat{\mathbf{p}}_T^{1,*} = N_T^{1,*}/N_T$ says that there are at least $N_T^{1,*}$ infected symptomatic individuals in the population. Therefore, it makes sense to bound p_1 as follows:

$$N_T^{1,*}/N \le p_1 \le \hat{\mathbf{p}}_T^{1,*}. \tag{21}$$

By the maximum entropy principle,¹⁷ the corrected estimator of p_1 is taken to be the expectation of a uniform distribution over $(N_T^{1,*}/N, \hat{\mathbf{p}}_T^{1,*})$. Formally, let U be a uniform distribution over the interval $(\hat{\mathbf{p}}_T^{1,*} \frac{N_T}{N}, \hat{\mathbf{p}}_T^{1,*})$. The corrected estimator of p_1 is defined as

$$\hat{p}_1 := E(U) = \frac{\hat{\mathbf{p}}_T^{1,*}}{2} \left(\frac{N_T}{N} + 1 \right). \tag{22}$$

From this point on, we proceed analogously to Section 3.1.1, replacing p_1 with \hat{p}_1 in Equations (16) to (18), to obtain

$$\hat{p}_{1}^{(1)} := \hat{p}_{1} \frac{N_{T}^{1,1}}{N_{T}^{1,*}} = \hat{p}_{1} \frac{\hat{\mathbf{p}}_{T}^{1,1}}{\hat{\mathbf{p}}_{T}^{1,*}},\tag{23}$$

0970288, 2023, 2,6, Downloaded from https://onlinelibrary.wiely.com/doi/10.1002/sim.9885, Wiley Online Library on [0907/2024]. See the Terms and Conditions (https://onlinelibrary.wiely.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Licensea

$$\hat{p}_{0}^{(1)} := \hat{p}_{0} \frac{N_{T}^{0,1}}{N_{T}^{0,*}} = (1 - \hat{p}_{1}) \frac{N_{T}^{0,1}}{N_{T}^{0,*}} = (1 - \hat{p}_{1}) \frac{\hat{\mathbf{p}}_{T}^{0,1}}{\hat{\mathbf{p}}_{T}^{0,*}}, \tag{24}$$

$$\hat{p}^{(1)} := \hat{p}_1^{(1)} + \hat{p}_0^{(1)}. \tag{25}$$

Algorithm 2 summarizes the procedure to produce estimators that correct the sampling bias when there is no testing error and not all symptomatic individuals are sampled.

Algorithm 2. Corrected estimator of prevalence without errors and not all symptomatic individuals sampled

1. For $\hat{\mathbf{p}}_{T}^{1,*} = \hat{\mathbf{p}}_{T}^{1,0} + \hat{\mathbf{p}}_{T}^{1,1}$, take

$$\hat{p}_1 = \frac{\hat{\mathbf{p}}_T^{1,*}}{2} \left(\frac{N_T}{N} + 1 \right).$$

2. Make

$$\hat{p}_1^{(1)} = \hat{p}_1 \frac{\hat{\mathbf{p}}_T^{1,1}}{\hat{\mathbf{p}}_T^{1,*}}.$$

- 3. Take $\hat{p}_0^{(1)} = \frac{\hat{\mathbf{p}}_T^{0,1}}{\hat{\mathbf{p}}_0^{0,*}} (1 \hat{p}_1)$, where $\hat{\mathbf{p}}_T^{0,*} = \hat{\mathbf{p}}_T^{0,0} + \hat{\mathbf{p}}_T^{0,1}$.
- 4. The estimated total prevalence is: $\hat{p}^{(1)} = \hat{p}_0^{(1)} + \hat{p}_1^{(1)}$.

Analogously to the previous asymptotic result, Theorem 3 shows that the naïve and corrected estimators are asymptotically normal. However, once again we find that the naïve estimator is not a consistent estimator of the true prevalence. Moreover, the corrected estimator will only be consistent if $E(\hat{p}_s) = p_s$, for s = 0, 1.

Theorem 3. Suppose $N \to \infty$ in such a way that $p_1 = N_1/N$ is kept fixed. Suppose additionally that, for s = 0, 1, there exists $\rho_s \in [0, 1]$, such that

$$\hat{p}_s \xrightarrow{p} \rho_s \tag{26}$$

as $N \to \infty$, where " $\stackrel{p}{\longrightarrow}$ " refers to convergence in probability. Then

$$N^{1/2} \left(\hat{\mathbf{p}}_T^{*,1} - \frac{p_0^{(1)} q(I_0) + p_1^{(1)} q(I_1)}{p_0 q(I_0) + p_1 q(I_1)} \right) \xrightarrow{\mathcal{L}} Z_2, \tag{27}$$

$$N^{1/2} \left(\hat{p}^{(1)} - \frac{\rho_0}{p_0} p_0^{(1)} - \frac{\rho_1}{p_1} p_1^{(1)} \right) \xrightarrow{\mathcal{L}} Z_3, \tag{28}$$

as $N \to \infty$, where $Z_2 \sim \mathcal{N}(0, V_{11} + V_{12})$, $Z_3 \sim \mathcal{N}(0, V_{13} + V_{14})$ are normally distributed random variables, and V_{11} , V_{12} , V_{13} , and V_{14} are given in the Appendix.

4 | WITH TESTING ERRORS

When testing errors are considered, the naïve estimators have an additional source of bias. Using (10), in this section, we present first the explicit form of the naïve estimators in the presence of sampling bias and testing errors stratified by

0970288, 2023, 2,6, Downloaded from https://onlinelibrary.wiely.com/doi/10.1002/sim.9885, Wiley Online Library on [0907/2024]. See the Terms and Conditions (https://onlinelibrary.wiely.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Licenses

symptoms. This is the most general form of naïve estimator considered in this article. As a corollary, unstratified errors are also considered. With naïve estimators in this general form, we then present their respective corrections.

Proposition 2. Let α_0 and β_0 be the false positive and false negative rate for asymptomatic individuals, respectively, and let α_1 and β_1 be the false positive and false negative rate for symptomatic individuals, respectively. The naïve estimators thus become:

$$\tilde{\mathbf{p}}_{T}^{0,0} = (1 - \check{\alpha}_{0})\hat{\mathbf{p}}_{T}^{0,0} + \check{\beta}_{0}\hat{\mathbf{p}}_{T}^{0,1},
\tilde{\mathbf{p}}_{T}^{0,1} = \check{\alpha}_{0}\hat{\mathbf{p}}_{T}^{0,0} + (1 - \check{\beta}_{0})\hat{\mathbf{p}}_{T}^{0,1},
\tilde{\mathbf{p}}_{T}^{1,0} = (1 - \check{\alpha}_{1})\hat{\mathbf{p}}_{T}^{1,0} + \check{\beta}_{1}\hat{\mathbf{p}}_{T}^{1,1},
\tilde{\mathbf{p}}_{T}^{1,1} = \check{\alpha}_{1}\hat{\mathbf{p}}_{T}^{1,0} + (1 - \check{\beta}_{1})\hat{\mathbf{p}}_{T}^{1,1},$$
(29)

where $\check{\alpha}_s$ and $\check{\beta}_s$, the proportion of false positives and false negatives in the sample, for s=0,1, approximate α_s and β_s respectively.

Analogously to previous definitions, let $\tilde{\mathbf{p}}_T^{s,*} = \tilde{\mathbf{p}}_T^{s,0} + \tilde{\mathbf{p}}_T^{s,1}$ and $\tilde{\mathbf{p}}_T^{*,i} = \tilde{\mathbf{p}}_T^{0,i} + \tilde{\mathbf{p}}_T^{1,i}$.

Corollary 1. If errors are not stratified, for s = 0, 1, the estimators of Proposition 2,

$$\tilde{\mathbf{p}}_{T}^{s,0} = (1 - \check{\alpha}_{s})\hat{\mathbf{p}}_{T}^{s,0} + \check{\beta}_{s}\hat{\mathbf{p}}_{T}^{s,1},$$

$$\tilde{\mathbf{p}}_{T}^{s,1} = \check{\alpha}_{s}\hat{\mathbf{p}}_{T}^{s,0} + (1 - \check{\beta}_{s})\hat{\mathbf{p}}_{T}^{s,1},$$
(30)

are such that $\check{\alpha}_s$ and $\check{\beta}_s$, the proportion of false positives and negatives in the sample for symptom class s, approximate the probabilities α and β of false positives and false negatives, respectively, independently of s.

Remark 2. The right-hand side of (29) and (30) contains the contribution to the naïve estimator by each group in the sample weighted by the probability of their errors. However, in either case, the proportions observed by practitioner are the tilde terms $\mathbf{\tilde{p}}_{T}^{s,i}$ in the left-hand side. The hat terms in the right-hand side, corresponding to (10) are unknown to him.

4.1 | Correction of testing errors

According to Remark 2, when testing errors are considered, estimators that correct them are necessary before applying the correction to sampling bias. This section presents such estimators.

Proposition 3. For s = 0, 1, assume $\hat{\alpha}_s$ and $\hat{\beta}_s$ are estimators of α_s and β_s , respectively, obtained from different data, satisfying also that they are independent of $\check{\alpha}_s$ and $\check{\beta}_s$. Then the estimators with correction for testing errors

$$\overline{p}_{T}^{s,0} = \hat{\mathbf{p}}_{T}^{s,*} \cdot \frac{\tilde{\mathbf{p}}_{T}^{s,0}/\hat{\mathbf{p}}_{T}^{s,*} - \hat{\beta}_{s}}{1 - \hat{\alpha}_{s} - \hat{\beta}_{s}} = \frac{\tilde{\mathbf{p}}_{T}^{s,0} - \hat{\beta}_{s}\tilde{\mathbf{p}}_{T}^{s,*}}{1 - \hat{\alpha}_{s} - \hat{\beta}_{s}},$$

$$\overline{p}_{T}^{s,1} = \hat{\mathbf{p}}_{T}^{s,*} \cdot \frac{\tilde{\mathbf{p}}_{T}^{s,1}/\hat{\mathbf{p}}_{T}^{s,*} - \hat{\alpha}_{s}}{1 - \hat{\alpha}_{s} - \hat{\beta}_{s}} = \frac{\tilde{\mathbf{p}}_{T}^{s,1} - \hat{\alpha}_{s}\tilde{\mathbf{p}}_{T}^{s,*}}{1 - \hat{\alpha}_{s} - \hat{\beta}_{s}}$$

$$(31)$$

approximate $\hat{\mathbf{p}}_T^{s,0}$ and $\hat{\mathbf{p}}_T^{s,1}$ respectively. Moreover, $\overline{p}_T^{s,0}/\overline{p}_T^{s,*}$ and $\overline{p}_T^{s,1}/\overline{p}_T^{s,*}$ are consistent estimators of $p_s^{(0)}/p_s$ and $p_s^{(1)}/p_s$, if $\hat{\alpha}_s$ and $\hat{\beta}_s$ are in turn consistent estimators of α_s and β_s , respectively.

4.1.1 | All symptomatic group is sampled

When all the symptomatic group is sampled, if $\check{\alpha}_s$ and $\hat{\alpha}_s$ are unbiased estimators of α_s , and $\check{\beta}_s$ and $\hat{\beta}_s$ are unbiased estimators of β_s , for s=0,1, following Algorithm 1, we obtain Algorithm 3 substituting $\hat{\mathbf{p}}_T^{s,i}$ by $\overline{p}_T^{s,i}$.

Algorithm 3. Corrected estimator of prevalence with errors and all symptomatic individuals sampled

1. For s = 0, 1, make

$$\begin{split} \bar{p}_T^{s,0} &= \frac{\tilde{\mathbf{p}}_T^{s,0} - \hat{\beta}_s \tilde{\mathbf{p}}_T^{s,*}}{1 - \hat{\alpha}_s - \hat{\beta}_s}, \\ \bar{p}_T^{s,1} &= \frac{\hat{\alpha}_s \tilde{\mathbf{p}}_T^{s,*} - \tilde{\mathbf{p}}_T^{s,1}}{\hat{\alpha}_s + \hat{\beta}_s - 1}. \end{split}$$

2. If $\bar{\alpha}_s = E\hat{\alpha}_s$ equals $E\check{\alpha}_s = \alpha_s$ and if $\bar{\beta}_s = E\hat{\beta}_s$ equals $E\check{\beta}_s = \beta_s$, for $\bar{p}_T^{1,*} = \bar{p}_T^{1,0} + \bar{p}_T^{1,1}$, take

$$\hat{p}_1 = \bar{p}_T^{1,*} \frac{N_T}{N}.$$

3. Make

$$\hat{p}_1^{(1)} = \hat{p}_1 \frac{\bar{p}_T^{1,1}}{\bar{p}_T^{1,*}}.$$

- 4. Take $\hat{p}_0^{(1)} = \frac{\bar{p}_T^{0,1}}{\bar{p}_0^{0,*}} (1 \hat{p}_1)$, where $\bar{p}_T^{0,*} = \bar{p}_T^{0,0} + \bar{p}_T^{0,1}$.
- 5. The estimated total prevalence is: $\hat{p}^{(1)} = \hat{p}_0^{(1)} + \hat{p}_1^{(1)}$.

4.1.2 | Not all the symptomatic group is sampled

Analogously to Section 3.1.2, replace $\hat{p}_T^{s,i}$ with $\overline{p}_T^{s,1}$ in Equations (22) to (25), to obtain

$$\hat{p}_1 := \frac{\bar{p}_T^{1,*}}{2} \left(\frac{N_T}{N} + 1 \right), \tag{32}$$

$$\hat{p}_1^{(1)} := \hat{p}_1 \frac{\overline{p}_T^{1,1}}{\overline{p}_T^{1,*}},\tag{33}$$

$$\hat{p}_0^{(1)} := \left(1 - \hat{p}_1\right) \frac{\overline{p}_T^{0,1}}{\overline{p}_T^{0,*}},\tag{34}$$

$$\hat{p}^{(1)} := \hat{p}_1^{(1)} + \hat{p}_0^{(1)}. \tag{35}$$

This information can be used to generate Algorithm 4. Theorem 4 summarizes the asymptotic behavior of the estimators involved in this section.

Theorem 4. Suppose the conditions of Theorem 3 hold, and that the estimators prevalences $\overline{p}_T^{s,0}/\overline{p}_T^{s,*}$ and $\overline{p}_T^{s,1}/\overline{p}_T^{s,*}$ converge to

$$\overline{\rho}_{s}^{(0)} = \frac{\beta_{s} - \overline{\beta}_{s}}{1 - \overline{\alpha}_{s} - \overline{\beta}_{s}} + \frac{1 - \alpha_{s} - \beta_{s}}{(1 - \overline{\alpha}_{s} - \overline{\beta}_{s})} \cdot \frac{p_{s}^{(0)}}{p_{s}}, \tag{36}$$

$$\overline{\rho}_{s}^{(1)} = \frac{\alpha_{s} - \overline{\alpha}_{s}}{1 - \overline{\alpha}_{s} - \overline{\beta}_{s}} + \frac{1 - \alpha_{s} - \beta_{s}}{(1 - \overline{\alpha}_{s} - \overline{\beta}_{s})} \cdot \frac{p_{s}^{(1)}}{p_{s}}$$
(37)

Algorithm 4. Corrected estimator of prevalence with errors and not all symptomatic individuals sampled

1. For s = 0, 1, make

$$\begin{split} \bar{p}_T^{s,0} &= \frac{\tilde{\mathbf{p}}_T^{s,0} - \hat{\beta}_s \tilde{\mathbf{p}}_T^{s,*}}{1 - \hat{\alpha}_s - \hat{\beta}_s}, \\ \bar{p}_T^{s,1} &= \frac{\hat{\alpha}_s \tilde{\mathbf{p}}_T^{s,*} - \tilde{\mathbf{p}}_T^{s,1}}{\hat{\alpha}_s + \hat{\beta}_s - 1}. \end{split}$$

2. If $\bar{\alpha}_s = E\hat{\alpha}_s$ equals $E\check{\alpha}_s = \alpha_s$ and if $\bar{\beta}_s = E\hat{\beta}_s$ equals $E\check{\beta}_s = \beta_s$, for $\bar{p}_T^{1,*} = \bar{p}_T^{1,0} + \bar{p}_T^{1,1}$, take

$$\hat{p}_1 = \frac{\bar{p}_T^{1,*}}{2} \left(\frac{N_T}{N} + 1 \right).$$

3. Make

$$\hat{p}_1^{(1)} = \hat{p}_1 \frac{\bar{p}_T^{1,1}}{\bar{p}_T^{1,*}}.$$

- 4. Take $\hat{p}_0^{(1)} = \frac{\bar{p}_D^{0,1}}{\bar{p}_D^{0,*}} (1 \hat{p}_1)$, where $\bar{p}_T^{0,*} = \bar{p}_T^{0,0} + \bar{p}_T^{0,1}$.
- 5. The estimated total prevalence is: $\hat{p}^{(1)} = \hat{p}_0^{(1)} + \hat{p}_1^{(1)}$.

as $N \to \infty$. Suppose further that the proportion of individuals wrongly classified in the sample $\check{\alpha}_s$, $\check{\beta}_s$, s=0,1 are independent. Then

$$N^{1/2} \left(\overline{p}_T^{*,1} - \frac{p_0 q(I_0) \overline{\rho}_0^{(1)} + p_1 q(I_1) \overline{\rho}_1^{(1)}}{q} \right) \xrightarrow{\mathcal{L}} Z_4, \tag{38}$$

$$N^{1/2} \left[\hat{p}^{(1)} - \left(\rho_0 \overline{\rho}_0^{(1)} + \rho_1 \overline{\rho}_1^{(1)} \right) \right] \xrightarrow{\mathcal{L}} Z_5, \tag{39}$$

as $N \to \infty$, where $\hat{p}^{(1)}$ is defined as in (35), $Z_4 \sim \mathcal{N}(0, V_5 + V_6)$, $Z_5 \sim \mathcal{N}(0, V_7 + V_8)$ are normally distributed random variables, and V_5 , V_6 , V_7 , and V_8 are defined in the Appendix.

Remark 3. If stratification is ignored, throughout all this section just take $\alpha = \alpha_0 = \alpha_1$ and $\hat{\alpha} = \hat{\alpha}_0 = \hat{\alpha}_1$. On the other hand, $\check{\alpha}_0$ and $\check{\alpha}_1$ are still distinct quantities, since they correspond to the observed positive error rates of stratum s = 0 and s = 1 respectively. Then do analogously with the beta terms to obtain β , $\hat{\beta}$, $\check{\beta}_1$, and $\check{\beta}_2$.

5 | DATA FROM THE ISRAELI MINISTRY OF HEALTH

In what follows, COVID-19 data from the Israeli Ministry of Health is considered. The Ministry of Health publicly released data for individuals tested for COVID-19 via a PCR assay from a nasal swab sample collected between March 22, 2020 and April 7, 2020. The dataset contains information on the test date, test result, clinical symptoms, gender of the individual, known contact with an infected individual and a binary indicator of whether the individual was 60 years of age or older. Symptoms include cough, fever, sore throat, shortness of breath and headache. For the purposes of illustrating the methodology, we will consider this as the population, consisting of 99 232 tested individuals, among whom 1862 were symptomatic (have shortness of breath or have at least three of four symptoms: cough, fever, sore throat, and headache) and 97 370 were asymptomatic. Among the total tested individuals, it was possible to identify 8393 infections through PCR testing. Among the individuals who tested positive, 1754 were symptomatic. The characteristics of the data set are presented in Table 2.

	Positive	Negative	Total
Symptomatic	1754	108	1862
Asymptomatic	6639	90 731	97 370
Total	8393	90 839	99 232

TABLE 3 Real proportions under stratified errors with different error combinations.

Error Comb.	Symptoms	Disease	Non-disease	Total
Comb. 1	Sympt.	1704	158	1862
	Asympt.	24719	72 651	97 370
	Total	26 423	72 809	99 232
Comb. 2	Sympt.	1686	176	1862
	Asympt.	15 646	81 724	97 370
	Total	17 332	81 900	99 232
Comb. 3	Asympt.	1707	155	1862
	Asympt.	20 116	77 254	97 370
	Total	21 823	77 409	99 232

Error rates will be stratified by symptoms. Thus, let α_0 and α_1 be the false positive rate for asymptomatic and symptomatic individuals, respectively, and β_0 and β_1 , the false negative rate for asymptomatic and symptomatic individuals, respectively. Although we do not have exact numbers of stratified false-positive and false-negative rates for PCR tests, a public report from UK Government Office for Science in 2020 indicated that the median false positive rate in the UK's COVID-19 RT-PCR testing program is 2.3% with IQR of 0.8% to 4.0%. Moreover, Arévalo-Rodríguez et al stated that after they collected information among all patients from 34 studies, the summary estimate of the false-negative rate was 13% with range of 1.8% to 58%. For the purpose of investigating how the estimates change with different assumed error rates, we will compare different combinations of error rates within the reasonable ranges according to the literature we found. The combinations of error rates in the population are assumed to have the values in (40). The actual number of individuals inside each group can be found in Table 3 after correcting Table 2 for these errors.

Comb. 1
$$\alpha_0 = 1\%$$
, $\alpha_1 = 3\%$, $\beta_0 = 20\%$, $\beta_1 = 2\%$;
Comb. 2 $\alpha_0 = 1\%$, $\alpha_1 = 4\%$, $\beta_0 = 10\%$, $\beta_1 = 2\%$;
Comb. 3 $\alpha_0 = 2\%$, $\alpha_1 = 3\%$, $\beta_0 = 15\%$, $\beta_1 = 5\%$. (40)

The real prevalence with **Comb. 1** is then

$$p^{(1)} = (1704 + 24719)/99232 = 0.266, (41)$$

and prevalence among the asymptomatic is 24719/97370 = 0.254.

The real prevalence with **Comb. 2** is then

$$p^{(1)} = (1686 + 15646)/99232 = 0.175, (42)$$

and prevalence among the asymptomatic is 15646/97370 = 0.161.

The real prevalence with **Comb. 3** is then

$$p^{(1)} = (1707 + 20\ 116)/99\ 232 = 0.220, (43)$$

and prevalence among the asymptomatic is 20116/97370 = 0.207.

TABLE 4 Stratified observed sample for 75% symptomatic and 25% asymptomatic, with all symptomatic individuals sampled.

Error Comb.	Symptoms	Positive	Negative	Total
Comb. 1	Symptomatics	1641	221	1862
	Asymptomatics	122	499	621
	Total	1763	720	2483
Comb. 2	Symptomatics	1642	220	1862
	Asymptomatics	94	527	621
	Total	1736	747	2483
Comb. 3	Symptomatics	1541	321	1862
	Asymptomatics	117	503	620
	Total	1658	824	2482

Additionally, in (44) we assume some sample error rates for the combinations (40). We emphasize that the actual values of (44) are "known unknowns" to the practitioner,²¹ and it is precisely their effect what needs to be corrected.

Comb. 1
$$\check{\alpha}_0 = 0.7\%$$
, $\check{\alpha}_1 = 3.5\%$, $\check{\beta}_0 = 25\%$, $\check{\beta}_1 = 4\%$;
Comb. 2 $\check{\alpha}_0 = 1.2\%$, $\check{\alpha}_1 = 3.6\%$, $\check{\beta}_0 = 12\%$, $\check{\beta}_1 = 3\%$;
Comb. 3 $\check{\alpha}_0 = 2.5\%$, $\check{\alpha}_1 = 2.9\%$, $\check{\beta}_0 = 18\%$, $\check{\beta}_1 = 10\%$. (44)

Finally, active information (defined in Appendix C) is used to compare how well Algorithms 3 and 4 and other proposed estimators in the literature are doing with respect to the real prevalence. The best estimator will be the one with active information \hat{I}^+ closer to 0. The competitors will be the method proposed by Díaz-Pachón and Rao, which assumes all symptomatic individuals are sampled, correcting only for sample bias and ignoring testing errors; ¹⁰ Diggle's Bayesian approach, which corrects for imperfect testing but ignores sampling bias; ²² and the Rogan-Gladen estimate, a frequentist method that only corrects for testing errors too. ²³ Neither of the competitors corrects for sampling bias and testing errors at the same time. In fact, as much as we searched, we could not find a methodology that simultaneously corrects for imperfect testing and sampling bias; this will be reflected in the analysis. All of the results are presented in Table 8.

Each of the following protocols presents a table with the sample results. These correspond to the observations given by (29). These values, the population size, and the estimated error rates from a different study ($\hat{\alpha}_s$ and $\hat{\beta}_s$) will be the input of an R program, for which the code is available at https://github.com/kalilizhou/BiasCorrection.git, with the four algorithms proposed in this article. The program thus obtains the correction. As a simplifying assumption in the remaining of this section, we take $\hat{\alpha}_s = \alpha_s$ and $\hat{\beta}_s = \beta_s$, for s = 0, 1.

Sampling Protocol 1: In the first scenario, all symptomatic individuals are sampled, as considered by Díaz-Pachón and Rao. ¹⁰ The sample consists of 2483 individuals. Among these, 1862 (75%) are symptomatic and 621 (25%) are asymptomatic. The sample error rates are taken from (44). The observed sampling results, corresponding to (29), are given by Table 4.

According to Table 4, for instance with **Comb. 1**, the naïve estimator is $\tilde{\mathbf{p}}_T^{1,*} = 1763/2483 \approx 0.71$. Since all the symptomatic group is sample, we use Algorithm 3, which produces the corrected estimator $\hat{p}^{(1)} = 0.248$. Table 8 presents these results as well as those of the other methods.

In this case, Diggle's correction was not implemented because it involves combinations in its logarithm that are difficult to approximate when the sample is moderately large. Under the assumption of sampling all symptomatic individuals, Díaz-Rao works very well, and RGE performs as poorly as the naïve estimators. Our Algorithm 3 is the best correction to the naïve estimate, producing the closest-to-zero actinfo. The corrected estimators of prevalence obtained from Algorithm 3 almost equal the real prevalence for all combinations of testing errors.

For the next protocols, the assumption that all symptomatic individuals were sampled is removed, which implies that the Diaz-Rao correction cannot be assessed and Algorithm 4 is followed.

TABLE 5 Stratified observed sample for 75% symptomatic and 25% asymptomatic (not all symptomatic individuals sampled).

Error Comb.	Symptoms	Positive	Negative	Total
Comb. 1	Sympt.	132	18	150
	Asympt.	10	40	50
	Total	142	58	200
Comb. 2	Sympt.	132	18	150
	Asympt.	8	42	50
	Total	140	60	200
Comb. 3	Sympt.	125	25	150
	Asympt.	9	41	50
	Total	134	66	200

TABLE 6 Stratified observed sample for 50% symptomatic and 50% asymptomatic.

Error Comb.	Symptoms	Positive	Negative	Total
Comb. 1	Sympt.	89	11	100
	Asympt.	19	81	100
	Total	108	92	200
Comb. 2	Sympt.	89	11	100
	Asympt.	15	85	100
	Total	104	96	200
Comb. 3	Sympt.	83	17	100
	Asympt.	19	81	100
	Total	102	98	200

Sampling Protocol 2: The sample consists of 200 individuals. Among these, 150 (75%) are symptomatic and 50 (25%) are asymptomatic. For both the symptomatic and asymptomatic groups, the sampling proportions are taken according to Table 3. Table 5 shows the observed sample. The summary of results under different methods is shown in Table 8.

Table 8 shows that, in this sampling scenario, our Algorithm 4 still has the best performance. In fact, Diggle's and Rogan-Gladen's estimates do as poorly as the naïve estimate. Algorithm 4 beats its competitors because it is the only one that corrects for sampling bias, whereas the other two only correct for testing errors. Notice that the additional information of Protocol 1 (knowing that all symptomatic individuals were sampled), in comparison to Protocol 2, greatly improves the performance of the correction, as reflected by the active information.

Sampling Protocol 3: In this scenario there are 100 symptomatic and 100 asymptomatic individuals. Again, the proportions inside each group were taken from Table 3. The observed sample is given in Table 6. After correcting the estimates, the summary of results under different methods for this sampling protocol is also presented in Table 8.

Compared to Protocol 2, Protocol 3 has less sampling bias. Therefore, all the methods perform better than in the previous scenario. But Algorithm 4 still works better than competitors.

Sampling Protocol 4: This sample is truly random, with $N_T = 200$, and it is obtained from Table 3. The observed sample is presented in Table 7. The results of the different methods for this scenario are presented in Table 8.

In this scenario, without sampling bias, all estimates perform extremely well, with Rogan-Gladen's frequentist estimates being optimal. Algorithm 4 works quite well. The results of these two methods are very close to the real prevalence.

TABLE 7 Stratified observed random sample of size 200.

Error Comb.	Symptoms	Positive	Negative	Total
Comb. 1	Sympt.	3	1	4
	Asympt.	39	157	196
	Total	42	158	200
Comb. 2	Sympt.	3	1	4
	Asympt.	30	166	196
	Total	33	167	200
Comb. 3	Sympt.	3	1	4
	Asympt.	37	159	196
	Total	40	160	200

TABLE 8 Comparison among methods of all sampling protocols with different combinations of stratified error rates.

Corrected estimates $\hat{\pmb{p}}^{(1)}$ (active information \hat{I}^+)								
Error Comb.	Real prevalence	Sampling protocol	Naïve	Díaz-Rao	Diggle	Rogan-Gladen	Algorithms 3 and 4	
Comb. 1	p = 0.266	1	0.710 (0.981)	0.212 (-0.229)	-	0.729 (1.009)	0.248 (-0.069)	
		2	0.710 (0.981)	-	0.784 (1.081)	0.729 (1.009)	0.486 (0.602)	
		3	0.540 (0.707)	-	0.592 (0.800)	0.564 (0.751)	0.398 (0.401)	
		4	0.210 (-0.237)	-	0.220 (-0.190)	0.249 (-0.066)	0.244 (-0.086)	
Comb. 2	p = 0.175	1	0.699 (1.387)	0.167 (-0.045)	-	0.711 (1.402)	0.173 (-0.011)	
		2	0.700 (1.388)	-	0.772 (1.484)	0.712 (1.403)	0.441 (0.926)	
		3	0.520 (1.091)	-	0.570 (1.181)	0.530 (1.109)	0.344 (0.679)	
		4	0.165 (-0.057)	-	0.168 (-0.041)	0.173 (-0.014)	0.168 (-0.047)	
Comb. 3	p = 0.220	1	0.668 (1.111)	0.204 (-0.076)	-	0.701 (1.159)	0.215 (-0.019)	
		2	0.670 (1.114)	-	0.738 (1.210)	0.703 (1.161)	0.448 (0.713)	
		3	0.510 (0.841)	-	0.558 (0.931)	0.537 (0.892)	0.371 (0.524)	
		4	0.200 (-0.095)	-	0.208 (-0.056)	0.215 (-0.024)	0.209 (-0.050)	

Note: Bold values represent the closest value to p in the second column.

Table 8 summarizes the results of all protocols and all combinations of testing errors. Our proposed Algorithms 3 (first row in each combination) and 4 (rows 2-4 in each combination) are always the best in Sampling Protocols 1-3 and perform as well as the naïve estimator under random sampling (Protocol 4), according to the actinfo assessment. The result shows the promising ability of our proposed algorithms to correct for both sampling bias and testing errors in prevalence estimation.

6 | DISCUSSION

Timely and accurate prevalence estimation of a disease is one of the most fundamental concepts in epidemiology and its importance is because it provides a measure of disease burden in a population at a particular point in time. It can also be part of a compendium of measures used to inform public health prevention policies to help slow the spread of disease through the population. To provide prevalence estimates that are reliable and generalizable, the sample must be comprehensive enough to capture all relevant subpopulations in the general population and as mentioned, for a number of diseases this can be challenging because many of these sub-populations can be hard-to-reach. Thus, sampling bias

corrections are needed. Interestingly, this article has presented new methodology where biased samples result due to over-sampling of symptomatic individuals. Such biased samples are here shown to be inconsistent in terms of not converging to the true proportion of infected individuals. In addition, Algorithms 1–4 go further and present corrections both for sampling bias and testing errors. Such corrections either eliminate bias completely (Algorithms 1 and 3), or reduce it substantially (Algorithms 2 and 4) when testing error rates are known or can be estimated consistently. However, the methodology generalizes easily regardless of how the biased samples resulted.

A limitation of our study is that we do not estimate error rates directly from our sample, but take the estimator from a previous independent sample. If this is not the case, then at least under the random sampling situation, prevalence can still be estimated using a Bayesian approach described by Diggle.²² This naturally results in increased variability of the prevalence estimate and relies on a reasonable prior distribution being elicited for the prevalence.

Sample pooling has also been proposed as an efficient way to estimate population prevalence because if the disease prevalence is low, then little information is accrued from individual tests.²⁴ This is sometimes called group testing. However, this implicitly assumes random sampling of pools which is clearly not the case considered here.

Another approach is to use population seroprevalence complex surveys.^{25,26} While inherently much more difficult to conduct and analyze, these can also suffer from non-ignorable non-response which can lead to biased estimates of prevalence. Indeed, biased sampling can be more generally cast within a missing data framework and the impact of different missing data mechanisms has been studied.²⁷

For some diseases it is becoming more common to use administrative data to estimate disease prevalence since for many countries these data cover large proportions of the population. Examples include Canada, Denmark, and Italy among others. This requires some effort to properly assemble these data sources,²⁸ but they have to date not proven as useful for emerging diseases like COVID-19 where surveillance studies dominated the earlier days of the pandemic.

In First-World countries, particularly in urban areas, testing practices seem to be well-described by oversampling of symptomatic individuals, sometimes even testing the whole group in a subpopulation, as it is the case with COVID-19 testing in universities and industries. A possible extension, however, is to consider the opposite situation in which the symptomatic group is under-sampled, producing an estimator that is biased because it underestimates prevalence. Such scenario is certainly relevant for COVID-19 too in several Third World countries, and even in difficult-to-reach subpopulations of First World countries.

AUTHOR CONTRIBUTIONS

Daniel Andrés Díaz-Pachón and J. Sunil Rao conceptualized the methodology framework and the paper. Lili Zhou, Daniel Andrés Díaz-Pachón, and Ola Hössjer developed the methodology details. Lili Zhou ran the examples and produced the R code. J. Sunil Rao and Chen Zhao ran the simulations. Daniel Andrés Díaz-Pachón, J. Sunil Rao, and Ola Hössjer reviewed and proofread the paper.

ACKNOWLEDGEMENTS

The authors thank the comments from the two anonymous reviewers that greatly improved the quality and readability of this manuscript.

FUNDING INFORMATION

Daniel Andrés Díaz-Pachón, Chen Zhao, and J. Sunil Rao acknowledge the support of the Copeland Foundation Award 2022 from the Department of Public Health Sciences at the University of Miami.

CONFLICT OF INTEREST STATEMENT

The authors declare no potential conflict of interests.

DATA AVAILABILITY STATEMENT

Code to implement Algorithms 1–4 is available at https://github.com/kalilizhou/BiasCorrection.git. The data used in Section 5 is publicly available at https://github.com/nshomron/COVIDpred.

ORCID

Daniel Andrés Díaz-Pachón https://orcid.org/0000-0001-6281-1720

REFERENCES

- 1. Tan S, Makela S, Heller D, et al. A Bayesian evidence synthesis approach to estimate disease prevalence in hard-to-reach populations: hepatitis C in new York City. *Epidemics*. 2018;23:96-109. doi:10.1016/j.epidem.2018.01.002
- 2. Alleva G, Arbia G, Falorsi PD, Zuliani A. A sample approach to the estimation of the critical parameters of the SARS-CoV-2 epidemics: an operational design with a focus on the Italian Health System. Technical Report. Rome: University of Sapienza; 2020.
- 3. Hellewell J, Abbott S, Gimma A, et al. Feasibility of controlling COVID-19 outbreaks by isolation of cases and contacts. *Lancet Glob Health*. 2020;8(4):e488-e496. doi:10.1016/S2214-109X(20)30074-7
- 4. Mancastroppa M, Castellano C, Vezzani A, Burioni R. Stochastic sampling effects favor manual over digital contact tracing. *Nat Commun.* 2021;12:1919. doi:10.1038/s41467-021-22082-7
- 5. Bengio Y, Janda R, Yu YW, et al. The need for privacy with public digital contact tracing during the COVID-19 pandemic. *Lancet Digit Health*. 2020;2(7):e342-e344. doi:10.1016/S2589-7500(20)30133-3
- 6. Andrews I, Kasy M. Identification of and correction for publication bias. Am Econ Rev. 2019;109(8):2766-2794. doi:10.1257/aer.20180310
- Kundu R, Shi X, Morrison J, Mukherjee B. A framework for understanding selection bias in real-world healthcare data; 2023. doi:10.48550/arXiv.2304.04652
- 8. Cortes C, Mohri M, Riley M, Rostamizadeh A. Sample selection bias correction theory. In: Freund Y, Györfi L, Turán G, Zeugmann T, eds. *Algorithmic Learning Theory*. Heidelberg: Springer; 2008:38-53.
- 9. Liu AC, Scholtus S, Waal TD. Correcting selection bias in big data by pseudo-weighting. *J Surv Stat Methodol*. 2022;smac029. doi:10.1093/jssam/smac029
- 10. Díaz-Pachón DA, Rao JS. A simple correction for COVID-19 sampling bias. J Theor Biol. 2021;512:110556. doi:10.1016/j.jtbi.2020.110556
- 11. Barbier J. High-dimensional inference: a statistical mechanics perspective. Ithaca Viaggio Nella Sci. 2020;XVI:99-137.
- 12. Hössjer O, Díaz-Pachón DA, Rao JS. A formal framework for knowledge acquisition: going beyond machine learning. *Entropy*. 2022;24(10):1469. doi:10.3390/e24101469
- 13. Shannon CE. A mathematical theory of communication. Bell Syst Tech J. 1948;27(3):379-423.
- Díaz-Pachón DA, Hössjer O, Marks RJ II. Is cosmological tuning fine or coarse? J Cosmol Astropart Phys. 2021;2021(07):020. doi:10.1088/1475-7516/2021/07/020
- 15. Díaz-Pachón DA, Hössjer O, Marks RJ II. Sometimes size does not matter. Found Phys. 2023;53:1. doi:10.1007/s10701-022-00650-1
- 16. Jaynes ET. Information theory and statistical mechanics. Phys Rev. 1957;106(4):620-630. doi:10.1103/PhysRev.106.620
- 17. Díaz-Pachón DA, Marks RJ II. Generalized active information: extensions to unbounded domains. *BIO-Complex*. 2020;2020(3):1-6. doi:10.5048/BIO-C.2020.3
- Zoabi Y, Deri-Rozov S, Shomron N. Machine learning-based prediction of COVID-19 diagnosis based on symptoms. NPJ Digit Med. 2021;4:3. doi:10.1038/s41746-020-00372-6
- 19. Scientific Advisory Group for Emergencies, Government Office for Science. GOS: impact of false positives and negatives, 3 June 2020. Technical Report. London: Government Office for Science; 2020.
- 20. Arevalo-Rodríguez I, Buitrago-García D, Simancas-Racines D, et al. False-negative results of initial RT-PCR assays for COVID-19: a systematic review. *PLoS One.* 2020;15(12):e0242958. doi:10.1371/journal.pone.0242958
- 21. Haug S, Marks RJ II, Dembski WA. Exponential contingency explosion: implications for artificial general intelligence. *IEEE Trans Syst Man Cybern Syst.* 2022;52(5):2800-2808. doi:10.1109/TSMC.2021.3056669
- 22. Diggle PJ. Estimating prevalence using an imperfect test. Epidemiol Res Int. 2011;2011:608719. doi:10.1155/2011/608719
- 23. Rogan WJ, Gladen B. Estimating prevalence from the results of a screening test. *Am J Epidemiol*. 1978;107(1):71-76. doi:10.1093/oxfordjournals.aje.a112510
- 24. Brynildsrud O. COVID-19 prevalence estimation by random sampling in population: optimal sample pooling under varying assumptions about true prevalence. *BMC Med Res Methodol.* 2020;20:196. doi:10.1186/s12874-020-01081-0
- 25. Carabaña JM. Datos de encuesta Para estimar la prevalencia de COVID-19. Un estudio piloto en Madrid capital. *Rev Esp Salud Publica*. 2020;94:e202011159.
- 26. Franceschi VB, Santos AS, Glaeser AB, et al. Population-based prevalence surveys during the COVID-19 pandemic: a systematic review. *Rev Med Virol.* 2021;31(4):e2200. doi:10.1002/rmv.2200
- 27. Hössjer O, Díaz-Pachón DA, Chen Z, Rao JS. An information theoretic approach to prevalence estimation and missing data; 2023. doi:10.48550/arXiv.2206.05120
- 28. Ward MM. Estimating disease prevalence and incidence using administrative data: some assembly required. *J Rheumatol.* 2013;40(8):1241-1243. doi:10.3899/jrheum.130675
- 29. Dembski WA, Marks RJ II. Bernoulli's principle of insufficient reason and conservation of information in computer search. *Proceedings of the 2009 IEEE International Conference on Systems, Man, and Cybernetics.* New York: IEEE; 2009:2647-2652. doi:10.1109/ICSMC.2009.5346119
- 30. Dembski WA, Marks RJ II. Conservation of information in search: measuring the cost of success. *IEEE Trans Syst Man Cybern A Syst Hum*. 2009;5(5):1051-1061. doi:10.1109/TSMCA.2009.2025027

0970233, 26, Downloaded from https://onlinelibrary.wiely.com/doi/10.1002/sim.9885, Wiley Online Library on [0907/2024]. See the Terms and Conditions (https://onlinelibrary.wiely.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons License

- 32. Montañez GD. The famine of forte: few search problems greatly favor your algorithm. 2017 IEEE International Conference on Systems, Man, and Cybernetics (SMC). New York: IEEE; 2017:477-482. doi:10.1109/SMC.2017.8122651
- 33. Montañez GD. A unified model of complex specified information. BIO-Complex. 2018;2018(4):1-26. doi:10.5048/BIO-C.2018.4
- 34. Wolpert DH, MacReady WG. No free lunch theorems for optimization. IEEE Trans Evol Comput. 1997;1(1):67-82. doi:10.1109/4235.585893
- 35. Díaz-Pachón DA, Sáenz JP, Rao JS, Dazard JE. Mode hunting through active information. *Appl Stoch Models Bus Ind.* 2019;35(2):376-393. doi:10.1002/asmb.2430
- 36. Liu T, Díaz-Pachón DA, Rao JS, Dazard JE. High dimensional mode hunting using pettiest component analysis. *IEEE Trans Pattern Anal Mach Intell*. 2023;45(4):4637-4649. doi:10.1109/TPAMI.2022.3195462
- 37. Díaz-Pachón DA, Marks RJ II. Active information requirements for fixation on the Wright-Fisher model of population genetics. BIO-Complex. 2020;2020(4):1-6. doi:10.5048/BIO-C.2020.4
- 38. Cover TM, Thomas JA. Elements of Information Theory. 2nd ed. New York: Wiley; 2006.
- 39. Díaz-Pachón DA, Sáenz JP, Rao JS. Hypothesis testing with active information. Stat Probab Lett. 2020;161:108742. doi:10.1016/j.spl.2020.108742
- 40. Díaz-Pachón DA, Hössjer O. Assessing, testing and estimating the amount of fine-tuning by means of active information. *Entropy*. 2022;24(10):1323. doi:10.3390/e24101323

How to cite this article: Zhou L, Díaz-Pachón DA, Zhao C, Rao JS, Hössjer O. Correcting prevalence estimation for biased sampling with testing errors. *Statistics in Medicine*. 2023;42(26):4713-4737. doi: 10.1002/sim.9885

APPENDIX A. GENERAL PROOFS

Proof of Proposition 1.

$$\hat{\mathbf{p}}_{T}^{s,i} = f_{S^{*}|T} \left(I_{s}^{(i)} | T = 1 \right) \\
= \frac{P \left[T = 1 | S^{*} = I_{s}^{(i)} \right]}{P[T = 1]} f_{S^{*}} \left(I_{s}^{(i)} \right) \\
= \frac{q \left(I_{s}^{(i)} \right)}{P[T = 1]} f_{S^{*}} \left(I_{s}^{(i)} \right) \\
= \frac{N_{T}^{s,i} / N_{s}^{(i)}}{N_{T} / N} \frac{N_{s}^{(i)}}{N} \\
= \frac{N_{T}^{s,i}}{N_{T}}.$$
(A1)

Remember that terms without subindex T are here population values, whereas terms with the subindex T are sample values. Notice in the fourth and fifth steps that $N_s^{(i)}$ cancels. Therefore, all the remaining information about $I_s^{(i)}$ comes from $N_T^{s,i}$, the number of tested (sampled) individuals with symptoms s and infectious status i. Now, from the third equality, this value is seen to come from $q(I_s^{(i)})$, the sampling probability of the group $I_s^{(i)}$. Therefore, all knowledge of $I_s^{(i)}$ comes from whatever knowledge we have about the sample mechanism $q(I_s^{(i)})$.

Proof of Theorem 2.

$$\begin{split} q^*(I_0) < q^*(I_1) &\Leftrightarrow \frac{N_T^{0,*}}{N_0} < \frac{N_T^{1,*}}{N_1} \\ &\Leftrightarrow \frac{N_1/N}{N_0/N} < \frac{N_T^{1,*}/N_T}{N_T^{0,*}/N_T} \end{split}$$

$$\Leftrightarrow \hat{\mathbf{p}}_T^{1,*} > p_1,$$

where the fourth step used that $p_0 + p_1 = 1 = \hat{\mathbf{p}}_T^{0,*} + \hat{\mathbf{p}}_T^{1,*}$.

Proof of Proposition 2. This proposition follows directly from the definition of testing errors. Consider for instance the first equation of (29). It stipulates that whereas the proportion of sampled individuals with s=i=0 is $\hat{\mathbf{p}}_T^{0,0}$, the reported fraction $\tilde{\mathbf{p}}_T^{0,0}$ of sampled individuals with s=i=0 differs from $\hat{\mathbf{p}}_T^{0,0}$ by an amount $\check{\beta}_0\hat{\mathbf{p}}_T^{0,1}$ – $\check{\alpha}_0\hat{\mathbf{p}}_T^{0,0}$, where the first term $\check{\beta}_0\hat{\mathbf{p}}_T^{0,1}$ is the fraction of (0, 1)-individuals wrongly classified as (0, 0), whereas the second term $\check{\alpha}_0\hat{\mathbf{p}}_T^{0,0}$ is the fraction of (0, 0)-individuals wrongly classified as (0, 1). The other three equations of (29) are motivated analogously.

Proof of Proposition 3. It follows from Equations (29) and (31) that $\bar{p}_T^{s,*} = \hat{\mathbf{p}}_T^{s,*} = \hat{\mathbf{p}}_T^{s,*} = \hat{\mathbf{p}}_T^{s,*}$ for s = 0, 1. This, and another application of (29) and (31) gives

$$\frac{\overline{p}_T^{s,0}}{\overline{p}_T^{s,*}} = \frac{\tilde{\mathbf{p}}_T^{s,0}/\hat{\mathbf{p}}_T^{s,*} - \hat{\beta}_s}{1 - \hat{\alpha}_s - \hat{\beta}_s} = \frac{\hat{\mathbf{p}}_T^{s,0}}{\hat{\mathbf{p}}_T^{s,*}} \cdot \frac{1 - \check{\alpha}_s - \check{\beta}_s}{1 - \hat{\alpha}_s - \hat{\beta}_s} + \frac{\check{\beta}_s - \hat{\beta}_s}{1 - \hat{\alpha}_s - \hat{\beta}_s}.$$
(A2)

Since $q(I_s^{(0)}) = q(I_s^{(1)})$ by (14), it follows that $\hat{\mathbf{p}}_T^{s,0}/\hat{\mathbf{p}}_T^{s,*}$ is a consistent estimator of $p_s^{(0)}/p_s$, and by assumption, $\hat{\alpha}_s$ and $\hat{\beta}_s$ are consistent estimators of α_s and β_s respectively. Moreover, Lemma 2 below implies that $\check{\alpha}_s$ and $\check{\beta}_s$ are consistent estimators of α_s and β_s as well. From this and (A2) it follows that $\overline{p}_T^{s,0}/\overline{p}_T^{s,*}$ is a consistent estimator of $p_s^{(0)}/p_s$. The fact that $\overline{p}_T^{s,1}/\overline{p}_T^{s,*}$ is a consistent estimator of $p_s^{(1)}/p_s$ is proved in the same way.

APPENDIX B. ASYMPTOTICS

As a preparation, we prove the following lemma that will be used as assumption in the main result of this section:

Lemma 1.
$$q(I_s^{(0)}) = q(I_s^{(1)}) = q(I_s)$$
.

Proof. The first equality is obtained by assumption (14). In order to prove the second equality, we use that a randomly chosen individual from I_s belongs to $I_s^{(i)}$ with probability $p_s^{(i)}/p_s$ for i=0,1. Conditioning on which subcohort of I_s the individual belongs to, it follows that

$$q(I_s) = \frac{p_s^{(0)}}{p_s} q(I_s^{(0)}) + \frac{p_s^{(1)}}{p_s} q(I_s^{(1)}).$$

Therefore, since $q(I_s)$ is a weighted average of $q(I_s^{(0)})$ and $q(I_s^{(1)})$, the second equality of Lemma 1 follows from the first one.

Hössjer et al (2023) proved a couple of theorems that we will use to prove the asymptotic results for the estimators discussed in this article.²⁷ We present them here for completeness, fitting their notation to ours.

Theorem 5 (Theorem 1 of Hössjer et al (2023)). Let $N \to \infty$ in such a way that N_1/N is always fixed, that Lemma 1 holds for fixed $q(I_0)$ and $q(I_1)$, and that there exists $\rho_s \in [0,1]$ such that $\hat{p}_s \xrightarrow{p} \rho_s$, as $N \to \infty$, for s = 0, 1. Then

$$N^{1/2} \left(\hat{\mathbf{p}}_{T}^{*,1} - \frac{p_{0}^{(1)}q(I_{0}) + p_{1}^{(1)}q(I_{1})}{q} \right) \xrightarrow{\mathcal{L}} \mathcal{N}(0, V_{1} + V_{2}),$$

$$N^{1/2} \left(\hat{p}^{(1)} - \left(\frac{\rho_{0}p_{0}^{(1)}}{p_{0}} + \frac{\rho_{1}p_{1}^{(1)}}{p_{1}} \right) \right) \xrightarrow{\mathcal{L}} \mathcal{N}(0, V_{3} + V_{4}),$$

$$N^{1/2} (\hat{p}_{s} - \rho_{s}) \xrightarrow{\mathcal{L}} \mathcal{N}(0, C_{00})$$
(B1)

as $N \to \infty$, where q is defined in (8),

$$\begin{split} V_1 &= \sum_{s=0}^1 \frac{p_s q(I_s) \left[1 - q(I_s)\right] \frac{p_s^{(1)}}{p_s} \left(1 - \frac{p_s^{(1)}}{p_s}\right)}{(p_0 q(I_0) + p_1 q(I_1))^2}, \\ V_2 &= \sum_{s=0}^1 \frac{p_s q(I_s) \left[1 - q(I_s)\right] \left(\frac{p_s^{(1)}}{p_s} - \frac{p_0^{(1)} q(I_0) + p_1^{(1)} q(I_1)}{p_0 q(I_0) + p_1 q(I_1)}\right)^2}{(p_0 q(I_0) + p_1 q(I_1))^2} \\ V_3 &= \sum_{s=0}^1 \rho_s^2 p_s^{-1} \frac{1 - q(I_s)}{q(I_s)} \frac{p_s^{(1)}}{p_s} \left(1 - \frac{p_s^{(1)}}{p_s}\right), \\ V_4 &= C_{00} \left(\frac{p_0^{(1)}}{p_0} - \frac{p_1^{(1)}}{p_1}\right)^2, \end{split}$$

and

$$C_{00} = \frac{(1+q)^2 B_{11} + \left(\frac{p_1 q(I_1)}{q}\right)^2 \left[p_0 q(I_0)(1-q(I_0)) + p_1 q(I_1)(1-q(I_1))\right] + 2(1+q) \frac{p_1 q(I_1)}{q} \Sigma_{\pi p 1}}{4},$$

$$B_{11} = \frac{p_1 q(I_1) \left\{q^2 \left[1-q(I_1)\right] - 2q p_1 q(I_1) \left[1-q(I_1)\right] + p_1 q(I_1) p_0 q(I_0) \left[1-q(I_0)\right] + p_1^2 q^2 (I_1) \left[1-q(I_1)\right]\right\}}{q^4},$$

$$\Sigma_{\pi p 1} = \frac{p_1 q(I_1) \left[1-q(I_1)\right]}{q} - \frac{p_1 q(I_1) \left\{p_0 q(I_0) \left[1-q(I_0)\right] + p_1 q(I_1) \left[1-q(I_1)\right]\right\}}{q^2}.$$

Proof of Theorem 1. Since $\hat{p}_1 = p_1$ and $\hat{p}_0 = p_0 = 1 - p_1$ are known, convergence in probability $\hat{p}_s \xrightarrow{p} \rho_s$ trivially holds with $\rho_s = p_s$, whereas the asymptotic weak limit of $N^{1/2}(\hat{p}_s - p_s)$ in (B1) degenerates to a one point distribution at 0 ($C_{00} = 0$). Since $q(I_1) = 1$, V_1 , V_2 , and V_3 , in Theorem 5 are readily simplified to

$$V_{01} = \frac{p_0 q(I_0) \left[1 - q(I_0)\right] \frac{p_0^{(1)}}{p_0} \left(1 - \frac{p_0^{(1)}}{p_0}\right)}{(p_0 q(I_0) + p_1)^2},$$
(B2)

$$V_{02} = \frac{p_0 q(I_0) \left[1 - q(I_0)\right] \left(\frac{p_0^{(1)}}{p_0} - \frac{p_0^{(1)} q(I_0) + p_1^{(1)}}{p_0 q(I_0) + p_1}\right)^2}{\left(p_0 q(I_0) + p_1\right)^2},$$
(B3)

$$V_{03} = p_0 \frac{1 - q(I_0)}{q(I_0)} \frac{p_0^{(1)}}{p_0} \left(1 - \frac{p_0^{(1)}}{p_0} \right),$$
(B4)

whereas $V_4 = 0$ since $C_{00} = 0$.

0970258, 2023, 26, Downloaded from https://onlinelibrary.wiley.com/doi/10.1002/sim.9885, Wiley Online Library on [09/07/2024]. See the Terms and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Licesease and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Licesease and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Licesease and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Licesease and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Licesease and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Licesease and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons and the conditions of the con

0970238, 2023, 2,6, Downloaded from https://onlinelibrary.wiely.com/doi/10.1002/sim 9885, Wiley Online Library on [0907/2024]. See the Terms and Conditions (https://onlinelibrary.wiely.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Licensea

Proof of Theorem 3. It is obtained directly from Theorem 5 by substituting V_1 , V_2 , V_3 , and V_4 for V_{11} , V_{12} , V_{13} , and V_{14} , respectively, where

$$V_{11} = \sum_{s=0}^{1} \frac{p_s q(I_s) \left[1 - q(I_s)\right] \frac{p_s^{(1)}}{p_s} \left(1 - \frac{p_s^{(1)}}{p_s}\right)}{\left(p_0 q(I_0) + p_1 q(I_1)\right)^2},$$
(B5)

$$V_{12} = \sum_{s=0}^{1} \frac{p_s q(I_s) \left[1 - q(I_s)\right] \left(\frac{p_s^{(1)}}{p_s} - \frac{p_0^{(1)} q(I_0) + p_1^{(1)} q(I_1)}{p_0 q(I_0) + p_1 q(I_1)}\right)^2}{(p_0 q(I_0) + p_1 q(I_1))^2},$$
(B6)

$$V_{13} = \sum_{s=0}^{1} \rho_s^2 p_s^{-1} \frac{1 - q(I_s)}{q(I_s)} \frac{p_s^{(1)}}{p_s} \left(1 - \frac{p_s^{(1)}}{p_s} \right), \tag{B7}$$

$$V_{14} = C_{00} \left(\frac{p_0^{(1)}}{p_0} - \frac{p_1^{(1)}}{p_1} \right)^2.$$
 (B8)

B.1 Asymptotics with testing errors

Before proving the asymptotic results with testing errors, some previous assumptions and results are used. First, we assume that the existing estimators of error rates are asymptotically normal. That is, there exists $\overline{\alpha} = (\overline{\alpha}_0, \overline{\alpha}_1)$ and $\overline{\beta} = (\overline{\beta}_0, \overline{\beta}_1)$ such that, for $\hat{\alpha} = (\hat{\alpha}_0, \hat{\alpha}_1)$ and $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)$,

$$N^{1/2} \left(\hat{\boldsymbol{\alpha}} - \overline{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}} - \overline{\boldsymbol{\beta}} \right) \xrightarrow{\mathcal{L}} \mathcal{N} \left(0, \begin{pmatrix} \Omega_{\alpha\alpha} & \Omega_{\alpha\beta} \\ \Omega_{\alpha\beta}^T & \Omega_{\beta\beta} \end{pmatrix} \right)$$
(B9)

as $N \to \infty$, where each of the terms in the variance-covariance matrix is a 2 × 2 matrix, and

$$\Omega_{\alpha\alpha} = (\Omega_{\alpha\alpha rs})_{r,s=0,1}
\Omega_{\beta\beta} = (\Omega_{\beta\beta rs})_{r,s=0,1}
\Omega_{\alpha\beta} = (\Omega_{\alpha\beta rs})_{r,s=0,1}.$$
(B10)

We also assume that $\check{\alpha}_0$, $\check{\alpha}_1$, $\check{\beta}_0$, and $\check{\beta}_1$ are all independent. Moreover, the following lemma will be used:

Lemma 2 (Lemma 2 from Hössjer et al (2023)). *The proportions of false positive and negatives in the sample satisfy*

$$N^{1/2}(\check{\alpha}_s - \alpha_s) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \Sigma_{\alpha\alpha s}),$$
 (B11)

and

$$N^{1/2}(\check{\beta}_s - \beta_s) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \Sigma_{\beta\beta s}),$$
 (B12)

respectively as $N \to \infty$, with

$$\begin{split} \Sigma_{\alpha\alpha s} &= \alpha_s (1 - \alpha_s) q / \left[p_s q(I_s) \left(1 - p_s^{(1)} / p_s \right) \right], \\ \Sigma_{\beta\beta s} &= \beta_s (1 - \beta_s) q p_s / \left[p_s q(I_s) p_s^{(1)} \right]. \end{split}$$

After transforming the notation of Hössjer et al (2023) to our notation, the asymptotic results of Theorem 4 are a direct consequence of the following result from Hössjer et al (2023):

Theorem 6 (Theorem 2 of Hössjer et al (2023)). Suppose the conditions of Theorem 3 hold, and additionally that the estimators $\overline{p}_T^{s,0}/\overline{p}_T^{s,*}$ and $\overline{p}_T^{s,1}/\overline{p}_T^{s,*}$ of prevalences of unaffected and affected in symptom group s converge to

$$\overline{\rho}_s^{(0)} = \frac{\beta_s - \overline{\beta}_s}{1 - \overline{\alpha}_s - \overline{\beta}_s} + \frac{1 - \alpha_s - \beta_s}{(1 - \overline{\alpha}_s - \overline{\beta}_s)} \cdot \frac{p_s^{(0)}}{p_s},\tag{B13}$$

$$\overline{\rho}_s^{(1)} = \frac{\alpha_s - \overline{\alpha}_s}{1 - \overline{\alpha}_s - \overline{\beta}_s} + \frac{1 - \alpha_s - \beta_s}{(1 - \overline{\alpha}_s - \overline{\beta}_s)} \cdot \frac{p_s^{(1)}}{p_s},\tag{B14}$$

respectively as $N \to \infty$. Suppose further that the proportion of individuals wrongly classified in the sample $\check{\alpha}_s$, $\check{\beta}_s$, s = 0, 1 are independent. Then

$$N^{1/2} \left(\overline{p}_T^{*,1} - \frac{p_0 q(I_0) \overline{\rho}_0^{(1)} + p_1 q(I_1) \overline{\rho}_1^{(1)}}{q} \right) \xrightarrow{\mathcal{L}} \mathcal{N}(0, V_5 + V_6), \tag{B15}$$

$$N^{1/2} \left[\hat{p}^{(1)} - \left(\rho_0 \overline{\rho}_0^{(1)} + \rho_1 \overline{\rho}_1^{(1)} \right) \right] \xrightarrow{\mathcal{L}} \mathcal{N}(0, V_7 + V_8)$$
 (B16)

as $N \to \infty$, where $\hat{p}^{(1)}$ is defined as in (35),

$$V_{5} = \sum_{r,s} \frac{p_{r}q(I_{r})p_{s}q(I_{s})}{q^{2}} \bar{A}_{rs},$$

$$V_{6} = \sum_{r,s} \bar{\rho}_{r}^{(1)} \bar{\rho}_{s}^{(1)} B_{rs},$$

$$V_{7} = \sum_{r,s} \rho_{r} \rho_{s} \bar{A}_{rs},$$

$$V_{8} = \sum_{r,s} \bar{\rho}_{r}^{(1)} \bar{\rho}_{s}^{(1)} C_{rs},$$

whereas

$$B_{11} = B_{00} = -B_{01} = -B_{10},$$

 $C_{00} = C_{11} - C_{01} = -C_{10},$

are defined in Theorem 5. Moreover,

$$\bar{A}_{ss} = K_{s1}^2 A_{ss} + K_{s2}^2 \Sigma_{\alpha \alpha s} + K_{s3}^2 \Sigma_{\beta \beta s}
+ K_{s4}^2 \Omega_{\alpha \alpha ss} + K_{s5}^2 \Omega_{\beta \beta ss} + 2K_{s4} K_{s5} \Omega_{\alpha \beta ss},$$
(B17)

and

$$\bar{A}_{rs} = K_{r4}K_{s4}\Omega_{\alpha\alpha rs} + K_{r5}K_{s5}\Omega_{\beta\beta rs} + K_{r4}K_{s5}\Omega_{\alpha\beta rs} + K_{r5}K_{s4}\Omega_{\alpha\beta sr},$$
(B18)

when $r \neq s$, with

$$K_{s1} = (1 - \alpha_s - \beta_s)/K_s,$$

$$K_{s2} = \left(1 - \frac{p_s^{(1)}}{p_s}\right)/K_s,$$

$$K_{s3} = -p_s^{(1)}/(p_sK_s),$$

$$K_{s4} = \left[\alpha_s + \overline{\beta}_s - 1 + \frac{p_s^{(1)}(1 - \alpha_s - \beta_s)}{p_s}\right]/K_s^2,$$

$$K_{s5} = \left[\alpha_s - \overline{\alpha}_s + \frac{p_s^{(1)}(1 - \alpha_s - \beta_s)}{p_s}\right]/K_s^2,$$

$$K_s = 1 - \overline{\alpha}_s - \overline{\beta}_s,$$
(B19)

and finally

$$A_{ss} = \frac{(1 - q(I_s))\frac{p_s^{(1)}}{p_s} \left(1 - \frac{p_s^{(1)}}{p_s}\right)}{p_s q(I_s)}.$$
 (B20)

Notice that, if $\overline{\alpha}_s = \alpha_s$ and $\overline{\beta}_s = \beta_s$ (that is, if the error rates are estimated consistently), then $\overline{p}_T^{s,i}/p_T^{s,*}$ is a consistent estimator of $\overline{\rho}_s^{(i)} = p_s^{(i)}/p_s$ for i = 0, 1. In particular, $\overline{p}_T^{s,1}/p_T^{s,*}$ is a consistent estimator of the prevalence $\overline{\rho}_s^{(1)} = p_s^{(i)}/p_s$ among the individuals of symptom group s.

APPENDIX C. ACTIVE INFORMATION: THE INDEX

Active information (actinfo) was introduced in search problems to quantify the amount of Shannon information introduced by the programmer in a search problem.²⁹⁻³¹ In machine learning, it has been used to show that no algorithm performs well for a large class of problems, in agreement with the so-called No Free Lunch Theorems.³²⁻³⁴ It has also been used for mode hunting,^{35,36} and to compare neutral to non-neutral models in population genetics.³⁷

We now use active information to analyze the bias. Through the eyes of actinfo, the bias is formally seen as the addition (if the parameter is overestimated) or subtraction (if the parameter is underestimated) of relevant information in the estimation of the parameter. Formally, active information is defined as

$$\hat{I}^{+} = \log(\hat{p}^{(1)}/p^{(1)}),\tag{C1}$$

where the logarithm is taken to be in base e, so that information is measured in nats. Thus defined, active information measures the amount of Shannon information of the estimator $\hat{p}^{(1)}$ to the true proportion $p^{(1)}$, and it is the quantity that is averaged in the Kullback-Leibler divergence.³⁸ That is, if the true proportion is overestimated, the active information will be positive and large; if the true proportion is underestimated, the active information will be negative; and if the true proportion is accurately estimated, the active information will be around zero.^{39,40} Because of Theorem 4, we interpret (C1) as an approximation of $I^+ = \log[(\rho_0 \overline{\rho}_0^{(1)} + \rho_1 \overline{\rho}_1^{(1)})/p^{(1)}]$.

APPENDIX D. SIMULATION

This section uses simulation to analyze the asymptotic behavior of the corrected estimator. The population has the following features:

- the proportion of positive cases with symptoms $p_1^{(1)}$ is 15%,
- the proportion of negative cases with symptoms $p_1^{(0)}$ is 5%,
- the proportion of positive cases without symptoms $p_0^{(1)}$ is 5%,
- and proportion of negative cases without symptoms $p_0^{(0)}$ is 75%.

Thus, the prevalence is $p_0^{(1)} + p_1^{(1)} = 20\%$, the proportion of symptomatic individuals in the population is $p_1^{(1)} + p_1^{(0)} = 20\%$, so the proportion of asymptomatic in the population is $p_0^{(1)} + p_0^{(0)} = 80\%$.

D.1 Correction without testing errors

We will run the simulation for multiple proportions of asymptomatic individuals getting sampled and will compare the results using a boxplot. The true prevalence will be known in the simulation, allowing us to evaluate the accuracy of our estimators.

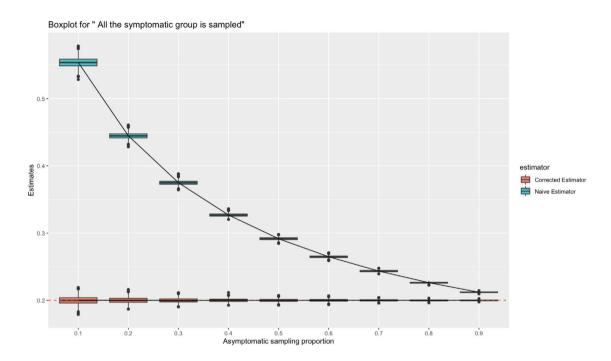


FIGURE D1 All the symptomatic group is sampled.

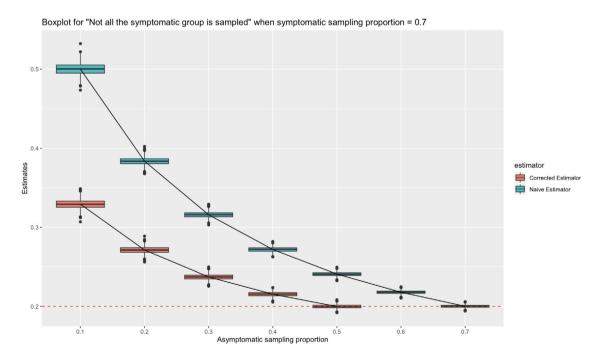


FIGURE D2 Not all the symptomatic group is sampled.

10970258, 2023, 26, Downloaded from https://onlinelibrary.wiley.com/doi/10.1002/sim.9885, Wiley Online Library on [09/07/2024]. See the Terms and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Liceseare (online Library) on [09/07/2024]. See the Terms and Conditions (https://onlinelibrary.wiley.com/terms-and-conditions) on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons Liceseare (online Library) on [09/07/2024].

TABLE D1 Estimated proportion of symptomatic in the population \hat{p}_1 (E(U)).

$\hat{\mathbf{p}}_{T}^{0,*}$	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
E(U)	0.3881	0.3031	0.2538	0.2221	0.1995	0.1829	0.1701	0.1596	0.1513

D.1.1 All the symptomatic group is sampled

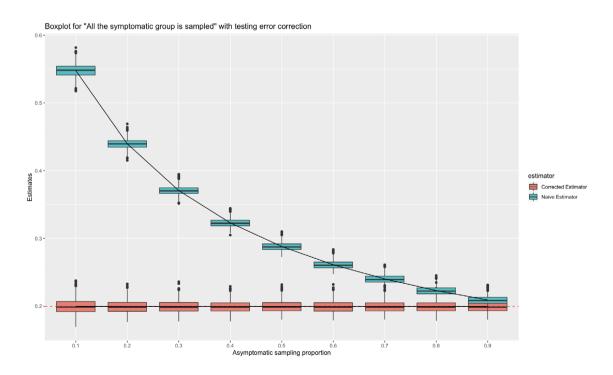
Assuming a population size of 10 000, we initially sampled all symptomatic individuals. However, we increased the sample size by including more asymptomatic individuals as we changed the proportion of those getting tested. This resulted in both the corrected and naive estimators approaching the true value, and the variance of the estimators decreasing, as shown in Figure D1.

D.2 Not all the symptomatic group is sampled

In this scenario, we assume that only 70% of symptomatic individuals underwent testing. From Figure D2, as the number of asymptomatic individuals in the total testing sample size increases, we observe that the corrected estimator converges to the true value faster than the naïve estimator. Thus we see once again that the testing error correction improves the accuracy of prevalence estimation. By Table D1, it can be seen that our estimated symptomatic rate in population \hat{p}_1 decreases as the number of asymptomatic individuals in the sample increases.

APPENDIX E. CORRECTION OF TESTING ERRORS

In the previous section, we explored the simulation of the correction of sampling error for a population without considering testing errors. In this section, we extend our analysis including testing error for asymptomatic and symptomatic individuals separately. Specifically, we will model the false positive and false negative rates for both groups in testing using normal distributions. The false positive rate for asymptomatic individuals $\check{\alpha}_0$ is assumed to follow a normal distribution with mean 0.01 and variance 0.0001, while the false negative rate for asymptomatic individuals \mathring{p}_0 follows a normal distribution with mean 0.2 and variance 0.0001. Similarly, the false positive rate for symptomatic individuals α_1 is assumed to follow a normal distribution with mean 0.03 and variance 0.0001, while the false negative rate for symptomatic individuals $\mathring{\beta}_1$ follows a normal distribution with mean 0.02 and variance 0.0001. The real value of the parameters is assumed to be



All the symptomatic group is sampled with testing error.

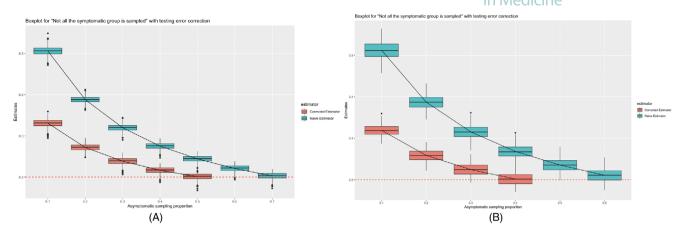


FIGURE E2 Not all the symptomatic group is sampled with testing error. (A) Unbiased correction parameters. (B) Biased correction parameters.

the mean of these distributions. We will consider two scenarios: one in which all symptomatic individuals are sampled, and another in which not all symptomatic individuals are sampled.

E.1 All symptomatic group is sampled

In this simulation, we assume that all symptomatic individuals are sampled for the testing group. We will adjust the proportion of asymptomatic individuals getting sampled from 0.1 to 0.9 to observe the effect of testing error correction on prevalence estimation. Based on the description in the previous section, we use the $\hat{\alpha}_0$, $\hat{\alpha}_1$, $\hat{\beta}_0$, and $\hat{\beta}_1$ obtained from other study as parameters for testing error correction. Here, we assume that $\hat{\alpha}_0$, $\hat{\alpha}_1$, $\hat{\beta}_0$, $\hat{\beta}_1$ follow a uniform distribution with a mean of the true values $\alpha_0 = 0.01$, $\alpha_1 = 0.03$, $\beta_0 = 0.2$, and $\beta_1 = 0.02$ in the simulation study.

From Figure E1, we expect to see that our corrected estimators are very close to the true value, while the naive estimator is approaching the true value as the proportion of asymptomatic individuals increases. Due to the additional variability introduced by testing error, we observe a larger variance of the corrected estimators compared to Algorithm 1.

E.2 Not all the symptomatic group is sampled

We also simulated an scenario where not all symptomatic individuals in the population were sampled for testing, accounting for testing error. Specifically, we assumed that 70% of symptomatic individuals in the population would go for a test.

Here we consider two parameter settings for the correction parameters obtained from other studies: the first one is unbiased, where $\hat{\alpha}_0$, $\hat{\alpha}_1$, $\hat{\beta}_0$, and $\hat{\beta}_1$ follow a uniform distribution with mean 0.01, 0.03, 0.2, and 0.02 (the true values), respectively; the second one is biased, where $\hat{\alpha}_0$, $\hat{\alpha}_1$, $\hat{\beta}_0$, and $\hat{\beta}_1$ follow a uniform distribution with mean 0.05, 0.1, 0.3, and 0.1, respectively. From the results in Figure E2A,B, it can be observed that the estimates from both settings converge to the true values as the proportion of asymptomatic increases in the sample, but the corrected estimate from the second setting has a larger variance.