

A Bayesian Approach to Account for Misclassification in Prevalence and
Trend Estimation

Supplemental Appendix

Martijn van Hasselt, Christopher R. Bollinger, and Jeremy Bray

A. DERIVATION OF IDENTIFIED SETS

A.1. Bounds for Cases I-III

Recall that $p_t = 0$ for all t . Here we present details on the derivation of the prevalence and trend bounds for Case I-III.

Case I. We assume first that $p_t = 0$ for all t , and $q_t = q^*$ is constant over time. From equation (1) it follows that $\pi_t = \mu_t/(1 - q^*)$. Since $\pi_t \leq 1$, the rate of false negative reporting satisfies $0 \leq q^* \leq 1 - \mu_t$ for $t = 1, \dots, T$. Defining $M = \max_s \mu_s$, we get $0 \leq q^* \leq 1 - M$. Together with the expression for π_t , this yields the bounds in (2). While the prevalence bounds can be used to calculate bounds on $\Delta\pi_{t,j}$, we find sharper bounds from the fact that $\Delta\pi_{t,j} = \Delta\mu_{t,j}/(1 - q^*)$. The sign of the derivative of the true trend with respect to q^* depends on the direction of the observed trend. If $\Delta\mu_{t,j} \geq 0$, the lower bound on $\Delta\pi_{t,j}$ is attained when $q^* = 0$, whereas the upper bound is attained when $q^* = 1 - M$. The situation is reversed when $\Delta\mu_{t,j} < 0$. Substituting these values of q^* into the expression for $\Delta\pi_{t,j}$ yields the bounds in (3).

Case II. In this case, $p_t = 0$ for all t , and $q_1 \leq q_2 \leq \dots \leq q_T$. As before, $0 \leq q_t \leq 1 - \mu_t$. Combined with the fact that q_t is non-decreasing, it follows that $0 \leq q_t \leq 1 - M_t^+$, where $M_t^+ = \max_{s \geq t} \mu_s$. This yields the bounds in (4). To derive the trend bounds, note that with time-varying q_t we have

$$\Delta\pi_{t,j} = \frac{\mu_{t+j}}{1 - q_{t+j}} - \frac{\mu_t}{1 - q_t}. \quad (18)$$

From (18) it is immediate that an upper bound for $\Delta\pi_{t,j}$ is obtained when $q_{t+j} = 1 - M_{t+j}^+$ and $q_t = 0$. This yields the upper bounds in (5). The lower bound for $\Delta\pi_{t,j}$ is attained when q_{t+j} is minimal and q_t is maximal, subject to the restriction that q_t is non-decreasing in t . Thus, at the lower bound we must have $q_t = q_{t+j} = q^*$ and

$$\Delta\pi_{t,j} = \Delta\mu_{t,j}/(1 - q^*). \quad (19)$$

Minimizing this depends on the sign of $\Delta\mu_{t,j}$: when the observed trend is non-negative, the lower bound is attained for $q^* = 0$. Substituting this into (19) yields the first lower bound in (5). When

$\Delta\mu_{t,j} < 0$, the lower bound is attained when q^* is maximal, subject to the restrictions

$$q^* = q_t \leq 1 - M_t^+, \quad q^* = q_{t+j} \leq 1 - M_{t+j}^+.$$

Since, by definition, M_t^+ is non-increasing in t , we find the lower bound at $q^* = 1 - M_t^+$. Substitution into (19) yields the second lower bound of (5).

Case III. In this case, $p_t = 0$ for all t , and $q_1 \geq q_2 \geq \dots \geq q_T$. Combined with $q_t \leq 1 - \mu_t$ for all t , we find $0 \leq q_t \leq 1 - M_t^-$, where $M_t^- = \max_{s \leq t} \mu_s$. This yields the prevalence bounds in (6). From (18), the lower bound on the trend is attained at $q_{t+j} = 0$ and $q_t = 1 - M_t^-$. This yields the lower bounds in (7). The upper bound is attained when q_{t+j} is maximal and q_t is minimal, subject to the restriction $q_t \geq q_{t+j}$. Thus, at the upper bound $q_t = q_{t+j} = q^*$ and the true trend is given by (19). Similar to the previous case, the sign of the derivative of $\Delta\pi_{t,j}$ depends on $\Delta\mu_{t,j}$. If the observed trend is non-negative, the derivative is non-negative and the upper bound for the true trend is attained when q^* is maximal, subject to

$$q^* = q_t \leq 1 - M_t^-, \quad q^* = q_{t+j} \leq 1 - M_{t+j}^-.$$

Since, by definition, M_t^- is non-decreasing in t , we find the upper bound at $q^* = 1 - M_{t+j}^-$. Substitution into (19) yields the first upper bound in (7). Alternatively, when $\Delta\mu_{t,j} < 0$, the derivative of $\Delta\pi_{t,j}$ with respect to q^* is negative and the upper bound for the true trend is attained at $q^* = 0$. Substituting into (19) yields the second upper bound in (7).

A.2. Bounds for Case IV

Lower Bound for the Trend

Recall that $a = 1 - x$ and $b = 1 + x$. From equation (1) it follows that, for a given value of q , the lower bound on $\Delta\pi_{t,j}$ is given by

$$\Delta\pi_{t,j}^L = \frac{\mu_{t+j}}{1 - aq} - \frac{\mu_t}{1 - bq}.$$

Minimizing this expression over q yields the (unconditional) lower bound on the trend. We have

$$\begin{aligned}\frac{d(\Delta\pi_{t,j}^L)}{dq} &= \frac{a\mu_{t+j}}{(1-aq)^2} - \frac{b\mu_t}{(1-bq)^2} \\ &\leq \frac{a\mu_{t+j}}{(1-aq)^2} - \frac{a\mu_t}{(1-bq)^2} \\ &\leq \frac{a\mu_{t+j}}{(1-aq)^2} - \frac{a\mu_t}{(1-aq)^2} \\ &\leq \frac{a\Delta\mu_t}{(1-aq)^2}.\end{aligned}$$

If $\Delta\mu_{t,j} < 0$, the derivative is negative, so that $\Delta\pi_{t,j}^L$ is minimized at $q = (1 - M)/b$:

$$\Delta\pi_{t,j}^L = \frac{\mu_{t+j}}{1 - (a/b)(1 - M)} - \frac{\mu_t}{M}. \quad (20)$$

Suppose now that $\Delta\mu_{t,j} \geq 0$. Define the function

$$f(q) = abq^2(b\mu_{t+j} - a\mu_t) - 2qab\Delta\mu_{t,j} + a\mu_{t+j} - b\mu_t.$$

Simple algebra shows that

$$\begin{aligned}\frac{d(\Delta\pi_{t,j}^L)}{dq} = 0 &\Leftrightarrow f(q) = 0, \\ \frac{d(\Delta\pi_{t,j}^L)}{dq} > 0 &\Leftrightarrow f(q) > 0, \\ \frac{d(\Delta\pi_{t,j}^L)}{dq} < 0 &\Leftrightarrow f(q) < 0.\end{aligned}$$

Note also that $\Delta\mu_{t,j} \geq 0$ implies that $b\mu_{t+j} - a\mu_t > 0$ and the function $f(q)$ has a minimum. We first consider the roots of $f(q)$, say q_1 and q_2 :

$$q_1 = \frac{\sqrt{a\mu_{t+j}} - \sqrt{b\mu_t}}{b\sqrt{a\mu_{t+j}} - a\sqrt{b\mu_t}}, \quad q_2 = \frac{\sqrt{a\mu_{t+j}} + \sqrt{b\mu_t}}{b\sqrt{a\mu_{t+j}} + a\sqrt{b\mu_t}}.$$

It is easy to show that $bq_2 > 1$, so that $q_2 > (1 - M)/b$. The root q_1 may be positive or negative, depending on the sign of its numerator.⁶

We now consider two cases. Suppose first that $a\mu_{t+j} < b\mu_t$. Then $q_1 < 0$ and $d(\Delta\pi_{t,j}^L)/dq < 0$

⁶Since $b\sqrt{a\mu_{t+j}} - a\sqrt{b\mu_t} = \sqrt{b}\sqrt{ab\mu_{t+j}} - \sqrt{a}\sqrt{ab\mu_t} \geq 0$ when $\Delta\mu_{t,j} \geq 0$.

for $q \in [0, (1 - M)/b]$. It follows that $\Delta\pi_{t,j}^L$ is minimized at $q = (1 - M)/b$, leading to the lower trend bound in (20). Conversely, suppose that $a\mu_{t+j} \geq b\mu_t$, so that $q_1 \geq 0$. We then further consider two subcases:

- (i) $q_1 > (1 - M)/b$. In this case, $d(\Delta\pi_{t,j}^L)/dq > 0$ for $q \in [0, (1 - M)/b]$ and the lower bound is minimized at $q = 0$:

$$\Delta\pi_{t,j}^L = \Delta\mu_{t,j}.$$

- (ii) $q_1 \leq (1 - M)/b$. In this case, the lower bound is attained at either $q = 0$ or $q = (1 - M)/b$:

$$\Delta\pi_{t,j}^L = \min \left\{ \Delta\mu_{t,j}, \frac{\mu_{t+j}}{1 - (a/b)(1 - M)} - \frac{\mu_t}{M} \right\}.$$

We can now summarize the lower trend bound as follows.

$$\Delta\pi_{t,j}^L = \begin{cases} \frac{\mu_{t+j}}{1 - (a/b)(1 - M)} - \frac{\mu_t}{M} & \text{if } a\mu_{t+j} < b\mu_t, \\ \Delta\mu_{t,j} & \text{if } a\mu_{t+j} \geq b\mu_t, \quad M > 1 - bq_1, \\ \min \left\{ \Delta\mu_{t,j}, \frac{\mu_{t+j}}{1 - (a/b)(1 - M)} - \frac{\mu_t}{M} \right\} & \text{if } a\mu_{t+j} \geq b\mu_t, \quad M \leq 1 - bq_1. \end{cases}$$

Substituting the definition of q_1 into the inequalities for M yields (9).

Upper Bound for the Trend

From equation (1) it follows that, for a given value of q , the upper bound on $\Delta\pi_{t,j}$ is given by

$$\Delta\pi_{t,j}^U = \frac{\mu_{t+j}}{1 - bq} - \frac{\mu_t}{1 - aq}.$$

Maximizing this expression over q yields the (unconditional) upper bound on the trend. We have

$$\begin{aligned} \frac{d(\Delta\pi_{t,j}^U)}{dq} &= \frac{b\mu_{t+j}}{(1 - bq)^2} - \frac{a\mu_t}{(1 - aq)^2} \\ &\geq \frac{b\mu_{t+j}}{(1 - bq)^2} - \frac{a\mu_t}{(1 - bq)^2} \\ &\geq \frac{b\mu_{t+j}}{(1 - bq)^2} - \frac{b\mu_t}{(1 - bq)^2} \\ &\geq \frac{b\Delta\mu_{t,j}}{(1 - bq)^2}. \end{aligned}$$

If $\Delta\mu_{t,j} > 0$, the derivative is positive and $\Delta\pi_{t,j}^U$ is maximized at $q = (1 - M)/b$:

$$\Delta\pi_{t,j}^U = \frac{\mu_{t+j}}{M} - \frac{\mu_t}{1 - (a/b)(1 - M)}. \quad (21)$$

Suppose now that $\Delta\mu_{t,j} \leq 0$. Define the function

$$h(q) = abq^2(a\mu_{t+j} - b\mu_t) - 2qab\Delta\mu_{t,j} + b\mu_{t+j} - a\mu_t.$$

Simple algebra shows that

$$\begin{aligned} \frac{d(\Delta\pi_{t,j}^U)}{dq} = 0 &\Leftrightarrow h(q) = 0, \\ \frac{d(\Delta\pi_{t,j}^U)}{dq} > 0 &\Leftrightarrow h(q) > 0, \\ \frac{d(\Delta\pi_{t,j}^U)}{dq} < 0 &\Leftrightarrow h(q) < 0. \end{aligned}$$

Note also that $\Delta\mu_{t,j} \leq 0$ implies that $a\mu_{t+j} - b\mu_t < 0$ and the function $h(q)$ has a maximum. We first consider the roots of $h(q)$, say q_3 and q_4 :

$$q_3 = \frac{\sqrt{a\mu_t} - \sqrt{b\mu_{t+j}}}{b\sqrt{a\mu_t} - a\sqrt{b\mu_{t+j}}}, \quad q_4 = \frac{\sqrt{a\mu_t} + \sqrt{b\mu_{t+j}}}{b\sqrt{a\mu_t} + a\sqrt{b\mu_{t+j}}}.$$

It is easy to show that $bq_4 > 1$, so that $q_4 > (1 - M)/b$. The root q_3 may be positive or negative, depending on the sign of its numerator.

We now consider two cases. Suppose first that $a\mu_t < b\mu_{t+j}$. Then $q_3 < 0$ and $d(\Delta\pi_{t,j}^U)/dq > 0$ for all $q \in [0, (1 - M)/b]$. It follows that $\Delta\pi_{t,j}^U$ is maximized at $q = (1 - M)/b$, leading to the upper trend bound in (21). Conversely, suppose that $a\mu_t \geq b\mu_{t+j}$, so that $q_3 \geq 0$. We then consider two subcases:

- (i) $q_3 > (1 - M)/b$. Then $d(\Delta\pi_{t,j}^U)/dq < 0$ for all $q \in [0, (1 - M)/b]$ and the upper bound is maximized at $q = 0$:

$$\Delta\pi_{t,j}^U = \Delta\mu_{t,j}.$$

(ii) $q_3 \leq (1 - M)/b$. In this case, the upper bound is attained at either $q = 0$ or $q = (1 - M)/b$:

$$\Delta\pi_{t,j}^U = \max \left\{ \Delta\mu_{t,j}, \frac{\mu_{t+j}}{M} - \frac{\mu_t}{1 - (a/b)(1 - M)} \right\}.$$

We can now summarize the upper trend bound as follows.

$$\Delta\pi_{t,j}^U = \begin{cases} \frac{\mu_{t+j}}{M} - \frac{\mu_t}{1 - (a/b)(1 - M)} & \text{if } b\mu_{t+j} > a\mu_t, \\ \Delta\mu_{t,j} & \text{if } b\mu_{t+j} \leq a\mu_t, \quad M > 1 - bq_3, \\ \max \left\{ \Delta\mu_{t,j}, \frac{\mu_{t+j}}{M} - \frac{\mu_t}{1 - (a/b)(1 - M)} \right\} & \text{if } b\mu_{t+j} \leq a\mu_t, \quad M \leq 1 - bq_3. \end{cases}$$

Substituting the definition of q_3 into the inequalities for M yields (10).

A.3. Bounds for Case V

Lower Bound for the Trend

For given values of p and \bar{q} , the lower bound for $\Delta\pi_{t,j}$ was given in (12) and repeated here:

$$\Delta\pi_{t,j}^L = \frac{\mu_{t+j} - p}{1 - p - a\bar{q}} - \frac{\mu_t - p}{1 - p - b\bar{q}}. \quad (22)$$

It follows that

$$\frac{d(\Delta\pi_{t,j}^L)}{dq} = \frac{a(\mu_{t+j} - p)}{(1 - p - a\bar{q})^2} - \frac{b(\mu_t - p)}{(1 - p - b\bar{q})^2} \leq \frac{a\Delta\mu_{t,j}}{(1 - p - a\bar{q})^2}.$$

When $\Delta\mu_{t,j} < 0$, the lower bound is minimized at $\bar{q} = (1 - M)/b$ and

$$\Delta\pi_{t,j}^L = \frac{\mu_{t+j} - p}{c - p} - \frac{\mu_t - p}{M - p}, \quad (23)$$

where $c = 1 - (a/b)(1 - M)$. It remains to minimize (23) with respect to p . Since $0 \leq p \leq m$, the minimum occurs either at $p = 0$, $p = m$ or possibly at an interior solution p_L^* , where the derivative

of (23) with respect to p is zero.⁷ Specifically,

$$p_L^* = \frac{c\sqrt{M - \mu_t} - M\sqrt{c - \mu_{t+j}}}{\sqrt{M - \mu_t} - \sqrt{c - \mu_{t+j}}},$$

and it remains to check whether $p_L^* \in [0, m]$. If $p_L^* \leq 0$, the lower trend bound is attained at $p = 0$; if $p_L^* \geq m$, the lower bound is attained at $p = m$, and if $0 < p_L^* < m$, the lower bound is attained at $p = p_L^*$. This leads to the following lower trend bound when $\Delta\mu_{t,j} < 0$:

$$\Delta\pi_{t,j}^L = \begin{cases} \frac{\mu_{t+j}}{c} - \frac{\mu_t}{M} & \text{if } p_L^* \leq 0, \\ \frac{\mu_{t+j} - p_L^*}{c - p_L^*} - \frac{\mu_t - p_L^*}{M - p_L^*} & \text{if } 0 < p_L^* < m, \\ \frac{\mu_{t+j} - m}{c - m} - \frac{\mu_t - m}{M - m} & \text{if } p_L^* \geq m. \end{cases}$$

Rewriting the restrictions on p_L^* in terms of μ_t and μ_{t+j} yields (13). On the other hand, when $\Delta\mu_{t,j} \geq 0$, the lower trend bound is obtained by minimizing (22) with respect to (p, \bar{q}) , subject to $0 \leq p \leq m$ and $0 \leq \bar{q} \leq (1 - M)/b$, but the solution cannot be easily characterized.

Upper Bound for Trend

For given values of p and \bar{q} , the upper bound for $\Delta\pi_{t,j}$ was given in (12) and is repeated here:

$$\Delta\pi_{t,j}^U = \frac{\mu_{t+j} - p}{1 - p - b\bar{q}} - \frac{\mu_t - p}{1 - p - a\bar{q}}. \quad (24)$$

It follows that

$$\frac{d(\Delta\pi_{t,j}^U)}{d\bar{q}} = \frac{b(\mu_{t+j} - p)}{(1 - p - b\bar{q})^2} - \frac{a(\mu_t - p)}{(1 - p - a\bar{q})^2} \geq \frac{b\Delta\mu_{t,j}}{(1 - p - b\bar{q})^2}.$$

Therefore, when $\Delta\mu_{t,j} > 0$, the upper bound is maximized at $\bar{q} = (1 - M)/b$ and

$$\Delta\pi_{t,j}^U = \frac{\mu_{t+j} - p}{M - p} - \frac{\mu_t - p}{c - p}. \quad (25)$$

It remains to maximize (25) with respect to p . As before, the maximum occurs either at $p = 0$, $p = m$ or possibly at an interior solution p^* that sets the derivative of (25) equal to zero.⁸ It can

⁷There are two solutions to the first-order conditions; it can be shown that one of these necessarily lies outside the interval $[0, m]$.

⁸As before, we discard one of the solutions to the first-order conditions, because it lies outside the interval $[0, m]$.

be shown that the candidate solution is given by

$$p_U^* = \frac{M\sqrt{c - \mu_t} - c\sqrt{M - \mu_{t+j}}}{\sqrt{c - \mu_t} - \sqrt{M - \mu_{t+j}}}.$$

When $p_U^* \leq 0$, the upper trend bound is attained at $p = 0$; if $p_U^* \geq m$, the upper bound is attained at $p = m$, and if $0 < p_U^* < m$, the upper bound is attained at $p = p_U^*$. This leads to the following upper bound on the trend, when $\Delta\mu_{t,j} > 0$:

$$\Delta\pi_{t,j}^U = \begin{cases} \frac{\mu_{t+j}}{M} - \frac{\mu_t}{c} & \text{if } p_U^* \leq 0, \\ \frac{\mu_{t+j} - p_U^*}{M - p_U^*} - \frac{\mu_t - p_U^*}{c - p_U^*} & \text{if } 0 < p_U^* < m, \\ \frac{\mu_{t+j} - m}{M - m} - \frac{\mu_t - m}{c - m} & \text{if } p_U^* \geq m. \end{cases}$$

Rewriting the restrictions on p_U^* in terms of μ_t and μ_{t+j} yields the bounds in (14)

On the other hand, when $\Delta\mu_{t,j} \leq 0$, the upper trend bound is obtained by maximizing (24) with respect to (p, \bar{q}) , subject to $0 \leq p \leq m$ and $0 \leq \bar{q} \leq (1 - M)/b$, but the solution cannot be easily characterized.

B. IDENTIFIED SETS AND HPD INTERVALS FOR THE PREVALENCE

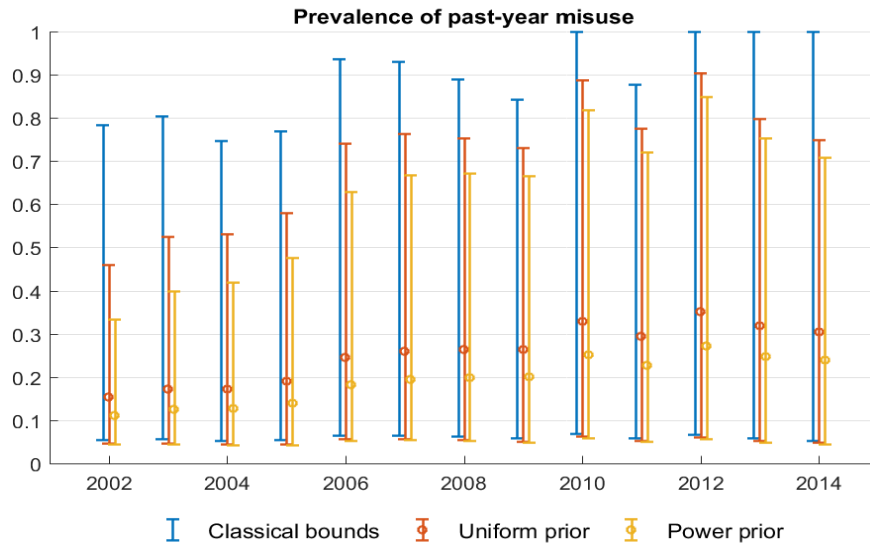


Figure B.1: classical bounds and 95% HPD intervals, q_t non-decreasing (Case II)

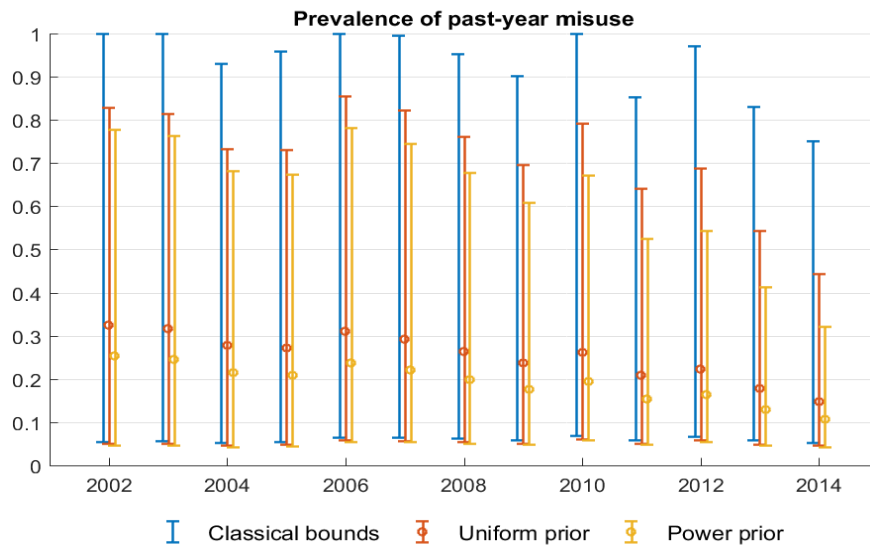


Figure B.2: classical bounds and 95% HPD intervals, q_t non-increasing (Case III)

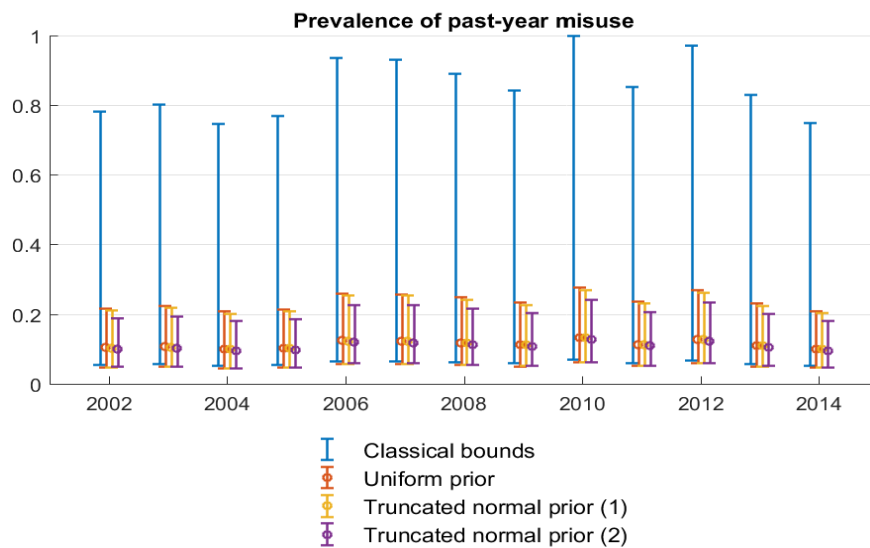


Figure B.3: classical bounds and 95% HPD intervals, uniform prior for \bar{q} (Case IV)

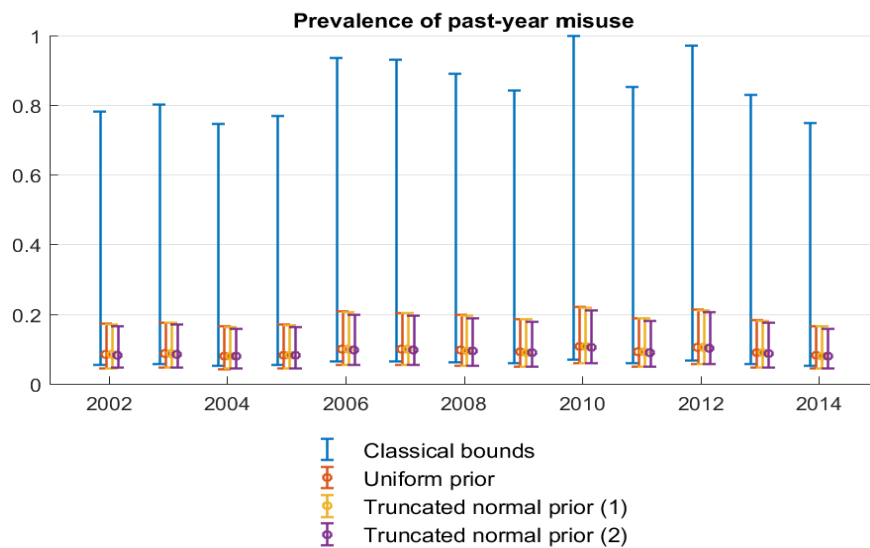


Figure B.4: Classical bounds and 95% HPD intervals, power prior for \bar{q} (Case IV)

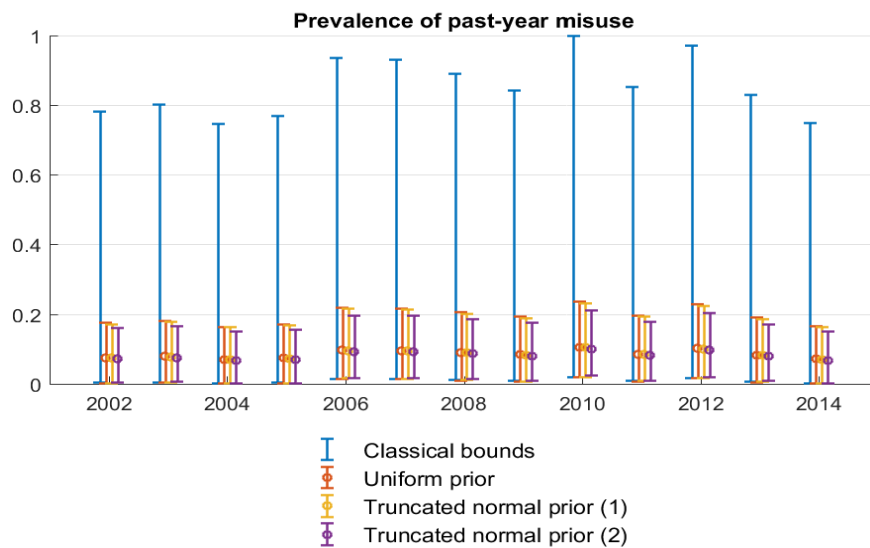


Figure B.5: classical bounds and 95% HPD intervals, uniform prior for \bar{q} (Case V)

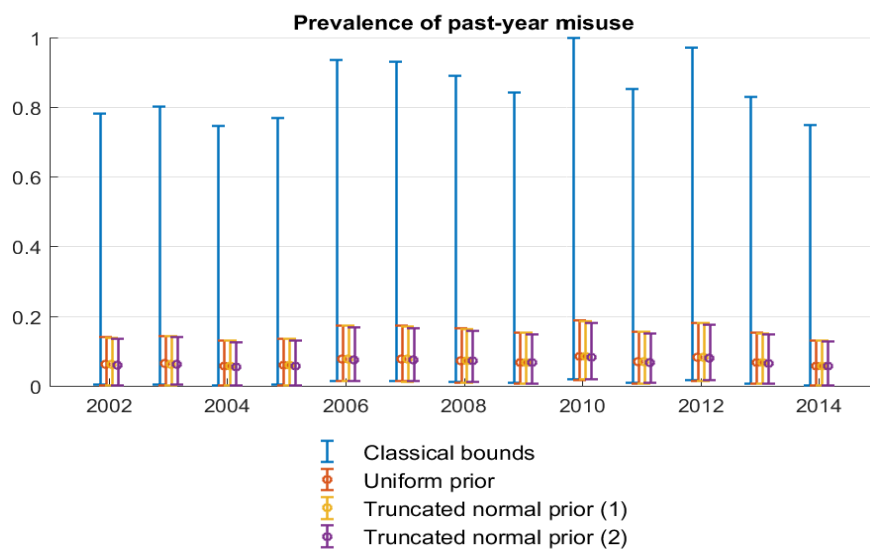


Figure B.6: Classical bounds and 95% HPD intervals, power prior for \bar{q} (Case V)

C. CHANGE IN AVERAGE PREVALENCE BETWEEN 2006-2009 AND 2010-2012

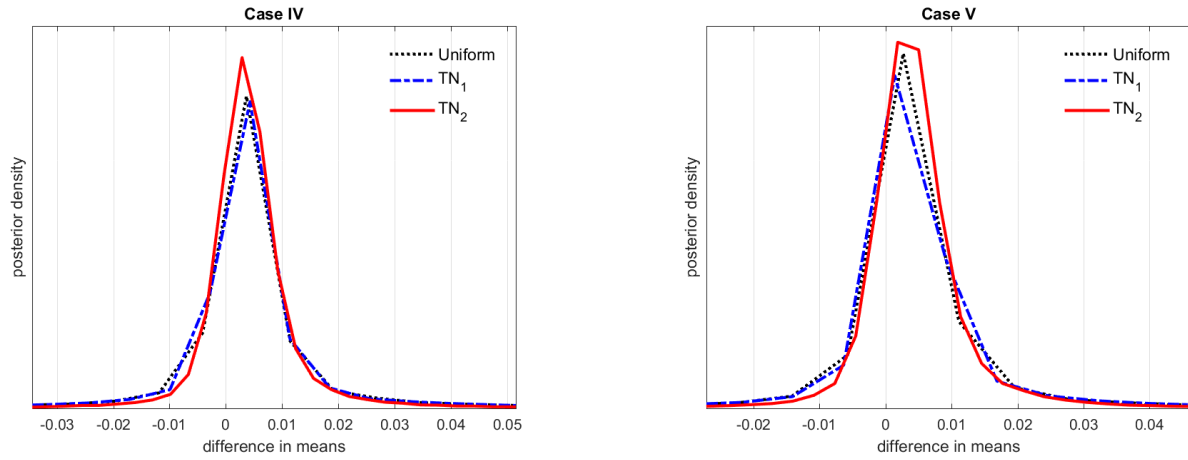


Figure C.1: posterior of difference in average prevalence between the periods 2006-2009 and 2010-2012, Case IV (left) and Case V (right).

Case	Prior	mean	std. dev.	2.5%	50%	97.5%	95% HPD	$P(+)$
IV	uniform	0.0003	0.0249	-0.0466	0.0002	0.0474	[-0.0470,0.0469]	0.5124
	TN_1	0.0004	0.0232	-0.0426	0.0002	0.0443	[-0.0420,0.0448]	0.5120
	TN_2	0.0003	0.0113	-0.0216	0.0003	0.0226	[-0.0215,0.0227]	0.5161
V	uniform	0.0004	0.0226	-0.0392	0.0003	0.0409	[-0.0394,0.0405]	0.5160
	TN_1	0.0003	0.0213	-0.0369	0.0003	0.0384	[-0.0375,0.0376]	0.5150
	TN_2	0.0004	0.0106	-0.0201	0.0003	0.0213	[-0.0198,0.0216]	0.5192

Table C.1: posterior summary of difference in average prevalence between the periods 2006-2009 and 2010-2012. $P(+)$ is the probability of an increase in average prevalence.

D. SUBGROUP ANALYSIS

Our main empirical analysis focuses on prescription opioid misuse among white men, ages 26 to 49 years old. Here, we consider misuse patterns for a select number of different subgroups. First, although our initial focus on middle-aged white men was motivated by prior research showing this to be a highly impacted population (Case and Deaton, 2015), we now compare this group to the population of 26-49 year olds who are *not* white men, based on recent research showing declines in life expectancy across a broader range of demographic groups, but still primarily concentrated among middle-aged individuals. Second, it is sometimes thought that the opioid misuse epidemic has been particularly devastating in rural areas (Keyes et al., 2014), yet more recent research suggests that non-rural populations may actually be at higher risk (Altekruse et al., 2020). The public-use NSDUH data classifies individuals as living in either (1) a large metro area, (2) a small metro area, or (3) a non-metro area. We combine the first two groups into a single metro group and compare 26-49 year olds in the metro group to 26-49 year olds in the non-metro group.

Section 4.3 considered five cases and, within each case, used several prior distributions. To focus on the between-group comparisons and to keep the analysis more concise here, we restrict ourselves to Case V and Assumption C-V. Specifically, we use a power prior for the base rate of false negative reporting (\bar{q}) and allow for a maximum deviation of 25% from that rate in each time period. The prior on those deviations is a normal distribution with mean 1, standard deviation 0.25, truncated to the interval $[0.75, 1.25]$. We believe that in this application, our prior is appealing for at least two reasons. First, false positive reporting is unlikely, but it may occur when survey respondents misinterpret the misuse question. The prior allows for that possibility. Second, while false-negative reporting has long been a concern in substance use surveys, little is known about the changes in these misreports, and their direction, over time. Our prior allows for substantial variation over time, yet does not impose any monotone trend.

A strength of our approach is the simplicity of inference for any summary statistic that can be calculated from the series $\{\pi_t; t = 1, \dots, T\}$. Simulated draws from the posterior of these parameters can be directly used to obtain random draws from the posterior of the statistic. As in our analysis in Section 4.3, $(\bar{\pi}_1 - \bar{\pi}_0)$ represents the difference in mean prevalence between the periods 2006-2009 and 2002-2005. Here we introduce another common summary parameter for this type of exercise:

slopes of the least squares trend line for π_t . A common practice is to use the repeated cross sections to estimate a trend line (either for the entire period, or for subperiods). Ignoring classification error for a moment, the estimate is obtained from the OLS regression of Y_{it} (the indicator for misuse by individual i in period t) on t . This type of analysis typically considers T fixed and inference should be based on the sampling variation from each repeated cross section. Rather than estimating a structural parameter, however, it should be viewed as a summary of the series $\{\pi_t; t = 1, \dots, T\}$. In our analysis, we examine the trend lines for the periods 2002-2007 ($\gamma_{02 \rightarrow 07}$) and 2008-2014 ($\gamma_{08 \rightarrow 14}$). Specifically, we calculate

$$\gamma_S = \frac{\frac{1}{T_S} \sum_{t \in S} t \pi_t - \left[\frac{1}{T_S} \sum_{t \in S} \pi_t \right] \bar{t}_S}{\frac{1}{T_S} \sum_{t \in S} t^2 - (\bar{t}_S)^2},$$

where S denotes the period, T_S the number of observations in S , and \bar{t}_S the average of t over S .⁹ While the actual trend in prevalence, if any, is not necessarily linear, we still calculate the linear projection slope as a useful summary across multiple time periods.

Figure D.1 compares the prevalence (left panel) and the one-year change in prevalence (right panel) between 26-49 year old white men and all others in the same age group. The HPD intervals are wider for white men, which reflects the smaller sample size. The posterior means suggest that white men experienced an increase in prevalence between 2002 and 2007, a fairly stable prevalence between 2008 and 2012, and, perhaps, a decrease in 2013 and 2014. For the comparison group, however, prevalence appears stable and fluctuates around 5% throughout the entire sample period. A similar picture emerges from the posterior summary statistics in Table D.1. For white men, the posterior mean of $(\bar{\pi}_1 - \bar{\pi}_0)$ suggests that the average prevalence in 2006-2009 was about 2 percentage points higher than during 2002-2005. The posterior of $\gamma_{02 \rightarrow 07}$ suggests an overall upward trend in prevalence between 2002 and 2007, with the 2.5% quantile exceeding zero and the HPD interval concentrated on positive values. In the comparison group, the posterior of all three parameters shows no clear evidence of any substantial changes or trends in prevalence.

Figure D.2 compares individuals 26-49 years old living in metro versus non-metro areas. The posterior mean of the prevalence of the non-metro group is sometimes higher and sometimes lower than in the metro group. No clear pattern emerges. The same can be said for the one-year change

⁹In the absence of classification error, we could regress Y_{it} on t , using data from period S . If T_S is fixed and the cross-sectional sample sizes n_t increase, this estimate will converge to γ_S .

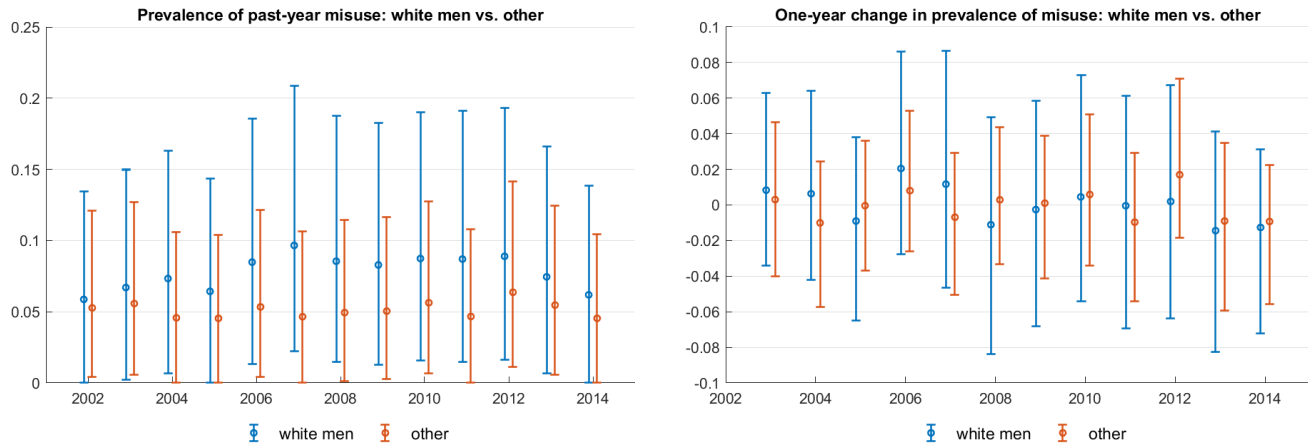


Figure D.1: posterior means and 95% HPD intervals for prevalence (left) and one-year change in prevalence (right), white men versus others.

population	parameter	mean	std. dev.	2.5%	50%	97.5%	95% HPD
white men (26-49)	$\bar{\pi}_1 - \bar{\pi}_0$	0.0216	0.0203	0.0042	0.0173	0.0738	[-0.0008,0.0632]
	$\gamma_{02 \rightarrow 07}$	0.0067	0.0066	0.0009	0.0053	0.0232	[-0.0007,0.0201]
	$\gamma_{08 \rightarrow 14}$	-0.0031	0.0049	-0.0141	-0.0024	0.0026	[-0.0128,0.0034]
other (26-49)	$\bar{\pi}_1 - \bar{\pi}_0$	0.0000	0.0137	-0.0218	0.0001	0.0217	[-0.0220,0.0213]
	$\gamma_{02 \rightarrow 07}$	-0.0011	0.0047	-0.0094	-0.0008	0.0049	[-0.0092,0.0051]
	$\gamma_{08 \rightarrow 14}$	0.0001	0.0036	-0.0053	0.0001	0.0059	[-0.0053,0.0060]

Table D.1: posterior comparison of white men versus other. The parameters are the difference in average prevalence between 2006-2009 and 2002-2005 ($\bar{\pi}_1 - \bar{\pi}_0$), and the linear projection slopes for prevalence in the periods 2002-2007 ($\gamma_{02 \rightarrow 07}$) and 2008-2014 ($\gamma_{08 \rightarrow 14}$).

in prevalence. From Table D.2, the posterior mean of ($\bar{\pi}_1 - \bar{\pi}_0$) is positive and almost twice as large for non-metro compared to metro, suggesting that non-metro areas experienced a larger increase in the prevalence of misuse. This is confirmed by comparing the posteriors of $\gamma_{02 \rightarrow 07}$ between the two groups. The left panel of Figure D.2 shows, however, that this might be driven by an outlier in the prevalence in non-metro areas in 2007.

Although the subgroup analyses we have presented here were intended to illustrate our proposed methods, they nonetheless provide some initial evidence for potential between-group differences and some important implications for policy. We find evidence that racial differences in prevalence may be declining, primarily caused by a decreasing trend among whites. We also find indications that the early focus on rural communities may have missed important impacts in urban communities. Our Bayesian approach shows no consistent difference in the prevalence between metro and non-

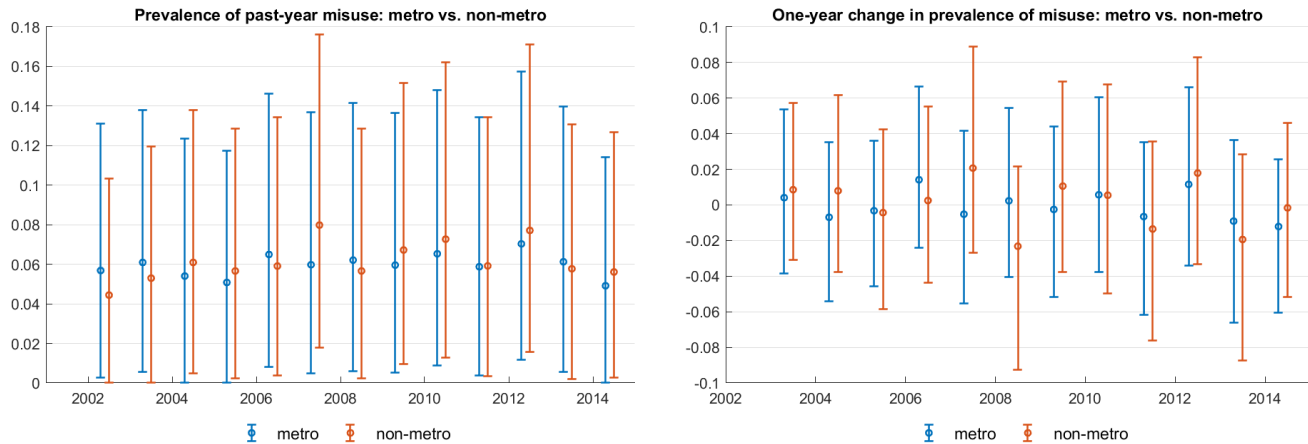


Figure D.2: posterior means and 95% HPD intervals for prevalence (left) and one-year change in prevalence (right), metro versus non-metro.

population	parameter	mean	std. dev.	2.5%	50%	97.5%	95% HPD
metro	$\bar{\pi}_1 - \bar{\pi}_0$	0.0059	0.0152	-0.0131	0.0045	0.0371	[-0.0151,0.0345]
	$\gamma_{02 \rightarrow 07}$	0.0007	0.0051	-0.0069	0.0005	0.0095	[-0.0077,0.0087]
	$\gamma_{08 \rightarrow 14}$	-0.0011	0.0039	-0.0086	-0.0008	0.0043	[-0.0082,0.0047]
non-metro	$\bar{\pi}_1 - \bar{\pi}_0$	0.0119	0.0169	-0.0064	0.0096	0.0512	[-0.0099,0.0460]
	$\gamma_{02 \rightarrow 07}$	0.0055	0.0060	-0.0002	0.0044	0.0199	[-0.0012,0.0175]
	$\gamma_{08 \rightarrow 14}$	-0.0006	0.0042	-0.0083	-0.0005	0.0061	[-0.0078,0.0065]

Table D.2: posterior comparison of metro versus non-metro. The parameters are the difference in average prevalence between 2006-2009 and 2002-2005 ($\bar{\pi}_1 - \bar{\pi}_0$), and the linear projection slopes for prevalence in the periods 2002-2007 (γ_{02-07}) and 2008-2014 (γ_{08-14}).

metro areas, although we found some weak evidence of a greater increase in the prevalence rate in non-metro areas over time.

While the priors that we used were the same across different groups, the posteriors were not. This is true because the data are informative about the bounds of the identified sets. Differences in these sets are revealed, which in turn affects the posteriors. Thus, despite the lack of formal identification, we learn something about heterogeneity within the larger population. Of course, this does not provide a formal statistical test for heterogeneity. While we use a parametric prior distribution, our Bayesian analysis is otherwise non-parametric. This aligns with the classical bounding literature, which has focused on the non-parametric identification and estimation of bounds.

Conducting subgroup analyses provides an informal way to see how covariates affect posterior

estimates of the prevalence. It is less clear how to assess that impact formally without using a parametric, model-based approach (e.g., a probit or logit model for prevalence). A parametric model would impose additional distributional assumptions (relative to our approach) and would also require the researcher to specify priors for each additional parameter. Without such a model, it is not clear how to proceed. For example, how can we test whether the prevalence of misuse in 2007 was higher for white men compared to others (see Figure D.1), when the prevalence in both groups is partially identified, and the identified sets and support of the posterior have substantial overlap?

It is equally difficult to test for differences in misreporting behavior between groups, or to determine how covariates might affect the probabilities of false negatives and false positives. It is reasonable to expect that misreporting rates can differ substantially between groups. In this section we have used the same prior for the misreporting probabilities for all groups. If reliable information on differential misreporting is available, the researcher can use it to specify group-specific priors. Given the non-parametric nature of our model, and the fact that misreporting rates are also partially identified, our approach does not allow us to determine how misreporting varies between groups. To assess the latter, one could specify a parametric misreporting model as in Nguimkeu et al. (2019). In our empirical context, this model would combine a bivariate probit model for actual misuse and misreporting with the assumption of no false positives ($p = 0$).¹⁰ Joint normality of the errors and an exclusion restriction between the misuse and misreporting equations are sufficient to identify the model parameters. While the use of a parametric, model-based approach can be very useful in practice, we do not pursue it further here, since it departs from the non-parametric Bayesian approach advocated in this paper.

¹⁰Nguimkeu et al. (2019) build on Poirier (1980), who first proposed this *partial observability* model.

E. EXTENSION TO REGRESSION MODELING

Here, we briefly discuss how our Bayesian approach may be adapted in the context of a regression model where the misclassified binary indicator appears as an explanatory variable. Using X_{it}^* and X_{it} (instead of Y_{it}^* and Y_{it}) to denote the true and observed values, respectively, of the binary indicator, and letting Y_{it} denote the dependent variable of interest, consider the model

$$Y_{it} = \alpha + X_{it}^* \beta + W_{it}' \gamma + U_{it}, \quad (26)$$

where W_{it} is a vector of covariates that are observed without error, and U_{it} is a residual. To simplify the discussion, we assume that X_{it}^* and W_{it} are exogenous and $E(U_{it}|X_{it}^*, W_{it}) = 0$. In an empirical application, Y_{it} can represent various health or labor market outcomes, and we are interested in estimating the association between opioid misuse and such outcomes. If the data are cross-sectional, it is well-known that β is partially identified. The upper and lower bounds for β are derived in Bollinger (1996) and can be expressed as functions of the first and second moments of $(W_{it}', X_{it}, Y_{it})$. This allows us to use the approach of Bollinger and Van Hasselt (2017b) to conduct inference about β . In this approach, we combine a nonparametric posterior for ϕ with a user-specified prior for the misclassification probabilities to simulate a sample from the posterior of β .

The situation is notably different in a repeated cross section. Using the results of Bollinger (1996), it can be shown that in time period $t = 1, \dots, T$,

$$\beta = \frac{(1 - p_t - q_t)\tilde{c}_t}{\tilde{s}_t - q_t(\mu_t - p_t) - p_t(1 - \mu_t)}, \quad (27)$$

where \tilde{c}_t is the covariance between \tilde{Y}_{it} and \tilde{X}_{it} , \tilde{s}_t is the variance of \tilde{X}_{it} , and $(\tilde{X}_{it}, \tilde{Y}_{it})$ are the residuals of (X_{it}, Y_{it}) after projecting out W_{it} . The (functions of) moments $(\tilde{c}_t, \tilde{s}_t, \mu_t)$ are identified and estimable, whereas (p_t, q_t) are unknown and give rise to partial identification of β . The assumption that β in equation (27) is fixed over time may suggest that β is identified if T is large enough, but this is not the case. An additional time period of data yields an additional parameter restriction through (27), but also introduces two additional parameters (p_t, q_t) . Hence, the problem of partial identification remains.

Equation (27) suggests a multi-step way to sample from the posterior of β . It is similar to the approaches described in Section 3.3 and in Bollinger and Van Hasselt (2017b). First, the parameters $(\tilde{c}_t, \tilde{s}_t, \mu_t)$ are functions of sample moments and defined through a set of moment conditions. The Bayesian bootstrap of Rubin (1981) allows us to generate draws from their nonparametric posterior (Chamberlain and Imbens, 2003). Second, we generate a draw from the conditional prior of (p_t, q_t) , given $(\tilde{c}_t, \tilde{s}_t, \mu_t)$.¹¹ Finally, a draw from the posterior of β is obtained by calculating the right-hand side of (27). Specifying a conditional prior for (p_t, q_t) for $t = 1, \dots, T$ is more complicated, however, than in Section 3.2. Equation (27) implies restrictions over time on the sequence (p_t, q_t) , which have to be incorporated in the prior. For example, when $T = 2$, there are four unknown probabilities. From equation (27),

$$\frac{(1 - p_1 - q_1)\tilde{c}_1}{\tilde{s}_1 - q_1(\mu_1 - p_1) - p_1(1 - \mu_1)} = \frac{(1 - p_2 - q_2)\tilde{c}_2}{\tilde{s}_2 - q_2(\mu_2 - p_2) - p_2(1 - \mu_2)},$$

so that only three of the probabilities are unrestricted. More generally, with T time periods, the joint prior on the misclassification probabilities is supported on a $(T + 1)$ -dimensional subspace.

Another important application that our analysis facilitates is the case where estimates of prevalence are included in a regression setting at a higher level of aggregation, such as counties or states. For example,

$$Y_{it} = \alpha + \pi_{it}\beta + W'_{it}\gamma + U_{it}.$$

Here, the prevalence π_{it} is unobserved. Typically, estimates μ_{it} from repeated cross sections of the appropriate geography are used but these are now biased for π_{it} . If individual-level data is available for each state or county, our approach can be used as a first step to generate the posterior of π_{it} . The second step involves estimating the regression model, and there are a number of different ways to proceed. For example, we can generate draws from the posterior of π_{it} as a data augmentation step in a Gibbs sampler to estimate α , β and γ (Tanner and Wong, 1987). Alternatively, we could use the posterior of π_{it} to implement a multiple (stochastic) imputation approach. This would yield multiple sets of estimates for the regression parameters and would account for the uncertainty surrounding the unobserved π_{it} .

¹¹The probabilities (p_t, q_t) are also partially identified and bounded by (functions of) the first and second moments of $(W'_{it}, X_{it}, Y_{it})$.

There are many applications where mismeasured indicator variables, or the estimates derived from them, are useful. In cases where the parameters of interest can be readily written as functions of the series $\{\pi_t; t = 1, \dots, T\}$, our approach can be implemented. In other situations, our approach may provide informative priors on unknown parameters or unobserved data. Finally, additional information about the relationship between the mismeasured Y_{it} and other variables may improve estimation or tighten bounds. We leave a more detailed investigation of these issues for future work.