

Sensitivity Analysis for Unmeasured Confounding in Meta-Analyses

Maya B. Mathur^{1,2*} and Tyler J. VanderWeele^{1,3}

¹ Department of Biostatistics, Harvard T. H. Chan School of Public Health, Boston, MA, USA

²Quantitative Sciences Unit, Stanford University, Palo Alto, CA, USA

³Department of Epidemiology, Harvard T. H. Chan School of Public Health, Boston, MA, USA

July 31, 2017

*Corresponding author:

mmathur@stanford.edu

Quantitative Sciences Unit (c/o Inna Sayfer)

1070 Arastradero Road

Palo Alto, CA

94305

Summary

Random-effects meta-analyses of observational studies can produce biased estimates if the synthesized studies are subject to unmeasured confounding. We propose sensitivity analyses quantifying the extent to which unmeasured confounding of specified magnitude could reduce to below a certain threshold the proportion of true effect sizes that are scientifically meaningful. We also develop converse methods to estimate the strength of confounding capable of reducing the proportion of scientifically meaningful true effects to below a chosen threshold. These methods apply when a “bias factor” is assumed to be normally distributed across studies or is assessed across a range of fixed values. Our estimators are derived using recently proposed sharp bounds on confounding bias within a single study that do not make assumptions regarding the unmeasured confounders themselves or the functional form of their relationships to the exposure and outcome of interest. We provide an R package, `ConfoundedMeta`, and a freely available online graphical user interface that compute point estimates and inference and produce plots for conducting such sensitivity analyses. These methods facilitate principled use of random-effects meta-analyses of observational studies to assess the strength of causal evidence for a hypothesis.

Key words: Bias; Confounding; Meta-analysis; Observational studies; Sensitivity analysis

1. INTRODUCTION

Meta-analyses can be indispensable for assessing the overall strength of evidence for a hypothesis and for precisely estimating effect sizes through aggregation of estimates. However, conclusions drawn from meta-analyses are only as reliable as the synthesized studies themselves; systematic bias in the meta-analyzed studies typically produces bias in the pooled point estimate (Egger et al., 1998). A common source of bias is unmeasured confounding (Shrier et al., 2007), which is our focus in this paper. When eliminating such bias by restricting attention to well-designed randomized studies is infeasible because the exposure cannot be randomized, an attractive option is to conduct sensitivity analyses assessing the extent to which unmeasured confounding of varying magnitudes could have compromised the results of the meta-analysis.

Existing sensitivity analyses for confounding bias or other internal biases in meta-analysis estimate a bias-corrected pooled point estimate by directly incorporating one or more bias parameters in the likelihood and placing a Bayesian prior on the distribution of these parameters (McCandless, 2012; Welton et al., 2009). An alternative frequentist approach models bias as additive or multiplicative within each study and then uses subjective assessment to elicit study-specific bias parameters (Turner et al., 2009). Although useful, these approaches typically require strong assumptions on the nature of unmeasured confounding (for example, requiring a single binary confounder), rely on the arbitrary specification of additive or multiplicative effects of bias, or require study-level estimates rather than only meta-analytic pooled estimates. Furthermore, the specified bias parameters do not necessarily lead to precise practical interpretations.

An alternative approach is to analytically bound the effect of unmeasured confounding on the results of a meta-analysis. To this end, bounding methods are currently available for point estimates of individual studies. We focus on sharp bounds derived by Ding & VanderWeele (2016) because of their generality and freedom from assumptions regarding the nature of the unmeasured confounders or the functional forms of their relationships with the exposure of interest and outcome. This approach subsumes several earlier approaches (Cornfield et al., 1959; Flanders & Khoury, 1990; Schlesselman, 1978) and relies on only two simple sensitivity parameters representing the strength of association of the unmeasured confounders with, firstly, the exposure and, secondly, the outcome.

The present paper extends these analytic bounds for single studies to the meta-analytic setting. Using standard estimates from a random-effects meta-analysis and intuitively interpretable sensitivity parameters on the magnitude of confounding, these results enable inference about the size of the true, unconfounded effects in a potentially heterogeneous population of studies. That is, we can select a minimum threshold of scientific importance for the magnitude of the true effect in any given study. If sensitivity analysis for unmeasured confounding indicates that too few studies in the meta-analysis have a true effect stronger than this threshold, then arguably the results of the meta-analysis are not robust to confounding, and scientifically meaningful causal conclusions are not warranted despite the observed point estimate. To this end, we develop estimators that answer the questions: “What proportion of studies would have a true effect size stronger than q in the presence of unmeasured confounding of a specified strength?” and “How severe would unmeasured confounding need to be to reduce to less than r the proportion of studies with true effect size stronger than q ?”. This approach to sensitivity analysis is essentially a meta-analytic extension of a recently

proposed metric (the “E-value”) that quantifies, for a single study, the minimum confounding bias capable of reducing the true effect to a chosen threshold (VanderWeele & Ding, 2017). We provide and demonstrate use of an R package (`ConfoundedMeta`) and a free, interactive online user interface for conducting such analyses and creating plots.

2. EXISTING BOUNDS ON CONFOUNDING BIAS IN A SINGLE STUDY

Ding & VanderWeele (2016) developed bounds for a single study as follows. Let X denote a binary exposure, Y a binary outcome, Z a vector of measured confounders, and U one or more unmeasured confounders. Let:

$$RR_{XY|z}^c = \frac{P(Y = 1 | X = 1, Z = z)}{P(Y = 1 | X = 0, Z = z)}$$

be the confounded relative risk (RR) of Y for $X = 1$ versus $X = 0$ conditional or stratified on the measured confounders $Z = z$.

Let its true, unconfounded counterpart standardized to the population be:

$$RR_{XY|z}^t = \frac{\sum_u P(Y = 1 | X = 1, Z = z, U = u) P(U = u | Z = z)}{\sum_u P(Y | X = 0, Z = z, U = u) P(U = u | Z = z)}$$

(Throughout, we use the term “true” as a synonym for “unconfounded” or “causal” when referring to both sample and population quantities. Also, henceforth, we condition implicitly on $Z = z$, dropping the explicit notation for brevity.)

Let $RR_{Xu} = P(U = u | X = 1) / P(U = u | X = 0)$. Define the first sensitivity parameter as $RR_{XU} = \max_u (RR_{Xu})$; that is, the maximal relative risk of $U = u$ for $X = 1$ versus $X = 0$ across strata of U . (If U is binary, this is just the relative risk relating X and U .) Next, for each stratum x of X , define a relative risk of Y on U , maximized across all possible

contrasts of U :

$$RR_{UY|X=x} = \frac{\max_u P(Y = 1|X = x, U = u)}{\min_u P(Y = 1|X = x, U = u)}, x \in \{0, 1\}$$

Define the second sensitivity parameter as $RR_{UY} = \max(RR_{UY|X=0}, RR_{UY|X=1})$. That is, considering both strata of X , it is the largest of the maximal relative risks of Y on U conditional on X . Then, Ding & VanderWeele (2016) showed that a sharp bound for the true effect is:

$$RR_{XY}^t \geq RR_{XY}^c / \frac{RR_{XU} \cdot RR_{UY}}{RR_{XU} + RR_{UY} - 1} \quad (2.1)$$

where we will refer to the “bias factor” $\frac{RR_{XU} \cdot RR_{UY}}{RR_{XU} + RR_{UY} - 1}$ as B .

If the two sensitivity parameters are equal ($RR_{XU} = RR_{UY}$), then to produce a bias factor B , each must exceed $B + \sqrt{B^2 - B}$ (Ding & VanderWeele, 2016). Thus, a useful transformation of B is the “confounding strength scale”, g , which is the minimum size of RR_{XU} and RR_{UY} under the assumption that they are equal:

$$g = B + \sqrt{B^2 - B} \quad \Leftrightarrow \quad B = \frac{g^2}{2g - 1} \quad (2.2)$$

If $RR_{XY}^c < 1$ (henceforth the “apparently preventive case”), then Equation (2.1) becomes (Ding & VanderWeele, 2016):

$$RR_{XY}^t \leq RR_{XY}^c \cdot \frac{RR_{XU}^* \cdot RR_{UY}}{RR_{XU}^* + RR_{UY} - 1}$$

where $RR_{XU}^* = \max_u (RR_{Xu}^{-1})$, i.e., the maximum of the inverse relative risks, rather than the relative risks themselves. Thus, B remains ≥ 1 , and we have $RR_{XY}^t \geq RR_{XY}^c$.

Although these results hold for multiple confounders, in the development to follow, we will use a single, categorical unmeasured confounder for clarity. However, all results can easily be

interpreted without assumptions on the type of exposure and unmeasured confounders, for instance by interpreting the relative risks defined above as “mean ratios” (Ding & VanderWeele, 2016).

3. RANDOM-EFFECTS META-ANALYSIS SETTING

In this paper, we use the aforementioned analytic bounds to derive counterparts for the random-effects meta-analysis model with the standard DerSimonian-Laird point estimate. This model assumes that each of k studies measures a potentially unique effect size M , such that $M \sim_{iid} N(\mu, V)$ for a grand mean μ and variance V . Let y_i be the point estimate of the i^{th} study and σ_i^2 the within-study variance (with the latter assumed fixed and known).

Analysis proceeds by first estimating V via one of many possible estimators, denoted τ^2 (Veroniki et al., 2015), then estimating μ via a weighted mean defined as:

$$\hat{y}_R = \frac{\sum_{i=1}^k w_i y_i}{\sum_{i=1}^k w_i}$$

The weights are inversely proportional to the total variance of each study (a sum of the between-study variance and the within-study variance), such that $w_i = 1/(\tau^2 + \sigma_i^2)$.

4. MAIN RESULTS

Consider k studies measuring relative risks with confounded population effect sizes on the log- RR scale, denoted M^c , such that $M^c \sim N(\mu^c, V^c)$. (Other outcome measures are considered briefly in the Discussion.) Let the corresponding true effects be M^t with expectation μ^t and variance V^t . Let \hat{y}_R^c be the standard inverse-variance-weighted random effects point estimate and τ_c^2 be a heterogeneity estimate, both computed from the confounded data.

Consider the bias factor on the log scale, $B^* = \log\left(\frac{RR_{XU} \cdot RR_{UY}}{RR_{XU} + RR_{UY} - 1}\right)$, and allow it to vary across studies under the assumption that $B^* \sim N(\mu_{B^*}, \sigma_{B^*}^2)$ independently of M^t . That is, we assume that the bias factor is independent of the true effects but not the confounded effects: naturally, studies with larger bias factors will tend to obtain larger effect sizes. The normality assumption on the bias factor holds approximately if, for example, its components (RR_{XU} and RR_{UY}) are identically and independently normal with relatively small variance (Web Appendix). We now develop three estimators enabling sensitivity analyses.

4.1. Proportion of studies with large effect sizes as a function of the bias factor

For an **apparently causative relative risk** ($\hat{y}_R^c > 0$, or equivalently the confounded pooled RR is greater than 1), define $p(q) = P(M^t > q)$ for any threshold q , i.e., the proportion of studies with true effect sizes larger than q . Then a consistent estimator of $p(q)$ is:

$$\hat{p}(q) = 1 - \Phi\left(\frac{q + \mu_{B^*} - \hat{y}_R^c}{\sqrt{\tau_c^2 - \sigma_{B^*}^2}}\right), \quad \tau_c^2 > \sigma_{B^*}^2$$

where Φ denotes the standard normal cumulative distribution function. In the special case in which the bias factor is fixed to μ_{B^*} across all studies, the same formula applies with $\sigma_{B^*}^2 = 0$.

Many common choices of heterogeneity estimators, τ_c^2 , are asymptotically independent of \hat{y}_R^c (Web Appendix), an assumption used for all standard errors in the main text. Results relaxing this assumption appear throughout the Web Appendix. An application of the delta method thus yields an approximate standard error:

$$\widehat{\text{SE}}(\hat{p}(q)) \approx \sqrt{\frac{\widehat{\text{Var}}(\hat{y}_R^c)}{\tau_c^2 - \sigma_{B^*}^2} + \frac{\widehat{\text{Var}}(\tau_c^2)(q + \mu_{B^*} - \hat{y}_R^c)^2}{4(\tau_c^2 - \sigma_{B^*}^2)^3}} \cdot \phi\left(\frac{q + \mu_{B^*} - \hat{y}_R^c}{\sqrt{\tau_c^2 - \sigma_{B^*}^2}}\right)$$

where ϕ denotes the standard normal density function. (If $\tau_c^2 \leq \sigma_{B^*}^2$, leaving one of the denominators undefined, this indicates that there is so little observed heterogeneity in the

confounded effect sizes that, given the specified bias distribution, V^t is estimated to be less than 0. Therefore, attention should be limited to a range of values of $\sigma_{B^*}^2$ such that $\tau_c^2 > \sigma_{B^*}^2$.)

For an **apparently preventative relative risk** ($\hat{y}_R^c < 0$ or the confounded pooled RR is less than 1), define instead $p(q) = P(M^t < q)$, i.e., the proportion of studies with true effect sizes less than q . Then a consistent estimator is:

$$\hat{p}(q) = \Phi \left(\frac{q - \mu_{B^*} - \hat{y}_R^c}{\sqrt{\tau_c^2 - \sigma_{B^*}^2}} \right), \tau_c^2 > \sigma_{B^*}^2$$

with approximate standard error:

$$\widehat{\text{SE}}(\hat{p}(q)) = \sqrt{\frac{\widehat{\text{Var}}(\hat{y}_R^c)}{\tau_c^2 - \sigma_{B^*}^2} + \frac{\widehat{\text{Var}}(\tau_c^2)(q - \mu_{B^*} - \hat{y}_R^c)^2}{4(\tau_c^2 - \sigma_{B^*}^2)^3}} \cdot \phi \left(\frac{q - \mu_{B^*} - \hat{y}_R^c}{\sqrt{\tau_c^2 - \sigma_{B^*}^2}} \right) \quad (4.1)$$

Because $\hat{p}(q)$ is monotonic in $\sigma_{B^*}^2$, the homogeneous bias case (i.e., $\sigma_{B^*}^2 = 0$) provides either an upper or lower bound on $\hat{p}(q)$ (Table 1). We later return to the practical utility of these results.

4.2. Bias factor required to reduce proportion of large effect sizes to a threshold

Conversely, we might consider the minimum common bias factor (on the RR scale) capable of reducing to less than r the proportion of studies with true effect exceeding q . We accordingly define $T(r, q) = B : P(M^t > q) = r$ to be this quantity, with B taken to be constant across studies. (Note that taking B to be constant does not necessarily imply that the unmeasured confounders themselves are identical across studies.) Then for an **apparently causative relative risk**, a consistent estimator for the the minimum common bias capable of reducing to less than r the proportion of studies with effects surpassing q is:

$$\hat{T}(r, q) = \exp \left\{ \Phi^{-1}(1 - r) \sqrt{\tau_c^2} - q + \hat{y}_R^c \right\} \quad (4.2)$$

with approximate standard error:

$$\widehat{\text{SE}}\left(\widehat{T}(r, q)\right) = \exp\left\{\sqrt{\widehat{\tau}_c^2}\left(\Phi^{-1}(1-r)\right) - q + \widehat{y}_R^c\right\} \sqrt{\widehat{\text{Var}}\left(\widehat{y}_R^c\right) + \frac{\widehat{\text{Var}}\left(\tau_c^2\right)\left(\Phi^{-1}(1-r)\right)^2}{4\tau_c^2}} \quad (4.3)$$

For an **apparently preventive relative risk**, we can instead consider the minimum common bias factor (on the RR scale) capable of reducing to less than r the proportion of studies with true effect less than q , thus defining $T(r, q) = B : P(M^t < q) = r$. Then a consistent estimator is:

$$\widehat{T}(r, q) = \exp\left\{q - \widehat{y}_R^c - \Phi^{-1}(r)\sqrt{\widehat{\tau}_c^2}\right\} \quad (4.4)$$

with approximate standard error:

$$\widehat{\text{SE}}\left(\widehat{T}(r, q)\right) = \exp\left\{q - \widehat{y}_R^c - \sqrt{\widehat{\tau}_c^2}\left(\Phi^{-1}(r)\right)\right\} \sqrt{\widehat{\text{Var}}\left(\widehat{y}_R^c\right) + \frac{\widehat{\text{Var}}\left(\tau_c^2\right)\left(\Phi^{-1}(r)\right)^2}{4\tau_c^2}} \quad (4.5)$$

4.3. Confounding strength required to reduce proportion of large effect sizes to a threshold

Under the assumption that the two components of the common bias factor are equal as in Equation 2.2, such that $g = RR_{XU} = RR_{UY}$, the bias can alternatively be parameterized on the confounding strength scale. Consider the minimum confounding strength required to lower to less than r the proportion of studies with true effect exceeding q and accordingly define $G(r, q) = g : P(M^t > q) = r$. For both the **apparently causative and the apparently preventive cases**, an application of Equation 2.2 yields:

$$\widehat{G}(r, q) = \widehat{T}(r, q) + \sqrt{\left(\widehat{T}(r, q)\right)^2 - \widehat{T}(r, q)} \quad (4.6)$$

with approximate standard error:

$$\widehat{\text{SE}}\left(\widehat{G}(r, q)\right) = \widehat{\text{SE}}\left(\widehat{T}(r, q)\right) \cdot \left(1 + \frac{2\widehat{T}(r, q) - 1}{2\sqrt{\widehat{T}(r, q)^2 - \widehat{T}(r, q)}}\right)$$

5. PRACTICAL USE AND INTERPRETATION

The estimators $\widehat{p}(q)$, $\widehat{T}(r, q)$, and $\widehat{G}(r, q)$ enable several types of sensitivity analysis. Firstly, $\widehat{p}(q)$ can be computed over a range of values of μ_{B^*} and $\sigma_{B^*}^2$. If $\widehat{p}(q)$ remains large for even large values of μ_{B^*} , this indicates that even if the influence of unmeasured confounding were substantial, a large proportion of studies nevertheless would have true effects of scientifically meaningful magnitudes. Similarly, $\widehat{T}(r, q)$ and $\widehat{G}(r, q)$ can be computed for r representing a “large enough” proportion of studies to warrant scientific interest; large values would again lead to the conclusion that results of the meta-analysis are relatively robust to unmeasured confounding. For example, by choosing $q = \log(1.10)$ and $r = 0.20$ and computing $\widehat{T}(r, q) = 2.50$ (equivalently, $\widehat{G}(r, q) = 4.44$), one might conclude: “The results of this meta-analysis are relatively robust to unmeasured confounding, insofar as a bias factor of 2.50 on the relative risk scale (e.g., a confounder associated with the exposure and outcome by risk ratios of 4.44 each) in each study would be capable of reducing to less than 20% the proportion of studies with true relative risks greater than 1.10, but weaker confounding could not do so.” On the other hand, small values of $\widehat{p}(q)$, $\widehat{T}(r, q)$, and $\widehat{G}(r, q)$ indicate that only weak unmeasured confounding would be required to reduce the effects to a scientifically unimportant level; the meta-analysis would therefore not warrant strong scientific conclusions regarding causation.

A general guideline might be to use $q = \log 1.10$ for an apparently causative relative risk or $q = \log 0.90$ for an apparently preventive relative risk. When the number of studies, k ,

is large (for example, ≥ 10), one might require at least 10% of studies ($r = 0.10$) to have effect sizes above q for results to be of scientific interest. For $k < 10$, one might select a higher threshold, such as $r = 0.20$ (thus requiring at least 20% of studies to have effects more extreme than, for example, $\log 1.10$). Of course, these guidelines can and should be adapted based on the substantive application. Furthermore, note that the amount of bias that would be considered “implausible” must be determined with attention to the design quality of the synthesized studies: a large bias factor may be plausible for a set of studies with poor confounding control and with high potential for unmeasured confounding, but not for a set of better-designed studies in which the measured covariates already provide good control of confounding.

Sensitivity analyses based on $\widehat{p}(q)$ should be reported for a wide range of values for μ_{B^*} and with $\sigma_{B^*}^2$ ranging from 0 to somewhat less than τ_c^2 . The bounds achieved when $\sigma_{B^*}^2 = 0$ (Table 1) can provide useful conservative analyses. For example, for $\widehat{y}_R^c > 0$ and $q > \widehat{\mu}^t$, the $\sigma_{B^*}^2 = 0$ case provides an upper bound on $\widehat{p}(q)$. When concluding that results are not robust to unmeasured confounding, the analysis with $\sigma_{B^*}^2 = 0$ is therefore conservative in that fewer true effect sizes would surpass q under heterogeneous bias. For example, if we calculated $\widehat{T}(r = 0.20, q = \log 1.10) = 1.20$, then an analysis like this would yield conclusions such as: “The results of this meta-analysis are relatively sensitive to unmeasured confounding. Even a bias factor as small as 1.20 in each study would reduce to less than 20% the proportion of studies with true relative risks greater than 1.10, and if the bias in fact varied across studies, then even fewer studies would surpass this effect size threshold.”

6. SOFTWARE AND APPLIED EXAMPLE

The present methods are implemented in an R package, `ConfoundedMeta`, which produces point estimates and inference for sensitivity analyses, tables across a user-specified grid of sensitivity parameters, and various plots. Descriptions of each function are provided in the Web Appendix and standard R documentation. A graphical user interface implementing the main functions is freely available (https://mmathur.shinyapps.io/meta_gui_2/).

We illustrate the package’s basic capabilities using an existing meta-analysis assessing, among several outcomes, the association of high versus low daily intake of soy protein with breast cancer risk among women (Trock et al., 2006). The analysis comprised 20 observational studies that varied in their degree of adjustment for suspected confounders, such as age, body mass index (BMI), and other risk factors. To obtain τ_c^2 and $\widehat{\text{Var}}(\tau_c^2)$ (which were not reported), we obtained study-level summary measures as reported in a table from Trock et al. (2006), approximating odds ratios with risk ratios given the rare outcome. This process is automated in the function `ConfoundedMeta::scrape_meta`. We estimated $\widehat{y}_R^c = \log 0.82$, $\widehat{\text{SE}}(\widehat{y}_R^c) = 8.8 \times 10^{-2}$ via the Hartung & Knapp (2001) adjustment (whose advantages were demonstrated by IntHout et al. (2014)), $\tau_c^2 = 0.10$ via the Paule & Mandel (1982) method, and $\widehat{\text{SE}}(\tau_c^2) = 5.0 \times 10^{-2}$.

Figure 1 (produced by `ConfoundedMeta::sens_plot`) displays the estimated proportion of studies with true relative risks < 0.90 as a function of either the bias factor or the confounding strength, holding constant $\sigma_{B^*}^2 = 0.01$. Table 2 (produced by `ConfoundedMeta::sens_table`) displays $\widehat{T}(r, q)$ and $\widehat{G}(r, q)$ across a grid of values for r and q . For example, only a bias factor exceeding 1.63 on the relative risk scale (equivalently, confounding association strengths of

2.64) could reduce to less than 10% the proportion of studies with true relative risks < 0.90 . However, variable bias across studies would reduce this proportion, and the confidence interval is wide.

7. SIMULATION STUDY

We assessed finite-sample performance of inference on $\widehat{p}(q)$ in a simple simulation study. While fixing the mean and variance of the true effects to $\mu^t = \log 1.4$ and $V^t = 0.15$ and the bias parameters to $\mu_{B^*} = \log 1.6$ and $\sigma_{B^*}^2 = 0.01$, we varied the number of studies ($k \in \{15, 25, 50, 200\}$) and the average sample size N within each study ($E[N] \in \{300, 500, 1000\}$). The fixed parameters were chosen to minimize artifacts from discarding pathological samples with $\tau_c^2 < \sigma_{B^*}^2$ or with truncated outcome probabilities due to extreme values of RR_{XY}^c . We ran 1000 simulations for each possible combination of k and $E[N]$, primarily assessing coverage of nominal 95% confidence intervals and secondarily assessing their precision (total width) and bias in $\widehat{p}(q)$.

For each study, we drew $N \sim \text{Unif}(150, 2E[N] - 150)$, using 150 as a minimum sample size to prevent model convergence failures, and drew the study's true effect size as $M^t \sim N(\mu^t, V^t)$. We simulated data for each subject under a model with a binary exposure ($X \sim \text{Bern}(0.5)$), a single binary unmeasured confounder, and a binary outcome. We set the two bias components equal to one another ($g = RR_{XU} = RR_{UY}$) and fixed $P(U = 1|X = 1) = 1$, allowing closed-form computation of:

$$P(U = 1|X = 0) = \frac{\exp(M^t)[1 + (g - 1)] - \exp(M^c)}{(g - 1)\exp(M^c)}$$

as in Ding & VanderWeele (2016). Within each stratum $X = x$, we simulated $U \sim$

Bern($P(U = 1|X = x)$). We simulated outcomes as $Y \sim \text{Bern}(\exp\{\log 0.05 + \log(g)U + M^t X\})$. Finally, we computed effect sizes and fit the random-effects model using the `metafor` package in R (Viechtbauer et al., 2010), estimating τ_c^2 per Paule & Mandel (1982) and $\widehat{\text{Var}}(\widehat{y}_R^c)$ with the Hartung & Knapp (2001) adjustment.

Results (Table 3) indicated approximately nominal performance for all combinations of k and $E[N]$, with precision appearing to depend more strongly on k than $E[N]$. As expected theoretically, $\widehat{p}(q)$ was approximately unbiased.

8. DISCUSSION

This paper has developed sensitivity analyses for unmeasured confounding in a random-effects meta-analysis of a relative risk outcome measure. Specifically, we have presented estimators for the proportion, $\widehat{p}(q)$, of studies with true effect sizes surpassing a threshold and for the minimum bias, $\widehat{T}(r, q)$, or confounding association strength, $\widehat{G}(r, q)$, in all studies that would be required to reduce to a threshold the proportion of studies with effect sizes less than q . Such analyses quantify the amount of confounding bias in terms of intuitively tractable sensitivity parameters. Computation of $\widehat{p}(q)$ uses two sensitivity parameters, namely the mean and variance across studies of a joint bias factor on the log-relative risk scale. Estimators $\widehat{T}(r, q)$ and $\widehat{G}(r, q)$ make reference to, and provide conclusions for, a single sensitivity parameter, chosen as either the common joint bias factor across studies or the strength of confounding associations on the relative risk scale. These methods assume that the bias factor is normally distributed or fixed across studies, but do not make further assumptions regarding the nature of unmeasured confounding.

Assessing sensitivity to unmeasured confounding is particularly important in meta-analyses of observational studies, where a central goal is to assess the current quality of evidence and to inform future research directions. If a well-designed meta-analysis yields a low value of $\widehat{T}(r, q)$ or $\widehat{G}(r, q)$ and thus is relatively sensitive to unmeasured confounding, this indicates that future research on the topic should prioritize randomized trials or designs and data collection that reduce unmeasured confounding. On the other hand, individual studies measuring moderate effect sizes with relatively wide confidence intervals may not, when considered individually, appear highly robust to unmeasured confounding; however, a meta-analysis aggregating their results may nevertheless suggest that a substantial proportion of the true effects are above a threshold of scientific importance even in the presence of some unmeasured confounding. Thus, conclusions of the meta-analysis may in fact be robust to moderate degrees of unmeasured confounding.

We focused on relative risk outcomes because of their frequency in biomedical meta-analyses and their mathematical tractability, which allows closed-form solutions with the introduction of only one assumption (on the distribution of the bias factor). To allow application of the present methods, an odds ratio outcome can be approximated as a relative risk if the outcome is rare. If the outcome is not rare, the odds ratio can be approximately converted to a relative risk by taking its square root; provided that the outcome probabilities are between 0.2 and 0.8, this transformation is always within 25% of the true relative risk (VanderWeele, in press). Comparable sensitivity analyses for other types of outcomes, such as mean differences, would require study-level summary measures (for example, of within-group means and variances) and in some cases would yield closed-form solutions only at the price of more stringent assumptions. Under the assumption of an underlying binary outcome with

high prevalence, such measures could be converted to log-odds ratios (Borenstein et al., 2009) and then to relative risks (VanderWeele, in press) as described above (see VanderWeele & Ding (2017)). It is important to note that, in circumstances discussed elsewhere (Tang, 2000; Thorlund et al., 2011), relative risk outcomes can produce biased meta-analytic estimates. When such biases in pooled point estimates or heterogeneity estimators are likely, sensitivity analyses will also be biased.

We operationalized “robustness to unmeasured confounding” as the proportion of true effects surpassing a threshold, an approach that focuses on the upper tail (for an apparently causative RR_{XY}^c) of the distribution of true effect sizes. Potentially, under substantial heterogeneity, a high proportion of true effect sizes could satisfy, for example, $RR_{XY}^t > 1.10$ while, simultaneously, a non-negligible proportion could be comparably strong in the opposite direction ($RR_{XY}^t < 0.90$). Such situations are intrinsic to the meta-analysis of heterogeneous effects, and in such settings, we recommend reporting the proportion of effect sizes below a symmetric threshold on the opposite side of the null (e.g., $\log 0.80$ if $q = \log 1.20$) both for the confounded distribution of effect sizes and for the distribution adjusted based on chosen bias parameters. For example, a meta-analysis that is potentially subject to unmeasured confounding and that estimates $\hat{y}_R^c = \log 1.15$ and $\tau_c^2 = 0.10$ would indicate that 45% of the effects RR_{XY}^c surpass 1.20, while 13% are less than 0.80. For a common $B^* = \log 1.10$ (equivalently, $g = 1.43$), we find that $\left(1 - \Phi\left(\frac{\log 1.20 - \log 1.15 + \log 1.10}{\sqrt{0.10}}\right)\right) \cdot 100\% = 33\%$ of the true effects surpass $RR_{XY}^c = 1.20$, while 20% are less than $RR_{XY}^c = 0.80$. More generally, random-effects meta-analyses could report the estimated proportion of effects above the null or above a specific threshold (along with a confidence interval for this proportion) as a continuous summary measure to supplement the standard pooled estimate and inference.

Together, these reporting practices could facilitate overall assessment of evidence strength and robustness to unmeasured confounding under effect heterogeneity.

The proposed sensitivity analyses in theory require only standard summary measures from a meta-analysis (namely, the estimated pooled effect and a heterogeneity estimator to compute point estimates, along with their estimated variances to compute inference), rather than study-level data. However, in practice, we find that reporting of τ_c^2 and $\widehat{\text{Var}}(\tau_c^2)$ is sporadic in the biomedical literature. Besides their utility for conducting sensitivity analyses, we consider τ_c^2 and $\widehat{\text{Var}}(\tau_c^2)$ to be inherently valuable to the scientific interpretation of heterogeneous effects. We therefore recommend that they be reported routinely for random-effects meta-analyses, even when related measures, such as the proportion of total variance attributable to effect heterogeneity (I^2), are also reported. To enable sensitivity analyses of existing meta-analyses that do not report the needed summary measures, the package `ConfoundedMeta` helps automate the process of obtaining and drawing inferences from study-level data from a published forest plot or table. The user can then simply fit a random-effects model of choice to obtain the required summary measures.

Our framework assumes that the bias factor is normally distributed or taken to be fixed across studies. Normality is approximately justified if, for example, RR_{XU} and RR_{UY} are approximately identically and independently normal with relatively small variance. Since RR_{UY} is in fact a maximum over strata of X and the range of U , future work could potentially consider an extreme-value distribution for this component, but such a specification would appear to require a computational, rather than closed-form, approach. Perhaps a more useful, conservative approach to assessing sensitivity to bias that may be highly skewed is to report $\widehat{T}(r, q)$ and $\widehat{G}(r, q)$ for a wide range of fixed values B^* , including those much larger than a

plausible mean.

An alternative sensitivity analysis approach would be to directly apply existing analytic bounds (Ding & VanderWeele, 2016) to each individual study in order to compute the proportion of studies with effect sizes more extreme than q given a particular bias factor. This has the downside of requiring access to study-level summary measures (rather than pooled estimates). Moreover, the confidence interval of each study may be relatively wide, such that no individual study appears robust to unmeasured confounding, while nevertheless a meta-analytic estimate that takes into account the distribution of effects may in fact indicate that some of these effects are likely robust. One could also alternatively conduct sensitivity analyses on the pooled point estimate itself, but such an approach is naïve to heterogeneity: when the true effects are highly variable, a non-negligible proportion of large true effects may remain even with the introduction of enough bias to attenuate the pooled estimate to a scientifically unimportant level.

In summary, our results have shown that sensitivity analyses for unmeasured confounding in meta-analyses can be conducted easily by extending results for individual studies. These methods are straightforward to implement through either our R package `ConfoundedMeta` or graphical user interface and ultimately help inform principled causal conclusions from meta-analyses.

REPRODUCIBILITY

All code required to reproduce the applied example and simulation study is publicly available (<https://osf.io/2r3gm/>).

SUPPLEMENTARY MATERIALS

Web Appendices referenced in Sections 4 and 6 are available with this paper at the *Biometrics* website on Wiley Online Library.

ACKNOWLEDGMENTS

This research was supported by National Defense Science and Engineering Graduate Fellowship 32 CFR 168a and NIH grant ES017876.

REFERENCES

- Borenstein, M., Hedges, L. V., Higgins, J., & Rothstein, H. R. (2009). *Introduction to meta-analysis*. Wiley Online Library.
- Cornfield, J., Haenszel, W., Hammond, E. C., Lilienfeld, A. M., Shimkin, M. B., & Wynder, E. L. (1959). Smoking and lung cancer: recent evidence and a discussion of some questions. *Journal of the National Cancer Institute*, *22*, 173–203.
- Ding, P., & VanderWeele, T. J. (2016). Sensitivity analysis without assumptions. *Epidemiology*, *27*(3), 368.
- Egger, M., Schneider, M., & Smith, G. D. (1998). Spurious precision? Meta-analysis of observational studies. *BMJ*, *316*(7125), 140.
- Flanders, W. D., & Khoury, M. J. (1990). Indirect assessment of confounding: graphic description and limits on effect of adjusting for covariates. *Epidemiology*, *1*(3), 239–246.

- Hartung, J., & Knapp, G. (2001). On tests of the overall treatment effect in meta-analysis with normally distributed responses. *Statistics in Medicine*, *20*(12), 1771–1782.
- IntHout, J., Ioannidis, J. P., & Borm, G. F. (2014). The Hartung-Knapp-Sidik-Jonkman method for random effects meta-analysis is straightforward and considerably outperforms the standard DerSimonian-Laird method. *BMC Medical Research Methodology*, *14*(1), 1.
- McCandless, L. C. (2012). Meta-analysis of observational studies with unmeasured confounders. *The International Journal of Biostatistics*, *8*(2), 368.
- Paule, R. C., & Mandel, J. (1982). Consensus values and weighting factors. *Journal of Research of the National Bureau of Standards*, *87*(5), 377–385.
- Schlesselman, J. J. (1978). Assessing effects of confounding variables. *American Journal of Epidemiology*, *108*(1), 3–8.
- Shrier, I., Boivin, J.-F., Steele, R. J., Platt, R. W., Furlan, A., Kakuma, R., . . . Rossignol, M. (2007). Should meta-analyses of interventions include observational studies in addition to randomized controlled trials? A critical examination of underlying principles. *American Journal of Epidemiology*, *166*(10), 1203–1209.
- Tang, J.-L. (2000). Weighting bias in meta-analysis of binary outcomes. *Journal of Clinical Epidemiology*, *53*(11), 1130–1136.
- Thorlund, K., Imberger, G., Walsh, M., Chu, R., Gluud, C., Wetterslev, J., . . . Thabane, L. (2011). The number of patients and events required to limit the risk of overestimation of intervention effects in meta-analysis: a simulation study. *PLoS One*, *6*(10), e25491.

- Trock, B. J., Hilakivi-Clarke, L., & Clarke, R. (2006). Meta-analysis of soy intake and breast cancer risk. *Journal of the National Cancer Institute*, *98*(7), 459–471.
- Turner, R. M., Spiegelhalter, D. J., Smith, G., & Thompson, S. G. (2009). Bias modelling in evidence synthesis. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, *172*(1), 21–47.
- VanderWeele, T., & Ding, P. (2017). Sensitivity analysis in observational research: introducing the E-value. *Annals of Internal Medicine*, doi: 10.7326/M16-2607.
- VanderWeele, T. J. (in press). On a square-root transformation of the odds ratio for a common outcome. *Epidemiology*.
- Veroniki, A. A., Jackson, D., Viechtbauer, W., Bender, R., Bowden, J., Knapp, G., ... Salanti, G. (2015). Methods to estimate the between-study variance and its uncertainty in meta-analysis. *Research Synthesis Methods*.
- Viechtbauer, W., et al. (2010). Conducting meta-analyses in R with the metafor package. *Journal of Statistical Software*, *36*(3), 1–48.
- Welton, N., Ades, A., Carlin, J., Altman, D., & Sterne, J. (2009). Models for potentially biased evidence in meta-analysis using empirically based priors. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, *172*(1), 119–136.

TABLES AND FIGURES

Table 1: *Bounds on $\hat{p}(q)$ provided by homogeneous bias with an apparently causative or preventive pooled effect. $\hat{\mu}^t$ estimates μ^t and is equal to $\hat{y}_R^c - \mu_{B^*}$ for $\hat{y}_R^c > 0$ or $\hat{y}_R^c + \mu_{B^*}$ for $\hat{y}_R^c < 0$.*

	$q > \hat{\mu}^t$	$q < \hat{\mu}^t$
$\hat{y}_R^c > 0$	Upper bound	Lower bound
$\hat{y}_R^c < 0$	Lower bound	Upper bound

Figure 1: *Impact of varying degrees of unmeasured confounding bias on proportion of true relative risks < 0.90*

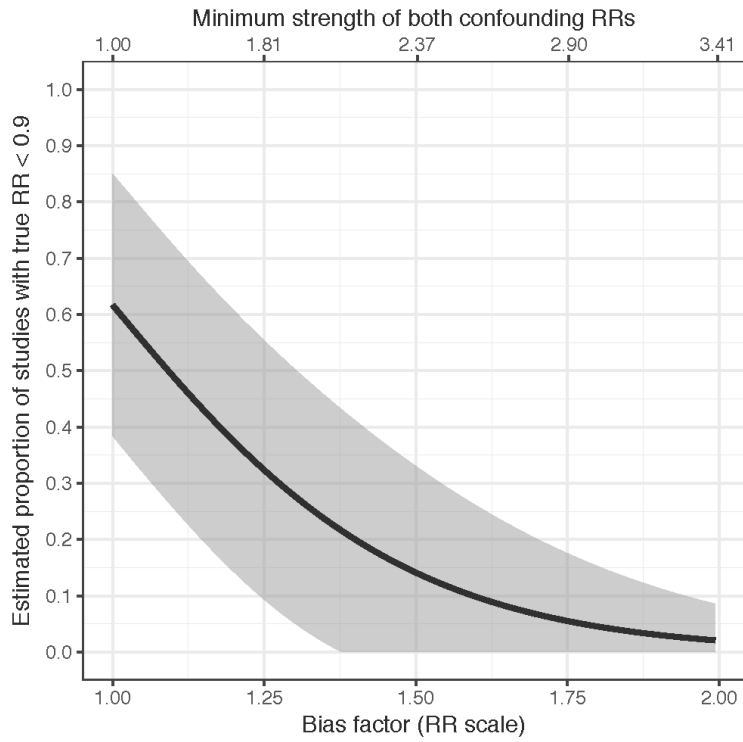


Table 2: $\hat{T}(r, q)$ and $\hat{G}(r, q)$ (in parentheses) for varying r and q . Blank cells indicate combinations for which no bias would be required.

		q	
r	0.70	0.80	0.90
0.1	1.27 (1.85)	1.45 (2.25)	1.63 (2.64)
0.2	1.10 (1.44)	1.26 (1.84)	1.42 (2.19)
0.3		1.14 (1.55)	1.29 (1.89)
0.4		1.05 (1.28)	1.18 (1.64)
0.5			1.09 (1.41)

Table 3: Point estimate bias, 95% confidence interval (CI) coverage, and 95% CI width for varying numbers of studies (k) and mean sample sizes within each study (Mean N).

k	Mean N	\hat{p} bias	CI coverage	CI width
15	300	0.030	0.968	0.572
25	300	0.034	0.976	0.452
50	300	0.031	0.967	0.315
200	300	0.028	0.929	0.154
15	500	0.022	0.967	0.524
25	500	0.022	0.977	0.408
50	500	0.025	0.974	0.283
200	500	0.024	0.934	0.140
15	1000	0.018	0.976	0.479
25	1000	0.016	0.976	0.370
50	1000	0.018	0.969	0.259
200	1000	0.015	0.970	0.129