

Mitigating Publication Bias Using Bayesian Stacking

Thomas A. Gibson

January 20, 2022

1 Introduction

Results from a meta-analysis may be skewed and unreliable in the presence of publication bias, where the publication or non-publication of a study depends on the statistical significance or magnitude of its results (Rothstein et al., 2006). Statistical methods for publication bias have been designed for sensitivity analysis, testing for the presence/magnitude of publication bias, and calculating bias-corrected parameter estimates. Most methods are either based on the funnel plot or selection models.

Methods based on the funnel plot (Light and Pillemer, 1984) – a scatterplot of effect sizes against their standard errors – inspect the plot’s asymmetry to test or correct for bias. Say we have S studies indexed by $i = 1, \dots, S$ and that y_i and $v_i = s_i^2$ are their estimated effect sizes and sampling variances. A popular non-parametric test for publication bias (Begg and Mazumdar, 1994) measures the rank correlation between standardized observed effect sizes y_i^* and the effect sizes’

variances v_i , where

$$\begin{aligned} y_i^* &= (y_i - \bar{y}) / (v_i^*)^{1/2} \\ \bar{y} &= \left(\sum_j (v_j^{-1}) / y_j \right) / \left(\sum_j v_j^{-1} \right) \\ v_i^* &= v_i - \left(\sum_j v_j^{-1} \right)^{-1}. \end{aligned}$$

Begg and Mazumdar (1994) measure the correlation between pairs (y_i^*, v_i) with Kendall’s tau, where a symmetric funnel plot would have correlation near zero. Egger’s test (Egger et al., 1997) fits a linear regression of observed standard normal deviates $\text{SND}_i = y_i / s_i$ against the inverse standard errors $1/s_i$, i.e. $\text{SND}_i = \alpha + \beta \times (1/s_i)$ with the null hypothesis $H_0 : \alpha = 0$. Other regression methods (Macaskill et al., 2001; Rücker et al., 2008; Thompson and Sharp, 1999; Peters et al., 2006) are similar to Egger’s test and use regression weights or transformations to improve upon Egger’s test in the presence of heterogeneity or for dichotomous outcomes (Jin et al., 2015). Lin and Chu (2018) develops a measure for the severity of publication bias based on the skewness of standardized deviates. The trim-and-fill method (Duval and Tweedie, 2000) calculates an adjusted mean effect estimate in a series of steps, by 1) estimating the number of missing studies k_0 , 2) “trimming” the smaller studies that are causing funnel plot asymmetry, 3) estimating the true mean effect with the remaining studies, and 4) replacing trimmed studies and their missing counterparts and re-estimating the mean effect and its variance. However, Duval and Tweedie (2000) recommend using trim-and-fill as a sensitivity analysis based on the *potential* number of missing studies, with general guidelines for sensitivity analysis given in Shi and Lin (2019). The trim-and-fill method is the only funnel plot-based method that offers an adjusted mean estimate, and it is not recommended if there is heterogeneity in study effects (Jin et al., 2015).

A second class of methods is based on *selection models*, first described in Hedges

(1984). Let Y be a random variable representing effect sizes for all studies in a population and let Θ be the parameters determining the sampling density $f(y; \Theta)$. Selection models assume a biased sampling scheme where the probability of a study being observed (published) is represented by a weight function $w(y; \lambda)$, where λ , a scalar or vector parameter, determines how certain studies may be more or less likely to be published. The weighted density for observed effect y_i is then

$$f^*(y_i; \Theta, \lambda) = \frac{f(y_i; \Theta)w(y_i; \lambda)}{\int f(y; \Theta)w(y; \lambda)dy} \quad (1.1)$$

and the likelihood function for all observed studies is

$$L(\Theta, \lambda) = \prod_{i=1}^S f^*(y_i; \Theta, \lambda). \quad (1.2)$$

Some models explicitly model the probability of publication for individual studies as a function of their p -values (Iyengar and Greenhouse, 1988; Hedges, 1992; Givens et al., 1997; Vevea and Hedges, 1995) or as a function of both the effect size and standard error (Copas, 1999; Copas and Shi, 2000, 2001). Hedges (1992) introduced stepped weight functions by dividing the unit interval $[0, 1]$ into K segments with $K - 1$ cut points, where studies that have p -values in different segments have different probabilities of publication. We refer to stepped selection functions by the number of cut points, i.e. a 1-step selection function might have a single cut point at $p = 0.05$, or a 2-step function might have cuts at $p = 0.05, 0.10$. Earlier selection models were recommended for bias-corrected effect size estimates, and were later recommended only for sensitivity analyses because of identifiability issues in smaller meta-analyses (Veeva and Woods, 2005; Jin et al., 2015). Sensitivity analyses use a grid representing varying levels of publication bias and estimate the mean effect under each assumed scenario. If results do not change much under an assumption of severe publication bias they are robust, and if results do change under an assumption of mild publication bias they are sensitive. Bayesian implementations of the

Copas selection model (Mavridis et al., 2013; Bai et al., 2020) have again allowed for estimation of mean effect sizes.

Recent approaches to mitigating publication bias have used Bayesian model averaged meta-analysis (BMA-MA) to consider a set of potential selection functions. Guan and Vandekerckhove (2016) considers four different selection functions based on p -values, including a no-bias model, an extreme-bias model where studies with p -values $p > \alpha$ are never published, a 1-step function where studies with $p > \alpha$ are published with some probability $\pi < 1$, and a model where the probability of publication decreases exponentially with p . Guan and Vandekerckhove (2016) only implement the models in a fixed-effects framework. Maier et al. (2020) evaluates a set of 12 models, using a $2 \times 2 \times 2$ factorial design with fixed/random effects, a true null/alternative hypothesis, and the presence/absence of publication bias. Maier et al. (2020) fit two-step and three-step selection functions based on p -values when publication bias is assumed, where the probability of publication changes at $p = 0.05$ (one-step) or at both $p = 0.05$ and $p = 0.10$ (two-step).

Bayesian model averaging (BMA) effectively assumes that one of the considered models is the “true” model, which is called the \mathcal{M} -closed setting (Bernardo and Smith, 2009). BMA does not perform as well under the \mathcal{M} -complete or \mathcal{M} -open settings, where the true data generating mechanism is too complex to implement or to put into a probabilistic framework (Bernardo and Smith, 2009; Le and Clarke, 2017). Multiple issues arise for BMA in these settings, including (a) the need to specify prior model probabilities, which makes little sense when we know the true model is not in our list, and (b) the model weights from BMA will converge to 1 for the model “closest” to the true model in terms of Kullback-Leibler divergence, and 0 for all others (Clyde and Iversen, 2013). Bayesian stacking of predictive distributions (Yao et al., 2018, 2021) is a method that outperforms Bayesian model

averaging in the \mathcal{M} -complete and \mathcal{M} -open settings and avoids issues (a) and (b) above. If we have K candidate models M_1, \dots, M_K , Yao et al. (2018) solve for model weights $w = (w_1, \dots, w_K)$ under the constraint $w \in \mathcal{S}_1^K$ where $\mathcal{S}_1^K = \{w \in [0, 1]^K : \sum_{k=1}^K w_k = 1\}$. They do this by solving

$$(\hat{w}_1, \dots, \hat{w}_K) = \max_{w \in \mathcal{S}_1^K} \frac{1}{S} \sum_{i=1}^n \log \sum_{k=1}^K w_k p(y_i | y_{-i}, M_k) \quad (1.3)$$

where $p(y_i | y_{-i}, M_k)$ is the leave-one-out (LOO) posterior predictive density with the i th data point left out evaluated at y_i . It would be computationally costly to refit each model M_k S times, so LOO densities $p(y_i | y_{-i}, M_k)$ are approximated using Pareto-smoothed importance sampling (Vehtari et al., 2017). We go through Bayesian stacking more thoroughly in Section 2.3.

Given that the true data generating mechanism for publication bias is likely much too complex to be specified in a simple selection model, we propose using Bayesian stacking to mitigate publication bias by fitting multiple Bayesian selection models and stacking over them. Assumed patterns of publication bias that poorly predict the observed data with LOO cross validation will be given little weight. We propose stacking over multiple types of models, including step functions (Vevea and Hedges, 1995) and Bayesian Copas selection models Mavridis et al. (2013); Bai et al. (2020). Section 2 describes relevant selection models for publication bias, Bayesian stacking, and how to implement Bayesian stacking of selection models. We then describe and summarize a simulation study in Section 3. The purpose of the simulation is to compare a stacked estimate of the mean effect size to estimates from individual selection models when the true data generator is not one of the fitted selection models. We apply the stacked model to a real dataset on the effectiveness of antidepressants in Section 4

2 Methods

We are doing a meta-analysis and have S studies indexed by $i = 1, \dots, S$, and each study provides an estimated effect y_i and an associated standard error s_i . Assuming estimates y_i are normally distributed, we calculate study i 's 2-sided p -value as $p_i = 2 \times (1 - \Phi(|y_i|/s_i))$ where $\Phi(\cdot)$ is the standard normal cumulative distribution function. The data model for each selection method is

$$y_i = \theta_i + s_i \epsilon_i \quad (2.1)$$

$$\theta_i | \theta, \tau^2 \sim N(\theta, \tau^2) \quad (2.2)$$

where θ_i represent random study effects normally distributed around global mean θ with variance τ^2 , and $\epsilon_i \sim N(0, 1)$. For stepped weight functions we combine s_i^2 and τ^2 into a single residual so that

$$y_i | \theta, \tau^2 \sim N(\theta, s_i^2 + \tau^2). \quad (2.3)$$

2.1 Selection models based on p-values

We define a stepped weight function $w(\cdot)$ similar to Vevea and Hedges (1995) and Vevea and Woods (2005), where $w(p)$ represents the probability that a study with p -value p is observed. Dividing the unit interval into K sub-intervals with descending endpoints a_k in which the weight function $w(p)$ is constant, we have that

$$w(p_i) = \begin{cases} \omega_1 & \text{if } a_1 < p_i < 1 \\ \omega_j & \text{if } a_j < p_i < a_{j-1} \\ \omega_K & \text{if } 0 < p_i < a_{K-1} \end{cases} \quad (2.4)$$

where $a_0 = 1$ and $a_K = 0$. Say $\boldsymbol{\omega} = (\omega_1, \dots, \omega_K)$. We then have that the likelihood contribution for each observation y_i given θ , τ^2 and weight function $w(\cdot)$ is

$$f(y_i|\theta, \tau^2, \boldsymbol{\omega}) = \frac{\phi(y_i; \theta, \tau^2 + s_i^2) \times w(p_i)}{\int \phi(x; \theta, \tau^2 + s_i^2) \times w(p(x, s_i)) dx}, \quad (2.5)$$

where $\phi(x; a, b)$ represents a normal probability density function with mean a and variance b , and $p(x, s_i) = 1 - \Phi(x/s_i)$ for one-sided p -values and $p(x, s_i) = 2(1 - \Phi(|x|/s_i))$ for two-sided p -values. Maier et al. (2020) place a “cumulative-Dirichlet” prior distribution on the weights $\boldsymbol{\omega}$, which effectively means placing a symmetric Dirichlet prior on an auxiliary parameter $\tilde{\boldsymbol{\omega}} \in (0, 1)^K$ and taking the cumulative sum

$$\begin{aligned} \tilde{\boldsymbol{\omega}} &\sim \text{Dirichlet}(\text{rep}(1, K)) \\ \boldsymbol{\omega} &= \text{cumulative-sum}(\tilde{\boldsymbol{\omega}}). \end{aligned} \quad (2.6)$$

This restricts $\boldsymbol{\omega}$ so that the K intervals in (2.4) have increasing probability of publication with decreasing p -values, and $\omega_K = 1$. The symmetric Dirichlet prior on $\tilde{\boldsymbol{\omega}}$ leads to prior means $(\frac{1}{K}, \frac{2}{K}, \dots, 1)$ for $\boldsymbol{\omega}$. Restricting $\omega_K = 1$ means each other ω_j represents the probability of publication for a study in interval j relative to the probability of publication in the lowest p -value interval. We consider a range of possible structures for the weight function $w(p)$ by varying both the number of intervals K and the choice of cut-points a_j for both one-sided and two-sided p -values.

2.2 Copas selection model

The Copas selection model (Copas, 1999; Copas and Shi, 2000, 2001) models the selection mechanism for publication as a function of study effect y_i and associated standard error s_i . We introduce a latent variable z_i modeled as

$$z_i = \gamma_0 + \frac{\gamma_1}{s_i} + \delta_i, \quad (2.7)$$

where z_i models the publication process such that study i is selected (published) only if $z_i > 0$. The parameter γ_0 controls the baseline probability of publication, and γ_1 defines the relationship between the observed standard deviation s_i and the probability of publication. Usually γ_1 is assumed to be positive, so that studies with smaller standard errors are more likely to be published. The random effects (ϵ_i, δ_i) are modeled as bivariate normal

$$\begin{pmatrix} \epsilon_i \\ \delta_i \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right) \quad (2.8)$$

where $\text{corr}(\epsilon_i, \delta_i) = \rho$ measures how the probability of selection changes with observed effect sizes.

Interpretation of the parameters $(\gamma_0, \gamma_1, \rho)$ can be difficult. The marginal probability that a study with standard error s_i is published is

$$P(z_i > 0 | s_i) = \Phi(\gamma_0 + \frac{\gamma_1}{s_i}).$$

Thus, if γ_0 is large and positive then all studies are published with high probability regardless of the value of s_i . We restrict γ_1 to be positive under the assumption that larger studies are more likely to be published for various reasons (e.g. more funding, quality of writing, etc.). Larger values of γ_1 lead to larger differences in publication probabilities for studies with different standard errors. The correlation parameter ρ is the main driver in how unadjusted estimates of θ are biased from the truth. If $\rho = 0$, then the selection process does not depend on observed effect sizes and unadjusted estimates are unbiased. Positive values of ρ indicate that observed effects y_i influence the selection process such that larger values of y_i (relative to their true mean θ_i) are being selected for, while negative values of ρ would show selection favoring larger negative values of y_i .

We consider two Bayesian adaptations of the Copas model (Mavridis et al., 2013;

Bai et al., 2020), which put prior distributions on all parameters including γ_0 and γ_1 . Mavridis et al. (2013) instead places priors on the lower and upper bounds for the probability of publication, P_{low} and P_{high} , as

$$\begin{aligned} P_{\text{low}} &\sim \text{Uniform}(L_1, L_2) \\ P_{\text{high}} &\sim \text{Uniform}(U_1, U_2), \end{aligned} \tag{2.9}$$

where (L_1, L_2) and (U_1, U_2) represent plausible ranges for the probability of publication for the studies with the largest and smallest standard errors, respectively. They then transform $(P_{\text{low}}, P_{\text{high}})$ to (γ_0, γ_1) with a 1-to-1 transformation using

$$\begin{aligned} \gamma_0 + \frac{\gamma_1}{s_{\text{max}}} &= \Phi^{-1}(P_{\text{low}}) \\ \gamma_0 + \frac{\gamma_1}{s_{\text{min}}} &= \Phi^{-1}(P_{\text{high}}) \end{aligned} \tag{2.10}$$

where s_{min} and s_{max} are the smallest and largest observed standard errors in the sample of S studies. Bai et al. (2020) place priors directly on γ_0 and γ_1 as

$$\begin{aligned} \gamma_0 &\sim \text{Uniform}(-2, 2) \\ \gamma_1 &\sim \text{Uniform}(0, s_{\text{max}}). \end{aligned} \tag{2.11}$$

They reason that this range of values allows for selection probabilities between 2.5% and 99.7% by restricting most of the mass for latent variables z_i to the range $(-2, 3)$. Priors (2.11) are meant to be default prior distributions, while (2.9) may require more problem-specific tuning, and the two priors lead to surprisingly different posterior distributions for the mean parameter θ .

2.3 Bayesian stacking

We use \mathcal{M} -open to refer to the setting in which our list of candidate models does not include the true data generating mechanism (Bernardo and Smith, 2009). Bayesian stacking (Yao et al., 2018) is an alternative to Bayesian model averaging that has superior performance in the \mathcal{M} -open setting. If we have K candidate models M_k

indexed by $k = 1, \dots, K$, the goal is to find the set of optimal weights $w \in \mathcal{S}_1^K$, $\mathcal{S}_1^K = \{w \in [0, 1]^K : \sum_{k=1}^K w_k = 1\}$, that maximizes a score S comparing the predictive distributions $p_k(\tilde{y}|y, M_k)$ to the true distribution $p_T(\tilde{y}|y)$. Formally, Yao et al. (2018) define the stacking problem as

$$\max_{w \in \mathcal{S}_1^K} \left(\sum_{k=1}^K w_k p(\cdot|y, M_k), p_T(\cdot|y) \right) \quad (2.12)$$

where the first argument P in $S(P, Q)$ is a weighted sum of model-specific posterior predictive distributions. Equation (2.12) is intractable as written, so Yao et al. (2018) replace the “true” predictive distribution $p_T(\tilde{y}|y)$ with observed values y_i , and replace model k ’s predictive distribution $p(\tilde{y}|y, M_k)$ with its corresponding leave-one-out (LOO) predictive distributions $\hat{p}_{k,-i}(y_i) = \int p(y_i|\theta_k, M_k)p(\theta_k|y_{-i}, M_k)d\theta_k$, where θ_k are the parameters in model k and subscript $-i$ denotes the data y with observation i left out. The authors recommend the logarithmic scoring rule, which reduces the stacking problem to solving for weights w with

$$(\hat{w}_1, \dots, \hat{w}_K) = \arg \max_{w \in \mathcal{S}_1^K} \frac{1}{n} \sum_{i=1}^n \log \sum_{k=1}^K w_k \hat{p}_{k,-i}(y_i) \quad (2.13)$$

via optimization. After solving for stacking weights \hat{w} , the *stacked posterior distribution* of a common parameter θ can be obtained by taking $\hat{w}_k \times T$ samples from each model k ’s posterior distribution $p(\theta|y)$ for a total of T posterior samples.

2.3.1 Pareto smoothed importance sampling

It would be computationally costly to refit each model n times to obtain LOO distributions $p_k(y_i|y_{-i}, M_k)$. Instead, Yao et al. (2018) use Pareto smoothed importance sampling (PSIS) (Vehtari et al., 2017) to obtain approximations. If we have draws $\theta_k^{(t)}$ from the full posterior $p(\theta_k|y, M_k)$, we can calculate importance ratios as

$$r_{i,k}^{(t)} = \frac{1}{p(y_i|\theta_k^{(t)}, M_k)} \propto \frac{p(\theta_k^{(t)}|y_{-i}, M_k)}{p(\theta_k^{(t)}|y, M_k)} \quad (2.14)$$

and use $r_{i,k}^{(t)}$ to generate model k 's importance sampling leave-one-out (IS-LOO) distribution

$$p_k(\tilde{y}_i|y_{-i}) \approx \frac{\sum_{t=1}^T r_{i,k}^{(t)} p(\tilde{y}_i|\theta_k^{(t)}, M_k)}{\sum_{t=1}^T r_{i,k}^{(t)}}, \quad (2.15)$$

which we evaluate at the left out data point y_i to obtain the IS-LOO predictive density $p_k(y_i|y_{-i})$.

This approximation may be unstable because the importance ratios $r_{i,k}^{(t)}$ can have high or infinite variance. Vehtari et al. (2017) mitigate this issue by fitting a generalized Pareto distribution to the upper tail (top 20%) of importance ratios, which returns smoothed importance weights $w_{i,k}^{(t)}$ to be used in place of the original importance ratios $r_{i,k}^{(t)}$ in (2.15). The Pareto-smoothed importance sampling leave-one-out (PSIS-LOO) expected log-predictive density (elpd) for point i in model k can be calculated as

$$\text{elpd}_{i,k} = \log \left(\frac{\sum_{t=1}^T w_{i,k}^{(t)} p(y_i|\theta_k^{(t)}, M_k)}{\sum_{t=1}^T w_{i,k}^{(t)}} \right). \quad (2.16)$$

We use the R package 'loo' (Vehtari et al., 2020) to implement PSIS-LOO, which requires only a $T \times n$ matrix of posterior pointwise log-likelihood samples. The loo package also solves for stacking weights w when given an $n \times K$ matrix of pointwise elpd estimates with one column for each model.

The estimated shape parameter \hat{k} from the generalized Pareto distribution can be used as a diagnostic for the reliability of PSIS-LOO approximations, where $\hat{k} > 0.7$ signals a potentially unreliable approximation. Vehtari et al. (2017) recommend sampling directly from $p_{k,-i}(y_i)$ (exact LOO) for problematic data points if the number of problematic y_i is small. Exact LOO requires fully refitting model k with the leave-one-out dataset y_{-i} and calculating the log-likelihood for the held out data point y_i . The exact LOO elpd for point i can then be combined with PSIS-LOO estimates for other data y_{-i} to solve for stacking weights w .

2.4 Stacking selection models for publication bias

It would be naive to think that either the stepped selection functions in Section 2.1 or the Copas models in Section 2.2 represent the true data generating mechanism for publication bias. As Bayesian stacking is designed to perform well in the event that our model list does not contain the true model, we propose stacking over both Copas models (Mavridis et al., 2013; Bai et al., 2020) and a variety of stepped selection functions to obtain a more robust posterior distribution for the mean parameter θ .

To fit the Copas models we rewrite model (??) - (??) as

$$\begin{pmatrix} y_i \\ z_i \end{pmatrix} \sim N \left(\begin{pmatrix} \theta \\ u_i \end{pmatrix}, \begin{pmatrix} \tau^2 + s_i^2 & \rho s_i \\ \rho s_i & 1 \end{pmatrix} \right) \mathbb{1}_{z_i > 0}, \quad (2.17)$$

which shows that we can first sample z_i from a truncated normal $z_i \sim N(u_i, 1) \mathbb{1}_{z_i > 0}$, and then sample $y_i | z_i \sim N(E[y_i | z_i], \text{Var}[y_i | z_i])$ where $E[y_i | z_i] = \theta + \rho s_i(z_i - u_i)$ and $\text{Var}[y_i | z_i] = \tau^2 + s_i^2(1 - \rho^2)$ (Mavridis et al., 2013). We also need to extract the log-likelihood for each observation i in order to stack models. The model construction (2.17) leads to a simple form for the log-likelihood

$$\begin{aligned} L(\theta, \tau^2, \rho, \gamma_0, \gamma_1) &= \sum_{i=1}^S \log[p(y_i | z_i > 0, s_i)] \\ &= \sum_{i=1}^S \log \left[\frac{p(z_i > 0 | y_i, s_i) f(y_i)}{p(z_i > 0 | s_i)} \right] \\ &= \sum_{i=1}^S \log(\phi(y_i; \theta, \tau^2)) - \log \Phi(u_i) + \log \Phi(v_i) \end{aligned} \quad (2.18)$$

where $\phi(y_i; \theta, \tau^2)$ represents the density at the point y_i of a normal distribution with mean θ and variance τ^2 , $\Phi(\cdot)$ represents the standard normal cumulative density

function, and

$$\begin{aligned} u_i &= \gamma_0 + \frac{\gamma_1}{s_i} \\ v_i &= \frac{u_i + \tilde{\rho}_i \frac{y_i - \theta}{\sqrt{\tau^2 + s_i^2}}}{\sqrt{(1 - \tilde{\rho}_i^2)}} \\ \tilde{\rho}_i &= \frac{s_i}{(\tau^2 + s_i^2)^{1/2}} \rho. \end{aligned}$$

While Bai et al. (2020) use default prior distributions (2.11) for the parameters γ_0 and γ_1 , Mavridis et al. (2013) instead advise researchers to use expert elicitation or historical data to specify (L_1, L_2) and (U_1, U_2) in (2.9) as plausible bounds for the lower and upper probabilities of publication. To avoid the need for strictly informative prior values, we specify $(L_1, L_2) = (0, 0.5)$ and $(U_1, U_2) = (0.5, 1)$, meaning we believe the lower bound for publication probability is between 0-50%, and the upper bound is between 50-100%. Mavridis et al. (2013) gives τ a half-normal prior $\tau \sim N(0, 10^2) \mathbb{1}_{\tau > 0}$, while Bai et al. (2020) uses a half-Cauchy prior $\tau \sim \text{Cauchy}(0, 1) \mathbb{1}_{\tau > 0}$. We fit the two Copas models in JAGS (Plummer et al., 2003).

As a default for the stepped models described in Section 2.1, we fit six different models. These include

1. Two-sided selection, p -value cutoff at 0.05
2. Two-sided selection, p -value cutoffs at 0.01 and 0.10
3. One-sided selection, p -value cutoff at 0.025
4. One-sided selection, p -value cutoffs at 0.025 and 0.5
5. One-sided selection, p -value cutoffs at 0.025 and 0.10
6. One-sided selection, p -value cutoffs at 0.005 and 0.05

We incorporate more one-sided than two-sided selection models under an assumption that one-sided selection is more likely in practice, as journals and researchers may select for results that match a specific directional effect. We fit the stepped selection models in Stan (Gelman et al., 2015) because of the ability to code the custom probability distribution (2.5), which also makes extraction of the log-likelihood matrix simple. For a given dataset we fit both Copas models and the six stepped models for a total of $K = 8$ models to stack over.

After optimizing stacking weights w , we sample from the stacked posterior distribution for θ by taking $w_k \times T$ samples with replacement from the posterior distribution of θ for model each model k and combining samples.

3 Simulation

The aim of the simulation is to assess how well the stacked model described in Section 2.4 *stacks up* against the individual models <please help me find a way to keep the pun in here> in estimating the true mean effect θ . We simulate data for meta-analyses using a $2 \times 2 \times 2 \times 4$ factorial design with

- overall mean effect $\theta = 0.1$ or 0.5
- two complicated selection functions
- a moderate or extreme level of selection bias for each selection function
- a small, medium, large, or very large amount of studies per meta-analysis

for 32 simulation scenarios. For each simulation scenario $j = 1, \dots, 24$ we set $\tau = 0.2$ and study-specific standard errors s_i , $i = 1, \dots, \tilde{S}_j$, are distributed as $\text{Uniform}(0.1, 0.8)$, where \tilde{S}_j is an *initial* number of studies per meta-analysis in scenario j . Study-specific true mean effects θ_i are distributed as $\theta_i \sim N(\theta_j, \tau^2)$, $i = 1, \dots, \tilde{S}_j$. We then

sample $y_i \sim N(\theta_i, s_i^2)$ as observed study effects. Each study i is assigned a probability of publication α_i based on one of two selection functions. The selection functions are chosen deliberately such that none of the stepped selection models or Copas models can individually capture the true selection mechanism. Whether or not each study i is included in the meta-analysis is then determined by an indicator $B_i \sim \text{Bernoulli}(\alpha_i)$, so that the number of studies S per meta-analysis is $S = \sum_{i=1}^{\tilde{S}_j} B_i$, which varies and is less than \tilde{S}_j . We choose \tilde{S}_j such that an average of $\sim 10, 20, 40$ and 80 studies survive the selection process and enter the meta-analysis for the small, medium, large, and very large simulation scenarios respectively.

We evaluate model performance using bias, 95% credible interval (CI) coverage, and root-mean squared error (RMSE).

3.1 Selection function 1

The first two selection functions are moderate (M) and extreme (E) versions of the same functional structure and are denoted f_{1M} and f_{1E} . The two functions decrease with one-sided p -values where $p_i = 1 - \Phi(y_i/s_i)$, and they have three steps at $p = 0.005, 0.2, 0.5$ with exponential decay between the three steps. Both functions are constant at $f_1(p) = 1$ for $p \in [0, 0.005)$. They then decrease exponentially for $p \in [0.005, 0.2)$ with $f_{1M}(p) = \exp(-0.5p)$ and $f_{1E}(p) = \exp(-2p)$. There is then a step at $p = 0.2$, and $f_{1M}(p) = \exp(-1p)$ and $f_{1E}(p) = \exp(-4p)$ for $p \in [0.2, 0.5)$. The functions are constant for $p \in [0.5, 1]$, with $f_{1M}(p) = 0.5$ and $f_{1E}(p) = 0.1$. See Figure 1 for a graph of both functions.

The proportion of studies to survive the selection mechanism depends on both the severity of selection (moderate or extreme) and the true mean θ . We sample an initial number of studies \tilde{S}_j , with \tilde{S}_j chosen so that an *average* of $S = 10, 20, 40$, or 80 studies survive and are included in the meta-analysis.

3.1.1 Results from simulation 1

We simulated $K = 200$ iterations for each of 16 scenarios. Figure 2 shows (absolute) bias for each model and simulation scenario. Left/right panels represent $\theta = 0.1$ or 0.5 , respectively, and top/bottom panels represent extreme or moderate selection. Stacking shows low bias for each sample size and combination of θ and selection severity, often having smaller bias than any individual model. Naturally, the standard model without any correction for publication bias has the largest bias in every scenario. The scenario with $\theta = 0.1$ and an extreme selection yields the largest bias across the board, with no model able to come close to reproducing the true mean.

RMSE for each scenario is shown in Figure 3. Stacking performs particularly well with extreme selection as sample sizes increase. The stacked posterior distribution of θ often has higher variance than most individual models, so even as it shows small biases it can be outperformed in terms of RMSE. However, models that have lower RMSE than stacking tend to have 95% coverage probabilities drop considerably as sample sizes increase, as shown in Figure 4. Stacking maintains coverage near the 95% nominal level in all scenarios except a) extreme selection with small true mean θ , and b) moderate selection with small θ *and* average sample size of 80. In (a), no model shows coverage probabilities near the nominal level, although stacking is among the closest for each sample size. For (b), all models except one-sided selection with steps at $p = 0.025, 0.5$ have coverage probabilities drop below 0.9. Overall, stacking tends to have low bias and 95% coverage probabilities closest to the nominal level.

4 Data analyses

We illustrate Bayesian stacking of selection models for publication bias on three datasets previously analyzed in the meta-analysis literature. In each example we fit the default $K = 8$ selection models and stack them to obtain the stacked posterior distribution, which we compare with results from the standard meta-analysis model.

4.1 Second-hand smoke and lung cancer

Hackshaw et al. (1997) analyzed a set of 37 studies measuring the effects of second-hand smoke on the likelihood of developing lung cancer. Individual studies measured the relative risk of lung cancer for women living with spouses who smoked vs women living with spouses who do not smoke. The authors performed a standard random-effects meta-analysis, finding a pooled relative risk (RR) of 1.24 (95% CI 1.13-1.36). The dataset from Hackshaw et al. (1997) has been used to illustrate publication bias models in the past (???)

We use log-relative risk as the main endpoint θ . Summaries for θ under each model are shown in Table 1, along with the standard model and stacked model. Three models contributed to the stacked posterior distribution: the Mavridis Copas model (weight = 0.63), one-sided stepped selection model with steps at $p = 0.025, 0.5$ (weight = .243), and the one-sided stepped selection model with steps at $p = 0.025, 0.1$ (weight = .127). Figure 5 shows posterior distributions for θ for each model. The stacked mean of θ is 0.108, with a 95% posterior interval of (-0.036, 0.255), which transforms to a mean relative risk of 1.11 (0.96, 1.29). Compared to the original results from Hackshaw et al. (1997), the stacked model estimates a mean increased risk of 11% rather than the 24% originally reported, with a much wider range of plausible values including the null value of 1.

4.2 Gender effects in grant proposals

Bornmann et al. (2007) compared the odds of a successful grant proposal for grants written by men compared to women using a dataset of 66 peer review procedures from 21 studies. Each study generally reported on one type of award and multiple peer review procedures for that award (e.g. TMR Marie Curie Fellowship for chemistry, engineering, mathematics, earth sciences, economics, physics and life sciences). Bornmann et al. (2007) fit a random effects meta-analysis with the log-odds ratio (logOR) as the main endpoint θ . The empirical Bayes estimate for the logOR of grant acceptance for men compared to women was 0.07 (95% CI 0.01-0.13), indicating a significant effect in favor of men. A funnel plot of study-specific effect sizes against their standard errors shows potential evidence of publication bias, so we fit the stacking procedure to the set of 66 results.

Table 2 shows model-specific point estimates and 95% intervals. The two models yielding the highest stacking weights were the Mavridis Copas model (weight = .345) and the Bai Copas model (weight = .655). Figure 6 shows posterior distributions of θ for the standard model, stacked model, and the two models contributing to the stacked model. The stacked posterior mean logOR was 0.051 with a 95% posterior interval of (-0.021, 0.119), still indicating a trend of grant proposals favoring men, but the 95% interval includes the null value of 0.

4.3 Recidivism and cognitive behavioral therapy

Landenberger and Lipsey (2005) analyzed a collection of 58 studies measuring how cognitive behavioral therapy (CBT) interventions are associated with recidivism in both adult and juvenile offenders. The authors fit a random effects meta-analysis model and report a mean odds ratio of 1.53 (logOR = 0.425) with $p < 0.001$. A

funnel plot of the 58 studies shows heavy asymmetry with skew towards positive results (favoring the CBT intervention).

Table 3 shows posterior summaries for the mean logOR for each model. The three models contributing to the stacked posterior were Mavridis Copas (weight = 0.142), Bai Copas (weight = 0.548), and one-sided selection with steps at $p = 0.025, 0.5$ (weight = 0.311). The stacked posterior mean logOR was 0.282 (95% CI 0.030, 0.469). Figure 7 shows posterior distributions for the standard model, stacked model, and the 3 models contributing to the stacked posterior. The stacked posterior logOR is shifted towards zero and is much more diffuse than the standard model.

5 Discussion

References

- Bai, R., Lin, L., Boland, M. R., and Chen, Y. (2020). A robust bayesian copas selection model for quantifying and correcting publication bias. *arXiv preprint arXiv:2005.02930*.
- Begg, C. B. and Mazumdar, M. (1994). Operating characteristics of a rank correlation test for publication bias. *Biometrics* pages 1088–1101.
- Bernardo, J. M. and Smith, A. F. (2009). *Bayesian theory*, volume 405. John Wiley & Sons.
- Bornmann, L., Mutz, R., and Daniel, H.-D. (2007). Gender differences in grant peer review: A meta-analysis. *Journal of Informetrics* **1**, 226–238.
- Clyde, M. and Iversen, E. S. (2013). Bayesian model averaging in the m-open framework. *Bayesian theory and applications* **14**, 483–498.

- Copas, J. (1999). What works?: selectivity models and meta-analysis. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* **162**, 95–109.
- Copas, J. and Shi, J. Q. (2000). Meta-analysis, funnel plots and sensitivity analysis. *Biostatistics* **1**, 247–262.
- Copas, J. and Shi, J. Q. (2001). A sensitivity analysis for publication bias in systematic reviews. *Statistical methods in medical research* **10**, 251–265.
- Duval, S. and Tweedie, R. (2000). Trim and fill: a simple funnel-plot-based method of testing and adjusting for publication bias in meta-analysis. *Biometrics* **56**, 455–463.
- Egger, M., Smith, G. D., Schneider, M., and Minder, C. (1997). Bias in meta-analysis detected by a simple, graphical test. *Bmj* **315**, 629–634.
- Gelman, A., Lee, D., and Guo, J. (2015). Stan: A probabilistic programming language for bayesian inference and optimization. *Journal of Educational and Behavioral Statistics* **40**, 530–543.
- Givens, G. H., Smith, D., and Tweedie, R. (1997). Publication bias in meta-analysis: a bayesian data-augmentation approach to account for issues exemplified in the passive smoking debate. *Statistical Science* **12**, 221–250.
- Guan, M. and Vandekerckhove, J. (2016). A bayesian approach to mitigation of publication bias. *Psychonomic bulletin & review* **23**, 74–86.
- Hackshaw, A. K., Law, M. R., and Wald, N. J. (1997). The accumulated evidence on lung cancer and environmental tobacco smoke. *Bmj* **315**, 980–988.

- Hedges, L. V. (1984). Estimation of effect size under nonrandom sampling: The effects of censoring studies yielding statistically insignificant mean differences. *Journal of Educational Statistics* **9**, 61–85.
- Hedges, L. V. (1992). Modeling publication selection effects in meta-analysis. *Statistical Science* **7**, 246–255.
- Iyengar, S. and Greenhouse, J. B. (1988). Selection models and the file drawer problem. *Statistical Science* pages 109–117.
- Jin, Z.-C., Zhou, X.-H., and He, J. (2015). Statistical methods for dealing with publication bias in meta-analysis. *Statistics in medicine* **34**, 343–360.
- Landenberger, N. A. and Lipsey, M. W. (2005). The positive effects of cognitive-behavioral programs for offenders: A meta-analysis of factors associated with effective treatment. *Journal of experimental criminology* **1**, 451–476.
- Le, T. and Clarke, B. (2017). A bayes interpretation of stacking for m-complete and m-open settings. *Bayesian Analysis* **12**, 807–829.
- Light, R. J. and Pillemer, D. B. (1984). Summing up: the science of reviewing research.
- Lin, L. and Chu, H. (2018). Quantifying publication bias in meta-analysis. *Biometrics* **74**, 785–794.
- Macaskill, P., Walter, S. D., and Irwig, L. (2001). A comparison of methods to detect publication bias in meta-analysis. *Statistics in medicine* **20**, 641–654.
- Maier, M., Bartoš, F., and Wagenmakers, E.-J. (2020). Robust Bayesian meta-analysis: Addressing publication bias with model-averaging.

- Mavridis, D., Sutton, A., Cipriani, A., and Salanti, G. (2013). A fully bayesian application of the copas selection model for publication bias extended to network meta-analysis. *Statistics in medicine* **32**, 51–66.
- Peters, J. L., Sutton, A. J., Jones, D. R., Abrams, K. R., and Rushton, L. (2006). Comparison of two methods to detect publication bias in meta-analysis. *Jama* **295**, 676–680.
- Plummer, M. et al. (2003). Jags: A program for analysis of bayesian graphical models using gibbs sampling. In *Proceedings of the 3rd international workshop on distributed statistical computing*, volume 124, pages 1–10. Vienna, Austria.
- Rothstein, H. R., Sutton, A. J., and Borenstein, M. (2006). *Publication bias in meta-analysis: Prevention, assessment and adjustments*. John Wiley & Sons.
- Rücker, G., Schwarzer, G., and Carpenter, J. (2008). Arcsine test for publication bias in meta-analyses with binary outcomes. *Statistics in Medicine* **27**, 746–763.
- Shi, L. and Lin, L. (2019). The trim-and-fill method for publication bias: practical guidelines and recommendations based on a large database of meta-analyses. *Medicine* **98**,.
- Thompson, S. G. and Sharp, S. J. (1999). Explaining heterogeneity in meta-analysis: a comparison of methods. *Statistics in medicine* **18**, 2693–2708.
- Vehtari, A., Gabry, J., Magnusson, M., Yao, Y., Bürkner, P.-C., Paananen, T., and Gelman, A. (2020). loo: Efficient leave-one-out cross-validation and waic for bayesian models. R package version 2.4.1.
- Vehtari, A., Gelman, A., and Gabry, J. (2017). Practical bayesian model evaluation

- using leave-one-out cross-validation and waic. *Statistics and computing* **27**, 1413–1432.
- Vevea, J. L. and Hedges, L. V. (1995). A general linear model for estimating effect size in the presence of publication bias. *Psychometrika* **60**, 419–435.
- Vevea, J. L. and Woods, C. M. (2005). Publication bias in research synthesis: sensitivity analysis using a priori weight functions. *Psychological methods* **10**, 428.
- Yao, Y., Pirš, G., Vehtari, A., and Gelman, A. (2021). Bayesian hierarchical stacking. *arXiv preprint arXiv:2101.08954* .
- Yao, Y., Vehtari, A., Simpson, D., and Gelman, A. (2018). Using stacking to average bayesian predictive distributions (with discussion). *Bayesian Analysis* **13**, 917–1007.

List of Figures

1	Selection mechanism 1: declining selection probabilities with increasing one-sided p-values. There are steps at $p = .2$ and $p = .5$ and exponential decay between cut-points. Selection probability is constant after $p = 0.5$	25
2	Absolute value of bias from 200 simulation replications using Selection mechanism 1. Left and right panels have $\theta = 0.1$ and $\theta = 0.5$, respectively. Top and bottom panels have extreme or moderate selection severity. The blue dashed line represents the stacked model.	26
3	RMSE from 200 simulation replications using Selection mechanism 1. Left and right panels have $\theta = 0.1$ and $\theta = 0.5$, respectively. Top and bottom panels have extreme or moderate selection severity. The blue dashed line represents the stacked model.	27
4	Proportion of 95% intervals covering the true mean θ from 200 simulation replications using Selection mechanism 1. Left and right panels have $\theta = 0.1$ and $\theta = 0.5$, respectively. Top and bottom panels have extreme or moderate selection severity. The blue dashed line represents the stacked model.	28
5	Posterior distributions of θ for each selection model using lung cancer data from Hackshaw et al. (1997). Colored lines show models where stacking weight was at least 0.01 (1%). Black line is the standard meta-analysis. Dashed line shows stacked posterior distribution.	29
6	Posterior distributions of θ for each selection model using grant application data from Bornmann et al. (2007). Colored lines show models where stacking weight was at least 0.01 (1%). Black line is the standard meta-analysis. Dashed line is stacked posterior distribution of θ	30
7	Posterior distributions of θ for each selection model using recidivism data from Landenberger and Lipsey (2005). Colored lines show models where stacking weight was at least 0.01 (1%). Black line is the standard meta-analysis. Dashed line shows stacked posterior distribution.	31

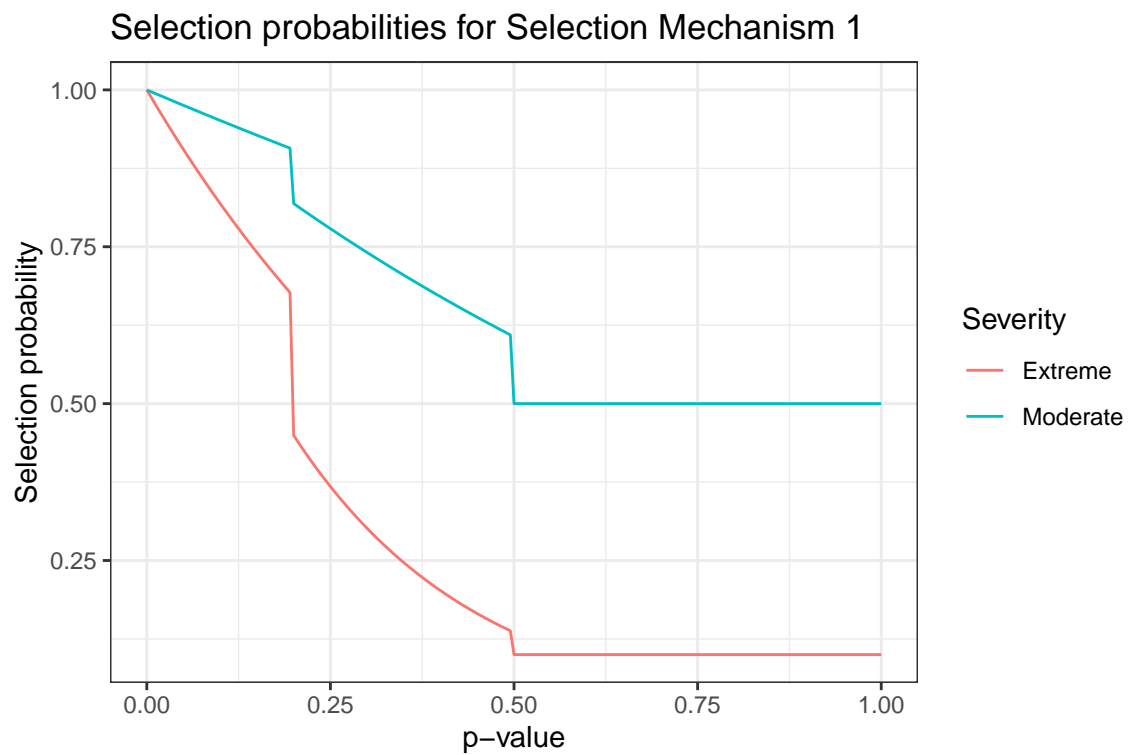


Figure 1: Selection mechanism 1: declining selection probabilities with increasing one-sided p-values. There are steps at $p = .2$ and $p = .5$ and exponential decay between cut-points. Selection probability is constant after $p = 0.5$.

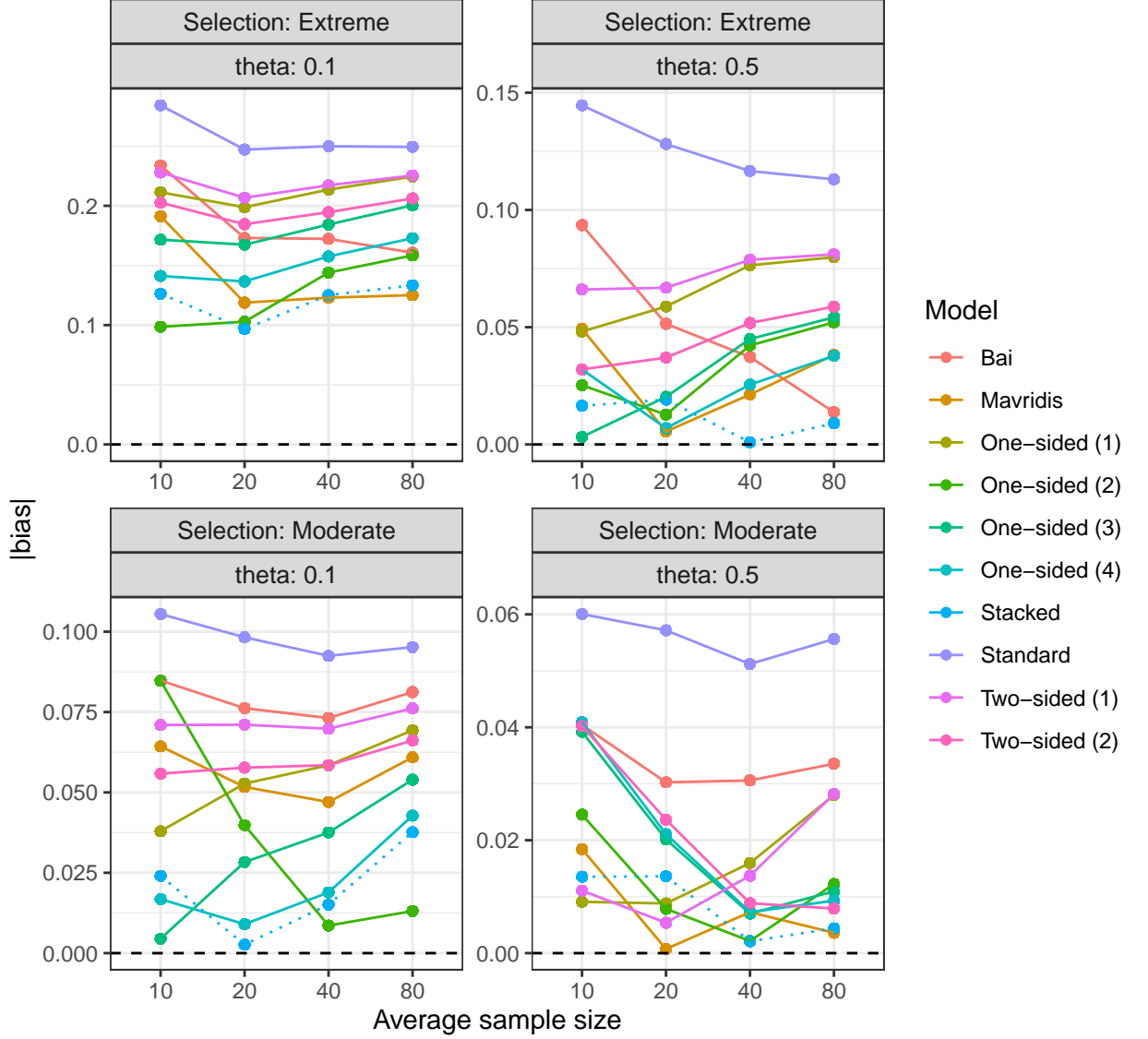


Figure 2: Absolute value of bias from 200 simulation replications using Selection mechanism 1.

Left and right panels have $\theta = 0.1$ and $\theta = 0.5$, respectively.

Top and bottom panels have extreme or moderate selection severity.

The blue dashed line represents the stacked model.

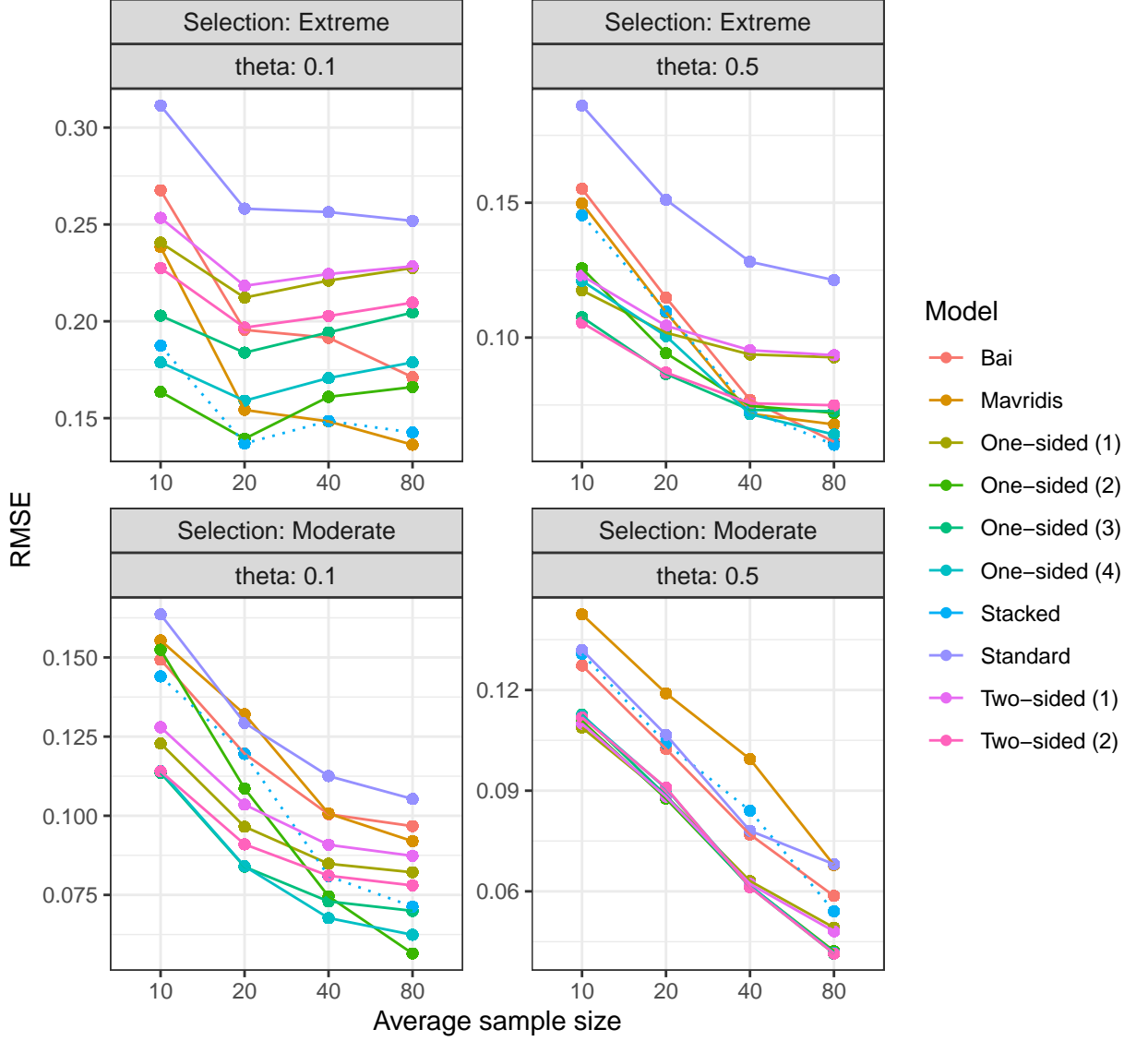


Figure 3: RMSE from 200 simulation replications using Selection mechanism 1. Left and right panels have $\theta = 0.1$ and $\theta = 0.5$, respectively. Top and bottom panels have extreme or moderate selection severity. The blue dashed line represents the stacked model.

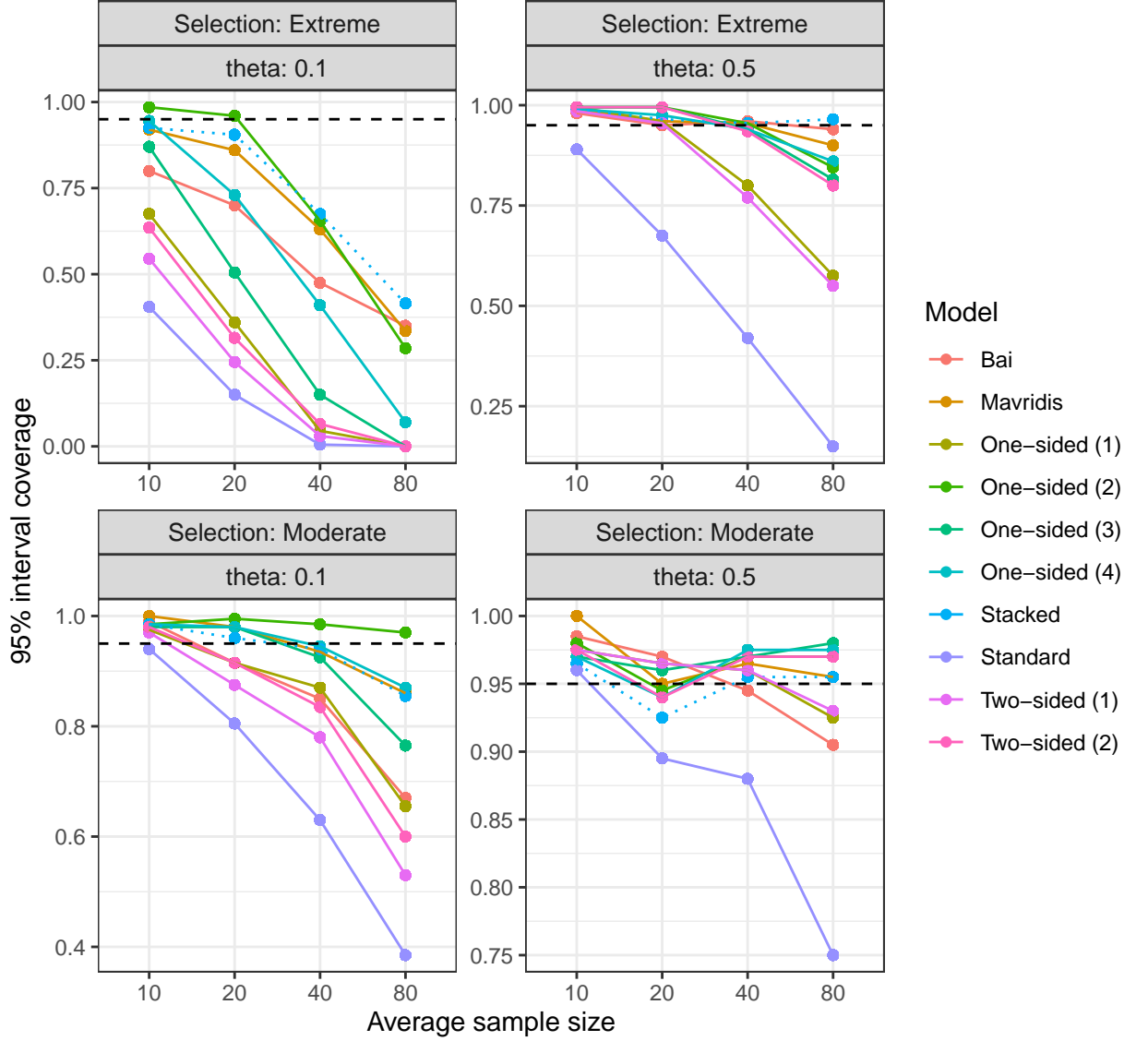


Figure 4: Proportion of 95% intervals covering the true mean θ from 200 simulation replications using Selection mechanism 1. Left and right panels have $\theta = 0.1$ and $\theta = 0.5$, respectively. Top and bottom panels have extreme or moderate selection severity. The blue dashed line represents the stacked model.

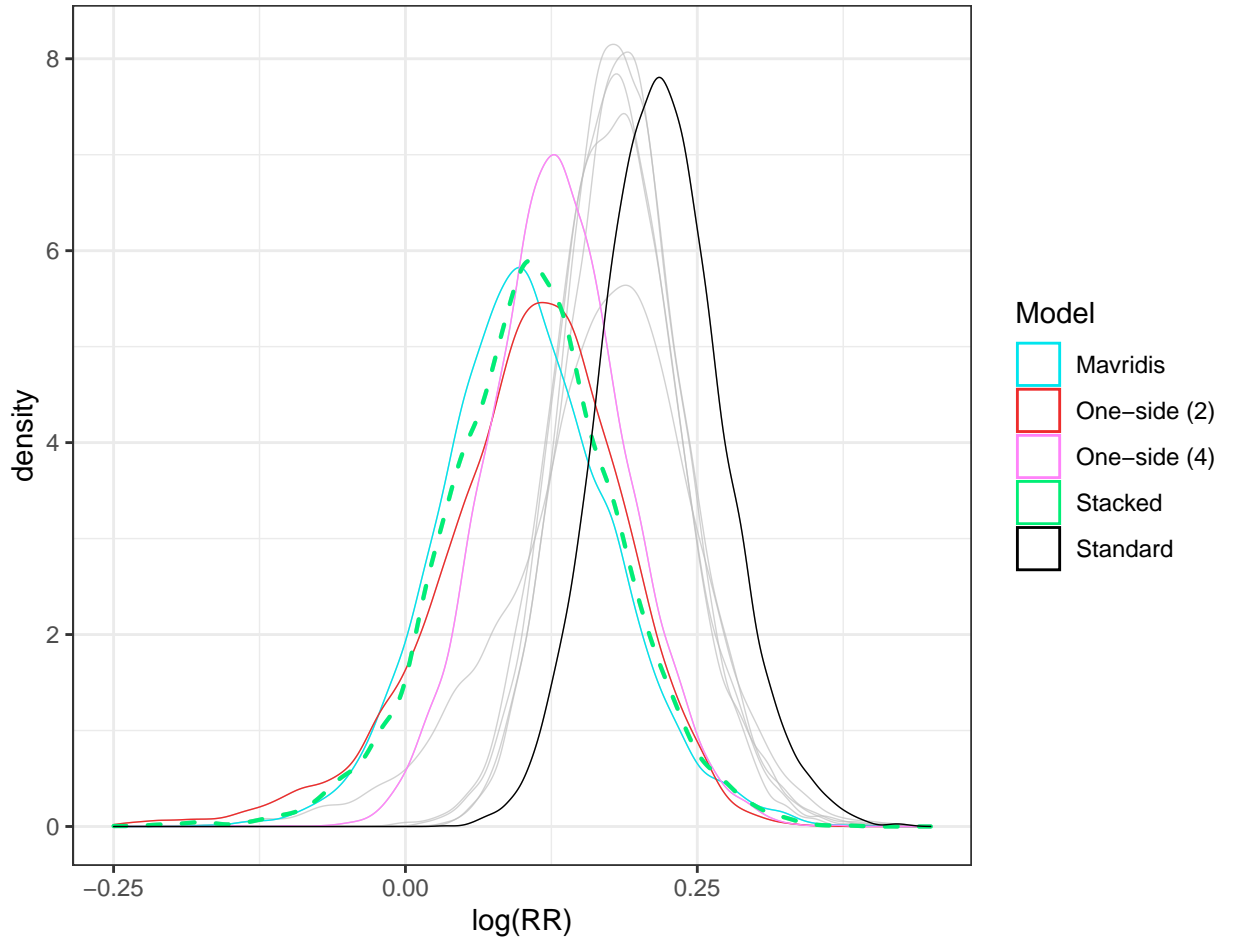


Figure 5: Posterior distributions of θ for each selection model using lung cancer data from Hackshaw et al. (1997). Colored lines show models where stacking weight was at least 0.01 (1%). Black line is the standard meta-analysis. Dashed line shows stacked posterior distribution.

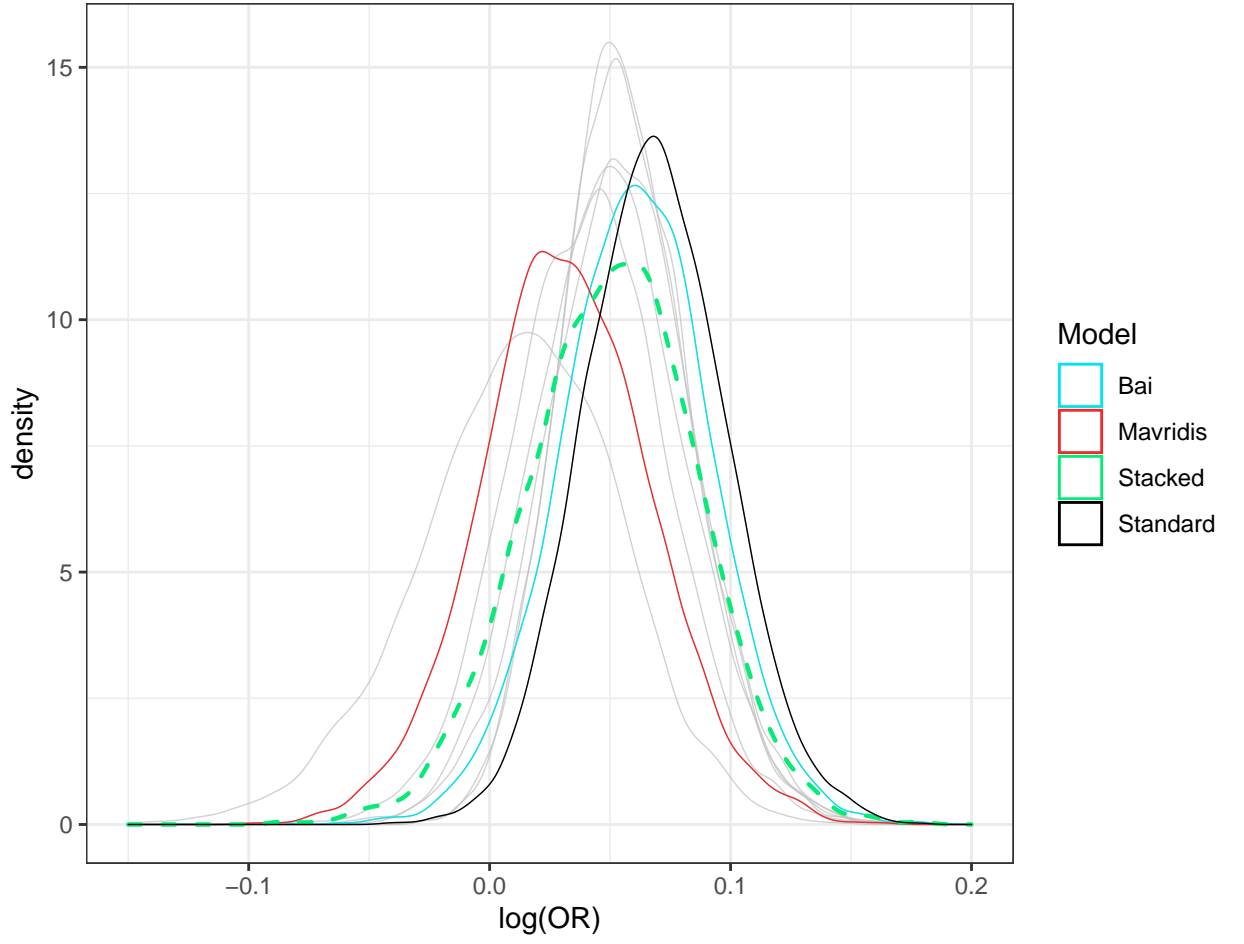


Figure 6: Posterior distributions of θ for each selection model using grant application data from Bornmann et al. (2007). Colored lines show models where stacking weight was at least 0.01 (1%). Black line is the standard meta-analysis. Dashed line is stacked posterior distribution of θ .

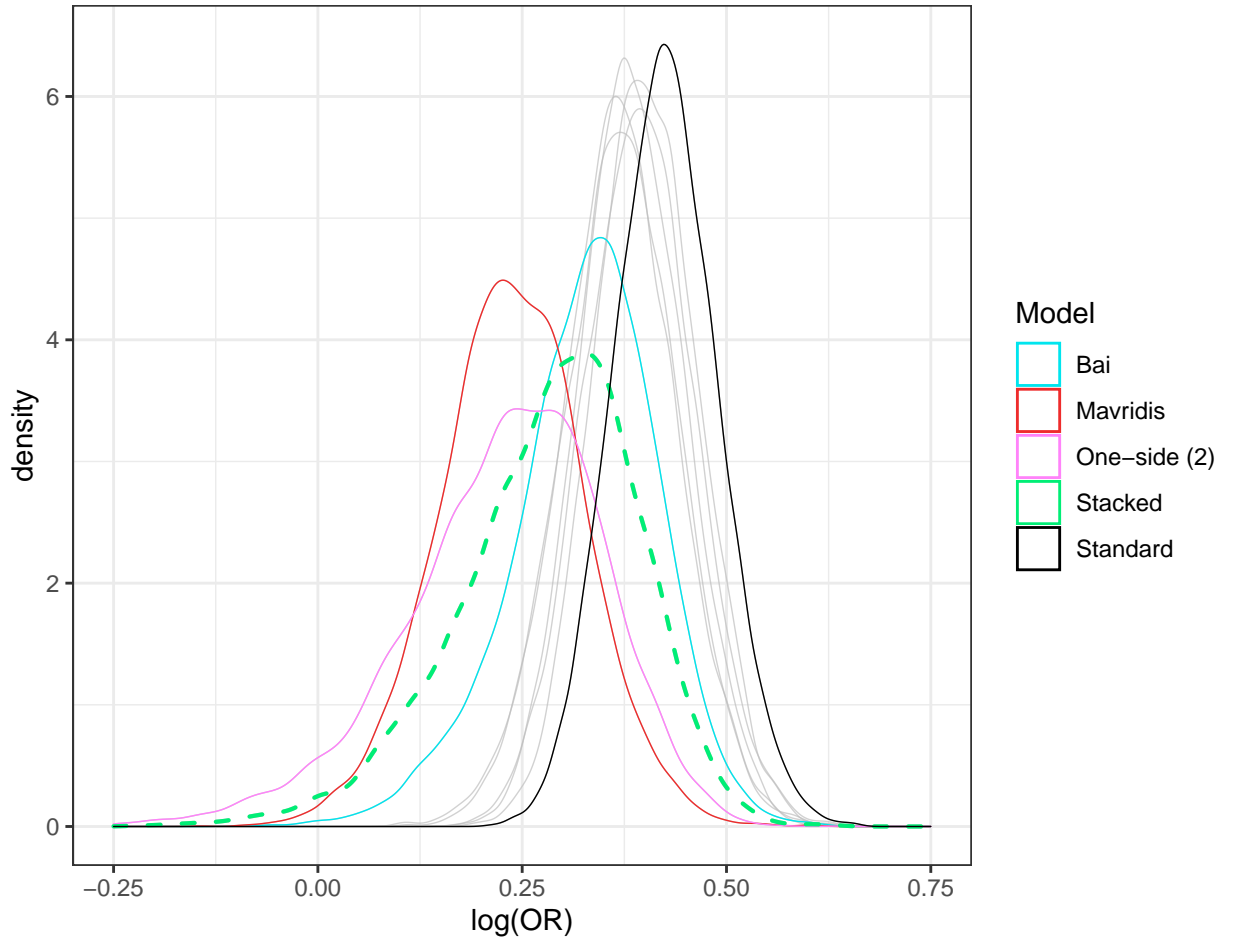


Figure 7: Posterior distributions of θ for each selection model using recidivism data from Landenberger and Lipsey (2005). Colored lines show models where stacking weight was at least 0.01 (1%). Black line is the standard meta-analysis. Dashed line shows stacked posterior distribution.

List of Tables

1	Posterior summaries for each model from numerical example 1 using data from Hackshaw et al. (1997). The stacked model has a drastically different posterior distribution for θ than the standard meta-analysis, with a mean closer to 0 and larger tails. Models contributing to the stacked posterior are the Mavridis Copas model and one-sided step models with steps at (.025, .5) and (.025, .1)	33
2	Posterior summaries for each model from numerical example 2 using data from Bornmann et al. (2007). While the standard analysis yields a 95% posterior credible interval (CI) excluding zero, the stacked posterior shifts the mean towards zero and the posterior CI includes zero. Models contributing to the stack are the Bai and Mavridis Copas models.	34
3	Posterior summaries for each model from numerical example 3 using data from Landenberger and Lipsey (2005). The stacked model heavily shifts the posterior distribution of θ towards zero. Models contributing to the stack are the Mavridis and Bai Copas models, and the one-sided step selection model with steps at (.025, .5).	35

Model	Mean	SD	2.5%	97.5%	Stacking Weight
Standard	0.219	0.052	0.122	0.327	—
Stacked	0.108	0.074	-0.044	0.253	—
Mavridis	0.103	0.073	-0.036	0.255	0.630
Bai	0.167	0.081	-0.016	0.309	0.000
Two-side (1)	0.190	0.052	0.093	0.297	0.000
Two-side (2)	0.186	0.049	0.093	0.286	0.000
One-side (1)	0.182	0.054	0.081	0.296	0.000
One-side (2)	0.105	0.082	-0.084	0.245	0.243
One-side (3)	0.183	0.053	0.085	0.294	0.000
One-side (4)	0.131	0.059	0.018	0.251	0.127

Table 1: Posterior summaries for each model from numerical example 1 using data from Hackshaw et al. (1997). The stacked model has a drastically different posterior distribution for θ than the standard meta-analysis, with a mean closer to 0 and larger tails. Models contributing to the stacked posterior are the Mavridis Copas model and one-sided step models with steps at (.025, .5) and (.025, .1)

Model	Mean	SD	2.5%	97.5%	Stacking Weight
Standard	0.069	0.030	0.012	0.128	—
Stacked	0.051	0.036	-0.021	0.119	—
Mavridis	0.032	0.035	-0.038	0.104	0.345
Bai	0.060	0.032	-0.003	0.122	0.655
Two-side (1)	0.057	0.027	0.007	0.112	0.000
Two-side (2)	0.056	0.027	0.007	0.112	0.000
One-side (1)	0.054	0.031	-0.007	0.115	0.000
One-side (2)	0.012	0.042	-0.075	0.090	0.000
One-side (3)	0.050	0.031	-0.010	0.111	0.000
One-side (4)	0.041	0.032	-0.023	0.104	0.000

Table 2: Posterior summaries for each model from numerical example 2 using data from Bornmann et al. (2007). While the standard analysis yields a 95% posterior credible interval (CI) excluding zero, the stacked posterior shifts the mean towards zero and the posterior CI includes zero. Models contributing to the stack are the Bai and Mavridis Copas models.

Model	Mean	SD	2.5%	97.5%	Stacking Weight
Standard	0.425	0.063	0.303	0.552	—
Stacked	0.282	0.112	0.030	0.469	—
Mavridis	0.236	0.089	0.058	0.408	0.142
Bai	0.326	0.090	0.124	0.485	0.548
Two-side (1)	0.402	0.063	0.283	0.525	0.000
Two-side (2)	0.385	0.064	0.264	0.515	0.000
One-side (1)	0.393	0.067	0.261	0.528	0.000
One-side (2)	0.227	0.123	-0.055	0.431	0.311
One-side (3)	0.369	0.069	0.231	0.507	0.000
One-side (4)	0.369	0.071	0.225	0.507	0.000

Table 3: Posterior summaries for each model from numerical example 3 using data from Landenberger and Lipsey (2005). The stacked model heavily shifts the posterior distribution of θ towards zero. Models contributing to the stack are the Mavridis and Bai Copas models, and the one-sided step selection model with steps at (.025, .5).