Select relevant moderators in meta-regression using Bayesian penalization

Caspar J. Van Lissa[1,2] & Sara van Erp[1]

[1] Utrecht University, dept. Methodology & Statistics

[2] Open Science Community Utrecht

Author Note

Add complete departmental affiliations for each author here. Each new line herein must be indented, like this line.

Enter author note here.

The authors made the following contributions. Caspar J. Van Lissa: Conceptualization, Writing - Original Draft Preparation, Programming front-end; Sara van Erp: Writing - Contributions and feedback, Programming back-end.

Correspondence concerning this article should be addressed to Caspar J. Van Lissa, Padualaan 14, 3584CH Utrecht, The Netherlands. E-mail: c.j.vanlissa@uu.nl

14                                            Abstract

15   One or two sentences providing a **basic introduction** to the field, comprehensible to a

16   scientist in any discipline.

17         Two to three sentences of **more detailed background**, comprehensible to scientists

18   in related disciplines.

19         One sentence clearly stating the **general problem** being addressed by this particular

20   study.

21         One sentence summarizing the main result (with the words "**here we show**" or their

22   equivalent).

23         Two or three sentences explaining what the **main result** reveals in direct comparison

24   to what was thought to be the case previously, or how the main result adds to previous

25   knowledge.

26         One or two sentences to put the results into a more **general context**.

27         Two or three sentences to provide a **broader perspective**, readily comprehensible to

28   a scientist in any discipline.

29         *Keywords:* keywords

30   Word count: X

Select relevant moderators in meta-regression using Bayesian penalization

Skeleton lasso/pema paper 1.) What is Meta-analysis? 2.) What is meta-regression and how does it complement meta-analysis? a.) Introduce Moderators b.) Study Heterogeneity + Random Sampling Error (and their difference)

2.5.) Fixed vs. random effects a.) Shortcomings of fixed effect models 3.) Shortcomings of current meta regressions w.r.t. estimating coefficients and heterogeneity: a.) Small sample size / overfitting b.) Non-normal data 4.) Various methods to estimate heterogeneity (and coefficients) a.) The use of WLS and REML 5.) Intro to Frequentist linear methods/Bayesian methods and Random forests, along with their (dis-)advantages: a.) Rma: uses WLS for estimation b.) MetaForest (Random Effects): Uses Random Forest Algorithm c.) Lasso Pema: uses penalized Lasso d.) Horseshoe Pema: uses horseshoe priors 6.) Goal of the current study 7.) Means of attaining goal and evaluation of performance: a.) simulation study b.) algorithmic performance c.) design factors d.) Impact of design factors on algorithim performance e.) Hypotheses of algorithmic performances

colour coding: I colour coded the text as to know from which file the text is copied GREEN = derived from 'Thesis_Metaforest' BLUE = Thesis_lasso BLACK = internship_report RED = Inserted myself

## Introduction

Meta-analysis is a quantitative form of evidence synthesis, whereby effect sizes from multiple similar studies are aggregated. In its simplest form, this aggregation consists of the computation of a summary effect as a weighted average of the observed effect sizes. This average is weighted to account for the fact that some observed effect sizes are assumed to be more informative about the underlying population effect. Each effect size is assigned a weight that determines how influential it is in calculating the summary effect. This weight is based on specific assumptions; for example, the *fixed effect* model assumes that

all observed effect sizes reflect one underlying true population effect size. This assumption is well-suited to the situation where effect sizes from close replication studies are meta-analyzed (**higgins_re-evaluation_2009?**, fabrigar_conceptualizing_2016, maxwell_is_2015). The *random effects* model, by contrast, assumes that population effect sizes follow a normal distribution. Each observed effect size provides information about the mean and standard deviation of this distribution of population effect sizes. This assumption is more appropriate when studies are conceptually similar and differences between them are random (**higgins_re-evaluation_2009?**, fabrigar_conceptualizing_2016, maxwell_is_2015).

Not all heterogeneity in effect sizes is random, however. Quantifiable between-study differences may introduce systematic heterogeneity. Such between-study differences are known as "moderators." For example, if studies have been replicated in Europe and the Americas, this difference can be captured by a binary moderator called "continent." Alternatively, if studies have used different dosages of the same drug, this may be captured by a continuous moderator called "dosage." Systematic heterogeneity in the observed effect sizes can be accounted for using *meta-regression* (Viechtbauer & López-López, 2015). This technique provides estimates of the effect of one or more study characteristics on the overall effect size, as well as of the overall effect size and residual heterogeneity after controlling for their influence.

One common application of meta-analysis is to summarize existing bodies of literature. In such situations, the number of moderators is often relatively high because similar research questions have been studied in different laboratories, using different methods, instruments, and samples. Each of these between-study differences could be coded as a moderator, and some of these moderators may explain systematic heterogeneity.

It is theoretically possible to account for the influence of multiple moderators using meta-regression. However, like any regression-based approach, meta-regression requires a

82 relatively high number of cases (studies) per parameter obtain sufficient power to examine

83 heterogeneity. In practice the number of available studies is often too low to examine

84 heterogeneity reliably (Riley, Higgins, & Deeks, 2011). At the same time, there are many

85 potential sources of heterogeneity, as similar research questions are studied in different

86 laboratories, using different methods, instruments, and samples. This leads to a problem

87 known as the "curse of dimensionality": the number of candidate moderators is large

88 relative to the number of cases in the data. Such cases do not fit comfortably into the

89 classic meta-analysis paradigm, which, like any regression-based approach, requires a high

90 number of cases per parameter. Between-studies thus presents a non-trivial challenge to

91 data aggregation using classic meta-analytic methods. At the same time, it also offers an

92 unexploited opportunity to learn which differences between studies have an impact on the

93 effect size found, if adequate exploratory techniques can be developed.

94 Addressing the curse of dimensionality necessitates *variable selection*: the selection of

95 a smaller subset of relevant moderators from a larger number of candidate moderators.

96 One way to perform variable selection is by relying on theory. However, in many fields of

97 science, theories exist at the individual level of analysis (e.g., in social science, at the level

98 of individual people). These theories do not necessarily generalize to the study level of

99 analysis. Using theories at the individual level for moderator selection at the study level

100 amounts to committing the ecological fallacy: generalizing inferences across levels of

101 analysis (**jargowskyEcologicalFallacy2004?**). To illustrate what theory at the study

102 level of analysis might look like, consider the so-called *decline effect*. It is a phenomenon

103 whereby effect sizes in a particular tranche of the literature seem to diminish over time

104 (**schoolerUnpublishedResultsHide2011?**). It has been theorized that the decline effect

105 can been attributed to regression to the mean: A finding initially draws attention from the

106 research community because an anomalously large effect size has been published, and

107 subsequent replications find smaller effect sizes. Based on the decline effect, we might thus

108 expect the variable "year of publication" to be a relevant moderator of study effect sizes.

109 Note that this prediction is valid even if year is orthogonal to the outcome of interest

110 within each study. Until more theory about the drivers of between-study heterogeneity is

111 developed, however, this approach will have limited utility for variable selection.

112     An alternative solution is to rely on statistical methods for variable selection. This is

113 a focal issue in the discipline of machine learning

114 (**hastieElementsStatisticalLearning2009?**). One technique that facilitates variable

115 selection is *regularization*: shrinking model parameters towards zero, such that only larger

116 parameters remain. Although this technique biases the parameter estimates, it also reduces

117 their variance, which has the advantage of producing more generalizable results that make

118 better predictions for new data (see **hastieElementsStatisticalLearning2009?**). This

119 paper introduces *Bayesian regularized meta-regression* (BRMA), an algorithm that uses

120 Bayesian estimation with regularizing priors to perform variable selection in meta-analysis.

121 The algorithm is implemented in the function `brma()` in the R-package `pema`.

## Statistical underpinnings

123     To understand how BRMA estimates the relevant parameters and performs variable

124 selection, it is instructional to first review the statistical underpinnings of the

125 aforementioned classic approaches to meta-analysis. First is the fixed-effect model, which

126 assumes that each observed effect size $T_i$ is an estimate of an underlying true effect size $\Theta$

127 (**hedgesFixedRandomeffectsModels1998?**). The only cause of heterogeneity in

128 observed effect sizes is presumed to be effect size-specific sampling variance, $v_i$, which is

129 treated as known, and computed as the square of the standard error of the effect size.

130 Thus, for a collection of $k$ studies, the observed effects sizes of individual studies $i$ (for $i =$

131 1,2, . . . $k$) are given by:

$$T_i = \Theta + \epsilon_i \tag{1}$$

$$\text{where } \epsilon_i \sim N(0, v_i) \tag{2}$$

132    Under the fixed effect model, the estimated population effect size $\hat{\theta}$ is obtained by

133  computing a weighted average of the observed effect sizes. If sampling error is assumed to

134  be the only source of variance in observed effect size, then it follows that studies with

135  smaller standard errors estimate the underlying true effect size more precisely. The

136  fixed-effect weights are thus simply the reciprocal of the sampling variance, $w_i = \frac{1}{v_i}$. The

137  estimate of the true effect is a weighted average across observed effect sizes:

$$\hat{\theta} = \frac{\sum_{i=1}^{k} w_i T_i}{\sum_{i=1}^{k} w_i} \tag{3}$$

138   ->

139    Whereas the fixed-effect model assumes that only one true population effect exists,

140  the random-effects model assumes that true effects may vary for unknown reasons, and

141  thus follow a (normal) distribution of their own (Hedges & Vevea, 1998). This

142  heterogeneity of the true effects is represented by their variance, $\tau^2$. The random effect

143  model thus assumes that the heterogeneity in observed effects can be decomposed into

144  sampling error and between-studies heterogeneity, resulting in the following equation for

145  the observed effect sizes:

$$T_i = \Theta + \zeta_i + \epsilon_i \tag{4}$$

$$\text{where } \zeta_i \sim N(0, tau^2) \tag{5}$$

$$\text{and } \epsilon_i \sim N(0, v_i) \tag{6}$$

146    In this model, $\Theta$ is the mean of the distribution of true effect sizes, and $\tau^2$ is its

147  variance, which can be interpreted as the variance between studies.

¹⁴⁸ If the true effect sizes follow a distribution, then even less precise studies (with larger

¹⁴⁹ sampling errors) may provide some information about this distribution. Like fixed-effect

¹⁵⁰ weights, random effects weights are still influenced by sampling error, but this influence is

¹⁵¹ attenuated by the estimated variance of the true effect sizes. The random-effects weights

¹⁵² are thus given by $w_i = \frac{1}{v_i + \hat{\tau}^2}$. It is important to note that, whereas the sampling error for

¹⁵³ each individual effect size is treated as known, between-study heterogeneity $\tau^2$ must be

¹⁵⁴ estimated. This estimate is represented by $\hat{\tau}^2$.

¹⁵⁵ **Meta-regression.** The random effects model assumes that causes of heterogeneity

¹⁵⁶ in the true effect sizes are unknown, and that their influence is random. Oftentimes,

¹⁵⁷ however, there are systematic sources of heterogeneity in true effect sizes. These

¹⁵⁸ between-study differences can be coded as moderators, and their influence can be

¹⁵⁹ estimated and controlled for using meta-regression. Meta-regression with $p$ moderators can

¹⁶⁰ be expressed with the following equation, where $x_{1...p}$ represent the moderators, and $\beta_{1...p}$

¹⁶¹ the regression coefficients:

$$T_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p + \zeta_i + \epsilon_i \tag{7}$$

$$\tag{8}$$

¹⁶² Note that $\beta_0$ represents the intercept of the distribution of true effect sizes after

¹⁶³ controlling for the moderators and the error term $\zeta_i$ represents residual between-studies

¹⁶⁴ heterogeneity. This term is still included because unexplained heterogeneity often remains

¹⁶⁵ after accounting for the moderators (Thompson & Sharp, 1999). This is a mixed-effects

¹⁶⁶ model; the intercept and effects of moderators are treated as fixed and the residual

¹⁶⁷ heterogeneity as random (Viechtbauer & López-López, 2015).

¹⁶⁸ To solve this model, the regression coefficients and residual heterogeneity must be

¹⁶⁹ estimated simultaneously. Numerous methods have been proposed to estimate

¹⁷⁰ meta-regression models, the most commonly used of which is restricted maximum

171 likelihood (REML). REML is an iterative method, meaning it performs the same

172 calculations repeatedly, updating the estimated regression coefficients and residual

173 heterogeneity until these estimates stabilize. In contrast to a regularized analysis

174 technique, this estimator produces low bias, which means that the average value of the

175 estimated regression coefficients and residual heterogeneity is close to their true values

176 (Panityakul et al., 2013; Hardy & Thompson, 1996).

## Regularized regression

178     Regularized regression biases parameter estimates towards zero by including a

179 shrinkage penalty in the estimation process. Before examining the Bayesian case, we will

180 explain the principle using frequentist OLS regression as an example. OLS regression

181 estimates the model parameters by minimizing the Residual Sum of Squares (RSS) of the

182 dependent variable, which is given by:

$$RSS = \sum_{i=1}^{n} (y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij})^2$$

183 The resulting parameter estimates are those that give the best predictions of the dependent

184 variable in the present dataset. Penalized regression, by contrast, adds a penalty term to

185 the quantity to be minimized. One commonly used penalty is the L1-norm of the

186 regression coefficients, or LASSO penalty, which corresponds to the sum of their absolute

187 values. This gives the penalized residual sum of squares:

$$PRSS = RSS + \lambda \sum_{j=1}^{p} |\beta_j|$$

188     Where $\lambda$ is a tuning parameter that determines how influential the penalty term will

189 be. If $\lambda$ is zero, the shrinkage penalty has no impact at all and the penalized regression will

190 produce the OLS estimates. If $\lambda \to \infty$, all coefficient shrink towards zero, producing the

191 null model. Because the penalty term is a function of the regression coefficients, the

192  optimizer has an incentive to keep the regression coefficients as small as possible. Note that

193  theLASSO penalty is but one example of a shrinkage penalty; other penalties exist, with

194  some unique properties.

## Bayesian estimation

196      Numerous methods have been proposed to accurately estimate the residual

197  heterogeneity, including the Hedges (HE), DerSimonian–Laird/Method of Moments (DL),

198  Sidik and Jonkman (SJ), Maximum Likelihood (ML), Restricted Maximum Likelihood

199  (REML), and Empirical Bayes (EB) method. These methods are mostly divided into two

200  groups: closed-form or non-iterative methods and iterative methods. The main difference

201  between these groups is that the closed form group uses a predetermined number of steps

202  to provide an estimation for the residual heterogeneity, whereas the iterative methods run

203  multiple iteration, as the name suggests, to converge to a solution when a specific criterion

204  is met. It is important to note that some iterative methods do not produce a solution when

205  they fail to converge after a predetermined amount of iteration.

206      In our scenario we are especially interested in an estimator which performs well under

207  the condition of a relative low number of studies. The Restricted Maximum Likelihood

208  (REML) seems to produce the lowest bias under this condition and is therefore preferred

209  (Panityakul et al., 2013; Hardy & Thompson, 1996). The REML is an iterative method

210  and needs a starting estimation of $\tau^2$ to start, usually it gets estimated by one of the

211  non-iterative methods (Viechtbauer & López-López, 2015). Besides the starting value of

212  $\tau^2$, it needs in every iteration an estimation of the regression coefficients of the moderators.

213  These are typically estimated by using the Weighted Least Squares (WLS) method. This is

214  a variation of the Ordinary Least Squares (OLS), but in the case of meta-analysis it is

215  necessary to assess weights to the coefficients. In systematic reviews large variation in

216  standard errors is often observed, which will result in large heteroscedasticity in the

217  estimation of the effects (Stanley & Doucouliagos, 2017). The addition of weights is a way

<sub>218</sub> to adjust for this heteroscedasticity. The weights are formulated as presented in equation

<sub>219</sub> (5).

<sub>220</sub>     The usage of a WLS method to estimate the regression coefficient may be

<sub>221</sub> problematic in the situation where a lot of moderators are measured without their specific

<sub>222</sub> effects, when the amount of studies is low and when moderators are dichotomous. The use

<sub>223</sub> of a least squares method will cause problems with the prediction accuracy and the model

<sub>224</sub> interpretability (James, Witten, Hastie, & Tibshirani, 2013). In the situation where a lot of

<sub>225</sub> moderators are measured and blindly included in the model, it may as well be the case that

<sub>226</sub> variables are included that are in fact not associated with the response. Including

<sub>227</sub> irrelevant variables in the model lowers the interpretability of the model (James et al.,

<sub>228</sub> 2013). An approach is necessary that automatically excludes the variables that are

<sub>229</sub> irrelevant i.e. performs variable selection. As explained before, in meta-analysis it is often

<sub>230</sub> the case that the number of moderators closely approaches or even exceeds the number of

<sub>231</sub> studies included in the analysis. A least squares method will display a lot variability in the

<sub>232</sub> fit when the number of variables is not much smaller than the number of studies (James et

<sub>233</sub> al., 2013). This means that the least squares method over fits the data and loses its power

<sub>234</sub> to be generalizable to future observations. When the number of variables exceeds the

<sub>235</sub> number of studies, the least squares method fails to produce one unique estimate and the

<sub>236</sub> method should not be used at all.

<sub>237</sub>     However, a least squares method could still be somewhat valuable in some situations.

<sub>238</sub> It is extremely suitable to estimate a linear relationship. In the case of dichotomous

<sub>239</sub> moderators, the relationship is always perfectly linear. A powerful non-linear estimation

<sub>240</sub> tool is in the situation of dichotomous moderators unnecessary and would not perform

<sub>241</sub> better at all. Whenever a non-linear relation gets fitted on data with an underlying linear

<sub>242</sub> relation, it will cause problems when this fit gets used for the prediction of future data.

<sub>243</sub> Given the various arguments, this paper provides an approach to tackle this problem of the

<sub>244</sub> least squares methods whilst still making use of a linear method. The weighted least

245 squares are replaced with the so-called LASSO regression for the estimation of the

246 regression coefficients. This algorithm shrinks or penalizes the regression coefficients and

247 performs variable selection (James et al., 2013; Hesterberg, Choi, Meier, & Fraley, 2008).

248 −>

249 **Algorithms and simulation + goal study** The goal of the present study was to

250 test whether a ME-MRA model with the lasso algorithm is able to outperform the

251 ME-MRA with least squares regression. More specifically, if the lasso is able to outperform

252 the least squares when in situation where the amount studies included in the analysis is

253 fairly low. To test this, two different algorithms are used; one called the rma, which makes

254 use of the WLS regression, and the lma, which makes use of a regularized lasso regression.

255 To test the lma and rma algorithms on the performance criteria, a simulation study is

256 performed. A simulation of the data is preferred over the use of real data. Simulated data

257 can be shaped to such an extent that it will have the all desired characteristics to test the

258 performance of the algorithm. Besides that, if simulated correctly, it will not have any

259 systematic errors or noise due to underlying models and it is more cost efficient.

260 Simulations are useful for evaluation of new methods like MetaForest and for the

261 comparison with alternative methods like metaCART and the classic approaches.

262 **Performance Criteria** The algorithms are evaluated on three different performance

263 criteria: The algorithms' predictive performance, their ability to estimate the residual

264 heterogeneity and their ability to detect and select the right moderators. The predictive

265 performance of the algorithms is defined by how well the algorithm is able to predict future

266 data. The algorithms have to estimate a model on a "training" dataset and then use this

267 model to see how well it fits on a second "testing" dataset. This is operationalized as the

268 cross-validated $R^2_{cv}$ (Van Lissa, 2017). The $R^2_{cv}$ is calculated using the fraction of variance

269 explained by the model on the testing dataset, relative to faction of variance explained by

270 the mean of the testing dataset. The mean of the testing dataset is the best prediction for

271 the testing data when there is no model present (van Lissa, 2017). The calculation of $R^2_{cv}$ is

272  expressed by the following equation:

273  $$R_{cv}^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y}_i)^2} \quad (9)$$

274  With $n$ being the number of studies in the testing dataset, $\hat{y}_i$ being the estimation for

275  study $i$, and $\bar{y}_i$ being the mean of the training dataset.

276  The ability of the algorithms to estimate the residual heterogeneity is by simply

277  taking the value of $\tau^2$ which the algorithm produces. The true value of the residual

278  heterogeneity is subtracted of the estimated value, solely to make the values more

279  interpretable. This means that a correct estimation of the residual heterogeneity will be

280  expressed by a value which exactly or close to zero. The residual heterogeneity is used as a

281  performance criterion because it is suspected that the lma model might not always be able

282  to predict residual heterogeneity correctly.

283  Variable selection is defined in terms of the algorithms ability to accredit positive

284  variable importance values to relevant moderators. Variable importance measures capture

285  the relative contribution of various moderators.

286  **Design factors** In the simulation study, meta analytic datasets will be simulated.

287  These datasets consist of two separate sub-datasets, a training- and a testing dataset. Both

288  sub-datasets will have the same characteristics with the exception of the number of studies

289  included. Certain characteristics of the sub-datasets will be manipulated to test how well

290  the algorithms perform under certain conditions. For each combination of characteristics,

291  or design factors, 100 datasets will be simulated. The design factors that will be

292  manipulated are the number of studies in the training data $k$ (22, 40 and 80), the average

293  within-study sample size $\bar{n}$ (40, 100 and 200), the population effect size $\beta$ (.2, .5 and .8)

294  and the residual heterogeneity $\tau^2$ (.01, .04 and .1). All the datasets will contain 20

295  moderators of which 10 are relevant and 10 are irrelevant. The moderators are binary and

296  follow $\sim \mathcal{B}]\nabla\backslash(0.5)$, which corresponds to an equal chance of being either one or zero. The

297  dependent variable $y_i$ represented by a $Hedges'g$. This is an estimator which takes the

standardized mean difference between a treatment and control group and is commonly

used in meta-analysis (Van Lissa, 2017). The true effect size $\theta_i$ is sampled out of a normal

distribution. The mean is computed by the assessing the values of the coefficients $\beta_j$, with

the values of the moderators and with the residual heterogeneity $\tau^2$ (Van Lissa, 2017).

This is in line with the calculation of $\theta_i$ represented in equation (6). The sampling error $\sigma_i^2$

is formed by varying the sizes of the samples of each study. The sample sizes $n_i \sim \mathcal{N}(\bar{n}, \frac{\bar{n}}{3}$

(Van Lissa, 2017).

Data were simulated using the random-effects model, based on four models: (A)

Main effect of one moderator, $\mu_i = \beta_1 x_{1i}$ (B) Two-way interaction,

$\mu_i = \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{1i} x_{2i}$ (E) Non-linear, cubic relationship, $\mu_i = \beta_1 x_{1i} + \beta_2 x_{1i}^2 + \beta_3 x_{1i}^3$

(F) Exponential relationship, $\mu_i = \beta_1 e^{x_{1i}}$

**Impact of design factors and hypotheses** These design factors are chosen on

purpose, because they are hypothesized to have an influence on the predictive performance

of the algorithms. The effect of the design factors ought to be either positive or negative on

the data. This means that some factor should, by increasing, make the data easier to be

analyzed, or make it more difficult to analyze. The amount of studies included in the

training data $k$ has a positive influence on the variance explained by the different

algorithms. This is due to the fact that there are simply more data points available to fit a

model on. The lma algorithm should be superior on the low value of $k$ over the rma

algorithm. The effect size $\beta$ has a positive impact on the ability of the algorithms to explain

variance. It can be hypothesized that the lma performs better at lower values of $\beta$ because

it is better equipped to detect and select variables when even when the amount of signal is

low. The residual heterogeneity $\tau^2$ should have a negative influence on the interpretability

of the data. Differences between the two algorithms could be present, but it remains

unclear which would perform better. The lma might perform better when the amount of

signal in the data is low or the noise is high, but it is also suspected to overestimate the

amount of heterogeneity and this could worsen if the $\tau^2$ increases. The $\bar{n}$ greatly influences

<sub>325</sub> the quality of the data. Higher values of within-study sample sizes reduce the sampling

<sub>326</sub> error. This will lead to a better prediction by the algorithms. In conclusion: higher values

<sub>327</sub> of $k$, $\beta$ and $\bar{n}$ will increase the quality of the data, where higher values of $\tau^2$ decrease the

<sub>328</sub> quality of the data. The lma is suspected to perform significantly better when the quality

<sub>329</sub> of the data is low, especially when the amount of studies in the sample is low, with the

<sub>330</sub> exception of the performance of the lma on the estimation of the residual heterogeneity

<sub>331</sub> # Results

<sub>332</sub> There were 3888 condition combinations in total. The algorithms ran on 100 different

<sub>333</sub> datasets on each of those, leaving 388800 cases to analyze. There were 20 of those cases for

<sub>334</sub> which the RMA algorithm had missing values on all metrics. Closer inspection showed that

<sub>335</sub> both the cubic and exponential model each contributed ten times to the missing values and

<sub>336</sub> only when 2, 3 or 6 moderators were taken up in the model. However, since the missing

<sub>337</sub> cases make up only 0.005% of the data, the cases were chosen to be omitted from further

<sub>338</sub> analysis entirely. Another observation is that the two-way interaction model only had

<sub>339</sub> results when the number of moderators in the model was either 3, 4 or 7, while the other

<sub>340</sub> models only had results when there were 2, 3 or 6 moderators. It is therefor more

<sub>341</sub> challenging to compare and interpret the effect the moderators had on the performance

<sub>342</sub> criteria between the two-way interaction and the other models.

<sub>343</sub> **Predictive performance**

<sub>344</sub> Predictive performance was operationalized by calculating the $R^2_{test}$ and $MSE_{test}$ for

<sub>345</sub> every combination of design factors, in further test denoted as $R^2$ and $MSE$ respectively.

<sub>346</sub> The densities for the $R^2_{test}$ and $MSE_{test}$ values were skewed however, which is why it was

<sub>347</sub> chosen to use the median $R^2$ as the metric for predictive performance, rather than the

<sub>348</sub> mean. The spread of the metrics was described using the Mean Absolute Deviation [MAD],

<sub>349</sub> rather than the standard deviation. It was found that the Horseshoe, Lasso and RMA

350   algorithm performed similarly overall, $R^2_{Hs} = 0.51 \pm 0.36$, $MSE_{Hs} = 0.21 \pm 0.18$ ;

351   $R^2_{Lasso} = 0.50 \pm 0.37$, $MSE_{Lasso} = 0.21 \pm 0.19$; $R^2_{RMA} = 0.50 \pm 0.37$,

352   $MSE_{RMA} = 0.22 \pm 0.23$ . The MetaForest algorithm performed worst on $R^2$:

353   $R^2_{Mf} = 0.35 \pm 0.38$, $MSE_{Mf} = 0.22 \pm 0.19$

354   To determine the effect of the design factors on $R^2$ for all algorithms, four separate

355   ANOVA's were performed; one per algorithm. The effect size $\eta^2$ per condition per

356   algorithm, including for all two-way interactions can be found in table 1. Do note that the

357   ANOVA's were performed with the normality assumption violated. The estimates serve

358   mostly as a guidance, rather than an absolute result.

359   Not too surprisingly, It was found that the true effect size $\beta$ had the largest effect on

360   $R^2$ for all algorithms. As $\beta$ increased, the performance of all algorithms increased as well.

361   However, $\beta$ did interact with the model that was estimated, being either a linear, two-way

362   interaction, cubic or an exponential model. A graphical representation of the interaction is

363   shown in image 1. When the exponential and cubic model were estimated, the increase of

364   $R^2$ slowed for every higher value of $\beta$. For the estimation of the two-way interaction model

365   the increase only slowed when $\beta$ went from 0.5 to 0.8, while the steepest increase for the

366   linear model estimation was when $\beta$ increased from 0.2 to 0.5. Also noteworthy is that $R^2$

367   stagnated during estimation of the cubic model as $\beta$ went up to 0.5, while for the other

368   models $R^2$ did keep increasing. There was little difference in median $R^2$ between Horsehoe,

369   Lasso and RMA, while MetaForest performed worst.

370   The second largest marginal effect was that of the estimated model. All algorithms

371   had the highest $R^2$ under the cubic model, followed by a similar performance on the

372   two-way interaction and exponential models. All algorithms performed worst for the linear

373   model. There again was little difference in performance between the Pema algorithms and

374   RMA, although MetaForest performed worst. Image 2A shows the relationship.

375   There also was a moderate interaction effect between the estimated model and the

376 amount of skewness of the input data $\alpha$, especially for the Pema algorithms. Again, Pema

377 and RMA algorithms performed best, followed by MetaForest in all conditions. Most

378 obvious to note is that all algorithms, except for when the linear model was estimated,

379 generally performed better on more skewed data, although the algorithms did perform

380 worse during estimation of the two-way interaction model when $\alpha$ went from 5 to 10.

381 Overall performance was best during estimation of the cubic model, but the performance

382 difference between estimated models decreased as $\alpha$ increased. Image 3 shows the

383 relationship.

384         The true residual heterogeneity $\tau^2$ had a negative linear relationship for all

385 algorithms on $R^2$. That is, as $\tau^2$ increased, $R^2$ decreased. Image 2B shows the relationship.

386         The mean sample size per study $\bar{n}$ also had a moderate effect. For all algorithms the

387 effect of $\bar{n}$ was positively linearly related with $R^2$. Again, the Pema and RMA algorithms

388 performed better than MetaForest. Image 2C shows the relationship.

389         An especially large effect was found for the number of studies used in the training

390 data $\kappa$ for MetaForest, while this effect was substantially smaller for RMA, Lasso and

391 Horseshoe. The relationship is positively linear for all algorithms, but the slope is

392 especially steep for MetaForest. Image 2D shows the relationship.

393         Finally, the number of moderators did not have a big effect for the Pema algorithms,

394 while for RMA and MetaForest the effect was more noticeable. The relationship is shown

395 in image 7. The relationship is generally negative with more moderators meaning worse

396 performance, although an increase can be observed as the number of moderators increase

397 from 4 to 6 for all algorithms except MetaForest. This increase in performance for

398 MetaForest appears when the number of moderators go from 3 to 4.

### Estimating residual heterogeneity

The ability of the algorithms to correctly estimate $\tau^2$ was operationalized by subtracting the true value for $\tau^2$ from the $\tau^2$ estimated by the algorithms. Again, the median and Mean Absolute Deviation were used as metrics for performance. The RMA algorithm showed the best results, $\Delta\tau^2_{RMA} = 0.02 \pm 0.06$, followed by the MetaForest algorithm $\Delta\tau^2_{Mf} = 0.09 \pm 0.13$. The Pema algorithms performed worst $\Delta\tau^2_{Lasso} = 0.23 \pm 0.18$; $\Delta\tau^2_{Hs} = 0.23 \pm 0.17$. The finding that all medians are positive implies that all algorithms have a bigger tendency to overestimate $\tau^2$ than underestimate it. One comment to make is that uncertainty of the estimates generally increased as $\Delta\tau^2$ also increased. This implies that there was more variation in performance as median performance worsened.

To determine the effect of the design factors on $\Delta\tau^2$ for all algorithms, four separate ANOVA's were performed. The effect size $\eta^2$ per condition per algorithm, including $\eta^2$ for all two-way interactions can be found in table 2. Again, the assumption of normality was violated.

The biggest predictor on the correct estimation of $\tau^2$ was the estimated model. This was mainly because the algorithms overestimated $\tau^2$ most when the model contained cubic terms. Image 5A shows the marginal relationship of the estimated model on $\Delta\tau^2$ . It becomes more clear why this overestimation occurred when showing the interaction between $\beta$ and the model estimated on $\Delta\tau^2$, shown in image 6. First note the general trend that during estimation of all models $\tau^2$ got more overestimated as $\beta$ increased, except during estimation of the linear model, where the effect of $\beta$ on $\Delta\tau^2$ was close to zero, except for MetaForest. However, note the scales for the y-axes. While estimating the two-way interaction, linear and exponential model, $\Delta\tau^2$ stayed well within a confined interval. However, the algorithms severely overestimated $\tau^2$ when the model contained cubic terms. Especially MetaForest overestimated $\tau^2$ substantially when $\beta = 0.8$ and the

estimated model is cubic: $\Delta\tau^2_{MF} = 2.92 \pm 2.37$. The other algorithms also had a $\Delta\tau^2 > 1$ in these conditions, but the results were not as severe. Interestingly, the Pema algorithms even outperformed the RMA algorithm in these conditions.

The marginal effects of $\beta$ on $\Delta\tau^2$ are shown in image 5B. MetaForest was affected most by the increase in $\beta$, but in general performed better than the Pema algorithms when $\beta < 0.8$. The RMA algorithm performed best overall.

The marginal effect of $\alpha$ on $\Delta\tau^2$ was rather minimal, although there was a slight decrease in $\Delta\tau^2$ as $\alpha$ increased. However, the decrease is more explicit when the interaction of $\alpha$ with the estimated model is added. Image 7 shows this interaction. The algorithms were rather unaffected by $\alpha$ for the linear model, and a small decrease in $\Delta\tau^2$ as $\alpha$ increased can be seen in the exponential model. When the two-way interaction model was estimated however, the algorithms benefitted as $\alpha$ increased, while for the cubic model, $\Delta\tau^2$ first increased as $\alpha$ increased from 0 to 2, but decreased as $\alpha$ increased from 2 to 10. RMA performed best, followed by MetaForest. The Pema algorithms performed similarly, but worst.

The effect of the true $\tau^2$ on $\Delta\tau^2$ was rather unnoticeable for the RMA and MetaForest algorithms. The tendency for the Pema algorithms on the other hand, was to overestimate $\tau^2$ more as the true $\tau^2$ increased. Image 5C shows the marginal relationship.

The effect of the number of moderators on $\Delta\tau^2$ was not that large either. A small increase in $\Delta\tau^2$ can be seen in the RMA and MetaForest algorithm as the number of moderators increased which was not found for the Pema algorithms. However, A small note is that MetaForest did substantially increase in $\Delta\tau^2$ as more moderators were added and the estimated model is cubic. Image 8 shows the interaction between the number of moderators and the estimated model.

$\kappa$ only had a substantial effect for MetaForest; the $\Delta\tau^2$ decreased quite rapidly if $\kappa$ increased, especially when the cubic model was estimated. For the other algorithms,

decreasing $\kappa$ had little to no effect on correctly estimating the residual heterogeneity. Image 9 shows the interaction of $\kappa$ with the estimated model.

Finally, the average number of observations in the studies did not have a substantial effect on $\Delta\tau^2$. Image 5D shows the marginal relationship.

**Variable selection**

To determine the extent to which the algorithms could perform variable selection correctly, the proportion true positives $[TP]$ and true negatives $[TN]$ were calculated. The $TP$ and $TN$ reflect how well the algorithms accredit importance to relevant moderators and discredit importance to irrelevant moderators respectively. It should be noted that $TP$ could only take on values 1 or 0 per simulated iteration, because in all models where $\beta > 0$, only one moderator was simulated to be relevant. As $\beta$ increased, the already relevant moderator increased in relevance, rather than spreading the relevance over the other moderators. $TN$ had a bigger range and could take on values dependent on how many moderators were taken up in the model. E.g. when $\beta = 0$ and $n_{mods} = n$, $n + 1$ different proportions were possible, $\frac{0}{n}$ up until $\frac{n}{n}$.

There were no differences in variables selected by Highest Density Intervals or Confidence Intervals for both Lasso and Horshoe and so for both algorithms it did not matter which interval type was analyzed. It was found that MetaForest had the highest proportion true positives: $TP_{Mf} = 0.98$, closely followed by RMA: $TP_{RMA} = 0.96$. Horshoe performed slightly better than Lasso; $TP_{Hs} = 0.91$; $TP_{Lasso} = 0.89$. As for $TN$, it was found that the pema algorithms performed best: $TN_{Hs}$ and $TN_{Lasso} = 0.93$, followed by RMA: $TN_{RMA} = 0.89$. MetaForest performed worst by a large margin: $TN_{Mf} = 0.50$. The Mean Absolute Deviation for all algorithms on both $TP$ and $TN$ was 0, except for MetaForests performance on $TN$, where the Mean Absolute Deviation was 0.44.

Perfomance on $TP$ and $TN$ were very high for all algorithms, with all mean

476  proportions, except MetaForests performance on $TN$, exceeding .89. This implies that

477  MetaForest had issues excluding irrelevant moderators from the models. Plots were

478  inspected to determine the effects of the design factors on the proportions. While

479  inspecting the plots, there were found to be little marginal effects of the design factors on

480  $TN$, while $TP$ was more affected.

481      Firstly, $\kappa$ only had a positive effect on $TP$. As $\kappa$ increased, $TP$ also increased.

482  MetaForest had the highest $TP$, followed by RMA and, lastly, Horsehoe and Lasso. The

483  increase in $TP$ for higher values of $\kappa$ was steeper for the Pema algorithms, however. There

484  was also an interaction of $\kappa$ with the estimated model shown in image 10. During

485  estimation of the linear and two-way interaction model, the relationship of $\kappa$ looked

486  relatively linear. At $\kappa = 20$ the $TP$ was relatively low for the algorithms, compared to the

487  cubic model where $TP$ starts at .98 and converges to 1 as $\kappa$ increased. This latter

488  relationship was also found for the exponential model, although the $TP$ at $\kappa = 20$ was

489  lower.

490      $\bar{n}$ had a positive and roughly linear relationship with $TP$ for all algorithms.

491  MetaForest performed best, followed by RMA, while Horseshoe and Lasso performed

492  worst. Image 11A shows the relationship.

493      $\beta$ had an interaction effect with the model estimated on $TN$. Only during estimation

494  of the two-way interaction model, $TN$ decreased as $\beta$ increased, For the other models, $TN$

495  remained stable. This could be because the interaction model was only fitted when there

496  were 3,4 or 7 moderators, while for the other models, only 2,3 or 6 moderators were used.

497  Image 12 shows the relationship. The effect of $\beta$ on $TP$ was positive; as $\beta$ increased, $TP$

498  increased too.

499      The true $\tau^2$ had a negative effect on $TP$, while $\alpha$ seemed to have little effect on both

500  $TP$ and $TN$. Image 11B shows the marginal relationship of $\tau^2$ on $TP$.

501      Finally, there was an effect of number of moderators on $TN$, but only for the two-way

interaction model. The $TN$ increased as the number of moderators did. Image 13 shows the interaction. This relationship was reversed for $TP$ and was found during estimation of all models, i.e. $TP$ decreased as the number of moderators increased. Image 14 shows the relationship.

## Discussion

# References

507