

1 Simultaneous SNP selection and adjustment for  
2 population structure in high dimensional prediction  
3 models

4 Sahir R Bhatnagar<sup>1,2</sup>, Yi Yang<sup>4</sup>, Tianyuan Lu<sup>2</sup>, Erwin Schurr<sup>6</sup>,  
5 JC Loredano-Osti<sup>7</sup>, Marie Forest<sup>2</sup>, Karim Oualkacha<sup>3</sup>, and  
6 Celia MT Greenwood<sup>1,2,5</sup>

7 <sup>1</sup>Department of Epidemiology, Biostatistics and Occupational Health,  
8 McGill University

9 <sup>2</sup>Lady Davis Institute, Jewish General Hospital, Montréal, QC

10 <sup>3</sup>Département de Mathématiques, Université de Québec À Montréal

11 <sup>4</sup>Department of Mathematics and Statistics, McGill University

12 <sup>5</sup>Departments of Oncology and Human Genetics, McGill University

13 <sup>6</sup>Department of Medicine, McGill University

14 <sup>7</sup>Department of Mathematics and Statistics, Memorial University

15 July 12, 2019

16 **Abstract**

17 Complex traits are known to be influenced by a combination of environmental fac-

18 tors and rare and common genetic variants. However, detection of such multivariate  
19 associations can be compromised by low statistical power and confounding by popu-  
20 lation structure. Linear mixed effects models (LMM) can account for correlations due  
21 to relatedness but have not been applicable in high-dimensional (HD) settings where  
22 the number of fixed effect predictors greatly exceeds the number of samples. False  
23 positives or false negatives can result from two-stage approaches, where the residuals  
24 estimated from a null model adjusted for the subjects' relationship structure are sub-  
25 sequently used as the response in a standard penalized regression model. To overcome  
26 these challenges, we develop a general penalized LMM framework called `ggmix` for  
27 simultaneous SNP selection and adjustment for population structure in high dimen-  
28 sional prediction models. Our method can accommodate several sparsity-inducing  
29 penalties such as the lasso, elastic net and group lasso, and also readily handles prior  
30 annotation information in the form of weights. We develop a blockwise coordinate  
31 descent algorithm which is highly scalable, computationally efficient and has theo-  
32 retical guarantees of convergence. Through simulations and two real data examples,  
33 we show that `ggmix` leads to better sensitivity and specificity compared to the two-  
34 stage approach or principal component adjustment with better prediction accuracy.  
35 `ggmix` can be used to construct polygenic risk scores and select instrumental variables  
36 in Mendelian randomization studies. Our algorithms are available in an R package  
37 (<https://github.com/greenwoodlab/ggmix>).

## 38 1 Author Summary

39 This work addresses a recurring challenge in the analysis and interpretation of genetic as-  
40 sociation studies: which genetic variants can best predict and are independently associated  
41 with a given phenotype in the presence of population structure? Not controlling confound-  
42 ing due to geographic population structure, family and/or cryptic relatedness can lead to  
43 spurious associations. Much of the existing research has therefore focused on modeling the

44 association between a phenotype and a single genetic variant in a linear mixed model with  
45 a random effect. However, this univariate approach may miss true associations due to the  
46 stringent significance thresholds required to reduce the number of false positives and also  
47 ignores the correlations between markers. We propose an alternative method for fitting  
48 high-dimensional multivariable models, which selects SNPs that are independently associ-  
49 ated with the phenotype while also accounting for population structure. We provide an  
50 efficient implementation of our algorithm and show through simulation studies and real data  
51 examples that our method outperforms existing methods in terms of prediction accuracy  
52 and controlling the false discovery rate.

## 53 2 Introduction

54 Genome-wide association studies (GWAS) have become the standard method for analyzing  
55 genetic datasets owing to their success in identifying thousands of genetic variants associated  
56 with complex diseases (<https://www.genome.gov/gwastudies/>). Despite these impressive  
57 findings, the discovered markers have only been able to explain a small proportion of the  
58 phenotypic variance; this is known as the missing heritability problem [1]. One plausible  
59 reason is that there are many causal variants that each explain a small amount of variation  
60 with small effect sizes [2]. Methods such GWAS, which test each variant or single nucleotide  
61 polymorphism (SNP) independently, may miss these true associations due to the stringent  
62 significance thresholds required to reduce the number of false positives [1]. Another major  
63 issue to overcome is that of confounding due to geographic population structure, family  
64 and/or cryptic relatedness which can lead to spurious associations [3]. For example, there  
65 may be subpopulations within a study that differ with respect to their genotype frequencies  
66 at a particular locus due to geographical location or their ancestry. This heterogeneity in  
67 genotype frequency can cause correlations with other loci and consequently mimic the signal  
68 of association even though there is no biological association [4, 5]. Studies that separate

69 their sample by ethnicity to address this confounding suffer from a loss in statistical power  
70 due to the drop in sample size.

71 To address the first problem, multivariable regression methods have been proposed which  
72 simultaneously fit many SNPs in a single model [6, 7]. Indeed, the power to detect an  
73 association for a given SNP may be increased when other causal SNPs have been accounted  
74 for. Conversely, a stronger signal from a causal SNP may weaken false signals when modeled  
75 jointly [6].

76 Solutions for confounding by population structure have also received significant attention in  
77 the literature [8, 9, 10, 11]. There are two main approaches to account for the relatedness  
78 between subjects: 1) the principal component (PC) adjustment method and 2) the linear  
79 mixed model (LMM). The PC adjustment method includes the top PCs of genome-wide  
80 SNP genotypes as additional covariates in the model [12]. The LMM uses an estimated  
81 covariance matrix from the individuals' genotypes and includes this information in the form  
82 of a random effect [3].

83 While these problems have been addressed in isolation, there has been relatively little  
84 progress towards addressing them jointly at a large scale. Region-based tests of association  
85 have been developed where a linear combination of  $p$  variants is regressed on the response  
86 variable in a mixed model framework [13]. In case-control data, a stepwise logistic-regression  
87 procedure was used to evaluate the relative importance of variants within a small genetic  
88 region [14]. These methods however are not applicable in the high-dimensional setting, i.e.,  
89 when the number of variables  $p$  is much larger than the sample size  $n$ , as is often the case in  
90 genetic studies where millions of variants are measured on thousands of individuals.

91 There has been recent interest in using penalized linear mixed models, which place a con-  
92 straint on the magnitude of the effect sizes while controlling for confounding factors such as  
93 population structure. For example, the LMM-lasso [15] places a Laplace prior on all main  
94 effects while the adaptive mixed lasso [16] uses the  $L_1$  penalty [17] with adaptively chosen

95 weights [18] to allow for differential shrinkage amongst the variables in the model. Another  
96 method applied a combination of both the lasso and group lasso penalties in order to select  
97 variants within a gene most associated with the response [19]. However, these methods are  
98 normally performed in two steps. First, the variance components are estimated once from  
99 a LMM with a single random effect. These LMMs normally use the estimated covariance  
100 matrix from the individuals' genotypes to account for the relatedness but assumes no SNP  
101 main effects (i.e. a null model). The residuals from this null model with a single random  
102 effect can be treated as independent observations because the relatedness has been effec-  
103 tively removed from the original response. In the second step, these residuals are used as the  
104 response in any high-dimensional model that assumes uncorrelated errors. This approach  
105 has both computational and practical advantages since existing penalized regression soft-  
106 ware such as `glmnet` [20] and `gglasso` [21], which assume independent observations, can be  
107 applied directly to the residuals. However, recent work has shown that there can be a loss in  
108 power if a causal variant is included in the calculation of the covariance matrix as its effect  
109 will have been removed in the first step [13, 22].

110 In this paper we develop a general penalized LMM framework called `ggmix` that simul-  
111 taneously selects variables and estimates their effects, accounting for between-individual  
112 correlations. Our method can accommodate several sparsity inducing penalties such as the  
113 lasso [17], elastic net [23] and group lasso [24]. `ggmix` also readily handles prior annotation  
114 information in the form of a penalty factor, which can be useful, for example, when dealing  
115 with rare variants. We develop a blockwise coordinate descent algorithm which is highly  
116 scalable and has theoretical guarantees of convergence to a stationary point. All of our  
117 algorithms are implemented in the `ggmix` R package hosted on GitHub with extensive docu-  
118 mentation (<https://github.com/greenwoodlab/ggmix>). We provide a brief demonstration  
119 of the `ggmix` package in Appendix C.

120 The rest of the paper is organized as follows. In Section 3, we compare the performance

121 of our proposed approach and demonstrate the scenarios where it can be advantageous to  
122 use over existing methods through simulation studies and two real data analyses. This is  
123 followed by a discussion of our results, some limitations and future directions in Section 4.  
124 Section 5 describes the `ggmix` model, the optimization procedure and the algorithm used to  
125 fit it.

## 126 **3 Results**

127 In this section we demonstrate the performance of `ggmix` in a simulation study and two real  
128 data applications.

### 129 **3.1 Simulation Study**

130 We evaluated the performance of `ggmix` in a variety of simulated scenarios. For each simu-  
131 lation scenario we compared `ggmix` to the `lasso` and the `twostep` method. For the `lasso`,  
132 we included the top 10 principal components from the simulated genotypes used to calcu-  
133 late the kinship matrix as unpenalized predictors in the design matrix. For the `twostep`  
134 method, we first fitted an intercept only model with a single random effect using the average  
135 information restricted maximum likelihood (AIREML) algorithm [25] as implemented in the  
136 `gaston` R package [26]. The residuals from this model were then used as the response in a  
137 regular `lasso` model. Note that in the `twostep` method, we removed the kinship effect in  
138 the first step and therefore did not need to make any further adjustments when fitting the  
139 penalized model. We fitted the `lasso` using the default settings and `standardize=FALSE`  
140 in the `glmnet` package [20]. For other parameters in our simulation study, we defined the  
141 following quantities:

- 142 •  $n$ : sample size

- 143 •  $c$ : percentage of causal SNPs
- 144 •  $\beta$ : true effect size vector of length  $p_{fixed}$
- 145 •  $S_0 = \{j; (\beta)_j \neq 0\}$  the index of the true active set with cardinality  $|S_0| = c \times p_{fixed}$
- 146 •  $\mathbf{X}^{(fixed)}$ :  $n \times p_{fixed}$  matrix of SNPs that were included as fixed effects in the model
- 147 •  $\mathbf{X}^{(causal)}$ :  $n \times |S_0|$  matrix of SNPs that were truly associated with the simulated phe-  
148 nototype, where  $\mathbf{X}^{(causal)} \subseteq \mathbf{X}^{(fixed)}$
- 149 •  $\mathbf{X}^{(other)}$ :  $n \times p_{other}$  matrix of SNPs that were used in the construction of the kinship  
150 matrix. Some of these  $\mathbf{X}^{(other)}$  SNPs, in conjunction with some of the SNPs in  $\mathbf{X}^{(fixed)}$   
151 were used in construction of the kinship matrix. We altered the balance between these  
152 two contributors and with the proportion of causal SNPs used to calculate kinship
- 153 •  $\mathbf{X}^{(kinship)}$ :  $n \times k$  matrix of SNPs used to construct the kinship matrix

154 We simulated data from the model

$$\mathbf{Y} = \mathbf{X}^{(fixed)}\beta + \mathbf{P} + \boldsymbol{\varepsilon} \quad (1)$$

155 where  $\mathbf{P} \sim \mathcal{N}(0, \eta\sigma^2\boldsymbol{\Phi})$  is the polygenic effect and  $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, (1 - \eta)\sigma^2\mathbf{I})$  is the error term.  
156 Here,  $\boldsymbol{\Phi}_{n \times n}$  is the covariance matrix calculated from  $\mathbf{X}^{(kinship)}$ ,  $\mathbf{I}_{n \times n}$  is the identity matrix  
157 and parameters  $\sigma^2$  and  $\eta \in [0, 1]$  determine how the variance is divided between  $\mathbf{P}$  and  
158  $\boldsymbol{\varepsilon}$ . The values of the parameters that we used were as follows: narrow sense heritability  
159  $\eta = \{0.1, 0.3\}$ , number of fixed effects  $p_{fixed} = 5,000$ , number of SNPs used to calculate the  
160 kinship matrix  $k = 10,000$ , percentage of causal SNPs  $c = \{0\%, 1\%\}$  and  $\sigma^2 = 1$ . In addition  
161 to these parameters, we also varied the amount of overlap between the causal SNPs and the  
162 SNPs used to generate the kinship matrix. We considered two main scenarios:

1. None of the causal SNPs are included in the calculation of the kinship matrix:

$$\mathbf{X}^{(kinship)} = \left[ \mathbf{X}^{(other)} \right]$$

2. All the causal SNPs are included in the calculation of the kinship matrix:

$$\mathbf{X}^{(kinship)} = \left[ \mathbf{X}^{(other)}; \mathbf{X}^{(causal)} \right].$$

163 Both kinship matrices were meant to contrast the model behavior when the causal SNPs  
164 are included in both the main effects and random effects versus when the causal SNPs are  
165 only included in the main effects. These scenarios are motivated by the current standard of  
166 practice in GWAS where the candidate marker is excluded from the calculation of the kinship  
167 matrix [8]. This approach becomes much more difficult to apply in large-scale multivariable  
168 models where there is likely to be overlap between the variables in the design matrix and  
169 kinship matrix. We simulated random genotypes from the BN-PSD admixture model with  
170 1D geography and 10 subpopulations using the `bnpsd` package [27, 28]. In Figure 1, we plot  
171 the estimated kinship matrix from a single simulated dataset in the form of a heatmap where  
172 a darker color indicates a closer genetic relationship.

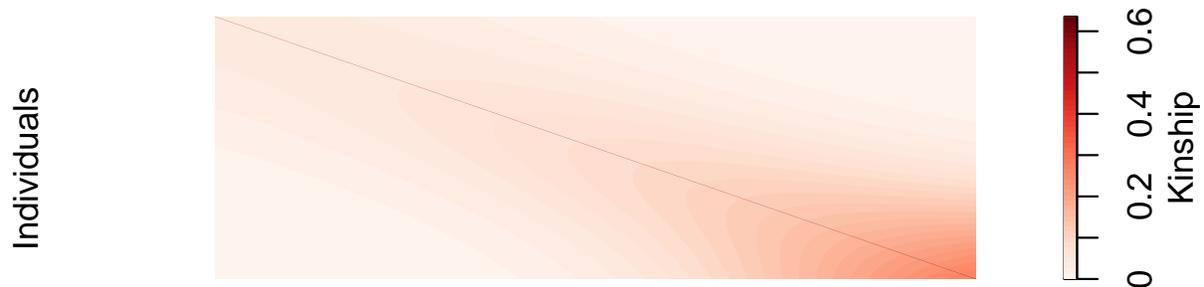


Figure 1: Example of an empirical kinship matrix used in simulation studies. This scenario models a 1D geography with extensive admixture.

173 In Figure 2 we plot the first two principal component scores calculated from the simulated  
174 genotypes used to calculate the kinship matrix in Figure 1, and color each point by sub-  
175 population membership. We can see that the PCs can identify the subpopulations which  
176 is why including them as additional covariates in a regression model has been considered a  
177 reasonable approach to control for confounding.

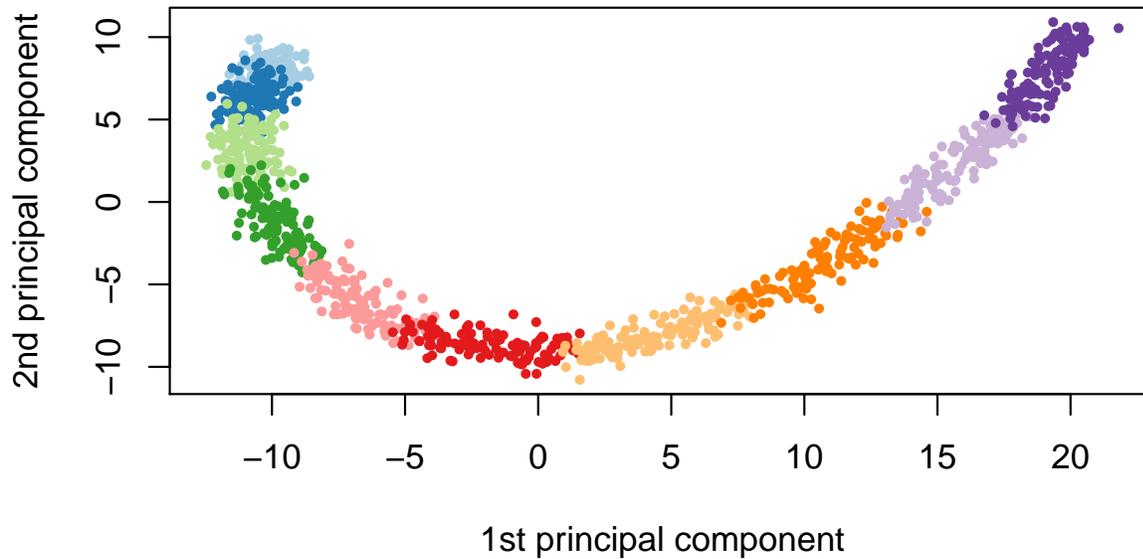


Figure 2: First two principal component scores of the genotype data used to estimate the kinship matrix where each color represents one of the 10 simulated subpopulations.

178 Using this set-up, we randomly partitioned 1000 simulated observations into 80% for training  
 179 and 20% for testing. The training set was used to fit the model and select the optimal tuning  
 180 parameter only, and the resulting model was evaluated on the test set. Let  $\hat{\lambda}$  be the esti-  
 181 mated value of the optimal regularization parameter,  $\hat{\beta}_{\hat{\lambda}}$  the estimate of  $\beta$  at regularization  
 182 parameter  $\hat{\lambda}$ , and  $\hat{S}_{\hat{\lambda}} = \{j; (\hat{\beta}_{\hat{\lambda}})_j \neq 0\}$  the index of the set of non-zero estimated coefficients.  
 183 We evaluated the methods based on correct sparsity defined as  $\frac{1}{p} \sum_{j=1}^p A_j$ , where

$$A_j = \begin{cases} 1 & \text{if } (\hat{\beta}_{\hat{\lambda}})_j = (\beta)_j = 0 \\ 1 & \text{if } (\hat{\beta}_{\hat{\lambda}})_j \neq 0, (\beta)_j \neq 0 \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

184 We also compared the test set prediction error based on the refitted unpenalized estimates

185 for each selected model, the estimation error ( $\|\widehat{\beta} - \beta\|_2^2$ ), true positive rate ( $|\widehat{S}_\lambda \cap S_0|/|S_0|$ ), false  
186 positive rate ( $|\widehat{S}_\lambda \setminus S_0|/|j \notin S_0|$ ), and the variance components ( $\eta, \sigma^2$ ) for the polygenic random  
187 effect and error term.

188 In Figure 3, we present the results for the scenario with 1% causal SNPs ( $c = 0.01$ ) which were  
189 all used in the calculation of the kinship matrix and true heritability  $\eta = 10\%$ . The complete  
190 simulation results are shown in supplementary Section B. We see that `ggmix` outperformed  
191 both the `twostep` and `lasso` in terms of correct sparsity and estimation error (Figure 3  
192 panels A and B). This was true regardless of true heritability and whether the causal SNPs  
193 were included in the calculation of the kinship matrix (Figures B.1, B.8, B.2 and B.9). Across  
194 all simulation scenarios, `ggmix` had the smallest root mean squared prediction error (RMSE)  
195 on the test set while also producing the most parsimonious models (Figures 3 panel B, B.3  
196 and B.13). Both the `lasso` and `twostep` had on average, slightly higher true positive rate  
197 compared to `ggmix` but came at the cost of a higher false positive rate (Figures 3 panel D, B.4  
198 and B.10). Both the `twostep` and `ggmix` overestimated the heritability though `ggmix` was  
199 closer to the true value (Figure 3 panel E). When none of the causal SNPs were in the  
200 kinship, both methods tended to overestimate the truth when  $\eta = 10\%$  and underestimate  
201 when  $\eta = 30\%$  (Figure B.11). Across all simulation scenarios `ggmix` was able to (on average)  
202 correctly estimate the error variance (Figures 3 panel F, B.6 and B.12). The `lasso` tended  
203 to overestimate  $\sigma^2$  in the null model while the `twostep` overestimated  $\sigma^2$  when none of the  
204 causal SNPs were in the kinship matrix.

205 Overall, we observed that variable selection results and RMSE for `ggmix` were similar regard-  
206 less of whether the causal SNPs were in the kinship matrix or not. This result is encouraging  
207 since in practice the kinship matrix is constructed from a random sample of SNPs across the  
208 genome, some of which are likely to be causal, particularly in polygenic traits. `ggmix` had  
209 very good Type 1 and II error control, while both the `lasso` and `twostep` had a very high  
210 false positive rate in all simulation scenarios. In particular, our simulation results show that

211 the principal component adjustment method may not be the best approach to control for  
 212 confounding by population structure, particularly when variable selection is of interest.

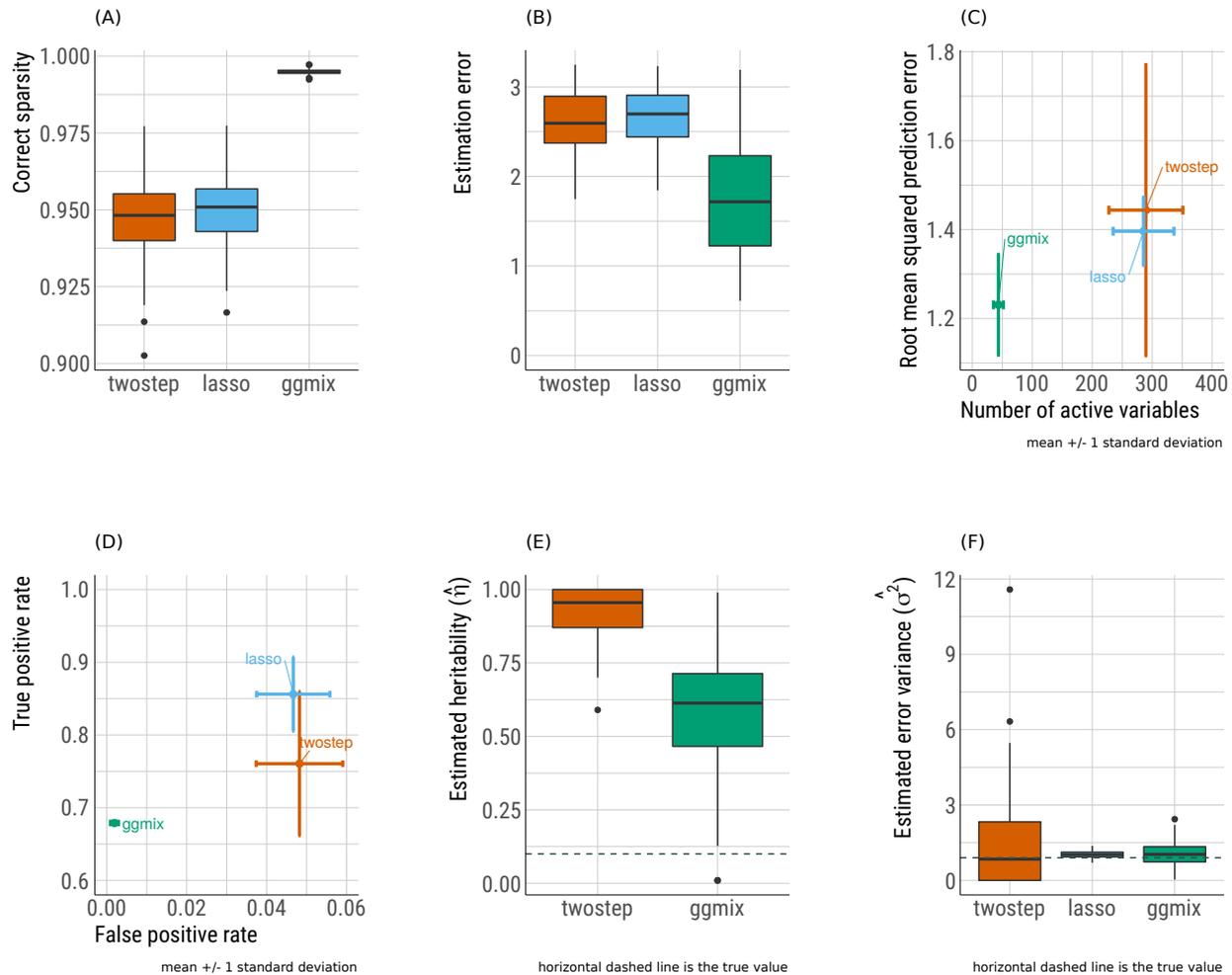


Figure 3: Results from 200 replications for the scenario with 1% causal SNPs ( $c = 0.01$ ) which are all used in the calculation of the kinship matrix and true heritability  $\eta = 10\%$ . (A) Correct sparsity as defined by Equation (2). (B) Estimation error defined as the squared distance between the estimated and true effect sizes (C) Root mean squared prediction error on the test set as a function of the number of selected variables. (D) True positive vs. false positive rate. (E) Heritability ( $\eta$ ) for **twostep** is estimated as  $\sigma_g^2 / (\sigma_g^2 + \sigma_e^2)$  from an intercept only LMM with a single random effect where  $\sigma_g^2$  and  $\sigma_e^2$  are the variance components for the random effect and error term, respectively.  $\eta$  is explicitly modeled in **ggmix**. There is no positive way to calculate  $\eta$  for the **lasso** since we are using a PC adjustment. (F) Error variance ( $\sigma^2$ ) for **twostep** is estimated from an intercept only LMM with a single random effect and is modeled explicitly in **ggmix**. For the **lasso** we use  $\frac{1}{n-|\hat{S}_\lambda|} \left\| \mathbf{Y} - \mathbf{X} \hat{\beta}_\lambda \right\|_2^2$  [29] as an estimator for  $\sigma^2$ .

## 213 3.2 Real Data Applications

214 Two datasets with contrasting features are used to illustrate the potential advantages of  
215 `ggmix` over existing approaches such as PC adjustment in a `lasso` regression. In one dataset,  
216 family structure induces low levels of correlation and sparsity in signals. In the second, a  
217 dataset involving mouse crosses, correlations are extremely strong and can confound sig-  
218 nals.

### 219 3.2.1 GAW20

220 In the most recent Genetic Analysis Workshop 20 (GAW20), the causal modeling group in-  
221 vestigated causal relationships between DNA methylation (exposure) within some genes and  
222 the change in high-density lipoproteins  $\Delta$ HDL (outcome) using Mendelian randomization  
223 (MR) [30]. Penalized regression methods were used to select SNPs strongly associated with  
224 the exposure in order to be used as an instrumental variable (IV) [31, 32]. However, since  
225 GAW20 data consisted of families, `twostep` methods were used which could have resulted  
226 in a large number of false positives or false negatives. `ggmix` is an alternative approach that  
227 could be used for selecting the IV while accounting for the family structure of the data.

228 We applied `ggmix` to all 200 GAW20 simulation datasets, each of 679 observations, and com-  
229 pared its performance to the `twostep` and `lasso` methods. Using a FaST-LMM (Factored  
230 Spectrally Transformed Linear Mixed Model) [33], we validated the effect of rs9661059 on  
231 blood lipid trait to be significant (genome-wide  $p = 6.29 \times 10^{-9}$ ). Though several other SNPs  
232 are also associated with the phenotype, these associations are probably mediated by CpG-  
233 SNP interaction pairs and do not reach statistical significance. Therefore, to avoid ambiguity,  
234 we only focused on chromosome 1 containing 51,104 SNPs where rs9661059 resides. Given  
235 that population admixture in the GAW20 data is likely, we estimated the population kinship  
236 using REAP [34] after decomposing population compositions using ADMIXTURE [35]. We  
237 supplied the estimated kinship matrix directly to `ggmix`. For both the `lasso` and `twostep`

238 methods, we adopted the same strategies as described in our simulation study in section 3.1,  
239 supplying the same kinship matrix estimated by REAP.

240 On each simulated replicate, we calibrated the methods so that they could be easily compared  
241 by fixing the true positive rate to 1 and then minimizing the false positive rate. Hence, the  
242 selected SNP, rs9661059, is likely to be the true positive for each method, and non-causal  
243 SNPs are excluded to the greatest extent. All of the three methods precisely choose the  
244 correct predictor without any false positives in more than half of the replicates, given the  
245 strong causal signal. When some false positives are selected, `ggmix` performs comparably  
246 to `twostep`, and the `lasso` tends to select more false positives (Figure 4). In terms of  
247 phenotype prediction, we observed that `ggmix` outperforms the `twostep` method without  
248 requiring more SNPs, while it achieves roughly the same prediction accuracy as `lasso` but  
249 with fewer non-causal SNPs (Figure 4).

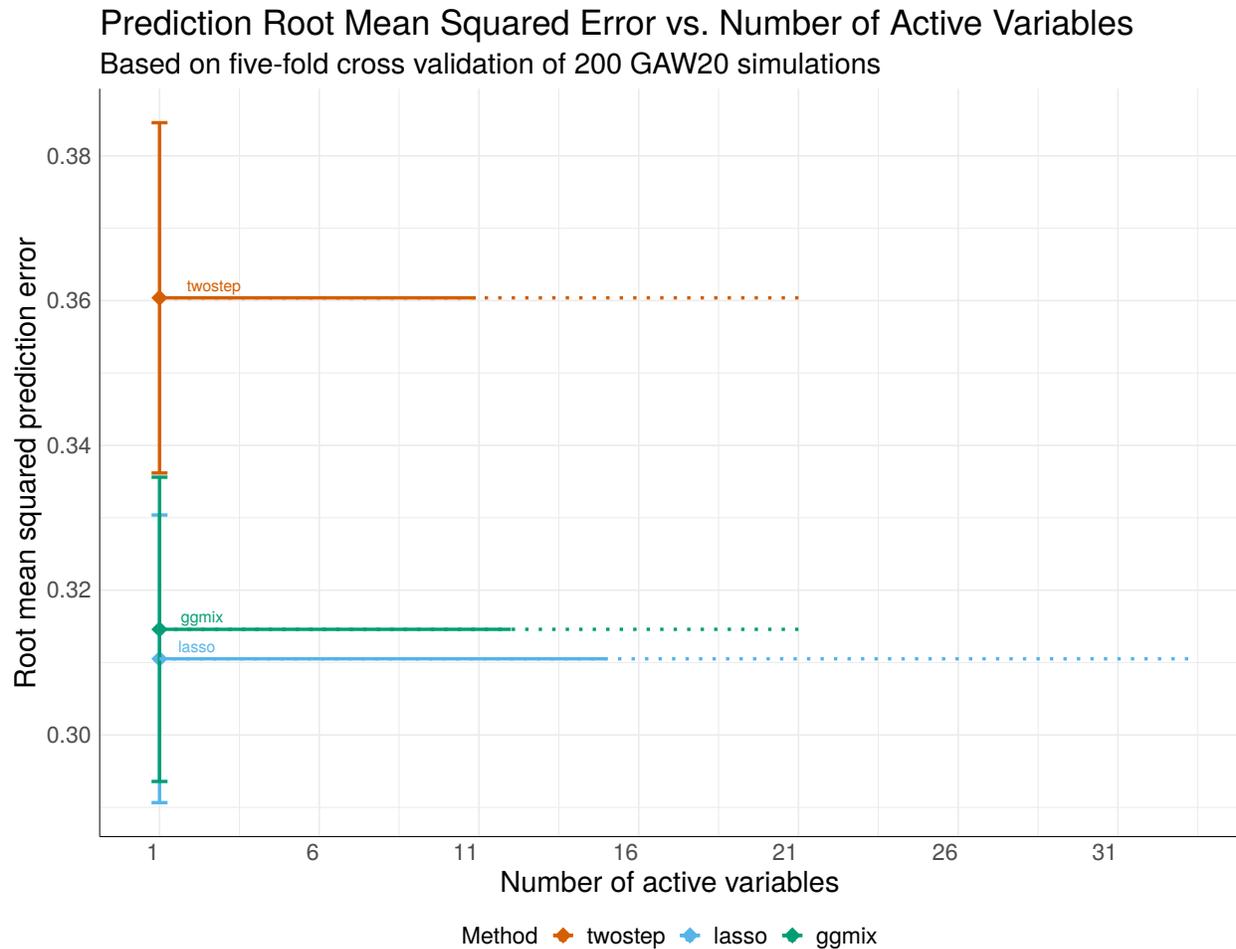


Figure 4: Mean  $\pm 1$  standard deviation of root mean squared error vs. number of active variables used by each method on the GAW20 data. Diamonds represent median number of active variables and the corresponding root mean square error. Horizontal solid lines span from median to the 90th percentile; Horizontal dotted lines span from the 90th percentile to the 95th percentile.

### 250 3.2.2 Mouse Crosses and Sensitivity to Mycobacterial Infection

251 Mouse inbred strains of genetically identical individuals are extensively used in research.  
252 Crosses of different inbred strains are useful for various studies of heritability focusing on  
253 either observable phenotypes or molecular mechanisms, and in particular, recombinant con-  
254 genic strains have been an extremely useful resource for many years [36]. However, ignor-  
255 ing complex genetic relationships in association studies can lead to inflated false positives

256 in genetic association studies when different inbred strains and their crosses are investi-  
257 gated [37, 38, 39]. Therefore, a previous study developed and implemented a mixed model  
258 to find loci associated with mouse sensitivity to mycobacterial infection [40]. The random  
259 effects in the model captured complex correlations between the recombinant congenic mouse  
260 strains based on the proportion of the DNA shared identical by descent. Through a se-  
261 ries of mixed model fits at each marker, new loci that impact growth of mycobacteria on  
262 chromosome 1 and chromosome 11 were identified.

263 Here we show that `ggmix` can identify these loci, as well as potentially others, in a single  
264 analysis. We reanalyzed the growth permissiveness in the spleen, as measured by colony  
265 forming units (CFUs), 6 weeks after infection from *Mycobacterium bovis* Bacille Calmette-  
266 Guerin (BCG) Russia strain as reported in [40].

267 By taking the consensus between the “main model” and the “conditional model” of the  
268 original study, we regarded markers D1Mit435 on chromosome 1 and D11Mit119 on chromo-  
269 some 11 as two true positive loci. Similar to the strategy used when analyzing the GAW20  
270 data, we optimized models by tuning the penalty factor such that these two loci are picked  
271 up, while the number of other active loci is minimized. To evaluate robustness of different  
272 models, we bootstrapped the 189-sample dataset and repeated the analysis 200 times. We  
273 directly estimated the kinship between mice using genotypes at 625 microsatellite markers.  
274 The estimated kinship entered directly into `ggmix` and `twostep`. For the `lasso`, we calcu-  
275 lated and included the first 10 principal components of the estimated kinship. Significant  
276 markers are defined as those captured in at least half of the bootstrap replicates, and in  
277 which the corresponding method successfully captures both pre-selected true positives with  
278 a penalty factor minimizing the number of active loci (Figure 5).

279 We demonstrate that `ggmix` recognizes the true associations more robustly than `twostep`  
280 and `lasso`. In almost all (99%) bootstrap replicates, `ggmix` is able to capture both true  
281 positives, while `twostep` failed in 19% of the replicates and `lasso` failed in 56% of the

282 replicates by missing of at least one of the two true positives (Figure 5). We also identified  
283 several other loci that might also be associated with susceptibility to mycobacterial infection  
284 (Table 1). Among these new potentially-associated markers, D2Mit156 was found to play a  
285 role in control of parasite numbers of *Leishmania tropica* in lymph nodes [41]. This locus is  
286 considered significant by our definition for both `ggmix` and `lasso`. An earlier study identified  
287 a parent-of-origin effect at D17Mit221 on CD4M levels [42]. This effect was more visible in  
288 crosses than in parental strains. In addition, D14Mit131, selected only by `ggmix`, was found  
289 to have a 9% loss of heterozygosity in hybrids of two inbred mouse strains [43], indicating the  
290 potential presence of putative suppressor genes pertaining to immune surveillance and tumor  
291 progression [44]. This result might also suggest association with anti-bacterial responses yet  
292 to be discovered.

### 3 RESULTS

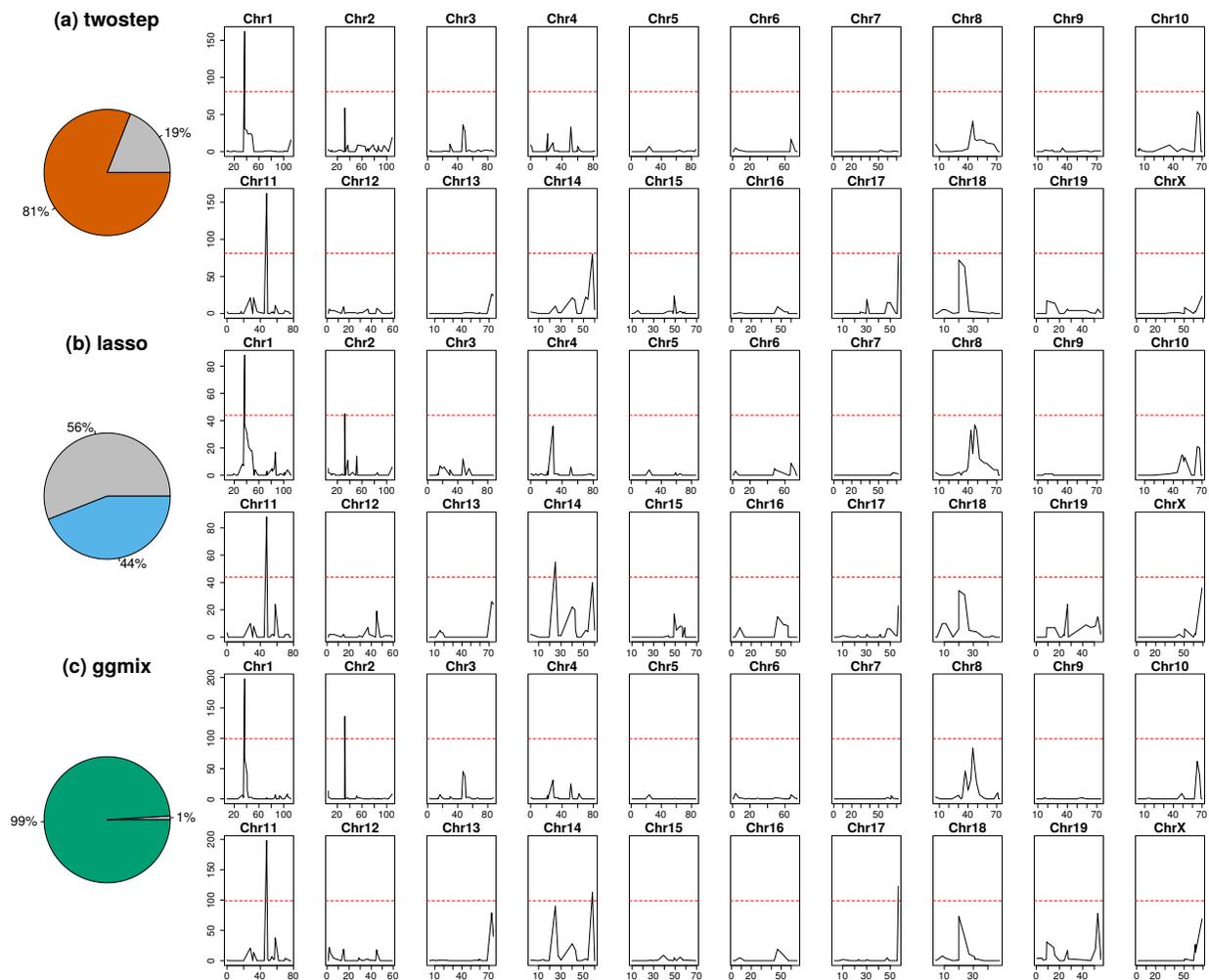


Figure 5: Comparison of model performance on the mouse cross data. Pie charts depict model robustness where grey areas denote bootstrap replicates on which the corresponding model is unable to capture both true positives using any penalty factor, whereas colored areas denote successful replicates. Chromosome-based signals record in how many successful replicates the corresponding loci are picked up by the corresponding optimized model. Red dashed lines delineate  $p$  value thresholds.

Table 1: Additional loci significantly associated with mouse susceptibility to myobacterial infection, after excluding two true positives. Loci needed to be identified in at least 50% of the successful bootstrap replicates that captured both true positive loci.

Method	Marker	Position in cM	Position in bp
<code>twostep</code>	N/A	N/A	N/A
<code>lasso</code>	D2Mit156	Chr2:31.66	Chr2:57081653-57081799
	D14Mit155	Chr14:31.52	Chr14:59828398-59828596
<code>ggmix</code>	D2Mit156	Chr2:31.66	Chr2:57081653-57081799
	D14Mit131	Chr14:63.59	Chr14:120006565-120006669
	D17Mit221	Chr17:59.77	Chr17:90087704-90087842

293

294

## 4 Discussion

295

296

297

298

299

300

301

302

303

304

305

306

307

We have developed a general penalized LMM framework called `ggmix` which simultaneously selects SNPs and adjusts for population structure in high dimensional prediction models. Through an extensive simulation study and two real data analyses, we show that the current approaches of PC adjustment and two-stage procedures are not necessarily sufficient to control for confounding by population structure leading to a high number of false positives or false negatives. Furthermore, `ggmix` showed improved prediction performance with a more parsimonious model compared to both the `lasso` and `twostep`. Our proposed method has excellent Type 1 error control and is robust to the inclusion of causal SNPs in the kinship matrix. Many methods for single-SNP analyses avoid this “proximal contamination” [8] by using a leave-one-chromosome-out scheme [45], i.e., construct the kinship matrix using all chromosomes except the one on which the marker being tested is located. However, this approach is not possible if we want to model many SNPs (across many chromosomes) jointly. We also demonstrated `ggmix` using two examples that mimic many experimental designs in

308 genetics. In the GAW20 example, we showed that while all methods were able to select  
309 the strongest causal SNP, `ggmix` did so with the least amount of false positives while also  
310 maintaining good predictive ability. In the mouse crosses example, we showed that `ggmix` is  
311 robust to perturbations in the data using a bootstrap analysis. Indeed, `ggmix` was able to  
312 consistently select the true positives across bootstrap replicates, while `twostep` failed in 19%  
313 of the replicates and `lasso` failed in 56% of the replicates by missing of at least one of the  
314 two true positives. Our re-analysis of the data also lead to some potentially new findings,  
315 not found by existing methods, that may warrant further study.

316 We emphasize here that previously developed methods such as the LMM-lasso [15] use a two-  
317 stage fitting procedure without any convergence details. From a practical point of view, there  
318 is currently no implementation that provides a principled way of determining the sequence  
319 of tuning parameters to fit, nor a procedure that automatically selects the optimal value of  
320 the tuning parameter. To our knowledge, we are the first to develop a coordinate gradient  
321 descent (CGD) algorithm in the specific context of fitting a penalized LMM for population  
322 structure correction with theoretical guarantees of convergence. Furthermore, we develop  
323 a principled method for automatic tuning parameter selection and provide an easy-to-use  
324 software implementation in order to promote wider uptake of these more complex methods  
325 by applied practitioners.

326 Although we derive a CGD algorithm for the  $\ell_1$  penalty, our approach can also be easily  
327 extended to other penalties such as the elastic net and group lasso with the same guarantees  
328 of convergence. A limitation of `ggmix` is that it first requires computing the covariance ma-  
329 trix with a computation time of  $\mathcal{O}(n^2k)$  followed by a spectral decomposition of this matrix  
330 in  $\mathcal{O}(n^3)$  time where  $k$  is the number of SNP genotypes used to construct the covariance  
331 matrix. This computation becomes prohibitive for large cohorts such as the UK Biobank [46]  
332 which have collected genetic information on half a million individuals. When the matrix of  
333 genotypes used to construct the covariance matrix is low rank, there are additional computa-

334 tional speedups that can be implemented. While this has been developed for the univariate  
335 case [8], to our knowledge, this has not been explored in the multivariable case. We are cur-  
336 rently developing a low rank version of the penalized LMM developed here, which reduces  
337 the time complexity from  $\mathcal{O}(n^2k)$  to  $\mathcal{O}(nk^2)$ .

338 There are other applications in which our method could be used as well. For example, there  
339 has been a renewed interest in polygenic risk scores (PRS) which aim to predict complex  
340 diseases from genotypes. `ggmix` could be used to build a PRS with the distinct advantage  
341 of modeling SNPs jointly, allowing for main effects as well as interactions to be accounted  
342 for. Based on our results, `ggmix` has the potential to produce more robust and parsimonious  
343 models than the `lasso` with better predictive accuracy. Our method is also suitable for fine  
344 mapping SNP association signals in genomic regions, where the goal is to pinpoint individual  
345 variants most likely to impact the underlying biological mechanisms of disease [47].

## 346 5 Materials and Methods

### 347 5.1 Model Set-up

348 Let  $i = 1, \dots, N$  be a grouping index,  $j = 1, \dots, n_i$  the observation index within a group  
349 and  $N_T = \sum_{i=1}^N n_i$  the total number of observations. For each group let  $\mathbf{y}_i = (y_1, \dots, y_{n_i})$  be  
350 the observed vector of responses or phenotypes,  $\mathbf{X}_i$  an  $n_i \times (p + 1)$  design matrix (with  
351 the column of 1s for the intercept),  $\mathbf{b}_i$  a group-specific random effect vector of length  
352  $n_i$  and  $\boldsymbol{\varepsilon}_i = (\varepsilon_{i1}, \dots, \varepsilon_{in_i})$  the individual error terms. Denote the stacked vectors  $\mathbf{Y} =$   
353  $(\mathbf{y}_1, \dots, \mathbf{y}_N)^T \in \mathbb{R}^{N_T \times 1}$ ,  $\mathbf{b} = (\mathbf{b}_1, \dots, \mathbf{b}_N)^T \in \mathbb{R}^{N_T \times 1}$ ,  $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_N)^T \in \mathbb{R}^{N_T \times 1}$ , and the  
354 stacked matrix  
355  $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_N)^T \in \mathbb{R}^{N_T \times (p+1)}$ . Furthermore, let  $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^T \in \mathbb{R}^{(p+1) \times 1}$  be a  
356 vector of fixed effects regression coefficients corresponding to  $\mathbf{X}$ . We consider the following

357 linear mixed model with a single random effect [48]:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{b} + \boldsymbol{\varepsilon} \quad (3)$$

358 where the random effect  $\mathbf{b}$  and the error variance  $\boldsymbol{\varepsilon}$  are assigned the distributions

$$\mathbf{b} \sim \mathcal{N}(0, \eta\sigma^2\boldsymbol{\Phi}) \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, (1 - \eta)\sigma^2\mathbf{I}) \quad (4)$$

359 Here,  $\boldsymbol{\Phi}_{N_T \times N_T}$  is a known positive semi-definite and symmetric covariance or kinship ma-  
 360 trix calculated from SNPs sampled across the genome,  $\mathbf{I}_{N_T \times N_T}$  is the identity matrix and  
 361 parameters  $\sigma^2$  and  $\eta \in [0, 1]$  determine how the variance is divided between  $\mathbf{b}$  and  $\boldsymbol{\varepsilon}$ . Note  
 362 that  $\eta$  is also the narrow-sense heritability ( $h^2$ ), defined as the proportion of phenotypic  
 363 variance attributable to the additive genetic factors [1]. The joint density of  $\mathbf{Y}$  is therefore  
 364 multivariate normal:

$$\mathbf{Y} | (\boldsymbol{\beta}, \eta, \sigma^2) \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \eta\sigma^2\boldsymbol{\Phi} + (1 - \eta)\sigma^2\mathbf{I}) \quad (5)$$

365 The LMM-Lasso method [15] considers an alternative but equivalent parameterization given  
 366 by:

$$\mathbf{Y} | (\boldsymbol{\beta}, \delta, \sigma_g^2) \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma_g^2(\boldsymbol{\Phi} + \delta\mathbf{I})) \quad (6)$$

where  $\delta = \sigma_e^2/\sigma_g^2$ ,  $\sigma_g^2$  is the genetic variance and  $\sigma_e^2$  is the residual variance. We instead  
 consider the parameterization in (5) since maximization is easier over the compact set  $\eta \in$   
 $[0, 1]$  than over the unbounded interval  $\delta \in [0, \infty)$  [48]. We define the complete parameter  
 vector as  $\boldsymbol{\Theta} := (\boldsymbol{\beta}, \eta, \sigma^2)$ . The negative log-likelihood for (5) is given by

$$-\ell(\boldsymbol{\Theta}) \propto \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \log(\det(\mathbf{V})) + \frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \quad (7)$$

367 where  $\mathbf{V} = \eta\Phi + (1 - \eta)\mathbf{I}$  and  $\det(\mathbf{V})$  is the determinant of  $\mathbf{V}$ .

Let  $\Phi = \mathbf{U}\mathbf{D}\mathbf{U}^T$  be the eigen (spectral) decomposition of the kinship matrix  $\Phi$ , where  $\mathbf{U}_{N_T \times N_T}$  is an orthonormal matrix of eigenvectors (i.e.  $\mathbf{U}\mathbf{U}^T = \mathbf{I}$ ) and  $\mathbf{D}_{N_T \times N_T}$  is a diagonal matrix of eigenvalues  $\Lambda_i$ .  $\mathbf{V}$  can then be further simplified [48]

$$\begin{aligned}
 \mathbf{V} &= \eta\Phi + (1 - \eta)\mathbf{I} \\
 &= \eta\mathbf{U}\mathbf{D}\mathbf{U}^T + (1 - \eta)\mathbf{U}\mathbf{I}\mathbf{U}^T \\
 &= \mathbf{U}\eta\mathbf{D}\mathbf{U}^T + \mathbf{U}(1 - \eta)\mathbf{I}\mathbf{U}^T \\
 &= \mathbf{U}(\eta\mathbf{D} + (1 - \eta)\mathbf{I})\mathbf{U}^T \\
 &= \mathbf{U}\tilde{\mathbf{D}}\mathbf{U}^T
 \end{aligned} \tag{8}$$

where

$$\tilde{\mathbf{D}} = \eta\mathbf{D} + (1 - \eta)\mathbf{I} \tag{9}$$

$$\begin{aligned}
 &= \eta \begin{bmatrix} \Lambda_1 & & & \\ & \Lambda_2 & & \\ & & \ddots & \\ & & & \Lambda_{N_T} \end{bmatrix} + (1 - \eta) \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \\
 &= \begin{bmatrix} 1 + \eta(\Lambda_1 - 1) & & & \\ & 1 + \eta(\Lambda_2 - 1) & & \\ & & \ddots & \\ & & & 1 + \eta(\Lambda_{N_T} - 1) \end{bmatrix} \\
 &= \text{diag} \{1 + \eta(\Lambda_1 - 1), 1 + \eta(\Lambda_2 - 1), \dots, 1 + \eta(\Lambda_{N_T} - 1)\}
 \end{aligned} \tag{10}$$

Since (9) is a diagonal matrix, its inverse is also a diagonal matrix:

$$\tilde{\mathbf{D}}^{-1} = \text{diag} \left\{ \frac{1}{1 + \eta(\Lambda_1 - 1)}, \frac{1}{1 + \eta(\Lambda_2 - 1)}, \dots, \frac{1}{1 + \eta(\Lambda_{N_T} - 1)} \right\} \quad (11)$$

From (8) and (10),  $\log(\det(\mathbf{V}))$  simplifies to

$$\begin{aligned} \log(\det(\mathbf{V})) &= \log \left( \det(\mathbf{U}) \det(\tilde{\mathbf{D}}) \det(\mathbf{U}^T) \right) \\ &= \log \left\{ \prod_{i=1}^{N_T} (1 + \eta(\Lambda_i - 1)) \right\} \\ &= \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) \end{aligned} \quad (12)$$

since  $\det(\mathbf{U}) = 1$ . It also follows from (8) that

$$\begin{aligned} \mathbf{V}^{-1} &= (\mathbf{U}\tilde{\mathbf{D}}\mathbf{U}^T)^{-1} \\ &= (\mathbf{U}^T)^{-1} (\tilde{\mathbf{D}})^{-1} \mathbf{U}^{-1} \\ &= \mathbf{U}\tilde{\mathbf{D}}^{-1}\mathbf{U}^T \end{aligned} \quad (13)$$

since for an orthonormal matrix  $\mathbf{U}^{-1} = \mathbf{U}^T$ . Substituting (11), (12) and (13) into (7) the negative log-likelihood becomes

$$-\ell(\Theta) \propto \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\beta)^T \mathbf{U}\tilde{\mathbf{D}}^{-1}\mathbf{U}^T (\mathbf{Y} - \mathbf{X}\beta) \quad (14)$$

$$\begin{aligned} &= \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^2} (\mathbf{U}^T\mathbf{Y} - \mathbf{U}^T\mathbf{X}\beta)^T \tilde{\mathbf{D}}^{-1} (\mathbf{U}^T\mathbf{Y} - \mathbf{U}^T\mathbf{X}\beta) \\ &= \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^2} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\beta)^T \tilde{\mathbf{D}}^{-1} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\beta) \\ &= \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^2} \sum_{i=1}^{N_T} \frac{\left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1}\beta_j \right)^2}{1 + \eta(\Lambda_i - 1)} \end{aligned} \quad (15)$$

368 where  $\tilde{\mathbf{Y}} = \mathbf{U}^T \mathbf{Y}$ ,  $\tilde{\mathbf{X}} = \mathbf{U}^T \mathbf{X}$ ,  $\tilde{Y}_i$  denotes the  $i^{\text{th}}$  element of  $\tilde{\mathbf{Y}}$ ,  $\tilde{X}_{ij}$  is the  $i, j^{\text{th}}$  entry of  $\tilde{\mathbf{X}}$   
369 and  $\mathbf{1}$  is a column vector of  $N_T$  ones.

## 370 5.2 Penalized Maximum Likelihood Estimator

371 We define the  $p + 3$  length vector of parameters  $\Theta := (\Theta_0, \Theta_1, \dots, \Theta_{p+1}, \Theta_{p+2}, \Theta_{p+3}) =$   
372  $(\beta, \eta, \sigma^2)$  where  $\beta \in \mathbb{R}^{p+1}$ ,  $\eta \in [0, 1]$ ,  $\sigma^2 > 0$ . In what follows,  $p + 2$  and  $p + 3$  are the indices  
373 in  $\Theta$  for  $\eta$  and  $\sigma^2$ , respectively. In light of our goals to select variables associated with the  
374 response in high-dimensional data, we propose to place a constraint on the magnitude of  
375 the regression coefficients. This can be achieved by adding a penalty term to the likelihood  
376 function (15). The penalty term is a necessary constraint because in our applications, the  
377 sample size is much smaller than the number of predictors. We define the following objective  
378 function:

$$Q_\lambda(\Theta) = f(\Theta) + \lambda \sum_{j \neq 0} v_j P_j(\beta_j) \quad (16)$$

379 where  $f(\Theta) := -\ell(\Theta)$  is defined in (15),  $P_j(\cdot)$  is a penalty term on the fixed regression  
380 coefficients  $\beta_1, \dots, \beta_{p+1}$  (we do not penalize the intercept) controlled by the nonnegative  
381 regularization parameter  $\lambda$ , and  $v_j$  is the penalty factor for  $j^{\text{th}}$  covariate. These penalty  
382 factors serve as a way of allowing parameters to be penalized differently. Note that we do  
383 not penalize  $\eta$  or  $\sigma^2$ . An estimate of the regression parameters  $\hat{\Theta}_\lambda$  is obtained by

$$\hat{\Theta}_\lambda = \arg \min_{\Theta} Q_\lambda(\Theta) \quad (17)$$

384 This is the general set-up for our model. In Section 5.3 we provide more specific details on  
385 how we solve (17).

### 386 5.3 Computational Algorithm

387 We use a general purpose block coordinate gradient descent algorithm (CGD) [49] to solve (17).  
 388 At each iteration, we cycle through the coordinates and minimize the objective function with  
 389 respect to one coordinate only. For continuously differentiable  $f(\cdot)$  and convex and block-  
 390 separable  $P(\cdot)$  (i.e.  $P(\boldsymbol{\beta}) = \sum_i P_i(\beta_i)$ ), Tseng and Yun [49] show that the solution gener-  
 391 ated by the CGD method is a stationary point of  $Q_\lambda(\cdot)$  if the coordinates are updated in a  
 392 Gauss-Seidel manner i.e.  $Q_\lambda(\cdot)$  is minimized with respect to one parameter while holding  
 393 all others fixed. The CGD algorithm has been successfully applied in fixed effects models  
 394 (e.g. [50], [20]) and linear mixed models with an  $\ell_1$  penalty [51]. In the next section we  
 395 provide some brief details about Algorithm 1. A more thorough treatment of the algorithm  
 396 is given in Appendix A.

---

**Algorithm 1:** Block Coordinate Gradient Descent

---

Set the iteration counter  $k \leftarrow 0$ , initial values for the parameter vector  $\Theta^{(0)}$  and convergence threshold  $\epsilon$ ;  
**for**  $\lambda \in \{\lambda_{max}, \dots, \lambda_{min}\}$  **do**  
     **repeat**  
          $\boldsymbol{\beta}^{(k+1)} \leftarrow \arg \min_{\boldsymbol{\beta}} Q_\lambda(\boldsymbol{\beta}, \eta^{(k)}, \sigma^2^{(k)})$   
          $\eta^{(k+1)} \leftarrow \arg \min_{\eta} Q_\lambda(\boldsymbol{\beta}^{(k+1)}, \eta, \sigma^2^{(k)})$   
          $\sigma^2^{(k+1)} \leftarrow \arg \min_{\sigma^2} Q_\lambda(\boldsymbol{\beta}^{(k+1)}, \eta^{(k+1)}, \sigma^2)$   
          $k \leftarrow k + 1$   
     **until** convergence criterion is satisfied:  $\|\Theta^{(k+1)} - \Theta^{(k)}\|_2 < \epsilon$ ;  
**end**

---

397 **5.3.1 Updates for the  $\beta$  parameter**

398 Recall that the part of the objective function that depends on  $\beta$  has the form

$$Q_\lambda(\Theta) = \frac{1}{2} \sum_{i=1}^{N_T} w_i \left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2 + \lambda \sum_{j=1}^p v_j |\beta_j| \quad (18)$$

399 where

$$w_i := \frac{1}{\sigma^2 (1 + \eta(\Lambda_i - 1))} \quad (19)$$

Conditional on  $\eta^{(k)}$  and  $\sigma^{2(k)}$ , it can be shown that the solution for  $\beta_j$ ,  $j = 1, \dots, p$  is given by

$$\beta_j^{(k+1)} \leftarrow \frac{\mathcal{S}_\lambda \left( \sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) \right)}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \quad (20)$$

where  $\mathcal{S}_\lambda(x)$  is the soft-thresholding operator

$$\mathcal{S}_\lambda(x) = \text{sign}(x)(|x| - \lambda)_+$$

400  $\text{sign}(x)$  is the signum function

$$\text{sign}(x) = \begin{cases} -1 & x < 0 \\ 0 & x = 0 \\ 1 & x > 0 \end{cases}$$

401 and  $(x)_+ = \max(x, 0)$ . We provide the full derivation in Appendix [A.1.2](#).

### 402 5.3.2 Updates for the $\eta$ paramter

403 Given  $\beta^{(k+1)}$  and  $\sigma^2^{(k)}$ , solving for  $\eta^{(k+1)}$  becomes a univariate optimization problem:

$$\eta^{(k+1)} \leftarrow \arg \min_{\eta} \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^2^{(k)}} \sum_{i=1}^{N_T} \frac{\left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j^{(k+1)}\right)^2}{1 + \eta(\Lambda_i - 1)} \quad (21)$$

404 We use a bound constrained optimization algorithm [52] implemented in the `optim` function  
405 in R and set the lower and upper bounds to be 0.01 and 0.99, respectively.

### 406 5.3.3 Updates for the $\sigma^2$ parameter

407 Conditional on  $\beta^{(k+1)}$  and  $\eta^{(k+1)}$ ,  $\sigma^2^{(k+1)}$  can be solved for using the following equation:

$$\sigma^2^{(k+1)} \leftarrow \arg \min_{\sigma^2} \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^{N_T} \frac{\left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j\right)^2}{1 + \eta(\Lambda_i - 1)} \quad (22)$$

There exists an analytic solution for (22) given by:

$$\sigma^2^{(k+1)} \leftarrow \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{\left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j^{(k+1)}\right)^2}{1 + \eta^{(k+1)}(\Lambda_i - 1)} \quad (23)$$

### 408 5.3.4 Regularization path

409 In this section we describe how determine the sequence of tuning parameters  $\lambda$  at which to  
410 fit the model. Recall that our objective function has the form

$$Q_{\lambda}(\Theta) = \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2} \sum_{i=1}^{N_T} w_i \left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j\right)^2 + \lambda \sum_{j=1}^p v_j |\beta_j| \quad (24)$$

411 The Karush-Kuhn-Tucker (KKT) optimality conditions for (24) are given by:

$$\begin{aligned}
 \frac{\partial}{\partial \beta_1, \dots, \beta_p} Q_\lambda(\Theta) &= \mathbf{0}_p \\
 \frac{\partial}{\partial \beta_0} Q_\lambda(\Theta) &= 0 \\
 \frac{\partial}{\partial \eta} Q_\lambda(\Theta) &= 0 \\
 \frac{\partial}{\partial \sigma^2} Q_\lambda(\Theta) &= 0
 \end{aligned}
 \tag{25}$$

412 The equations in (25) are equivalent to

$$\begin{aligned}
 \sum_{i=1}^{N_T} w_i \tilde{X}_{i1} \left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right) &= 0 \\
 \frac{1}{v_j} \sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right) &= \lambda \gamma_j, \\
 \gamma_j \in \begin{cases} \text{sign}(\hat{\beta}_j) & \text{if } \hat{\beta}_j \neq 0 \\ [-1, 1] & \text{if } \hat{\beta}_j = 0 \end{cases}, & \text{for } j = 1, \dots, p \\
 \frac{1}{2} \sum_{i=1}^{N_T} \frac{\Lambda_i - 1}{1 + \eta(\Lambda_i - 1)} \left( 1 - \frac{\left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2}{\sigma^2(1 + \eta(\Lambda_i - 1))} \right) &= 0 \\
 \sigma^2 - \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{\left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2}{1 + \eta(\Lambda_i - 1)} &= 0
 \end{aligned}
 \tag{26}$$

413 where  $w_i$  is given by (19),  $\tilde{\mathbf{X}}_{-1}^T$  is  $\tilde{\mathbf{X}}^T$  with the first column removed,  $\tilde{\mathbf{X}}_1^T$  is the first column  
414 of  $\tilde{\mathbf{X}}^T$ , and  $\gamma \in \mathbb{R}^p$  is the subgradient function of the  $\ell_1$  norm evaluated at  $(\hat{\beta}_1, \dots, \hat{\beta}_p)$ .  
415 Therefore  $\hat{\Theta}$  is a solution in (17) if and only if  $\hat{\Theta}$  satisfies (26) for some  $\gamma$ . We can determine  
416 a decreasing sequence of tuning parameters by starting at a maximal value for  $\lambda = \lambda_{max}$   
417 for which  $\hat{\beta}_j = 0$  for  $j = 1, \dots, p$ . In this case, the KKT conditions in (26) are equivalent

418 to

$$\begin{aligned}
 \frac{1}{v_j} \sum_{i=1}^{N_T} \left| w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \tilde{X}_{i1} \beta_0 \right) \right| &\leq \lambda, \quad \forall j = 1, \dots, p \\
 \beta_0 &= \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{i1} \tilde{Y}_i}{\sum_{i=1}^{N_T} w_i \tilde{X}_{i1}^2} \\
 \frac{1}{2} \sum_{i=1}^{N_T} \frac{\Lambda_i - 1}{1 + \eta(\Lambda_i - 1)} \left( 1 - \frac{\left( \tilde{Y}_i - \tilde{X}_{i1} \beta_0 \right)^2}{\sigma^2 (1 + \eta(\Lambda_i - 1))} \right) &= 0 \\
 \sigma^2 &= \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{\left( \tilde{Y}_i - \tilde{X}_{i1} \beta_0 \right)^2}{1 + \eta(\Lambda_i - 1)}
 \end{aligned} \tag{27}$$

419 We can solve the KKT system of equations in (27) (with a numerical solution for  $\eta$ ) in order  
 420 to have an explicit form of the stationary point  $\hat{\Theta}_0 = \{\hat{\beta}_0, \mathbf{0}_p, \hat{\eta}, \hat{\sigma}^2\}$ . Once we have  $\hat{\Theta}_0$ , we  
 421 can solve for the smallest value of  $\lambda$  such that the entire vector  $(\hat{\beta}_1, \dots, \hat{\beta}_p)$  is 0:

$$\lambda_{max} = \max_j \left\{ \left| \frac{1}{v_j} \sum_{i=1}^{N_T} \hat{w}_i \tilde{X}_{ij} \left( \tilde{Y}_i - \tilde{X}_{i1} \hat{\beta}_0 \right) \right| \right\}, \quad j = 1, \dots, p \tag{28}$$

422 Following Friedman et al. [20], we choose  $\tau \lambda_{max}$  to be the smallest value of tuning parameters  
 423  $\lambda_{min}$ , and construct a sequence of  $K$  values decreasing from  $\lambda_{max}$  to  $\lambda_{min}$  on the log scale.  
 424 The defaults are set to  $K = 100$ ,  $\tau = 0.01$  if  $n < p$  and  $\tau = 0.001$  if  $n \geq p$ .

### 425 5.3.5 Warm Starts

426 The way in which we have derived the sequence of tuning parameters using the KKT con-  
 427 ditions, allows us to implement warm starts. That is, the solution  $\hat{\Theta}$  for  $\lambda_k$  is used as the  
 428 initial value  $\Theta^{(0)}$  for  $\lambda_{k+1}$ . This strategy leads to computational speedups and has been  
 429 implemented in the `ggmix` R package.

### 430 5.3.6 Prediction of the random effects

431 We use an empirical Bayes approach (e.g. [53]) to predict the random effects  $\mathbf{b}$ . Let the  
432 maximum a posteriori (MAP) estimate be defined as

$$\hat{\mathbf{b}} = \arg \max_{\mathbf{b}} f(\mathbf{b}|\mathbf{Y}, \boldsymbol{\beta}, \eta, \sigma^2) \quad (29)$$

where, by using Bayes rule,  $f(\mathbf{b}|\mathbf{Y}, \boldsymbol{\beta}, \eta, \sigma^2)$  can be expressed as

$$\begin{aligned} f(\mathbf{b}|\mathbf{Y}, \boldsymbol{\beta}, \eta, \sigma^2) &= \frac{f(\mathbf{Y}|\mathbf{b}, \boldsymbol{\beta}, \eta, \sigma^2)\pi(\mathbf{b}|\eta, \sigma^2)}{f(\mathbf{Y}|\boldsymbol{\beta}, \eta, \sigma^2)} \\ &\propto f(\mathbf{Y}|\mathbf{b}, \boldsymbol{\beta}, \eta, \sigma^2)\pi(\mathbf{b}|\eta, \sigma^2) \\ &\propto \exp \left\{ -\frac{1}{2\sigma^2}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b})^T \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b}) - \frac{1}{2\eta\sigma^2} \mathbf{b}^T \boldsymbol{\Phi}^{-1} \mathbf{b} \right\} \\ &= \exp \left\{ -\frac{1}{2\sigma^2} \left[ (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b})^T \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b}) + \frac{1}{\eta} \mathbf{b}^T \boldsymbol{\Phi}^{-1} \mathbf{b} \right] \right\} \quad (30) \end{aligned}$$

Solving for (29) is equivalent to minimizing the exponent in (30):

$$\hat{\mathbf{b}} = \arg \min_{\mathbf{b}} \left\{ (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b})^T \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b}) + \frac{1}{\eta} \mathbf{b}^T \boldsymbol{\Phi}^{-1} \mathbf{b} \right\} \quad (31)$$

Taking the derivative of (31) with respect to  $\mathbf{b}$  and setting it to 0 we get:

$$\begin{aligned}
 0 &= -2\mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{b}) + \frac{2}{\eta}\boldsymbol{\Phi}^{-1}\mathbf{b} \\
 &= -\mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) + \left(\mathbf{V}^{-1} + \frac{1}{\eta}\boldsymbol{\Phi}^{-1}\right)\mathbf{b} \\
 \hat{\mathbf{b}} &= \left(\mathbf{V}^{-1} + \frac{1}{\hat{\eta}}\boldsymbol{\Phi}^{-1}\right)^{-1}\mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \\
 &= \left(\mathbf{U}\tilde{\mathbf{D}}^{-1}\mathbf{U}^T + \frac{1}{\hat{\eta}}\mathbf{U}\mathbf{D}^{-1}\mathbf{U}^T\right)^{-1}\mathbf{U}\tilde{\mathbf{D}}^{-1}\mathbf{U}^T(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \\
 &= \left(\mathbf{U}\left[\tilde{\mathbf{D}}^{-1} + \frac{1}{\hat{\eta}}\mathbf{D}^{-1}\right]\mathbf{U}^T\right)^{-1}\mathbf{U}\tilde{\mathbf{D}}^{-1}(\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\hat{\boldsymbol{\beta}}) \\
 &= \mathbf{U}\left[\tilde{\mathbf{D}}^{-1} + \frac{1}{\hat{\eta}}\mathbf{D}^{-1}\right]^{-1}\mathbf{U}^T\mathbf{U}\tilde{\mathbf{D}}^{-1}(\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\hat{\boldsymbol{\beta}})
 \end{aligned}$$

433 where  $\mathbf{V}^{-1}$  is given by (13), and  $(\hat{\boldsymbol{\beta}}, \hat{\eta})$  are the estimates obtained from Algorithm 1.

### 434 5.3.7 Phenotype prediction

435 Here we describe the method used for predicting the unobserved phenotype  $\mathbf{Y}^*$  in a set of  
 436 individuals with predictor set  $\mathbf{X}^*$  that were not used in the model training e.g. a testing  
 437 set. Let  $q$  denote the number of observations in the testing set and  $N - q$  the number of  
 438 observations in the training set. We assume that a `ggmix` model has been fit on a set of  
 439 training individuals with observed phenotype  $\mathbf{Y}$  and predictor set  $\mathbf{X}$ . We further assume  
 440 that  $\mathbf{Y}$  and  $\mathbf{Y}^*$  are jointly multivariate Normal:

$$\begin{bmatrix} \mathbf{Y}^* \\ \mathbf{Y} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \boldsymbol{\mu}_{1(q \times 1)} \\ \boldsymbol{\mu}_{2(N-q) \times 1} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11(q \times q)} & \boldsymbol{\Sigma}_{12_{q \times (N-q)}} \\ \boldsymbol{\Sigma}_{21_{(N-q) \times q}} & \boldsymbol{\Sigma}_{22_{(N-q) \times (N-q)}} \end{bmatrix} \right) \quad (32)$$

441 Then, from standard multivariate Normal theory, the conditional distribution  $\mathbf{Y}^* | \mathbf{Y}, \eta, \sigma^2, \boldsymbol{\beta}, \mathbf{X}, \mathbf{X}^*$   
 442 is  $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$  where

$$\boldsymbol{\mu}^* = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{Y} - \boldsymbol{\mu}_2) \quad (33)$$

$$\boldsymbol{\Sigma}^* = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21} \quad (34)$$

443 The phenotype prediction is thus given by:

$$\boldsymbol{\mu}_{q \times 1}^* = \mathbf{X}^*\boldsymbol{\beta} + \frac{1}{\sigma^2}\boldsymbol{\Sigma}_{12}\mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \quad (35)$$

$$= \mathbf{X}^*\boldsymbol{\beta} + \frac{1}{\sigma^2}\boldsymbol{\Sigma}_{12}\mathbf{U}\tilde{\mathbf{D}}^{-1}\mathbf{U}^T(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \quad (36)$$

$$= \mathbf{X}^*\boldsymbol{\beta} + \frac{1}{\sigma^2}\boldsymbol{\Sigma}_{12}\mathbf{U}\tilde{\mathbf{D}}^{-1}(\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}) \quad (37)$$

$$= \mathbf{X}^*\boldsymbol{\beta} + \frac{1}{\sigma^2}\eta\sigma^2\boldsymbol{\Phi}^*\mathbf{U}\tilde{\mathbf{D}}^{-1}(\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}) \quad (38)$$

$$= \mathbf{X}^*\boldsymbol{\beta} + \eta\boldsymbol{\Phi}^*\mathbf{U}\tilde{\mathbf{D}}^{-1}(\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}) \quad (39)$$

444 where  $\boldsymbol{\Phi}^*$  is the  $q \times (N - q)$  covariance matrix between the testing and training individu-  
445 als.

### 446 5.3.8 Choice of the optimal tuning parameter

447 In order to choose the optimal value of the tuning parameter  $\lambda$ , we use the generalized  
448 information criterion [54] (GIC):

$$GIC_\lambda = -2\ell(\hat{\boldsymbol{\beta}}, \hat{\sigma}^2, \hat{\eta}) + a_n \cdot \hat{df}_\lambda \quad (40)$$

449 where  $\hat{df}_\lambda$  is the number of non-zero elements in  $\hat{\boldsymbol{\beta}}_\lambda$  [55] plus two (representing the variance  
450 parameters  $\eta$  and  $\sigma^2$ ). Several authors have used this criterion for variable selection in mixed  
451 models with  $a_n = \log N_T$  [51, 56], which corresponds to the BIC. We instead choose the high-

## 5 MATERIALS AND METHODS

---

452 dimensional BIC [57] given by  $a_n = \log(\log(N_T)) * \log(p)$ . This is the default choice in our  
453 `ggmix` R package, though the interface is flexible to allow the user to select their choice of  
454  $a_n$ .

## 455 Availability of data and material

- 456 1. The GAW20 data is freely available upon request from <https://www.gaworkshop.org/data-sets>.
- 457
- 458 2. Mouse cross data is available from ES upon request.
- 459 3. The entire simulation study is reproducible. Source code available at <https://github.com/sahirbhatnagar/ggmix/tree/pgen/simulation>. This includes scripts for `ggmix`,
- 460 `lasso` and `twostep` methods.
- 461
- 462 4. The R package `ggmix` is freely available from GitHub at <https://github.com/greenwoodlab/ggmix>.
- 463
- 464 5. A website describing how to use the package is available at <https://sahirbhatnagar.com/ggmix/>.
- 465

## 466 Competing interests

467 The authors declare that they have no competing interests.

## 468 Author's contributions

469 SRB, KO, YY and CMTG conceived the idea. SRB developed the algorithms, software  
470 and simulation study. TL completed the real data analysis. ES and JCLO provided data  
471 and interpretations. SRB, TL and CMTG wrote a draft of the manuscript then all authors  
472 edited, read and approved the final manuscript.

## 473 Acknowledgements

474 SRB was supported by the Ludmer Centre for Neuroinformatics and Mental Health and  
475 the Canadian Institutes for Health Research PJT 148620. This research was enabled in  
476 part by support provided by Calcul Québec ([www.calculquebec.ca](http://www.calculquebec.ca)) and Compute Canada  
477 ([www.computeCanada.ca](http://www.computeCanada.ca)). The funders had no role in study design, data collection and  
478 analysis, decision to publish, or preparation of the manuscript.

## 479 Supporting Information

480 Contains the following sections:

481 A **Block Coordinate Descent Algorithm** - a detailed description of the algorithm  
482 used to fit our `ggmix` model

483 B **Additional Simulation Results** - complete simulation results

484 C **ggmix Package Showcase** - a vignette describing how to use our `ggmix` R package

## 485 References

486 [1] Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorff LA, Hunter DJ, et al. Finding  
487 the missing heritability of complex diseases. *Nature*. 2009;461(7265):747. 3, 22

488 [2] Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, et al. Common  
489 SNPs explain a large proportion of the heritability for human height. *Nature genetics*.  
490 2010;42(7):565. 3

491 [3] Astle W, Balding DJ, et al. Population structure and cryptic relatedness in genetic  
492 association studies. *Statistical Science*. 2009;24(4):451–471. 3, 4

## REFERENCES

- 493 [4] Song M, Hao W, Storey JD. Testing for genetic associations in arbitrarily structured  
494 populations. *Nature genetics*. 2015;47(5):550–554. 3
- 495 [5] Marchini J, Cardon LR, Phillips MS, Donnelly P. The effects of human population  
496 structure on large genetic association studies. *Nature genetics*. 2004;36(5):512. 3
- 497 [6] Hoggart CJ, Whittaker JC, De Iorio M, Balding DJ. Simultaneous analysis of all SNPs in  
498 genome-wide and re-sequencing association studies. *PLoS genetics*. 2008;4(7):e1000130.  
499 4
- 500 [7] Li J, Das K, Fu G, Li R, Wu R. The Bayesian lasso for genome-wide association studies.  
501 *Bioinformatics*. 2010;27(4):516–523. 4
- 502 [8] Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, Heckerman D. FaST linear  
503 mixed models for genome-wide association studies. *Nature methods*. 2011;8(10):833–  
504 835. 4, 8, 19, 21
- 505 [9] Kang HM, Sul JH, Zaitlen NA, Kong Sy, Freimer NB, Sabatti C, et al. Variance  
506 component model to account for sample structure in genome-wide association studies.  
507 *Nature genetics*. 2010;42(4):348. 4
- 508 [10] Yu J, Pressoir G, Briggs WH, Bi IV, Yamasaki M, Doebley JF, et al. A unified mixed-  
509 model method for association mapping that accounts for multiple levels of relatedness.  
510 *Nature genetics*. 2006;38(2):203. 4
- 511 [11] Eu-Ahsunthornwattana J, Miller EN, Fakiola M, Jeronimo SM, Blackwell JM, Cordell  
512 HJ, et al. Comparison of methods to account for relatedness in genome-wide association  
513 studies with family-based data. *PLoS Genet*. 2014;10(7):e1004445. 4
- 514 [12] Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Princi-  
515 pal components analysis corrects for stratification in genome-wide association studies.  
516 *Nature genetics*. 2006;38(8):904. 4

## REFERENCES

- 517 [13] Oualkacha K, Dastani Z, Li R, Cingolani PE, Spector TD, Hammond CJ, et al. Ad-  
518 justed sequence kernel association test for rare variants controlling for cryptic and family  
519 relatedness. *Genetic epidemiology*. 2013;37(4):366–376. 4, 5
- 520 [14] Cordell HJ, Clayton DG. A unified stepwise regression procedure for evaluating the  
521 relative effects of polymorphisms within a gene using case/control or family data:  
522 application to HLA in type 1 diabetes. *The American Journal of Human Genetics*.  
523 2002;70(1):124–141. 4
- 524 [15] Rakitsch B, Lippert C, Stegle O, Borgwardt K. A Lasso multi-marker mixed  
525 model for association mapping with population structure correction. *Bioinformatics*.  
526 2013;29(2):206–214. 4, 20, 22
- 527 [16] Wang D, Eskridge KM, Crossa J. Identifying QTLs and epistasis in structured plant  
528 populations using adaptive mixed LASSO. *Journal of agricultural, biological, and en-  
529 vironmental statistics*. 2011;16(2):170–184. 4
- 530 [17] Tibshirani R. Regression shrinkage and selection via the lasso. *Journal of the Royal  
531 Statistical Society Series B (Methodological)*. 1996;p. 267–288. 4, 5
- 532 [18] Zou H. The adaptive lasso and its oracle properties. *Journal of the American statistical  
533 association*. 2006;101(476):1418–1429. 5
- 534 [19] Ding X, Su S, Nandakumar K, Wang X, Fardo DW. A 2-step penalized regression  
535 method for family-based next-generation sequencing association studies. In: *BMC pro-  
536 ceedings*. vol. 8. BioMed Central; 2014. p. S25. 5
- 537 [20] Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear models  
538 via coordinate descent. *Journal of statistical software*. 2010;33(1):1. 5, 6, 26, 30, 43
- 539 [21] Yang Y, Zou H. A fast unified algorithm for solving group-lasso penalize learning  
540 problems. *Statistics and Computing*. 2015;25(6):1129–1141. 5

## REFERENCES

- 541 [22] Yang J, Zaitlen NA, Goddard ME, Visscher PM, Price AL. Advantages and pitfalls in  
542 the application of mixed-model association methods. *Nature genetics*. 2014;46(2):100.  
543 5
- 544 [23] Zou H, Hastie T. Regularization and variable selection via the elastic net. *Journal of*  
545 *the Royal Statistical Society: Series B (Statistical Methodology)*. 2005;67(2):301–320.  
546 5
- 547 [24] Yuan M, Lin Y. Model selection and estimation in regression with grouped vari-  
548 ables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*.  
549 2006;68(1):49–67. 5
- 550 [25] Gilmour AR, Thompson R, Cullis BR. Average information REML: an efficient algo-  
551 rithm for variance parameter estimation in linear mixed models. *Biometrics*. 1995;p.  
552 1440–1450. 6
- 553 [26] Dandine-Roulland C. *gaston: Genetic Data Handling (QC, GRM, LD, PCA) and*  
554 *Linear Mixed Models*; 2018. R package version 1.5.3. Available from: [https:](https://CRAN.R-project.org/package=gaston)  
555 [//CRAN.R-project.org/package=gaston](https://CRAN.R-project.org/package=gaston). 6
- 556 [27] Ochoa A, Storey JD. FST and kinship for arbitrary population structures I: Generalized  
557 definitions. *bioRxiv*. 2016;. 8
- 558 [28] Ochoa A, Storey JD. FST and kinship for arbitrary population structures II: Method  
559 of moments estimators. *bioRxiv*. 2016;. 8
- 560 [29] Reid S, Tibshirani R, Friedman J. A study of error variance estimation in lasso regres-  
561 sion. *Statistica Sinica*. 2016;p. 35–67. 12
- 562 [30] Davey Smith G, Ebrahim S. Mendelian randomization: can genetic epidemiology con-  
563 tribute to understanding environmental determinants of disease? *International journal*  
564 *of epidemiology*. 2003;32(1):1–22. 13

## REFERENCES

- 565 [31] Cherlin S, Howey RA, Cordell HJ. Using penalized regression to predict phenotype  
566 from SNP data. In: BMC proceedings. vol. 12. BioMed Central; 2018. p. 38. 13
- 567 [32] Zhou W, Lo SH. Analysis of genotype by methylation interactions through sparsity-  
568 inducing regularized regression. In: BMC proceedings. vol. 12. BioMed Central; 2018.  
569 p. 40. 13
- 570 [33] Howey RA, Cordell HJ. Application of Bayesian networks to GAW20 genetic and blood  
571 lipid data. In: BMC proceedings. vol. 12. BioMed Central; 2018. p. 19. 13
- 572 [34] Thornton T, Tang H, Hoffmann TJ, Ochs-Balcom HM, Caan BJ, Risch N. Esti-  
573 mating kinship in admixed populations. *The American Journal of Human Genetics*.  
574 2012;91(1):122–138. 13
- 575 [35] Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in  
576 unrelated individuals. *Genome research*. 2009;19(9):1655–1664. 13
- 577 [36] Fortin A, Diez E, Rochefort D, Laroche L, Malo D, Rouleau GA, et al. Recombinant  
578 congenic strains derived from A/J and C57BL/6J: a tool for genetic dissection of com-  
579 plex traits. *Genomics*. 2001;74(1):21–35. 15
- 580 [37] Bennett BJ, Farber CR, Orozco L, Kang HM, Ghazalpour A, Siemers N, et al. A  
581 high-resolution association mapping panel for the dissection of complex traits in mice.  
582 *Genome research*. 2010;20(2):281–290. 16
- 583 [38] Flint J, Eskin E. Genome-wide association studies in mice. *Nature Reviews Genetics*.  
584 2012;13(11):807. 16
- 585 [39] Cheng R, Lim JE, Samocha KE, Sokoloff G, Abney M, Skol AD, et al. Genome-wide  
586 association studies and the problem of relatedness among advanced intercross lines and  
587 other highly recombinant populations. *Genetics*. 2010;185(3):1033–1044. 16

## REFERENCES

- 588 [40] Di Pietrantonio T, Hernandez C, Girard M, Verville A, Orlova M, Belley A, et al.  
589 Strain-specific differences in the genetic control of two closely related mycobacteria.  
590 PLoS pathogens. 2010;6(10):e1001169. 16
- 591 [41] Sohrabi Y, Havelková H, Kobets T, Šíma M, Volkova V, Grekov I, et al. Mapping the  
592 Genes for Susceptibility and Response to *Leishmania tropica* in Mouse. PLoS neglected  
593 tropical diseases. 2013;7(7):e2282. 17
- 594 [42] Jackson AU, Fornés A, Galecki A, Miller RA, Burke DT. Multiple-trait quantitative  
595 trait loci analysis using a large mouse sibship. Genetics. 1999;151(2):785–795. 17
- 596 [43] Stern1 M, Benavides F, Klingelberger E, Conti2 C. Allelotype analysis of chemi-  
597 cally induced squamous cell carcinomas in F1 hybrids of two inbred mouse strains with  
598 different susceptibility to tumor progression. Carcinogenesis. 2000;21(7):1297–1301. 17
- 599 [44] Lasko D, Cavenee W, Nordenskjöld M. Loss of constitutional heterozygosity in human  
600 cancer. Annual review of genetics. 1991;25(1):281–314. 17
- 601 [45] Loh PR, Tucker G, Bulik-Sullivan BK, Vilhjalmsón BJ, Finucane HK, Salem RM, et al.  
602 Efficient Bayesian mixed-model analysis increases association power in large cohorts.  
603 Nature genetics. 2015;47(3):284. 19
- 604 [46] Allen N, Sudlow C, Downey P, Peakman T, Danesh J, Elliott P, et al. UK Biobank:  
605 Current status and what it means for epidemiology. Health Policy and Technology.  
606 2012;1(3):123–126. 20
- 607 [47] Spain SL, Barrett JC. Strategies for fine-mapping complex traits. Human molecular  
608 genetics. 2015;24(R1):R111–R119. 21
- 609 [48] Pirinen M, Donnelly P, Spencer CC, et al. Efficient computation with a linear mixed  
610 model on large-scale data sets with applications to genetic studies. The Annals of  
611 Applied Statistics. 2013;7(1):369–390. 22, 23

- 612 [49] Tseng P, Yun S. A coordinate gradient descent method for nonsmooth separable mini-  
613 mization. *Mathematical Programming*. 2009;117(1):387–423. 26, 43, 46
- 614 [50] Meier L, Van De Geer S, Bühlmann P. The group lasso for logistic regression. *Journal*  
615 *of the Royal Statistical Society: Series B (Statistical Methodology)*. 2008;70(1):53–71.  
616 26, 43
- 617 [51] Schelldorfer J, Bühlmann P, DE G, VAN S. Estimation for High-Dimensional Lin-  
618 ear Mixed-Effects Models Using L1-Penalization. *Scandinavian Journal of Statistics*.  
619 2011;38(2):197–214. 26, 33, 43
- 620 [52] Byrd RH, Lu P, Nocedal J, Zhu C. A limited memory algorithm for bound constrained  
621 optimization. *SIAM Journal on Scientific Computing*. 1995;16(5):1190–1208. 28
- 622 [53] Wakefield J. Bayesian and frequentist regression methods. Springer Science & Business  
623 Media; 2013. 31
- 624 [54] Nishii R. Asymptotic properties of criteria for selection of variables in multiple regres-  
625 sion. *The Annals of Statistics*. 1984;p. 758–765. 33
- 626 [55] Zou H, Hastie T, Tibshirani R, et al. On the degrees of freedom of the lasso. *The*  
627 *Annals of Statistics*. 2007;35(5):2173–2192. 33
- 628 [56] Bondell HD, Krishna A, Ghosh SK. Joint Variable Selection for Fixed and Random  
629 Effects in Linear Mixed-Effects Models. *Biometrics*. 2010;66(4):1069–1077. 33
- 630 [57] Fan Y, Tang CY. Tuning parameter selection in high dimensional penalized likeli-  
631 hood. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*.  
632 2013;75(3):531–552. 34
- 633 [58] Xie Y. *Dynamic Documents with R and knitr*. vol. 29. CRC Press; 2015. 64

## 634 A Block Coordinate Descent Algorithm

635 We use a general purpose block coordinate descent algorithm (CGD) [49] to solve (17). At  
 636 each iteration, the algorithm approximates the negative log-likelihood  $f(\cdot)$  in  $Q_\lambda(\cdot)$  by a  
 637 strictly convex quadratic function and then applies block coordinate descent to generate a  
 638 descent direction followed by an inexact line search along this direction [49]. For continuously  
 639 differentiable  $f(\cdot)$  and convex and block-separable  $P(\cdot)$  (i.e.  $P(\beta) = \sum_i P_i(\beta_i)$ ), [49] show  
 640 that the solution generated by the CGD method is a stationary point of  $Q_\lambda(\cdot)$  if the coor-  
 641 dinates are updated in a Gauss-Seidel manner i.e.  $Q_\lambda(\cdot)$  is minimized with respect to one  
 642 parameter while holding all others fixed. The CGD algorithm can thus be run in parallel and  
 643 therefore suited for large  $p$  settings. It has been successfully applied in fixed effects models  
 644 (e.g. [50], [20]) and [51] for mixed models with an  $\ell_1$  penalty. Following Tseng and Yun [49],  
 645 the CGD algorithm is given by Algorithm 2.

646 The Armijo rule is defined as follows [49]:

Choose  $\alpha_{init}^{(k)} > 0$  and let  $\alpha^{(k)}$  be the largest element of  $\{\alpha_{init}^{(k)} \delta^r\}_{r=0,1,2,\dots}$  satisfying

$$Q_\lambda(\Theta_j^{(k)} + \alpha^{(k)} d^{(k)}) \leq Q_\lambda(\Theta_j^{(k)}) + \alpha^{(k)} \varrho \Delta^{(k)} \quad (45)$$

where  $0 < \delta < 1$ ,  $0 < \varrho < 1$ ,  $0 \leq \gamma < 1$  and

$$\Delta^{(k)} := \nabla f(\Theta_j^{(k)}) d^{(k)} + \gamma (d^{(k)})^2 H_{jj}^{(k)} + \lambda P(\Theta_j^{(k)} + d^{(k)}) - \lambda P(\Theta_j^{(k)}) \quad (46)$$

648 Common choices for the constants are  $\delta = 0.1$ ,  $\varrho = 0.001$ ,  $\gamma = 0$ ,  $\alpha_{init}^{(k)} = 1$  for all  $k$  [51].

649 Below we detail the specifics of Algorithm 2 for the  $\ell_1$  penalty.

## A BLOCK COORDINATE DESCENT ALGORITHM

---

**Algorithm 2:** Coordinate Gradient Descent Algorithm to solve (17)

---

Set the iteration counter  $k \leftarrow 0$  and choose initial values for the parameter vector

$\Theta^{(0)}$ ;

**repeat**

Approximate the Hessian  $\nabla^2 f(\Theta^{(k)})$  by a symmetric matrix  $H^{(k)}$ :

$$H^{(k)} = \text{diag} \left[ \min \left\{ \max \left\{ \left[ \nabla^2 f(\Theta^{(k)}) \right]_{jj}, c_{min} \right\}, c_{max} \right\} \right]_{j=1, \dots, p} \quad (41)$$

**for**  $j = 1, \dots, p$  **do**

Solve the descent direction  $d^{(k)} := d_{H^{(k)}}(\Theta_j^{(k)})$ ;

**if**  $\Theta_j^{(k)} \in \{\beta_1, \dots, \beta_p\}$  **then**

$$d_{H^{(k)}}(\Theta_j^{(k)}) \leftarrow \arg \min_d \left\{ \nabla f(\Theta_j^{(k)})d + \frac{1}{2}d^2 H_{jj}^{(k)} + \lambda P(\Theta_j^{(k)} + d) \right\} \quad (42)$$

**end**

**end**

Choose a stepsize;

$$\alpha_j^{(k)} \leftarrow \text{line search given by the Armijo rule}$$

Update;

$$\widehat{\Theta}_j^{(k+1)} \leftarrow \widehat{\Theta}_j^{(k)} + \alpha_j^{(k)} d^{(k)}$$

Update;

$$\widehat{\eta}^{(k+1)} \leftarrow \arg \min_{\eta} \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^{2(k)}} \sum_{i=1}^{N_T} \frac{\left( \widetilde{Y}_i - \sum_{j=0}^p \widetilde{X}_{ij+1} \beta_j^{(k+1)} \right)^2}{1 + \eta(\Lambda_i - 1)} \quad (43)$$

Update;

$$\widehat{\sigma}^{2(k+1)} \leftarrow \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{\left( \widetilde{Y}_i - \sum_{j=0}^p \widetilde{X}_{ij+1} \beta_j^{(k+1)} \right)^2}{1 + \eta^{(k+1)}(\Lambda_i - 1)} \quad (44)$$

$k \leftarrow k + 1$

**until** convergence criterion is satisfied;

---

## 650 A.1 $\ell_1$ penalty

651 The objective function is given by

$$Q_\lambda(\Theta) = f(\Theta) + \lambda|\beta| \quad (47)$$

### 652 A.1.1 Descent Direction

653 For simplicity, we remove the iteration counter ( $k$ ) from the derivation below.

654 For  $\Theta_j^{(k)} \in \{\beta_1, \dots, \beta_p\}$ , let

$$d_H(\Theta_j) = \arg \min_d G(d) \quad (48)$$

655 where

$$G(d) = \nabla f(\Theta_j)d + \frac{1}{2}d^2 H_{jj} + \lambda|\Theta_j + d|$$

656 Since  $G(d)$  is not differentiable at  $-\Theta_j$ , we calculate the subdifferential  $\partial G(d)$  and search

657 for  $d$  with  $0 \in \partial G(d)$ :

$$\partial G(d) = \nabla f(\Theta_j) + dH_{jj} + \lambda u \quad (49)$$

658 where

$$u = \begin{cases} 1 & \text{if } d > -\Theta_j \\ -1 & \text{if } d < -\Theta_j \\ [-1, 1] & \text{if } d = -\Theta_j \end{cases} \quad (50)$$

659 We consider each of the three cases in (49) below

1.  $d > -\Theta_j$

$$\begin{aligned} \partial G(d) &= \nabla f(\Theta_j) + dH_{jj} + \lambda = 0 \\ d &= \frac{-(\nabla f(\Theta_j) + \lambda)}{H_{jj}} \end{aligned}$$

## A BLOCK COORDINATE DESCENT ALGORITHM

---

Since  $\lambda > 0$  and  $H_{jj} > 0$ , we have

$$\frac{-(\nabla f(\Theta_j) - \lambda)}{H_{jj}} > \frac{-(\nabla f(\Theta_j) + \lambda)}{H_{jj}} = d \stackrel{\text{def}}{>} -\Theta_j$$

The solution can be written compactly as

$$d = \text{mid} \left\{ \frac{-(\nabla f(\Theta_j) - \lambda)}{H_{jj}}, -\Theta_j, \frac{-(\nabla f(\Theta_j) + \lambda)}{H_{jj}} \right\}$$

660 where  $\text{mid} \{a, b, c\}$  denotes the median (mid-point) of  $a, b, c$  [49].

2.  $d < -\Theta_j$

$$\begin{aligned} \partial G(d) &= \nabla f(\Theta_j) + dH_{jj} - \lambda = 0 \\ d &= \frac{-(\nabla f(\Theta_j) - \lambda)}{H_{jj}} \end{aligned}$$

Since  $\lambda > 0$  and  $H_{jj} > 0$ , we have

$$\frac{-(\nabla f(\Theta_j) + \lambda)}{H_{jj}} < \frac{-(\nabla f(\Theta_j) - \lambda)}{H_{jj}} = d \stackrel{\text{def}}{<} -\Theta_j$$

Again, the solution can be written compactly as

$$d = \text{mid} \left\{ \frac{-(\nabla f(\Theta_j) - \lambda)}{H_{jj}}, -\Theta_j, \frac{-(\nabla f(\Theta_j) + \lambda)}{H_{jj}} \right\}$$

3.  $d_j = -\Theta_j$

There exists  $u \in [-1, 1]$  such that

$$\begin{aligned} \partial G(d) &= \nabla f(\Theta_j) + dH_{jj} + \lambda u = 0 \\ d &= \frac{-(\nabla f(\Theta_j) + \lambda u)}{H_{jj}} \end{aligned}$$

## A BLOCK COORDINATE DESCENT ALGORITHM

For  $-1 \leq u \leq 1$ ,  $\lambda > 0$  and  $H_{jj} > 0$  we have

$$\frac{-(\nabla f(\Theta_j) + \lambda)}{H_{jj}} \leq d \stackrel{\text{def}}{=} -\Theta_j \leq \frac{-(\nabla f(\Theta_j) - \lambda)}{H_{jj}}$$

The solution can again be written compactly as

$$d = \text{mid} \left\{ \frac{-(\nabla f(\Theta_j) - \lambda)}{H_{jj}}, -\Theta_j, \frac{-(\nabla f(\Theta_j) + \lambda)}{H_{jj}} \right\}$$

661 We see all three cases lead to the same solution for (48). Therefore the descent direction for  
 662  $\Theta_j^{(k)} \in \{\beta_1, \dots, \beta_p\}$  for the  $\ell_1$  penalty is given by

$$d = \text{mid} \left\{ \frac{-(\nabla f(\beta_j) - \lambda)}{H_{jj}}, -\beta_j, \frac{-(\nabla f(\beta_j) + \lambda)}{H_{jj}} \right\} \quad (51)$$

### 663 A.1.2 Solution for the $\beta$ parameter

664 If the Hessian  $\nabla^2 f(\Theta^{(k)}) > 0$  then  $H^{(k)}$  defined in (41) is equal to  $\nabla^2 f(\Theta^{(k)})$ . Using  $\alpha_{init} = 1$ ,  
 665 the largest element of  $\left\{ \alpha_{init}^{(k)} \delta^r \right\}_{r=0,1,2,\dots}$  satisfying the Armijo Rule inequality is reached for  
 666  $\alpha^{(k)} = \alpha_{init}^{(k)} \delta^0 = 1$ . The Armijo rule update for the  $\beta$  parameter is then given by

$$\beta_j^{(k+1)} \leftarrow \beta_j^{(k)} + d^{(k)}, \quad j = 1, \dots, p \quad (52)$$

667 Substituting the descent direction given by (51) into (52) we get

$$\beta_j^{(k+1)} = \text{mid} \left\{ \beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) - \lambda)}{H_{jj}}, 0, \beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) + \lambda)}{H_{jj}} \right\} \quad (53)$$

668 We can further simplify this expression. Let

$$w_i := \frac{1}{\sigma^2 (1 + \eta(\Lambda_i - 1))} \quad (54)$$

## A BLOCK COORDINATE DESCENT ALGORITHM

669 .

Re-write the part depending on  $\beta$  of the negative log-likelihood in (15) as

$$g(\beta^{(k)}) = \frac{1}{2} \sum_{i=1}^{N_T} w_i \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} - \tilde{X}_{ij} \beta_j^{(k)} \right)^2 \quad (55)$$

The gradient and Hessian are given by

$$\nabla f(\beta_j^{(k)}) := \frac{\partial}{\partial \beta_j^{(k)}} g(\beta^{(k)}) = - \sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} - \tilde{X}_{ij} \beta_j^{(k)} \right) \quad (56)$$

$$H_{jj} := \frac{\partial^2}{\partial \beta_j^{(k)2}} g(\beta^{(k)}) = \sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2 \quad (57)$$

Substituting (56) and (57) into  $\beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) - \lambda)}{H_{jj}}$

$$\begin{aligned} & \beta_j^{(k)} + \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} - \tilde{X}_{ij} \beta_j^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \\ &= \beta_j^{(k)} + \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} - \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2 \beta_j^{(k)}}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \\ &= \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \end{aligned} \quad (58)$$

Similarly, substituting (56) and (57) in  $\beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) + \lambda)}{H_{jj}}$  we get

$$\frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) - \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \quad (59)$$

## A BLOCK COORDINATE DESCENT ALGORITHM

Finally, substituting (58) and (59) into (53) we get

$$\begin{aligned} \beta_j^{(k+1)} &= \text{mid} \left\{ \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) - \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2}, 0, \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \right\} \\ &= \frac{\mathcal{S}_\lambda \left( \sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) \right)}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \end{aligned} \quad (60)$$

Where  $\mathcal{S}_\lambda(x)$  is the soft-thresholding operator

$$\mathcal{S}_\lambda(x) = \text{sign}(x)(|x| - \lambda)_+$$

$\text{sign}(x)$  is the signum function

$$\text{sign}(x) = \begin{cases} -1 & x < 0 \\ 0 & x = 0 \\ 1 & x > 0 \end{cases}$$

670 and  $(x)_+ = \max(x, 0)$ .

## B ADDITIONAL SIMULATION RESULTS

### 671 B Additional Simulation Results

#### 672 B.1 Null Model ( $c = 0$ )

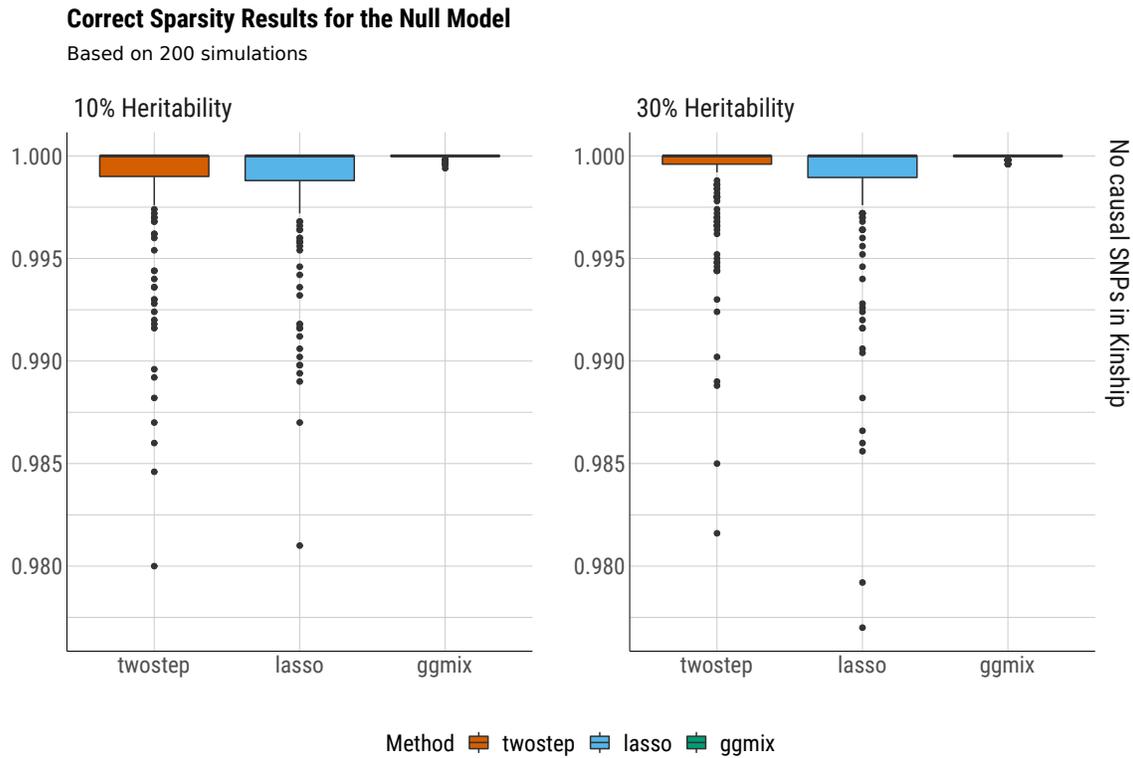


Figure B.1: Boxplots of the correct sparsity from 200 replications by the true heritability  $\eta = \{10\%, 30\%\}$ .

## B ADDITIONAL SIMULATION RESULTS

### Estimation Error Results for the Null Model

Based on 200 simulations

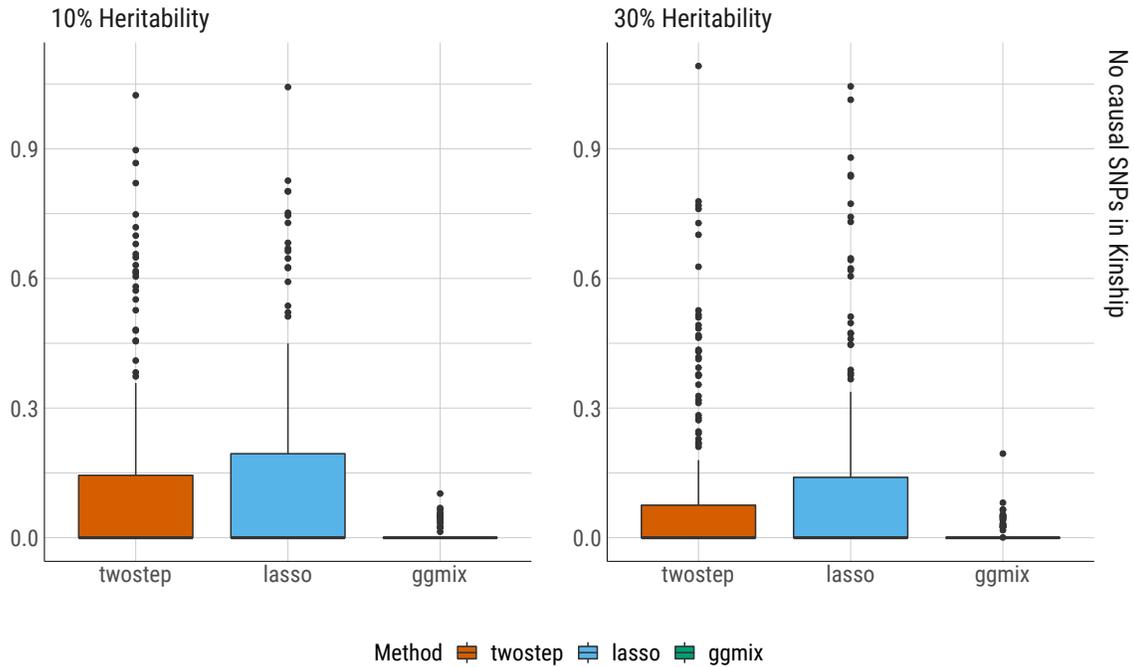
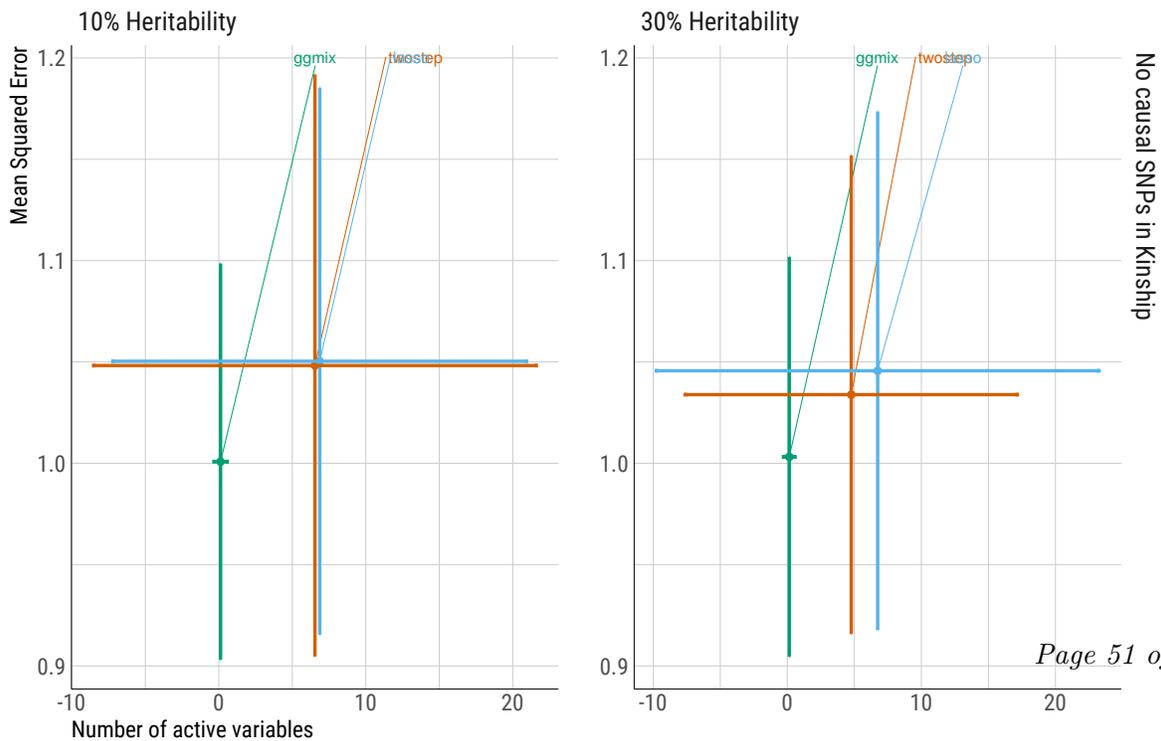


Figure B.2: Boxplots of the estimation error from 200 replications by the true heritability  $\eta = \{10\%, 30\%\}$ .

### Mean Squared Error vs. Number of Active Variable (Mean +/- 1 SD) for Null Model

Based on 200 simulations



## B ADDITIONAL SIMULATION RESULTS

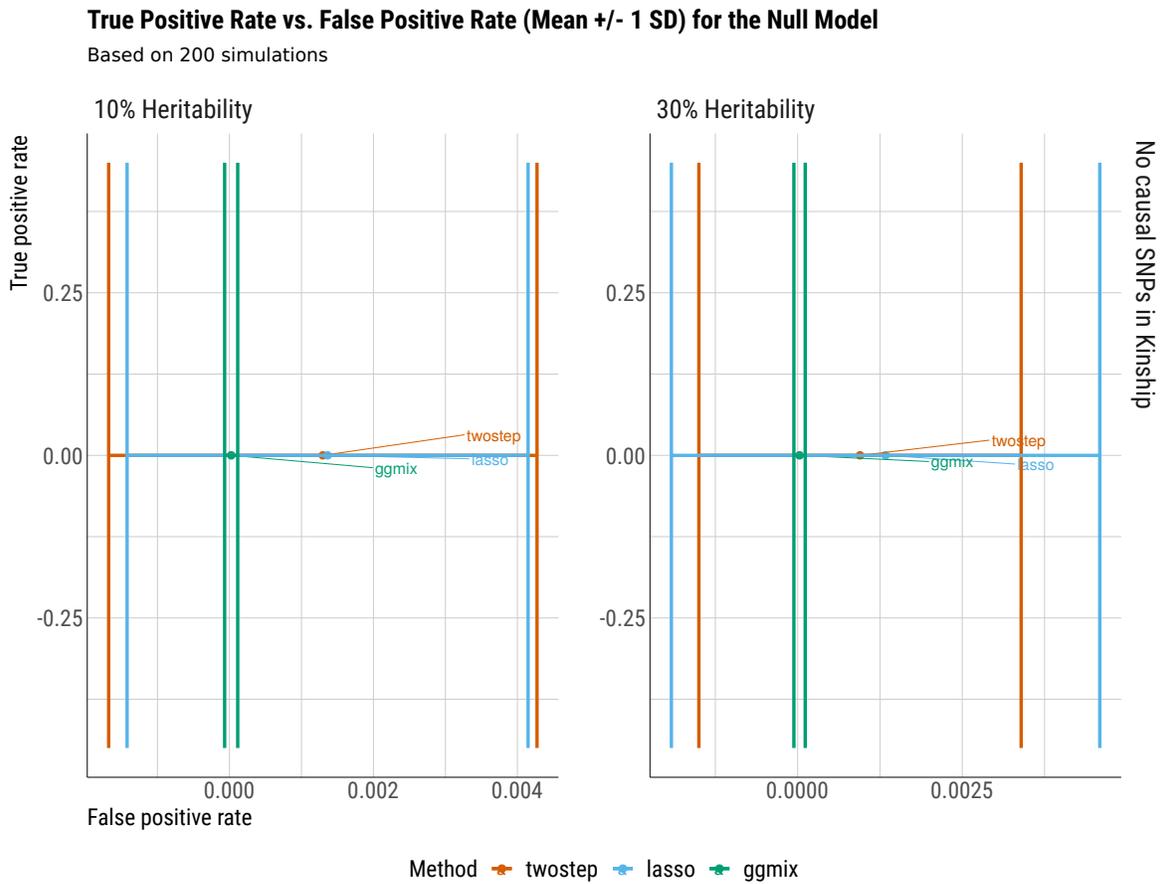


Figure B.4: Means  $\pm$  1 standard deviation of true positive rate vs. false positive rate from 200 replications by the true heritability  $\eta = \{10\%, 30\%\}$ .

## B ADDITIONAL SIMULATION RESULTS

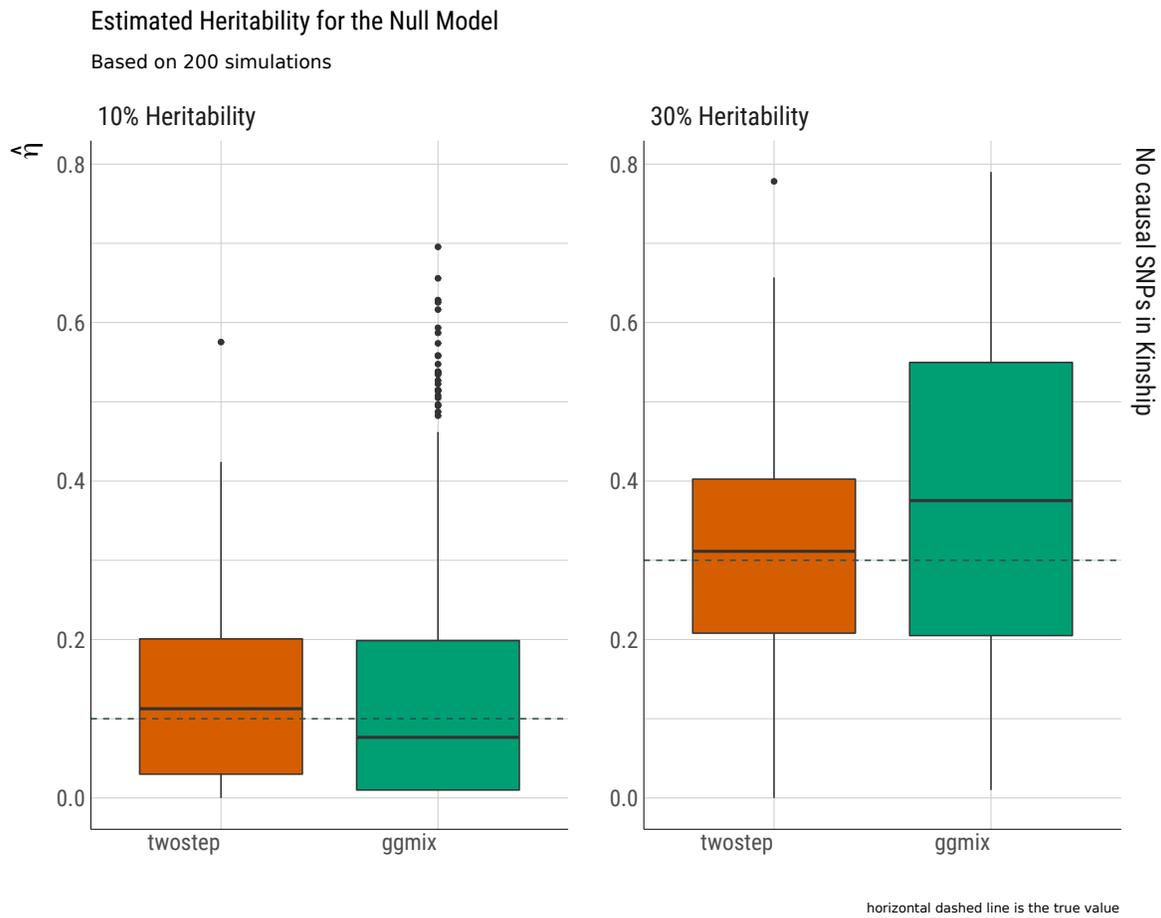


Figure B.5: Boxplots of the heritability estimate  $\hat{\eta}$  from 200 simulations by the true heritability  $\eta = \{10\%, 30\%\}$ .

## B ADDITIONAL SIMULATION RESULTS

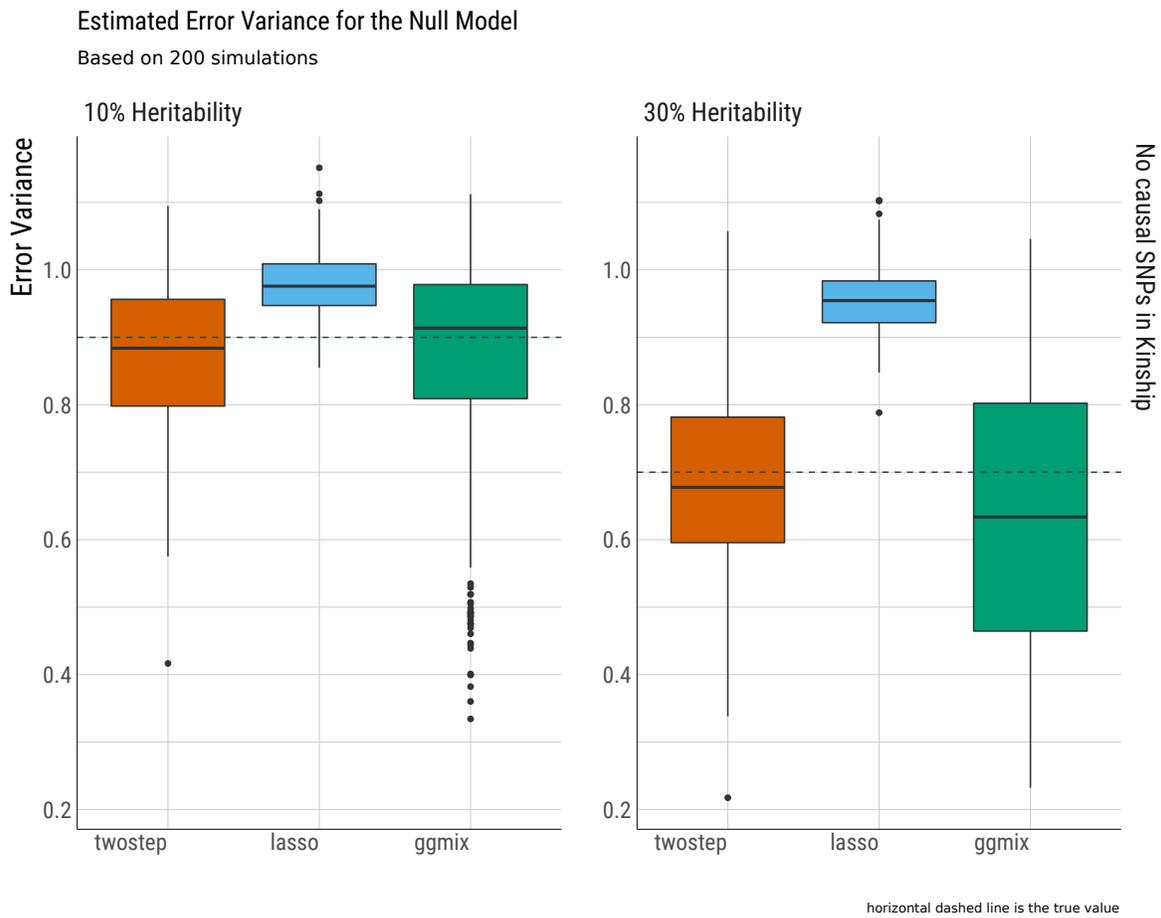


Figure B.6: Boxplots of the estimated error variance from 200 replications by the true heritability  $\eta = \{10\%, 30\%\}$ .

## B ADDITIONAL SIMULATION RESULTS

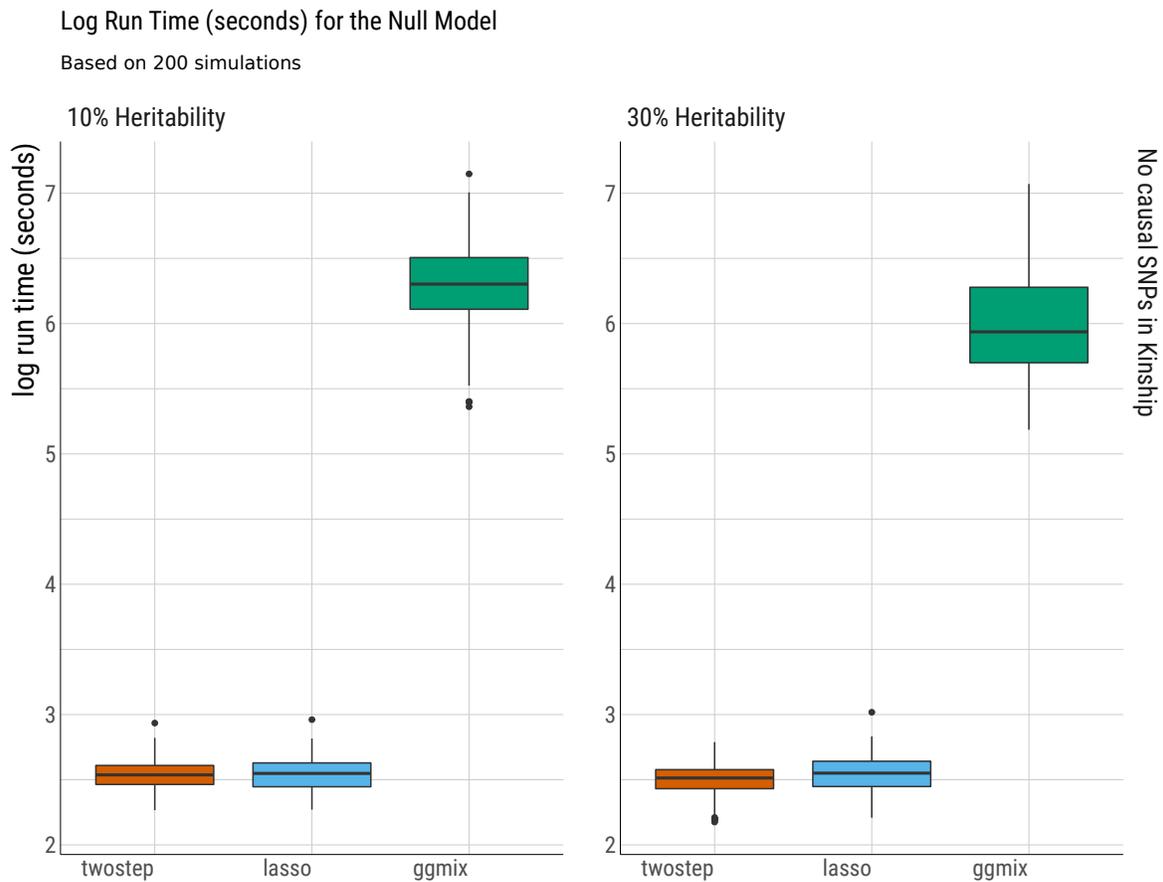


Figure B.7: Run time (in log seconds) for null model for `twostep`, `lasso` and `ggmix`.

### 673 B.2 1% of SNPs are Causal ( $c = 0.01$ )

## B ADDITIONAL SIMULATION RESULTS

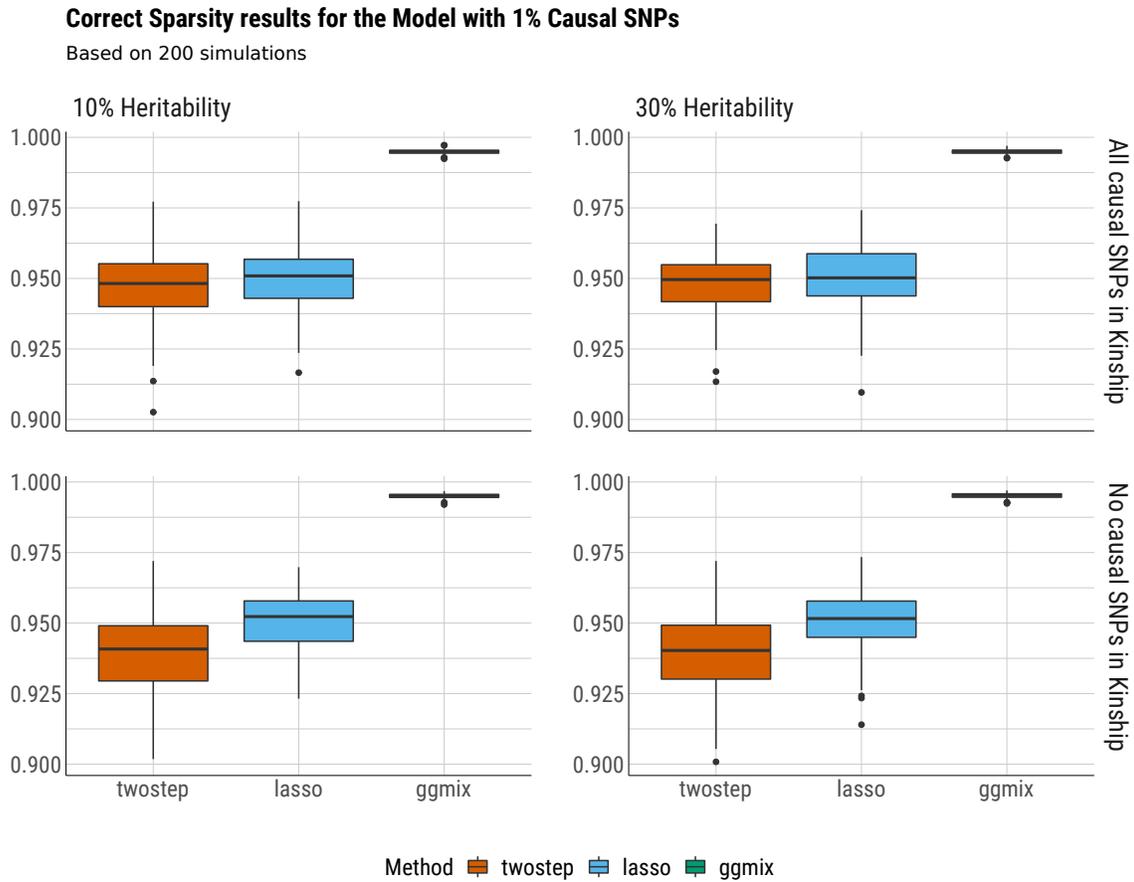


Figure B.8: Boxplots of the correct sparsity from 200 replications by the true heritability  $\eta = \{10\%, 30\%\}$  and number of causal SNPs that were included in the calculation of the kinship matrix for the model with 1% causal SNPs ( $c = 0.01$ ).

## B ADDITIONAL SIMULATION RESULTS

### Estimation Error results for the Model with 1% Causal SNPs

Based on 200 simulations

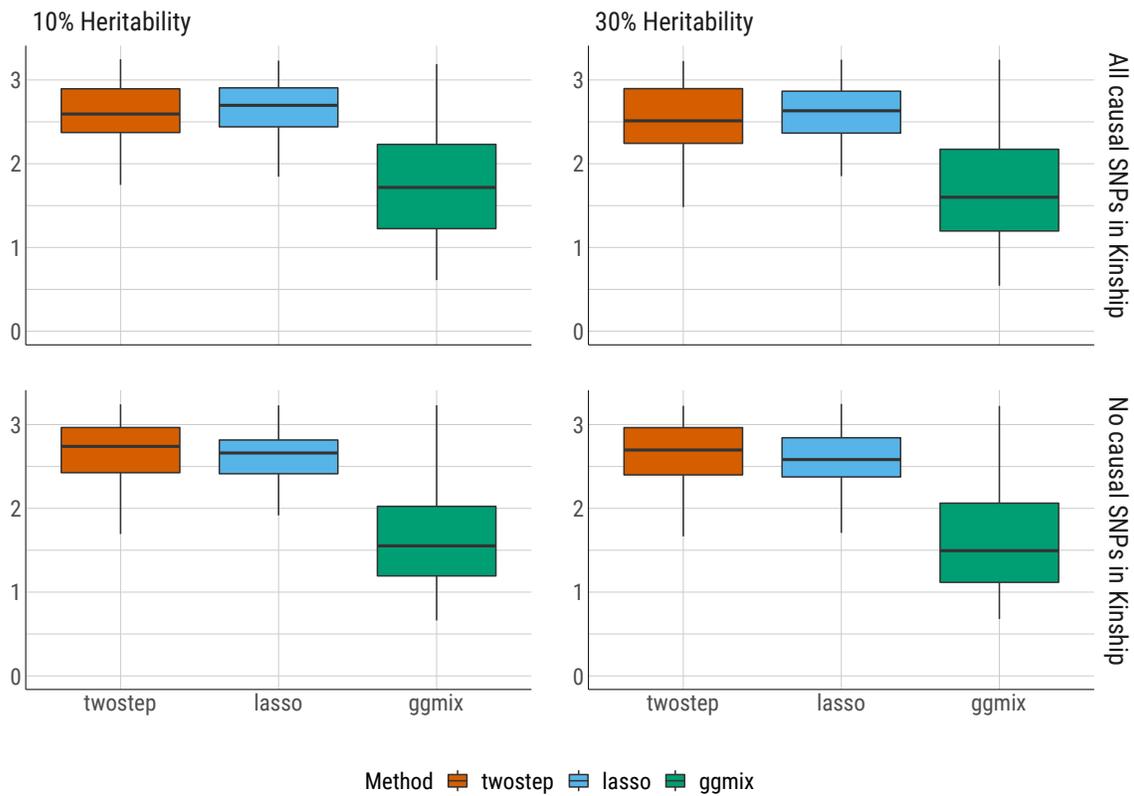


Figure B.9: Boxplots of the estimation error from 200 replications by the true heritability  $\eta = \{10\%, 30\%\}$  and number of causal SNPs that were included in the calculation of the kinship matrix for the model with 1% causal SNPs ( $c = 0.01$ ).

## B ADDITIONAL SIMULATION RESULTS

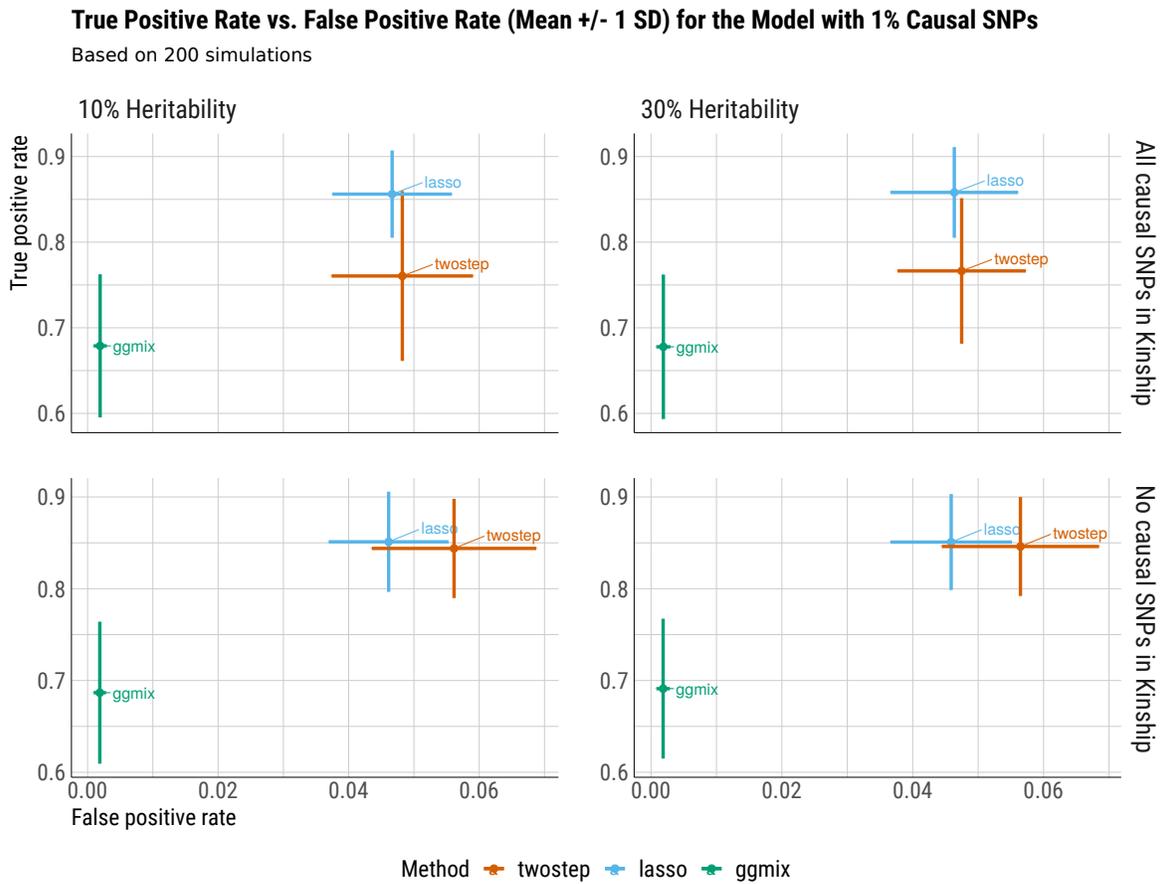


Figure B.10: Means  $\pm 1$  standard deviation of true positive rate vs. false positive rate from 200 replications by the true heritability  $\eta = \{10\%, 30\%\}$  and number of causal SNPs that were included in the calculation of the kinship matrix for the model with 1% causal SNPs ( $c = 0.01$ ).

## B ADDITIONAL SIMULATION RESULTS

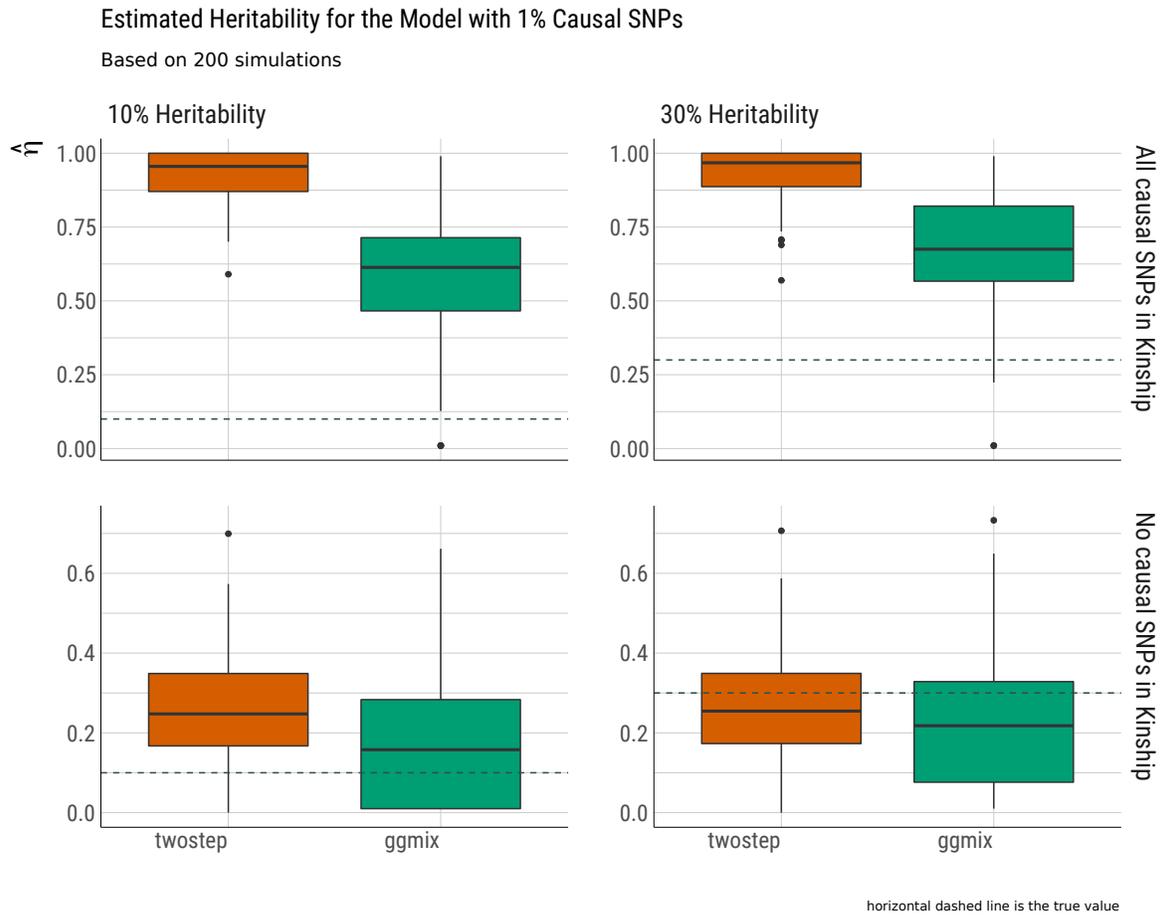


Figure B.11: Boxplots of the heritability estimate  $\hat{\eta}$  from 200 replications by the true heritability  $\eta = \{10\%, 30\%\}$  and number of causal SNPs that were included in the calculation of the kinship matrix for the model with 1% causal SNPs ( $c = 0.01$ ).

## B ADDITIONAL SIMULATION RESULTS

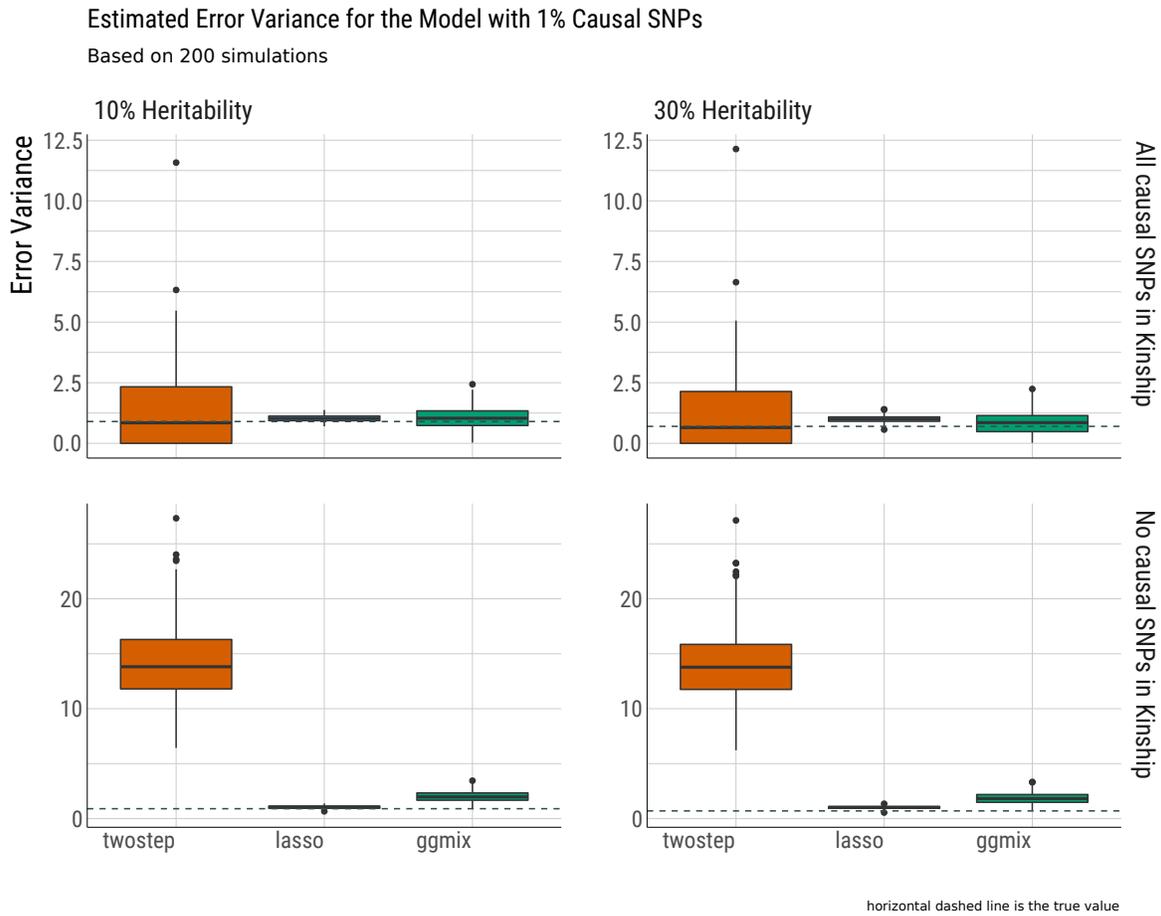


Figure B.12: Boxplots of the estimated error variance from 200 replications by the true heritability  $\eta = \{10\%, 30\%\}$  and number of causal SNPs that were included in the calculation of the kinship matrix for the model with 1% causal SNPs ( $c = 0.01$ ).

## B ADDITIONAL SIMULATION RESULTS

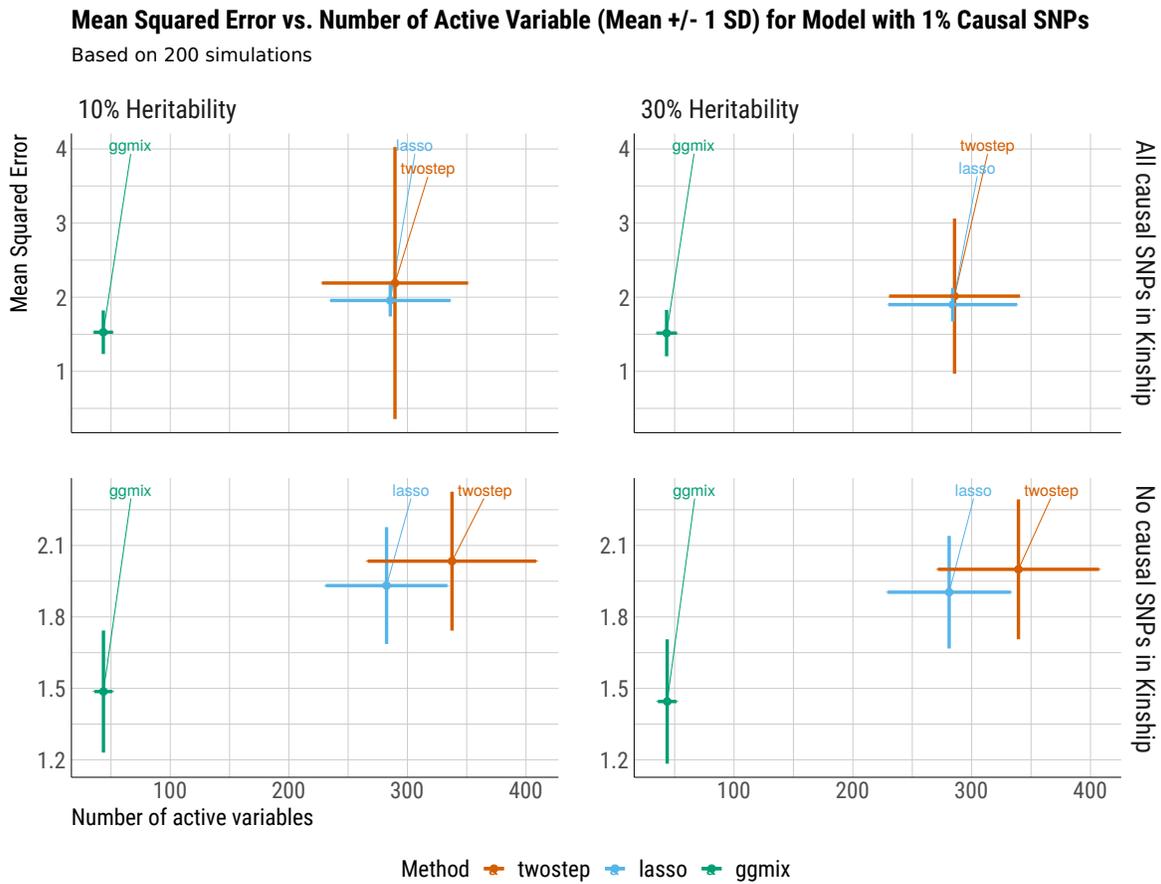


Figure B.13: Root mean squared prediction error on the test set vs. the number of active variables from 200 replications by the true heritability  $\eta = \{10\%, 30\%\}$  and number of causal SNPs that were included in the calculation of the kinship matrix for the model with 1% causal SNPs ( $c = 0.01$ ).

## B ADDITIONAL SIMULATION RESULTS

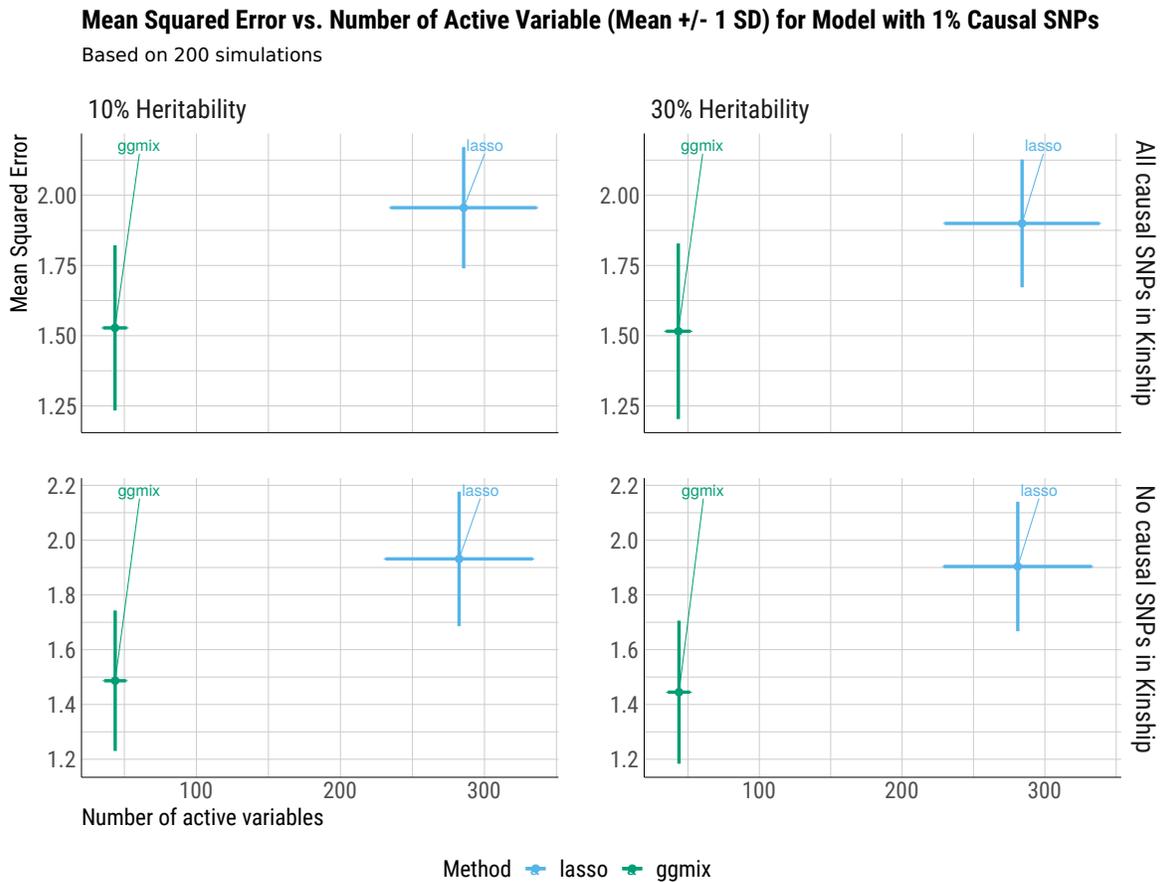


Figure B.14: Mean squared error vs number of active variables results from 200 replications by the true heritability  $\eta = \{10\%, 30\%\}$  and number of causal SNPs that were included in the calculation of the kinship matrix for the model with 1% causal SNPs ( $c = 0.01$ ), for 1% causal SNPs for ggmix and lasso only.

## B ADDITIONAL SIMULATION RESULTS

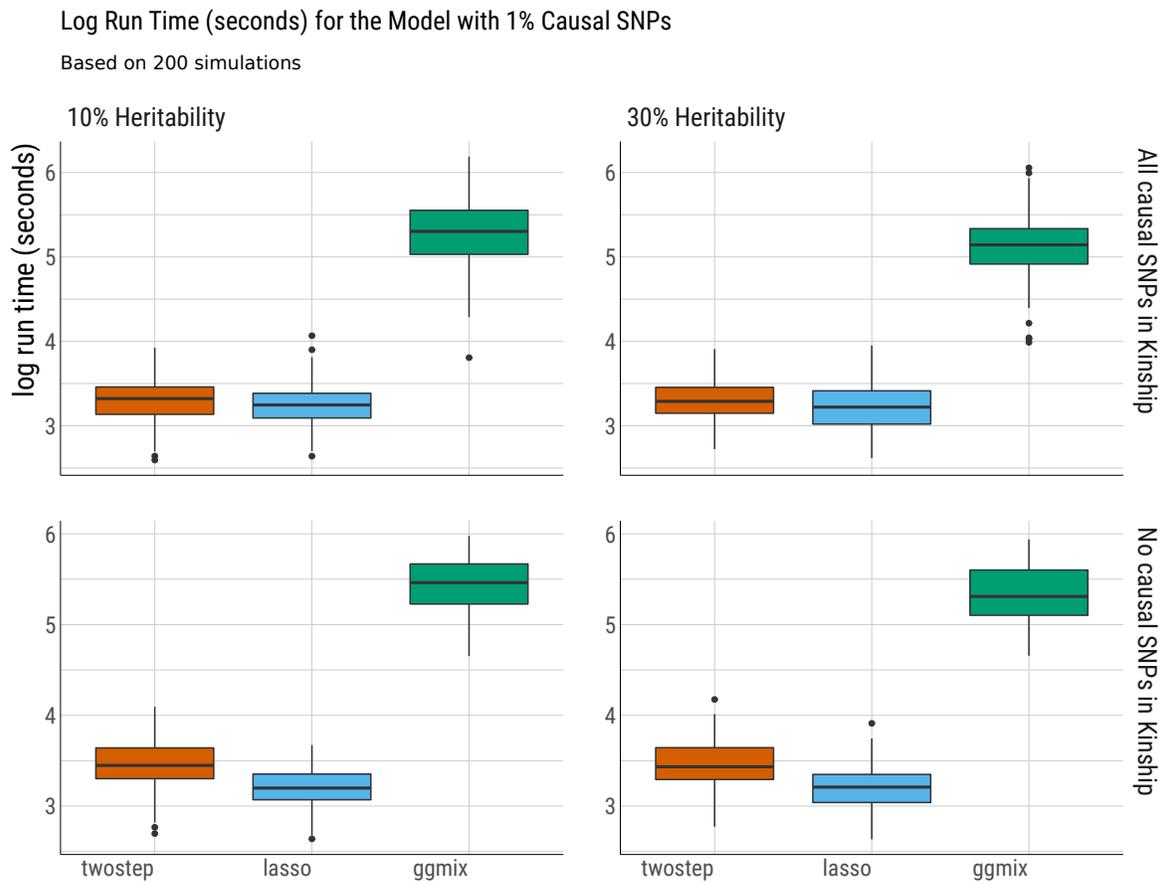


Figure B.15: Run time (in log seconds) from 200 replications by the true heritability  $\eta = \{10\%, 30\%\}$  and number of causal SNPs that were included in the calculation of the kinship matrix for the model with 1% causal SNPs ( $c = 0.01$ ).

## 674 C ggmix Package Showcase

675 In this section we briefly introduce the freely available and open source `ggmix` package in R.  
676 More comprehensive documentation is available at <https://sahirbhatnagar.com/ggmix>.  
677 Note that this entire section is reproducible; the code and text are combined in an `.Rnw`<sup>1</sup> file  
678 and compiled using `knitr` [58].

### 679 C.1 Installation

680 The package can be installed from [GitHub](#) via

```
install.packages("pacman")  
pacman::p_load_gh('sahirbhatnagar/ggmix')
```

681 To showcase the main functions in `ggmix`, we will use the simulated data which ships with  
682 the package and can be loaded via:

```
library(ggmix)  
data("admixed")  
names(admixed)  
  
## [1] "y"           "x"           "causal"  
## [4] "beta"       "kin"         "Xkinship"  
## [7] "not_causal" "causal_positive" "causal_negative"  
## [10] "x_lasso"
```

683 For details on how this data was simulated, see `help(admixed)`.

684 There are three basic inputs that `ggmix` needs:

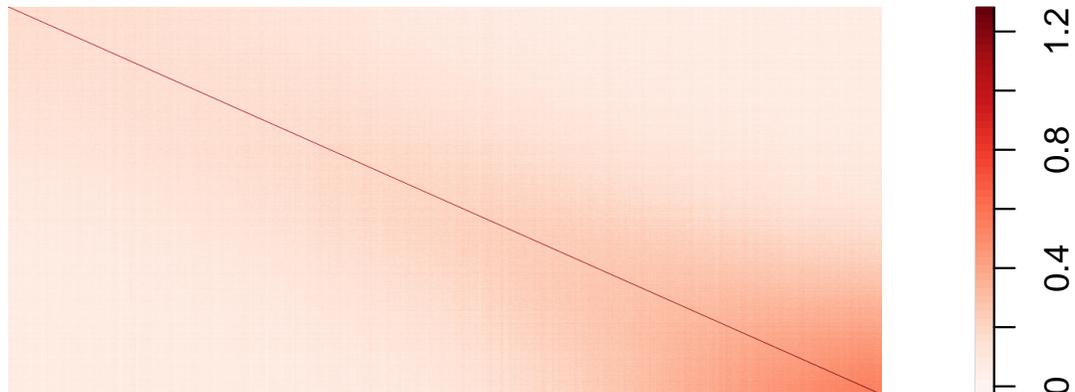
- 685 1.  $Y$ : a continuous response variable
- 686 2.  $X$ : a matrix of covariates of dimension  $N \times p$  where  $N$  is the sample size and  $p$  is the  
687 number of covariates
- 688 3.  $\Phi$ : a kinship matrix

---

<sup>1</sup>scripts available at <https://github.com/sahirbhatnagar/ggmix/tree/pgen/manuscript>

689 We can visualize the kinship matrix in the admixed data using the popkin package:

```
# need to install the package if you don't have it
# pacman::p_load_gh('StoreyLab/popkin')
popkin::plotPopkin(admixed$kin)
```



690

## 691 C.2 Fit the linear mixed model with Lasso Penalty

692 We will use the most basic call to the main function of this package, which is called `ggmix`.

693 This function will by default fit a  $L_1$  penalized linear mixed model (LMM) for 100 distinct

694 values of the tuning parameter  $\lambda$ . It will choose its own sequence:

```
fit <- ggmix(x = admixed$x, y = admixed$y, kinship = admixed$kin)
```

## C GGMIX PACKAGE SHOWCASE

```
names(fit)

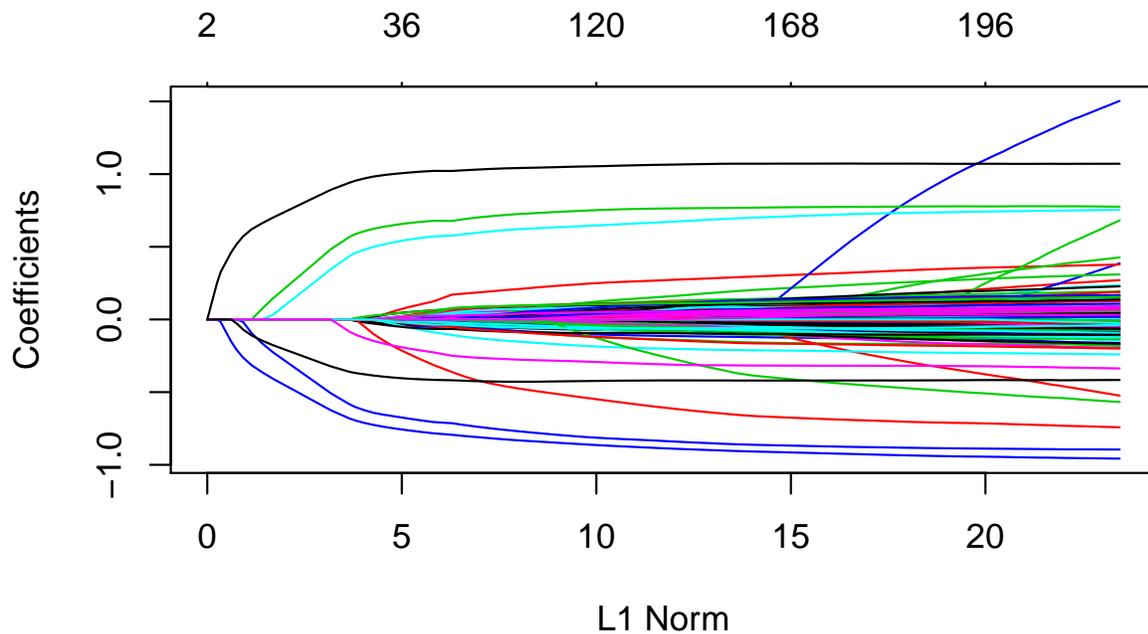
## [1] "result"      "ggmix_object" "n_design"     "p_design"
## [5] "lambda"     "coef"         "b0"           "beta"
## [9] "df"         "eta"          "sigma2"       "nlambda"
## [13] "cov_names"   "call"

class(fit)

## [1] "lassofullrank" "ggmix_fit"
```

695 We can see the solution path for each variable by calling the `plot` method for objects of  
696 class `ggmix_fit`:

```
plot(fit)
```



697

698 We can also get the coefficients for given value(s) of lambda using the `coef` method for  
699 objects of class `ggmix_fit`:

```
# only the first 5 coefficients printed here for brevity
```

```
coef(fit, s = c(0.1,0.02))[1:5, ]  
  
## 5 x 2 Matrix of class "dgeMatrix"  
##           1           2  
## (Intercept) -0.3824525 -0.030224599  
## X62          0.0000000  0.000000000  
## X185          0.0000000  0.001444518  
## X371          0.0000000  0.009513475  
## X420          0.0000000  0.000000000
```

700 Here,  $\mathbf{s}$  specifies the value(s) of  $\lambda$  at which the extraction is made. The function uses linear  
701 interpolation to make predictions for values of  $\mathbf{s}$  that do not coincide with the lambda  
702 sequence used in the fitting algorithm.

703 We can also get predictions ( $X\hat{\beta}$ ) using the `predict` method for objects of class `ggmix_fit`:

```
# need to provide x to the predict function  
# predict for the first 5 subjects  
predict(fit, s = c(0.1,0.02), newx = admixed$x[1:5,])  
  
##           1           2  
## id1 -1.19165061 -1.3123392  
## id2 -0.02913052  0.3885923  
## id3 -2.00084875 -2.6460043  
## id4 -0.37255277 -0.9542463  
## id5 -1.03967831 -2.1377268
```

### 704 C.3 Find the Optimal Value of the Tuning Parameter

705 We use the Generalized Information Criterion (GIC) to select the optimal value for  $\lambda$ . The  
706 default is  $a_n = \log(\log(n)) * \log(p)$  which corresponds to a high-dimensional BIC (HD-  
707 BIC):

```
# pass the fitted object from ggmix to the gic function:
```

## C GGMIX PACKAGE SHOWCASE

```
hdbic <- gic(fit)
class(hdbic)

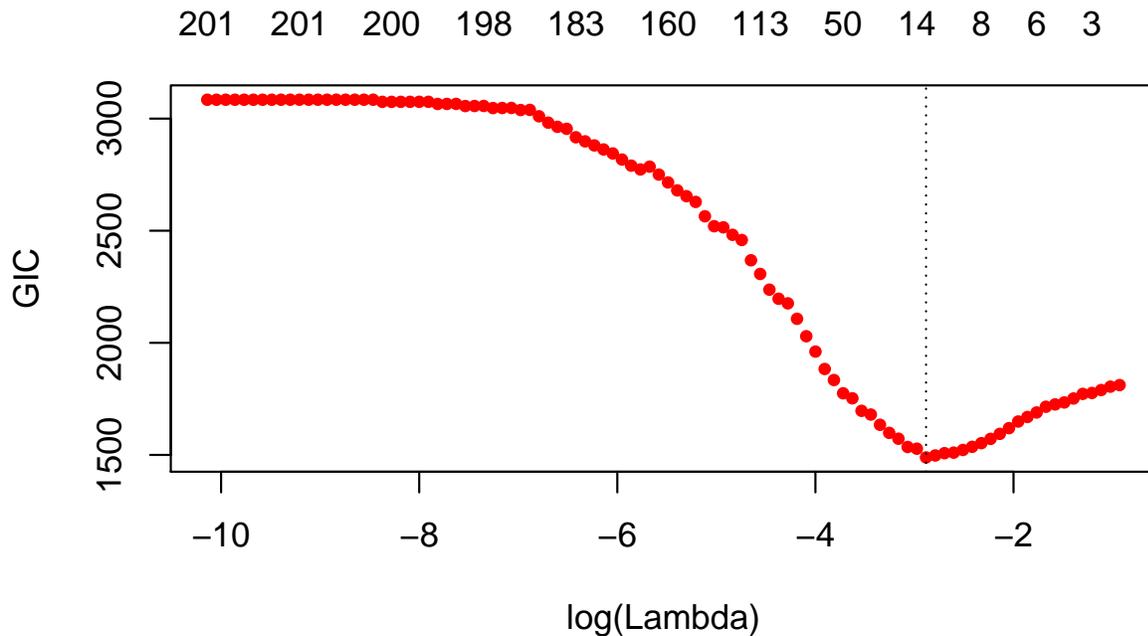
## [1] "ggmix_gic"      "lassofullrank" "ggmix_fit"

# we can also fit the BIC by specifying the an argument
bicfit <- gic(fit, an = log(length(admixed$y)))
```

708 We can plot the HDBIC values against  $\log(\lambda)$  using the `plot` method for objects of class

709 `ggmix_gic`:

```
plot(hdbic)
```



710

711 The optimal value for  $\lambda$  according to the HDBIC, i.e., the  $\lambda$  that leads to the minimum HDBIC

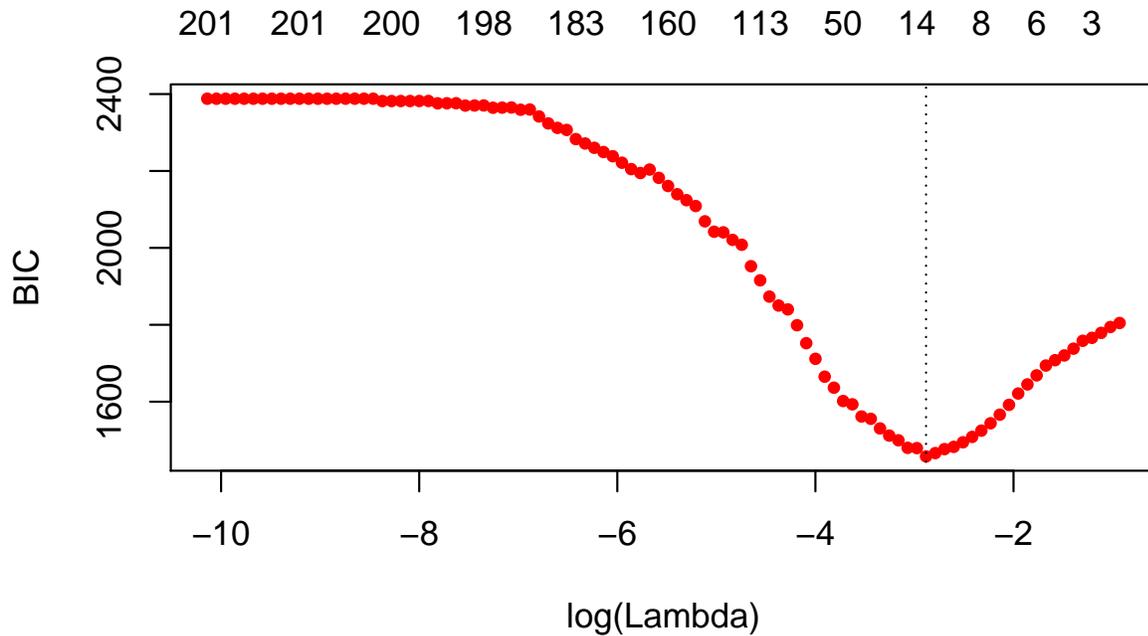
712 is:

```
hdbic[["lambda.min"]]

## [1] 0.05596623
```

713 We can also plot the BIC results:

```
plot(bicfit, ylab = "BIC")
```



714

```
bicfit[["lambda.min"]]  
## [1] 0.05596623
```

#### 715 C.4 Get Coefficients Corresponding to Optimal Model

716 We can use the object outputted by the `gic` function to extract the coefficients corresponding  
717 to the selected model using the `coef` method for objects of class `ggmix_gic`:

```
coef(hdbic)[1:5, , drop = FALSE]  
  
## 5 x 1 sparse Matrix of class "dgCMatrix"  
##           1  
## (Intercept) -0.2668419  
## X62         .  
## X185         .  
## X371         .  
## X420         .
```

718 We can also extract just the nonzero coefficients which also provide the estimated variance

719 components  $\eta$  and  $\sigma^2$ :

```
coef(hdbic, type = "nonzero")  
  
##              1  
## (Intercept) -0.26684191  
## X336        -0.67986393  
## X7638        0.43403365  
## X1536        0.93994982  
## X1943        0.56600730  
## X2849       -0.58157979  
## X56         -0.08244685  
## X4106       -0.35939830  
## eta         0.26746240  
## sigma2     0.98694300
```

720 We can also make predictions from the `hdbic` object, which by default will use the model  
721 corresponding to the optimal tuning parameter:

```
predict(hdbic, newx = admixed$x[1:5,])  
  
##              1  
## id1 -1.3061041  
## id2  0.2991654  
## id3 -2.3453664  
## id4 -0.4486012  
## id5 -1.3895793
```

## 722 C.5 Extracting Random Effects

723 The user can compute the random effects using the provided `ranef` method for objects of  
724 class `ggmix_gic`. This command will compute the estimated random effects for each subject  
725 using the parameters of the selected model:

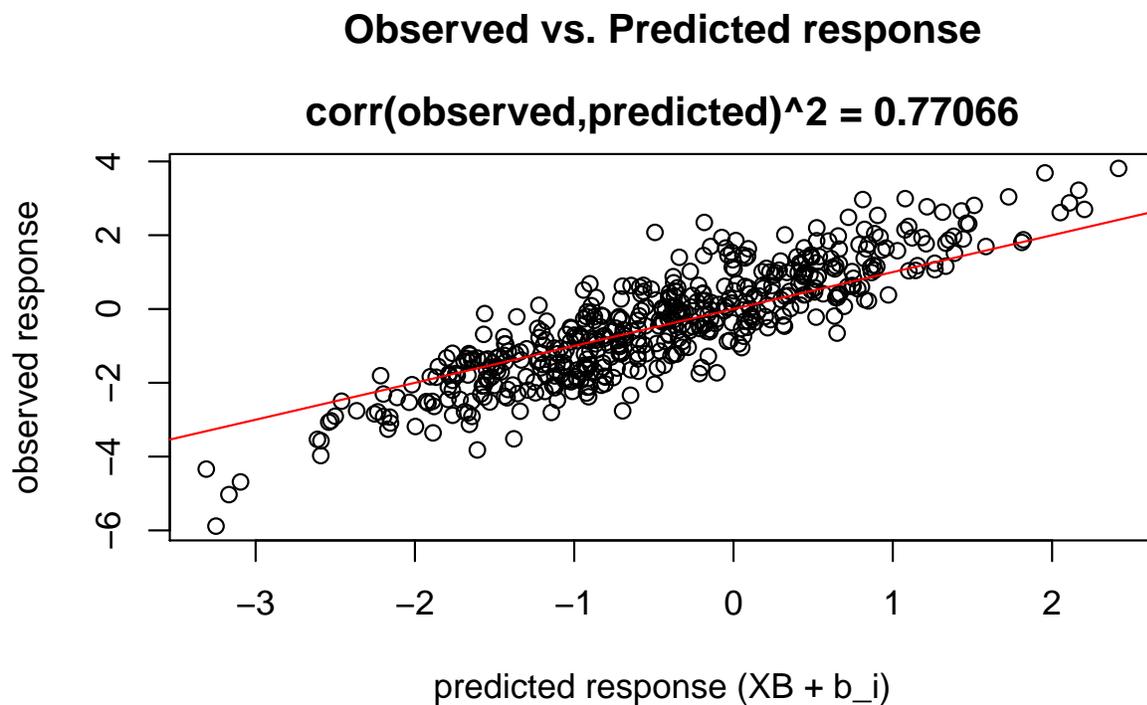
```
ranef(hdbic)[1:5]  
  
## [1] -0.02548691 -0.10011680  0.13020240 -0.30650997  0.16045768
```

## 726 C.6 Diagnostic Plots

727 We can also plot some standard diagnostic plots such as the observed vs. predicted response,  
728 QQ-plots of the residuals and random effects and the Tukey-Anscombe plot. These can be  
729 plotted using the `plot` method on a `ggmix_gic` object as shown below.

### 730 C.6.1 Observed vs. Predicted Response

```
plot(hdbic, type = "predicted", newx = admixed$x, newy = admixed$y)
```

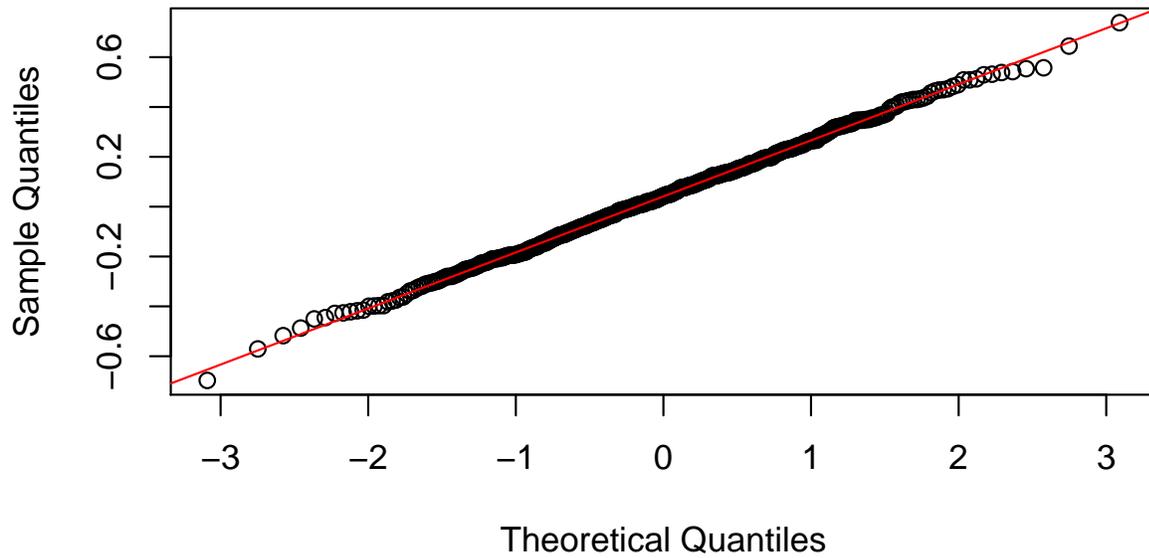


731

### 732 C.6.2 QQ-plots for Residuals and Random Effects

```
plot(hdbic, type = "QQranef", newx = admixed$x, newy = admixed$y)
```

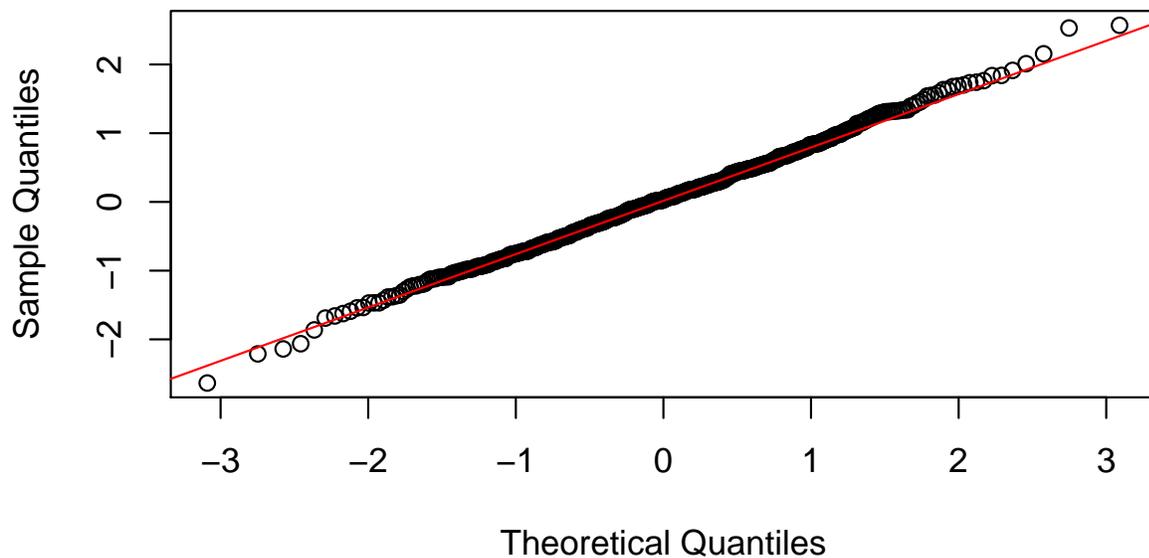
### QQ-Plot of the random effects at lambda = 0.06



733

```
plot(hdbic, type = "QQresid", newx = admixed$x, newy = admixed$y)
```

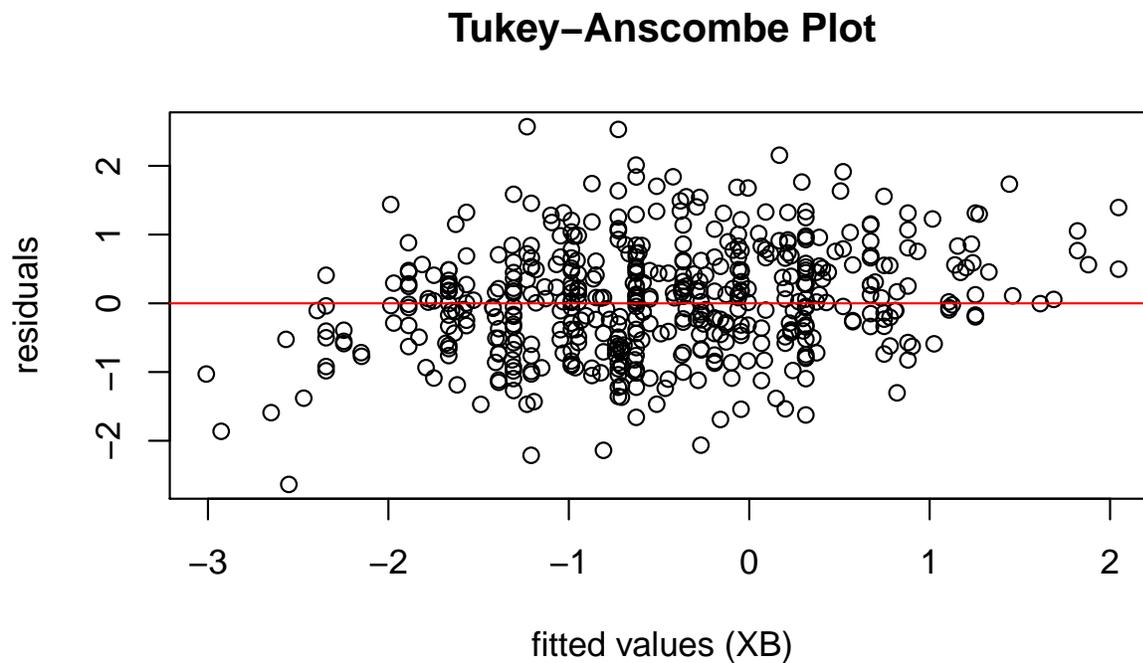
### QQ-Plot of the residuals at lambda = 0.06



734

735 C.6.3 Tukey-Anscombe Plot

```
plot(hdbic, type = "Tukey", newx = admixed$x, newy = admixed$y)
```



736