

# A General Framework for Variable Selection in Linear Mixed Models with Applications to Genetic Studies with Structured Populations

Sahir R Bhatnagar<sup>1,2</sup>, Karim Oualkacha<sup>3</sup>, Yi Yang<sup>4</sup>, Marie Forest<sup>2</sup>, and  
Celia MT Greenwood<sup>1,2,5</sup>

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

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

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

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

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

July 2, 2018

## Abstract

Complex traits are thought to be influenced by a combination of environmental factors and rare and common genetic variants. However, detection of such multivariate associations can be compromised by low statistical power and confounding by population structure. Linear mixed effect models (LMM) can account for correlations due to relatedness but are not applicable in high-dimensional (HD) settings where the number

of predictors greatly exceeds the number of samples. False negatives can result from two-stage approaches, where the residuals estimated from a null model adjusted for the subjects' relationship structure are subsequently used as the response in a standard penalized regression model. To overcome these challenges, we develop a general penalized LMM framework that simultaneously selects and estimates variables, accounting for between individual correlations, in one step. Our method can accommodate several sparsity inducing penalties such as the lasso, elastic net and group lasso, and also readily handles prior annotation information in the form of weights. We develop a groupwise-majorization descent algorithm which is highly scalable, computationally efficient and has theoretical guarantees of the convergence. Through simulations, we show that our method has better power over the two-stage approach, particularly for polygenic traits. We apply our method to identify SNPs that predict bone mineral density in the UK Biobank cohort. This approach can also be used to generate genetic risk scores and finding groups of predictors associated with the response, such as variants within a gene or pathway. Our algorithms are available in an R package (<https://github.com/sahirbhatnagar/ggmix>).

## 1 Introduction

Genome-wide association studies (GWAS) have become the standard method for analyzing genetic datasets owing to their success in identifying thousands of genetic variants associated with complex diseases (<https://www.genome.gov/gwastudies/>). Despite these impressive findings, the discovered markers have only been able to explain a small proportion of the phenotypic variance known as the missing heritability problem (1). One plausible explanation is that there are many causal variants that each explain a small amount of variation with small effect sizes (2). Methods such GWAS, which test each variant independently, are likely to miss these true associations due to the stringent significance thresholds required to reduce the number of false positives (1). Another major problem is that of confounding

due to geographic population structure, family and/or cryptic relatedness which can lead to spurious associations (3). For example, there may be subpopulations within a study that differ with respect to their genotype frequencies at a particular locus due to geographical location or their ancestry. This heterogeneity in genotype frequency can cause correlations with other loci and consequently mimic the signal of association even though there is no biological association (4, 5).

To address these two problems, several multivariable regression methods have been proposed that would simultaneously fit many SNPs in a single model while accounting for the population structure (4, 6, 7, 8). There are two main approaches to control for confounding by population structure: 1) the principal component (PC) adjustment method and 2) the linear mixed model (LMM). The PC adjustment method includes the top PCs of genome-wide single nucleotide polymorphism (SNP) genotypes as fixed effects while the LMM includes a random polygenic effect (3).

then talk about how this doesn't work when  $p > n$ . Talk about lasso motivations include fine mapping, mendel randomization,

## 1.1 Measures of Relatedness

"Although the exact genetic relationships between individuals in the samples are unknown, we could take advantage of the high-density genotype information to empirically estimate the level of relatedness between reportedly unrelated individuals (9)."

SNP genotypes can be coded as dummy variables with homozygotes being assigned a 0.0, heterozygotes being a 0.5, and opposite homozygotes being a 1.0 under an additive model or, for models involving dominance or recessive effects, with heterozygotes being assigned a 0.0 or 1.0, respectively. For the analyses we describe below, we assumed an additive model.

see <http://dalexander.bol.ucla.edu/preprints/admixture-preprint.pdf> for details

about confounding by population structure: “Cluster analysis directly seeks the ancestral clusters in the data, while principal component analysis (PCA) constructs low-dimensional projections of the data that explain the gross variation in marker genotypes, which in practice is the variation between populations”

In Table 2 we outline existing *multivariate* methods for genetic data containing related samples. MLMM (10) and LMM-Lasso (8) regress one trait (or phenotype) against multiple predictors (e.g. SNPs) while accounting for the population structure. Neither GCAT (4) nor QTCAT (11) use mixed models. GCAT uses an inverse regression approach coupled with logistic factor analysis. QTCAT doesn’t account for population structure, but instead searches for groups of highly correlated markers that are associated with the phenotype.

Note that there is confusion in the genetics literature on the meaning of multivariate linear mixed models. For example, GEMMA (12) is an association method for multiple traits against a single SNP, but is referred to in their paper as a “multivariate mixed model”. MTMM (13) is also an association method for multiple phenotypes against a single SNP but is referred to as a “multi-trait mixed model”.

See the review by (14) for comparison of *single* locus methods accounting for relatedness in GWAS with family data.

Table 1: Existing multivariate (multi-locus) methods for genetic data containing related samples.

Method	Software	Description
Multi-locus mixed-model (MLMM) (10)	<a href="#">R package on Github</a>	Approximate (2-step), stepwise mixed-model regression with forward inclusion and backward elimination. Since variance attributed to random polygenic term decreases when cofactors are added to the model, heritable variance estimate as a criterion to stop forward inclusion. Association testing.
LMM-Lasso (8)	<a href="#">Python code on Github</a>	Approximate (2-step), Laplacian shrinkage prior over the fixed effects. Optimize $\delta = \sigma_e^2 / \sigma_g^2$ . Stability selection used to assess significance
GCAT (4)	<a href="#">R package on Github</a>	Inverse regression approach where the association is tested by modeling genotype variation in terms of the trait plus model terms accounting for structure. The terms accounting for structure were based on the logistic factor analysis (15) approach
QTCAT (11)	<a href="#">R Package on Github</a>	Quantitative Trait Cluster Association Test. Do not account for population structure but instead search for clusters of highly correlated markers that are significantly associated to the phenotype using a hierarchical testing procedure for correlated covariates (16)

## 2 Penalized Mixed Models

### 2.1 Model Set-up

Let  $i = 1, \dots, N$  be the grouping index,  $j = 1, \dots, n_i$  the observation index within a group 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 the observed vector of responses,  $\mathbf{X}_i$  an  $n_i \times (p + 1)$  design matrix (with the column of 1s for the intercept),  $\mathbf{b}_i$  a group-specific random effect vector of length  $n_i$  and  $\boldsymbol{\varepsilon}_i = (\varepsilon_{i1}, \dots, \varepsilon_{in_i})$  the individual error terms. Furthermore, denote the stacked vectors  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)^T \in \mathbb{R}^{N_T \times 1}$ ,  $\mathbf{b} = (\mathbf{b}_1, \dots, \mathbf{b}_N)^T \in \mathbb{R}^{N_T \times 1}$ ,  $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_N)^T \in \mathbb{R}^{N_T \times 1}$ , and the stacked matrix  $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_N)^T \in \mathbb{R}^{N_T \times (p+1)}$ . Furthermore, let  $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^T \in \mathbb{R}^{(p+1) \times 1}$  a vector of fixed effects regression coefficients corresponding to  $\mathbf{X}$ . Following (17), we consider the following linear mixed model with a single random effect:

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

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

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

Here,  $\boldsymbol{\Phi}_{N_T \times N_T}$  is a known positive semi-definite and symmetric kinship matrix,  $\mathbf{I}_{N_T \times N_T}$  is the identity matrix and parameters  $\sigma^2$  and  $\eta \in [0, 1]$  determine how the variance is divided between  $\mathbf{b}$  and  $\boldsymbol{\varepsilon}$ . The joint density of  $\mathbf{Y}$  is multivariate normal:

$$\mathbf{Y} | (\boldsymbol{\beta}, \eta, \sigma^2) \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \eta\sigma^2\boldsymbol{\Phi} + (1 - \eta)\sigma^2\mathbf{I}) \quad (3)$$

Alternatively we may consider the parameterization in (18):

$$\mathbf{Y} | (\boldsymbol{\beta}, \delta, \sigma_g^2) \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma_g^2(\boldsymbol{\Phi} + \delta\mathbf{I})) \quad (4)$$

where  $\delta = \sigma_e^2/\sigma_g^2$ ,  $\sigma_g^2$  is the genetic variance and  $\sigma_e^2$  is the residual variance. (17) consider the parameterization in (3) since maximization is easier over the compact set  $\eta \in [0, 1]$  than over the unbounded interval  $\delta \in [0, \infty)$  as is done in (4) by (18).

Define the complete parameter vector  $\boldsymbol{\Theta} = (\boldsymbol{\beta}, \eta, \sigma^2)$ . The negative log-likelihood for (3) is given by

$$-\ell(\boldsymbol{\Theta}) \propto \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \log(\det(\mathbf{V})) + \frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \quad (5)$$

where  $\mathbf{V} = \eta\boldsymbol{\Phi} + (1 - \eta)\mathbf{I}$  and  $\det(\mathbf{V})$  is the determinant of  $\mathbf{V}$ . Let  $\boldsymbol{\Phi} = \mathbf{U}\mathbf{D}\mathbf{U}^T$  be the eigen (spectral) decomposition of the kinship matrix  $\boldsymbol{\Phi}$ , 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 (17)

$$\begin{aligned} \mathbf{V} &= \eta\boldsymbol{\Phi} + (1 - \eta)\mathbf{I} \\ &= \eta\mathbf{U}\mathbf{D}\mathbf{U}^T + (1 - \eta)\mathbf{U}\mathbf{I}\mathbf{U}^T \\ &= \mathbf{U}\eta\mathbf{D}\mathbf{U}^T + \mathbf{U}(1 - \eta)\mathbf{I}\mathbf{U}^T \\ &= \mathbf{U}(\eta\mathbf{D} + (1 - \eta)\mathbf{I})\mathbf{U}^T \\ &= \mathbf{U}\tilde{\mathbf{D}}\mathbf{U}^T \end{aligned} \quad (6)$$

where

$$\tilde{\mathbf{D}} = \eta \mathbf{D} + (1 - \eta) \mathbf{I} \quad (7)$$

$$\begin{aligned} &= \eta \begin{bmatrix} \Lambda_1 & & & \\ & \Lambda_2 & & \\ & & \ddots & \\ & & & \Lambda_{N_T} \end{bmatrix} + (1 - \eta) \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \\ &= \begin{bmatrix} 1 + \eta(\Lambda_1 - 1) & & & \\ & 1 + \eta(\Lambda_2 - 1) & & \\ & & \ddots & \\ & & & 1 + \eta(\Lambda_{N_T} - 1) \end{bmatrix} \\ &= \text{diag} \{1 + \eta(\Lambda_1 - 1), 1 + \eta(\Lambda_2 - 1), \dots, 1 + \eta(\Lambda_{N_T} - 1)\} \end{aligned} \quad (8)$$

Since (7) is a diagonal matrix, its inverse is also a diagonal matrix:

$$\tilde{\mathbf{D}}^{-1} = \text{diag} \left\{ \frac{1}{1 + \eta(\Lambda_1 - 1)}, \frac{1}{1 + \eta(\Lambda_2 - 1)}, \dots, \frac{1}{1 + \eta(\Lambda_{N_T} - 1)} \right\} \quad (9)$$

From (6) and (8),  $\log(\det(\mathbf{V}))$  simplifies to

$$\begin{aligned} \log(\det(\mathbf{V})) &= \log \left( \det(\mathbf{U}) \det(\tilde{\mathbf{D}}) \det(\mathbf{U}^T) \right) \\ &= \log \left\{ \prod_{i=1}^{N_T} (1 + \eta(\Lambda_i - 1)) \right\} \\ &= \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) \end{aligned} \quad (10)$$



since  $\det(\mathbf{U}) = 1$ . It also follows from (6) that

$$\begin{aligned}\mathbf{V}^{-1} &= (\mathbf{U}\tilde{\mathbf{D}}\mathbf{U}^T)^{-1} \\ &= (\mathbf{U}^T)^{-1} (\tilde{\mathbf{D}})^{-1} \mathbf{U}^{-1} \\ &= \mathbf{U}\tilde{\mathbf{D}}^{-1}\mathbf{U}^T\end{aligned}\tag{11}$$

since for an orthonormal matrix  $\mathbf{U}^{-1} = \mathbf{U}^T$ . Substituting (9), (10) and (11) into (5) the negative log-likelihood becomes

$$-\ell(\boldsymbol{\Theta}) \propto \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{U}\tilde{\mathbf{D}}^{-1}\mathbf{U}^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})\tag{12}$$

$$\begin{aligned}&= \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^2} (\mathbf{U}^T\mathbf{Y} - \mathbf{U}^T\mathbf{X}\boldsymbol{\beta})^T \tilde{\mathbf{D}}^{-1} (\mathbf{U}^T\mathbf{Y} - \mathbf{U}^T\mathbf{X}\boldsymbol{\beta}) \\ &= \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^2} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta})^T \tilde{\mathbf{D}}^{-1} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}) \\ &= \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^2} \sum_{i=1}^{N_T} \frac{(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1}\beta_j)^2}{1 + \eta(\Lambda_i - 1)}\end{aligned}\tag{13}$$

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  $\tilde{\mathbf{X}}$  and  $\mathbf{1}$  is a column vector of  $N_T$  ones.

## 2.2 Penalized Maximum Likelihood Estimator

We define the  $p + 3$  length vector of parameters  $\boldsymbol{\Theta} := (\Theta_0, \Theta_1, \dots, \Theta_{p+1}, \Theta_{p+2}, \Theta_{p+3}) = (\boldsymbol{\beta}, \eta, \sigma^2)$  where  $\boldsymbol{\beta} \in \mathbb{R}^{p+1}$ ,  $\eta \in [0, 1]$ ,  $\sigma^2 > 0$ . In what follows,  $p + 2$  and  $p + 3$  are the indices in  $\boldsymbol{\Theta}$  for  $\eta$  and  $\sigma^2$ , respectively. Define the objective function:

$$Q_\lambda(\boldsymbol{\Theta}) = f(\boldsymbol{\Theta}) + \lambda \sum_{j \neq 0} v_j P_j(\beta_j)\tag{14}$$

where  $f(\Theta) := -\ell(\Theta)$  is defined in (13),  $P_j(\cdot)$  is a penalty term on the fixed regression coefficients  $\beta_1, \dots, \beta_{p+1}$  (we do not penalize the intercept), controlled by the nonnegative regularization parameter  $\lambda$ , and  $v_j$  is the penalty factor for  $j$ th covariate. These penalty factors serve as a way of allowing parameters to be penalized differently. Note that we do not penalize  $\eta$  or  $\sigma^2$ . The penalty term is a necessary constraint because in our applications, the sample size is much smaller than the number of predictors. An estimate of the regression parameters  $\hat{\Theta}_\lambda$  is obtained by

$$\hat{\Theta}_\lambda = \arg \min_{\Theta} Q_\lambda(\Theta) \quad (15)$$

### 3 Computational Algorithm version 1

To solve for (15) we use a block relaxation technique (19) given by Algorithm 1

---

**Algorithm 1:** Block Relaxation Algorithm

---

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**

$$\beta^{(k+1)} \leftarrow \arg \min_{\beta} Q_\lambda \left( \beta, \eta^{(k)}, \sigma^{2(k)} \right)$$

$$\eta^{(k+1)} \leftarrow \arg \min_{\eta} Q_\lambda \left( \beta^{(k+1)}, \eta, \sigma^{2(k)} \right)$$

$$\sigma^{2(k+1)} \leftarrow \arg \min_{\sigma^2} Q_\lambda \left( \beta^{(k+1)}, \eta^{(k+1)}, \sigma^2 \right)$$

$k \leftarrow k + 1$

**until** convergence criterion is satisfied:  $\left\| \Theta^{(k+1)} - \Theta^{(k)} \right\|_2 < \epsilon$ ;

**end**

---

Below we discuss the specifics of Algorithm 1

### 3.1 Updates for the $\beta$ parameter

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

$$Q_\lambda(\Theta) = \frac{1}{2} \sum_{i=1}^{N_T} w_i \left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2 + \lambda \sum_{j=1}^p v_j |\beta_j| \quad (16)$$

where

$$w_i := \frac{1}{\sigma^2 (1 + \eta(\Lambda_i - 1))} \quad (17)$$

However `glmnet` solves the following problem:

$$\beta^{(k+1)} \leftarrow \arg \min_{\beta} \frac{1}{2 \sum_{i=1}^{N_T} \tilde{w}_i^{(k)}} \sum_{i=1}^{N_T} \tilde{w}_i^{(k)} \left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2 + \lambda \sum_{j=1}^p v_j |\beta_j| \quad (18)$$

where

$$\tilde{w}_i^{(k)} = N_T \cdot \frac{w_i^{(k)}}{\sum_{i=1}^{N_T} w_i^{(k)}} \quad (19)$$

Note that  $\sum_i \tilde{w}_i^{(k)} = N_T$ . We can simplify (18) to be:

$$\begin{aligned} \beta^{(k+1)} &\leftarrow \arg \min_{\beta} \frac{1}{2N_T} \sum_{i=1}^{N_T} N_T \cdot \frac{w_i^{(k)}}{\sum_{i=1}^{N_T} w_i^{(k)}} \left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2 + \lambda \sum_{j=1}^p v_j |\beta_j| \\ \beta^{(k+1)} &\leftarrow \arg \min_{\beta} \frac{1}{2 \sum_{i=1}^{N_T} w_i^{(k)}} \sum_{i=1}^{N_T} w_i^{(k)} \left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2 + \lambda \sum_{j=1}^p v_j |\beta_j| \end{aligned} \quad (20)$$

In order to make (16) to be in the form of (20), we must scale the lambda accordingly:

$$\beta^{(k+1)} \leftarrow \arg \min_{\beta} \frac{1}{2} \sum_{i=1}^{N_T} w_i^{(k)} \left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2 + \frac{\lambda}{\sum_{i=1}^{N_T} w_i^{(k)}} \sum_{j=1}^p v_j |\beta_j| \quad (21)$$

Conditional on  $\eta^{(k)}$  and  $\sigma^{2(k)}$ , it can be shown that the solution for  $\beta$  is a weighted lasso problem with observation weights given by (17).

The full derivation is given in Section 7.1. Therefore,  $\beta^{(k+1)}$  can be efficiently solved using the `glmnet` algorithm (20). Note that the rescaling of the weights to sum to  $N_T$  is what is being done in `glmnet`.

### 3.2 Updates for the $\eta$ paramter

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

Given  $\beta^{(k+1)}$  and  $\sigma^{2(k)}$ , solving for  $\eta^{(k+1)}$  becomes a univariate optimization problem. We use a bound constrained optimization algorithm (21) implemented in the `optim` function in R and set the lower and upper bounds to be 0 and 1, respectively.

### 3.3 Updates for the $\sigma^2$ parameter

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

Conditional on  $\beta^{(k+1)}$  and  $\eta^{(k+1)}$ , there exists an analytic solution for  $\sigma^{2(k+1)}$ :

$$\begin{aligned} \frac{\partial}{\partial \sigma^2} Q_{\lambda}(\Theta) &= \frac{N_T}{2\sigma^2} - \frac{1}{2\sigma^4} \sum_{i=1}^{N_T} \frac{\left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j^{(k+1)} \right)^2}{1 + \eta^{(k+1)}(\Lambda_i - 1)} = 0 \\ \sigma^{2(k+1)} &= \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{\left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j^{(k+1)} \right)^2}{1 + \eta^{(k+1)}(\Lambda_i - 1)} \end{aligned} \quad (24)$$

### 3.4 Regularization path

Recall that our objective function has the form

$$Q_\lambda(\Theta) = \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2} \sum_{i=1}^{N_T} w_i \left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2 + \lambda \sum_{j=1}^p v_j |\beta_j| \quad (25)$$

The Karush-Kuhn-Tucker (KKT) optimality conditions for (25) are given by:

$$\begin{aligned} \frac{\partial}{\partial \beta_1, \dots, \beta_p} Q_\lambda(\Theta) &= \mathbf{0}_p \\ \frac{\partial}{\partial \beta_0} Q_\lambda(\Theta) &= 0 \\ \frac{\partial}{\partial \eta} Q_\lambda(\Theta) &= 0 \\ \frac{\partial}{\partial \sigma^2} Q_\lambda(\Theta) &= 0 \end{aligned} \quad (26)$$

The equations in (26) are equivalent to

$$\begin{aligned} \sum_{i=1}^{N_T} w_i \tilde{X}_{i1} \left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right) &= 0 \\ \frac{1}{v_j} \sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right) &= \lambda \gamma_j, \\ \gamma_j &\in \begin{cases} \text{sign}(\hat{\beta}_j) & \text{if } \hat{\beta}_j \neq 0 \\ [-1, 1] & \text{if } \hat{\beta}_j = 0 \end{cases}, \quad \text{for } j = 1, \dots, p \\ \frac{1}{2} \sum_{i=1}^{N_T} \frac{\Lambda_i - 1}{1 + \eta(\Lambda_i - 1)} \left( 1 - \frac{\left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2}{\sigma^2(1 + \eta(\Lambda_i - 1))} \right) &= 0 \\ \sigma^2 - \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{\left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2}{1 + \eta(\Lambda_i - 1)} &= 0 \end{aligned} \quad (27)$$

where  $w_i$  is given by (17),  $\tilde{\mathbf{X}}_{-1}^T$  is  $\tilde{\mathbf{X}}^T$  with the first column removed,  $\tilde{\mathbf{X}}_1^T$  is the first column

of  $\tilde{\mathbf{X}}^T$ , and  $\boldsymbol{\gamma} \in \mathbb{R}^p$  is the subgradient function of the  $\ell_1$  norm evaluated at  $(\hat{\beta}_1, \dots, \hat{\beta}_p)$ . Therefore  $\hat{\boldsymbol{\Theta}}$  is a solution in (15) if and only if  $\hat{\boldsymbol{\Theta}}$  satisfies (27) for some  $\gamma$ .

we find the solution for the other parameters such that the KKT conditions are verified.  
page 17 of ss with learning

Therefore we can determine a decreasing sequence of tuning parameters by starting at a maximal value for  $\lambda = \lambda_{max}$  for which  $\hat{\beta}_j = 0$  for  $j = 1, \dots, p$ . In this case, the KKT conditions in (27) are equivalent to

$$\begin{aligned} \frac{1}{v_j} \sum_{i=1}^{N_T} \left| w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \tilde{X}_{i1} \beta_0 \right) \right| &\leq \lambda, \quad \forall j = 1, \dots, p \\ \beta_0 &= \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{i1} \tilde{Y}_i}{\sum_{i=1}^{N_T} w_i \tilde{X}_{i1}^2} \\ \frac{1}{2} \sum_{i=1}^{N_T} \frac{\Lambda_i - 1}{1 + \eta(\Lambda_i - 1)} \left( 1 - \frac{\left( \tilde{Y}_i - \tilde{X}_{i1} \beta_0 \right)^2}{\sigma^2 (1 + \eta(\Lambda_i - 1))} \right) &= 0 \\ \sigma^2 &= \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{\left( \tilde{Y}_i - \tilde{X}_{i1} \beta_0 \right)^2}{1 + \eta(\Lambda_i - 1)} \end{aligned} \quad (28)$$

We can solve the KKT system of equations in (28) (with a numerical solution for  $\eta$ ) in order to have an explicit form of the stationary point  $\hat{\boldsymbol{\Theta}}_0 = \{\hat{\beta}_0, \mathbf{0}_p, \hat{\eta}, \hat{\sigma}^2\}$ . Once we have  $\hat{\boldsymbol{\Theta}}_0$ , we can solve for the smallest value of  $\lambda$  such that the entire vector  $(\hat{\beta}_1, \dots, \hat{\beta}_p)$  is 0:

$$\lambda_{max} = \max_j \left\{ \left| \frac{1}{v_j} \sum_{i=1}^{N_T} \tilde{w}_i \tilde{X}_{ij} \left( \tilde{Y}_i - \tilde{X}_{i1} \hat{\beta}_0 \right) \right| \right\}, \quad j = 1, \dots, p \quad (29)$$

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

### 3.5 Warm Starts

The way in which we have derived the sequence of tuning parameters using the KKT conditions, allows us to implement warm starts. That is, the solution  $\hat{\Theta}$  for  $\lambda_k$  is used as the initial value  $\Theta^{(0)}$  for  $\lambda_{k+1}$ .

### 3.6 Prediction of the random effects

We use an empirical Bayes approach (e.g. (22)) to predict the random effects  $\mathbf{b}$ . Let the maximum a posteriori (MAP) estimate be defined as

$$\hat{\mathbf{b}} = \arg \max_{\mathbf{b}} f(\mathbf{b}|\mathbf{Y}, \boldsymbol{\beta}, \eta, \sigma^2) \quad (30)$$

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

$$\begin{aligned} f(\mathbf{b}|\mathbf{Y}, \boldsymbol{\beta}, \eta, \sigma^2) &= \frac{f(\mathbf{Y}|\mathbf{b}, \boldsymbol{\beta}, \eta, \sigma^2)\pi(\mathbf{b}|\eta, \sigma^2)}{f(\mathbf{Y}|\boldsymbol{\beta}, \eta, \sigma^2)} \\ &\propto f(\mathbf{Y}|\mathbf{b}, \boldsymbol{\beta}, \eta, \sigma^2)\pi(\mathbf{b}|\eta, \sigma^2) \\ &\propto \exp \left\{ -\frac{1}{2\sigma^2}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b})^T \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b}) - \frac{1}{2\eta\sigma^2} \mathbf{b}^T \boldsymbol{\Phi}^{-1} \mathbf{b} \right\} \\ &= \exp \left\{ -\frac{1}{2\sigma^2} \left[ (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b})^T \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b}) + \frac{1}{\eta} \mathbf{b}^T \boldsymbol{\Phi}^{-1} \mathbf{b} \right] \right\} \end{aligned} \quad (31)$$

Solving for (30) is equivalent to minimizing the exponent in (31):

$$\hat{\mathbf{b}} = \arg \min_{\mathbf{b}} \left\{ (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b})^T \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b}) + \frac{1}{\eta} \mathbf{b}^T \boldsymbol{\Phi}^{-1} \mathbf{b} \right\} \quad (32)$$

Taking the derivative of (32) with respect to  $\mathbf{b}$  and setting it to 0 we get:

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

where  $\mathbf{V}^{-1}$  is given by (11), and  $(\hat{\boldsymbol{\beta}}, \hat{\eta})$  are the estimates obtained from Algorithm 1.

### 3.7 Choice of the tuning parameter

We use the BIC:

$$BIC_{\lambda} = -2\ell(\hat{\boldsymbol{\beta}}, \hat{\sigma}^2, \hat{\eta}) + c \cdot \hat{df}_{\lambda} \quad (33)$$

where  $\hat{df}_{\lambda}$  is the number of non-zero elements in  $\hat{\boldsymbol{\beta}}_{\lambda}$  (23) plus two (representing the variance parameters  $\eta$  and  $\sigma^2$ ). Several authors have used this criterion for variable selection in mixed models with  $c = \log N_T$  (24, 25) and  $c = \log N$  (26) (where  $N$  is the number of groups). Other authors have proposed  $c = \log(\log(N_T)) * \log(N_T)$  (27).

## 4 Low rank similarity matrix

Let  $\mathbf{K} \in \mathbb{R}^{N_T \times k}$  be the matrix containing the  $k$  SNPs used to compute the factored kinship matrix  $\boldsymbol{\Phi}$  given by

$$\boldsymbol{\Phi} = \mathbf{K}\mathbf{K}^T \quad (34)$$



Furthermore, let  $\mathbf{K} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$  be the singular value decomposition (SVD) of  $\mathbf{K}$ . Plugging this into (34) we get

$$\begin{aligned}
 \Phi &= (\mathbf{U}\mathbf{\Lambda}\mathbf{V}^T) (\mathbf{U}\mathbf{\Lambda}\mathbf{V}^T)^T \\
 &= \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T\mathbf{V}\mathbf{\Lambda}\mathbf{U}^T \\
 &= \mathbf{U}\mathbf{\Lambda}\mathbf{\Lambda}\mathbf{U}^T \\
 &= \mathbf{U}\mathbf{\Sigma}\mathbf{U}^T,
 \end{aligned} \tag{35}$$

Therefore, the eigenvectors of  $\Phi$  are equal to the singular vectors of  $\mathbf{K}$  (denoted by  $\mathbf{U}$ ), and the eigenvalues of  $\Phi$  (denoted by the diagonal matrix  $\mathbf{\Sigma}$ ) are equal to the square of the singular values of  $\mathbf{K}$  (28). This allows us to bypass the explicit computation of the kinship matrix by directly applying SVD on the SNP matrix  $\mathbf{W}$ . (18) noted that the computational time for fitting the LMM can be reduced if the matrix  $\mathbf{K}$  is not full rank, i.e., when  $k < N_T$ . This is due to the fact that the matrix  $\mathbf{D}_{N_T \times N_T}$  contains  $k$  non-zero eigenvalues followed by  $N_T - k$  zeros on the diagonal. Let  $\mathbf{U} \equiv [\mathbf{U}_1 \ \mathbf{U}_2]$ , where  $\mathbf{U}_1 \in \mathbb{R}^{N_T \times k}$  and  $\mathbf{U}_2 \in \mathbb{R}^{N_T \times (N_T - k)}$  are the matrices of singular vectors corresponding to the  $k$  non-zero and  $N_T - k$  zero eigenvalues, respectively. Then (35) can be written as

$$\Phi = \mathbf{U}_1 \mathbf{\Sigma} \mathbf{U}_1^T \tag{36}$$

We now try to simplify the log-likelihood (12). Since there are  $N_T - k$  zero eigenvalues, the second term in (12) reduces to

$$\frac{1}{2} \left( \sum_{i=1}^k \log(1 + \eta(\Sigma_i - 1)) + (N_T - k) \log(1 - \eta) \right) \tag{37}$$

where  $\Sigma_i = \Lambda_i^2$ , and  $\Lambda_i$  is the  $i^{\text{th}}$  singular value of  $\mathbf{W}$ . Let  $a \equiv (\mathbf{Y} - \mathbf{X}\beta)$ . The third term

in (12) can be written as

$$\begin{aligned} \frac{1}{2\sigma^2} a^T [\eta \Phi + (1 - \eta) \mathbf{I}_n]^{-1} a &= \frac{1}{2\sigma^2} a^T [\eta \mathbf{U}_1 \Sigma_1 \mathbf{U}_1^T + (1 - \eta) \mathbf{I}_n]^{-1} a \\ &= \frac{1}{2\sigma^2} a^T [\mathbf{C} \mathbf{B} \mathbf{C}^T + \mathbf{A}]^{-1} a \end{aligned}$$

where

$$\mathbf{A} = (1 - \eta) \mathbf{I}_n$$

$$\mathbf{B} = \Sigma_1$$

$$\mathbf{C} = \sqrt{\eta} \mathbf{U}_1$$

$$\mathbf{C}^T = \sqrt{\eta} \mathbf{U}_1^T$$

Assuming  $\mathbf{C} \mathbf{B} \mathbf{C}^T + \mathbf{A}$  is non-singular, the inverse of  $[\mathbf{C} \mathbf{B} \mathbf{C}^T + \mathbf{A}]$  is given explicitly by the Woodbury formula (29)

$$(\mathbf{A} + \mathbf{C} \mathbf{B} \mathbf{C}^T)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{C} (\mathbf{B}^{-1} + \mathbf{C}^T \mathbf{A}^{-1} \mathbf{C})^{-1} \mathbf{C}^T \mathbf{A}^{-1} \quad (38)$$

Substituting the values for  $\mathbf{A}, \mathbf{B}$  and  $\mathbf{C}$  into (38) we get

$$\begin{aligned} (\mathbf{A} + \mathbf{C} \mathbf{B} \mathbf{C}^T)^{-1} &= \frac{1}{1 - \eta} \mathbf{I}_{N_T} - \frac{\sqrt{\eta}}{1 - \eta} \mathbf{I}_{N_T} \mathbf{U}_1 \left( \Sigma_1^{-1} + \frac{\eta}{1 - \eta} \mathbf{U}_1^T \mathbf{I}_{N_T} \mathbf{U}_1 \right)^{-1} \frac{\sqrt{\eta}}{1 - \eta} \mathbf{U}_1^T \mathbf{I}_{N_T} \\ &= \frac{1}{1 - \eta} \left[ \mathbf{I}_{N_T} - \frac{\eta}{1 - \eta} \mathbf{U}_1 \left( \Sigma_1^{-1} + \frac{\eta}{1 - \eta} \mathbf{I}_k \right)^{-1} \mathbf{U}_1^T \right] \\ &= \frac{1}{1 - \eta} \left[ \mathbf{I}_{N_T} - \frac{\eta}{1 - \eta} \mathbf{U}_1 \left( \frac{\eta}{1 - \eta} \left( \frac{1 - \eta}{\eta} \Sigma_1^{-1} + \mathbf{I}_k \right) \right)^{-1} \mathbf{U}_1^T \right] \\ &= \frac{1}{1 - \eta} \left[ \mathbf{I}_{N_T} - \mathbf{U}_1 \left( \frac{1 - \eta}{\eta} \Sigma_1^{-1} + \mathbf{I}_k \right)^{-1} \mathbf{U}_1^T \right] \end{aligned} \quad (39)$$

where we have used the following identities:  $\mathbf{I}_k = \mathbf{U}_1^T \mathbf{U}_1$ ,  $\mathbf{I}_{N_T - k} = \mathbf{U}_2^T \mathbf{U}_2$ .

Substituting (37) and (39) in (12) we obtain

$$\begin{aligned}
 -\ell(\boldsymbol{\Theta}) \propto & \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \left( \sum_{i=1}^k \log(1 + \eta(\Sigma_i - 1)) + (N_T - k) \log(1 - \eta) \right) + \\
 & \frac{1}{2} \left\{ (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \left[ \frac{1}{\sigma^2(1 - \eta)} \left( \mathbf{I}_{N_T} - \mathbf{U}_1 \left( \frac{1 - \eta}{\eta} \boldsymbol{\Sigma}_1^{-1} + \mathbf{I}_k \right)^{-1} \mathbf{U}_1^T \right) \right] (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \right\}
 \end{aligned} \tag{40}$$

## 5 Group Lasso with Low-rank Similarity Matrix

This section focuses on the part of the log-likelihood (40) that depends on  $\boldsymbol{\beta}$ .

### 5.1 Model

Only the third term of the log-likelihood (40) depends on  $\boldsymbol{\beta}$ :

$$\frac{1}{2} \left\{ (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \left[ \frac{1}{\sigma^2(1 - \eta)} \left( \mathbf{I}_{N_T} - \mathbf{U}_1 \left( \frac{1 - \eta}{\eta} \boldsymbol{\Sigma}_1^{-1} + \mathbf{I}_k \right)^{-1} \mathbf{U}_1^T \right) \right] (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \right\} \tag{41}$$

Equation (41) can be written more generally as

$$L(\boldsymbol{\beta} \mid \mathbf{D}) = \frac{1}{2} [\mathbf{Y} - \hat{\mathbf{Y}}]^T \mathbf{W} [\mathbf{Y} - \hat{\mathbf{Y}}]$$

where  $\hat{\mathbf{Y}} = \sum_{j=1}^p \beta_j X_j$ ,  $\mathbf{D}$  is the working data  $\{\mathbf{Y}, \mathbf{X}\}$ , and  $\mathbf{W}$  is an  $N_T \times N_T$  weight matrix given by

$$\mathbf{W} = \frac{1}{\sigma^2(1 - \eta)} \left( \mathbf{I}_{N_T} - \mathbf{U}_1 \left( \frac{1 - \eta}{\eta} \boldsymbol{\Sigma}_1^{-1} + \mathbf{I}_k \right)^{-1} \mathbf{U}_1^T \right) \tag{42}$$

Assume that we the predictors in the design matrix  $\mathbf{X} \in \mathbb{R}^{N_T \times p}$  belong to  $K$  groups and that the group membership is already defined such that  $(1, 2, \dots, p) = \bigcup_{k=1}^K I_k$  and the cardinality of index set  $I_k$  is  $p_k$ ,  $I_k \cap I_{k'} = \emptyset$  for  $k \neq k'$ ,  $1 \leq k, k' \leq K$ . Thus group  $k$  contains  $p_k$

predictors, which are  $x_j$ 's for  $j \in I_k$ , and  $1 \leq k \leq K$ . If an intercept is included, then  $I_1 = \{1\}$ . Given the group partition, we use  $\boldsymbol{\beta}_{(k)}$  to denote the segment of  $\boldsymbol{\beta}$  corresponding to group  $k$ . This notation is used for any  $p$ -dimensional vector. We consider the group lasso penalized estimator

$$\min_{\boldsymbol{\beta}} L(\boldsymbol{\beta} \mid \mathbf{D}) + \lambda \sum_{k=1}^K w_k \|\boldsymbol{\beta}_{(k)}\|_2, \quad (43)$$

The loss function  $L$  satisfies the quadratic majorization (QM) condition, since there exists a  $p \times p$  matrix  $\mathbf{H} = \mathbf{X}^\top \mathbf{W} \mathbf{X}$ , and  $\nabla L(\boldsymbol{\beta} \mid \mathbf{D}) = -\left(Y - \hat{Y}\right)^\top \mathbf{W} \mathbf{X}$ , which may only depend on the data  $\mathbf{D}$ , such that for all  $\boldsymbol{\beta}, \boldsymbol{\beta}^*$ ,

$$L(\boldsymbol{\beta} \mid \mathbf{D}) \leq L(\boldsymbol{\beta}^* \mid \mathbf{D}) + (\boldsymbol{\beta} - \boldsymbol{\beta}^*)^\top \nabla L(\boldsymbol{\beta}^* \mid \mathbf{D}) + \frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\beta}^*)^\top \mathbf{H}(\boldsymbol{\beta} - \boldsymbol{\beta}^*). \quad (44)$$

## 5.2 Algorithm

Noticing that the penalty term  $\sum_{k=1}^K w_k \|\boldsymbol{\beta}_{(k)}\|_2$  is separable with respect to the indices of the features  $k = 1, \dots, K$ , we can derive the *groupwise-majorization-descent* (GMD) algorithm for computing the solution of (43) when the loss function satisfies the QM condition. Let  $\tilde{\boldsymbol{\beta}}$  denote the current solution of  $\boldsymbol{\beta}$ . Without loss of generality, let us derive the GMD update of  $\tilde{\boldsymbol{\beta}}_{(k)}$ , the coefficients of group  $k$ . Define  $\mathbf{H}_k$  as the sub-matrix of  $\mathbf{H}$  corresponding to group  $k$ . For example, if group 2 is  $\{2, 4\}$  then  $\mathbf{H}_{(2)}$  is a  $2 \times 2$  matrix with

$$\mathbf{H}_{(2)} = \begin{bmatrix} h_{2,2} & h_{2,4} \\ h_{4,2} & h_{4,4} \end{bmatrix},$$

where  $h_{i,j}$  is the  $i, j$ th entry of the  $\mathbf{H}$  matrix. Write  $\boldsymbol{\beta}$  such that  $\boldsymbol{\beta}_{(k')} = \tilde{\boldsymbol{\beta}}_{(k')}$  for  $k' \neq k$ . Given  $\boldsymbol{\beta}_{(k')} = \tilde{\boldsymbol{\beta}}_{(k')}$  for  $k' \neq k$ , the optimal  $\boldsymbol{\beta}_{(k)}$  is defined as

$$\arg \min_{\boldsymbol{\beta}_{(k)}} L(\boldsymbol{\beta} \mid \mathbf{D}) + \lambda w_k \|\boldsymbol{\beta}_{(k)}\|_2. \quad (45)$$

Unfortunately, there is no closed form solution to (45) for a general loss function with general design matrix. We overcome the computational obstacle by taking advantage of the QM condition. From (44) we have

$$L(\boldsymbol{\beta} \mid \mathbf{D}) \leq L(\tilde{\boldsymbol{\beta}} \mid \mathbf{D}) + (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})^\top \nabla L(\tilde{\boldsymbol{\beta}} \mid \mathbf{D}) + \frac{1}{2}(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})^\top \mathbf{H}(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}).$$

Write  $U(\tilde{\boldsymbol{\beta}}) = -\nabla L(\tilde{\boldsymbol{\beta}} \mid \mathbf{D})$ . Using

$$\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}} = (\underbrace{0, \dots, 0}_{k-1}, \boldsymbol{\beta}_{(k)} - \tilde{\boldsymbol{\beta}}_{(k)}, \underbrace{0, \dots, 0}_{K-k}),$$

we can write

$$L(\boldsymbol{\beta} \mid \mathbf{D}) \leq L(\tilde{\boldsymbol{\beta}} \mid \mathbf{D}) - (\boldsymbol{\beta}_{(k)} - \tilde{\boldsymbol{\beta}}_{(k)})^\top U_{(k)} + \frac{1}{2}(\boldsymbol{\beta}_{(k)} - \tilde{\boldsymbol{\beta}}_{(k)})^\top \mathbf{H}_{(k)}(\boldsymbol{\beta}_{(k)} - \tilde{\boldsymbol{\beta}}_{(k)}). \quad (46)$$

where

$$U_{(k)} = \frac{\partial}{\partial \boldsymbol{\beta}_{(k)}} L_Q(\boldsymbol{\beta} \mid \mathbf{D}) = -\left(Y - \hat{Y}\right)^\top \mathbf{W} \mathbf{X}_{(k)}, \quad (47)$$

$$\mathbf{H}_{(k)} = \frac{\partial^2}{\partial \boldsymbol{\beta}_{(k)} \partial \boldsymbol{\beta}_{(k)}^\top} L_Q(\boldsymbol{\beta} \mid \mathbf{D}) = \mathbf{X}_{(k)}^\top \mathbf{W} \mathbf{X}_{(k)}. \quad (48)$$

Let  $\eta_k$  be the largest eigenvalue of  $\mathbf{H}_{(k)}$ . We set  $\gamma_k = (1 + \varepsilon^*)\eta_k$ , where  $\varepsilon^* = 10^{-6}$ . Then we can further relax the upper bound in (46) as

$$L(\boldsymbol{\beta} \mid \mathbf{D}) \leq L(\tilde{\boldsymbol{\beta}} \mid \mathbf{D}) - (\boldsymbol{\beta}^{(k)} - \tilde{\boldsymbol{\beta}}^{(k)})^\top U_{(k)} + \frac{1}{2}\gamma_k(\boldsymbol{\beta}^{(k)} - \tilde{\boldsymbol{\beta}}^{(k)})^\top (\boldsymbol{\beta}^{(k)} - \tilde{\boldsymbol{\beta}}^{(k)}). \quad (49)$$

It is important to note that the inequality strictly holds unless for  $\boldsymbol{\beta}^{(k)} = \tilde{\boldsymbol{\beta}}^{(k)}$ . Instead of

minimizing (45) we solve

$$\arg \min_{\tilde{\boldsymbol{\beta}}^{(k)}} L(\tilde{\boldsymbol{\beta}} \mid \mathbf{D}) - (\boldsymbol{\beta}^{(k)} - \tilde{\boldsymbol{\beta}}^{(k)})^\top U_{(k)} + \frac{1}{2} \gamma_k (\boldsymbol{\beta}^{(k)} - \tilde{\boldsymbol{\beta}}^{(k)})^\top (\boldsymbol{\beta}^{(k)} - \tilde{\boldsymbol{\beta}}^{(k)}) + \lambda w_k \|\boldsymbol{\beta}^{(k)}\|_2. \quad (50)$$

Denote by  $\tilde{\boldsymbol{\beta}}^{(k)}(\text{new})$  the solution to (50). It is straightforward to see that  $\tilde{\boldsymbol{\beta}}^{(k)}(\text{new})$  has a simple closed-form expression

$$\tilde{\boldsymbol{\beta}}^{(k)}(\text{new}) = \frac{1}{\gamma_k} \left( U_{(k)} + \gamma_k \tilde{\boldsymbol{\beta}}^{(k)} \right) \left( 1 - \frac{\lambda w_k}{\|U_{(k)} + \gamma_k \tilde{\boldsymbol{\beta}}^{(k)}\|_2} \right)_+. \quad (51)$$

Algorithm 2 summarizes the details of GMD.

---

**Algorithm 2:** The GMD algorithm for general group-lasso learning.

---

1. For  $k = 1, \dots, K$ , compute  $\gamma_k$ , the largest eigenvalue of  $\mathbf{H}^{(k)}$ .
  2. Initialize  $\tilde{\boldsymbol{\beta}}$ .
  3. Repeat the following cyclic groupwise updates until convergence:
    - for  $k = 1, \dots, K$ , do step (3.1)–(3.3)
      - 3.1 Compute  $U(\tilde{\boldsymbol{\beta}}) = -\nabla L(\tilde{\boldsymbol{\beta}} \mid \mathbf{D})$ .
      - 3.2 Compute  $\tilde{\boldsymbol{\beta}}^{(k)}(\text{new}) = \frac{1}{\gamma_k} \left( U_{(k)} + \gamma_k \tilde{\boldsymbol{\beta}}^{(k)} \right) \left( 1 - \frac{\lambda w_k}{\|U_{(k)} + \gamma_k \tilde{\boldsymbol{\beta}}^{(k)}\|_2} \right)_+.$
      - 3.3 Set  $\tilde{\boldsymbol{\beta}}^{(k)} = \tilde{\boldsymbol{\beta}}^{(k)}(\text{new})$ .
- 

### 5.3 Convergence

We can prove the strict descent property of GMD by using the MM principle (30, 31, 32).

Define

$$Q(\boldsymbol{\beta} \mid \mathbf{D}) = L(\tilde{\boldsymbol{\beta}} \mid \mathbf{D}) - (\boldsymbol{\beta}^{(k)} - \tilde{\boldsymbol{\beta}}^{(k)})^\top U_{(k)} + \frac{1}{2} \gamma_k (\boldsymbol{\beta}^{(k)} - \tilde{\boldsymbol{\beta}}^{(k)})^\top (\boldsymbol{\beta}^{(k)} - \tilde{\boldsymbol{\beta}}^{(k)}) + \lambda w_k \|\boldsymbol{\beta}^{(k)}\|_2. \quad (52)$$

Obviously,  $Q(\boldsymbol{\beta} \mid \mathbf{D}) = L(\boldsymbol{\beta} \mid \mathbf{D}) + \lambda w_k \|\boldsymbol{\beta}^{(k)}\|_2$  when  $\boldsymbol{\beta}^{(k)} = \tilde{\boldsymbol{\beta}}^{(k)}$  and (??) shows that  $Q(\boldsymbol{\beta} \mid \mathbf{D}) > L(\boldsymbol{\beta} \mid \mathbf{D}) + \lambda w_k \|\boldsymbol{\beta}^{(k)}\|_2$  when  $\boldsymbol{\beta}^{(k)} \neq \tilde{\boldsymbol{\beta}}^{(k)}$ . After updating  $\tilde{\boldsymbol{\beta}}^{(k)}$  using (??), we

have

$$\begin{aligned}
 L(\tilde{\boldsymbol{\beta}}^{(k)}(\text{new}) \mid \mathbf{D}) + \lambda w_k \|\tilde{\boldsymbol{\beta}}^{(k)}(\text{new})\|_2 &\leq Q(\tilde{\boldsymbol{\beta}}^{(k)}(\text{new}) \mid \mathbf{D}) \\
 &\leq Q(\tilde{\boldsymbol{\beta}} \mid \mathbf{D}) \\
 &= L(\tilde{\boldsymbol{\beta}} \mid \mathbf{D}) + \lambda w_k \|\tilde{\boldsymbol{\beta}}^{(k)}\|_2.
 \end{aligned}$$

Moreover, if  $\tilde{\boldsymbol{\beta}}^{(k)}(\text{new}) \neq \tilde{\boldsymbol{\beta}}^{(k)}$ , then the first inequality becomes

$$L(\tilde{\boldsymbol{\beta}}^{(k)}(\text{new}) \mid \mathbf{D}) + \lambda w_k \|\tilde{\boldsymbol{\beta}}^{(k)}(\text{new})\|_2 < Q(\tilde{\boldsymbol{\beta}}^{(k)}(\text{new}) \mid \mathbf{D}).$$

Therefore, the objective function is strictly decreased after updating all groups in a cycle, unless the solution does not change after each groupwise update. If this is the case, we can show that the solution must satisfy the KKT conditions, which means that the algorithm converges and finds the right answer. To see this, if  $\tilde{\boldsymbol{\beta}}^{(k)}(\text{new}) = \tilde{\boldsymbol{\beta}}^{(k)}$  for all  $k$ , then by the update formula (51) we have that for all  $k$

$$\tilde{\boldsymbol{\beta}}^{(k)} = \frac{1}{\gamma_k} \left( U^{(k)} + \gamma_k \tilde{\boldsymbol{\beta}}^{(k)} \right) \left( 1 - \frac{\lambda w_k}{\|U^{(k)} + \gamma_k \tilde{\boldsymbol{\beta}}^{(k)}\|_2} \right) \quad \text{if } \|U^{(k)} + \gamma_k \tilde{\boldsymbol{\beta}}^{(k)}\|_2 > \lambda w_k, \quad (53)$$

$$\tilde{\boldsymbol{\beta}}^{(k)} = \mathbf{0} \quad \text{if } \|U^{(k)} + \gamma_k \tilde{\boldsymbol{\beta}}^{(k)}\|_2 \leq \lambda w_k. \quad (54)$$

By straightforward algebra we obtain the KKT conditions:

$$\begin{aligned}
 -U^{(k)} + \lambda w_k \cdot \frac{\tilde{\boldsymbol{\beta}}^{(k)}}{\|\tilde{\boldsymbol{\beta}}^{(k)}\|_2} &= \mathbf{0} \quad \text{if } \tilde{\boldsymbol{\beta}}^{(k)} \neq \mathbf{0}, \\
 \|U^{(k)}\|_2 &\leq \lambda w_k \quad \text{if } \tilde{\boldsymbol{\beta}}^{(k)} = \mathbf{0},
 \end{aligned}$$

where  $k = 1, 2, \dots, K$ . Therefore, if the objective function stays unchanged after a cycle, the algorithm necessarily converges to the right answer.

## 5.4 Fitting Options and Algorithms

Recall  $\mathbf{K} \in \mathbb{R}^{N_T \times k}$  is the matrix containing the  $k$  SNPs used to compute the factored kinship matrix  $\Phi$ . The dimension of this matrix will determine the algorithm used as shown in the table below.

Table 2: Algorithm used based on dimension of  $\mathbf{K}$ .

Dimension of $\mathbf{K}$	lasso	group lasso
$N_T > k$	gcdnet (or degenerate gglasso)	gglasso (GMD Algorithm with weight matrix)
$N_T < k$	glmnet (Coordinate descent with observation weights)	gglasso (GMD Algorithm with observation weights)

## 6 Simulation Study

To assess the performance of penfam we used genotyped data from the UK Biobank cohort to maintain LD structure. We restricted our simulation study to 1st degree relatives defined by the KING estimate for kinship coefficients. We define the following quantities:

- $c$ : percentage of causal SNPs
- $\rho$ : linkage disequilibrium between two SNPs
- $\mathbf{X}^{(test)}$ :  $n \times 1000$  matrix of SNPs that have been randomly sampled across the genome, with sampling weights proportional to the size of each chromosome. These are the SNPs that will be included as fixed effects in our model.



- $\mathbf{X}^{(causal)}$ :  $n \times (c \times 1000)$  matrix of SNPs out of the SNPs included in the fixed effect model that will be truly associated with the simulated phenotype, where  $\mathbf{X}^{(causal)} \subseteq \mathbf{X}^{(test)}$
- $\mathbf{X}^{(other)}$ :  $n \times 4000$  matrix of SNPs that have been randomly sampled across the genome, with sampling weights proportional to the size of each chromosome. This matrix will be used in the construction of the kinship matrix. Some of these  $\mathbf{X}^{(other)}$  SNPs, in conjunction with some of the SNPs in  $\mathbf{X}^{(test)}$  will be used in construction of the kinship matrix. We will alter the balance between these two contributors and with the proportion of causal SNPs used to calculate kinship. The maximum LD between any two SNPs in  $\mathbf{X}^{(test)}$  and  $\mathbf{X}^{(other)}$  will be  $\rho$ .
- $\mathbf{X}^{(kinship)}$ :  $n \times k$  matrix of SNPs used to construct the kinship matrix.
- $\beta_j$ : effect size for the  $j^{th}$  SNP, simulated from a standard normal distribution for  $j = 1, \dots, (c \times 1000)$
- $Y^* = \sum_{j=1}^{c \times 1000} \beta_j \mathbf{X}_j^{(causal)}$
- $Y = Y^* + k \cdot \varepsilon$ , where the error term  $\varepsilon$  is generated from a standard normal distribution, and  $k$  is chosen such that the signal-to-noise ratio  $SNR = (Var(Y^*)/Var(\varepsilon))$  is 1

We will consider the following simulation scenarios. In each scenario we consider  $c = \{0.1, 0.5\}$  and  $\rho = \{0.1, 0.5, 0.9\}$ :

## Scenario 1

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

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

### Scenario 3

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

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

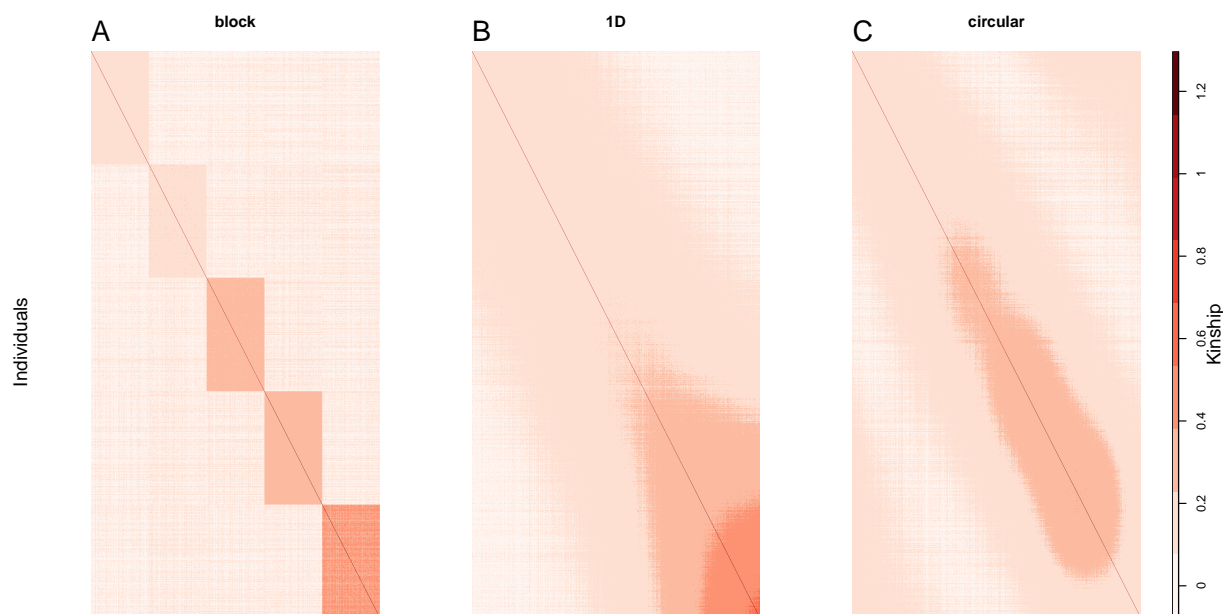


Figure 1: Empirical kinship matrices used in simulation studies.

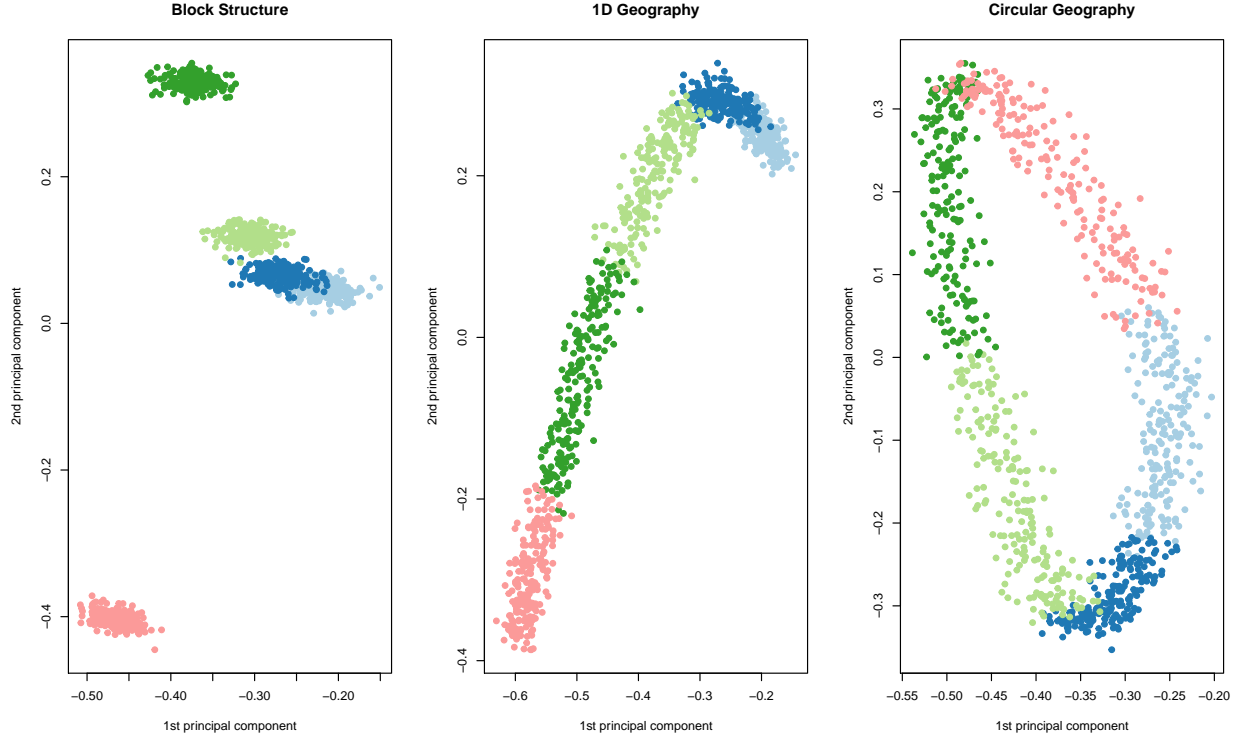


Figure 2: First two principal component scores of the kinship matrix where each color represents one of the 5 simulated subpopulations. The first panel corresponds to 5 independent subpopulations, the second corresponds to a 1 dimensional geographical structure and the third panel corresponds to a circular geography.

## 6.1 Results

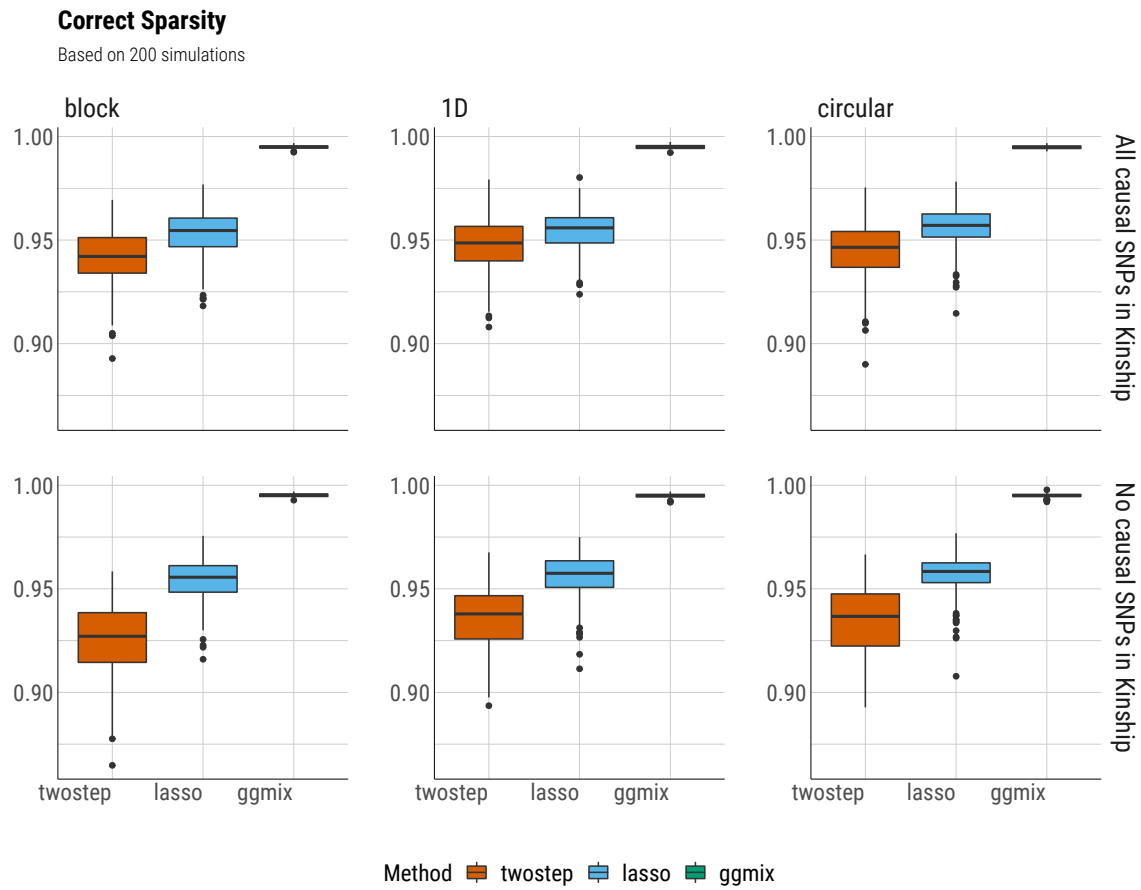


Figure 3: Boxplots of the correct sparsity from 200 simulations by kinship geography and number of causal SNPs that were included in the calculation of the kinship matrix.

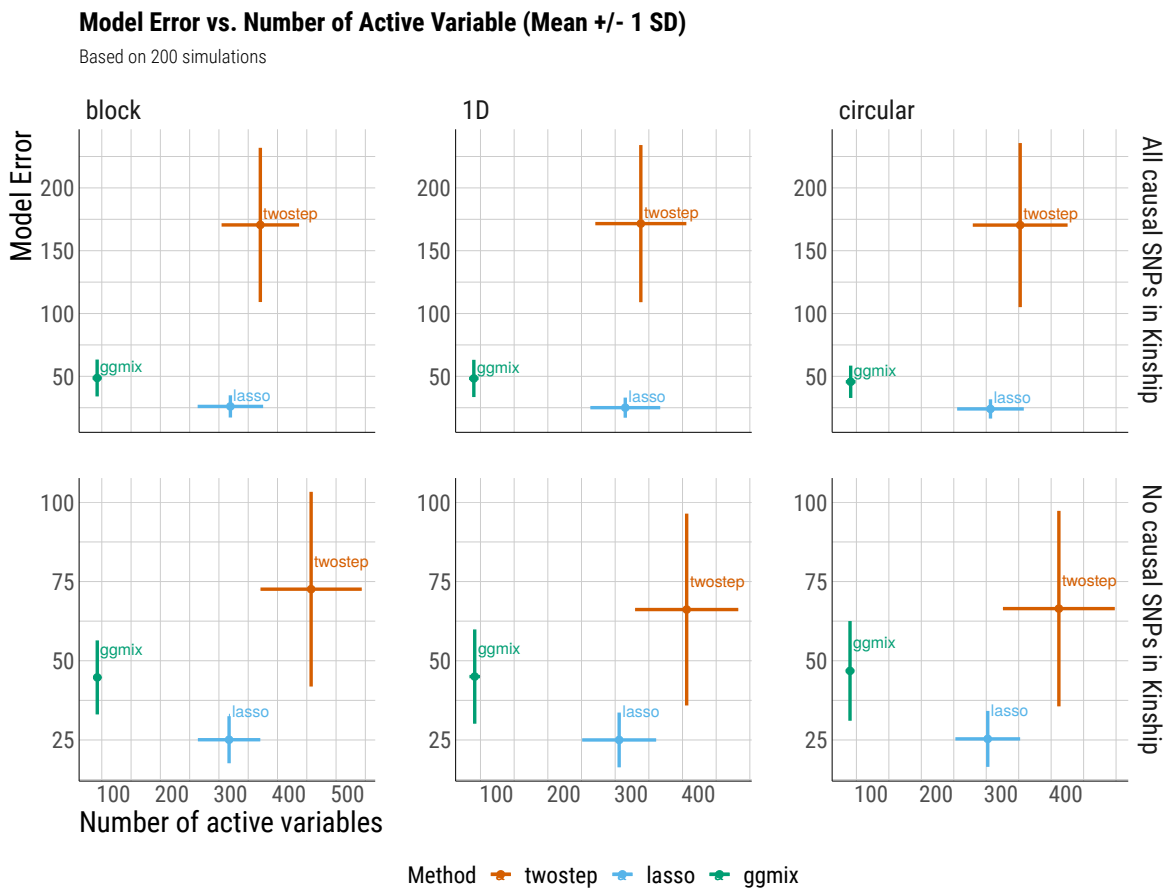


Figure 4: Model error vs number of active variables results.

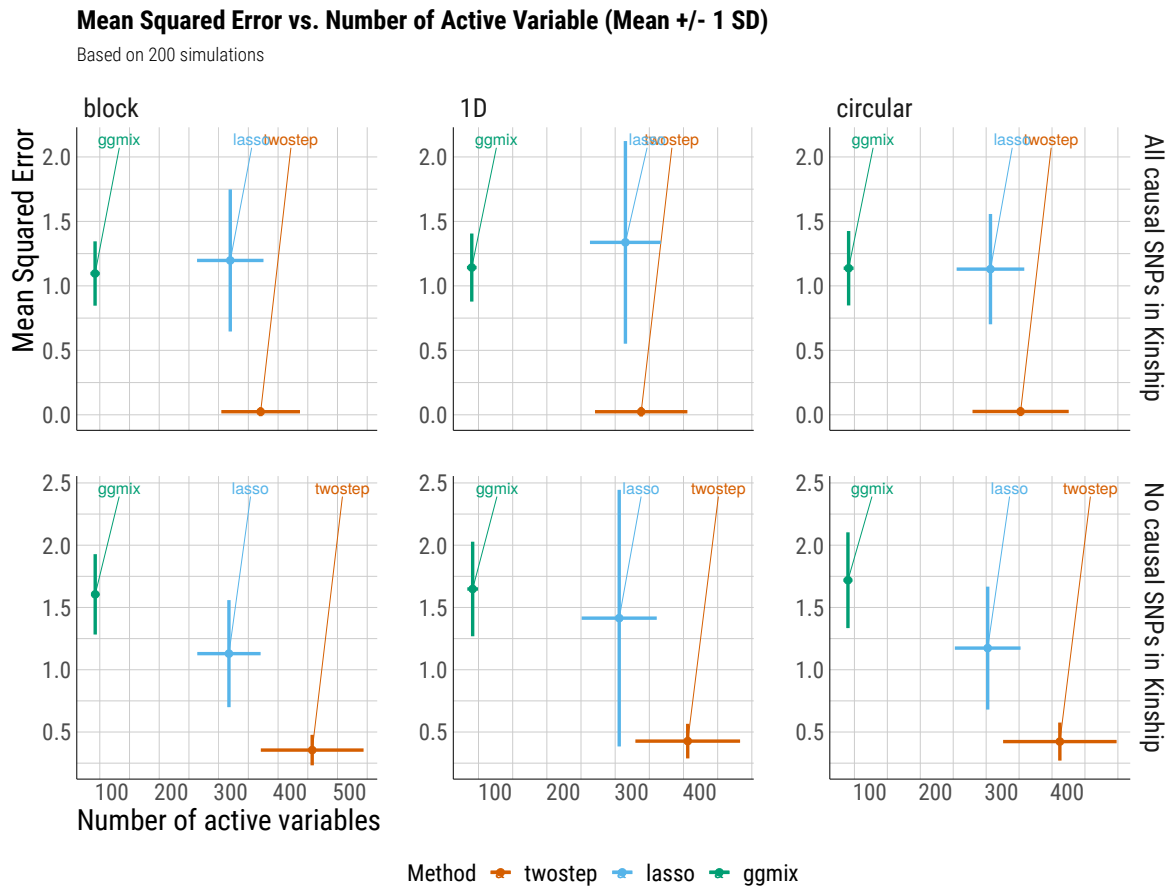


Figure 5: Mean squared error vs number of active variables results.

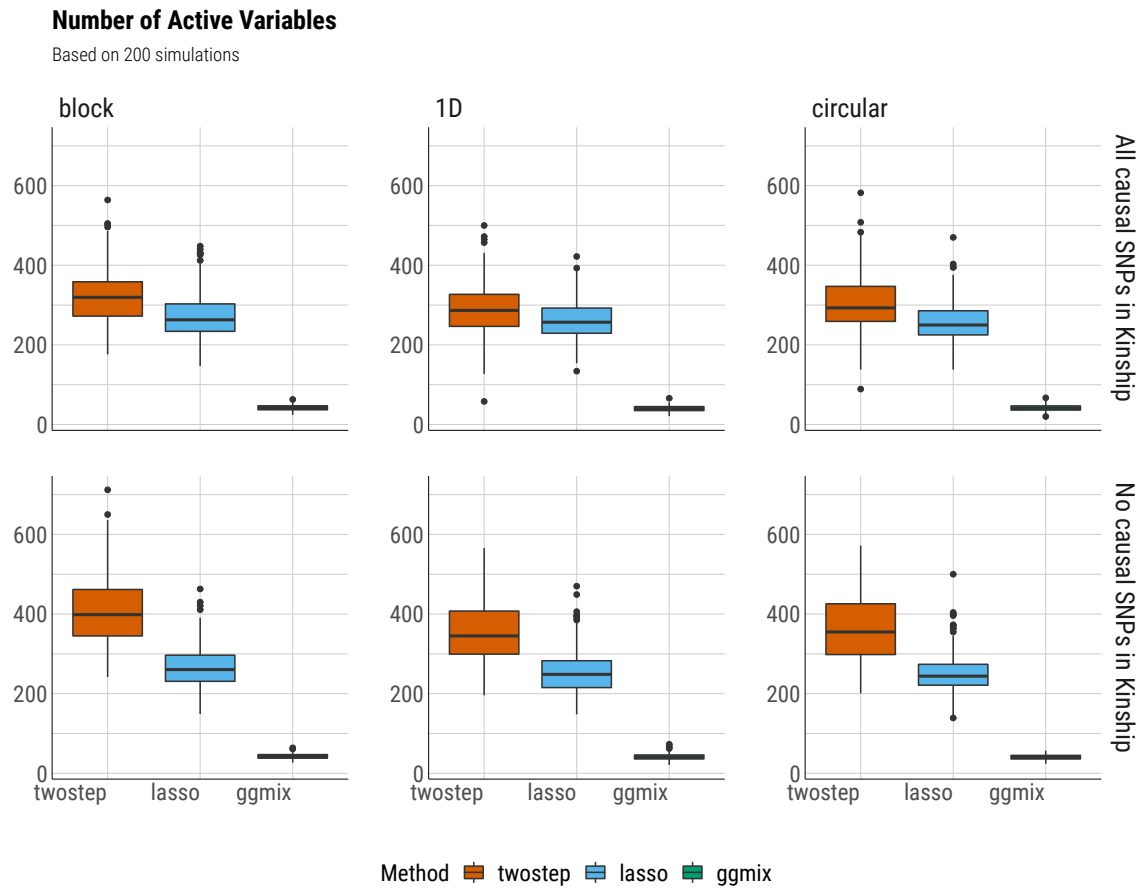


Figure 6: Boxplots of the number of active variables from 200 simulations by kinship geography and number of causal SNPs that were included in the calculation of the kinship matrix.

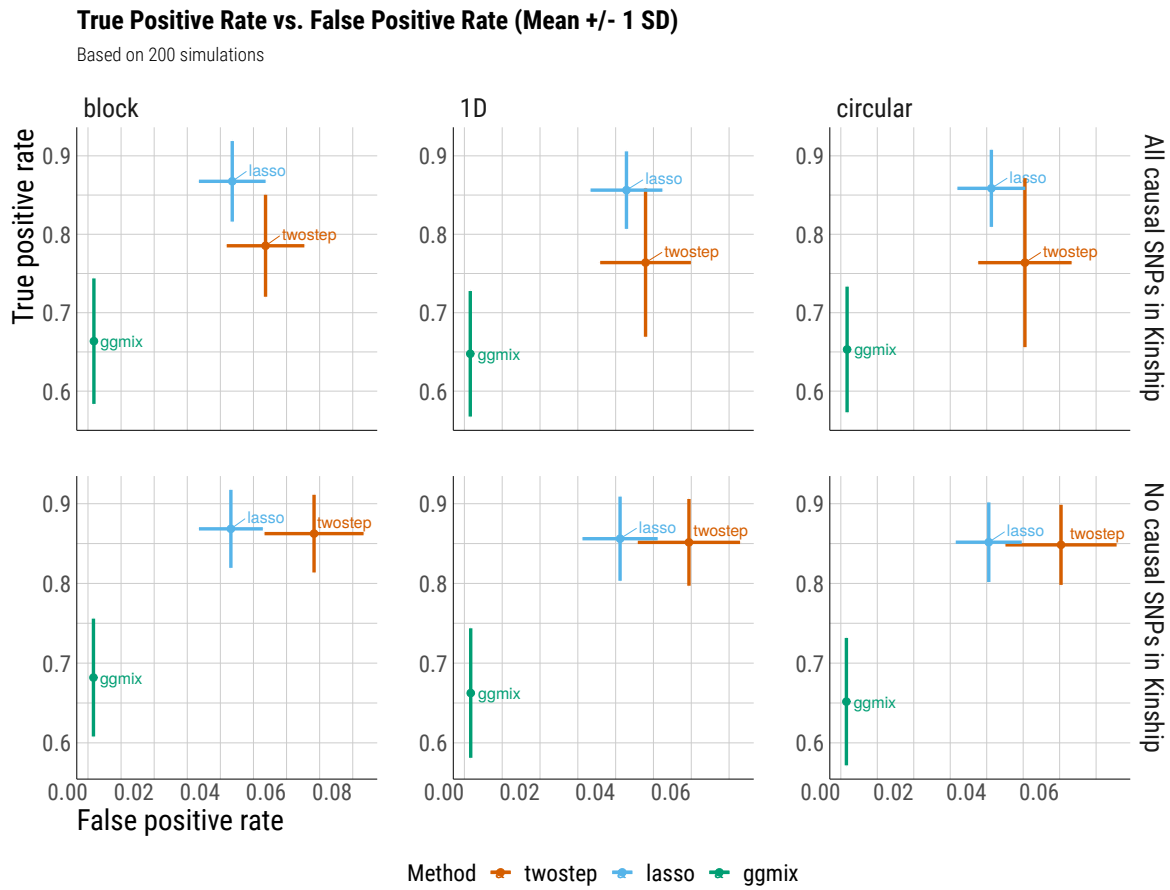


Figure 7: Means  $\pm$  1 standard deviation of true positive rate vs. false positive rate from 200 simulations by kinship geography and number of causal SNPs that were included in the calculation of the kinship matrix.



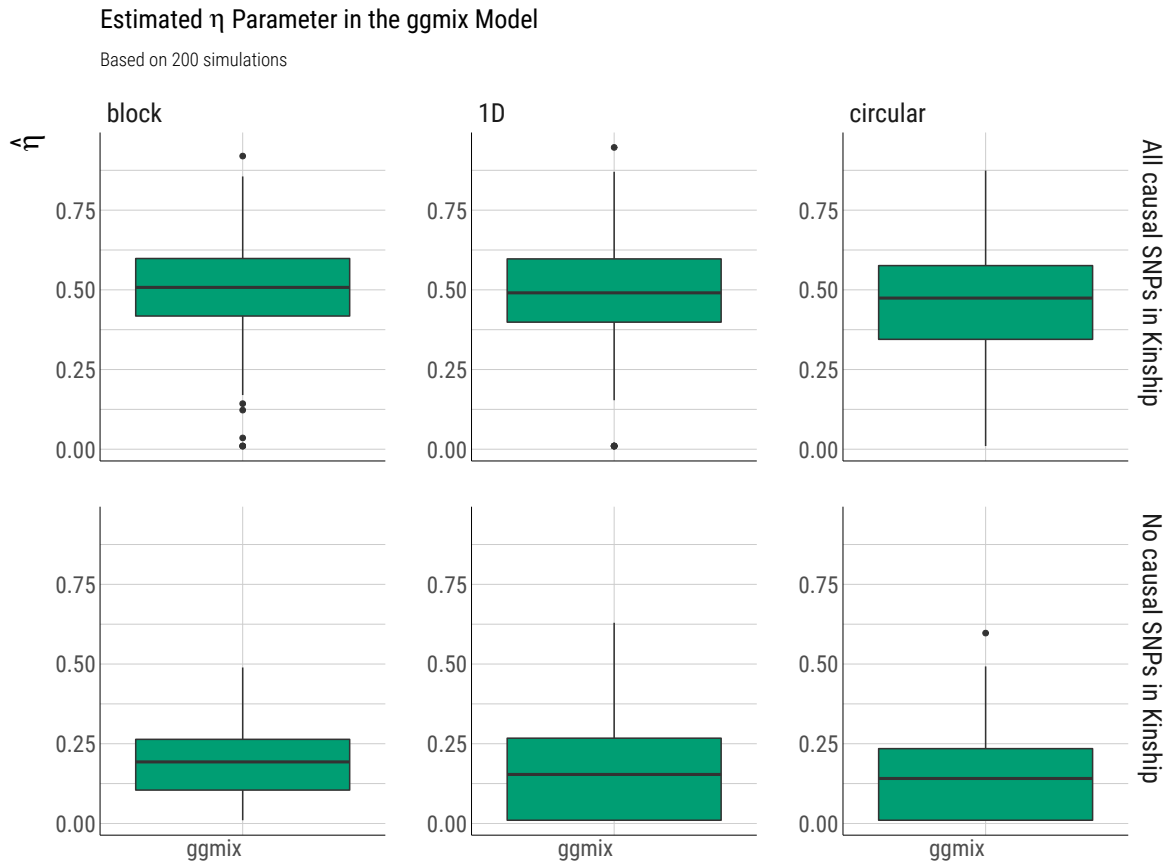


Figure 8: Boxplot of  $\hat{\eta}$  in the `ggmix` model from 200 simulations by kinship geography and number of causal SNPs that were included in the calculation of the kinship matrix.

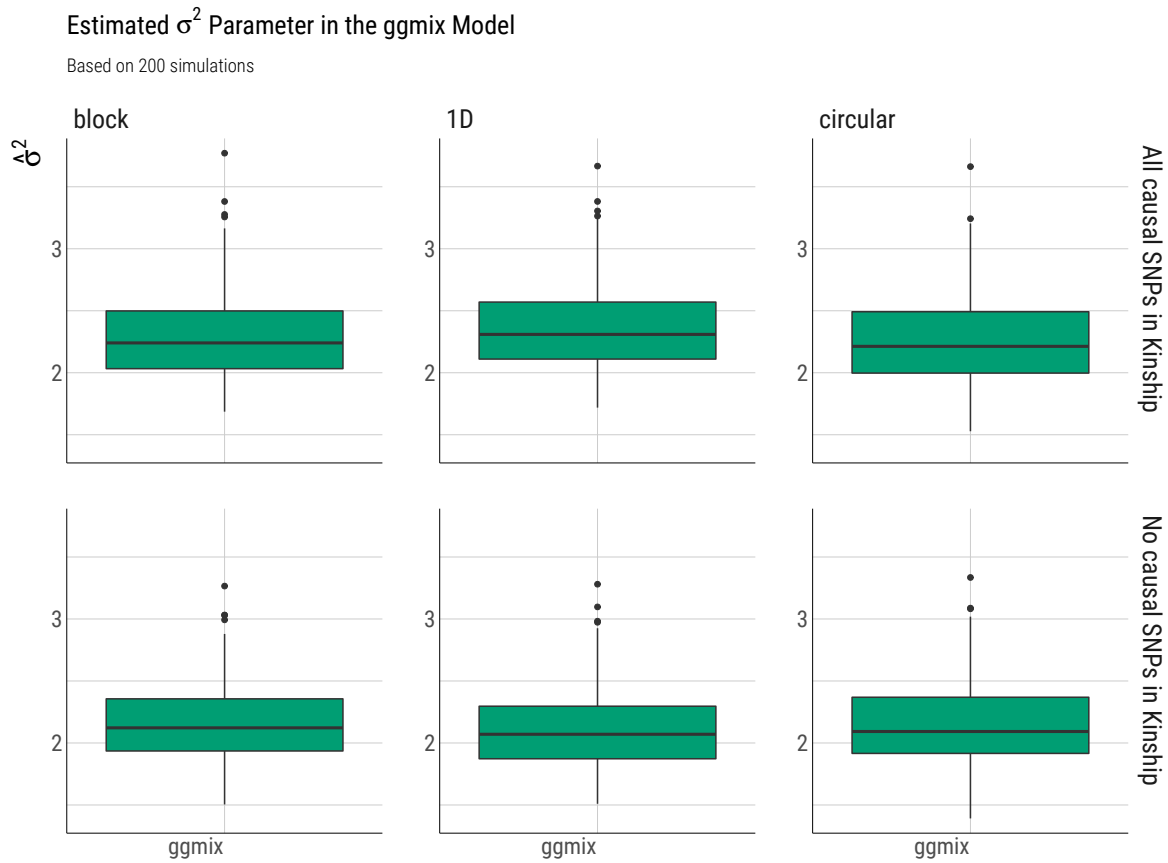


Figure 9: Boxplot of  $\hat{\eta}$  in the `ggmix` model from 200 simulations by kinship geography and number of causal SNPs that were included in the calculation of the kinship matrix.

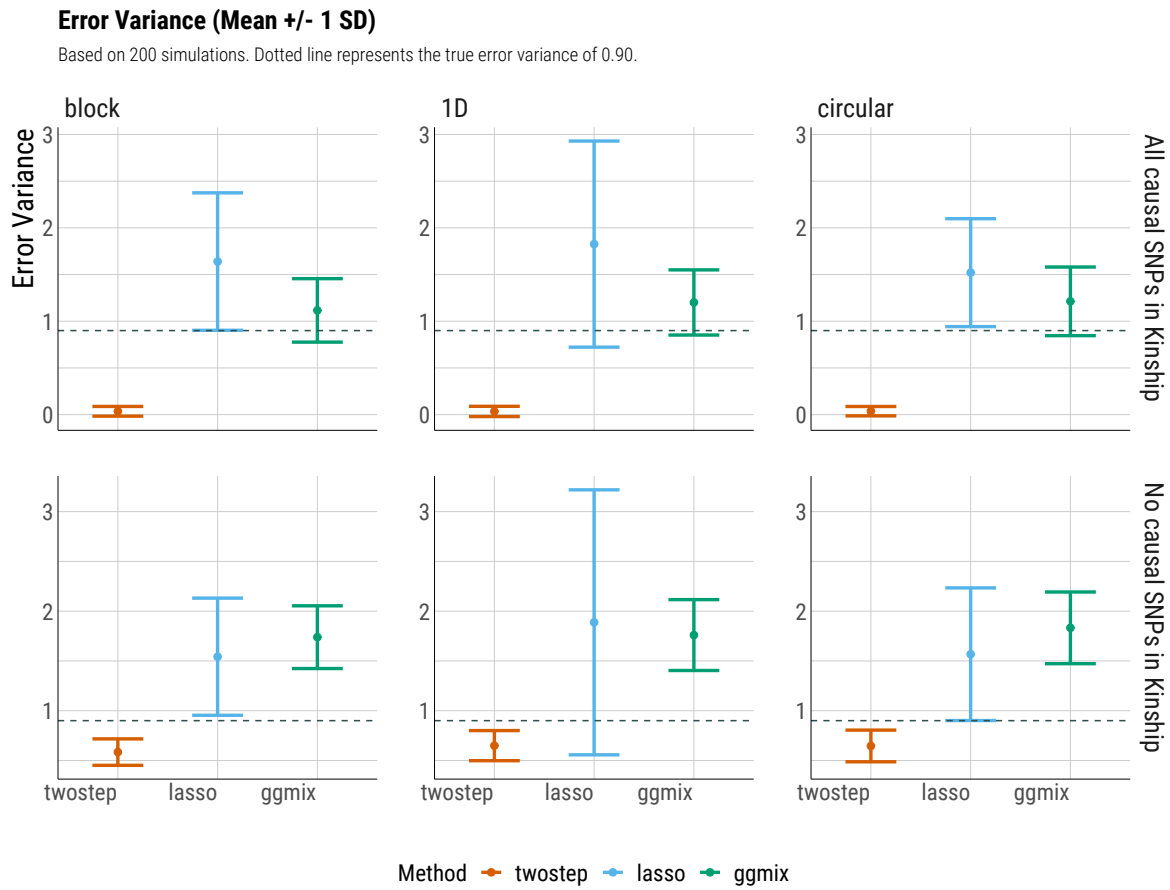


Figure 10: Means  $\pm 1$  standard deviation of error variance from 200 simulations by kinship geography and number of causal SNPs that were included in the calculation of the kinship matrix.

## 7 Computational Algorithm version 2

We use a general purpose block coordinate descent algorithm (CGD) (33) to solve (15). At each iteration, the algorithm approximates the negative log-likelihood  $f(\cdot)$  in  $Q_\lambda(\cdot)$  by a strictly convex quadratic function and then applies block coordinate descent to generate a decent direction followed by an inexact line search along this direction (33). For continuously differentiable  $f(\cdot)$  and convex and block-separable  $P(\cdot)$  (i.e.  $P(\boldsymbol{\beta}) = \sum_i P_i(\beta_i)$ ), (33) show that the solution generated by the CGD method is a stationary point of  $Q_\lambda(\cdot)$  if the coordinates are updated in a Gauss-Seidel manner i.e.  $Q_\lambda(\cdot)$  is minimized with respect to one parameter while holding all others fixed. The CGD algorithm can thus be run in parallel and therefore suited for large  $p$  settings. It has been successfully applied in fixed effects models (e.g. (34), (20)) and (25) for mixed models with an  $\ell_1$  penalty.

Following (33), the CGD algorithm is given in Algorithm 3.

---

**Algorithm 3:** Coordinate Gradient Descent Algorithm

---

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  $H^{(k)}$ :

$$H^{(k)} = \text{diag} \left[ \min \left\{ \max \left\{ \left[ \nabla^2 f(\Theta^{(k)}) \right]_{jj}, c_{\min} \right\}, c_{\max} \right\} \right]_{j=1, \dots, p+1} \quad (55)$$

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

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

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

$$d_{H^{(k)}}(\Theta_j^{(k)}) \leftarrow \arg \min_d \left\{ \nabla f(\Theta_j^{(k)})d + \frac{1}{2}d^2 H_{jj}^{(k)} + \lambda P(\Theta_j^{(k)} + d) \right\} \quad (56)$$

**end**

**if**  $\Theta_j^{(k)} \in \{\eta\}$  **then**

$$d_{H^{(k)}}(\Theta_j^{(k)}) \leftarrow -\nabla f(\Theta_j^{(k)})/H_{jj}^{(k)} \quad (57)$$

**end**

Choose a stepsize;

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

Update;

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

**end**

Update;

$$\hat{\sigma}^{2(k+1)} \leftarrow \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{([\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\hat{\boldsymbol{\beta}}^{(k+1)}]_i)^2}{1 + \hat{\eta}^{(k+1)}(\Lambda_i - 1)} \quad (58)$$

$k \leftarrow k + 1$

**until** convergence criterion is satisfied;

We note that conditional on  $\widehat{\boldsymbol{\beta}}$  and  $\widehat{\eta}$ , there exists an analytic solution for  $\widehat{\sigma^2}$ :

$$\begin{aligned}\frac{\partial}{\partial \sigma^2} f(\boldsymbol{\Theta}) &= \frac{N_T}{2\sigma^2} - \frac{1}{2\sigma^4} \sum_{i=1}^{N_T} \frac{([\widetilde{\mathbf{Y}} - \widetilde{\mathbf{X}}\boldsymbol{\beta}]_i)^2}{1 + \eta(\Lambda_i - 1)} = 0 \\ \widehat{\sigma^2} &= \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{([\widetilde{\mathbf{Y}} - \widetilde{\mathbf{X}}\widehat{\boldsymbol{\beta}}]_i)^2}{1 + \widehat{\eta}(\Lambda_i - 1)}\end{aligned}\quad (59)$$

The Armijo rule is defined as follows (33):

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

$$Q_\lambda(\Theta_j^{(k)} + \alpha^{(k)} d^{(k)}) \leq Q_\lambda(\Theta_j^{(k)}) + \alpha^{(k)} \varrho \Delta^{(k)} \quad (60)$$

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

$$\Delta^{(k)} := \nabla f(\Theta_j^{(k)}) d^{(k)} + \gamma (d^{(k)})^2 H_{jj}^{(k)} + \lambda P(\Theta_j^{(k)} + d^{(k)}) - \lambda P(\Theta^{(k)}) \quad (61)$$

Common choices for the constants are  $\delta = 0.1$ ,  $\varrho = 0.001$ ,  $\gamma = 0$ ,  $\alpha_{init}^{(k)} = 1$  for all  $k$  (25).

Below we detail the specifics of Algorithm 3 for different penalty functions  $P(\boldsymbol{\beta})$ .

## 7.1 $\ell_1$ penalty

The objective function is given by

$$Q_\lambda(\boldsymbol{\Theta}) = f(\boldsymbol{\Theta}) + \lambda |\boldsymbol{\beta}| \quad (62)$$

### 7.1.1 Descent Direction

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

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

$$d_H(\Theta_j) = \arg \min_d G(d) \quad (63)$$

where

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

Since  $G(d)$  is not differentiable at  $-\Theta_j$ , we calculate the subdifferential  $\partial G(d)$  and search for  $d$  with  $0 \in \partial G(d)$ :

$$\partial G(d) = \nabla f(\Theta_j) + dH_{jj} + \lambda u \quad (64)$$

where

$$u = \begin{cases} 1 & \text{if } d > -\Theta_j \\ -1 & \text{if } d < -\Theta_j \\ [-1, 1] & \text{if } d = -\Theta_j \end{cases} \quad (65)$$

We consider each of the three cases in (64) below

1.  $d > -\Theta_j$

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

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

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

where  $\text{mid} \{a, b, c\}$  denotes the median (mid-point) of  $a, b, c$  (33).

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

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

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

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

We see all three cases lead to the same solution for (63). Therefore the descent direction for  $\Theta_j^{(k)} \in \{\beta_1, \dots, \beta_p\}$  for the  $\ell_1$  penalty is given by

$$d = \text{mid} \left\{ \frac{-(\nabla f(\beta_j) - \lambda)}{H_{jj}}, -\beta_j, \frac{-(\nabla f(\beta_j) + \lambda)}{H_{jj}} \right\} \quad (66)$$

### 7.1.2 Solution for the $\beta$ parameter

If the Hessian  $\nabla^2 f(\Theta^{(k)}) > 0$  then  $H^{(k)}$  defined in (55) is equal to  $\nabla^2 f(\Theta^{(k)})$ . Using  $\alpha_{init} = 1$ , the largest element of  $\left\{ \alpha_{init}^{(k)} \delta^r \right\}_{r=0,1,2,\dots}$  satisfying the Armijo Rule inequality is reached for  $\alpha^{(k)} = \alpha_{init}^{(k)} \delta^0 = 1$ . The Armijo rule update for the  $\beta$  parameter is then given by

$$\beta_j^{(k+1)} \leftarrow \beta_j^{(k)} + d^{(k)}, \quad j = 1, \dots, p \quad (67)$$

Substituting the descent direction given by (66) into (67) we get

$$\beta_j^{(k+1)} = \text{mid} \left\{ \beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) - \lambda)}{H_{jj}}, 0, \beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) + \lambda)}{H_{jj}} \right\} \quad (68)$$

We can further simplify this expression. Let

$$w_i := \frac{1}{\sigma^2 (1 + \eta(\Lambda_i - 1))} \quad (69)$$

Re-write the part depending on  $\beta$  of the negative log-likelihood in (13) as

$$g(\beta^{(k)}) = \frac{1}{2} \sum_{i=1}^{N_T} w_i \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} - \tilde{X}_{ij} \beta_j^{(k)} \right)^2 \quad (70)$$

The gradient and Hessian are given by

$$\nabla f(\beta_j^{(k)}) := \frac{\partial}{\partial \beta_j^{(k)}} g(\beta^{(k)}) = - \sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} - \tilde{X}_{ij} \beta_j^{(k)} \right) \quad (71)$$

$$H_{jj} := \frac{\partial^2}{\partial \beta_j^{(k)2}} g(\beta^{(k)}) = \sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2 \quad (72)$$

Substituting (71) and (72) into  $\beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)})) - \lambda}{H_{jj}}$

$$\begin{aligned} & \beta_j^{(k)} + \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} - \tilde{X}_{ij} \beta_j^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \\ &= \beta_j^{(k)} + \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} - \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2 \beta_j^{(k)}}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \\ &= \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \end{aligned} \quad (73)$$

Similarly, substituting (71) and (72) in  $\beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)})) + \lambda}{H_{jj}}$  we get

$$\frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) - \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \quad (74)$$

Finally, substituting (73) and (74) into (68) we get

$$\begin{aligned}\beta_j^{(k+1)} &= \text{mid} \left\{ \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) - \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2}, 0, \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \right\} \\ &= \frac{\mathcal{S}_\lambda \left( \sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) \right)}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2}\end{aligned}\quad (75)$$

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

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

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

We note that the parameter update for  $\beta_j$  given by (75) takes the same form as the weighted updates of the `glmnet` algorithm (20) (Section 2.4, equation (10)) with  $\alpha = 1$ .

## References

- [1] Teri A Manolio, Francis S Collins, Nancy J Cox, David B Goldstein, Lucia A Hindorff, David J Hunter, Mark I McCarthy, Erin M Ramos, Lon R Cardon, Aravinda Chakravarti, et al. Finding the missing heritability of complex diseases. *Nature*, 461(7265):747, 2009. 2
- [2] Jian Yang, Beben Benyamin, Brian P McEvoy, Scott Gordon, Anjali K Henders, Dale R Nyholt, Pamela A Madden, Andrew C Heath, Nicholas G Martin, Grant W Montgomery, et al. Common snps explain a large proportion of the heritability for human height. *Nature genetics*, 42(7):565, 2010. 2
- [3] William Astle, David J Balding, et al. Population structure and cryptic relatedness in genetic association studies. *Statistical Science*, 24(4):451–471, 2009. 3
- [4] Minsun Song, Wei Hao, and John D Storey. Testing for genetic associations in arbitrarily structured populations. *Nature genetics*, 47(5):550–554, 2015. 3, 4, 5
- [5] Jonathan Marchini, Lon R Cardon, Michael S Phillips, and Peter Donnelly. The effects of human population structure on large genetic association studies. *Nature genetics*, 36(5):512, 2004. 3
- [6] Clive J Hoggart, John C Whittaker, Maria De Iorio, and David J Balding. Simultaneous analysis of all snps in genome-wide and re-sequencing association studies. *PLoS genetics*, 4(7):e1000130, 2008. 3
- [7] Dong Wang, Kent M Eskridge, and Jose Crossa. Identifying qtls and epistasis in structured plant populations using adaptive mixed lasso. *Journal of agricultural, biological, and environmental statistics*, 16(2):170–184, 2011. 3
- [8] Barbara Rakitsch, Christoph Lippert, Oliver Stegle, and Karsten Borgwardt. A lasso

- multi-marker mixed model for association mapping with population structure correction. *Bioinformatics*, 29(2):206–214, 2013. [3](#), [4](#), [5](#)
- [9] Hyun Min Kang, Jae Hoon Sul, Noah A Zaitlen, Sit-ye Kong, Nelson B Freimer, Chiara Sabatti, Eleazar Eskin, et al. Variance component model to account for sample structure in genome-wide association studies. *Nature genetics*, 42(4):348, 2010. [3](#)
- [10] Vincent Segura, Bjarni J Vilhjálmsson, Alexander Platt, Arthur Korte, Ümit Seren, Quan Long, and Magnus Nordborg. An efficient multi-locus mixed-model approach for genome-wide association studies in structured populations. *Nature genetics*, 44(7):825–830, 2012. [4](#), [5](#)
- [11] Jonas R Klasen, Elke Barbez, Lukas Meier, Nicolai Meinshausen, Peter Bühlmann, Maarten Koornneef, Wolfgang Busch, and Korbinian Schneeberger. A multi-marker association method for genome-wide association studies without the need for population structure correction. *Nature communications*, 7:13299, 2016. [4](#), [5](#)
- [12] Xiang Zhou and Matthew Stephens. Efficient multivariate linear mixed model algorithms for genome-wide association studies. *Nature methods*, 11(4):407–409, 2014. [4](#)
- [13] Arthur Korte, Bjarni J Vilhjálmsson, Vincent Segura, Alexander Platt, Quan Long, and Magnus Nordborg. A mixed-model approach for genome-wide association studies of correlated traits in structured populations. *Nature genetics*, 44(9):1066–1071, 2012. [4](#)
- [14] Jakris Eu-Ahsunthornwattana, E Nancy Miller, Michaela Fakiola, Selma MB Jeronimo, Jenefer M Blackwell, Heather J Cordell, Wellcome Trust Case Control Consortium 2, et al. Comparison of methods to account for relatedness in genome-wide association studies with family-based data. *PLoS Genet*, 10(7):e1004445, 2014. [4](#)
- [15] W. Hao, M. Song, and J. D. Storey. Probabilistic models of genetic variation in struc-

- tured populations applied to global human studies. *ArXiv e-prints*, December 2013. 5
- [16] Nicolai Meinshausen. Hierarchical testing of variable importance. *Biometrika*, pages 265–278, 2008. 5
- [17] Matti Pirinen, Peter Donnelly, Chris CA Spencer, et al. Efficient computation with a linear mixed model on large-scale data sets with applications to genetic studies. *The Annals of Applied Statistics*, 7(1):369–390, 2013. 6, 7
- [18] Christoph Lippert, Jennifer Listgarten, Ying Liu, Carl M Kadie, Robert I Davidson, and David Heckerman. Fast linear mixed models for genome-wide association studies. *Nature methods*, 8(10):833–835, 2011. 7, 17
- [19] Jan De Leeuw. Block-relaxation algorithms in statistics. In *Information systems and data analysis*, pages 308–324. Springer, 1994. 10
- [20] Jerome Friedman, Trevor Hastie, and Rob Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1):1, 2010. 12, 14, 36, 43
- [21] Richard H Byrd, Peihuang Lu, Jorge Nocedal, and Ciyu Zhu. A limited memory algorithm for bound constrained optimization. *SIAM Journal on Scientific Computing*, 16(5):1190–1208, 1995. 12
- [22] Jon Wakefield. *Bayesian and frequentist regression methods*. Springer Science & Business Media, 2013. 15
- [23] Hui Zou, Trevor Hastie, Robert Tibshirani, et al. On the “degrees of freedom” of the lasso. *The Annals of Statistics*, 35(5):2173–2192, 2007. 16
- [24] Howard D Bondell, Arun Krishna, and Sujit K Ghosh. Joint variable selection for fixed

- and random effects in linear mixed-effects models. *Biometrics*, 66(4):1069–1077, 2010. 16
- [25] Jürg Schelldorfer, Peter Bühlmann, GEER DE, and SARA VAN. Estimation for high-dimensional linear mixed-effects models using l1-penalization. *Scandinavian Journal of Statistics*, 38(2):197–214, 2011. 16, 36, 38
- [26] Joseph G Ibrahim, Hongtu Zhu, Ramon I Garcia, and Ruixin Guo. Fixed and random effects selection in mixed effects models. *Biometrics*, 67(2):495–503, 2011. 16
- [27] Hansheng Wang, Bo Li, and Chenlei Leng. Shrinkage tuning parameter selection with a diverging number of parameters. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(3):671–683, 2009. 16
- [28] Daniel P Berrar, Werner Dubitzky, Martin Granzow, et al. *A practical approach to microarray data analysis*. Springer, 2003. 17
- [29] Gene H Golub and Charles F Van Loan. *Matrix computations*, volume 3. JHU Press, 2012. 18
- [30] K. Lange, D. Hunter, and I. Yang. Optimization transfer using surrogate objective functions (with discussion). *Journal of Computational and Graphical Statistics*, 9:1–20, 2000. 22
- [31] D.R. Hunter and K. Lange. A tutorial on MM algorithms. *The American Statistician*, 58(1):30–37, 2004. 22
- [32] T. Wu and K. Lange. The MM alternative to EM. *Statistical Science*, 4:492–505, 2010. 22
- [33] Paul Tseng and Sangwoon Yun. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*, 117(1):387–423, 2009. 36, 37, 38, 40

- [34] Lukas Meier, Sara Van De Geer, and Peter Bühlmann. The group lasso for logistic regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(1):53–71, 2008. [36](#)



## A Algorithm Details

In this section we provide more specific details about the algorithms used to solve th objective function.

### A.1 title

## B title