

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

Sahir R Bhatnagar^{1,2}, Karim Oualkacha³, Yi Yang⁴, Marie Forest², and
Celia MT Greenwood^{1,2,5}

¹Department of Epidemiology, Biostatistics and Occupational Health, McGill
University

²Lady Davis Institute, Jewish General Hospital, Montréal, QC

³Département de Mathématiques, Université de Québec À Montréal

⁴Department of Mathematics and Statistics, McGill University

⁵Departments of Oncology and Human Genetics, McGill University

July 4, 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 or single nucleotide polymorphism (SNP) 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 issue to overcome 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 the first problem, multivariable regression methods have been proposed which simultaneously fit many SNPs in a single model (6, 7). Indeed, the power to detect an association for a given SNP may be increased when other causal SNPs have been accounted for. Conversely, a stronger signal from a causal SNP may weaken false signals when modeled jointly (6).

Confounding by population structure has also received significant attention in the literature (8, 9, 10, 11). There are two main approaches to account for the relatedness between subjects: 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 SNP genotypes as additional covariates in the model (12). The LMM uses an estimated covariance matrix from the individuals' genotypes and includes this information in the form of a random effect (3).

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

In light of this, there has been recent interest in penalized linear mixed models which place a constraint on the magnitude of the effect sizes while controlling for confounding influences such as population structure. For example, the LMM-lasso (15) places a Laplace prior on all main effects while the adaptive mixed lasso (16) uses the L_1 penalty (17) with adaptively chosen weights (18) to allow for differential shrinkage amongst the variables in the model. Another method applied a combination of both the lasso and group lasso penalties in order to select variants within a gene most associated with the response (19). One potential issue with these methods is that they are performed in two steps. First, the variance components are estimated once from a LMM with a single random effect that uses the estimated covariance matrix from the individuals' genotypes to account for the relatedness but assumes no SNP effects. In the second step, these are treated as known quantities by regressing the SNPs on the residuals from the first step, effectively treating the observations as independent. This approach has both computational and practical advantages since existing penalized regression software such as `glmnet` (20) and `gglasso` (21), which assume independent observations, can be applied directly to the residuals. However, recent work has shown that there can be a loss in power if a causal variant is included in the calculation of the covariance matrix as its effect will have been removed in the first step (13, 22). Another issue with the aforementioned methods is that they first require computing the covariance matrix with a computation time of $\mathcal{O}(n^2k)$ followed by a spectral decomposition of this matrix in $\mathcal{O}(n^3)$ time where k is the number of SNP genotypes used to construct the covariance matrix. These methods become prohibitive to use for large cohorts such as the UK Biobank (23) which have collected genetic information on half a million individuals. There is thus a need to develop newer methodologies that reflect the increasing size and genetic heterogeneity of the large cohort studies being assembled today.

In this paper we develop a general penalized LMM framework called `ggmix` 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 (17), elastic net (24) and group lasso (25). `ggmix` also readily handles prior annotation information in the form of a penalty factor, which can be useful for example when dealing with rare variants. We develop a blockwise coordinate descent algorithm which is highly scalable and has theoretical guarantees of convergence to a stationary point. When the matrix of genotypes used to construct the covariance matrix is low rank, there are additional computational speedups that can be implemented. While this has been developed for the univariate case (8), to our knowledge, this hasn't been explored in the multivariable case. The LMM-lasso paper mentions that this is possible but does not provide further details on how this can be implemented in a penalized mixed model framework. In the sequel, we develop a low rank version of the blockwise coordinate descent algorithm which reduces the time complexity from $\mathcal{O}(n^2k)$ to $\mathcal{O}(nk^2)$. All of our algorithms are implemented in the `ggmix` R package hosted on GitHub with extensive documentation (<http://sahirbhatnagar.com/ggmix/>). We provide a brief demonstration of the `ggmix` package in Appendix B.

The rest of the paper is organized as follows. Section 2 describes the `ggmix` model. Section 3 contains the optimization procedure and the algorithm used to fit the `ggmix` model. In Section 4, we compare the performance of our proposed approach and demonstrate the scenarios where it can be advantageous to use over existing methods through simulation studies. Section 5 contains some real data examples and Section 6 discusses some limitations and future directions.

2 Penalized Linear 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 or phenotypes, \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. Denote the stacked vectors $\mathbf{Y} = (\mathbf{y}_i, \dots, \mathbf{y}_N)^T \in \mathbb{R}^{N_T \times 1}$, $\mathbf{b} = (\mathbf{b}_i, \dots, \mathbf{b}_N)^T \in \mathbb{R}^{N_T \times 1}$, $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}_i, \dots, \boldsymbol{\varepsilon}_N)^T \in \mathbb{R}^{N_T \times 1}$, and the stacked matrix $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_N)^T \in \mathbb{R}^{N_T \times (p+1)}$. Furthermore, let $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^T \in \mathbb{R}^{(p+1) \times 1}$ be a vector of fixed effects regression coefficients corresponding to \mathbf{X} . We consider the following linear mixed model with a single random effect (28):

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

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 covariance or kinship matrix, $\mathbf{I}_{N_T \times N_T}$ is the identity matrix and parameters σ^2 and $\eta \in [0, 1]$ determine how the variance is divided between \mathbf{b} and $\boldsymbol{\varepsilon}$. Furthermore, η is also the narrow-sense heritability (h^2), defined as the proportion of phenotypic variance attributable to the additive genetic factors (1). The joint density of \mathbf{Y} is multivariate normal:

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

The LMM-Lasso method (15) considers an alternative parameterization given by:

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

where $\delta = \sigma_e^2/\sigma_g^2$, σ_g^2 is the genetic variance and σ_e^2 is the residual variance. We instead consider the parameterization in (3) since maximization is easier over the compact set $\eta \in [0, 1]$ than over the unbounded interval $\delta \in [0, \infty)$ (28). We define the complete parameter

vector as $\Theta := (\beta, \eta, \sigma^2)$. The negative log-likelihood for (3) is given by

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

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 Φ , 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 Λ_i . \mathbf{V} can then be further simplified (28)

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

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^{th} element of $\tilde{\mathbf{Y}}$, \tilde{X}_{ij} is the i, j^{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 η and σ^2 , respectively. In light of our goals to select variables associated with the response in high-dimensional data, we consider placing a constraint on the magnitude of the regression coefficients. This can be achieved by adding a penalty term to the likelihood

function (13). The penalty term is a necessary constraint because in our applications, the sample size is much smaller than the number of predictors. We define the following objective function:

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

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 λ , 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 η or σ^2 . An estimate of the regression parameters $\hat{\Theta}_\lambda$ is obtained by

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

This is the general set-up for our model. In Section 3 we provide more specific details on how we solve (15).

3 Computational Algorithm

We use a general purpose block coordinate gradient descent algorithm (CGD) (29) to solve (15). At each iteration, we cycle through the coordinates and minimize the objective function with respect to one coordinate only. For continuously differentiable $f(\cdot)$ and convex and block-separable $P(\cdot)$ (i.e. $P(\beta) = \sum_i P_i(\beta_i)$), Tseng and Yun (29) 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 has been successfully applied in fixed effects models (e.g. (30), (20)) and linear mixed models with an ℓ_1 penalty (31). In the next section we provide some brief details about Algorithm 1. A more thorough treatment of the algorithm 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 ϵ ;
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

3.1 Updates for the β parameter

Recall that the part of the objective function that depends on β 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

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

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

$$\beta_j^{(k+1)} \leftarrow \frac{\mathcal{S}_{\lambda} \left(\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left(\tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_{\ell}^{(k)} \right) \right)}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \quad (18)$$

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 provide the full derivation in Appendix [A.1.2](#).

3.2 Updates for the η paramter

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

$$\eta^{(k+1)} \leftarrow \arg \min_{\eta} \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^{2(k)}} \sum_{i=1}^{N_T} \frac{\left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j^{(k+1)} \right)^2}{1 + \eta(\Lambda_i - 1)} \quad (19)$$

We use a bound constrained optimization algorithm ([32](#)) implemented in the `optim` function in R and set the lower and upper bounds to be 0.01 and 0.99, respectively.

3.3 Updates for the σ^2 parameter

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

$$\sigma^{2(k+1)} \leftarrow \arg \min_{\sigma^2} \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^{N_T} \frac{\left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j^{(k+1)} \right)^2}{1 + \eta(\Lambda_i - 1)} \quad (20)$$

There exists an analytic solution for ([20](#)) given by:

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

3.4 Regularization path

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

$$Q_\lambda(\Theta) = \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2} \sum_{i=1}^{N_T} w_i \left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2 + \lambda \sum_{j=1}^p v_j |\beta_j| \quad (22)$$

The Karush-Kuhn-Tucker (KKT) optimality conditions for (22) are given by:

$$\begin{aligned} \frac{\partial}{\partial \beta_1, \dots, \beta_p} Q_\lambda(\Theta) &= \mathbf{0}_p \\ \frac{\partial}{\partial \beta_0} Q_\lambda(\Theta) &= 0 \\ \frac{\partial}{\partial \eta} Q_\lambda(\Theta) &= 0 \\ \frac{\partial}{\partial \sigma^2} Q_\lambda(\Theta) &= 0 \end{aligned} \quad (23)$$

The equations in (23) are equivalent to

$$\begin{aligned} \sum_{i=1}^{N_T} w_i \tilde{X}_{i1} \left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right) &= 0 \\ \frac{1}{v_j} \sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right) &= \lambda \gamma_j, \\ \gamma_j &\in \begin{cases} \text{sign}(\hat{\beta}_j) & \text{if } \hat{\beta}_j \neq 0 \\ [-1, 1] & \text{if } \hat{\beta}_j = 0 \end{cases}, \quad \text{for } j = 1, \dots, p \\ \frac{1}{2} \sum_{i=1}^{N_T} \frac{\Lambda_i - 1}{1 + \eta(\Lambda_i - 1)} \left(1 - \frac{\left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2}{\sigma^2 (1 + \eta(\Lambda_i - 1))} \right) &= 0 \\ \sigma^2 - \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{\left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2}{1 + \eta(\Lambda_i - 1)} &= 0 \end{aligned} \quad (24)$$

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 ℓ_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 (24) for some $\boldsymbol{\gamma}$. 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 (24) are equivalent to

$$\begin{aligned} \frac{1}{v_j} \sum_{i=1}^{N_T} \left| w_i \tilde{X}_{ij} \left(\tilde{Y}_i - \tilde{X}_{i1} \beta_0 \right) \right| &\leq \lambda, \quad \forall j = 1, \dots, p \\ \beta_0 &= \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{i1} \tilde{Y}_i}{\sum_{i=1}^{N_T} w_i \tilde{X}_{i1}^2} \\ \frac{1}{2} \sum_{i=1}^{N_T} \frac{\Lambda_i - 1}{1 + \eta(\Lambda_i - 1)} \left(1 - \frac{\left(\tilde{Y}_i - \tilde{X}_{i1} \beta_0 \right)^2}{\sigma^2 (1 + \eta(\Lambda_i - 1))} \right) &= 0 \\ \sigma^2 &= \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{\left(\tilde{Y}_i - \tilde{X}_{i1} \beta_0 \right)^2}{1 + \eta(\Lambda_i - 1)} \end{aligned} \quad (25)$$

We can solve the KKT system of equations in (25) (with a numerical solution for η) 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 λ such that the entire vector $(\hat{\beta}_1, \dots, \hat{\beta}_p)$ is 0:

$$\lambda_{max} = \max_j \left\{ \left| \frac{1}{v_j} \sum_{i=1}^{N_T} \hat{w}_i \tilde{X}_{ij} \left(\tilde{Y}_i - \tilde{X}_{i1} \hat{\beta}_0 \right) \right| \right\}, \quad j = 1, \dots, p \quad (26)$$

Following Friedman et al. (20), we choose $\tau \lambda_{max}$ to be the smallest value of tuning parameters λ_{min} , and construct a sequence of K values decreasing from λ_{max} to λ_{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 $\hat{\boldsymbol{\Theta}}$ for λ_k is used as the

initial value $\Theta^{(0)}$ for λ_{k+1} . This strategy leads to computational speedups and has been implemented in the `ggmix` R package.

3.6 Prediction of the random effects

We use an empirical Bayes approach (e.g. (33)) to predict the random effects \mathbf{b} . Let the maximum a posteriori (MAP) estimate be defined as

$$\hat{\mathbf{b}} = \arg \max_{\mathbf{b}} f(\mathbf{b}|\mathbf{Y}, \boldsymbol{\beta}, \eta, \sigma^2) \quad (27)$$

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

$$\begin{aligned} f(\mathbf{b}|\mathbf{Y}, \boldsymbol{\beta}, \eta, \sigma^2) &= \frac{f(\mathbf{Y}|\mathbf{b}, \boldsymbol{\beta}, \eta, \sigma^2)\pi(\mathbf{b}|\eta, \sigma^2)}{f(\mathbf{Y}|\boldsymbol{\beta}, \eta, \sigma^2)} \\ &\propto f(\mathbf{Y}|\mathbf{b}, \boldsymbol{\beta}, \eta, \sigma^2)\pi(\mathbf{b}|\eta, \sigma^2) \\ &\propto \exp \left\{ -\frac{1}{2\sigma^2}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b})^T \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b}) - \frac{1}{2\eta\sigma^2} \mathbf{b}^T \boldsymbol{\Phi}^{-1} \mathbf{b} \right\} \\ &= \exp \left\{ -\frac{1}{2\sigma^2} \left[(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b})^T \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b}) + \frac{1}{\eta} \mathbf{b}^T \boldsymbol{\Phi}^{-1} \mathbf{b} \right] \right\} \quad (28) \end{aligned}$$

Solving for (27) is equivalent to minimizing the exponent in (28):

$$\hat{\mathbf{b}} = \arg \min_{\mathbf{b}} \left\{ (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b})^T \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{b}) + \frac{1}{\eta} \mathbf{b}^T \boldsymbol{\Phi}^{-1} \mathbf{b} \right\} \quad (29)$$

Taking the derivative of (29) with respect to \mathbf{b} and setting it to 0 we get:

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

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 optimal tuning parameter

In order to choose the optimal value of the tuning parameter λ , we use the generalized information criterion (34) (GIC):

$$GIC_{\lambda} = -2\ell(\hat{\boldsymbol{\beta}}, \hat{\sigma}^2, \hat{\eta}) + a_n \cdot \hat{df}_{\lambda} \quad (30)$$

where \hat{df}_{λ} is the number of non-zero elements in $\hat{\boldsymbol{\beta}}_{\lambda}$ (35) plus two (representing the variance parameters η and σ^2). Several authors have used this criterion for variable selection in mixed models with $a_n = \log N_T$ (31, 36), which corresponds to the BIC. We instead choose the high-dimensional BIC (37) given by $a_n = \log(\log(N_T)) * \log(p)$. This is the default choice in our `ggmix` R package, though the interface is flexible to allow the user to select their choice of a_n .

4 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
- ρ : 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 ρ .
- $\mathbf{X}^{(kinship)}$: $n \times k$ matrix of SNPs used to construct the kinship matrix.
- β_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 ε 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 simulate data from

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{b} + \boldsymbol{\varepsilon} \quad (31)$$

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

$$\mathbf{b} \sim \mathcal{N}(0, \eta\boldsymbol{\Phi}) \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, (1 - \eta)\mathbf{I}) \quad (32)$$

and $\eta = 0.1$. $n = 1k, p = 5k, X_{kinship} = 10k$. 1% of the 5k SNPs are causal.

Scenario 1

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

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

Scenario 2

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

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

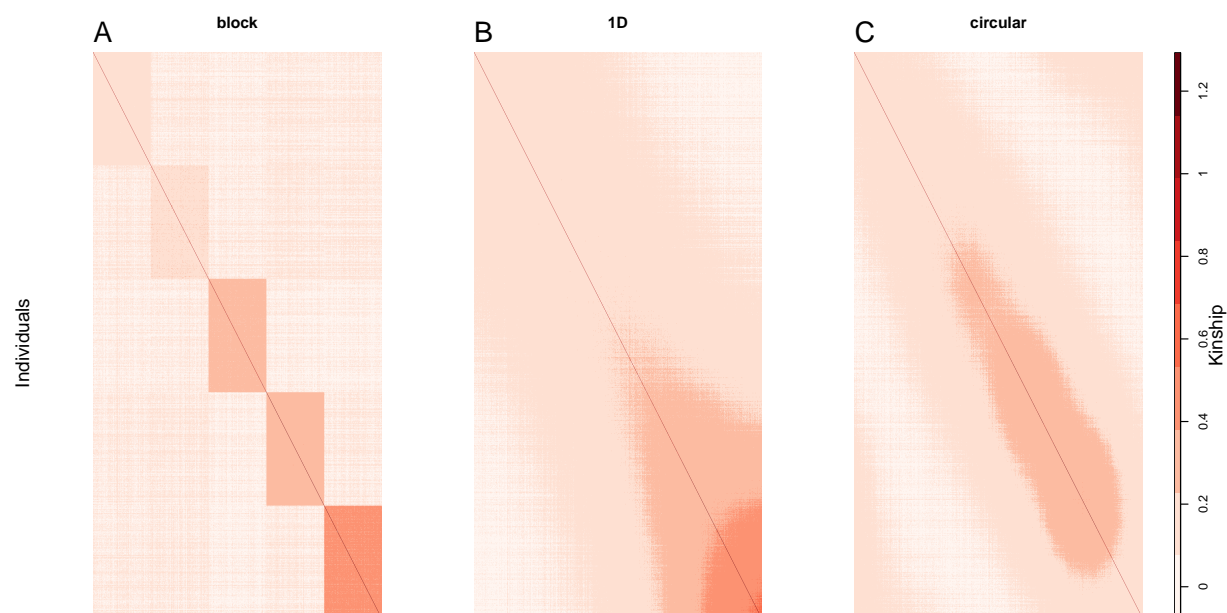


Figure 1: Empirical kinship matrices used in simulation studies.

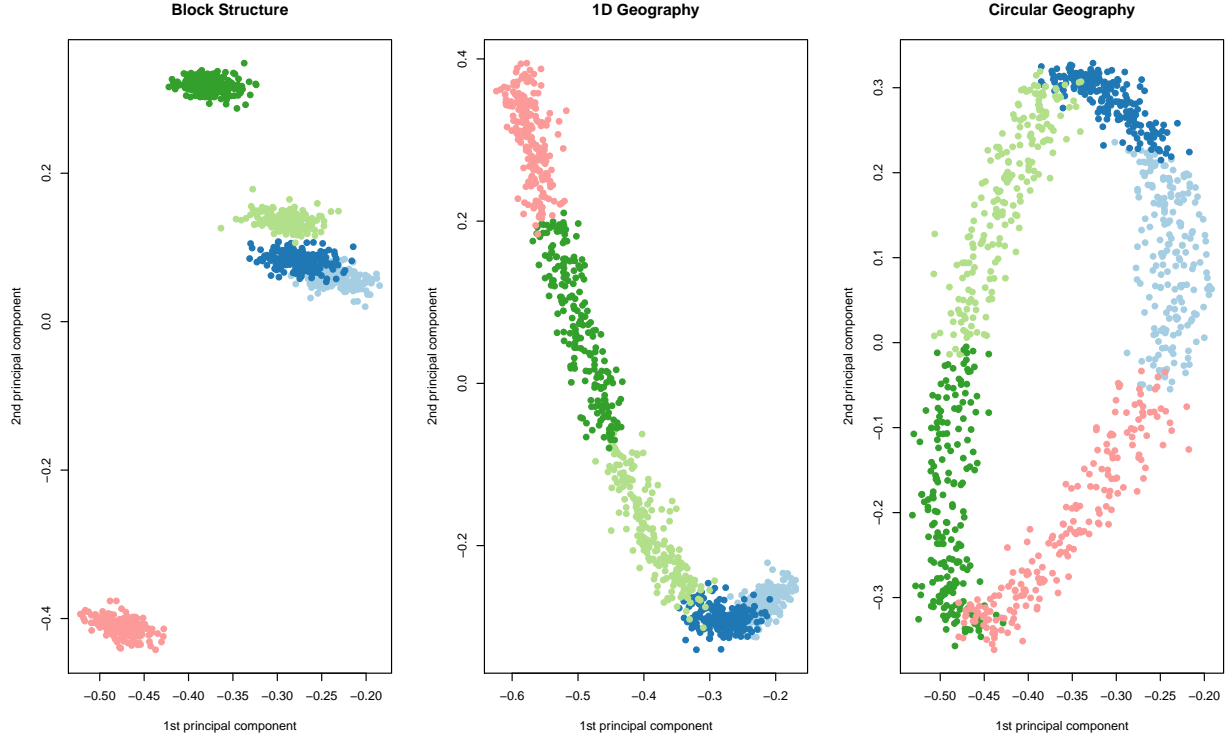


Figure 2: First two principal component scores of the kinship matrix where each color represents one of the 5 simulated subpopulations. The first panel corresponds to 5 independent subpopulations, the second corresponds to a 1 dimensional geographical structure and the third panel corresponds to a circular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

4.1 Results

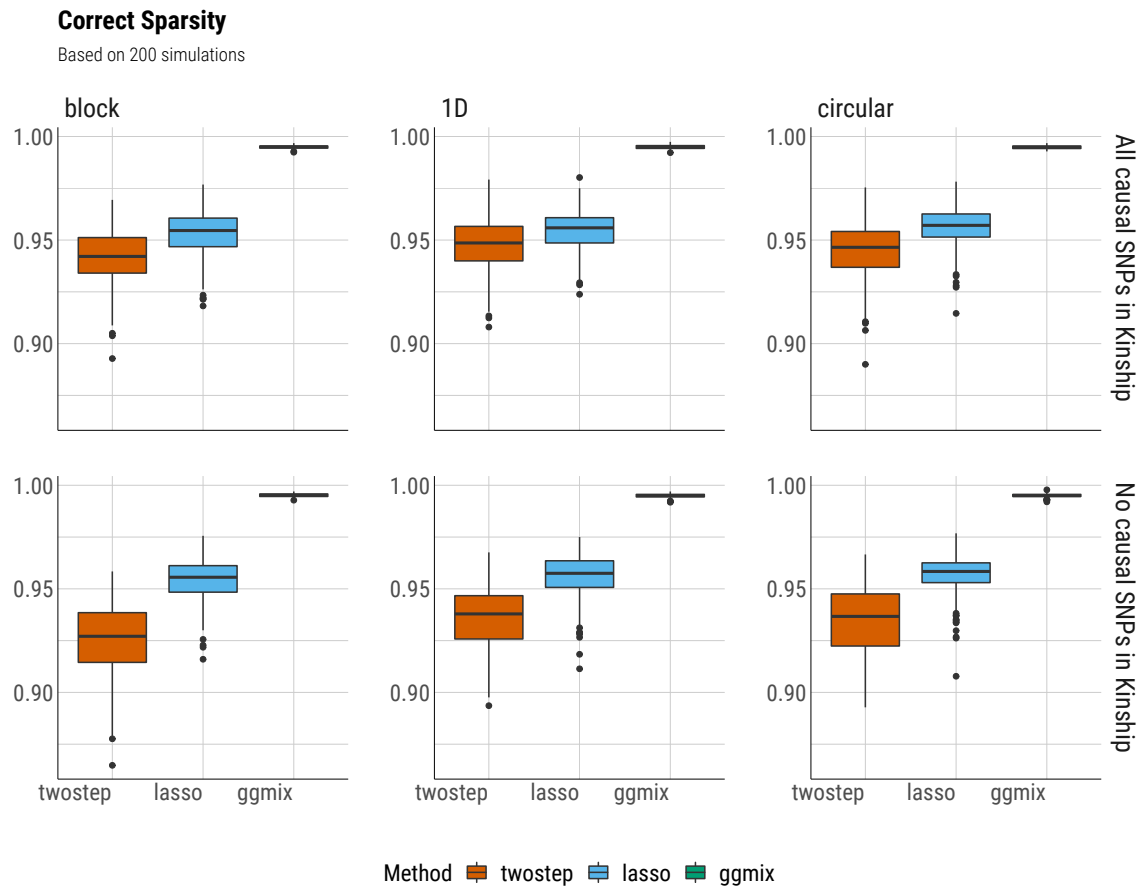


Figure 3: Boxplots of the correct sparsity from 200 simulations by kinship geography and number of causal SNPs that were included in the calculation of the kinship matrix.

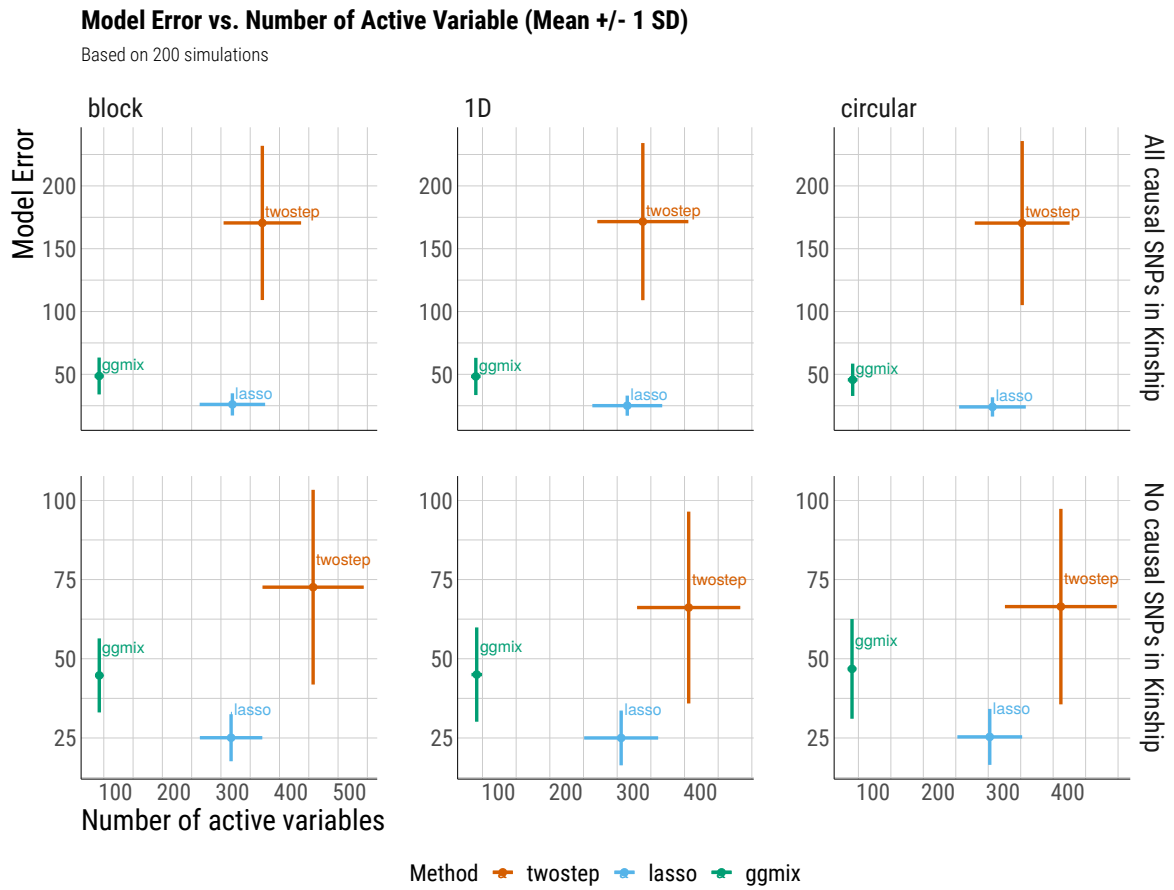


Figure 4: Model error vs number of active variables results.

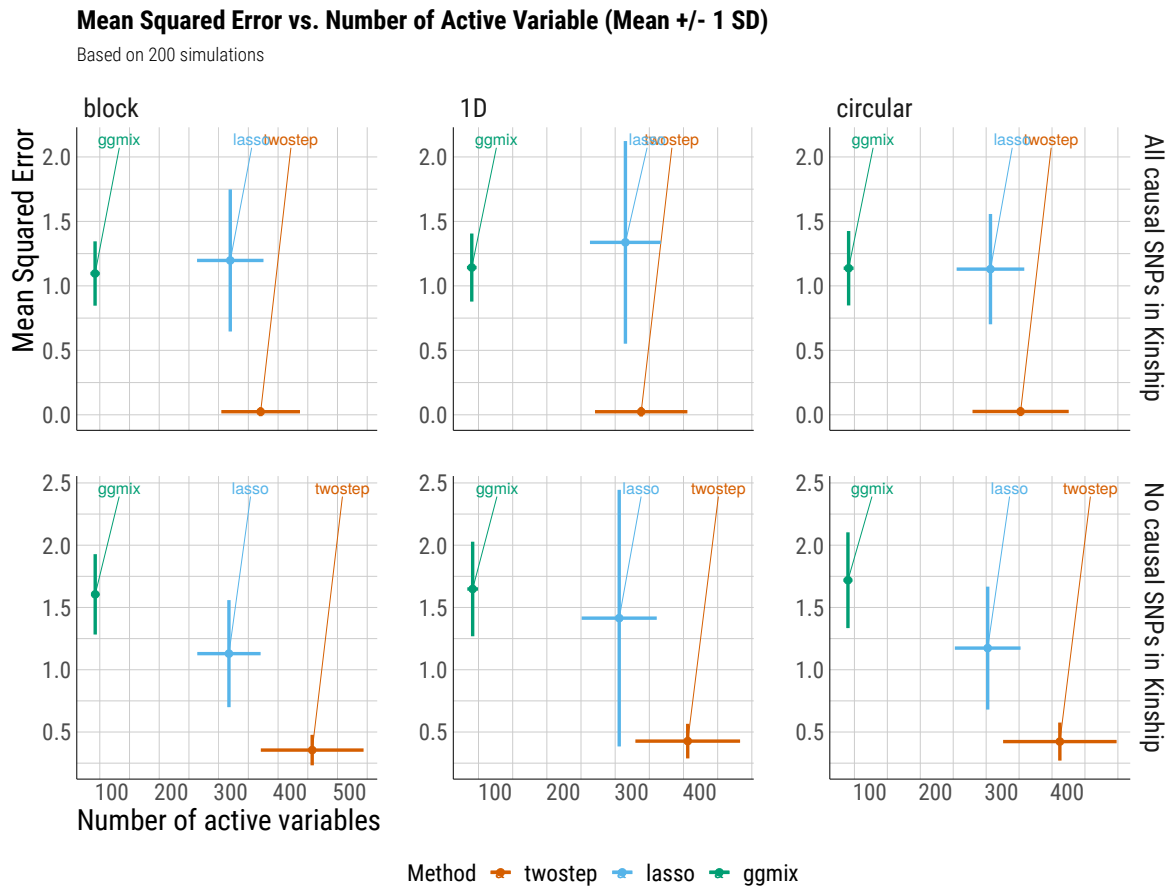


Figure 5: Mean squared error vs number of active variables results.

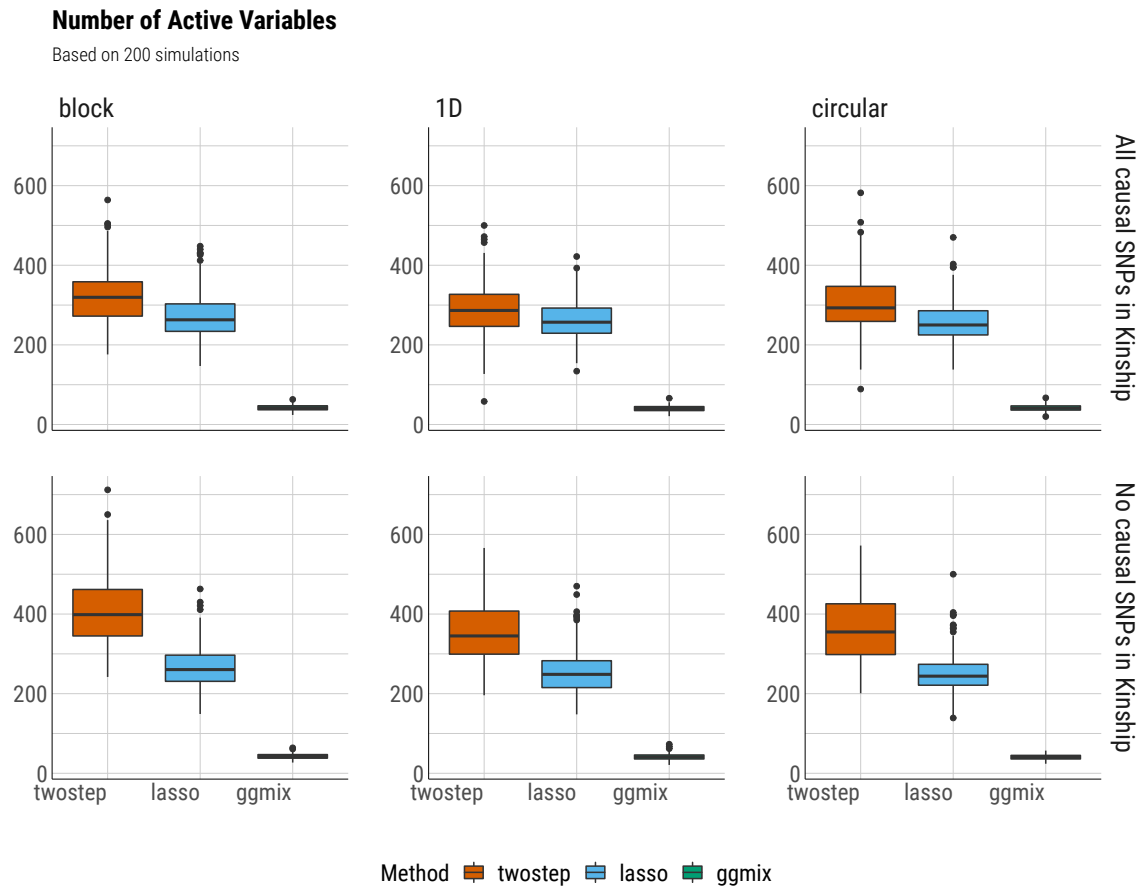


Figure 6: Boxplots of the number of active variables from 200 simulations by kinship geography and number of causal SNPs that were included in the calculation of the kinship matrix.

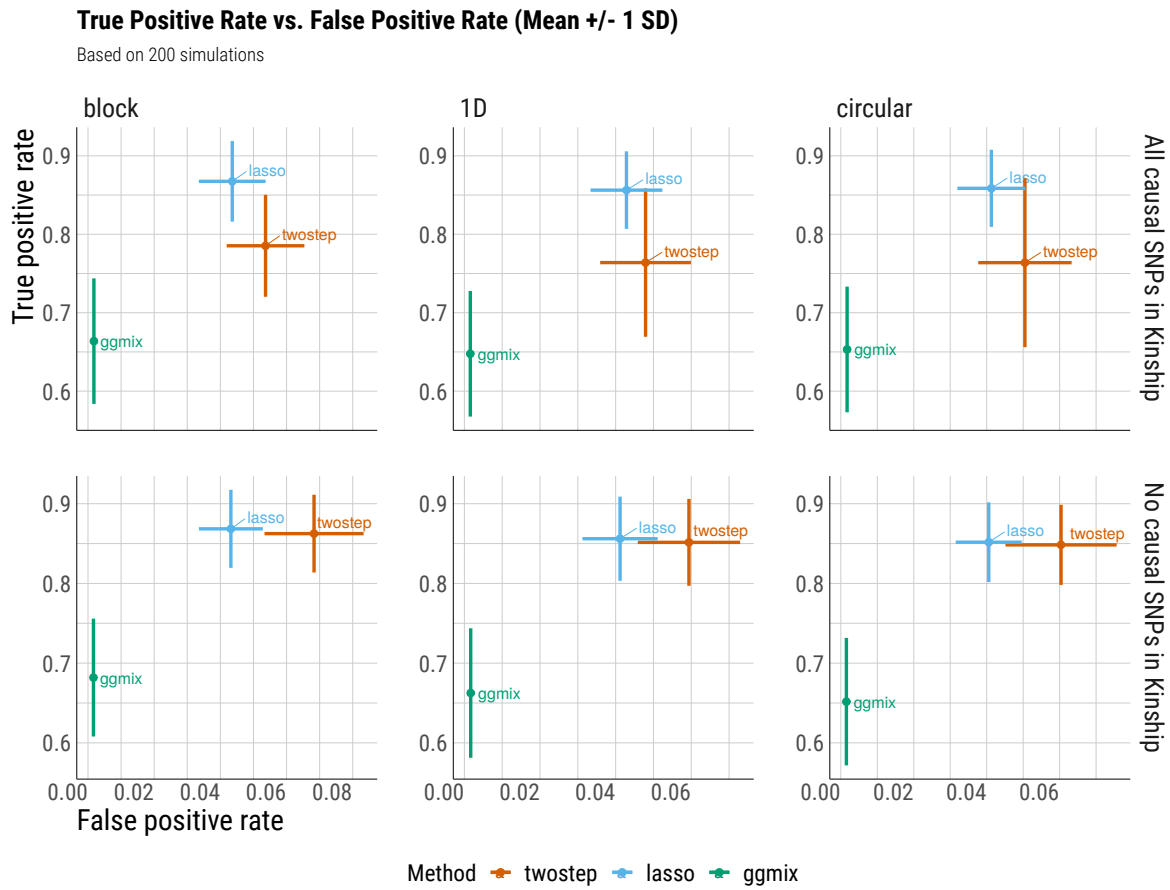


Figure 7: Means \pm 1 standard deviation of true positive rate vs. false positive rate from 200 simulations by kinship geography and number of causal SNPs that were included in the calculation of the kinship matrix.

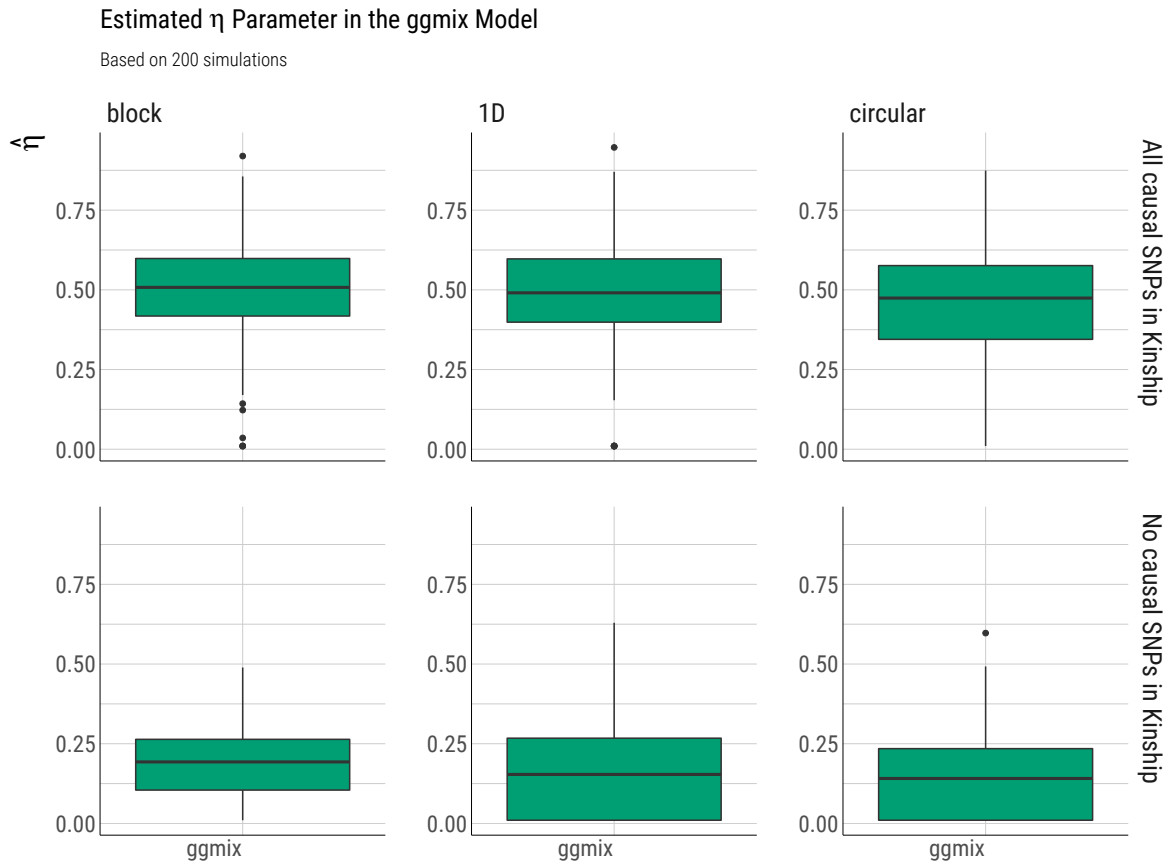


Figure 8: Boxplot of $\hat{\eta}$ in the `ggmix` model from 200 simulations by kinship geography and number of causal SNPs that were included in the calculation of the kinship matrix.

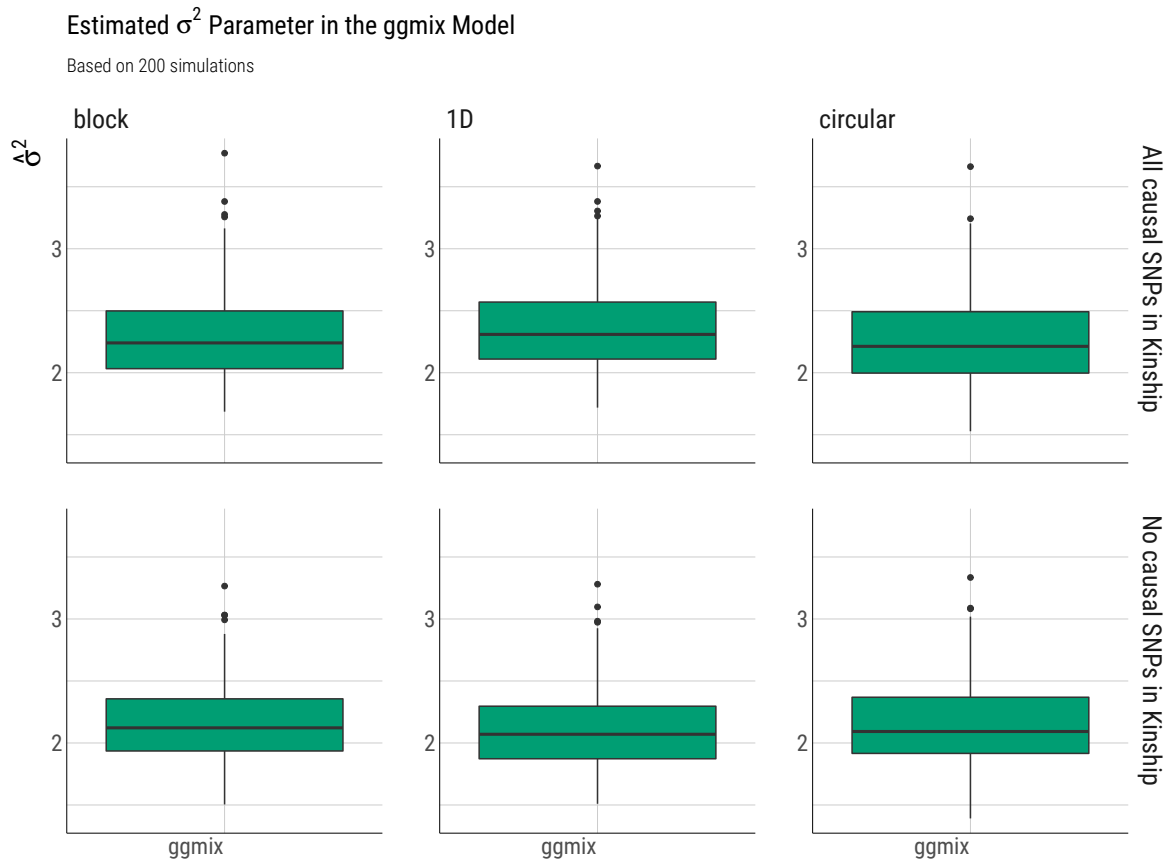


Figure 9: Boxplot of $\hat{\eta}$ in the `ggmix` model from 200 simulations by kinship geography and number of causal SNPs that were included in the calculation of the kinship matrix.

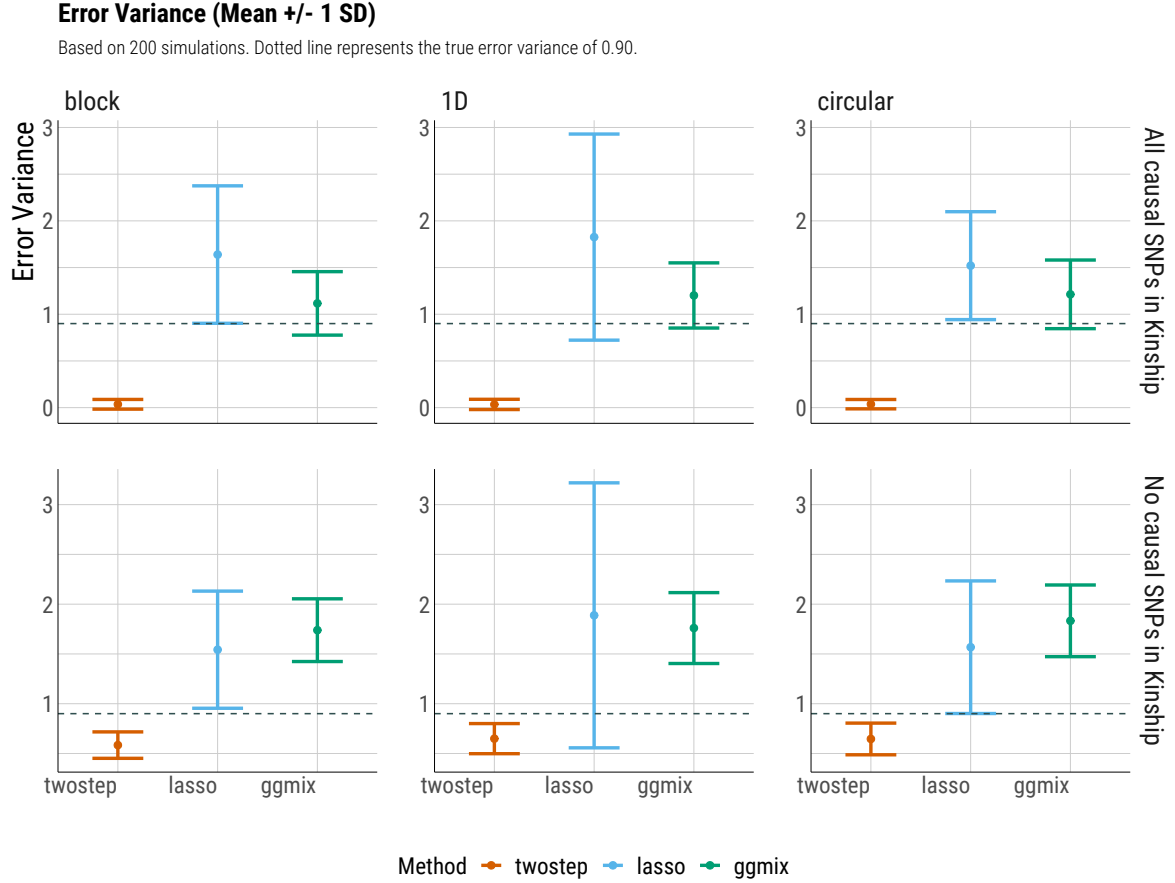


Figure 10: Means \pm 1 standard deviation of error variance from 200 simulations by kinship geography and number of causal SNPs that were included in the calculation of the kinship matrix.

5 Discussion

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 are

method has better power over the two-stage approach, particularly for polygenic traits.

While the predominant motivation for these methods has been association testing, we believe that there are other applications in which they can be used as well. For example, in the most recent Genetic Analysis Workshop 20 (GAW20), the causal modeling group investigated causal relationships between DNA methylation (exposure) within some genes and the change in high-density lipoproteins Δ HDL (outcome) using Mendelian randomization (MR) (26). Penalized regression methods could be used to select SNPs strongly associated with the exposure in order to be used as an instrumental variable (IV). However, since GAW20 data consisted of families, two step methods were used which could have resulted in a loss of power. `ggmix` is an alternative approach that could be used for selecting the IV while accounting for the familial structure of the data. Our method is also suitable for fine mapping SNP association signals in genomic regions, where the goal is to pinpoint individual variants most likely to impact the underlying biological mechanisms of disease (27).

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](#), [7](#)
- [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](#)
- [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] Jiahua Li, Kiranmoy Das, Guifang Fu, Runze Li, and Rongling Wu. The bayesian lasso for genome-wide association studies. *Bioinformatics*, 27(4):516–523, 2010. [3](#)
- [8] 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. [3](#), [5](#)

-
- [9] Hyun Min Kang, Jae Hoon Sul, Noah A Zaitlen, Sit-yee 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] Jianming Yu, Gael Pressoir, William H Briggs, Irie Vroh Bi, Masanori Yamasaki, John F Doebley, Michael D McMullen, Brandon S Gaut, Dahlia M Nielsen, James B Holland, et al. A unified mixed-model method for association mapping that accounts for multiple levels of relatedness. *Nature genetics*, 38(2):203, 2006. 3
- [11] 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. 3
- [12] Alkes L Price, Nick J Patterson, Robert M Plenge, Michael E Weinblatt, Nancy A Shadick, and David Reich. Principal components analysis corrects for stratification in genome-wide association studies. *Nature genetics*, 38(8):904, 2006. 3
- [13] Karim Oualkacha, Zari Dastani, Rui Li, Pablo E Cingolani, Timothy D Spector, Christopher J Hammond, J Brent Richards, Antonio Ciampi, and Celia MT Greenwood. Adjusted sequence kernel association test for rare variants controlling for cryptic and family relatedness. *Genetic epidemiology*, 37(4):366–376, 2013. 3, 4
- [14] Heather J Cordell and David G Clayton. A unified stepwise regression procedure for evaluating the relative effects of polymorphisms within a gene using case/control or family data: application to hla in type 1 diabetes. *The American Journal of Human Genetics*, 70(1):124–141, 2002. 3
- [15] 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. 4, 7

-
- [16] 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. [4](#)
- [17] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996. [4](#), [5](#)
- [18] Hui Zou. The adaptive lasso and its oracle properties. *Journal of the American statistical association*, 101(476):1418–1429, 2006. [4](#)
- [19] Xiuhua Ding, Shaoyong Su, Kannabiran Nandakumar, Xiaoling Wang, and David W Fardo. A 2-step penalized regression method for family-based next-generation sequencing association studies. In *BMC proceedings*, volume 8, page S25. BioMed Central, 2014. [4](#)
- [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. [4](#), [10](#), [14](#), [34](#), [40](#)
- [21] Yi Yang and Hui Zou. A fast unified algorithm for solving group-lasso penalize learning problems. *Statistics and Computing*, 25(6):1129–1141, 2015. [4](#)
- [22] Jian Yang, Noah A Zaitlen, Michael E Goddard, Peter M Visscher, and Alkes L Price. Advantages and pitfalls in the application of mixed-model association methods. *Nature genetics*, 46(2):100, 2014. [4](#)
- [23] Naomi Allen, Cathie Sudlow, Paul Downey, Tim Peakman, John Danesh, Paul Elliott, John Gallacher, Jane Green, Paul Matthews, Jill Pell, et al. Uk biobank: Current status and what it means for epidemiology. *Health Policy and Technology*, 1(3):123–126, 2012. [4](#)

-
- [24] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005. 5
- [25] Ming Yuan and Yi Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006. 5
- [26] George Davey Smith and Shah Ebrahim. mendelian randomization: can genetic epidemiology contribute to understanding environmental determinants of disease? *International journal of epidemiology*, 32(1):1–22, 2003. 5
- [27] Sarah L Spain and Jeffrey C Barrett. Strategies for fine-mapping complex traits. *Human molecular genetics*, 24(R1):R111–R119, 2015. 5
- [28] 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
- [29] Paul Tseng and Sangwoon Yun. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*, 117(1):387–423, 2009. 10, 34, 37
- [30] 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. 10, 34
- [31] 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. 10, 16, 34
- [32] Richard H Byrd, Peihuang Lu, Jorge Nocedal, and Ciyou Zhu. A limited memory

- algorithm for bound constrained optimization. *SIAM Journal on Scientific Computing*, 16(5):1190–1208, 1995. 12
- [33] Jon Wakefield. *Bayesian and frequentist regression methods*. Springer Science & Business Media, 2013. 15
- [34] Ryuei Nishii. Asymptotic properties of criteria for selection of variables in multiple regression. *The Annals of Statistics*, pages 758–765, 1984. 16
- [35] 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
- [36] 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
- [37] Yingying Fan and Cheng Yong Tang. Tuning parameter selection in high dimensional penalized likelihood. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 75(3):531–552, 2013. 16
- [38] Yihui Xie. *Dynamic Documents with R and knitr*, volume 29. CRC Press, 2015. 41

A Block Coordinate Descent Algorithm

We use a general purpose block coordinate descent algorithm (CGD) (29) 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 (29). For continuously differentiable $f(\cdot)$ and convex and block-separable $P(\cdot)$ (i.e. $P(\beta) = \sum_i P_i(\beta_i)$), (29) 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. (30), (20)) and (31) for mixed models with an ℓ_1 penalty. Following Tseng and Yun (29), the CGD algorithm is given by Algorithm 2.

The Armijo rule is defined as follows (29):

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

$$Q_\lambda(\Theta_j^{(k)} + \alpha^{(k)} d^{(k)}) \leq Q_\lambda(\Theta_j^{(k)}) + \alpha^{(k)} \varrho \Delta^{(k)} \quad (38)$$

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

$$\Delta^{(k)} := \nabla f(\Theta_j^{(k)}) d^{(k)} + \gamma (d^{(k)})^2 H_{jj}^{(k)} + \lambda P(\Theta_j^{(k)} + d^{(k)}) - \lambda P(\Theta_j^{(k)}) \quad (39)$$

Common choices for the constants are $\delta = 0.1$, $\varrho = 0.001$, $\gamma = 0$, $\alpha_{init}^{(k)} = 1$ for all k (31).

Below we detail the specifics of Algorithm 2 for the ℓ_1 penalty.

Algorithm 2: Coordinate Gradient Descent Algorithm to solve (15)

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 $H^{(k)}$:

$$H^{(k)} = \text{diag} \left[\min \left\{ \max \left\{ \left[\nabla^2 f(\Theta^{(k)}) \right]_{jj}, c_{\min} \right\}, c_{\max} \right\} \right]_{j=1, \dots, p+1} \quad (33)$$

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**

$$d_{H^{(k)}}(\Theta_j^{(k)}) \leftarrow \arg \min_d \left\{ \nabla f(\Theta_j^{(k)})d + \frac{1}{2}d^2 H_{jj}^{(k)} + \lambda P(\Theta_j^{(k)} + d) \right\} \quad (34)$$

end

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

$$d_{H^{(k)}}(\Theta_j^{(k)}) \leftarrow -\nabla f(\Theta_j^{(k)})/H_{jj}^{(k)} \quad (35)$$

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;

$$\hat{\eta}^{(k+1)} \leftarrow \arg \min_{\eta} \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^{2(k)}} \sum_{i=1}^{N_T} \frac{\left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j^{(k+1)} \right)^2}{1 + \eta(\Lambda_i - 1)} \quad (36)$$

Update;

$$\hat{\sigma}^{2(k+1)} \leftarrow \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{\left(\tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j^{(k+1)} \right)^2}{1 + \eta^{(k+1)}(\Lambda_i - 1)} \quad (37)$$

$k \leftarrow k + 1$

until convergence criterion is satisfied;

A.1 ℓ_1 penalty

The objective function is given by

$$Q_\lambda(\boldsymbol{\Theta}) = f(\boldsymbol{\Theta}) + \lambda|\boldsymbol{\beta}| \quad (40)$$

A.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

$$d_H(\Theta_j) = \arg \min_d G(d) \quad (41)$$

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 $0 \in \partial G(d)$:

$$\partial G(d) = \nabla f(\Theta_j) + dH_{jj} + \lambda u \quad (42)$$

where

$$u = \begin{cases} 1 & \text{if } d > -\Theta_j \\ -1 & \text{if } d < -\Theta_j \\ [-1, 1] & \text{if } d = -\Theta_j \end{cases} \quad (43)$$

We consider each of the three cases in (42) 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 (29).

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 (41). Therefore the descent direction for $\Theta_j^{(k)} \in \{\beta_1, \dots, \beta_p\}$ for the ℓ_1 penalty is given by

$$d = \text{mid} \left\{ \frac{-(\nabla f(\beta_j) - \lambda)}{H_{jj}}, -\beta_j, \frac{-(\nabla f(\beta_j) + \lambda)}{H_{jj}} \right\} \quad (44)$$

A.1.2 Solution for the β parameter

If the Hessian $\nabla^2 f(\Theta^{(k)}) > 0$ then $H^{(k)}$ defined in (33) 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 β parameter is then given by

$$\beta_j^{(k+1)} \leftarrow \beta_j^{(k)} + d^{(k)}, \quad j = 1, \dots, p \quad (45)$$

Substituting the descent direction given by (44) into (45) we get

$$\beta_j^{(k+1)} = \text{mid} \left\{ \beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) - \lambda)}{H_{jj}}, 0, \beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) + \lambda)}{H_{jj}} \right\} \quad (46)$$

We can further simplify this expression. Let

$$w_i := \frac{1}{\sigma^2 (1 + \eta(\Lambda_i - 1))} \quad (47)$$

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

$$g(\beta^{(k)}) = \frac{1}{2} \sum_{i=1}^{N_T} w_i \left(\tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} - \tilde{X}_{ij} \beta_j^{(k)} \right)^2 \quad (48)$$

The gradient and Hessian are given by

$$\nabla f(\beta_j^{(k)}) := \frac{\partial}{\partial \beta_j^{(k)}} g(\beta^{(k)}) = - \sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left(\tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} - \tilde{X}_{ij} \beta_j^{(k)} \right) \quad (49)$$

$$H_{jj} := \frac{\partial^2}{\partial \beta_j^{(k)2}} g(\beta^{(k)}) = \sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2 \quad (50)$$

Substituting (49) and (50) into $\beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) - \lambda)}{H_{jj}}$

$$\begin{aligned} & \beta_j^{(k)} + \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left(\tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} - \tilde{X}_{ij} \beta_j^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \\ &= \beta_j^{(k)} + \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left(\tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} - \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2 \beta_j^{(k)}}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \\ &= \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left(\tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \end{aligned} \quad (51)$$

Similarly, substituting (49) and (50) in $\beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) + \lambda)}{H_{jj}}$ we get

$$\frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left(\tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) - \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \quad (52)$$

Finally, substituting (51) and (52) into (46) we get

$$\begin{aligned}\beta_j^{(k+1)} &= \text{mid} \left\{ \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left(\tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) - \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2}, 0, \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left(\tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \right\} \\ &= \frac{\mathcal{S}_\lambda \left(\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left(\tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) \right)}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2}\end{aligned}\quad (53)$$

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 β_j given by (53) takes the same form as the weighted updates of the `glmnet` algorithm (20) (Section 2.4, equation (10)) with $\alpha = 1$.

B ggmix Package Showcase

In this section we briefly introduce the freely available and open source `ggmix` package in R. More comprehensive documentation is available at <https://sahirbhatnagar.com/ggmix>. Note that this entire section is reproducible; the code and text are combined in an `.Rnw`¹ file and compiled using `knitr` (38).

B.1 Installation

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

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

To showcase the main functions in `ggmix`, we will use the simulated data which ships with 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"
```

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

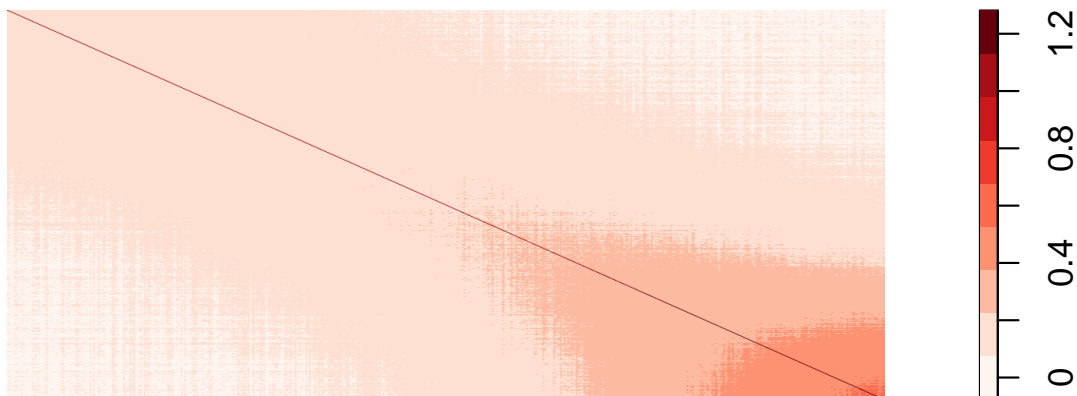
There are three basic inputs that `ggmix` needs:

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

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

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)
```



B.2 Fit the linear mixed model with Lasso Penalty

We will use the most basic call to the main function of this package, which is called `ggmix`. This function will by default fit a L_1 penalized linear mixed model (LMM) for 100 distinct values of the tuning parameter λ . It will choose its own sequence:

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

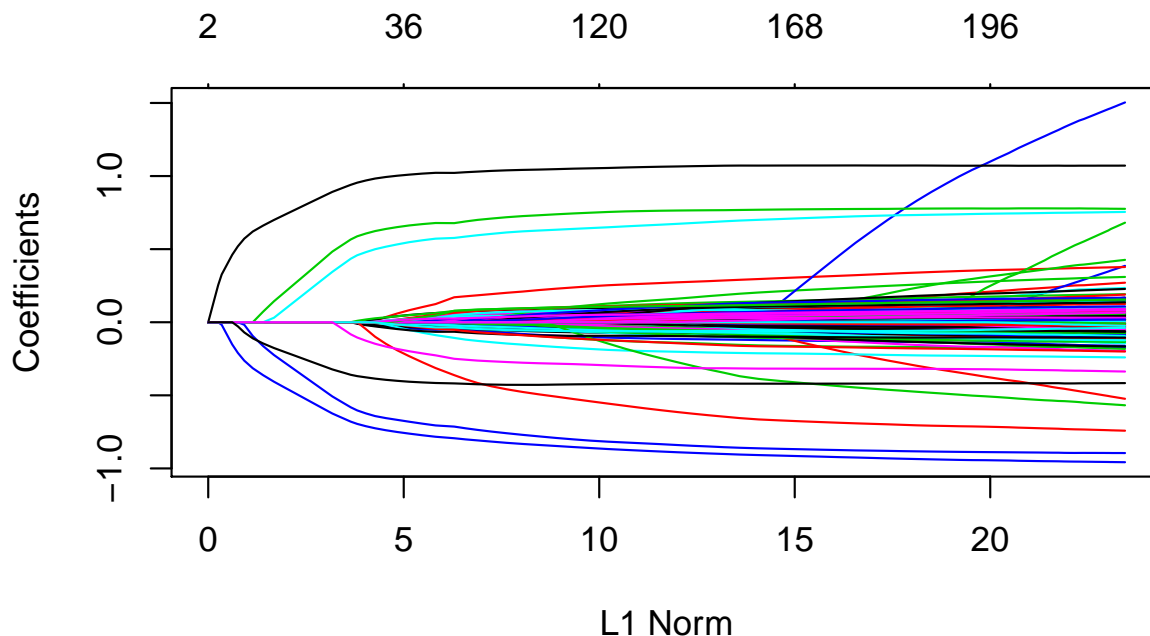
## [1] "result"      "ggmix_object" "n_design"     "p_design"
## [5] "lambda"      "coef"         "b0"           "beta"
## [9] "df"          "eta"          "sigma2"       "nlambda"
## [13] "cov_names"   "call"

class(fit)
```

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

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

```
plot(fit)
```



We can also get the coefficients for given value(s) of λ using the `coef` method for 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 -0.030227753
## X62          0.0000000  0.000000000
## X185          0.0000000  0.001444670
## X371          0.0000000  0.009513604
## X420          0.0000000  0.000000000
```

Here, `s` specifies the value(s) of λ at which the extraction is made. The function uses linear

interpolation to make predictions for values of \mathbf{s} that do not coincide with the lambda sequence used in the fitting algorithm.

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 -1.3123396
## id2 -0.02913052  0.3885921
## id3 -2.00084875 -2.6460045
## id4 -0.37255277 -0.9542455
## id5 -1.03967831 -2.1377274
```

B.3 Find the Optimal Value of the Tuning Parameter

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

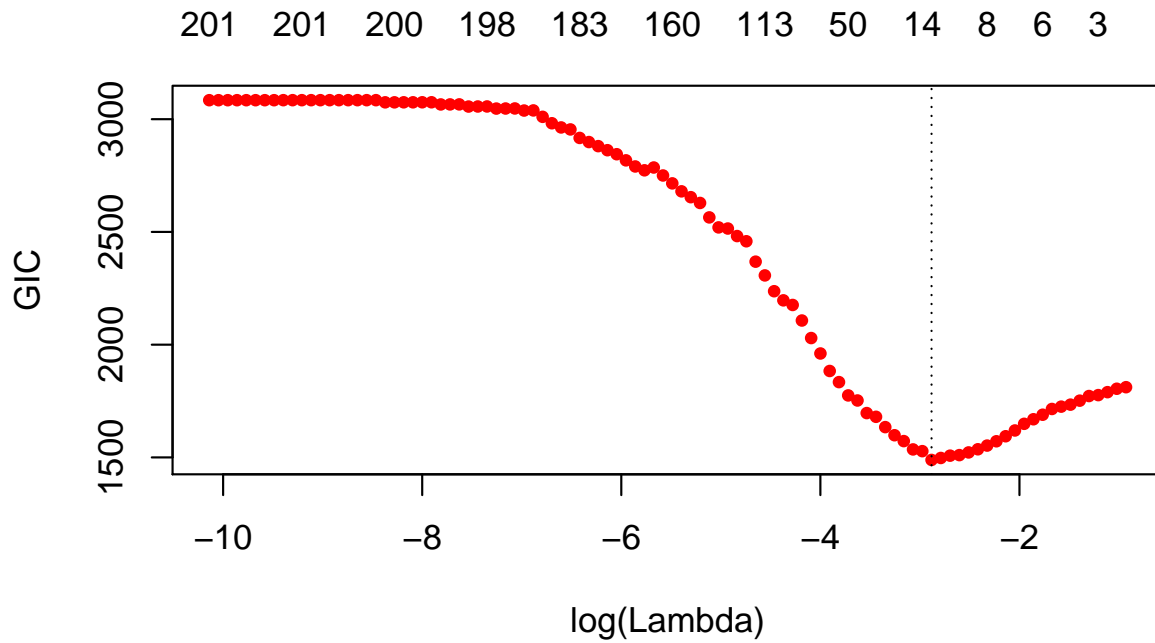
```
# pass the fitted object from ggmix to the gic function:
hdbic <- gic(fit)
class(hdbic)

## [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)))
```

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

```
plot(hdbic)
```

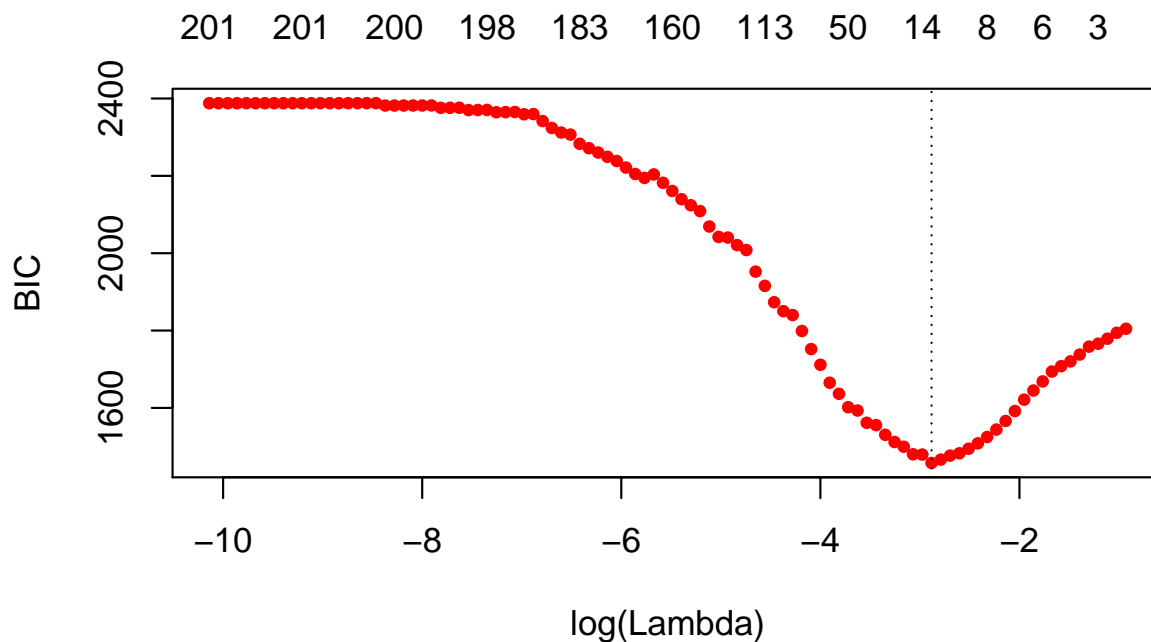


The optimal value for λ according to the HDBIC, i.e., the λ that leads to the minimum HDBIC is:

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

We can also plot the BIC results:

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



```
bicfit[["lambda.min"]]
```

```
## [1] 0.05596623
```

B.4 Get Coefficients Corresponding to Optimal Model

We can use the object outputted by the `gic` function to extract the coefficients corresponding 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         .
```

We can also extract just the nonzero coefficients which also provide the estimated variance components η and σ^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
```

We can also make predictions from the `hdbic` object, which by default will use the model 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
```

B.5 Extracting Random Effects

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

```
ranef(hdbic)[1:5]

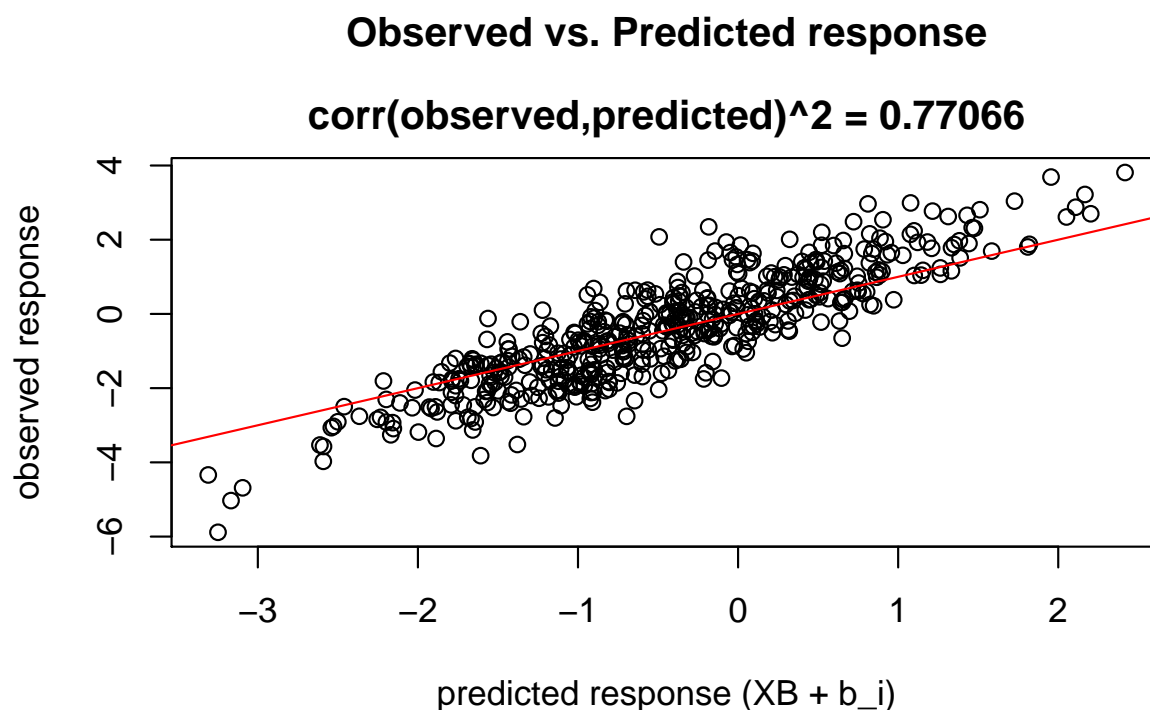
## [1] -0.02548691 -0.10011680  0.13020240 -0.30650997  0.16045768
```


B.6 Diagnostic Plots

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

B.6.1 Observed vs. Predicted Response

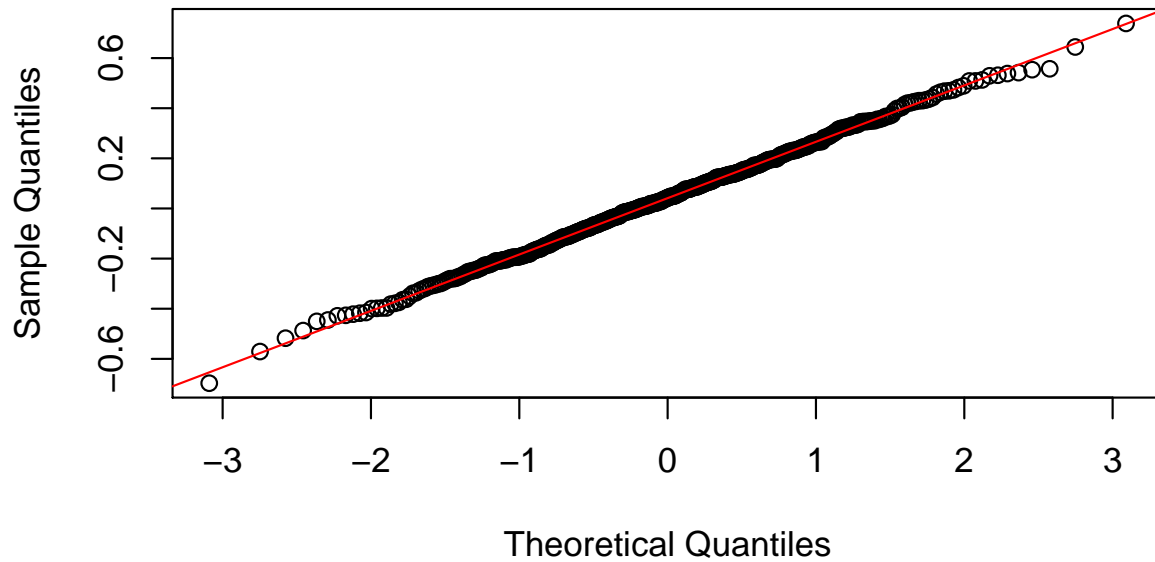
```
plot(hdbic, type = "predicted", newx = admixed$x, newy = admixed$y)
```



B.6.2 QQ-plots for Residuals and Random Effects

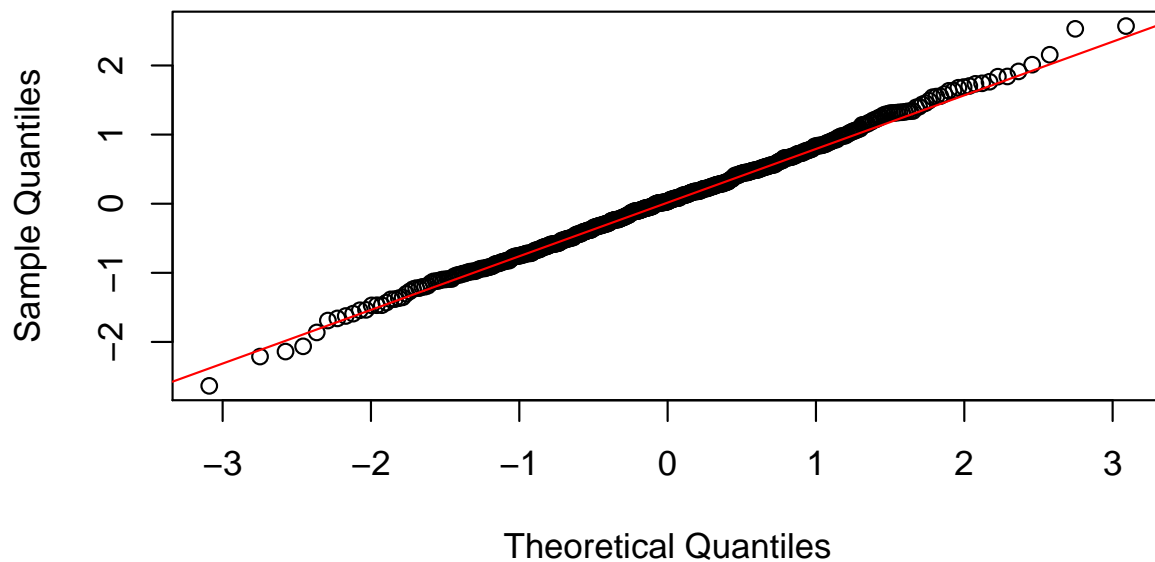
```
plot(hdbic, type = "QQranef", newx = admixed$x, newy = admixed$y)
```

QQ-Plot of the random effects at $\lambda = 0.06$



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

QQ-Plot of the residuals at $\lambda = 0.06$



B.6.3 Tukey-Anscombe Plot

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

