

---

# Simultaneous SNP selection and adjustment for population structure in high dimensional prediction models

Sahir R Bhatnagar<sup>1,2,\*</sup>, Yi Yang<sup>3</sup>, Tianyuan Lu<sup>4,5</sup>, Erwin Schurr<sup>6</sup>, JC Loredo-Osti<sup>7</sup>, Marie Forest<sup>8</sup>, Karim Oualkacha<sup>9</sup>, Celia MT Greenwood<sup>1,4,5,10,11</sup>

**1** Department of Epidemiology, Biostatistics and Occupational Health, McGill University, Montréal, Québec, Canada

**2** Department of Diagnostic Radiology, McGill University, Montréal, Québec, Canada

**3** Department of Mathematics and Statistics, McGill University, Montréal, Québec, Canada

**4** Quantitative Life Sciences, McGill University, Montréal, Québec, Canada

**5** Lady Davis Institute, Jewish General Hospital, Montréal, Québec, Canada

**6** Department of Medicine, McGill University, Montréal, Québec, Canada

**7** Department of Mathematics and Statistics, Memorial University, St. John's, Newfoundland, Canada

**8** École de Technologie Supérieure, Montréal, Québec, Canada

**9** Département de Mathématiques, UQÀM, Montréal, Québec, Canada

**10** Department of Oncology, McGill University, Montréal, Québec, Canada

**11** Department of Human Genetics, McGill University, Montréal, Québec, Canada

\* sahir.bhatnagar@mcgill.ca

## Supporting Information

Contains the following sections:

---

**A Block Coordinate Descent Algorithm** - a detailed description of the algorithm used to fit our `gmmix` model.

**B Additional Real Data Analysis Results** - supporting information for the GAW20 and UK Biobank analyses

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

---

## S1 Text

### Block Coordinate Descent Algorithm

#### Model Set-up

Let  $i = 1, \dots, N$  be a 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  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)^T \in \mathbb{R}^{N_T \times 1}$ ,  $\mathbf{b} = (\mathbf{b}_1, \dots, \mathbf{b}_N)^T \in \mathbb{R}^{N_T \times 1}$ ,  $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_N)^T \in \mathbb{R}^{N_T \times 1}$ , and the stacked matrix  $\mathbf{X} = (\mathbf{X}_1^T, \dots, \mathbf{X}_N^T) \in \mathbb{R}^{N_T \times (p+1)}$ . 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 [1]:

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

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

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

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

---

multivariate normal:

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

We consider the parameterization in (3) since maximization is easier over the compact set  $\eta \in [0, 1]$  than over the unbounded interval  $\delta \in [0, \infty)$  [1]. We define the complete parameter vector as  $\boldsymbol{\Theta} := (\boldsymbol{\beta}, \eta, \sigma^2)$ . The negative log-likelihood for (3) is given by

$$-\ell(\boldsymbol{\Theta}) \propto \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \log(\det(\mathbf{V})) + \frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \quad (4)$$

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  $\Lambda_i$ . In the main text we show that  $\mathbf{V}$  can then be further simplified to

$$\mathbf{V} = \mathbf{U}\tilde{\mathbf{D}}\mathbf{U}^T \quad (5)$$

where

$$\tilde{\mathbf{D}} = \text{diag}\{1 + \eta(\Lambda_1 - 1), 1 + \eta(\Lambda_2 - 1), \dots, 1 + \eta(\Lambda_{N_T} - 1)\} \quad (6)$$

Since (6) is a diagonal matrix, its inverse is also a diagonal matrix:

$$\tilde{\mathbf{D}}^{-1} = \text{diag}\left\{\frac{1}{1 + \eta(\Lambda_1 - 1)}, \frac{1}{1 + \eta(\Lambda_2 - 1)}, \dots, \frac{1}{1 + \eta(\Lambda_{N_T} - 1)}\right\} \quad (7)$$

From (5) and (6),  $\log(\det(\mathbf{V}))$  simplifies to

$$\begin{aligned}
\log(\det(\mathbf{V})) &= \log \left( \det(\mathbf{U}) \det(\tilde{\mathbf{D}}) \det(\mathbf{U}^T) \right) \\
&= \log \left\{ \prod_{i=1}^{N_T} (1 + \eta(\Lambda_i - 1)) \right\} \\
&= \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1))
\end{aligned} \tag{8}$$

since  $\det(\mathbf{U}) = 1$ . It also follows from (5) that

$$\begin{aligned}
\mathbf{V}^{-1} &= \left( \mathbf{U} \tilde{\mathbf{D}} \mathbf{U}^T \right)^{-1} \\
&= (\mathbf{U}^T)^{-1} \left( \tilde{\mathbf{D}} \right)^{-1} \mathbf{U}^{-1} \\
&= \mathbf{U} \tilde{\mathbf{D}}^{-1} \mathbf{U}^T
\end{aligned} \tag{9}$$

since for an orthonormal matrix  $\mathbf{U}^{-1} = \mathbf{U}^T$ . Substituting (7), (8) and (9) into (4) the negative log-likelihood becomes

$$-\ell(\Theta) = \frac{N_T}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^2} \sum_{i=1}^{N_T} \frac{\left( \tilde{Y}_i - \sum_{j=0}^p \tilde{X}_{ij+1} \beta_j \right)^2}{1 + \eta(\Lambda_i - 1)} \tag{10}$$

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.

## Penalized Maximum Likelihood Estimator

We define the  $p + 3$  length vector of parameters  $\Theta := (\Theta_0, \Theta_1, \dots, \Theta_{p+1}, \Theta_{p+2}, \Theta_{p+3}) = (\boldsymbol{\beta}, \eta, \sigma^2)$  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  $\Theta$  for  $\eta$  and  $\sigma^2$ , respectively. In light of our goals to select variables associated with the response in high-dimensional data, we propose to place a constraint on the magnitude of the regression coefficients. This can be achieved by adding a penalty term to the likelihood

---

function (10). 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:

$$Q_\lambda(\Theta) = f(\Theta) + \lambda \sum_{j \neq 0} v_j P_j(\beta_j) \quad (11)$$

where  $f(\Theta) := -\ell(\Theta)$  is defined in (10),  $P_j(\cdot)$  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  $\sigma^2$ . An estimate of the regression parameters  $\widehat{\Theta}_\lambda$  is obtained by

$$\widehat{\Theta}_\lambda = \arg \min_{\Theta} Q_\lambda(\Theta) \quad (12)$$

We use a general purpose block coordinate descent algorithm (CGD) [5] to solve (12). 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 decent to generate a decent direction followed by an inexact line search along this direction [5]. For continuously differentiable  $f(\cdot)$  and convex and block-separable  $P(\cdot)$  (i.e.  $P(\beta) = \sum_i P_i(\beta_i)$ ), [5] 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. [6], [7]) and [4] for mixed models with an  $\ell_1$  penalty. Following Tseng and Yun [5], the CGD algorithm is given by Algorithm 1.

---

**Algorithm 1:** Coordinate Gradient Descent Algorithm to solve (12)

---

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

$$\Theta^{(0)};$$

**repeat**

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

$$H^{(k)} = \text{diag} \left[ \min \left\{ \max \left\{ \left[ \nabla^2 f(\Theta^{(k)}) \right]_{jj}, c_{min} \right\} c_{max} \right\} \right]_{j=1,\dots,p} \quad (13)$$

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

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

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

$$d_{H^{(k)}}(\Theta_j^{(k)}) \leftarrow \arg \min_d \left\{ \nabla f(\Theta_j^{(k)})d + \frac{1}{2}d^2 H_{jj}^{(k)} + \lambda P(\Theta_j^{(k)} + d) \right\} \quad (14)$$

**end**

**end**

Choose a stepsize;

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

Update;

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

Update;

$$\widehat{\eta}^{(k+1)} \leftarrow \arg \min_{\eta} \frac{1}{2} \sum_{i=1}^{N_T} \log(1 + \eta(\Lambda_i - 1)) + \frac{1}{2\sigma^2(k)} \sum_{i=1}^{N_T} \frac{\left( \widetilde{Y}_i - \sum_{j=0}^p \widetilde{X}_{ij+1} \beta_j^{(k+1)} \right)^2}{1 + \eta(\Lambda_i - 1)} \quad (15)$$

Update;

$$\widehat{\sigma^2}^{(k+1)} \leftarrow \frac{1}{N_T} \sum_{i=1}^{N_T} \frac{\left( \widetilde{Y}_i - \sum_{j=0}^p \widetilde{X}_{ij+1} \beta_j^{(k+1)} \right)^2}{1 + \eta^{(k+1)}(\Lambda_i - 1)} \quad (16)$$

$$k \leftarrow k + 1$$

**until** convergence criterion is satisfied;

---

The Armijo rule is defined as follows [5]:

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

$$Q_\lambda(\Theta_j^{(k)} + \alpha^{(k)} d^{(k)}) \leq Q_\lambda(\Theta_j^{(k)}) + \alpha^{(k)} \varrho \Delta^{(k)} \quad (17)$$

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

$$\Delta^{(k)} := \nabla f(\Theta_j^{(k)}) d^{(k)} + \gamma (d^{(k)})^2 H_{jj}^{(k)} + \lambda P(\Theta_j^{(k)} + d^{(k)}) - \lambda P(\Theta_j^{(k)}) \quad (18)$$

Common choices for the constants are  $\delta = 0.1$ ,  $\varrho = 0.001$ ,  $\gamma = 0$ ,  $\alpha_{init}^{(k)} = 1$  for all  $k$  [4].

Below we detail the specifics of Algorithm 1 for the  $\ell_1$  penalty.

## $\ell_1$ penalty

The objective function is given by

$$Q_\lambda(\Theta) = f(\Theta) + \lambda |\beta| \quad (19)$$

## 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

$$d_H(\Theta_j) = \arg \min_d G(d) \quad (20)$$

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

$$\partial G(d) = \nabla f(\Theta_j) + dH_{jj} + \lambda u \quad (21)$$

where

$$u = \begin{cases} 1 & \text{if } d > -\Theta_j \\ -1 & \text{if } d < -\Theta_j \\ [-1, 1] & \text{if } d = \Theta_j \end{cases} \quad (22)$$

We consider each of the three cases in (21) 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  $a, b, c$  [5].

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 (20). Therefore the descent direction for  $\Theta_j^{(k)} \in \{\beta_1, \dots, \beta_p\}$  for the  $\ell_1$  penalty is given by

$$d = \text{mid} \left\{ \frac{-(\nabla f(\beta_j) - \lambda)}{H_{jj}}, -\beta_j, \frac{-(\nabla f(\beta_j) + \lambda)}{H_{jj}} \right\} \quad (23)$$

### Solution for the $\beta$ parameter

If the Hessian  $\nabla^2 f(\Theta^{(k)}) > 0$  then  $H^{(k)}$  defined in (13) is equal to  $\nabla^2 f(\Theta^{(k)})$ . Using  $\alpha_{init} = 1$ , the largest element of  $\{\alpha_{init}^{(k)} \delta^r\}_{r=0,1,2,\dots}$  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

$$\beta_j^{(k+1)} \leftarrow \beta_j^{(k)} + d^{(k)}, \quad j = 1, \dots, p \quad (24)$$

Substituting the descent direction given by (23) into (24) we get

$$\beta_j^{(k+1)} = \text{mid} \left\{ \beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) - \lambda)}{H_{jj}}, 0, \beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) + \lambda)}{H_{jj}} \right\} \quad (25)$$

We can further simplify this expression. Let

$$w_i := \frac{1}{\sigma^2 (1 + \eta(\Lambda_i - 1))} \quad (26)$$

Re-write the part depending on  $\beta$  of the negative log-likelihood in (10) as

$$g(\beta^{(k)}) = \frac{1}{2} \sum_{i=1}^{N_T} w_i \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} - \tilde{X}_{ij} \beta_j^{(k)} \right)^2 \quad (27)$$

The gradient and Hessian are given by

$$\nabla f(\beta_j^{(k)}) := \frac{\partial}{\partial \beta_j^{(k)}} g(\beta^{(k)}) = - \sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} - \tilde{X}_{ij} \beta_j^{(k)} \right) \quad (28)$$

$$H_{jj} := \frac{\partial^2}{\partial \beta_j^{(k)2}} g(\beta^{(k)}) = \sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2 \quad (29)$$

Substituting (28) and (29) into  $\beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) - \lambda)}{H_{jj}}$

$$\begin{aligned}
& \beta_j^{(k)} + \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} - \tilde{X}_{ij} \beta_j^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \\
&= \beta_j^{(k)} + \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} - \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2 \beta_j^{(k)}}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \\
&= \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \tag{30}
\end{aligned}$$

Similarly, substituting (28) and (29) in  $\beta_j^{(k)} + \frac{-(\nabla f(\beta_j^{(k)}) + \lambda)}{H_{jj}}$  we get

$$\frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) - \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \tag{31}$$

Finally, substituting (30) and (31) into (25) we get

$$\begin{aligned}
\beta_j^{(k+1)} &= \text{mid} \left\{ \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) - \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2}, 0, \frac{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) + \lambda}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \right\} \\
&= \frac{\mathcal{S}_\lambda \left( \sum_{i=1}^{N_T} w_i \tilde{X}_{ij} \left( \tilde{Y}_i - \sum_{\ell \neq j} \tilde{X}_{i\ell} \beta_\ell^{(k)} \right) \right)}{\sum_{i=1}^{N_T} w_i \tilde{X}_{ij}^2} \tag{32}
\end{aligned}$$

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

# 1 Additional Real Data Analysis Results

## 1.1 Distribution of SNPs used in UK Biobank analysis

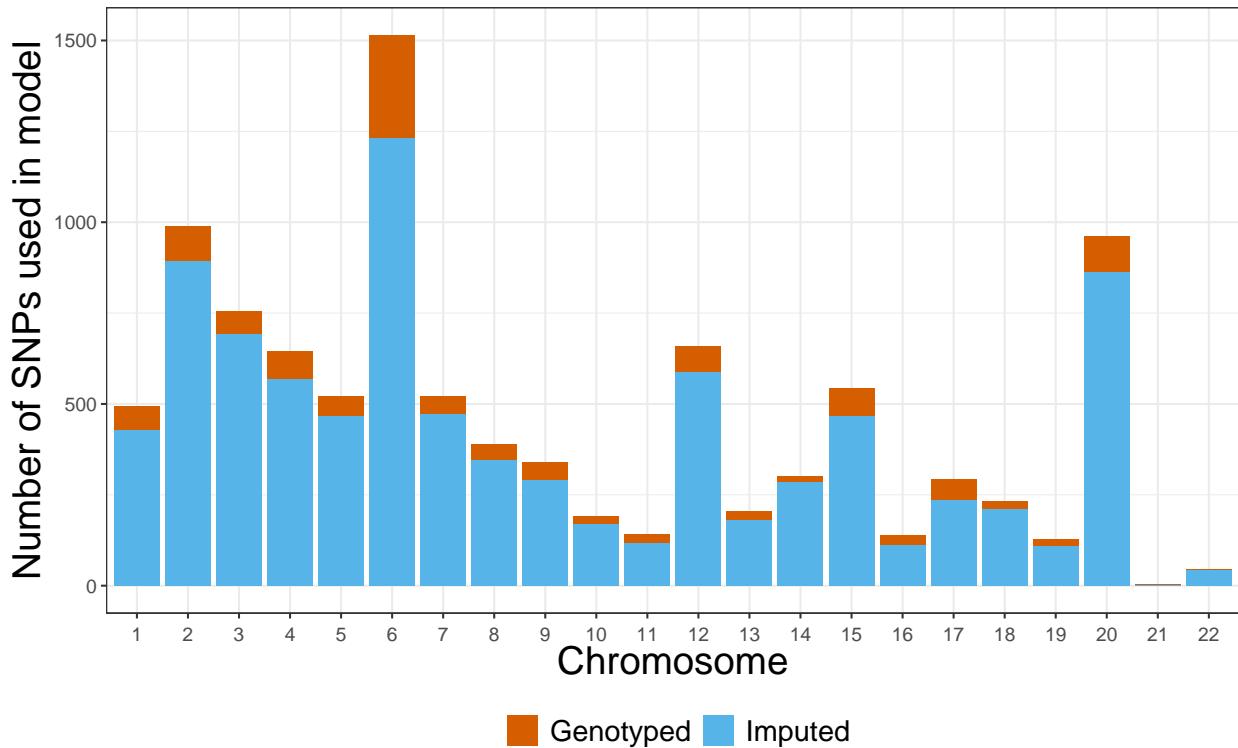


Figure 1.1: Distribution of SNPs used in UK Biobank analysis by chromosome and whether or not the SNP was imputed.

## 1.2 LD structure among the markers in the GAW20 and the mouse dataset

We illustrate the LD structure among the markers in the GAW20 dataset and the mouse dataset separately in Figures 1.2 and 1.3, respectively. In Figure 1.2, we show the pairwise  $r^2$  for 655 SNPs within a 1Mb-window around the causal SNP rs9661059 (indicated) that we focused on. The dotplot above the heatmap denotes  $r^2$  between each SNP and the causal SNP. It is clear that although strong correlation does exist between some SNPs, none of these nearby SNPs is correlated with the causal SNP. The only dot denoting an  $r^2 = 1$  represents the causal SNP itself.

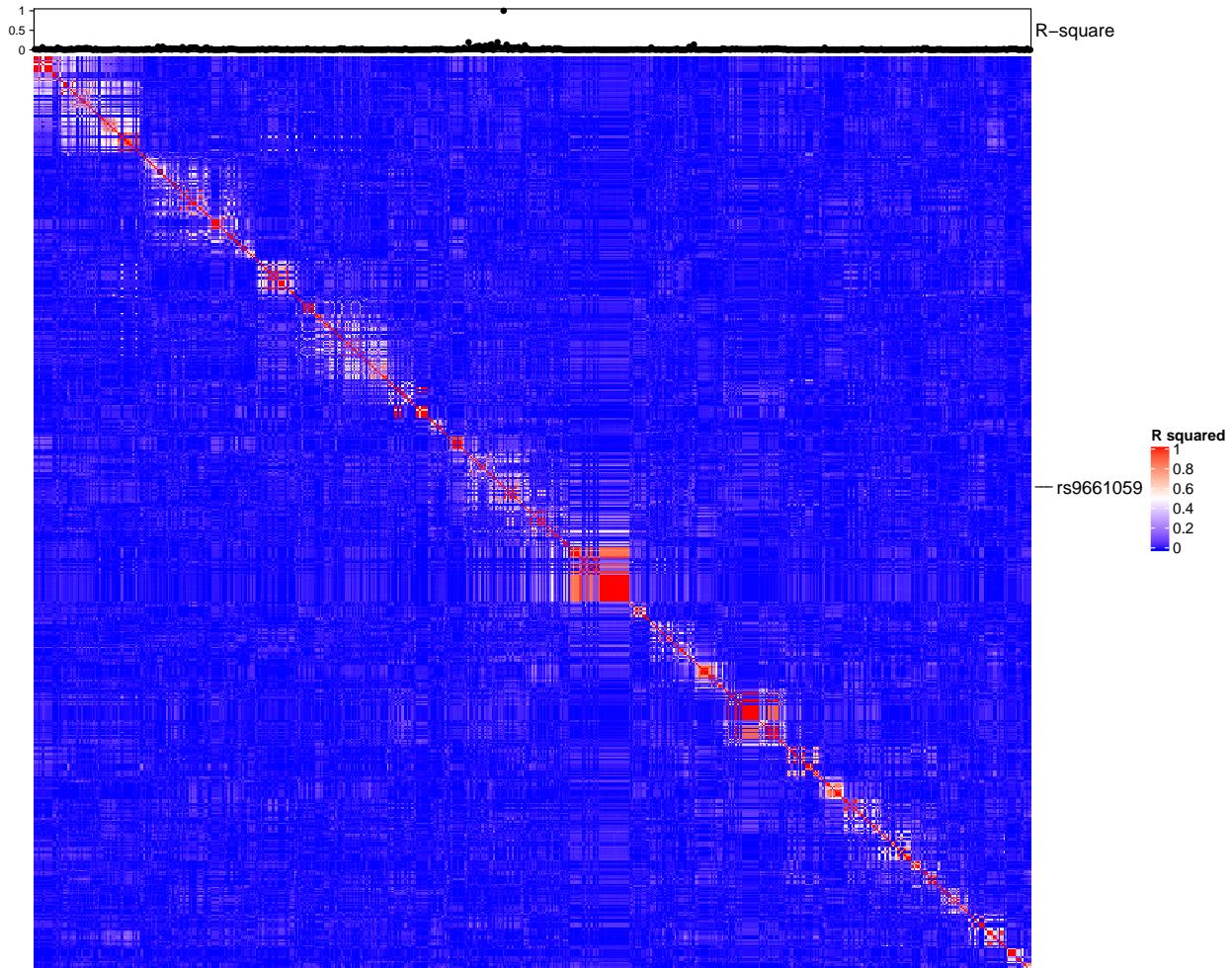


Figure 1.2: LD structure among the markers in the GAW20 dataset

In Figure 1.3, we show the pairwise  $r^2$  for all microsatellite markers in the mouse dataset. It is clear that many markers are considerably strongly correlated with each other, as we expected.

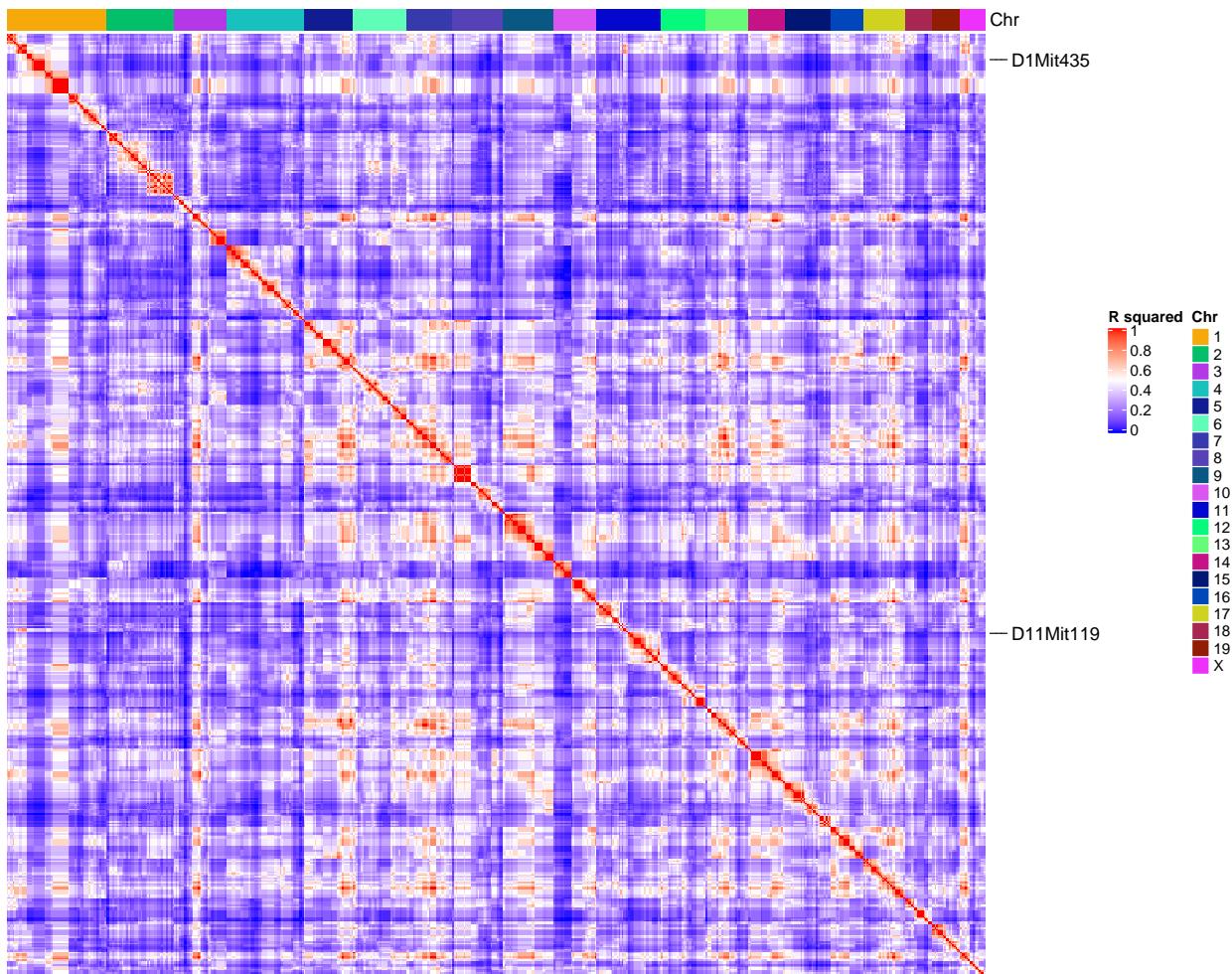


Figure 1.3: LD structure among the markers in the mouse dataset

## 2 ggmmix Package Showcase

In this section we briefly introduce the freely available and open source `ggmmix` package in R. More comprehensive documentation is available at <https://sahirbhatnagar.com/ggmmix>. Note that this entire section is reproducible; the code and text are combined in an `.Rnw`<sup>1</sup> file and compiled using `knitr` [8].

### 2.1 Installation

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

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

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

```
## library(ggmmix)
data("admixed")
names(admixed)

## [1] "ytrain"      "ytune"       "ytest"        "xtrain"
## [5] "xtune"       "xtest"       "xtrain_lasso" "xtune_lasso"
## [9] "xtest_lasso" "Xkinship"    "kin_train"    "kin_tune_train"
## [13] "kin_test_train" "mu_train"   "causal"       "beta"
## [17] "not_causal"   "kinship"    "coancestry"  "PC"
## [21] "subpops"
```

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

There are three basic inputs that `ggmmix` needs:

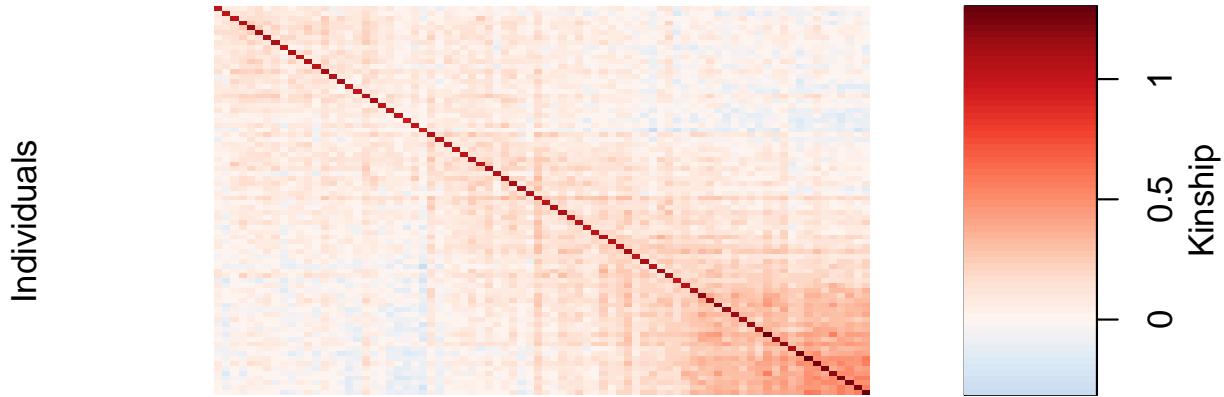
1.  $Y$ : a continuous response variable
2.  $X$ : a matrix of covariates of dimension  $N \times p$  where  $N$  is the sample size and  $p$  is the number of covariates
3.  $\Phi$ : a kinship matrix

---

<sup>1</sup>scripts available at <https://github.com/sahirbhatnagar/ggmmix/tree/pgen/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::plot_popkin(admixed$kin_train)
```



## 2.2 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  $\lambda$ . It will choose its own sequence:

```
fit <- ggmix(x = admixed$xtrain,
              y = admixed$ytrain,
              kinship = admixed$kin_train)

names(fit)

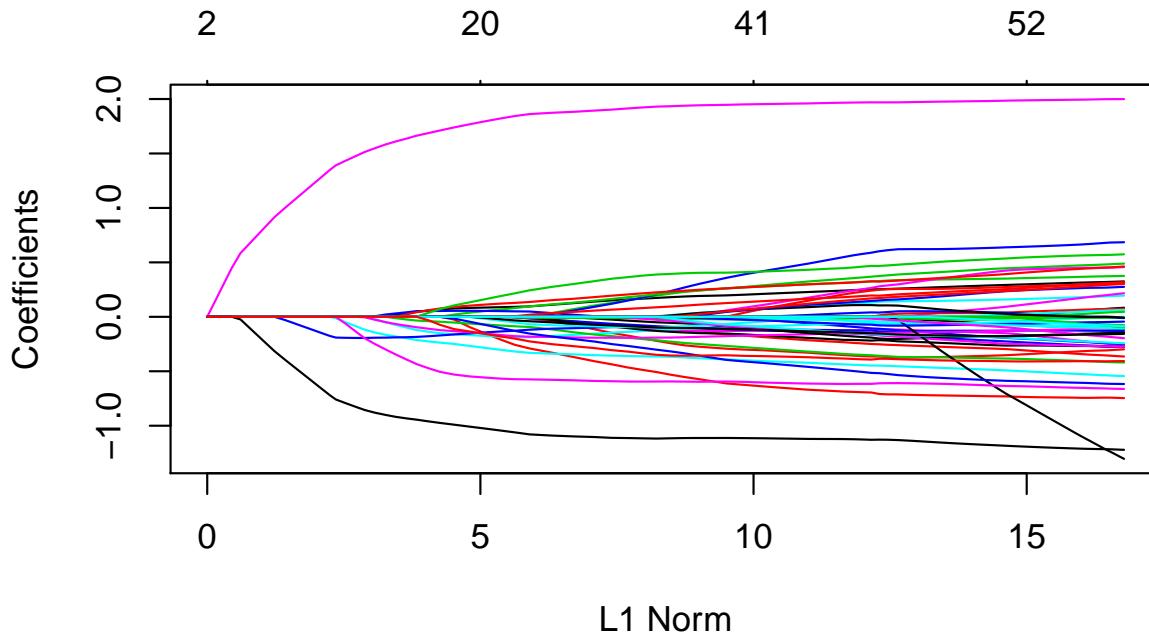
## [1] "result"      "ggmix_object" "n_design"    "p_design"    "lambda"
## [6] "coef"        "b0"          "beta"        "df"         "eta"
## [11] "sigma2"       "nlambda"     "cov_names"   "call"

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 lambda 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) -0.03715135  0.247105426
## X23        0.00000000  0.098030248
## X36        0.00000000 -0.013022250
## X38        0.00000000  0.005378361
## X40        0.00000000  0.004028934
```

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

interpolation to make predictions for values of  $s$  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$xtest[1:5,])

##          1          2
## id26  2.30208546  2.45597763
## id39  0.87334032  1.62931898
## id45 -0.12296837 -0.06075786
## id52 -0.03715135 -0.97519671
## id53 -0.21046107 -0.23151040
```

### 2.3 Find the Optimal Value of the Tuning Parameter

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

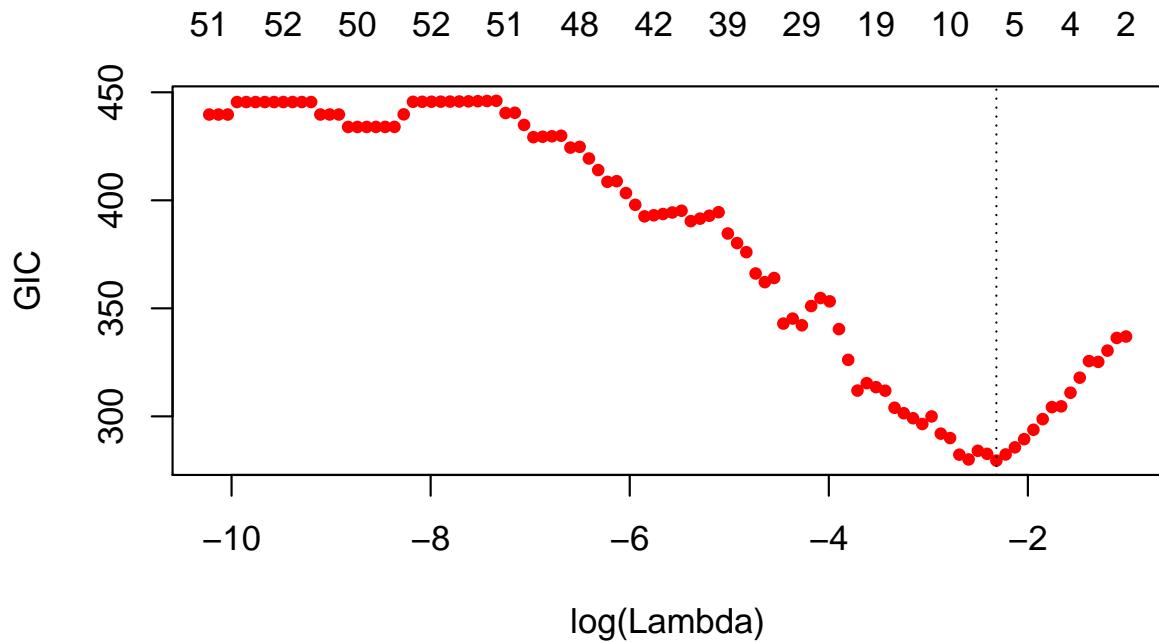
```
# pass the fitted object from ggmix to the gic function:
hdbic <- gic(fit)
class(hdbic)

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

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

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

```
plot(hdbic)
```

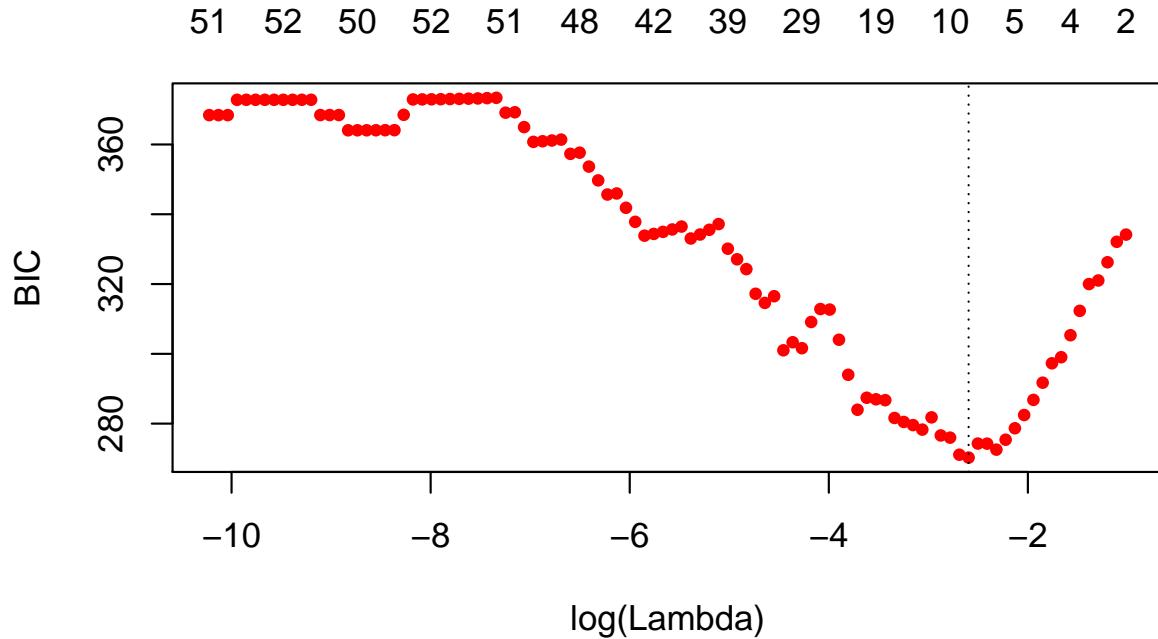


The optimal value for  $\lambda$  according to the HDBIC, i.e., the  $\lambda$  that leads to the minium HDBIC is:

```
hdbic[["lambda.min"]]
## [1] 0.09862269
```

We can also plot the BIC results:

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



```
bicfit[["lambda.min"]]
## [1] 0.07460445
```

## 2.4 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) -0.03660806
## X23         .
## X36         .
## X38         .
## X40         .
```

We can also extract just the nonzero coefficients which also provide the estimated variance components  $\eta$  and  $\sigma^2$ :

```
coef(hdbic, type = "nonzero")

##          1
## (Intercept) -0.03660806
## X302       -0.17607392
## X524        1.34951500
## X538       -0.72052613
## eta         0.99000000
## sigma2      1.60476289
```

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$xtest[1:5,])

##          1
## id26  2.31027410
## id39  0.86922183
## id45 -0.12814532
## id52 -0.03660806
## id53 -0.21268198
```

## 2.5 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]

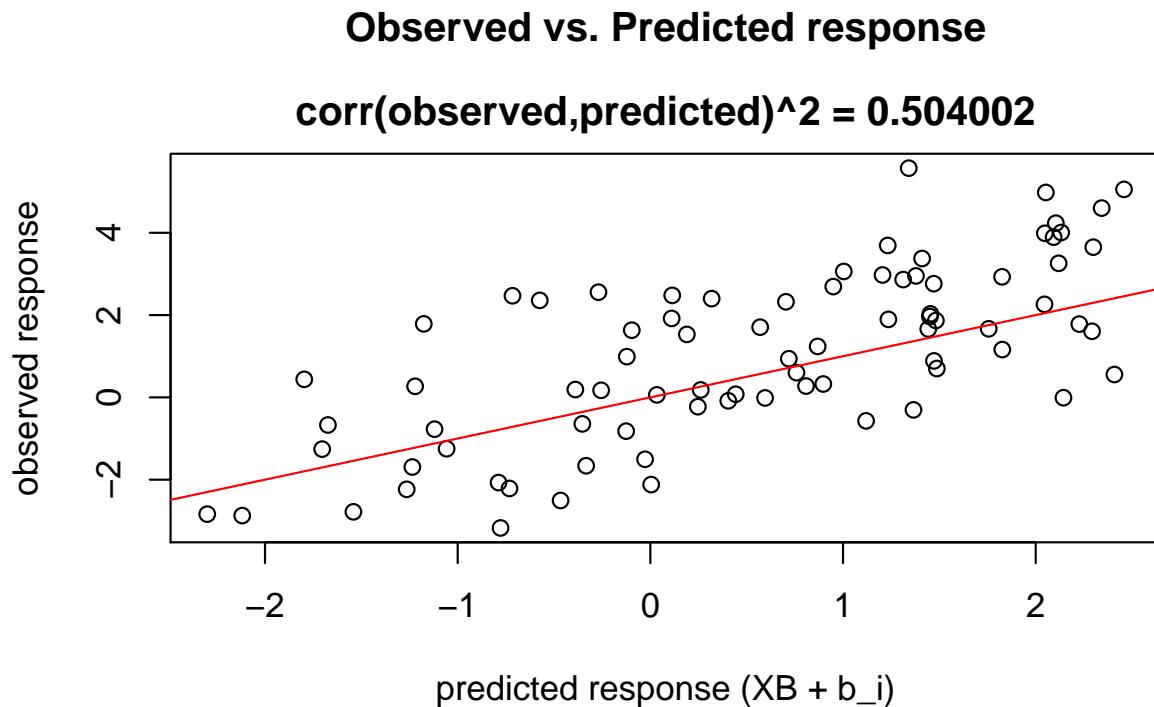
## [1] -2.4889655  1.1834200 -0.5641832 -0.9310334 -0.3458703
```

## 2.6 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.

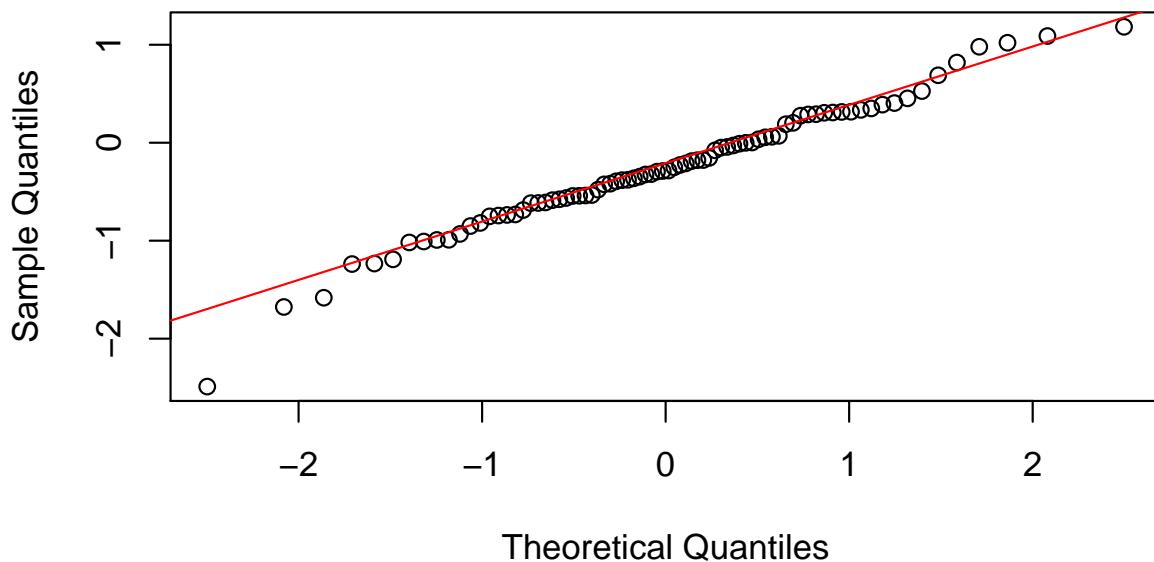
### 2.6.1 Observed vs. Predicted Response

```
plot(hdbic, type = "predicted", newx = admixed$xtrain, newy = admixed$ytrain)
```

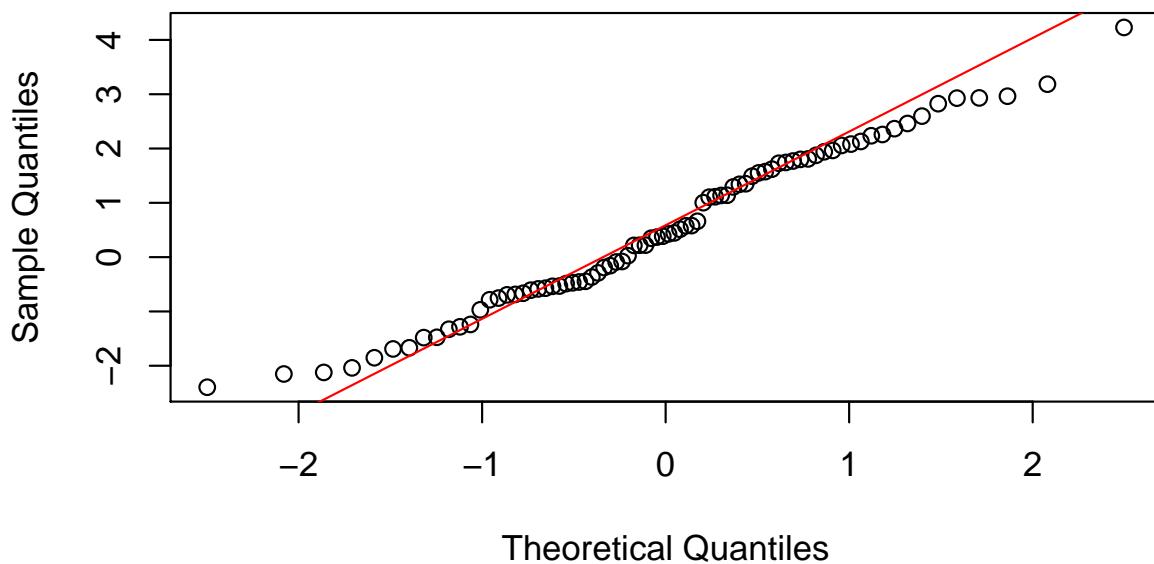


### 2.6.2 QQ-plots for Residuals and Random Effects

```
plot(hdbic, type = "QQranef", newx = admixed$xtrain, newy = admixed$ytrain)
```

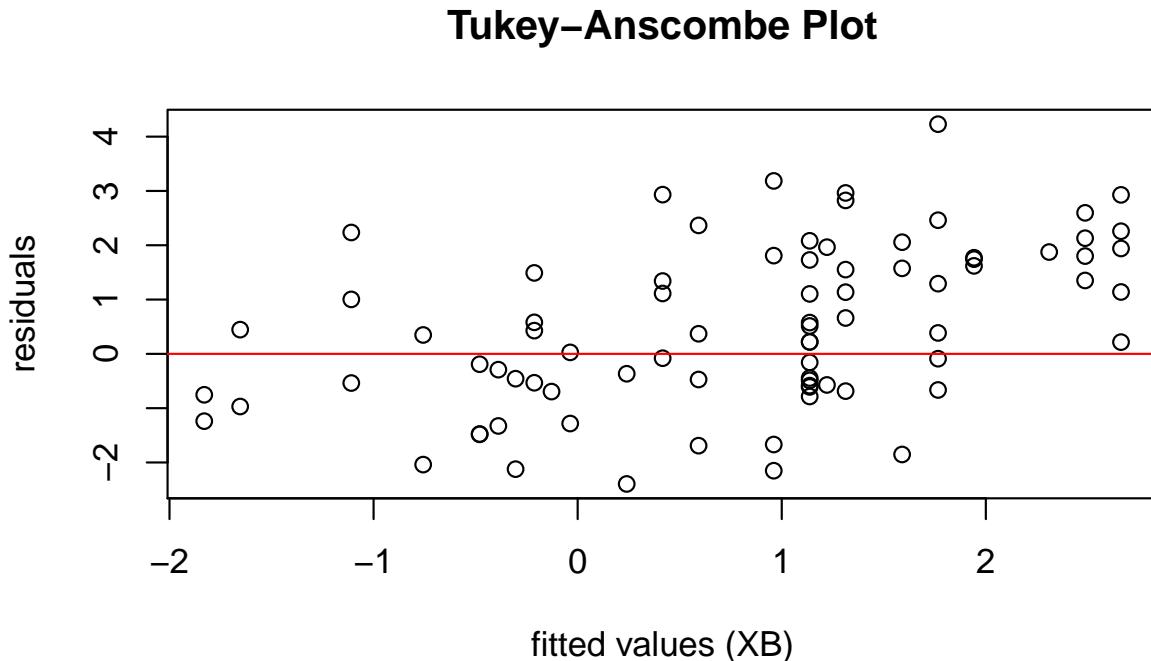
**QQ-Plot of the random effects at lambda = 0.10**

```
plot(hdbic, type = "QQresid", newx = admixed$xtrain, newy = admixed$ytrain)
```

**QQ-Plot of the residuals at lambda = 0.10**

### 2.6.3 Tukey-Anscombe Plot

```
plot(hdbic, type = "Tukey", newx = admixed$xtrain, newy = admixed$ytrain)
```



## References

- [1] Pirinen M, Donnelly P, Spencer CC, et al. Efficient computation with a linear mixed model on large-scale data sets with applications to genetic studies. *The Annals of Applied Statistics*. 2013;7(1):369–390. [3](#), [4](#)
- [2] Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorff LA, Hunter DJ, et al. Finding the missing heritability of complex diseases. *Nature*. 2009;461(7265):747. [3](#)
- [3] Rakitsch B, Lippert C, Stegle O, Borgwardt K. A Lasso multi-marker mixed model for association mapping with population structure correction. *Bioinformatics*. 2013;29(2):206–214.
- [4] Schelldorfer J, Bühlmann P, DE G, VAN S. Estimation for High-Dimensional Lin-

- ear Mixed-Effects Models Using L1-Penalization. Scandinavian Journal of Statistics. 2011;38(2):197–214. [6](#), [8](#)
- [5] Tseng P, Yun S. A coordinate gradient descent method for nonsmooth separable minimization. Mathematical Programming. 2009;117(1):387–423. [6](#), [9](#)
- [6] Meier L, Van De Geer S, Bühlmann P. The group lasso for logistic regression. Journal of the Royal Statistical Society: Series B (Statistical Methodology). 2008;70(1):53–71. [6](#)
- [7] Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear models via coordinate descent. Journal of statistical software. 2010;33(1):1. [6](#)
- [8] Xie Y. Dynamic Documents with R and knitr. vol. 29. CRC Press; 2015. [17](#)