

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

Joint work with

Sahir Bhatnagar (McGill), Yi Yang (McGill), Celia Greenwood (McGill)

BIRS Workshop 2018

Motivation

Genetic Analysis Workshop (GAW20, March 4-7, 2017, San Diego, US)



Genetic Analysis Workshop



[Home](#)

[About](#)

[GAW20](#)

[Register](#)

[Related Links](#)

[Contact](#)

GAW20: DATA SETS

Epigenetic and Pharmacogenomic Data

The data set for GAW20 draws on themes of pharmacogenomics and epigenetics, some of the most requested topics in a 2015 survey of the GAW mailing list. The GAW20 'real' data set includes metabolic syndrome diagnoses and HDL and triglyceride levels before and after treatment with fenofibrate as well as genome-wide methylation pre- and post-treatment and dense genome-wide SNPs from the [GOLDN project](#). For more detail on

¹GOLDEN project: Genetics of Lipid Lowering Drugs and Diet Network Study

Motivation

- ▶ Our contribution in GAW20

Investigating potential causal relationships between SNPs, DNA methylation and HDL

Lai Jiang^{1,2}, Kaiqiong Zhao^{1,2}, Kathleen Klein², Angelo J Canty⁵,
Karim Oualkacha³, Celia MT Greenwood^{*1,2,4}

- ▶ Our contribution in GAW20

Investigating potential causal relationships between SNPs, DNA methylation and HDL

Lai Jiang^{1,2}, Kaiqiong Zhao^{1,2}, Kathleen Klein², Angelo J Canty⁵,
Karim Oualkacha³, Celia MT Greenwood^{*1,2,4}

- ▶ Contribution of the Causal modelling group

Causal modeling in a multi-omics setting: insights from Genetic Analysis Workshop 20

Jonathan Auerbach*, Richard Howey*, Lai Jiang*, Anne Justice*, Liming Li*, Karim Oualkacha*,
Sergi Sayols-Baixeras*, Stella W. Aslibekyan†

*Contributed equally; listed in alphabetical order

Motivation

- ▶ Our contribution in GAW20 consisted of investigation of causal relationship between DNA methylation (exposure) within some genes and Δ HDL (outcome)

Motivation

- ▶ Our contribution in GAW20 consisted of investigation of causal relationship between DNA methylation (exposure) within some genes and Δ HDL (outcome)
- ▶ DNA methylation in these genes has been shown association with HDL

Motivation

- ▶ Our contribution in GAW20 consisted of investigation of causal relationship between DNA methylation (exposure) within some genes and Δ HDL (outcome)
- ▶ DNA methylation in these genes has been shown association with HDL
- ▶ We used Mendelian randomization to explore causal relationship
- ▶ We used SNPs around the analyzed genes as Instrumental Variables (IVs) to interrogate the causal relationship

Motivation

- ▶ Our contribution in GAW20 consisted of investigation of causal relationship between DNA methylation (exposure) within some genes and Δ HDL (outcome)
- ▶ DNA methylation in these genes has been shown association with HDL
- ▶ We used Mendelian randomization to explore causal relationship
- ▶ We used SNPs around the analyzed genes as Instrumental Variables (IVs) to interrogate the causal relationship



Challenges in GAW20 Data Sets

- ▶ GAW20 SNPs data was high-dimensional
- ▶ There was a need for data regularization in order to select SNPs strongly associated with the exposure
- ▶ Penalized LS regression can be used (Lasso, SCAD, MCP or Elastic net)

Challenges in GAW20 Data Sets

- ▶ GAW20 SNPs data was high-dimensional
- ▶ There was a need for data regularization in order to select SNPs strongly associated with the exposure
- ▶ Penalized LS regression can be used (Lasso, SCAD, MCP or Elastic net)
- ▶ But, data consists of families !
- ▶ In the GAW20, all regularized methods
 - ▶ either did not control for the family structure

Challenges in GAW20 Data Sets

- ▶ GAW20 SNPs data was high-dimensional
- ▶ There was a need for data regularization in order to select SNPs strongly associated with the exposure
- ▶ Penalized LS regression can be used (Lasso, SCAD, MCP or Elastic net)
- ▶ But, data consists of families !
- ▶ In the GAW20, all regularized methods
 - ▶ either did not control for the family structure
 - ▶ or used two-steps adjustment for the family structure (including our group)

Challenges in GAW20 Data Sets

- ▶ Two-steps adjustment:
 - ▶ Step 1 : uses LMM to adjust for subjects relationship

¹Oualkacha et al. Gene. Epi. (2013)

Challenges in GAW20 Data Sets

- ▶ Two-steps adjustment:
 - ▶ Step 1 : uses LMM to adjust for subjects relationship
 - ▶ Step 2 : uses residuals from Step 1 in variable-selection
LS-regression methods to select SNPs

¹Oualkacha et al. Gene. Epi. (2013)

Challenges in GAW20 Data Sets

- ▶ Two-steps adjustment:
 - ▶ Step 1 : uses LMM to adjust for subjects relationship
 - ▶ Step 2 : uses residuals from Step 1 in variable-selection LS-regression methods to select SNPs
- ▶ Two-steps procedure is a valid approach
- ▶ In association testing, (GRAMMAR) it is known to suffer from huge power loss ¹

¹Oualkacha et al. Gene. Epi. (2013)

Proposal

Aim:

We believe that performing variable selection and controlling for familial and/or hidden relationships simultaneously in high-dimensional settings, are likely to be of great interest to the genetic scientists community

Proposal

Aim:

We believe that performing variable selection and controlling for familial and/or hidden relationships simultaneously in high-dimensional settings, are likely to be of great interest to the genetic scientists community

Proposal:

We propose, `ggmix`, a two-in-one procedure which controls for structured populations and performs variable selection in Linear Mixed Models

Data and Model

- ▶ Phenotype: $\mathbf{Y} = (y_1, \dots, y_n) \in \mathbb{R}^n$
- ▶ SNPs: $\mathbf{X} = (\mathbf{X}_1; \dots; \mathbf{X}_n)^T \in \mathbb{R}^{n \times p}$, where $p \gg n$
- ▶ Twice the Kinship matrix or Realized Relationship matrix:
 $\Phi \in \mathbb{R}^{n \times n}$
- ▶ Regression Coefficients: $\beta = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$
- ▶ Polygenic random effect: $\mathbf{P} = (P_1, \dots, P_n) \in \mathbb{R}^n$
- ▶ Error: $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n) \in \mathbb{R}^n$

Data and Model

- ▶ Phenotype: $\mathbf{Y} = (y_1, \dots, y_n) \in \mathbb{R}^n$
- ▶ SNPs: $\mathbf{X} = (\mathbf{X}_1; \dots, \mathbf{X}_n)^T \in \mathbb{R}^{n \times p}$, where $p \gg n$
- ▶ Twice the Kinship matrix or Realized Relationship matrix:
 $\Phi \in \mathbb{R}^{n \times n}$
- ▶ Regression Coefficients: $\beta = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$
- ▶ Polygenic random effect: $\mathbf{P} = (P_1, \dots, P_n) \in \mathbb{R}^n$
- ▶ Error: $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n) \in \mathbb{R}^n$
- ▶ We consider the following LMM with a single random effect:

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{P} + \varepsilon$$

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

- ▶ σ^2 is the phenotype total variance
- ▶ $\eta \in [0, 1]$ is the phenotype heritability (narrow sense)
- ▶ $\mathbf{Y} | (\beta, \eta, \sigma^2) \sim \mathcal{N}(\mathbf{X}\beta, \eta\sigma^2\Phi + (1 - \eta)\sigma^2\mathbf{I})$

Likelihood

- ▶ The negative log-likelihood is given by

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

$$\mathbf{V} = \eta\Phi + (1 - \eta)\mathcal{I}$$

Likelihood

- ▶ The negative log-likelihood is given by

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

$$\mathbf{V} = \eta \Phi + (1 - \eta) \mathcal{I}$$

- ▶ Assume the spectral decomposition of Φ

$$\Phi = \mathbf{U} \mathbf{D} \mathbf{U}^T$$

- ▶ \mathbf{U} is an $n \times n$ orthogonal matrix and \mathbf{D} is an $n \times n$ diagonal matrix

Likelihood

- ▶ The negative log-likelihood is given by

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

$$\mathbf{V} = \eta \Phi + (1 - \eta) \mathcal{I}$$

- ▶ Assume the spectral decomposition of Φ

$$\Phi = \mathbf{U} \mathbf{D} \mathbf{U}^T$$

- ▶ \mathbf{U} is an $n \times n$ orthogonal matrix and \mathbf{D} is an $n \times n$ diagonal matrix
- ▶ One can write

$$\mathbf{V} = \mathbf{U}(\eta \mathbf{D} + (1 - \eta) \mathcal{I}) \mathbf{U}^T = \mathbf{U} \mathbf{W} \mathbf{U}^T$$

Likelihood

- ▶ The negative log-likelihood is given by

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

$$\mathbf{V} = \eta \Phi + (1 - \eta) \mathcal{I}$$

- ▶ Assume the spectral decomposition of Φ

$$\Phi = \mathbf{U} \mathbf{D} \mathbf{U}^T$$

- ▶ \mathbf{U} is an $n \times n$ orthogonal matrix and \mathbf{D} is an $n \times n$ diagonal matrix
- ▶ One can write

$$\mathbf{V} = \mathbf{U}(\eta \mathbf{D} + (1 - \eta) \mathcal{I}) \mathbf{U}^T = \mathbf{U} \mathbf{W} \mathbf{U}^T$$

with $\mathbf{W} = \text{diag}(w_i)_{i=1}^n$, $w_i = \eta \mathbf{D}_{ii} + (1 - \eta)$

Likelihood

- ▶ Projection of \mathbf{Y} (and columns of \mathbf{X}) into $\text{Span}(\mathbf{U})$ leads to a simplified correlation structure for the transformed data:
 $\tilde{\mathbf{Y}} = \mathbf{U}^\top \mathbf{Y}$
- ▶ $\tilde{\mathbf{Y}} | (\boldsymbol{\beta}, \eta, \sigma^2) \sim \mathcal{N}(\tilde{\mathbf{X}}\boldsymbol{\beta}, \sigma^2 \mathbf{W})$, with $\tilde{\mathbf{X}} = \mathbf{U}^\top \mathbf{X}$

Likelihood

- ▶ Projection of \mathbf{Y} (and columns of \mathbf{X}) into $\text{Span}(\mathbf{U})$ leads to a simplified correlation structure for the transformed data:

$$\tilde{\mathbf{Y}} = \mathbf{U}^T \mathbf{Y}$$

- ▶ $\tilde{\mathbf{Y}} | (\boldsymbol{\beta}, \eta, \sigma^2) \sim \mathcal{N}(\tilde{\mathbf{X}}\boldsymbol{\beta}, \sigma^2 \mathbf{W})$, with $\tilde{\mathbf{X}} = \mathbf{U}^T \mathbf{X}$
- ▶ The negative log-likelihood can then be expressed as

$$-\ell(\boldsymbol{\Theta}) \propto \frac{n}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^n \log(w_i) + \frac{1}{2\sigma^2} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta})^T \mathbf{W}^{-1} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta})$$

Likelihood

- ▶ Projection of \mathbf{Y} (and columns of \mathbf{X}) into $\text{Span}(\mathbf{U})$ leads to a simplified correlation structure for the transformed data:

$$\tilde{\mathbf{Y}} = \mathbf{U}^T \mathbf{Y}$$

- ▶ $\tilde{\mathbf{Y}} | (\boldsymbol{\beta}, \eta, \sigma^2) \sim \mathcal{N}(\tilde{\mathbf{X}}\boldsymbol{\beta}, \sigma^2 \mathbf{W})$, with $\tilde{\mathbf{X}} = \mathbf{U}^T \mathbf{X}$

- ▶ The negative log-likelihood can then be expressed as

$$-\ell(\boldsymbol{\Theta}) \propto \frac{n}{2} \log(\sigma^2) + \frac{1}{2} \sum_{i=1}^n \log(w_i) + \frac{1}{2\sigma^2} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta})^T \mathbf{W}^{-1} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta})$$

- ▶ For fixed σ^2 and η , solving for $\boldsymbol{\beta}$ is a weighted least squares problem

Penalized Maximum Likelihood Estimator

- ▶ Define the objective function:

$$Q_\lambda(\Theta) = -\ell(\Theta) + \lambda \sum_j p_j(\beta_j)$$

- ▶ $p_j(\cdot)$ is a penalty term on β_1, \dots, β_p
- ▶ An estimate of the model parameters $\hat{\Theta}_\lambda$ is obtained by

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

Block Relaxation (De Leeuw, 1994)

To solve for the optimization problem we use a block relaxation technique

Set $k \leftarrow 0$, initial values for the parameter vector $\Theta^{(0)}$ and ϵ ;

for $\lambda \in \{\lambda_{max}, \dots, \lambda_{min}\}$ **do**

repeat

$$\text{For } j = 1, \dots, p, \beta_j^{(k+1)} \leftarrow \arg \min_{\beta_j} Q_\lambda \left(\beta_{-j}^{(k)}, \eta^{(k)}, \sigma^2^{(k)} \right)$$

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

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

$$k \leftarrow k + 1$$

until *convergence criterion is satisfied:*

$$\|\Theta^{(k+1)} - \Theta^{(k)}\|_2 < \epsilon;$$

end

Algorithm 1: Block Relaxation Algorithm

Coordinate Gradient Descent Method

- ▶ We take advantage of smoothness of $\ell(\Theta)$
- ▶ We approximate $Q_\lambda(\Theta)$ by a strictly convex quadratic function (using gradient)
- ▶ We use CGD to calculate a descent direction
- ▶ To achieve the descent property for the objective function, we employ further line search

¹Tseng P& Yun S. Math. Program., Ser. B, (2009)

Coordinate Gradient Descent Method

- ▶ We take advantage of smoothness of $\ell(\Theta)$
- ▶ We approximate $Q_\lambda(\Theta)$ by a strictly convex quadratic function (using gradient)
- ▶ We use CGD to calculate a descent direction
- ▶ To achieve the descent property for the objective function, we employ further line search

Theorem [Convergence] ¹:

If $\{\Theta^{(k)}, k = 0, 1, 2, \dots\}$ is a sequence of iterates generated by the iteration map of Algorithm 1, then each cluster point (i.e. limit point) of $\{\Theta^{(k)}, k = 0, 1, 2, \dots\}$ is a stationary point of $Q_\lambda(\Theta)$

¹Tseng P& Yun S. Math. Program., Ser. B, (2009)

Choice of the tuning parameter

- ▶ We use the BIC:

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

- ▶ \widehat{df}_λ is the number of non-zero elements in $\hat{\beta}_\lambda$ plus two ¹
- ▶ Several authors ² have used this criterion for variable selection in mixed models with $c = \log n$
- ▶ Other authors ³ have proposed $c = \log(\log(n)) * \log(n)$

¹Zou et al. The Annals of Statistics, (2007)

²Bondell et al. Biometrics (2010)

³Wang et al. JRSS(Ser. B), (2009)

Simulation study

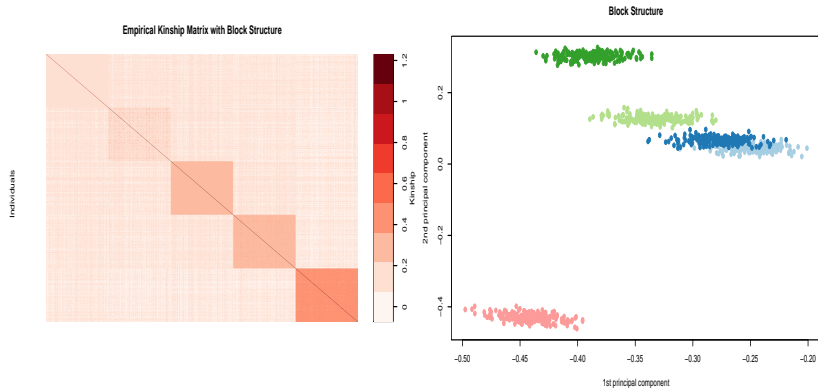
- ▶ We simulate genotypes from the BN-PSD Admixture Model
- ▶ a : percentage of causal SNPs
- ▶ $\mathbf{X}^{(test)}$: $n \times 5000$ matrix of SNPs randomly sampled across the genome
- ▶ $\mathbf{X}^{(causal)}$: $n \times (a * 5000)$ matrix of SNPs that are truly associated with the simulated phenotype, $\mathbf{X}^{(causal)} \subseteq \mathbf{X}^{(test)}$
- ▶ β_j : effect size for the j^{th} SNP, simulated from a $Uniform(0.3, 0.7)$ for $j = 1, \dots, (a * 5000)$
- ▶ $\mathbf{Y} | (\beta, \eta, \sigma^2) \sim \mathcal{N}(\mathbf{X}^{(causal)}\beta, \eta\sigma^2\Phi + (1 - \eta)\sigma^2\mathbf{I})$

¹<https://cran.r-project.org/package=bnpsd>

RRM/Kinship matrix construction

- ▶ $\mathbf{X}^{(other)}$: $n \times 10,000$ matrix of simulated SNPs
- ▶ $\mathbf{X}^{(kinship)}$: matrix of SNPs used to construct the RRM/Kinship matrix
 - ▶ Scenario 1: $\mathbf{X}^{(kinship)} = \mathbf{X}^{(other)} \leftarrow$ No overlap
 - ▶ Scenario 2: $\mathbf{X}^{(kinship)} = [\mathbf{X}^{(other)}, \mathbf{X}^{(causal)}] \leftarrow$ 100% overlap
- ▶ In each scenario we considered $a = 0, 0.01, \eta = 0.1, 0.5$ and $\sigma^2 = 1$

Empirical Kinship Matrix



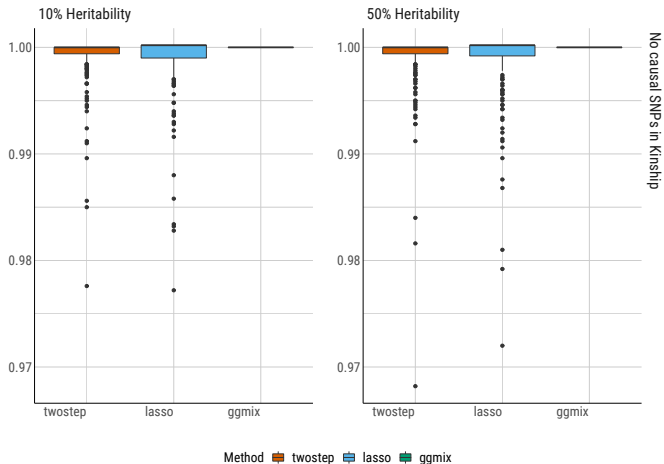
(a) Kinship Matrix

(b) 1st and 2nd PC

Correct Sparsity for Null Model

Correct Sparsity Results for the Null Model

Based on 200 simulations

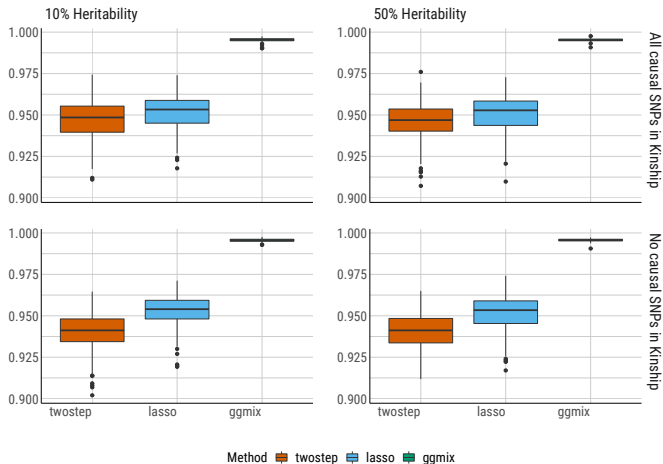


$\eta = 10\%$

Correct Sparsity for Model with 1% Causal SNPs

Correct Sparsity results for the Model with 1% Causal SNPs

Based on 200 simulations

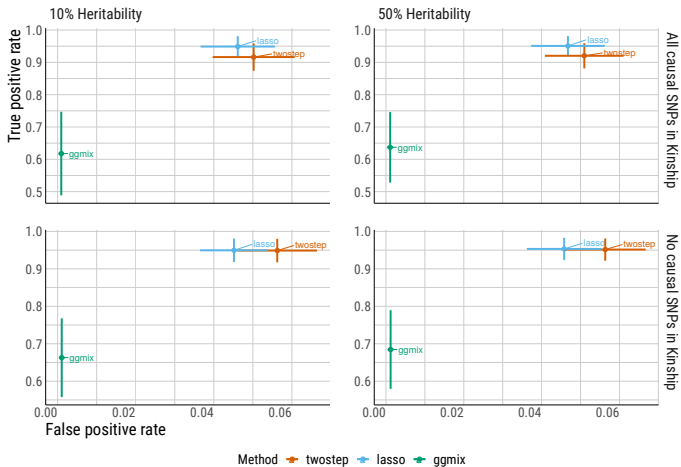


$\eta = 10\%$

True Positive vs. False Positive Rate

True Positive Rate vs. False Positive Rate (Mean \pm 1 SD) for the Model with 1% Causal SNPs

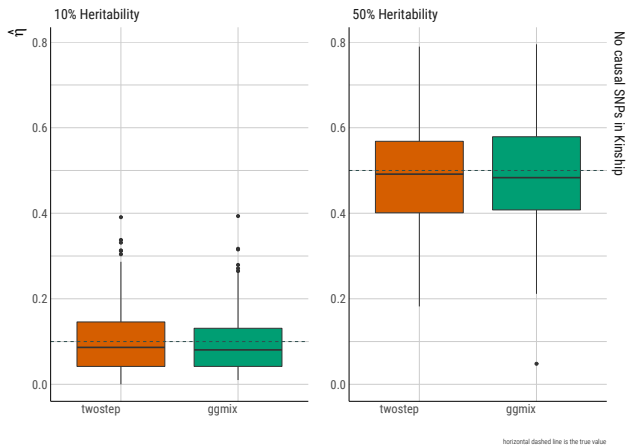
Based on 200 simulations



Heritability Estimates for Null Model

Estimated Heritability for the Null Model

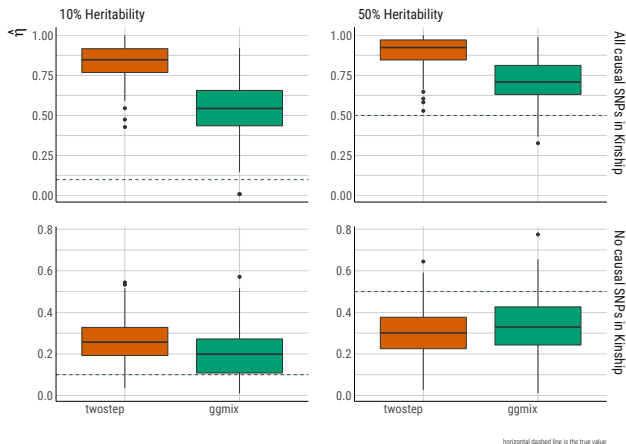
Based on 200 simulations



Heritability Estimates for Model with 1% Causal SNPs

Estimated Heritability for the Model with 1% Causal SNPs

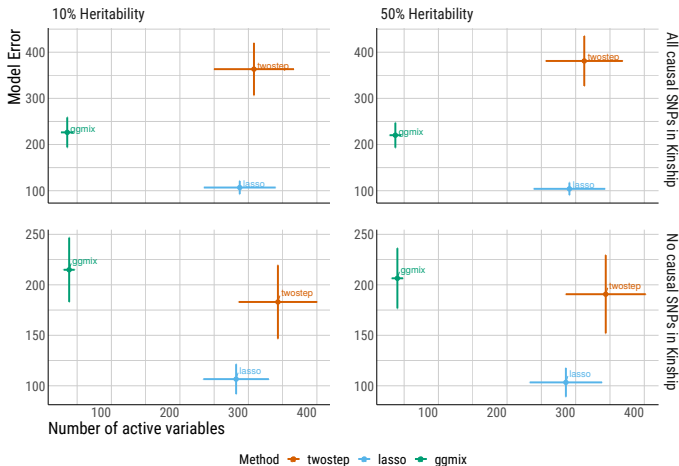
Based on 200 simulations



Model Error vs. Number Active for Model with 1% Causal SNPs

Model Error vs. Number of Active Variable (Mean +/- 1 SD) for Model with 1% Causal SNPs

Based on 200 simulations



Discussion/Future work

- ▶ We have presented a regularized mixed model, *ggmix*, for variable-selecting in the presence of confounding influences
- ▶ In some situations, prior information of the predictors (e.g. SNPs) groups structure is available

Discussion/Future work

- ▶ We have presented a regularized mixed model, *ggmix*, for variable-selecting in the presence of confounding influences
- ▶ In some situations, prior information of the predictors (e.g. SNPs) groups structure is available
- ▶ Theoretical development of group-Lasso in LMM is already done

Discussion/Future work

- ▶ We have presented a regularized mixed model, *ggmix*, for variable-selecting in the presence of confounding influences
- ▶ In some situations, prior information of the predictors (e.g. SNPs) groups structure is available
- ▶ Theoretical development of group-Lasso in LMM is already done

Discussion/Future work

- ▶ Capturing the subjects relationship using random effect requires VCs estimation
- ▶ Random effect modelling leads to a non-convex optimization problem
- ▶ Fixed effects models are good alternatives to random effects models for analysis of Longitudinal/Panel data ¹
- ▶ Capturing familial structure using a penalized FE model could be an interesting avenue to explore

¹Roger Koenker, JMA, (2004)

References

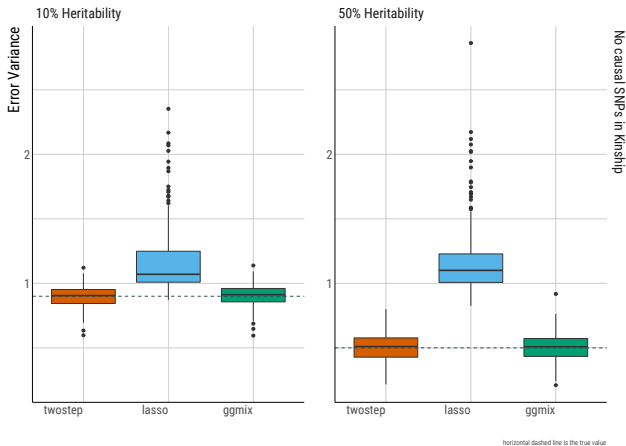
1. Christoph Lippert, Jennifer Listgarten, Ying Liu, Carl M Kadie, Robert I Davidson, and David Heckerman. Fast linear mixed models for genome-wide association studies. *Nature methods*, 8(10):833–835, 2011.
2. Matti Pirinen, Peter Donnelly, Chris CA Spencer, et al. Efficient computation with a linear mixed model on large-scale data sets with applications to genetic studies. *The Annals of Applied Statistics*, 7(1):369–390, 2013.
3. Paul Tseng and Sangwoon Yun. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*, 117(1):387–423, 2009.
4. Jerome Friedman, Trevor Hastie, and Rob Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1):1, 2010.

<http://sahirbhatnagar.com/ggmix/articles/introduction-to-ggmix.html>

Error Variance for Null Model

Estimated Error Variance for the Null Model

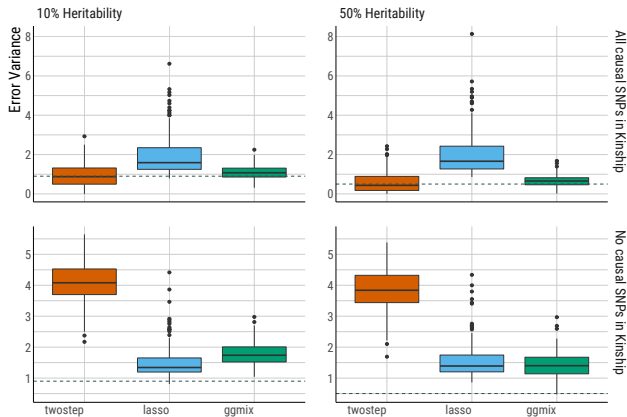
Based on 200 simulations



Error Variance for Model with 1% Causal SNPs

Estimated Error Variance for the Model with 1% Causal SNPs

Based on 200 simulations



horizontal dashed line is the true value