

A Sparse Group Lasso multi-marker mixed model for Genome-Wide association studies

Yingjie GUO, Chenxi WU*, Ao LI, Xiaoyan LIU, Alon KEINAN, Maozu GUO

April 24, 2017

Abstract

1 Introduction

2 Materials and Methods

2.1 Genotypes and phenotypes

2.2

2.3 Multivariate Linear mixed model

Our approach builds on linear mixed model, explaining the phenotype variability by a sum of individual genetic effects and random confounding effects. Suppose that m individuals of a phenotype are collected, a linear mixed model in association mapping is typically expressed as

$$\mathbf{y} = \underbrace{\sum_{j=1}^n \beta_j \mathbf{x}_j}_{\text{Fixed effect}} + \underbrace{\mathbf{Z}\mathbf{u}}_{\text{confounding}} + \underbrace{\Psi}_{\text{noise}} \quad (1)$$

Random effect

Where $\mathbf{y} = (y_1, \dots, y_m)$ is an $m \times 1$ vector of observed phenotypes, and $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_q)$ is an $m \times q$ matrix of fixed effects including SNPs and other confounding variables. $\boldsymbol{\beta}$ is a $q \times 1$ vector representing coefficients of the fixed effects. \mathbf{Z} is usually a designed matrix which gives group information of the samples. Here we treat each sample belongs to an unique group that \mathbf{Z} turns out to be an identity matrix. \mathbf{u} is the random effect of the

*Max-Planck Institute of Mathematics, Vivatsgasse 7, 53111 Bonn, Germany, wuchenxi2013@gmail.com

mixed model with $Var(\mathbf{u}) = \delta_g^2 \mathbf{K}$. To account for confounding by population structure, family structure, and cryptic relatedness in GWAS, \mathbf{K} can be reliably estimated from genetic markers, e.g., using the realized relationship matrix (RRM) which captures the overall genetic similarity between all pairs of samples.

The resulting mixed model is typically applied in the context of single candidate SNP, that is, restricting the sum in Equation (1) to a particular SNP while ignoring all others. This independent analysis can be compromised by complex genetic architectures with some genetic factors masking others.

As an alternative, we proposed an efficient approach to carry out joint inference in the model implied by Equation (1). Our approach assesses all the SNPs at the same time while accounting for their interdependencies and without making any assumptions on their ordering. To allow for applications to genome-wide SNP data, we place a XX prior over the fixed effects β_j , assigning zero-effect size to majority of SNPs via bi-level selection method that select non-zero effect SNPs both at the group level and within group level, which as done in the classical sparse group lasso.

We call this approach SGL-LMM (Sparse Group Lasso-Linear Mixed Model), as it combines the advantages of established LMM with sparse group Lasso regression. The resulting model allows for dissecting the explained phenotype variance into individual SNP effects and effects caused by population structure.

2.3.1 Sparse Group Lasso - Linear Mixed Model

Let $\mathbf{X} \in \mathbb{R}^{m \times q}$ denotes the matrix of fixed effects for m individuals. This matrix includes the column of 1s corresponding to the offset, the covariances, and all the SNPs to be considered. We model the phenotype for m individuals, $\mathbf{y} \in \mathbb{R}^{m \times 1}$ as the sum of fixed effects β_j of SNPs \mathbf{x}_j and random effect due to confounding influences \mathbf{u} [see Equation (1)]. The fixed effect terms are summed over genome-wide SNPs, where the great majority of SNPs have zero-effect size, that is, $\beta_j = 0$, which is achieved by a XX prior on all weights. The random effect \mathbf{u} is not observed directly. Instead, we assume that the distribution of \mathbf{u} is Gaussian with covariance \mathbf{K} , $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \sigma_g^2 \mathbf{K})$.

Assuming Gaussian noise, $\psi \sim \mathcal{N}(\mathbf{0}, \sigma_e^2 \mathbf{I})$, and marginalizing over the random vector \mathbf{u} , we can write down the conditional posterior distribution over the weight vector β :

$$p(\beta | \mathbf{y}, \mathbf{X}, \mathbf{K}, \sigma_g^2, \sigma_e^2, \lambda) \propto \underbrace{\left(\mathcal{N}(\mathbf{y} | \sum_{j=1}^q \beta_j \mathbf{x}_j, \sigma_g^2 \mathbf{K} + \sigma_e^2 \mathbf{I}) \right)}_{\text{marginallikelihood}} \underbrace{\pi(\beta)}_{\text{prior}} \quad (2)$$

With the prior of the form

$$\pi(\beta) \propto \exp(-\lambda \sum_{G \in \mathcal{G}} (\alpha_G \|\beta_G\|_2 + \|\beta_G\|_1)) \quad (3)$$

Here, suppose we are given M (possibly overlapping) groups $\mathcal{G} = \{G_1, G_2, \dots, G_M\}$, so that $G_i \subset \{1, 2, \dots, p\} \forall i$, of maximum size B . These groups contain sets of "similar" features, such as SNPs within a region of a gene. We assume that all but $k \ll M$ groups are identically zero. Among the active groups, we further assume that at most a fraction $\alpha \in (0, 1)$ of the coefficients per group are non zero. $\alpha_G > 0$ are constants that balance the tradeoff between the group norms and the ℓ_1 norm. λ denotes the sparsity hyperparameter of our prior, σ_g^2 is the magnitude of the random effect component, and σ_e^2 denotes the magnitude of the residual variance.

2.3.2 Parameter inference

Learning the hyperparameters $\theta = \{\sigma_g^2, \sigma_e^2\}$ and the weights β jointly is a hard non-convex optimization problems. Here, we propose a combination of fitting some of these parameters on the null model with the individual SNP effects excluded and reduction to a sparse overlapping group lasso problem.

We first optimize σ_g^2, σ_e^2 by maximum likelihood under the null model, ignoring the effect of individual SNPs. The analogous procedure is widely used in single-SNP mixed models and has been shown to yield near-identical results to an exact approach.

Evaluation of the log likelihood The log likelihood is parameterized by a weight vector β and the variances of the random components, σ_g^2 and σ_e^2 .

$$LL(\sigma_g^2, \sigma_e^2, \beta) = \log \mathcal{N}(\mathbf{y} | \mathbf{X}\beta, \sigma_g^2 \mathbf{K} + \sigma_e^2 \mathbf{I}) \quad (4)$$

To speed up the computation needed, we introducing $\delta \equiv \frac{\sigma_e^2}{\sigma_g^2}$, the covariance matrix becomes $\sigma_g^2(\mathbf{K} + \delta \mathbf{I})$, and the likelihood becomes a function of β, δ and σ_g^2 :

$$LL(\delta, \sigma_g^2, \beta) = \log \mathcal{N}(\mathbf{y} | \mathbf{X}\beta, \sigma_g^2(\mathbf{K} + \delta \mathbf{I})) \quad (5)$$

Using the formula for the multivariate Normal distribution, we obtain

$$LL(\delta, \sigma_g^2, \beta) = -\frac{1}{2} \left(m \log(2\pi\sigma_g^2) + \log(|\mathbf{K} + \delta \mathbf{I}|) + \frac{1}{\sigma_g^2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{K} + \delta \mathbf{I})^{-1} (\mathbf{y} - \mathbf{X}\beta) \right) \quad (6)$$

Since \mathbf{K} is symmetric matrix, letting $\mathbf{U} \mathbf{S} \mathbf{U}^T = \mathbf{K}$ be the spectral decomposition of \mathbf{K} , and noting that $\mathbf{I} = \mathbf{U} \mathbf{U}^T$, Equation (6) becomes

$$LL(\delta, \sigma_g^2, \beta) = -\frac{1}{2} \left(m \log(2\pi\sigma_g^2) + \log(|\mathbf{U} \mathbf{S} \mathbf{U}^T + \delta \mathbf{U} \mathbf{U}^T|) + \frac{1}{\sigma_g^2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{U} \mathbf{S} \mathbf{U}^T + \delta \mathbf{U} \mathbf{U}^T)^{-1} (\mathbf{y} - \mathbf{X}\beta) \right)$$

Next, we factor out \mathbf{U} and \mathbf{U}^T from the covariance of the Normal, so that it becomes the diagonal matrix $(\mathbf{S} + \delta \mathbf{I})$, obtaining

$$-\frac{1}{2} \left(m \log(2\pi\sigma_g^2) + \log(|\mathbf{U}(\mathbf{S} + \delta\mathbf{I})\mathbf{U}^T|) + \frac{1}{\sigma_g^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{U}(\mathbf{S} + \delta\mathbf{I})\mathbf{U}^T)^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right) \quad (7)$$

The determinant of $|\mathbf{U}(\mathbf{S} + \delta\mathbf{I})\mathbf{U}^T|$ can be written as $|(\mathbf{S} + \delta\mathbf{I})|$ using the properties that $|\mathbf{AB}| = |\mathbf{A}||\mathbf{B}|$, and that $|\mathbf{U}| = |\mathbf{U}^T| = 1$. The inverse of this part can be rewritten as $\mathbf{U}(\mathbf{S} + \delta\mathbf{I})^{-1}\mathbf{U}^T$ using the properties that $(\mathbf{AB}^{-1}) = \mathbf{B}^{-1}\mathbf{A}^{-1}$, that $\mathbf{U}^{-1} = \mathbf{U}^T$, and that $\mathbf{U}^{-T} = \mathbf{U}$. Thus, after additionally moving out \mathbf{U} from the covariance term so that it now acts as a rotation matrix on the inputs \mathbf{X} and targets \mathbf{y} , we obtain

$$\begin{aligned} & -\frac{1}{2} \left(m \log(2\pi\sigma_g^2) + \log(|\mathbf{U}||(\mathbf{S} + \delta\mathbf{I})||\mathbf{U}^T|) + \frac{1}{\sigma_g^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{U}(\mathbf{S} + \delta\mathbf{I})^{-1} \mathbf{U}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right) \\ &= -\frac{1}{2} \left(m \log(2\pi\sigma_g^2) + \log(|(\mathbf{S} + \delta\mathbf{I})|) + \frac{1}{\sigma_g^2} ((\mathbf{U}^T \mathbf{y}) - (\mathbf{U}^T \mathbf{X})\boldsymbol{\beta})^T (\mathbf{S} + \delta\mathbf{I})^{-1} ((\mathbf{U}^T \mathbf{y}) - (\mathbf{U}^T \mathbf{X})\boldsymbol{\beta}) \right) \end{aligned} \quad (8)$$

Finding the maximum likelihood genetic variance with null model We start by eliminating $\boldsymbol{\beta}$ from Equation (8), and set derivative with respect to σ_g^2 to zero. As the covariance matrix of the Normal distribution is now a diagonal matrix $(\mathbf{S} + \delta\mathbf{I})$, we could obtain

$$0 = -\frac{1}{2} \left(\frac{n}{\hat{\sigma}_g^2} - \frac{1}{\hat{\sigma}_g^4} \sum_{i=1}^m \frac{([U^T \mathbf{y}]_i)^2}{[\mathbf{S}]_{ii} + \delta} \right)$$

Multiplying both sides by $2\hat{\sigma}_g^4$ and solving for $\hat{\sigma}_g^2$, we get

$$\hat{\sigma}_g^2 = \frac{1}{n} \sum_{i=1}^m \frac{([U^T \mathbf{y}]_i)^2}{[\mathbf{S}]_{ii} + \delta}$$

This equation can be evaluated in $\mathcal{O}(m)$.

Optimization of δ Plugging in $\hat{\sigma}_g^2$ into Equation(8) with null model, the log likelihood becomes a function only of δ , $LL(\delta, \hat{\sigma}_g^2(\delta)) = LL(\delta)$:

$$LL(\delta) = -\frac{1}{2} \left(n \log(2\pi) + \sum_{i=1}^m \log([\mathbf{S}]_{ii} + \delta) + n + n \log \frac{1}{n} \left(\sum_{i=1}^m \frac{([U^T \mathbf{y}]_i)^2}{[\mathbf{S}]_{ii} + \delta} \right) \right)$$

We optimize this function of δ using a one-dimensional numerical optimizer to find the maximum likelihood value of δ , from which the maximum likelihood values of all the parameters can be directly computed.

Reduction to sparse overlapping group lasso problem Having fixed δ , we use the spectral decomposition of \mathbf{K} again to rotate our data such that the covariance matrix becomes isotropic

$$p(\boldsymbol{\beta}|\tilde{\mathbf{y}}, \tilde{\mathbf{X}}, \sigma_g^2) \propto \mathcal{N}(\tilde{\mathbf{y}}|\sum_{j=1}^q \beta_j \tilde{\mathbf{x}}_j, \sigma_g^2 \mathbf{I})\pi(\boldsymbol{\beta}) \quad (9)$$

Here, $\tilde{\mathbf{X}}$ denotes the rotated and rescaled genotypes, and $\tilde{\mathbf{y}}$ denotes the respective phenotypes

$$\begin{aligned} \tilde{\mathbf{X}} &= (\mathbf{S} + \delta \mathbf{I})^{-\frac{1}{2}} \mathbf{U}^T \mathbf{X} \\ \tilde{\mathbf{y}} &= (\mathbf{S} + \delta \mathbf{I})^{-\frac{1}{2}} \mathbf{U}^T \mathbf{y} \end{aligned}$$

Sparse group lasso problem Using the transformation above, we can obtain the Maximum-A-Posteriori parameter estimation of the most probable weights in Equation (9)

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{MAP} &= \underset{\boldsymbol{\beta}}{\operatorname{argmax}} p(\boldsymbol{\beta}|\tilde{\mathbf{y}}, \tilde{\mathbf{X}}, \sigma_g^2) = \underset{\boldsymbol{\beta}}{\operatorname{argmax}} p(\tilde{\mathbf{y}}|\boldsymbol{\beta}, \tilde{\mathbf{X}}, \sigma_g^2)\pi(\boldsymbol{\beta}) \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmax}} \log p(\tilde{\mathbf{y}}|\boldsymbol{\beta}, \tilde{\mathbf{X}}, \sigma_g^2)\pi(\boldsymbol{\beta}) \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmax}} \{ \log \pi(\boldsymbol{\beta}) + \log \prod_{j=1}^m p(\tilde{\mathbf{y}}_j|\boldsymbol{\beta}, \tilde{\mathbf{x}}^j, \sigma_g^2) \} \end{aligned} \quad (10)$$

The log likelihood is as follow:

$$\begin{aligned} \log \prod_{j=1}^m p(\tilde{\mathbf{y}}_j|\boldsymbol{\beta}, \tilde{\mathbf{x}}^j, \sigma_g^2) &= \sum_{j=1}^m \log p(\tilde{\mathbf{y}}_j|\boldsymbol{\beta}, \tilde{\mathbf{x}}^j, \sigma_g^2) \\ &= \sum_{j=1}^m \log \left(\frac{1}{(2\pi)^{1/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (\tilde{\mathbf{y}}_j - \boldsymbol{\beta}^T \tilde{\mathbf{x}}^j)^T \Sigma^{-1} (\tilde{\mathbf{y}}_j - \boldsymbol{\beta}^T \tilde{\mathbf{x}}^j) \right) \right) \end{aligned}$$

Where $|\Sigma|^{1/2} = \sigma_g^m$, and $\Sigma^{-1} = \frac{1}{\sigma_g^2} \mathbf{I}$

$$\log \prod_{j=1}^m p(\tilde{\mathbf{y}}_j|\boldsymbol{\beta}, \tilde{\mathbf{x}}^j, \sigma_g^2) =$$

2.4	Overview of geXGB
2.5	Simulation study
2.6	Competitive methods
3	Results
3.1	Results for the simulation study
3.2	Real data sets analysis
4	Discussion
5	Conflict of interest
6	Acknowledgements
	References