

PROCEEDINGS

Open Access

# Penalized multivariate linear mixed model for longitudinal genome-wide association studies

Jin Liu<sup>1†</sup>, Jian Huang<sup>2†</sup>, Shuangge Ma<sup>3,4\*</sup>

From Genetic Analysis Workshop 18  
Stevenson, WA, USA. 13-17 October 2012

## Abstract

We consider analysis of Genetic Analysis Workshop 18 data, which involves multiple longitudinal traits and dense genome-wide single-nucleotide polymorphism (SNP) markers. We use a multivariate linear mixed model to account for the covariance of random effects and multivariate residuals. We divide the SNPs into groups according to the genes they belong to and score them using weighted sum statistics. We propose a penalized approach for genetic variant selection at the gene level. The overall modeling and penalized selection method is referred to as the penalized multivariate linear mixed model. Cross-validation is used for tuning parameter selection. A resampling approach is adopted to evaluate the relative stability of the identified genes. Application to the Genetic Analysis Workshop 18 data shows that the proposed approach can effectively select markers associated with phenotypes at gene level.

## Background

The Genetic Analysis Workshop 18 (GAW18) data consists of multiple longitudinal traits and dense genome-wide single-nucleotide polymorphism (SNP) markers. A commonly used approach for identifying markers associated with traits is to conduct single-variant analysis and then adjust for multiple comparisons on each trait. However, for complex polygenic traits, single-variant analysis methods may not be appropriate as they fail to take into account the accumulated and/or joint effects of multiple genetic variants on the traits. In addition, analyzing each trait separately does not take into account the correlation among traits, and thus can be ineffective. To overcome these limitations, we developed a joint analysis approach referred to as the penalized multivariate linear mixed model (PMLMM). This approach takes into account covariance of both random effects and residuals and uses a group minimax concave penalty (MCP) approach [1] for variant selection at the gene level. A resampling approach is adopted to evaluate

the relative stability of the identified genes. Our analysis of the GAW18 data indicates that the proposed approach can effectively select markers associated with multiple traits at the gene level.

## Methods

Consider a genetic association study with longitudinal measurements on  $N$  subjects,  $p$  genetic variants, and  $q$  environmental exposure covariates. Here a genetic variant can be a single SNP marker or a score representing a group of SNPs. For subject  $i$ , suppose that there are  $n_i$  longitudinal measurements on  $m$  traits. Let  $Y_i$  be the  $n_i \times m$  trait matrix for subject  $i$ . Let  $Y$  be the  $n \times m$  trait matrix for all the  $N$  subjects, where  $n = \sum_{i=1}^N n_i$ . The transpose of  $Y$  is  $Y' = (Y'_1, \dots, Y'_N)$ . Let  $X_i$  be the  $n_i \times p$  matrix consisting of the genetic variant scores of subject  $i$ . Let  $Z_i$  be the  $n_i \times q$  covariate matrix. We center all the measurements to have sample means equal to zero. When  $m = 1$ , this setting simplifies to that in Schellldorfer et al [2].

Consider the multivariate linear mixed model

$$Y_i = X_i B + Z_i C_i + E_i, \quad i = 1, \dots, N, \quad (1)$$

where  $B$  is a  $p \times m$  matrix representing the effects of  $p$  genetic variants on  $m$  traits, and  $C_i$  is a  $q \times m$  matrix

\* Correspondence: shuangge.ma@yale.edu

† Contributed equally

<sup>3</sup>School of Public Health, Yale University, 60 College Street, New Haven, CT 06520, USA

Full list of author information is available at the end of the article

representing the subject specific effects of the covariates  $Z_i$  for the  $i^{\text{th}}$  subject. We treat  $C_i$  as random effects. Assume that (a)  $E_i \sim \text{MN}_{n_i \times m}(0, \sum_1 I_{n_i})$ , that is,  $E_i$  is row-independent with column covariance matrix  $\Sigma_1$ , and each  $E_i$  is independent for  $i = 1, \dots, N$ ; (b)  $C_i \sim \text{MN}_{q \times m}(0, \sum_2 D)$ , where  $\Sigma_2$  is the column covariance matrix and  $D$  is the row covariance matrix, and each  $C_i$  is independent for  $i = 1, \dots, N$ ; (c) each  $E_i$  and  $C_i$  is independent; and (d)  $\Sigma_1 = \Sigma_2 = \Sigma$ .

Then  $Z_i C_i + E_i \sim \text{MN}_{n_i \times m}(0, \sum_1 Z_i D Z_i' + I_{n_i})$  and  $Y \sim \text{MN}_{n \times m}(XB, \sum_1 V)$  where  $V = \text{Diag}(V_1, \dots, V_N)$  and  $V_i = Z_i D Z_i' + I_{n_i}$ , where  $I_{n_i}$  is an  $n_i \times n_i$  identity matrix. A more detailed description of this model can be found in Liu et al [3].

From Dawid [4], the negative log-likelihood function is:

$$-\ell(B, V, \Sigma) = \text{constant} + \frac{n}{2} \log |\Sigma| + \frac{m}{2} \log |V| + \frac{1}{2} \text{tr}(\sum_1^{-1} (y - XB) V^{-1} (y - XB)) \quad (2)$$

Hastie et al. [5] suggest using  $\hat{\Sigma} = y'y/n$  for estimating  $\Sigma$ . We estimate  $D$  by using the estimates from  $m$  univariate linear mixed models and subsequently get the estimate  $\hat{V}$  of  $V$  as  $\hat{V}_i = Z_i \hat{D} Z_i' + I_{n_i}$ . Given  $\hat{\Sigma}$  and  $\hat{V}$ , we can transform the negative likelihood function into a weighted least squares criterion for estimating  $B$ , which is  $\text{tr}(\sum_1^{-1} (y - XB) \hat{V}^{-1} (y - XB))$ . For variant selection, we adopt the group MCP approach [6]. The overall penalized objective function is

$$Q(B) = \text{tr}(\sum_1^{-1} (y - XB) \hat{V}^{-1} (y - XB)) + \sum_{j=1}^p \rho(\|B_j\|_2; \lambda, \gamma), \quad (3)$$

where  $B_j$  is the  $j^{\text{th}}$  row of  $B$  and  $\rho(t; \lambda, \gamma) = \lambda \int_0^{|t|} (1 - x/(\gamma\lambda))_+ dx$  is the MCP with tuning parameter  $\lambda$  and regularization parameter  $\gamma$  [7].

For computation, we use a group coordinate descent algorithm [1]. The group MCP involves a regularization parameter and a tuning parameter. Generally speaking, smaller values of  $\gamma$  are better at retaining the unbiasedness of the MCP penalty for large coefficients, but they have the risk of creating objective functions that have problems with nonconvexity [8], are difficult to optimize, and yield solutions that are discontinuous with respect to  $\lambda$ . Simulation studies by Breheny and Huang [8] suggest that  $\gamma = 6$  is a reasonable choice. Therefore, we fix it to be 6 in our analyses. We search for the optimal value of  $\lambda$  using 5-fold cross-validation.

## Results

The GAW18 data set consists of dense genome-wide markers with longitudinal measurements on systolic and

diastolic blood pressure (SBP and DBP) and other covariates. Other measurements include gender, age, year of examination, use of antihypertensive medications, and tobacco smoking at up to 4 time points. In this study, we analyze the 157 unrelated individuals using SBP and DBP as traits and other medical and demographic covariates as random effects. Gene annotations for SNP data are obtained from <http://www.scandb.org>. SNPs in each gene are scored using weighted sum statistics to generate gene-level measurements [9]. After quality control, we have the genetic scores of 10,400 genes for further analysis. SBP, DBP, and genetic scores are standardized to have zero means and unit variances. This procedure removes the estimation of intercepts and makes the genes comparable.

We apply the proposed PMLMM to identify genetic variants that are associated with both SBP and DBP at the gene level. As a benchmark, we also analyze each trait separately using a penalized linear mixed model (PLMM) approach. Table 1 shows the genes identified using PMLMM. Table 2 summarizes the overlaps of genes selected using the different approaches. Although there is overlap, PMLMM and PLMM identify significantly different sets of genes. We evaluate the relative stability of identification of each gene using a resampling approach and calculate the observed occurrence index (OOI) [10]. A larger value of OOI indicates that the corresponding identified gene is more stably identified. Table 1 also shows OOI results. The identified genes have reasonably high OOIs.

## Discussion

In this study, we analyze the GAW18 data and develop a PMLMM approach. A multivariate linear mixed model is used to model variance components among traits and longitudinal measurements. A penalization approach is adopted for variant selection. In the estimation procedure, it can be considered heuristic to use  $\hat{\Sigma}$  and  $\hat{V}$  as proposed. Assumptions (a) to (c) are standard in mixed models, but the assumption that  $\Sigma_1 = \Sigma_2$  may be restrictive. Because our study is to identify multitrait-associated markers at the gene level, the restriction on variance components does not affect the selection result significantly. We are currently developing a similar approach to update variance components with more relaxed assumptions on  $\Sigma_1$  and  $\Sigma_2$ . An iterative algorithm can be implemented to solve for  $B$ ,  $\Sigma$ , and  $V$ . In variant selection, our method is designed to search for genes associated with all the traits considered. When different sets of genetic variants are suspected to be associated with different phenotypes, the sparse group penalization approach [11] can be applied.

## Conclusions

We have presented a penalized multivariate linear mixed model (PMLMM) for detecting pleiotropic genetic

**Table 1 Genes identified by PMLMM: estimates for SBP and DBP, and OOI**

Gene	SBP	DBP	OOI	Gene	SBP	DBP	OOI
<i>MMEL1</i>	0.002	-0.002	0.333	<i>TMEM41B</i>	0.027	0.033	0.403
<i>CD52</i>	0.085	0.060	0.697	<i>ARNTL</i>	0.024	0.006	0.247
<i>DPH2</i>	0.071	-0.032	0.323	<i>SPTY2D1</i>	0.025	-0.007	0.507
<i>C8A</i>	0.018	0.032	0.563	<i>CHST1</i>	-0.008	-0.031	0.540
<i>DNAJB4</i>	-0.028	-0.022	0.333	<i>MRE11A</i>	-0.042	-0.007	0.623
<i>HS2ST1</i>	0.002	0.006	0.307	<i>ENOX1</i>	-0.068	-0.032	0.647
<i>PROK1</i>	0.006	0.010	0.373	<i>LOC100132760</i>	-0.041	0.048	0.693
<i>THBS3</i>	-0.004	0.001	0.337	<i>SPRY2</i>	-0.004	-0.005	0.297
<i>C1orf182</i>	2E-04	0.045	0.490	<i>GABRG3</i>	-0.027	0.006	0.573
<i>TGFBR2</i>	-0.033	-0.030	0.627	<i>THBS1</i>	0.023	0.028	0.353
<i>LOC100129194</i>	0.005	0.011	0.217	<i>CSPG4</i>	0.098	-0.012	0.880
<i>LMOD3</i>	-0.013	-0.028	0.493	<i>C15orf27</i>	0.047	-0.014	0.490
<i>LOC653712</i>	0.017	0.001	0.450	<i>LOC100128570</i>	0.026	0.019	0.283
<i>LAMP3</i>	0.034	0.008	0.627	<i>HOMER2</i>	0.024	0.013	0.250
<i>EIF2B5</i>	-0.014	-0.014	0.417	<i>ADAMTS17</i>	0.002	0.002	0.270
<i>EHHADH</i>	0.003	-0.052	0.677	<i>SLC16A11</i>	0.015	-0.004	0.170
<i>SFRS12</i>	0.005	0.007	0.290	<i>ALDH3A1</i>	0.002	0.021	0.657
<i>C5orf32</i>	0.019	-0.029	0.560	<i>FLJ44815</i>	0.019	-0.005	0.210
<i>ZNF346</i>	0.001	-0.024	0.553	<i>TANC2</i>	0.003	-0.001	0.187
<i>LOC100128901</i>	-0.069	0.006	0.627	<i>PDE6G</i>	-0.056	-0.012	0.377
<i>OGDH</i>	0.123	0.038	0.777	<i>C19orf6</i>	0.013	0.048	0.577
<i>NSUN5</i>	0.035	-0.014	0.660	<i>TMEM146</i>	-0.001	-0.063	0.550
<i>PPP1R3A</i>	-0.034	-0.041	0.453	<i>STX10</i>	4E-04	0.015	0.333
<i>MEST</i>	-2E-04	-3E-04	0.197	<i>RLN3</i>	0.024	-0.034	0.603
<i>NOM1</i>	0.029	-0.001	0.630	<i>CYP4F11</i>	0.018	0.005	0.127
<i>FLJ41200</i>	0.013	0.005	0.333	<i>LOC728326</i>	0.048	-0.040	0.547
<i>LRRC19</i>	-0.015	-0.040	0.480	<i>ZNF585A</i>	0.028	0.082	0.720
<i>CCIN</i>	0.006	-0.004	0.437	<i>SUPT5H</i>	0.020	-0.014	0.233
<i>LOC100130911</i>	0.004	0.024	0.467	<i>FLJ10490</i>	0.004	2E-04	0.327
<i>PTCH1</i>	0.004	-1E-04	0.300	<i>ZNF331</i>	0.027	0.003	0.330
<i>DFNB31</i>	0.057	-0.008	0.710	<i>BACE2</i>	-0.116	-0.074	0.940
<i>OR52D1</i>	0.023	-0.017	0.490	<i>KRTAP10-12</i>	0.002	0.003	0.233

**Table 2 Overlap of selected genes between PMLMM and PLMM**

	PMLMM	PLMM*	PLMM†
PMLMM	64	24	16
PLMM <sup>1</sup>		40	0
PLMM <sup>2</sup>			29

\*PLMM on SBP.

†PLMM on DBP.

associations among multiple traits in the presence of pedigree structures. The proposed approach combines the advantages of mixed models that allow for elegant correction for pedigree-based family data and integrative analysis of multiple traits. Compared with PLMM which considers one trait at a time, the proposed PMLMM can achieve better performance when the pleiotropic effect is appropriately modeled.

**Competing interests**

The authors declare that they have no competing interests.

**Authors' contributions**

All authors were involved in study design. JL conducted the numerical work. All authors were involved in manuscript preparation, and read and approved the final manuscript.

**Acknowledgements**

This study was supported by NIH grants CA142774, CA165923, and CA120988, the VA Cooperative Studies Program of the Department of Veterans Affairs, Office of Research and Development, and 2012LD001 from National Bureau of Statistics of China. The GAW18 whole genome sequence data were provided by the T2D-GENES Consortium, which is supported by NIH grants U01 DK085524, U01 DK085584, U01 DK085501, U01 DK085526, and U01 DK085545. The other genetic and phenotypic data for GAW18 were provided by the San Antonio Family Heart Study and San Antonio Family Diabetes/Gallbladder Study, which are supported by NIH grants P01 HL045222, R01 DK047482, and R01 DK053889. The Genetic Analysis Workshop is supported by NIH grant R01 GM031575.

This article has been published as part of BMC Proceedings Volume 8 Supplement 1, 2014: Genetic Analysis Workshop 18. The full contents of the supplement are available online at <http://www.biomedcentral.com/bmcproc/>

supplements/8/S1. Publication charges for this supplement were funded by the Texas Biomedical Research Institute.

#### Authors' details

<sup>1</sup>School of Public Health, University of Illinois at Chicago, 1601 W. Taylor Street, Chicago, IL 60612, USA. <sup>2</sup>Department of Statistics & Actuarial Science, Department of Biostatistics, University of Iowa, 241 Schaeffer Hall, Iowa City, IA 52242, USA. <sup>3</sup>School of Public Health, Yale University, 60 College Street, New Haven, CT 06520, USA. <sup>4</sup>VA Cooperative Studies Program Coordinating Center, West Haven, CT 06516, USA.

Published: 17 June 2014

#### References

1. Huang J, Wei F, Ma S: Semiparametric regression pursuit. *Stat Sin* 2012, **22**:1403-1426.
2. Schellhdorfer J, van de Geer S: Estimation for high-dimensional linear mixed-effects models using L1-penalization. *Scand Stat Theory Appl* 2011, **38**(2):197-214.
3. Liu J, Yang C, Shi X, Zhao H, Huang J, Ma S: A penalized multiple-trait mixed model for association mapping with population structure correction. *Technical Report (arXiv preprint arXiv:1305.4413)* 2013.
4. Dawid A: Some matrix-variate distribution theory: notational considerations and a Bayesian application. *Biometrika* 1981, **68**:265-274.
5. Hastie T, Tibshirani R, Friedman J: *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. New York, Springer-Verlag, 2009.
6. Liu J, Huang J, Ma S: Analysis of genome-wide association studies with multiple outcomes using penalization. *PLoS One* 2012, **7**:e51198.
7. Zhang CH: Nearly unbiased variable selection under minimax concave penalty. *Ann Stat* 2010, **38**:894-942.
8. Breheny P, Huang J: Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection. *Ann Appl Stat* 2011, **5**:232-253.
9. Madsen B, Browning S: A groupwise association test for rare mutations using a weighted sum statistic. *PLoS Genet* 2009, **5**:e1000384.
10. Huang J, Ma S: Variable selection in the accelerated failure time model via the bridge method. *Lifetime Data Anal* 2010, **16**:176-195.
11. Liu J, Huang J, Xie Y, Ma S: Sparse group penalized integrative analysis of multiple cancer prognosis datasets. *Genet Res (Camb)* 2013, **95**(2-3):68-77.

doi:10.1186/1753-6561-8-S1-S73

**Cite this article as:** Liu et al.: Penalized multivariate linear mixed model for longitudinal genome-wide association studies. *BMC Proceedings* 2014 **8**(Suppl 1):S73.

**Submit your next manuscript to BioMed Central  
and take full advantage of:**

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at  
www.biomedcentral.com/submit

