

Pleiotropic mapping for genome-wide association studies using group variable selection.

B. Liquet^{1,2} and K. Mengersen² and A. N. Pettitt² and M. Sutton²

¹ LMAP, Université de Pau et des Pays de L'Adour'

² ACEMS, QUT, Australia



Pleiotropic mapping for genome-wide association studies using group variable selection

- ▶ **Pleiotropy**: genetic variants which affect **multiple different complex diseases**
- ▶ **Example**: genetic variants which affect both **Breast and Thyroid cancer**.
- ▶ Results from GWAS suggest that complex diseases are often affected by many variants with small effects (**known as polygenicity**)
- ▶ **Aims**:
 - ▶ statistical method to **leverage pleiotropic effects**
 - ▶ **incorporate prior pathway knowledge** to increase statistical power and identify important risk variants.

Genomics Data: Wide Data, High Dimensional Data

Wide Data



Thousands / Millions of Variables

Hundreds of Samples

Screening and fdr,
Lasso, SVM, Stepwise

We have too many variables, prone to overfitting. Need to remove variable, or regularize, or both

- ▶ **Main constraint:** situation with $p > n$
- ▶ **Strong colinearity** among the variables.

Group structures within the data

- ▶ **Genomics**: genes within the same pathway have similar functions and act together in regulating a biological system.

↔ These genes can add up to have a larger effect

↔ can be detected as a group (i.e., at a pathway or gene set/module level).

Group structures within the data

- ▶ **Genomics**: genes within the same pathway have similar functions and act together in regulating a biological system.

↔ These genes can add up to have a larger effect

↔ can be detected as a group (i.e., at a pathway or gene set/module level).

We consider variables are divided into groups:

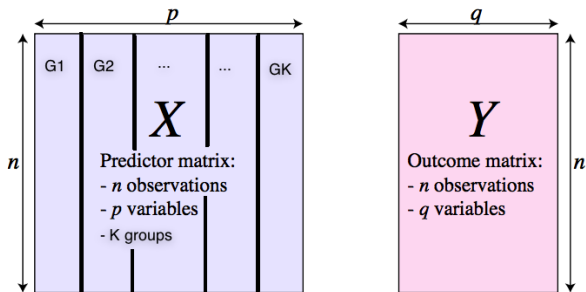
- ▶ Example p : SNPs grouped into K genes

$$\mathbf{X} = [\underbrace{SNP_1, \dots, SNP_k}_{gene_1} | \underbrace{SNP_{k+1}, SNP_{k+2}, \dots, SNP_h}_{gene_2} | \dots | \underbrace{SNP_{l+1}, \dots, SNP_p}_{gene_K}]$$

- ▶ Example p : genes grouped into K pathways/modules ($X_j = gene_j$)

$$\mathbf{X} = [\underbrace{X_1, X_2, \dots, X_k}_{M_1} | \underbrace{X_{k+1}, X_{k+2}, \dots, X_h}_{M_2} | \dots | \underbrace{X_{l+1}, X_{l+2}, \dots, X_p}_{M_K}]$$

Our contribution for Multivariate phenotypes



- ▶ Select **group variables** taking into account the data structures; **all the variables** within a group are selected otherwise none of them are selected
- ▶ Combine **both sparsity of groups and within each group**; only **relevant variables** within a group are selected

Our contribution for Multivariate phenotypes

Frequentist Approaches: Partial Least Squares (PLS)

- ▶ Sparse Group PLS : $\text{SNP} \subset \text{Gene}$ or $\text{Gene} \subset \text{Pathways}$

Liquet B., Lafaye de Micheaux P., Hejblum B. and Thiebaut R., (2016) *Group and Sparse Group Partial Least Square Approaches Applied in Genomics Context*. **Bioinformatics**, 32(1), 35–42.

- ▶ Sparse Group subgroup PLS : $\text{SNP} \subset \text{Gene} \subset \text{Pathways}$

M. Sutton, R. Thiebaut, and B. Liquet. (2018) *Sparse group subgroup Partial Least Squares with application to genomics data*. *Statistics in Medicine*.

Our contribution for Multivariate phenotypes

Frequentist Approaches: Partial Least Squares (PLS)

- ▶ Sparse Group PLS : SNP \subset Gene or Gene \subset Pathways

Liquet B., Lafaye de Micheaux P., Hejblum B. and Thiebaut R., (2016) *Group and Sparse Group Partial Least Square Approaches Applied in Genomics Context*. **Bioinformatics**, 32(1), 35–42.

- ▶ Sparse Group subgroup PLS : SNP \subset Gene \subset Pathways

M. Sutton, R. Thiebaut, and B. Liquet. (2018) *Sparse group subgroup Partial Least Squares with application to genomics data*. *Statistics in Medicine*.

Main ideas:

- ▶ combining L_1 and L_2 penalties into the optimization function
- ▶ **Sparse Group Penalties:**

$$\lambda_1 \sum_{g=1}^G \sqrt{p_g} \|\beta_g\|_2 + \lambda_2 \|\beta\|_1$$

Our contribution for Multivariate phenotypes

Bayesian Approaches: Multivariate regression model

- ▶ Bayesian group lasso model with spike and slab priors

Liquet, B., Mengersen, K., Pettitt, A. N. and Sutton, M. (2016). *Bayesian Variable Selection Regression Of Multivariate Responses For Group Data* Bayesian Analysis. Volume 12, Number 4 (2017), 1039-1067.

Main ideas:

- ▶ **spike and slab** priors providing variable selection at the **group level**.
- ▶ **hierarchical** spike and slab prior structure to select variables both at the **group level and within each group**.

Extention for Pleiotropy: Model for multiple GWAS

- ▶ Suppose we have data from K independent GWAS datasets, $\mathcal{D} = \mathcal{D}_1 \cup \dots \cup \mathcal{D}_K$, where $\mathcal{D}_k = (\{y_1, x_1\}, \dots, \{y_{n_k}, x_{n_k}\})$
- ▶ $y_{ik} \in \{0, 1\}$ denotes the phenotype of the k th study
- ▶ $x_{ik} \in \mathbb{R}^p$ is the vector with corresponding p SNPs.
- ▶ **Logistic regression model**

$$\text{Logit}(P(y_{ik} = 1|x_{ik})) = x_{ik}^T \beta_{\cdot k} \text{ for } k = 1, \dots, K,$$

- ▶ $\beta_{\cdot k} \in \mathbb{R}^p$ the regression coefficients for the k th GWAS.
- ▶ Let $\beta_j \in \mathbb{R}^K$, $j = 1, \dots, p$, the vector of K regression coefficients corresponding to the j th SNP over the K GWAS.

Group Structure

- ▶ SNPs can be partitioned into G groups (genes or Pathways)
- ▶ Let $\pi_g, g = 1, \dots, G$ the set of SNPs contained in the g th group with $p_g = |\pi_g|$.
- ▶ Matrix of all regression coefficients as $\mathbf{B} = (\beta_{\cdot 1}, \dots, \beta_{\cdot K}) = (\beta_{1 \cdot}, \dots, \beta_{p \cdot})^T$.

Frequentist Approach

- ▶ The log likelihood for the combined datasets:

$$p(\mathcal{D} | \mathbf{B}) = \sum_{k=1}^K L_k \quad \text{where } L_k \text{ Log-Likelihood of study } k$$

- ▶ The penalised likelihood estimate

$$\widehat{\mathbf{B}} = \underset{\mathbf{B} \in \mathbb{R}^{p \times K}}{\operatorname{argmin}} \left\{ - \sum_{k=1}^K L_k + \lambda_1 \|\mathbf{B}\|_{G_{2,1}} + \lambda_2 \|\mathbf{B}\|_{\ell_{2,1}} \right\} \quad (1)$$

- ▶ $G_{2,1}$ -norm penalty $\|\mathbf{B}\|_{G_{2,1}} = \sum_{g=1}^G \sqrt{\sum_{i \in \pi_g} \sum_{j=1}^K \beta_{ik}^2}$
- ▶ $\ell_{2,1}$ -norm penalty $\|\mathbf{B}\|_{\ell_{2,1}} = \sum_{i=1}^p \sqrt{\sum_{k=1}^K \beta_{ik}^2}$ respectively.
- ▶ The $G_{2,1}$ -norm fixes the group structure across studies and encourages sparsity at a gene level.
- ▶ The $\ell_{2,1}$ -norm which allows sparsity within a group.

Inference

- ▶ Inference using the **alternating direction method of multipliers algorithm** (ADMM).
- ▶ Novel approach for **identifying pleiotropic effects** as it accounts for gene specific and SNP specific effects using a variable selection approach.
- ▶ The method is only capable of producing a **point estimate of \mathbf{B}** and accurate estimation of the variance for these parameters is not easily given.

Bayesian Logistic regression with multivariate spike and slab prior: LogitMBGL-SS

- ▶ Let $\gamma = (\gamma_1, \dots, \gamma_p)^T$ indicate the association status for SNPs where $\gamma_j = 1$ indicates that the j th SNP is associated to all K traits.
- ▶ **Spike and slab prior** for the j th SNP $\beta_j \in \mathbb{R}^K$,

$$\beta_j \sim (1 - \gamma_j)\mathcal{N}_K(0, \tau_j^2 \mathbf{V}) + \gamma_j \delta_0(\beta_j)$$

$$\tau_j^2 \sim \text{Gamma}\left(\frac{K+1}{2}, \frac{\lambda}{2}\right),$$

$$\mathbf{V} \sim IW(d, \mathbf{Q}),$$

$$\gamma_j \sim \text{Bernolli}(\alpha_0)$$

$$\alpha_0 \sim \text{Beta}(a, b)$$

for $j = 1, \dots, p$, where $\delta_0(\beta_j)$ denotes a point mass at $\mathbf{0} \in \mathbb{R}^K$.

- ▶ Here, $\mathbf{V} \in \mathbb{R}^{K \times K}$ is a covariance matrix modeling the covariance of the SNP effect on the traits.

Extension

- ▶ Should perform well when the **SNPs are independent**.
- ▶ **GWAS datasets**: **strong correlations** that can occur between SNPs within the same gene.
- ▶ **Solution**: reparameterise the coefficients to handle the sparsity at a gene grouping level and individual feature level separately.
- ▶ $\tau \in \mathbb{R}^p$ to model individual sparsity
- ▶ $\mathbf{b}^{(g)} \in \mathbb{R}^{p_g K}$ with $\mathbf{b}^{(g)} = (\mathbf{b}_1^{(g)T}, \dots, \mathbf{b}_{p_g}^{(g)T})$ where $\mathbf{b}_j^{(g)} \in \mathbb{R}^K$ for group sparsity.

$$\beta_j = \tau_j \mathbf{b}_j^{(g)}, \quad \text{where } \tau_j \geq 0, \quad \text{for all } j \in \pi_g.$$

Bayesian Logistic regression using multivariate sparse group selection with spike and slab priors

$$\beta_{j\cdot} = \tau_j \mathbf{b}_j^{(g)}, \quad \text{where } \tau_j \geq 0, \quad \text{for all } j \in \pi_g.$$

We assume the following multivariate spike and slab

$$\mathbf{b}^{(g)} \sim (1 - \alpha_0) \mathcal{N}_{p_g}(\mathbf{0}, \mathbb{I}_{p_g} \otimes \mathbf{V}) + \alpha_0 \delta_0(\mathbf{b}^{(g)})$$

$$\tau_j \sim (1 - \alpha_1) \mathcal{N}^+(0, s^2) + \alpha_1 \delta_0(\tau_j),$$

$$\alpha_0 \sim \text{Beta}(a_1, a_2)$$

$$\alpha_1 \sim \text{Beta}(c_1, c_2)$$

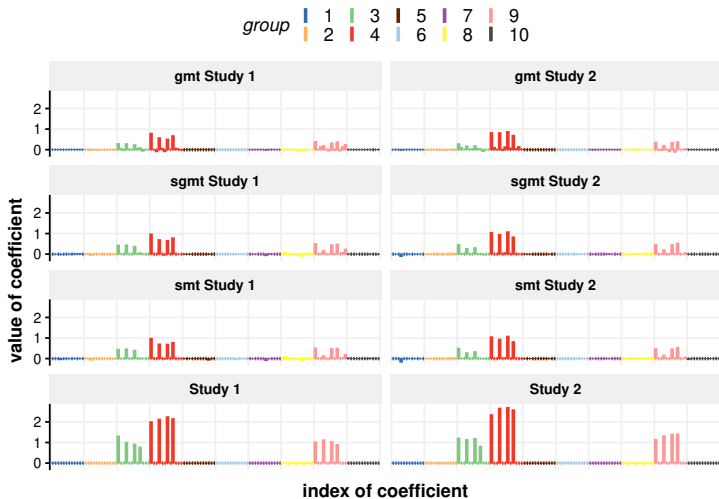
$$s^2 \sim \text{InvGamma}(1, t)$$

for $j \in \pi_g$ and $g = 1, \dots, G$

Signal recovery:

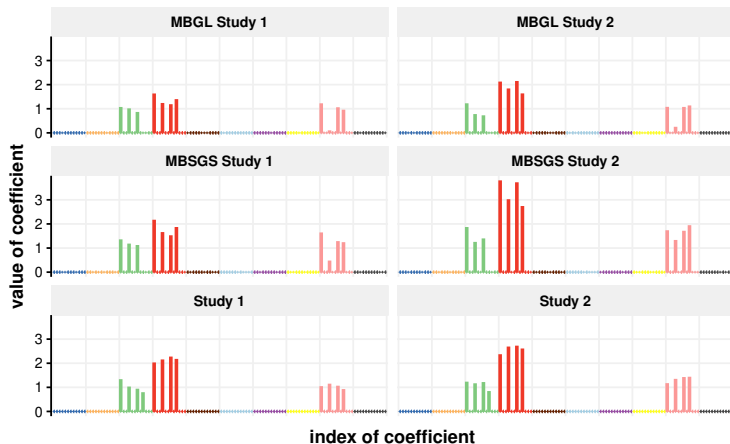
- (i) **GMT**: Grouped multi-task penalised logistic regression
($\lambda_1 > 0, \lambda_2 = 0$) using $G_{2,1}$ -norm
- (ii) **SMT**: Sparse multi-task penalised logistic regression
($\lambda_1 = 0, \lambda_2 > 0$) using $\ell_{2,1}$ -norm
- (iii) **SGMT**: Sparse group multi-task penalised logistic regression
($\lambda_1 > 0, \lambda_2 > 0$)
- (iv) **LOGITMBGL**: Bayesian logistic regression using multivariate group lasso with spike and slab prior
- (v) **LOGITMBSGS**: Bayesian logistic regression using multivariate sparse group selection with spike and slab prior

Results: Frequentist



Results: Bayesian

group 1 3 5 7 9
 2 4 6 8 10



Main Conclusion on the simulation studies

- ▶ The penalised approaches perform reasonably well in **variable selection** but the reconstructed signal is **underestimated**.
- ▶ In general, the penalised likelihood methods suffer in terms of **false negatives**, selecting more variables to be nonzero than the Bayesian methods.
- ▶ The Bayesian methods perform the best in terms of **signal recovery** measured by the ℓ_1 error and **variable selection** performance metrics.
- ▶ The penalised likelihood approaches are computationally efficient using **alternating direction method of multipliers algorithm**
- ▶ Simulation results suggest that when **computationally possible** the Bayesian estimators should be used.
- ▶ **The multivariate Bayesian sparse group selection with spike and slab prior** performed the best in terms of signal recovery.
- ▶ **The Bayesian method** provides a natural method for quantifying the **variability** of the estimated coefficients.

What Next ?

- ▶ Application on real data: case/control studies
 - ▶ Breast Cancer and Thyroide Cancer
 - ▶ Thyroide Cancer (482 case, 463 control)
 - ▶ Breast Cancer (1172 case, 1125 control)
 - ▶ 6677 SNPs from 618 genes from 10 non-overlapping gene pathways.

Acknowledgement

This work is funded by “ La ligue contre le Cancer”



It is part of our current project on

“Cross Cancer Genomic Investigation of Pleiotropy effect and GxE interactions at pathway level: application to breast and thyroid cancers”