Honglang Wang Depart. of Stat. & Prob. wangho16@msu.edu

**Omics** Data Integration

Statistical Genetics/Genomics Journal



### Summary and discussion of "Joint Analysis of SNP and Gene Expression Data in Genetic Association Studies of Complex Diseases"

Omics data integration @ Statistical Genetics/Genomics Journal Club @ MSU

**Abstract** Technological platforms have advanced to a stage where many biological entities, e.g., genes, transcripts, and proteins, can be measured on the whole genome scale, yielding massive high-throughput 'omics data, such as genetic, genomic, epigenetic, protemic, and metabolomic data. The 'omics era provides and unprecedented promise of understanding common complex diseases, developing strategies for disease risk assessment, early detection, and prevention and intervention, and personalized therapies. In genetical genomics studies, it is important to jointly analyze gene expression data and genetic variants in exploring their associations with complex traits. This is for the discussion of the paper (Huang et al., 2014).

# Contents

1	Preliminary         1.1       Variance component test in GLM         1.2       Mediation Analysis	<b>2</b> 2 3
2	Summary of the paper         2.1       Methodology         2.2       Implementation	<b>4</b> 4 7
3	Simulation         3.1       Linkage Disequilibrium         3.2       Simulation Results	<b>8</b> 8 9
4	Discussion         4.1       Another way         4.2       Research Direction	<b>18</b> 18 19

## 1 Preliminary

### 1.1 Variance component test in GLM

In this section, we will recall some useful results for the variance component testing in generalized linear model with mixed effects (Lin, 1997).

Consider the following generalized linear mixed model

$$\mathbb{E}(Y|\mathbf{b}) = \boldsymbol{\mu}^{b}, \operatorname{Var}(Y|\mathbf{b}) = \operatorname{diag}\{\phi a_{i}^{-1}v(\boldsymbol{\mu}_{i}^{b})\};\$$

$$g(\boldsymbol{\mu}^{b}) = \boldsymbol{\eta}^{b} = \mathbf{X}\boldsymbol{\alpha} + \mathbf{Z}\mathbf{b},$$

$$\mathbf{b} \sim F(\mathbf{0}, \mathbf{D}(\boldsymbol{\theta})),$$
(1)

where  $g(\cdot)$  is a monotonic differentiable link function, and the covariance matrix of the random effects has the property that  $\boldsymbol{\theta} = \mathbf{0}$  implies  $\mathbf{D}(\boldsymbol{\theta}) = \mathbf{0}$ .

1. First of all by GLM theory we have the conditional log-quasilikelihood of  $\alpha$  given  ${\bf b}$ 

$$l_i(\boldsymbol{\alpha}; \mathbf{b}) \propto \int_{y_i}^{\mu_i^b} \frac{a_i(Y_i - u)}{\phi v(u)} du.$$
 (2)

,

Note that for the logistic mixed model, we have the conditional log-likelihood, so we could just use it here.

2. Then we have the marginal (integrated) quasilikelihood of  $(\alpha, \theta)$ 

$$L(\boldsymbol{\alpha}, \boldsymbol{\theta}) = \int \prod_{i} \exp(l_i(\boldsymbol{\alpha}; \mathbf{b})) dF(\mathbf{b}; \boldsymbol{\theta}) = \int \exp(\sum_{i=1}^{n} l_i(\boldsymbol{\alpha}; \mathbf{b})) dF(\mathbf{b}; \boldsymbol{\theta}).$$
(3)

3. Approximate the above likelihood (3) by using Laplace method since its hard to calculate the integral:

$$\exp(\sum_{i=1}^{n} l_{i}(\boldsymbol{\alpha}; \mathbf{b})) \approx \exp\left\{\sum_{i=1}^{n} l_{i}(\boldsymbol{\alpha}; \mathbf{0})\right\} \left\{1 + \sum_{i=1}^{n} \frac{\partial l_{i}(\boldsymbol{\alpha}; \mathbf{0})}{\partial \eta_{i}} \mathbf{Z}_{i}^{\mathsf{T}} \mathbf{b} + \frac{1}{2} \mathbf{b}^{\mathsf{T}} \left[\left\{\sum_{i=1}^{n} \frac{\partial l_{i}(\boldsymbol{\alpha}; \mathbf{0})}{\partial \eta_{i}} \mathbf{Z}_{i}\right\} \left\{\sum_{i=1}^{n} \frac{\partial l_{i}(\boldsymbol{\alpha}; \mathbf{0})}{\partial \eta_{i}} \mathbf{Z}_{i}^{\mathsf{T}}\right\} + \sum_{i=1}^{n} \frac{\partial^{2} l_{i}(\boldsymbol{\alpha}; \mathbf{0})}{\partial \eta_{i}^{2}} \mathbf{Z}_{i} \mathbf{Z}_{i}^{\mathsf{T}}\right] \mathbf{b}\right\},$$

$$L(\boldsymbol{\alpha}, \boldsymbol{\theta}) = \mathbb{E}_{\mathbf{b}}\left\{\exp(\sum_{i=1}^{n} l_{i}(\boldsymbol{\alpha}; \mathbf{b}))\right\} \approx \exp\left\{\sum_{i=1}^{n} l_{i}(\boldsymbol{\alpha}; \mathbf{0})\right\} \left\{1 + \frac{1}{2} \operatorname{tr}\left(\left[\left\{\sum_{i=1}^{n} \frac{\partial l_{i}(\boldsymbol{\alpha}; \mathbf{0})}{\partial \eta_{i}} \mathbf{Z}_{i}\right\} \left\{\sum_{i=1}^{n} \frac{\partial l_{i}(\boldsymbol{\alpha}; \mathbf{0})}{\partial \eta_{i}} \mathbf{Z}_{i}^{\mathsf{T}}\right\} + \sum_{i=1}^{n} \frac{\partial^{2} l_{i}(\boldsymbol{\alpha}; \mathbf{0})}{\partial \eta_{i}^{2}} \mathbf{Z}_{i} \mathbf{Z}_{i}^{\mathsf{T}}\right] \mathbf{D}(\boldsymbol{\theta})\right)\right\}$$

$$l(\boldsymbol{\alpha}, \boldsymbol{\theta}) = \log L(\boldsymbol{\alpha}, \boldsymbol{\theta}) \approx \sum_{i=1}^{n} l_{i}(\boldsymbol{\alpha}; \mathbf{0}) + \frac{1}{2} \operatorname{tr}\left(\left[\left\{\sum_{i=1}^{n} \frac{\partial l_{i}(\boldsymbol{\alpha}; \mathbf{0})}{\partial \eta_{i}} \mathbf{Z}_{i}\right\} \left\{\sum_{i=1}^{n} \frac{\partial l_{i}(\boldsymbol{\alpha}; \mathbf{0})}{\partial \eta_{i}} \mathbf{Z}_{i}^{\mathsf{T}}\right\} + \sum_{i=1}^{n} \frac{\partial^{2} l_{i}(\boldsymbol{\alpha}; \mathbf{0})}{\partial \eta_{i}^{2}} \mathbf{Z}_{i} \mathbf{Z}_{i}^{\mathsf{T}}\right] \mathbf{D}(\boldsymbol{\theta})\right).$$

$$(4)$$

Omics Data Integration  $\bullet$  Honglang Wang

page 2 of 19

4. A global score statistic for testing  $H_0: \theta = 0$  is constructed as follows:

$$\chi_G^2 := U_{\theta}(\hat{\alpha})^{\mathsf{T}} \tilde{I}(\hat{\alpha})^{-1} U_{\theta}(\hat{\alpha})$$
(5)

where  $\hat{\boldsymbol{\alpha}}$  is the MLE estimator under the null hypothesis (which is the usual generalized linear model),  $U_{\boldsymbol{\theta}}(\boldsymbol{\alpha})$  is the gradient vector  $\partial l(\boldsymbol{\alpha}, \boldsymbol{\theta})/\partial \boldsymbol{\theta}$ , and  $\tilde{I}$  is the information matrix of  $\boldsymbol{\theta}$  under the null, which takes the form

$$\tilde{I} = I_{\theta\theta} - I_{\alpha\theta}^{\dagger} I_{\alpha\alpha}^{-1} I_{\alpha\theta},$$
  
with  $I_{\alpha\theta} = \mathbb{E}(\frac{\partial l}{\partial \alpha} \frac{\partial l}{\partial \theta^{\dagger}}), I_{\theta\theta} = \mathbb{E}(\frac{\partial l}{\partial \theta} \frac{\partial l}{\partial \theta^{\dagger}})$  and  $I_{\alpha\alpha} = \mathbb{E}(\frac{\partial l}{\partial \alpha} \frac{\partial l}{\partial \alpha^{\dagger}}).$ 

#### **1.2 Mediation Analysis**

Mediation analysis (VanderWeele and Vansteelandt, 2010) is just one type of causal inference in statistics (Rubin, 1990).



Figure 1: Example of mediation with exposure A, mediator M, outcome Y, and covariates C.

We will let A denote an exposure of interest, Y a dichotomous outcome, and M a potential mediator. We let C denote a set of baseline covariates not affected by the exposure. The relations among these variables are depicted in Figure (1). For example, A may denote genetic variants such as SNP, M m-RNA expression level, and Y cardiovascular disease. A question of interest may then be the extent to which the effect of genetic variants A on cardiovascular disease Y is mediated through m-RNA expression level M and the extent to which it is through other ways.

Let Y(a, m) be the potential outcome that would have been observed if A = a and M = m, and M(a) be the potential outcome of M had the A been set to a.

The direct effect of SNPs is the effect of the SNPs on the disease outcome that is not through gene expression, whereas the indirect effect is the effect of the SNPs on the disease outcome that is through the gene expression. We can define the direct effect (DE), the indirect effect (IE) and the total effect (TE) of the SNPs, respectively, on the log odds ratio (OR) scale as:

1. the total effect (TE), conditional on C = c, comparing exposure level a with  $a^*$ , is defined by

$$\log OR_{a,a^*|c}^{TE} = \text{logit}(\mathbb{P}(Y(a, M(a)) = 1)) - \text{logit}(\mathbb{P}(Y(a^*, M(a^*)) = 1)).$$

2. the direct effect (DE), conditional on C = c and M = M(a), comparing exposure level a with  $a^*$ , is defined by

$$\log OR_{a,a^*|c}^{DE}(a) = \text{logit}(\mathbb{P}(Y(a, M(a)) = 1)) - \text{logit}(\mathbb{P}(Y(a^*, M(a)) = 1)).$$

3. the indirect effect (IE), conditional on C = c and A = a, comparing exposure level a with  $a^*$ , is defined by

$$\log OR_{a,a^*|c}^{IE}(a) = \text{logit}(\mathbb{P}(Y(a, M(a)) = 1)) - \text{logit}(\mathbb{P}(Y(a, M(a^*)) = 1)).$$

Thus we have that

$$\log OR_{a,a^*|c}^{TE} = \log OR_{a,a^*|c}^{DE}(a) + \log OR_{a,a^*|c}^{IE}(a^*).$$

### 2 Summary of the paper

#### 2.1 Methodology

In this paper, the authors proposed to jointly model a set of SNPs within a gene, a gene expression, and disease status, where a logistic model is used to model the dependence of disease status on the SNP set and the gene expression, and a linear model is used for the dependence of the gene expression on the SNP-set, both adjusting for covariates. We are primarily interested in testing whether a gene, whose effects are captured by SNPs and/or gene expression, is associated with a disease phenotype.

A SNP-set and gene expression pair can be defined in multiple ways. For example, the SNP-set is the SNPs in a gene and the expression of the gene. Alternatively, one can choose the SNP-set as the eQTLs of the corresponding gene expression.

Consider the following model:

$$\operatorname{logit}\{\mathbb{P}(Y_i = 1 | \mathbf{S}_i, G_i, \mathbf{X}_i; \boldsymbol{\beta}_{\mathbf{S}}, \boldsymbol{\gamma})\} = \mathbf{X}_i^{\mathsf{T}} \boldsymbol{\alpha} + \mathbf{S}_i^{\mathsf{T}} \boldsymbol{\beta}_{\mathbf{S}} + G_i \boldsymbol{\beta}_G + \mathbf{C}_i^{\mathsf{T}} \boldsymbol{\gamma}, i = 1, 2, \cdots, n \quad (6)$$

where  $\boldsymbol{\alpha} \in \mathbb{R}^{q}, \boldsymbol{\beta}_{\mathbf{S}} \in \mathbb{R}^{p}, \boldsymbol{\gamma} \in \mathbb{R}^{p}, \beta_{G} \in \mathbb{R}, \mathbf{X}_{i}$  is the *q*-dim vector of covariates,  $\mathbf{S}_{i}$  is the *p*-dim vector of SNPs,  $G_{i} \in \mathbb{R}$  is the expression level,  $\mathbf{C}_{i} = G_{i}\mathbf{S}_{i}$  is the interaction, and  $Y_{i}$  is the dichotomous response. And since the SNPs in a gene might be large and some might be highly correlated (due to linkage disequilibrium), we assume  $\beta_{\mathbf{S},j} \stackrel{IID}{\sim} (0, \tau_{S})$  and  $\gamma_{j} \stackrel{IID}{\sim} (0, \tau_{I})$ . We are interested in the following global test problem

$$H_0: \tau_S = \tau_I = \beta_G = 0. \tag{7}$$

**Remark.** How to understand our hypothesis testing in terms of the mediation analysis? In order to do this, we have to introduce the model which states the relationship between  $G_i$  and  $(\mathbf{X}_i, \mathbf{S}_i)$ :

$$G_i = \mathbf{X}_i^{\mathsf{T}} \boldsymbol{\phi} + \mathbf{S}_i^{\mathsf{T}} \boldsymbol{\delta} + \varepsilon_i, \tag{8}$$

where  $\varepsilon_i \sim N(0, \sigma_G^2)$ . Now we are ready to calculate the TE, DE and IE by the following formula

$$\log OR_{\mathbf{s},\mathbf{s}^*|\mathbf{x}}^{TE} = (\mathbf{s} - \mathbf{s}^*)^{\mathsf{T}} \{ \boldsymbol{\beta}_S + \boldsymbol{\beta}_G \boldsymbol{\delta} + \boldsymbol{\gamma} (\mathbf{x}^{\mathsf{T}} \boldsymbol{\phi} + \mathbf{s}^{*\mathsf{T}} \boldsymbol{\delta} + \boldsymbol{\beta}_G \sigma_G^2) + \boldsymbol{\delta} \mathbf{s}^{\mathsf{T}} \boldsymbol{\gamma} \} + \frac{1}{2} \sigma_G^2 (\mathbf{s} + \mathbf{s}^*)^{\mathsf{T}} \boldsymbol{\gamma} (\mathbf{s} - \mathbf{s}^*)^{\mathsf{T}} \boldsymbol{\gamma}, \log OR_{\mathbf{s},\mathbf{s}^*|\mathbf{x}}^{DE} (\mathbf{s}) = (\mathbf{s} - \mathbf{s}^*)^{\mathsf{T}} \{ \boldsymbol{\beta}_S + \boldsymbol{\gamma} (\mathbf{x}^{\mathsf{T}} \boldsymbol{\phi} + \mathbf{s}^{*\mathsf{T}} \boldsymbol{\delta} + \boldsymbol{\beta}_G \sigma_G^2) \} + \frac{1}{2} \sigma_G^2 (\mathbf{s} + \mathbf{s}^*)^{\mathsf{T}} \boldsymbol{\gamma} (\mathbf{s} - \mathbf{s}^*)^{\mathsf{T}} \boldsymbol{\gamma}, \log OR_{\mathbf{s},\mathbf{s}^*|\mathbf{x}}^{IE} (\mathbf{s}^*) = (\mathbf{s} - \mathbf{s}^*)^{\mathsf{T}} \{ \boldsymbol{\beta}_S + \boldsymbol{\delta} \mathbf{s}^{\mathsf{T}} \boldsymbol{\gamma} \}.$$
(9)

Then we have the following explanation:

- 1. If  $\delta \neq \mathbf{0}$ , i.e. the gene expression is associated with the SNPs, then  $H_0 : \beta_{\mathbf{S}} = \mathbf{0}, \beta_G = 0, \gamma = \mathbf{0} \Leftrightarrow H_0 : DE = 0, IE = 0$ . The test for the joint effects of SNPs in a SNP set and a gene expression on disease, i.e. the total effect of a gene, is equivalent to a test for the total SNP effect on the disease.
- 2. If  $\boldsymbol{\delta} = \mathbf{0}$ , i.e. the gene expression is not associated with the SNPs, then  $H_0: \boldsymbol{\beta}_{\mathbf{S}} = \mathbf{0}, \boldsymbol{\beta}_G = 0, \boldsymbol{\gamma} = \mathbf{0}$  simply ties to evaluate the joint effect of SNP-set, gene expression and the interaction on disease.

By the above general methodology for the variance component testing, we have the following score statistics for  $\tau_S, \tau_I, \beta_G$ 

$$U_S = (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0)^{\mathsf{T}} \mathbf{S} \mathbf{S}^{\mathsf{T}} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0), U_I = (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0)^{\mathsf{T}} \mathbf{C} \mathbf{C}^{\mathsf{T}} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0), U_G = (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0)^{\mathsf{T}} \mathbf{G} \mathbf{G}^{\mathsf{T}} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0),$$
(10)

where  $\hat{\mu}_0$  is the MEL of the mean of **Y** under the null hypothesis  $H_0$  by using the standard method in GLM.

The authors proposed the following weighted sum of three scores as the test statistic for the null hypothesis:

$$Q = \frac{1}{n} \{ w_1 U_S + w_2 U_G + w_3 U_I \}, \tag{11}$$

where the weights  $w_1, w_2, w_3$  are chosen to be the inverse of the standard deviation of the  $U_S, U_G, U_I$ , which can be calculated explicitly (Lin, 1997) as follows

$$1/w_1^2 = \mathbf{1}^{\mathsf{T}}(\mathbf{SS}^{\mathsf{T}} \cdot \mathbf{K} \cdot \mathbf{SS}^{\mathsf{T}})\mathbf{1}, 1/w_2^2 = \mathbf{1}^{\mathsf{T}}(\mathbf{GG}^{\mathsf{T}} \cdot \mathbf{K} \cdot \mathbf{GG}^{\mathsf{T}})\mathbf{1}, 1/w_3^2 = \mathbf{1}^{\mathsf{T}}(\mathbf{CC}^{\mathsf{T}} \cdot \mathbf{K} \cdot \mathbf{CC}^{\mathsf{T}})\mathbf{1},$$
(12)

where  $\mathbf{K} = ((K_{ij}))$  with  $K_{ii} = -4\hat{\mu}_{0i}^4 + 8\hat{\mu}_{0i}^3 - 5\hat{\mu}_{0i}^2 + \hat{\mu}_{0i}$  and for  $i \neq j$ ,  $K_{ij} = 2[\hat{\mu}_{0i}(1 - \hat{\mu}_{0j})][\hat{\mu}_{0j}(1 - \hat{\mu}_{0j})]$ . Here  $\mathbf{A} \cdot \mathbf{B}$  denotes the component-wise multiplication.

Now the asymptotic distribution of our test statistic Q under the null hypothesis can be derived as follows.

$$Q = \frac{1}{n} \{ w_1 U_S + w_2 U_G + w_3 U_I \} = \frac{1}{n} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0)^{\mathsf{T}} \{ w_1 \mathbf{S} \mathbf{S}^{\mathsf{T}} + w_2 \mathbf{G} \mathbf{G}^{\mathsf{T}} + w_3 \mathbf{C} \mathbf{C}^{\mathsf{T}} \} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0)$$
$$= \frac{1}{n} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0)^{\mathsf{T}} (\sqrt{w_1} \mathbf{S}, \sqrt{w_2} \mathbf{G}, \sqrt{w_3} \mathbf{C}) (\sqrt{w_1} \mathbf{S}, \sqrt{w_2} \mathbf{G}, \sqrt{w_3} \mathbf{C})^{\mathsf{T}} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0)$$
$$= \frac{1}{n} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0)^{\mathsf{T}} \mathbf{V} \mathbf{V}^{\mathsf{T}} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0) = \| \frac{1}{\sqrt{n}} \mathbf{V}^{\mathsf{T}} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0) \|_2^2 = \| \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{V}_i (Y_i - \hat{\boldsymbol{\mu}}_{0i}) \|_2^2$$

where  $\mathbf{V}_{i} = (\sqrt{w_{1}}\mathbf{S}_{i}^{\mathsf{T}}, \sqrt{w_{2}}G_{i}, \sqrt{w_{3}}\mathbf{C}_{i}^{\mathsf{T}})^{\mathsf{T}} \in \mathbb{R}^{2p+1}$ . Note that for  $S_{U}(\boldsymbol{\theta}) := \frac{1}{\sqrt{n}}\sum_{i=1}^{n}\mathbf{U}_{i}(Y_{i}-\mu_{i})$  with  $\mathbf{U}_{i} = (\mathbf{X}_{i}^{\mathsf{T}}, \mathbf{V}_{i}^{\mathsf{T}}) \in \mathbb{R}^{q+2p+1}$ , where  $\boldsymbol{\theta} = (\boldsymbol{\alpha}^{\mathsf{T}}, \boldsymbol{\beta}_{\mathbf{S}}^{\mathsf{T}}, \beta_{G}, \boldsymbol{\gamma}^{\mathsf{T}})^{\mathsf{T}}$  with  $\boldsymbol{\theta}^{0} = (\boldsymbol{\alpha}^{0\mathsf{T}}, \mathbf{0}^{\mathsf{T}}, 0, \mathbf{0}^{\mathsf{T}})^{\mathsf{T}}$ , from the GLM theory, we have that (not rigorous)

$$S_U(\boldsymbol{\theta}^0) \stackrel{d}{\to} \mathcal{N}(\mathbf{0}, \mathbf{D})$$
 (13)

where  $\mathbf{D} = \begin{pmatrix} \mathbf{D}_{XX} & \mathbf{D}_{XV} \\ \mathbf{D}_{VX} & \mathbf{D}_{VV} \end{pmatrix} = n^{-1} \mathbf{U}^{\mathsf{T}} \mathbf{W} \mathbf{U}$  with  $\mathbf{W} = \text{diag}\{\mu_i(1-\mu_i)\}.$ 

By Taylor expansion, we have

$$\frac{1}{\sqrt{n}}\sum_{i=1}^{n}\mathbf{V}_{i}(Y_{i}-\hat{\mu}_{0i})\approx\mathbf{A}S_{U}(\boldsymbol{\theta}^{0})$$

where  $\mathbf{A} = (-\mathbf{D}_{XV}^{\mathsf{T}}\mathbf{D}_{XX}^{-1}, \mathbf{I}_{2p+1})$ , and then

$$\frac{1}{\sqrt{n}} \sum_{i=1}^{n} \mathbf{V}_{i}(Y_{i} - \hat{\mu}_{0i}) \xrightarrow{d|H_{0}} \mathcal{N}(\mathbf{0}, \mathbf{A}\mathbf{D}\mathbf{A}^{\mathsf{T}}),$$
(14)

which implies the following mixture of  $\chi^2$  asymptotic distribution

$$Q = \|\frac{1}{\sqrt{n}} \sum_{i=1}^{n} \mathbf{V}_{i} (Y_{i} - \hat{\mu}_{0i})\|_{2}^{2} \xrightarrow{d|H_{0}} \sum_{l=1}^{2p+1} (\mathbf{A}_{l}^{\mathsf{T}} \boldsymbol{\epsilon})^{2} := Q_{0}$$
(15)

where  $\mathbf{A}_l$  is the *l*-th row of  $\mathbf{A}$  and  $\boldsymbol{\epsilon} \sim \mathrm{N}(\mathbf{0}, \mathbf{D})$ .

And people use a scaled  $\chi^2$  distribution by matching the first two moments to approximate the mixture of  $\chi^2$  distribution, i.e.  $Q_0 \approx \kappa \chi^2_{\nu}$ , with  $\kappa = \operatorname{Var}(Q_0)/(2\mathbb{E}(Q_0))$  and  $\nu = 2[\mathbb{E}(Q_0)]^2 / \operatorname{Var}(Q_0).$ 

The paper want to demonstrate the following points:

- 1. For eQTL SNPs, i.e.  $\delta \neq 0$ , the null hypothesis is equivalent to no total genetic effect, which is defined above.
- 2. For non eQTL SNPs, i.e.  $\delta = 0$ , then there is no such equivalence.

3. For traditional SNP only genetic analysis, we use the following model

$$\operatorname{logit}\{\mathbb{P}(Y_i = 1 | \mathbf{S}_i, \mathbf{X}_i)\} = \mathbf{X}_i^{\mathsf{T}} \boldsymbol{\alpha}^* + \mathbf{S}_i^{\mathsf{T}} \boldsymbol{\beta}_{\mathbf{S}}^*, i = 1, 2, \cdots, n,$$
(16)

for testing  $H_0: \beta_{\mathbf{S}}^* = \mathbf{0}$ . If our integrative model is true without interaction, then we have the following random intercept true model

$$\operatorname{logit}\{\mathbb{P}(Y_i=1|\mathbf{S}_i,\mathbf{X}_i,\epsilon_i)\} = \mathbf{X}_i^{\mathsf{T}}(\boldsymbol{\alpha}+\beta_G\boldsymbol{\phi}) + \mathbf{S}_i^{\mathsf{T}}(\boldsymbol{\beta}_{\mathbf{S}}+\beta_G\boldsymbol{\delta}) + \beta_G\epsilon_i, i=1,2,\cdots,n,$$

which leads to by integrating the random intercept out

logit{
$$\mathbb{P}(Y_i = 1 | \mathbf{S}_i, \mathbf{X}_i)$$
}  $\approx c \Big\{ \mathbf{X}_i^{\mathsf{T}}(\boldsymbol{\alpha} + \beta_G \boldsymbol{\phi}) + \mathbf{S}_i^{\mathsf{T}}(\boldsymbol{\beta}_{\mathbf{S}} + \beta_G \boldsymbol{\delta}) \Big\}, i = 1, 2, \cdots, n.$ 

Thus testing  $H_0: \beta_{\mathbf{S}}^* = \mathbf{0}$  is approximately equivalent to testing for no total effect of the SNPs.

And if with interaction, it can be shown that the difference is that traditional approach will lose power.

### 2.2 Implementation

In order to raise the power for different alternatives such as only SNPs, (SNPs,gene expression), and (SNPs, gene expression, interaction), which correspond to the following three test statistics

$$Q_{SGC} = \frac{1}{n} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0)^{\mathsf{T}} \{ w_1 \mathbf{S} \mathbf{S}^{\mathsf{T}} + w_2 \mathbf{G} \mathbf{G}^{\mathsf{T}} + w_3 \mathbf{C} \mathbf{C}^{\mathsf{T}} \} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0),$$

$$Q_{SG} = \frac{1}{n} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0)^{\mathsf{T}} \{ w_1 \mathbf{S} \mathbf{S}^{\mathsf{T}} + w_2 \mathbf{G} \mathbf{G}^{\mathsf{T}} \} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0),$$

$$Q_S = \frac{1}{n} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0)^{\mathsf{T}} \{ w_1 \mathbf{S} \mathbf{S}^{\mathsf{T}} \} (\mathbf{Y} - \hat{\boldsymbol{\mu}}_0),$$
(17)

the authors proposed to use the minimum of the three p-values as the new test statistic. But the null distribution of the minimum p-values is hard to derive. So we resort to the score-based wild bootstrapping (Kline and Santos, 2012).

1. Generate the bootstrapped score:

$$\boldsymbol{\epsilon}^{*b} = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \mathbf{U}_i (Y_i - \hat{\mu}_{0i}) \mathcal{N}_i, b = 1, 2, \cdots, B,$$

where  $\mathcal{N}_i \stackrel{IID}{\sim} \mathcal{N}(0, 1)$ .

2. Produce three p-values: first notice that the three different  $Q_{0,SGC}, Q_{0,SG}, Q_{0,S}$ correspond to the three different matrix **A** with different  $\mathbf{V}_i = (\sqrt{w_1} \mathbf{S}_i^{\mathsf{T}}, \sqrt{w_2} G_i,$   $\sqrt{w_3} \mathbf{C}_i^{\mathsf{T}})^{\mathsf{T}}, (\sqrt{w_1} \mathbf{S}_i^{\mathsf{T}}, \sqrt{w_2} G_i)^{\mathsf{T}} \text{ and } \sqrt{w_1} \mathbf{S}_i \text{ respectively. And we denote these three different <math>\mathbf{A}$  by  $\mathbf{A}_{SGC}, \mathbf{A}_{SG}, \mathbf{A}_S$ . Then we have

$$Q_{0,SGC}^{*b} = \|\mathbf{A}_{SGC} \boldsymbol{\epsilon}^{*b}\|_{2}^{2},$$

$$Q_{0,SG}^{*b} = \|\mathbf{A}_{SG} \boldsymbol{\epsilon}^{*b}\|_{2}^{2},$$

$$Q_{0,S}^{*b} = \|\mathbf{A}_{S} \boldsymbol{\epsilon}^{*b}\|_{2}^{2},$$
(18)

which leads to the three different p-values

$$\hat{P}_{SGC} = \mathbb{P}^*(Q_{0,SGC}^* > Q_{SGC}), \hat{P}_{SG} = \mathbb{P}^*(Q_{0,SG}^* > Q_{SG}), \hat{P}_S = \mathbb{P}^*(Q_{0,S}^* > Q_S),$$
(19)

which can be estimated by  $\sum_{b=1}^{B} \mathbb{1}(Q_{0,SGC}^{*b} > Q_{SGC})/B, \sum_{b=1}^{B} \mathbb{1}(Q_{0,SG}^{*b} > Q_{SG})/B, \sum_{b=1}^{B} \mathbb{1}(Q_{0,S}^{*b} > Q_{SG})/B.$ 

3. The minimum p-value statistic is now  $\hat{P}_{\min} = \min\{\hat{P}_{SGC}, \hat{P}_{SG}, \hat{P}_{S}\}$ , whose null distribution can be approximated by the empirical distribution of

$$\{\hat{P}^{b}_{\min} = \min\{\mathbb{P}^{*}(Q^{*}_{0,SGC} > Q^{*b}_{0,SGC}), \mathbb{P}^{*}(Q^{*}_{0,SG} > Q^{*b}_{0,SG}), \mathbb{P}^{*}(Q^{*}_{0,S} > Q^{*b}_{0,S})\}\}_{b=1}^{B},$$

where the items such as  $\mathbb{P}^*(Q_{0,SGC}^* > Q_{0,SGC}^{*b})$  can be approximated by  $\sum_{b'=1}^{B} \mathbb{1}(Q_{0,SGC}^{*b'} > Q_{0,SGC}^{*b})/B$ .

4. The p-value for the minimum p-value statistic can be calculated via

$$\sum_{b=1}^{B} \mathbb{1}(\hat{P}_{\min} < \hat{P}_{\min}^{b})/B.$$

### 3 Simulation

#### 3.1 Linkage Disequilibrium

There are four types of Haplotype for two loci in Table 1. And the data structure is usually given by the Table 2.

	В	b	Total
А	$p_{AB}$	$p_{Ab}$	$p_A$
a	$p_{aB}$	$p_{ab}$	$p_a$
Total	$p_B$	$p_b$	1

Table 1: **Haplotype Frequencies.** For two loci, locus A and locus B, there are four haplotypes: AB, Ab, aB, ab.

	$BB(P_{BB})$	$Bb(P_{Bb})$	$bb(P_{bb})$
$AA(P_{AA})$	$n_{22}(p_{AB}^2)$	$n_{21}(2p_{AB}p_{Ab})$	$n_{20}(p_{Ab}^2)$
$Aa(P_{Aa})$	$n_{12}(2p_{AB}p_{aB})$	$n_{11}(2p_{AB}p_{ab}+2p_{Ab}p_{aB})$	$n_{10}(2p_{Ab}p_{ab})$
$aa(P_{aa})$	$n_{02}(p_{aB}^2)$	$n_{01}(2p_{aB}p_{ab})$	$n_{00}(p_{ab}^2)$

Table 2: Data Structure and expected genotype frequencies (assuming a random mating). For two loci, locus A and locus B, there are four haplotypes: AB, Ab, aB, ab.  $n = \sum_{i,j=0,1,2} n_{ij}$ .

Linkage Equilibrium (Expected for Distant Loci) means two loci A and B are independent, which is equivalent to the follows from the contingency table

$$p_{AB} = p_A p_B,$$

which implies the following three

$$p_{Ab} = p_A p_b, \ p_{aB} = p_a p_B, \ p_{ab} = p_a p_b.$$

Then Linkage Disequilibrium (Expected for Nearby Loci) means

$$p_{AB} \neq p_A p_B$$
.

And we have the coefficient of linkage disequilibrium (LD) between the two loci in the population,

$$D = p_{AB} - p_A p_B \tag{20}$$

which reflects the degree of LD. The larger D, the stronger LD.

Since D corresponds to the covariance between the loci, we have the following standardized version

$$r = D/\sqrt{p_A(1-p_A)p_B(1-p_B)}.$$
(21)

In summary, if we are given  $p_A, p_B$  and r, then we can calculate D, and then the haplotype frequencies such as  $p_{AB}$ . Then we will have the genotype frequencies such as  $p_{AaBB}$ . And then by the conditional arguments, we can simulate the data such as

$$P(Bb|aa) = \frac{P(aaBb)}{P(aa)} = \frac{2p_{aB}p_{ab}}{p_{aa}}.$$

### 3.2 Simulation Results

We tried the simulation for the following date generation:

We have q=1, i.e. with X only the intercept and p=10 SNPs p=10;n=c(100,200,500)

pAs=c(0.1,0.15,0.1,0.2,0.25,0.1,0.2,0.05,0.4,0.3)#MAF rs=c(0.1,0.8,0.5,0.4,-0.3,-0.2,0.2,0.15,0.4)#r k0=4#the fourth one is the only causal SNP

delta=c(0, 0.5, 1)
beta\_S=c(0, 0.1, 0.2, 0.3, 0.4, 0.5))
beta\_G=c(0, 0.2, 0.5)
gamma=c(0, 0.2, 0.5)

```
generate SNP-set with LD structure
```

```
FSNP=sample(2:0,n,replace=TRUE,c(pAs[1]^2,2*pAs[1]*(1-pAs[1]),(1-pAs[1])^2))
SNP_data=FSNP
for(j in 2:p)
   NSNP=genSNP(FSNP, rs[j-1], pAs[j-1],pAs[j]) FSNP=NSNP SNP_data=cbind(SNP_data,FSNP)
```

generate m-RNA expression data by the linear model without covariates
G=matrix(delta\*SNP\_data[,k0]+rnorm(n),ncol=1)

```
generate the response
```

```
eta=-0.2+beta_S*SNP_data[,k0]+beta_G*G+gamma*SNP_data[,k0]*G
p_response=apply(matrix(eta,ncol=1),1,function(x) exp(x)/(1+exp(x)))
Y=apply(matrix(p_response,ncol=1),1,function(x) rbinom(1,1,x))
```



Figure 2: Empirical Size and Power. We used the score-based bootstrapping here. N=500, B=500.



Figure 3: Empirical Size and Power. We used the score-based bootstrapping here. N=500, B=500.



Figure 4: Empirical Size and Power. We used the score-based bootstrapping here. N=500, B=500.



Figure 5: Empirical Size and Power. We used the score-based bootstrapping here. N=500, B=500.



Figure 6: Empirical Size and Power. We used the score-based bootstrapping here. N=500, B=500.



Figure 7: Empirical Size and Power. We used the score-based bootstrapping here. N=500, B=500.



Figure 8: Empirical Size and Power. We used the score-based bootstrapping here. N=500, B=500.

From the simulation results, we could see that

- 1. For  $\delta = 0$  case, in general, the model which is closest to the true model has the highest power.
- 2. For  $\delta \neq 0$  case, it will have more power to detect some particular alternative hypothesis that the genetic effect is indirectly through expression level. Look at Figure 5 and compare the Plot (5a) and Plot (5d).

### 4 Discussion

#### 4.1 Another way

For simplicity, we consider (note that our goal is testing the association between the outcome and a set of SNPs)

$$Y_i = \mathbf{G}_i^{\mathsf{T}} \boldsymbol{\beta}_G + \epsilon_{1i}, \tag{22}$$

$$\mathbf{G}_{i}^{\mathsf{T}}\boldsymbol{\beta}_{G} = \mathbf{S}_{i}^{\mathsf{T}}\boldsymbol{\alpha}_{S} + \epsilon_{2i} \tag{23}$$

and thus we have

$$Y_i = \mathbf{S}_i^{\mathsf{T}} \boldsymbol{\gamma}_S + \epsilon_i, \tag{24}$$

where  $\boldsymbol{\gamma}_S = \boldsymbol{\alpha}_S, \epsilon_i = \epsilon_{1i} + \epsilon_{2i} \sim N(0, \sigma_1^2 + \sigma_2^2)$ . The null hypothesis:

$$H_0: \boldsymbol{\alpha}_S = \mathbf{0}, i.e. \ \boldsymbol{\gamma}_S = \mathbf{0} \tag{25}$$

which is equivalent to no genetic effect. And this true model is reflecting our primary interest of alternative that  $\mathbf{S}_i$  is associated with  $Y_i$  through regulation of the expression of  $\mathbf{G}_i$  (Zhao et al., 2014).

If our integrative model is true, and the null hypothesis of no association between  $\mathbf{S}_i$ and  $Y_i$  is equivalent to  $\boldsymbol{\alpha}_S = \mathbf{0}$  in the integrative model and  $\boldsymbol{\gamma}_S = \mathbf{0}$  in the usual linear model. Now the question is which one is more powerful.

- 1. Usual linear model:  $\hat{\boldsymbol{\gamma}}_{S} = (\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{Y}$  with  $\operatorname{Var}(\hat{\boldsymbol{\gamma}}_{S}) = \boldsymbol{\Sigma}_{SS}^{-1}(\sigma_{1}^{2} + \sigma_{2}^{2})/n$  where  $\boldsymbol{\Sigma}_{SS} = \mathbb{E}(\mathbf{S}_{i}\mathbf{S}_{i}^{\mathsf{T}}).$
- 2. Integrative model: we first have  $\hat{\boldsymbol{\beta}}_{G} = (\mathbf{G}^{\mathsf{T}}\mathbf{G})^{-1}\mathbf{G}^{\mathsf{T}}\mathbf{Y}$  and then we have  $\hat{\boldsymbol{\alpha}}_{S} = (\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{G}\hat{\boldsymbol{\beta}}_{G}$ . So  $\operatorname{Var}(\hat{\boldsymbol{\alpha}}_{S}) = \operatorname{Var}(\hat{\boldsymbol{\alpha}}_{S} \boldsymbol{\alpha}_{S}) = \operatorname{Var}((\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{G}\hat{\boldsymbol{\beta}}_{G} \boldsymbol{\alpha}_{S}) = \operatorname{Var}((\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{G}\hat{\boldsymbol{\beta}}_{G} (\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{S}\boldsymbol{\alpha}_{S}) = \operatorname{Var}((\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{G}\hat{\boldsymbol{\beta}}_{G} (\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{G}\boldsymbol{\beta}_{G} + (\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{G}\boldsymbol{\beta}_{G} (\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{G}\boldsymbol{\beta}_{G} + (\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{G}\boldsymbol{\beta}_{G} (\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{G}\boldsymbol{\beta}_{G} + \operatorname{Var}((\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{G}\boldsymbol{\beta}_{G} (\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{G}\boldsymbol{\beta}_{G}) + \operatorname{Var}((\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{G}\boldsymbol{\beta}_{G} (\mathbf{S}^{\mathsf{T}}\mathbf{S})^{-1}\mathbf{S}^{\mathsf{T}}\mathbf{S}\boldsymbol{\alpha}_{S}) = \Sigma_{SS}^{-1}\Sigma_{SG}\Sigma_{GS}^{-1}\Sigma_{GS}\Sigma_{SS}^{-1}\sigma_{1}^{2}/n + \Sigma_{SS}^{-1}\sigma_{2}^{2}/n.$
- 3. Now for  $\mathbf{S}_i \in \mathbb{R}^p$ , we take p = 1 for illustration. So if  $\boldsymbol{\Sigma}_{SS}^{-1}\boldsymbol{\Sigma}_{SG}\boldsymbol{\Sigma}_{GG}^{-1}\boldsymbol{\Sigma}_{GS} < 1$ we will have  $\operatorname{Var}(\hat{\boldsymbol{\alpha}}_S) < \operatorname{Var}(\hat{\boldsymbol{\gamma}}_S)$ . If Gene expression are independent, then it's equivalent to  $\|\operatorname{Corr}(S_i, \mathbf{G}_i)\|_2^2 < 1$ . In other words, we gain the most power if the  $\mathbf{G}_i$  are weakly correlated with  $\mathbf{S}_i$ . This is sensible, because otherwise the expression

data would add little additional information. In the extreme case where they are perfectly correlated, our integrative analysis would be no different from a standard analysis. On the other hand, while the integrative approach has more relative power for weak correlations, its absolute power can be low if the correlations are two low, as in the extreme case where the correlation is zero we have  $\alpha_S = 0$ , i.e. null hypothesis holds. In the ideal setting, the correlations are weak but  $\alpha_S$  is still large, which only possible when  $\mathbf{G}_i$  is highly associated with  $Y_i$  so that  $\boldsymbol{\beta}_G$  is large.

### 4.2 Research Direction

- 1. How to generalize it to the nonlinear interaction situation?
- 2. Recently Haris et al. (2014) proposed how to fit the following interaction model

$$\operatorname{logit}\{\mathbb{P}(Y_i = 1\} = \mathbf{S}_i^{\mathsf{T}}\boldsymbol{\beta}_{\mathbf{S}} + \mathbf{G}_i^{\mathsf{T}}\boldsymbol{\beta}_G + \mathbf{S}_i^{\mathsf{T}}\mathbf{B}\mathbf{G}_i, i = 1, 2, \cdots, n$$
(26)

under the strong heredity condition that if an interaction term is included in the model then both of the corresponding main effects must be present or weak heredity that if an interaction term is included in the model then at least one of the corresponding main effects must be present. That is variable selection with the constraints.

## References

- Asad Haris, Daniela Witten, and Noah Simon. Convex modeling of interactions with strong heredity. arXiv preprint arXiv:1410.3517, 2014.
- Yen-Tsung Huang, Tyler J VanderWeele, and Xihong Lin. Joint analysis of snp and gene expression data in genetic association studies of complex diseases. *The annals of applied statistics*, 8(1):352, 2014.
- Patrick Kline and Andres Santos. A score based approach to wild bootstrap inference. Journal of Econometric Methods, 1(1):23–41, 2012.
- Xihong Lin. Variance component testing in generalised linear models with random effects. *Biometrika*, 84(2):309–326, 1997.
- Donald B Rubin. Formal mode of statistical inference for causal effects. Journal of Statistical Planning and Inference, 25(3):279–292, 1990.
- Tyler J VanderWeele and Stijn Vansteelandt. Odds ratios for mediation analysis for a dichotomous outcome. *American journal of epidemiology*, 172(12):1339–1348, 2010.
- Sihai D Zhao, T Tony Cai, and Hongzhe Li. More powerful genetic association testing via a new statistical framework for integrative genomics. *Biometrics*, 64:1–21, 2014. ISSN 1541-0420. doi: 10.1111/biom.12206. URL http://dx.doi.org/10.1111/biom. 12206.