

# Variance Components Estimation for Linear Mixed Models

Jing Ma

Public Health Sciences Division  
Fred Hutch Cancer Research Center

Banff Workshop  
February 5, 2019

The ACE model

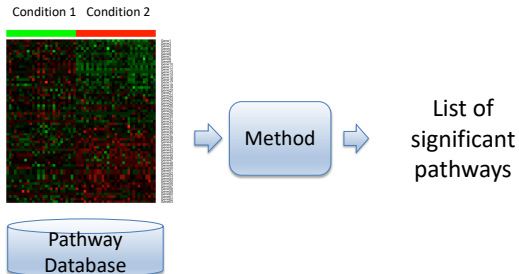
$$y = \mathbf{X}\beta + \gamma + c + e$$

$$\gamma \sim N(0, \sigma_\gamma^2 \mathbf{A}), \quad c \sim N(0, \sigma_c^2 \mathbf{C}), \quad e \sim N(0, \sigma_e^2 I_n)$$

- ▶  $y$ :  $n \times 1$  vector of quantitative traits
- ▶  $\mathbf{A}$ :  $n \times n$  genetic related matrix
- ▶  $\mathbf{C}$ :  $n \times n$  matrix for shared environment
- ▶  $h^2 = \sigma_\gamma^2 / (\sigma_\gamma^2 + \sigma_c^2 + \sigma_e^2)$

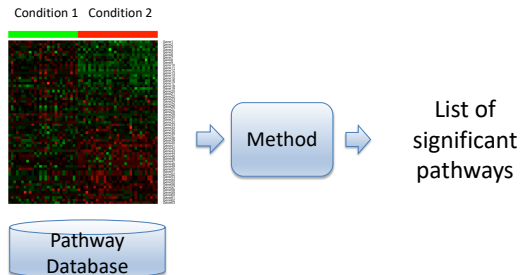
## Scientific Question

Whether a *genetic/metabolic pathway* is involved in responding to changes in environmental conditions or in specific cell functions.



## Scientific Question

Whether a *genetic/metabolic pathway* is involved in responding to changes in environmental conditions or in specific cell functions.



- ▶ Reduce the complexity; more explanatory power.

The NetGSA model<sup>1</sup> (for  $t = 1, 2, i = 1, \dots, n_t$ )

$$y_i^{(t)} = \mu^{(t)} + \gamma_i^{(t)} + e_i^{(t)}$$

$$\gamma^{(t)} \sim N(0, A_t)$$

$$e_i^{(t)} \sim N(0, \sigma_e^2 I_p)$$

- ▶  $y_i^{(t)}$ :  $p \times 1$  vector of observation for individual  $i$  in group  $t$
- ▶  $A_t^{-1}$ :  $p \times p$  network information matrix
- ▶ Test statistic for  $H_0 : \mu_G^{(1)} = \mu_G^{(2)}$  depends on the variance components.

---

<sup>1</sup> Ma J, et al. Bioinformatics. 2016

## Pros

- ▶ Statistically efficient

## Cons

- ▶ Need to invert  $n \times n$  matrices (or  $p \times p$  matrices in NetGSA) → computationally expensive, e.g.
  - ▶  $n > 100K$  in heritability estimation
  - ▶  $p \approx 3K$  in enrichment analysis

The residual after removing fixed effects is

$$\varepsilon = \mathbf{y} - \mathbf{X}\beta = \gamma + \mathbf{c} + \mathbf{e},$$

whose second moment is

$$\mathbb{E}(\varepsilon\varepsilon') = \sigma_\gamma^2 \mathbf{A} + \sigma_c^2 \mathbf{C} + \sigma_e^2 \mathbf{I}_n.$$

---

<sup>2</sup> Sofer T. Stat. Appl. Genet. Mol. Biol. 2017

The residual after removing fixed effects is

$$\varepsilon = \mathbf{y} - \mathbf{X}\beta = \gamma + \mathbf{c} + \mathbf{e},$$

whose second moment is

$$\mathbb{E}(\varepsilon\varepsilon') = \sigma_\gamma^2 \mathbf{A} + \sigma_c^2 \mathbf{C} + \sigma_e^2 \mathbf{I}_n.$$

Let  $\text{Vec}(\mathbf{A})$  denote the upper triangular part of a matrix  $\mathbf{A}$  including the diagonal. The HE method solves for  $\sigma_j^2$  by regressing

$$\tilde{\mathbf{Y}} = \text{Vec}(\hat{\varepsilon}\hat{\varepsilon}') \in \mathbb{R}^{n^*}, \quad n^* = \frac{n(n+1)}{2}$$

on the design matrix

$$\tilde{\mathbf{X}} = [\text{Vec}(\mathbf{I}_n), \text{Vec}(\mathbf{A}), \text{Vec}(\mathbf{C})] \in \mathbb{R}^{n^* \times 3}.$$

---

<sup>2</sup> Sofer T. Stat. Appl. Genet. Mol. Biol. 2017



## Pros

- ▶ Only need to invert a  $3 \times 3$  matrix

## Cons

- ▶ May get negative estimates
- ▶ Computational cost is  $O(n^2)$  in the ACE model and  $O(np^2)$  in NetGSA → can be inefficient if  $n$  and  $p$  are large

Negative estimates



Use non-negative least squares (NNLS)

NNLS solves for the variance components by minimizing

$$\frac{1}{n^*} \left\{ \theta' \tilde{X}' \tilde{X} \theta - 2\theta' \tilde{X}' \tilde{Y} \right\}, \quad \text{s.t. } \theta \geq 0,$$

where

$$\frac{1}{n^*} \tilde{X}' \tilde{X} = \frac{1}{n^*} \sum_{i=1}^{n^*} \tilde{X}'_i \tilde{X}_i,$$

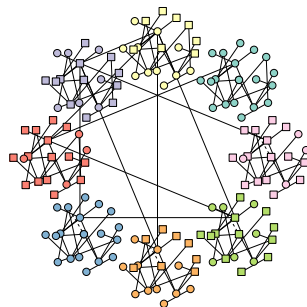
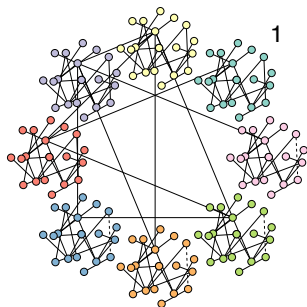
$$\frac{1}{n^*} \tilde{X}' \tilde{Y} = \frac{1}{n^*} \sum_{i=1}^{n^*} \tilde{X}'_i \tilde{Y}_i.$$

- ▶ NNLS depends only on the average products between rows of  $\tilde{X}$  and between rows of  $\tilde{X}$  and  $\tilde{Y}$ .
- ▶ We can approximate these values by subsampling rows of  $\tilde{X}$  and  $\tilde{Y}$  multiple times to get robust NNLS estimates.



- ▶ Classical CLT fails because entries in  $\tilde{Y}$  are dependent
- ▶ We evoke the CLT for weakly dependent processes to conclude consistency for the HE/REHE estimator

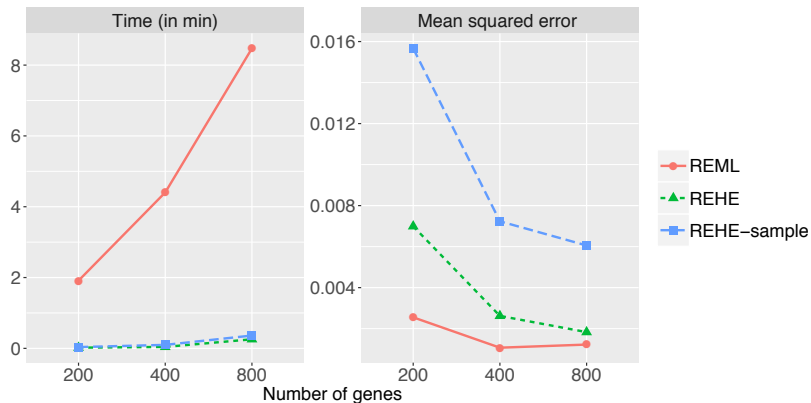
- ▶ 8 subnetworks with varying degrees of enrichment



● set 1 ● set 2 ● set 3 ● set 4 ● set 5 ● set 6 ● set 7 ● set 8

# Simulations: VC Estimation

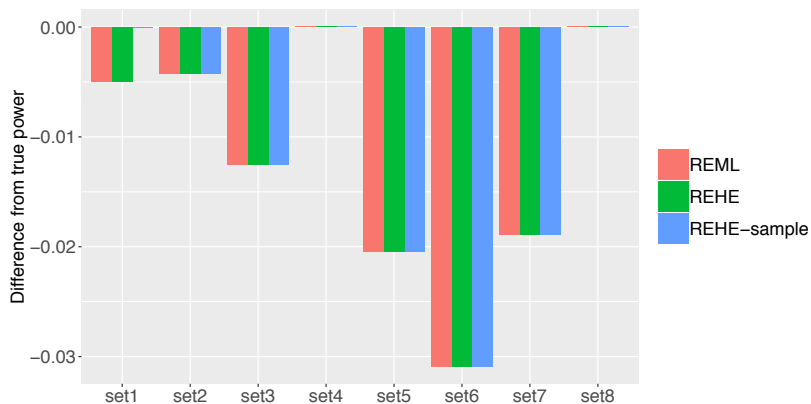
►  $n_1 = n_2 = 200, \sigma_\gamma^2 = \sigma_\epsilon^2 = 1$





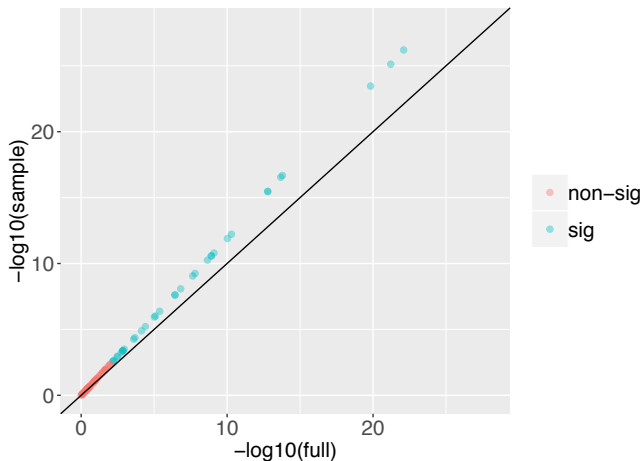
# Simulations: Power

- ▶  $p = 800, n_1 = n_2 = 200, \sigma_\gamma^2 = \sigma_\epsilon^2 = 1, \mu^{(2)} - \mu^{(1)} = 0.1$



- ▶ Gene expression from 160 normal vs 264 tumor samples
- ▶  $p = 2800$  genes with Entrez IDs
- ▶ Network topology information extracted from BioGrid
- ▶ Analysis of 96 KEGG signaling pathways

- ▶ Each dot represents one pathway; FDR  $p$ -value threshold at 0.01



The ACE model with  $\mathbf{X} \in \mathbb{R}^{n \times m}$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\gamma} + \mathbf{c} + \mathbf{e}.$$

Let  $L$  be the  $(n - m) \times n$  matrix with its rows spanning the kernel space of  $\mathbf{X}'$ .  
Then

$$\mathbb{E}[L \mathbf{y} \mathbf{y}' L'] = L(\sigma_{\boldsymbol{\gamma}}^2 \mathbf{A} + \sigma_{\mathbf{c}}^2 \mathbf{C} + \sigma_{\mathbf{e}}^2 I_n) L'.$$

The ACE model with  $\mathbf{X} \in \mathbb{R}^{n \times m}$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\gamma} + \mathbf{c} + \mathbf{e}.$$

Let  $L$  be the  $(n - m) \times n$  matrix with its rows spanning the kernel space of  $\mathbf{X}'$ .  
Then

$$\mathbb{E}[L \mathbf{y}\mathbf{y}' L'] = L(\sigma_{\boldsymbol{\gamma}}^2 \mathbf{A} + \sigma_{\mathbf{c}}^2 \mathbf{C} + \sigma_{\mathbf{e}}^2 I_n) L'.$$

- ▶  $\mathbf{y}\mathbf{y}'$  is a sample outer product derived from the Euclidean distance.

The ACE model with  $\mathbf{X} \in \mathbb{R}^{n \times m}$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\gamma} + \mathbf{c} + \mathbf{e}.$$

Let  $L$  be the  $(n - m) \times n$  matrix with its rows spanning the kernel space of  $\mathbf{X}'$ .  
Then

$$\mathbb{E}[L \mathbf{y} \mathbf{y}' L'] = L(\sigma_{\boldsymbol{\gamma}}^2 \mathbf{A} + \sigma_{\mathbf{c}}^2 \mathbf{C} + \sigma_{\mathbf{e}}^2 I_n) L'.$$

- ▶  $\mathbf{y} \mathbf{y}'$  is a sample outer product derived from the Euclidean distance.
- ▶ If we do not observe  $\mathbf{y}$  but have an outer product matrix  $M$  defined from a distance measure suitable for microbiome data, we can detect heritable microbial communities by

$$\mathbb{E}[L M L'] = L(\sigma_{\boldsymbol{\gamma}}^2 \mathbf{A} + \sigma_{\mathbf{c}}^2 \mathbf{C} + \sigma_{\mathbf{e}}^2 I_n) L'.$$

- ▶ Kun Yue (UW)
- ▶ Ali Shojaie (UW)