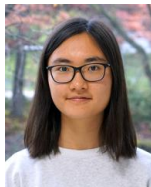


Statistical Methods for Analysis of Correlated Data

Jing Ma

Division of Public Health Sciences
Fred Hutchinson Cancer Center

Genetic Analysis of Mendelian and Complex Disorders Course
27 July 2022



Kun Yue



Mike Hellstern

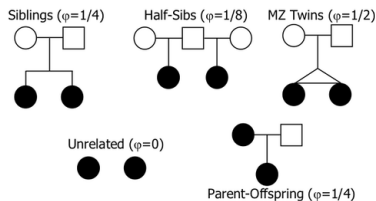


Ali Shojaie



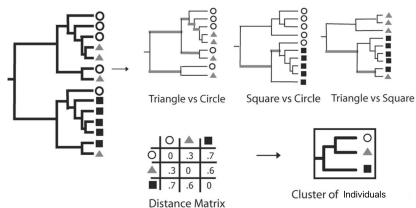
George Michailidis

Correlation among Samples



Genetic relatedness

Human genetics: kinship

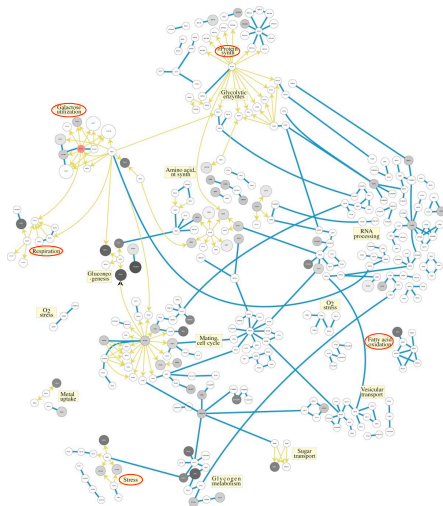


Phylogenetic relatedness¹

Microbiome: phylogenetic distances

¹Lozupone and Knight. *Appl Environ Microbiol.* '05

Correlation among Variables



- ▶ Nodes: genes
- ▶ Edges: protein → DNA and protein – protein
- ▶ Genes form functional modules

Fig: Integrated physical interaction network in yeast *Saccharomyces cerevisiae*².

²Ideker et al. *Science*. 01'

Genome-wide Association Analysis

Gene Set Analysis

Scientific Question: to identify associations of genotypes with phenotypes.

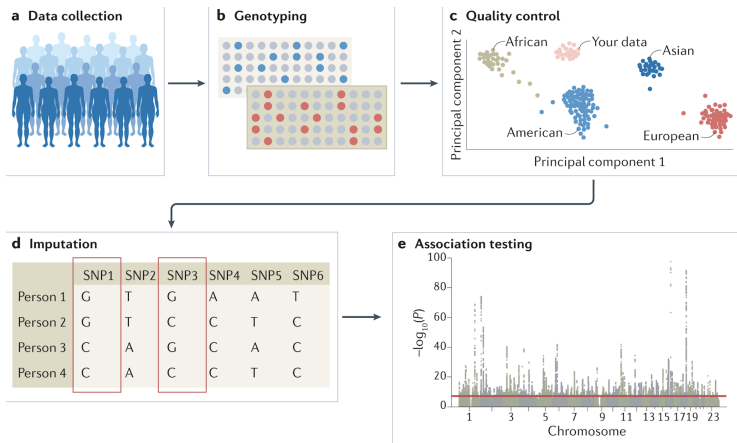


Fig: Steps of a GWAS experiment³.

³Uffelmann et al. *Nat Rev Methods Primers*. '21

Statistical model

$$y = W\alpha + X_s\beta_s + \gamma + \epsilon$$

$$\gamma \sim N(0, \sigma_\gamma^2 K)$$

$$\epsilon \sim N(0, \sigma_\epsilon^2 I_n)$$

This is a linear mixed model where

Statistical model

$$y = W\alpha + X_s\beta_s + \gamma + \epsilon$$

$$\gamma \sim N(0, \sigma_\gamma^2 K)$$

$$\epsilon \sim N(0, \sigma_\epsilon^2 I_n)$$

This is a linear mixed model where

- ▶ y : an $n \times 1$ vector of quantitative traits (e.g., red blood cell count)

Statistical model

$$y = W\alpha + X_s\beta_s + \gamma + \epsilon$$

$$\gamma \sim N(0, \sigma_\gamma^2 K)$$

$$\epsilon \sim N(0, \sigma_\epsilon^2 I_n)$$

This is a linear mixed model where

- ▶ y : an $n \times 1$ vector of quantitative traits (e.g., red blood cell count)
- ▶ W : a matrix of covariates (e.g. age, sex, ancestry)

Statistical model

$$y = W\alpha + X_s\beta_s + \gamma + \epsilon$$

$$\gamma \sim N(0, \sigma_\gamma^2 K)$$

$$\epsilon \sim N(0, \sigma_\epsilon^2 I_n)$$

This is a linear mixed model where

- ▶ y : an $n \times 1$ vector of quantitative traits (e.g., red blood cell count)
- ▶ W : a matrix of covariates (e.g. age, sex, ancestry)
- ▶ X_s : an $n \times 1$ vector of genotype values at SNP s

Statistical model

$$y = W\alpha + X_s\beta_s + \gamma + \epsilon$$

$$\gamma \sim N(0, \sigma_\gamma^2 K)$$

$$\epsilon \sim N(0, \sigma_\epsilon^2 I_n)$$

This is a linear mixed model where

- ▶ y : an $n \times 1$ vector of quantitative traits (e.g., red blood cell count)
- ▶ W : a matrix of covariates (e.g. age, sex, ancestry)
- ▶ X_s : an $n \times 1$ vector of genotype values at SNP s
- ▶ β_s : the strength of association between SNP s and y

Statistical model

$$y = W\alpha + X_s\beta_s + \gamma + \epsilon$$

$$\gamma \sim N(0, \sigma_\gamma^2 K)$$

$$\epsilon \sim N(0, \sigma_\epsilon^2 I_n)$$

This is a linear mixed model where

- ▶ y : an $n \times 1$ vector of quantitative traits (e.g., red blood cell count)
- ▶ W : a matrix of covariates (e.g. age, sex, ancestry)
- ▶ X_s : an $n \times 1$ vector of genotype values at SNP s
- ▶ β_s : the strength of association between SNP s and y
- ▶ γ : a random effect that captures the polygenic effect of other SNPs

Statistical model

$$y = W\alpha + X_s\beta_s + \gamma + \epsilon$$

$$\gamma \sim N(0, \sigma_\gamma^2 K)$$

$$\epsilon \sim N(0, \sigma_\epsilon^2 I_n)$$

This is a linear mixed model where

- ▶ y : an $n \times 1$ vector of quantitative traits (e.g., red blood cell count)
- ▶ W : a matrix of covariates (e.g. age, sex, ancestry)
- ▶ X_s : an $n \times 1$ vector of genotype values at SNP s
- ▶ β_s : the strength of association between SNP s and y
- ▶ γ : a random effect that captures the polygenic effect of other SNPs
- ▶ K : $n \times n$ kinship matrix

Input data: (W, X_s, y, K)

Association testing

$$H_0 : \beta_s = 0$$

Heritability estimation

$$h^2 = \frac{\sigma_\gamma^2}{\sigma_\gamma^2 + \sigma_\epsilon^2}$$

$$y = W\alpha + X_s\beta_s + u,$$

where $u \sim \mathcal{N}(0, \sigma_\gamma^2 V)$ and $V = K + \sigma_\epsilon^2 / \sigma_\gamma^2 I_n$.

$$y = W\alpha + X_s\beta_s + u,$$

where $u \sim \mathcal{N}(0, \sigma_\gamma^2 V)$ and $V = K + \sigma_\epsilon^2 / \sigma_\gamma^2 I_n$.

Generalized least squares

$$\begin{bmatrix} \hat{\alpha} \\ \hat{\beta}_s \end{bmatrix} = (X^T V^{-1} X)^{-1} X^T V^{-1} y$$

$$y = W\alpha + X_s\beta_s + u,$$

where $u \sim \mathcal{N}(0, \sigma_\gamma^2 V)$ and $V = K + \sigma_\epsilon^2 / \sigma_\gamma^2 I_n$.

Generalized least squares

$$\begin{bmatrix} \hat{\alpha} \\ \hat{\beta}_s \end{bmatrix} = (X^T V^{-1} X)^{-1} X^T V^{-1} y$$

Both analysis tasks require estimating the variance components!

Let $Z \in \{0, 1, 2\}^{n \times q}$ denote the remaining q SNPs (i.e. excluding X_s).

Let $Z \in \{0, 1, 2\}^{n \times q}$ denote the remaining q SNPs (i.e. excluding X_s).
Consider the model

$$y = W\alpha + X_s\beta_s + Z\tau + \epsilon$$

Let $Z \in \{0, 1, 2\}^{n \times q}$ denote the remaining q SNPs (i.e. excluding X_s).
Consider the model

$$y = W\alpha + X_s\beta_s + Z\tau + \epsilon$$

The coefficient β_s is the effect of SNP s on y after adjusting the effects introduced by other SNPs Z .

Let $Z \in \{0, 1, 2\}^{n \times q}$ denote the remaining q SNPs (i.e. excluding X_s).
Consider the model

$$y = W\alpha + X_s\beta_s + Z\tau + \epsilon$$

The coefficient β_s is the effect of SNP s on y after adjusting the effects introduced by other SNPs Z .

Fitting the model

$$\tau_j \sim \mathcal{N}(0, \sigma_\tau^2), \quad j = 1, \dots, q$$

Let $Z \in \{0, 1, 2\}^{n \times q}$ denote the remaining q SNPs (i.e. excluding X_s).
Consider the model

$$y = W\alpha + X_s\beta_s + Z\tau + \epsilon$$

The coefficient β_s is the effect of SNP s on y after adjusting the effects introduced by other SNPs Z .

Fitting the model

$$\tau_j \sim \mathcal{N}(0, \sigma_\gamma^2), \quad j = 1, \dots, q$$

Averaging over the distribution of τ_j 's, we obtain

$$y \sim \mathcal{N}(W\alpha + X_s\beta_s, \sigma_\gamma^2 ZZ^T + \sigma_\epsilon^2 I_n)$$

Let $Z \in \{0, 1, 2\}^{n \times q}$ denote the remaining q SNPs (i.e. excluding X_s).
Consider the model

$$y = W\alpha + X_s\beta_s + Z\tau + \epsilon$$

The coefficient β_s is the effect of SNP s on y after adjusting the effects introduced by other SNPs Z .

Fitting the model

$$\tau_j \sim \mathcal{N}(0, \sigma_\gamma^2), \quad j = 1, \dots, q$$

Averaging over the distribution of τ_j 's, we obtain

$$y \sim \mathcal{N}(W\alpha + X_s\beta_s, \sigma_\gamma^2 ZZ^T + \sigma_\epsilon^2 I_n)$$

The kinship $K = ZZ^T$ is a natural choice.

Maximum likelihood (null)

$$\max_{\sigma_\gamma^2, \sigma_\varepsilon^2/\sigma_\gamma^2} \left\{ -\frac{1}{2} \log |\sigma_\gamma^2 V| - \frac{1}{2} \sigma_\gamma^{-2} (y - W\hat{\alpha})^\top V^{-1} (y - W\hat{\alpha}) \right\}$$

Maximum likelihood (null)

$$\max_{\sigma_{\gamma}^2, \sigma_{\epsilon}^2 / \sigma_{\gamma}^2} \left\{ -\frac{1}{2} \log |\sigma_{\gamma}^2 V| - \frac{1}{2} \sigma_{\gamma}^{-2} (y - W\hat{\alpha})^T V^{-1} (y - W\hat{\alpha}) \right\}$$

Restricted maximum likelihood (REML)

$$\max_{\sigma_{\gamma}^2, \sigma_{\epsilon}^2 / \sigma_{\gamma}^2} \{ \text{likelihood of } L^T y \}$$

where $L^T W = 0$ and L^T has full row rank.

Maximum likelihood (null)

$$\max_{\sigma_{\gamma}^2, \sigma_{\epsilon}^2 / \sigma_{\gamma}^2} \left\{ -\frac{1}{2} \log |\sigma_{\gamma}^2 V| - \frac{1}{2} \sigma_{\gamma}^{-2} (y - W\hat{\alpha})^T V^{-1} (y - W\hat{\alpha}) \right\}$$

Restricted maximum likelihood (REML)

$$\max_{\sigma_{\gamma}^2, \sigma_{\epsilon}^2 / \sigma_{\gamma}^2} \{ \text{likelihood of } L^T y \}$$

where $L^T W = 0$ and L^T has full row rank.

☺ REML estimator has the **smallest variance** among all estimators

Maximum likelihood (null)

$$\max_{\sigma_\gamma^2, \sigma_\epsilon^2/\sigma_\gamma^2} \left\{ -\frac{1}{2} \log |\sigma_\gamma^2 V| - \frac{1}{2} \sigma_\gamma^{-2} (y - W\hat{\alpha})^\top V^{-1} (y - W\hat{\alpha}) \right\}$$

Restricted maximum likelihood (REML)

$$\max_{\sigma_\gamma^2, \sigma_\epsilon^2/\sigma_\gamma^2} \{ \text{likelihood of } L^\top y \}$$

where $L^\top W = 0$ and L^\top has full row rank.

- ☺ REML estimator has the **smallest variance** among all estimators
- ☹ REML is **computationally expensive**: need to invert $n \times n$ matrices where $n > 100K$ in large studies

Maximum likelihood (null)

$$\max_{\sigma_\gamma^2, \sigma_\epsilon^2/\sigma_\gamma^2} \left\{ -\frac{1}{2} \log |\sigma_\gamma^2 V| - \frac{1}{2} \sigma_\gamma^{-2} (y - W\hat{\alpha})^\top V^{-1} (y - W\hat{\alpha}) \right\}$$

Restricted maximum likelihood (REML)

$$\max_{\sigma_\gamma^2, \sigma_\epsilon^2/\sigma_\gamma^2} \{ \text{likelihood of } L^\top y \}$$

where $L^\top W = 0$ and L^\top has full row rank.

- ☺ REML estimator has the **smallest variance** among all estimators
- ☹ REML is **computationally expensive**: need to invert $n \times n$ matrices where $n > 100K$ in large studies

Need alternatives that can balance statistical and computational efficiency.

Assume no fixed effects for the moment. The model is

$$y = \gamma + \epsilon.$$

Assume no fixed effects for the moment. The model is

$$y = \gamma + \epsilon.$$

The second moment of y is

$$E(yy^T) = \sigma_\gamma^2 K + \sigma_\epsilon^2 I_n.$$

Assume no fixed effects for the moment. The model is

$$y = \gamma + \epsilon.$$

The second moment of y is

$$E(yy^T) = \sigma_\gamma^2 K + \sigma_\epsilon^2 I_n.$$

yy^T is a linear function of K and I_n !

Let $\text{vec}(K)$ denote the vectorization of K by stacking its columns. Let $n^* = n^2$ and

$$\tilde{Y} = \text{vec}(yy^T) \in \mathbb{R}^{n^*}, \quad \tilde{X} = [\text{vec}(I_n), \text{vec}(K)] \in \mathbb{R}^{n^* \times 2}.$$

⁴Haseman and Elston. *Behavior Genetics*. '72; Sofer T. *Stat. Appl. Genet. Mol. Biol.* '17

Let $\text{vec}(K)$ denote the vectorization of K by stacking its columns. Let $n^* = n^2$ and

$$\tilde{Y} = \text{vec}(yy^T) \in \mathbb{R}^{n^*}, \quad \tilde{X} = [\text{vec}(I_n), \text{vec}(K)] \in \mathbb{R}^{n^* \times 2}.$$

HE regression⁴ solves for σ_j^2 by minimizing

$$\frac{1}{n^*} (\tilde{Y} - \tilde{X}\sigma^2)^T (\tilde{Y} - \tilde{X}\sigma^2)$$

⁴Haseman and Elston. *Behavior Genetics*. '72; Sofer T. *Stat. Appl. Genet. Mol. Biol.* '17

Let $\text{vec}(K)$ denote the vectorization of K by stacking its columns. Let $n^* = n^2$ and

$$\tilde{Y} = \text{vec}(yy^T) \in \mathbb{R}^{n^*}, \quad \tilde{X} = [\text{vec}(I_n), \text{vec}(K)] \in \mathbb{R}^{n^* \times 2}.$$

HE regression⁴ solves for σ_j^2 by minimizing

$$\frac{1}{n^*} (\tilde{Y} - \tilde{X}\sigma^2)^T (\tilde{Y} - \tilde{X}\sigma^2)$$

- ☺ The HE estimator is **unbiased**.
- ☺ HE is **computationally efficient**: $O(dn^2)$ as opposed to $O(n^3)$ for REML

⁴Haseman and Elston. *Behavior Genetics*. '72; Sofer T. *Stat. Appl. Genet. Mol. Biol.* '17

Let $\text{vec}(K)$ denote the vectorization of K by stacking its columns. Let $n^* = n^2$ and

$$\tilde{Y} = \text{vec}(yy^T) \in \mathbb{R}^{n^*}, \quad \tilde{X} = [\text{vec}(I_n), \text{vec}(K)] \in \mathbb{R}^{n^* \times 2}.$$

HE regression⁴ solves for σ_j^2 by minimizing

$$\frac{1}{n^*} (\tilde{Y} - \tilde{X}\sigma^2)^T (\tilde{Y} - \tilde{X}\sigma^2)$$

- ☺ The HE estimator is **unbiased**.
- ☺ HE is **computationally efficient**: $O(dn^2)$ as opposed to $O(n^3)$ for REML
- ☹ **May get negative estimates**: truncation to zero?

⁴Haseman and Elston. *Behavior Genetics*. '72; Sofer T. *Stat. Appl. Genet. Mol. Biol.* '17

☺ Avoid negative estimates by **non-negative least squares (NNLS)**

☺ Avoid negative estimates by **non-negative least squares (NNLS)**

REHE solves for the variance components by minimizing

$$\frac{1}{n^*} (\tilde{Y} - \tilde{X}\sigma^2)^\top (\tilde{Y} - \tilde{X}\sigma^2) = \frac{1}{n^*} \left\{ (\sigma^2)^\top \tilde{X}^\top \tilde{X} \sigma^2 - 2(\sigma^2)^\top \tilde{X}^\top \tilde{Y} \right\},$$

subject to $\sigma^2 \geq 0$.

☺ Avoid negative estimates by **non-negative least squares (NNLS)**

REHE solves for the variance components by minimizing

$$\frac{1}{n^*} (\tilde{Y} - \tilde{X}\sigma^2)^\top (\tilde{Y} - \tilde{X}\sigma^2) = \frac{1}{n^*} \left\{ (\sigma^2)^\top \tilde{X}^\top \tilde{X} \sigma^2 - 2(\sigma^2)^\top \tilde{X}^\top \tilde{Y} \right\},$$

subject to $\sigma^2 \geq 0$.

☺ Global minimizer is guaranteed due to convexity.

- ☺ Avoid negative estimates by **non-negative least squares (NNLS)**

REHE solves for the variance components by minimizing

$$\frac{1}{n^*} (\tilde{Y} - \tilde{X}\sigma^2)^\top (\tilde{Y} - \tilde{X}\sigma^2) = \frac{1}{n^*} \left\{ (\sigma^2)^\top \tilde{X}^\top \tilde{X} \sigma^2 - 2(\sigma^2)^\top \tilde{X}^\top \tilde{Y} \right\},$$

subject to $\sigma^2 \geq 0$.

- ☺ Global minimizer is guaranteed due to convexity.
- ☺ Computational cost of REHE is comparable to HE, both faster than REML.

- ☺ Avoid negative estimates by **non-negative least squares (NNLS)**

REHE solves for the variance components by minimizing

$$\frac{1}{n^*} (\tilde{Y} - \tilde{X}\sigma^2)^\top (\tilde{Y} - \tilde{X}\sigma^2) = \frac{1}{n^*} \left\{ (\sigma^2)^\top \tilde{X}^\top \tilde{X} \sigma^2 - 2(\sigma^2)^\top \tilde{X}^\top \tilde{Y} \right\},$$

subject to $\sigma^2 \geq 0$.

- ☺ Global minimizer is guaranteed due to convexity.
- ☺ Computational cost of REHE is comparable to HE, both faster than REML.
- ☹ **May get zero estimates**



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

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

We can approximate these inner products by subsampling rows of \tilde{X} and \tilde{Y} .

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

We can approximate these inner products by subsampling rows of \tilde{X} and \tilde{Y} .

REHE with Resampling (reREHE)

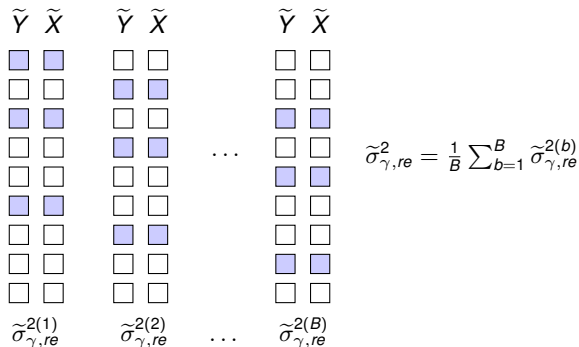
\tilde{Y}	\tilde{X}	\tilde{Y}	\tilde{X}	...	\tilde{Y}	\tilde{X}
						
						
				...		
						
						
						
						
$\tilde{\sigma}_{\gamma, re}^{2(1)}$		$\tilde{\sigma}_{\gamma, re}^{2(2)}$...	$\tilde{\sigma}_{\gamma, re}^{2(B)}$	

$$\tilde{\sigma}_{\gamma, re}^2 = \frac{1}{B} \sum_{b=1}^B \tilde{\sigma}_{\gamma, re}^{2(b)}$$

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

We can approximate these inner products by subsampling rows of \tilde{X} and \tilde{Y} .

REHE with Resampling (reREHE)



☺ reREHE estimates are **strictly positive** and can be faster to compute.

Suppose we have covariates. The null model

$$y = W\alpha + \gamma + \epsilon.$$

⁵ K^\dagger can be replaced by K when n is large.

Suppose we have covariates. The null model

$$y = W\alpha + \gamma + \epsilon.$$

Let $P^\perp = I_n - W(W^T W)^{-1} W^T$ denote the projection matrix onto the orthogonal complement of the column space of W . Let

$$y^\dagger = P^\perp y, \quad \gamma^\dagger = P^\perp \gamma, \quad \epsilon^\dagger = P^\perp \epsilon$$

⁵ K^\dagger can be replaced by K when n is large.

Suppose we have covariates. The null model

$$y = W\alpha + \gamma + \epsilon.$$

Let $P^\perp = I_n - W(W^\top W)^{-1}W^\top$ denote the projection matrix onto the orthogonal complement of the column space of W . Let

$$y^\dagger = P^\perp y, \quad \gamma^\dagger = P^\perp \gamma, \quad \epsilon^\dagger = P^\perp \epsilon$$

We obtain a new model with no covariates

$$y^\dagger = \gamma^\dagger + \epsilon^\dagger, \quad \gamma^\dagger \sim \mathcal{N}(0, \sigma_\gamma^2 K^\dagger)$$

where $K^\dagger = P^\perp K P^\perp$ ⁵.

⁵ K^\dagger can be replaced by K when n is large.

Constructing Confidence Intervals



⁶Can also construct quantile confidence interval

Parametric Bootstrap

- ▶ Compute REHE estimates $\tilde{\sigma}_\gamma^2, \tilde{\sigma}_\epsilon^2$ based on \tilde{Y}, K, I_n ;
- ▶ For $b = 1$ to B
 - ▶ Generate response vector $\tilde{Y}^{*(b)}$ from $\mathcal{N}(0, \tilde{\sigma}_\gamma^2 K + \tilde{\sigma}_\epsilon^2 I_n)$;
 - ▶ Compute REHE estimates $\tilde{\sigma}_\gamma^{2(b)}, \tilde{\sigma}_\epsilon^{2(b)}$, based on $\tilde{Y}^{*(b)}, K, I_n$;

⁶Can also construct quantile confidence interval

Parametric Bootstrap

- ▶ Compute REHE estimates $\tilde{\sigma}_\gamma^2, \tilde{\sigma}_\epsilon^2$ based on \tilde{Y}, K, I_n ;
- ▶ For $b = 1$ to B
 - ▶ Generate response vector $\tilde{Y}^{*(b)}$ from $\mathcal{N}(0, \tilde{\sigma}_\gamma^2 K + \tilde{\sigma}_\epsilon^2 I_n)$;
 - ▶ Compute REHE estimates $\tilde{\sigma}_\gamma^{2(b)}, \tilde{\sigma}_\epsilon^{2(b)}$, based on $\tilde{Y}^{*(b)}, K, I_n$;

Wald-type confidence interval⁶

$$\left[\tilde{\sigma}_\gamma^2 - z_{\alpha/2} \times \text{s.e.} \left(\tilde{\sigma}_\gamma^{2(b)} \right), \tilde{\sigma}_\gamma^2 + z_{\alpha/2} \times \text{s.e.} \left(\tilde{\sigma}_\gamma^{2(b)} \right) \right],$$

where $z_{\alpha/2}$ is the $(1 - \alpha/2)$ -th percentile of the standard normal distribution.

⁶Can also construct quantile confidence interval

- ▶ $n = 12,502$ after removing observations with missing values

- ▶ $n = 12,502$ after removing observations with missing values
- ▶ y : red blood cell count

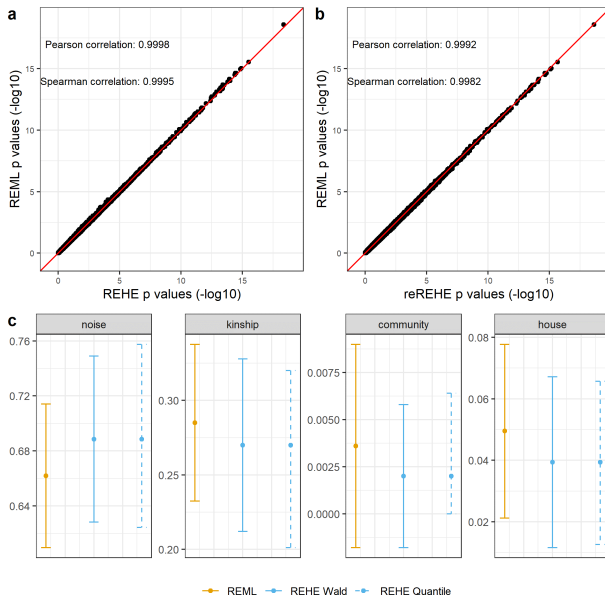
- ▶ $n = 12,502$ after removing observations with missing values
- ▶ y : red blood cell count
- ▶ X : 4,100,028 SNPs

- ▶ $n = 12,502$ after removing observations with missing values
- ▶ y : red blood cell count
- ▶ X : 4,100,028 SNPs
- ▶ Covariates W : age, sex, cigarette use, field center indicator, genetic subgroup indicator, ancestry, sampling weights

- ▶ $n = 12,502$ after removing observations with missing values
- ▶ y : red blood cell count
- ▶ X : 4,100,028 SNPs
- ▶ Covariates W : age, sex, cigarette use, field center indicator, genetic subgroup indicator, ancestry, sampling weights
- ▶ Variance components: genetic relatedness, membership of household, and membership of community group

- ▶ $n = 12,502$ after removing observations with missing values
- ▶ y : red blood cell count
- ▶ X : 4,100,028 SNPs
- ▶ Covariates W : age, sex, cigarette use, field center indicator, genetic subgroup indicator, ancestry, sampling weights
- ▶ Variance components: genetic relatedness, membership of household, and membership of community group
- ▶ Bonferroni correction for multiple testing

- ▶ $n = 12,502$ after removing observations with missing values
- ▶ y : red blood cell count
- ▶ X : 4,100,028 SNPs
- ▶ Covariates W : age, sex, cigarette use, field center indicator, genetic subgroup indicator, ancestry, sampling weights
- ▶ Variance components: genetic relatedness, membership of household, and membership of community group
- ▶ Bonferroni correction for multiple testing
- ▶ REHE took 2.4 min for estimation and 18 min for inference; REML 23.9 min



Synthetic data were generated from

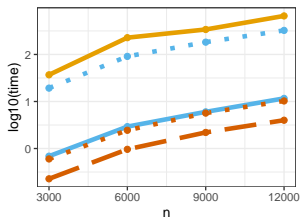
$$y = \sigma_0^2 I_n + \sigma_1^2 K_1,$$

where K_1 is a submatrix of the genetic relatedness matrix from HCHS/SOL.

- ▶ $n \in \{3,000, 6,000, 9,000, 12,000\}$
- ▶ $(\sigma_0^2, \sigma_1^2) \in \{(0.1, 0.1), (0.01, 0.1)\}$

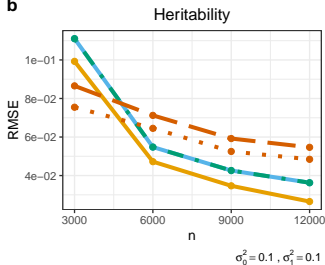
23% HE estimates were negative before truncation at zero
($n = 3000, \sigma_0^2 = 0.01$).

a



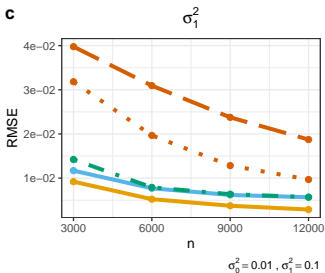
REML REHE est REHE CI reREHE 0.05 reREHE 0.1

b



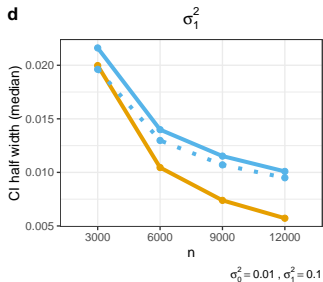
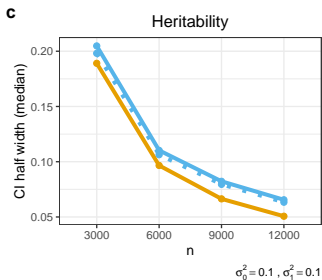
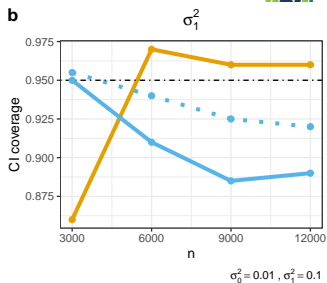
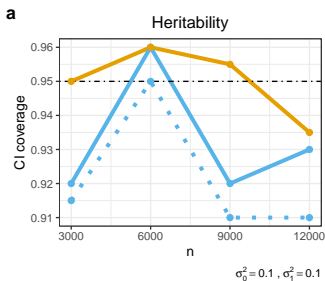
REML REHE HE reREHE 0.05 reREHE 0.1

c



REML REHE HE reREHE 0.05 reREHE 0.1

Confidence Interval Results



REML REHE Wald REHE Quantile

Genome-wide Association Analysis

Gene Set Analysis

Gene Set: a set of all SNPs located near a list of related genes.

Gene Set: a set of all SNPs located near a list of related genes.

Scientific Question: whether a *gene set* is associated with a trait.

Gene Set: a set of all SNPs located near a list of related genes.

Scientific Question: whether a *gene set* is associated with a trait.

Motivation: many biological processes are driven by mechanisms involving more than one SNP

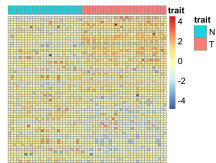
Gene Set: a set of all SNPs located near a list of related genes.

Scientific Question: whether a *gene set* is associated with a trait.

Motivation: many biological processes are driven by mechanisms involving more than one SNP

- ☺ Easy interpretation
- ☺ Fewer number of gene sets compared to number of genes/SNPs
- ☺ More power by pooling many weaker signals

Input



Pathway Database



Methods

GSEA

SPIA

DEGraph

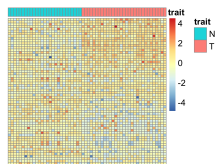
NetGSA

...

Output

List of significant pathway

Input



Pathway Database



Methods

GSEA

SPIA

DEGraph

NetGSA

...

Output

List of significant pathway

Pathway Database

KEGG, MSigDB, BioCarta, Reactome, MetaCyc, etc.

Motivation: genes are not independent

Motivation: genes are not independent

Most existing methods rely on curated interactions from pathway databases.

- ☹ Curated networks can be **incomplete** and/or **inaccurate**
- ☹ Curated networks lack **condition/disease-specific** alterations in interactions

Motivation: genes are not independent

Most existing methods rely on curated interactions from pathway databases.

- ☹ Curated networks can be **incomplete** and/or **inaccurate**
- ☹ Curated networks lack **condition/disease-specific** alterations in interactions

Which null hypothesis?

Motivation: genes are not independent

Most existing methods rely on curated interactions from pathway databases.

- ☹ Curated networks can be **incomplete** and/or **inaccurate**
- ☹ Curated networks lack **condition/disease-specific** alterations in interactions

Which null hypothesis?

- ▶ The genes in a given pathway are at most as differentially expressed as those outside the pathway (camera, PathNet).

Motivation: genes are not independent

Most existing methods rely on curated interactions from pathway databases.

- ☹ Curated networks can be **incomplete** and/or **inaccurate**
- ☹ Curated networks lack **condition/disease-specific** alterations in interactions

Which null hypothesis?

- ▶ The genes in a given pathway are at most as differentially expressed as those outside the pathway (camera, PathNet).
- ▶ The observed number of DE genes is just by chance and the DE genes are randomly located in the pathway (SPIA, Pathway-Express)

Motivation: genes are not independent

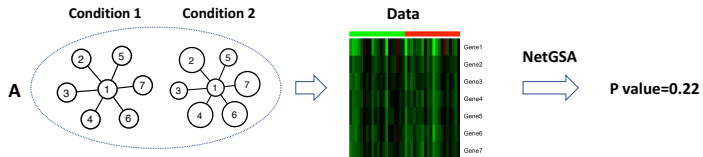
Most existing methods rely on curated interactions from pathway databases.

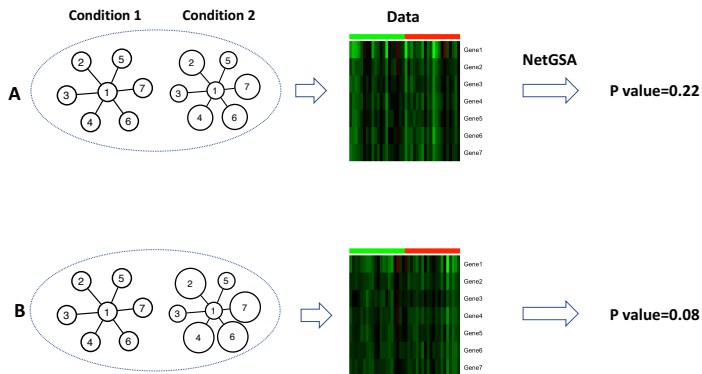
- ☹ Curated networks can be **incomplete** and/or **inaccurate**
- ☹ Curated networks lack **condition/disease-specific** alterations in interactions

Which null hypothesis?

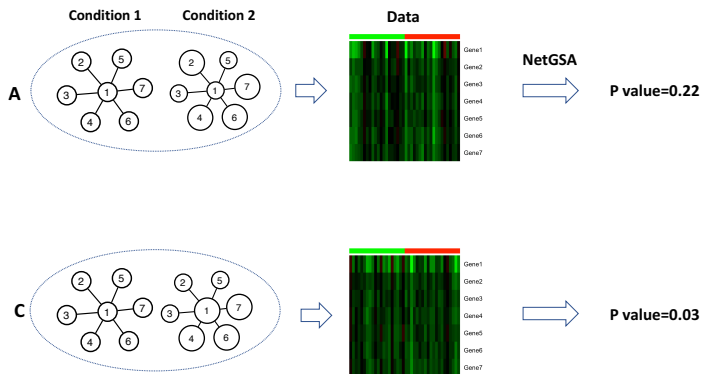
- ▶ The genes in a given pathway are at most as differentially expressed as those outside the pathway (camera, PathNet).
- ▶ The observed number of DE genes is just by chance and the DE genes are randomly located in the pathway (SPIA, Pathway-Express)
- ▶ **Self-contained null** (NetGSA, DEGraph and topologyGSA)

NetGSA - Toy Example

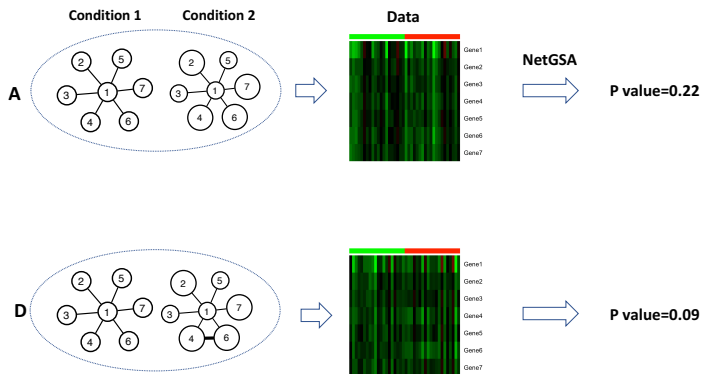




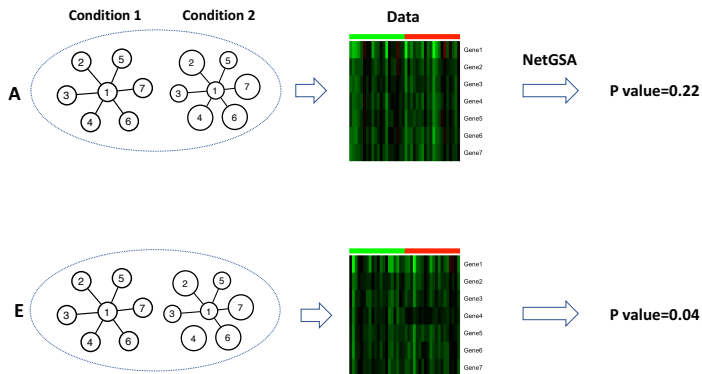
Nodes 2, 4, 6, 7 have larger changes in mean in case B than in case A.



Node 1 as opposed to node 2 has change in mean in case C.



There is an additional change in correlation between nodes 4 and 6 in case D.



There is an additional change in correlation between nodes 1 and 4 in case E.

What Drives Gene Set Significance



- ▶ Change in mean values of genes in the set
- ▶ Position of genes: hub genes are more important
- ▶ Change in gene-gene interaction

- ▶ Change in mean values of genes in the set
- ▶ Position of genes: hub genes are more important
- ▶ Change in gene-gene interaction

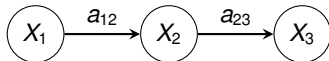
NetGSA captures all three factors!

Let $Y \in \mathbb{R}^p$ denote the expression values of p genes from an arbitrary sample. Suppose $Y = X + \epsilon$, where X is signal and ϵ is noise.

⁷Shojaie and Michailidis. *JCB*. '09

Let $Y \in \mathbb{R}^p$ denote the expression values of p genes from an arbitrary sample. Suppose $Y = X + \epsilon$, where X is signal and ϵ is noise.

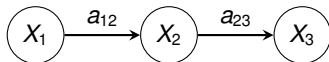
Assume the p genes are related via a network $A = (a_{ij})$ where a_{ij} denotes the strength of association between genes i and j .



⁷Shojaie and Michailidis. *JCB*. '09

Let $Y \in \mathbb{R}^p$ denote the expression values of p genes from an arbitrary sample. Suppose $Y = X + \epsilon$, where X is signal and ϵ is noise.

Assume the p genes are related via a network $A = (a_{ij})$ where a_{ij} denotes the strength of association between genes i and j .



We model X via the latent variable model⁷

$$X_1 = \gamma_1$$

$$X_2 = a_{12}X_1 + \gamma_2$$

$$X_3 = a_{23}X_2 + \gamma_3 = a_{12}a_{23}\gamma_1 + a_{23}\gamma_2 + \gamma_3$$

where $\gamma_j \sim \mathcal{N}(\mu_j, \sigma_\gamma^2)$ represents the baseline expression of gene j .

⁷Shojaie and Michailidis. *JCB*. '09

$$Y = \Lambda\gamma + \epsilon, \quad \gamma \sim \mathcal{N}(\mu, \sigma_\gamma^2 I_p), \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2 I_p)$$

where

$$\Lambda = \begin{pmatrix} 1 & 0 & 0 \\ a_{12} & 1 & 0 \\ a_{12}a_{23} & a_{23} & 1 \end{pmatrix}$$

is the **influence matrix** of the gene network $\Lambda = (I_p - A)^{-1}$.

$$Y = \Lambda\gamma + \epsilon, \quad \gamma \sim \mathcal{N}(\mu, \sigma_\gamma^2 I_p), \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2 I_p)$$

where

$$\Lambda = \begin{pmatrix} 1 & 0 & 0 \\ a_{12} & 1 & 0 \\ a_{12}a_{23} & a_{23} & 1 \end{pmatrix}$$

is the **influence matrix** of the gene network $\Lambda = (I_p - A)^{-1}$.

Statistical Inference

Given data Y_i ($i = 1, \dots, n$) and network A , test for a gene set G

$$H_0 : \mu_G^{(1)} = \mu_G^{(2)}$$

$$Y = \Lambda\gamma + \epsilon, \quad \gamma \sim \mathcal{N}(\mu, \sigma_\gamma^2 I_p), \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2 I_p)$$

where

$$\Lambda = \begin{pmatrix} 1 & 0 & 0 \\ a_{12} & 1 & 0 \\ a_{12}a_{23} & a_{23} & 1 \end{pmatrix}$$

is the **influence matrix** of the gene network $\Lambda = (I_p - A)^{-1}$.

Statistical Inference

Given data Y_i ($i = 1, \dots, n$) and network A , test for a gene set G

$$H_0 : \mu_G^{(1)} = \mu_G^{(2)}$$

or

$$H_0^{net} : (\Lambda^{(1)} \mu^{(1)})_G = (\Lambda^{(2)} \mu^{(2)})_G$$

A can be directed acyclic or undirected.

⁸Ma et al. *Bioinformatics*. '16

A can be directed acyclic or undirected.

A is weighted.

⁸Ma et al. *Bioinformatics*. '16

A can be directed acyclic or undirected.

A is weighted.

NetGSA infers the weights from data (**independent from Y**) using **graphical models**.

⁸Ma et al. *Bioinformatics*. '16

A can be directed acyclic or undirected.

A is weighted.

NetGSA infers the weights from data (independent from Y) using graphical models.

☺ Many RNA-seq data are available

⁸Ma et al. *Bioinformatics*. '16

A can be directed acyclic or undirected.

A is weighted.

NetGSA infers the weights from data (**independent from Y**) using **graphical models**.

☺ Many RNA-seq data are available

☺ Can use curated networks as **side information** to improve data-driven network inference⁸

⁸Ma et al. *Bioinformatics*. '16

$$A = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} & \begin{pmatrix} \cdot & ? & 1 & 0 & ? & 0 \\ ? & \cdot & ? & ? & 0 & ? \\ 1 & ? & \cdot & ? & 0 & 0 \\ 0 & ? & ? & \cdot & ? & 1 \\ ? & 0 & 0 & ? & \cdot & ? \\ 0 & ? & 0 & 1 & ? & \cdot \end{pmatrix} \end{matrix}$$

- 0: there is no interaction; 1: there is interaction; ?: unknown

$$A = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \end{matrix} \\ \begin{pmatrix} \cdot & ? & 1 & 0 & ? & 0 \\ ? & \cdot & ? & ? & 0 & ? \\ 1 & ? & \cdot & ? & 0 & 0 \\ 0 & ? & ? & \cdot & ? & 1 \\ ? & 0 & 0 & ? & \cdot & ? \\ 0 & ? & 0 & 1 & ? & \cdot \end{pmatrix} & \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} \end{matrix}$$

- ▶ **0**: there is no interaction; **1**: there is interaction; **?**: unknown
- ▶ Given data, we use **graphical models** to incorporate existing information using a **constrained optimization** framework.

$$A = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \end{matrix} \\ \begin{pmatrix} \cdot & ? & 1 & 0 & ? & 0 \\ ? & \cdot & ? & ? & 0 & ? \\ 1 & ? & \cdot & ? & 0 & 0 \\ 0 & ? & ? & \cdot & ? & 1 \\ ? & 0 & 0 & ? & \cdot & ? \\ 0 & ? & 0 & 1 & ? & \cdot \end{pmatrix} & \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} \end{matrix}$$

- ▶ 0: there is no interaction; 1: there is interaction; ?: unknown
- ▶ Given data, we use **graphical models** to incorporate existing information using a **constrained optimization** framework.
- ▶ Can **estimate novel interactions** and **validate existing information**.

$$A = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \end{matrix} \\ \begin{pmatrix} \cdot & ? & 1 & 0 & ? & 0 \\ ? & \cdot & ? & ? & 0 & ? \\ 1 & ? & \cdot & ? & 0 & 0 \\ 0 & ? & ? & \cdot & ? & 1 \\ ? & 0 & 0 & ? & \cdot & ? \\ 0 & ? & 0 & 1 & ? & \cdot \end{pmatrix} & \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} \end{matrix}$$

- ▶ **0**: there is no interaction; **1**: there is interaction; **?**: unknown
- ▶ Given data, we use **graphical models** to incorporate existing information using a **constrained optimization** framework.
- ▶ Can **estimate novel interactions** and **validate existing information**.
- ▶ Consistent estimation of network **requires fewer observations**, depending on the available external information.

⁹Hellstern et al. *PLoS Comp Bio.* '21

Partition large networks into smaller ones by estimating a block diagonal network.



This strategy improves computational speed with little loss in performance⁹.

⁹Hellstern et al. *PLoS Comp Bio.* '21

Incomplete Pathway Information



Pathway memberships may be unknown.

¹⁰Ma et al. *Bioinformatics*. '19

Pathway memberships may be unknown.

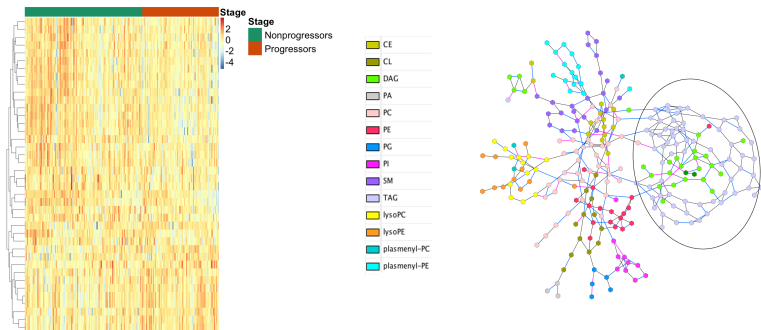


Fig: Inferred lipid interaction network in Chronic Kidney Disease progression

DNEA¹⁰ uses data to estimate the network topology, identify modules by consensus clustering of the network, and perform enrichment analysis.

¹⁰Ma et al. *Bioinformatics*. '19

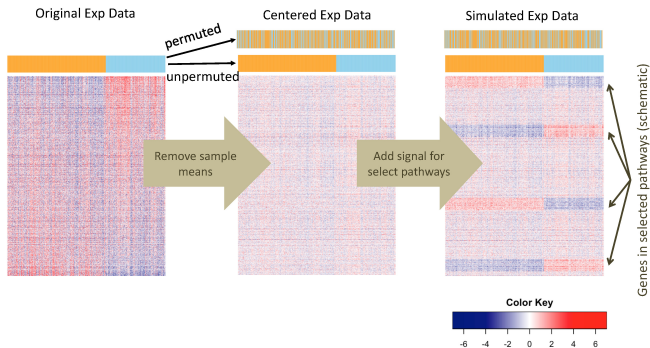
Competitive null:

- ▶ SPIA (Tarca et al. '09)
- ▶ camera (Wu and Smyth, '12)
- ▶ PathNet (Dutta, et al. '12)

Self-contained null:

- ▶ topologyGSA (Massa et al. '10)
- ▶ DEGraph (Jacob et al. '12)
- ▶ NetGSA (Ma et al. '16)

Synthetic data were generated from TCGA¹¹. $p = 2598$ genes; $n_1 = 403$ ER positive samples; $n_2 = 117$ ER negative samples.



Permuting the sample labels removes any difference in gene-gene correlation.

¹¹TCGA. *Nature*. '12

100 KEGG pathways (`graphite` R package).

Table 2 Average type I errors over multiple pathways, grouped by pathway sizes, for the TCGA breast cancer study [26].

Method	Pathway size	
	≤ 75	> 75
Pathway-Express	0*	0*
NetGSA	0.052	0.103
SPIA	0*	0*
topologyGSA	0.506	0.754
CAMERA	0.002	0.003
DEGraph	0.001	0.001
PathNet	0.048	0.057

* Under the self-contained null, the number of DE genes is zero. SPIA and Pathway-Express can not assess the impact of pathways that do not have any DE genes.

Clockwise from top left to bottom left: *Glucagon signaling pathway, AMPK signaling pathway, Insulin signaling pathway, and B cell receptor signaling pathway.*

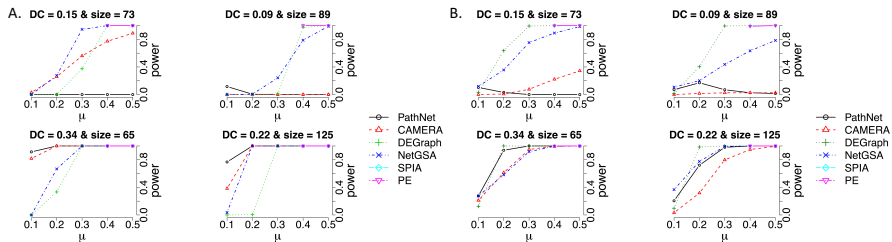


Fig: A: sample labels same as in TCGA; B: sample labels permuted.

Powers are averaged over multiple pathways that have similar proportion of affected genes.

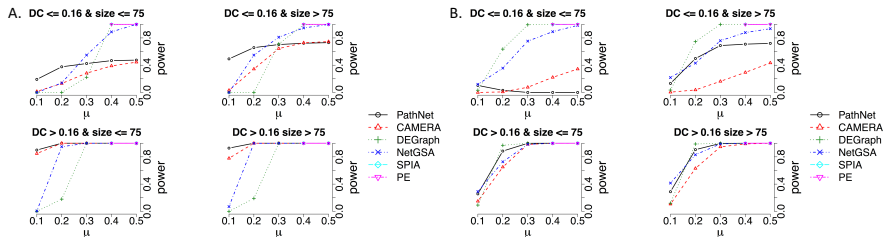
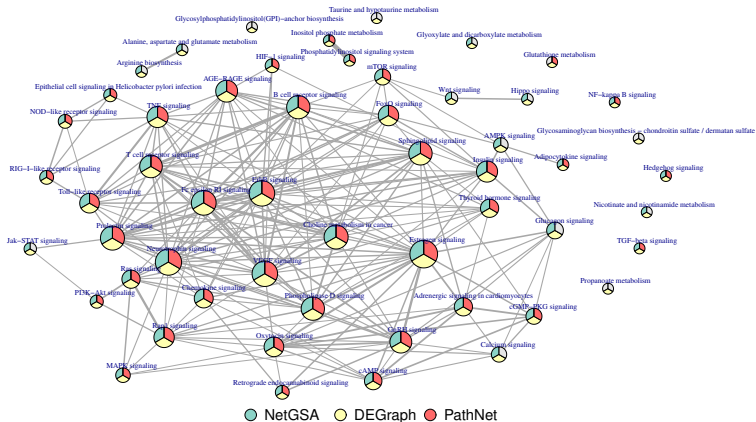
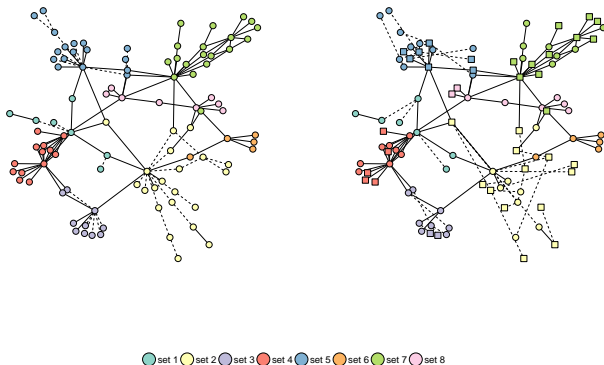


Fig: A: sample labels same as in TCGA; B: sample labels permuted.



- ▶ Nodes: pathways
- ▶ Edges: share of genes (top 5%)

Synthetic data were generated from a DREAM network with changes in network topology.



sets 1, 6: no change

sets 3, 8: 20% nodes with differential means

sets 4, 5: 40% nodes with differential means

sets 2, 7: 60% nodes with differential means

sets 1, 2, 3, 5: also have changes in topology

Table: Empirical powers averaged in 100 replications.

Method	1	2	3	4	5	6	7	8
NetGSA	0.08	0.89	0.96	0.14	0.99	0.02	0.94	0.03
DEGraph	0.18	1.00	1.00	0.49	1.00	0.06	0.62	0.31
true power	0.12	0.93	0.98	0.11	0.99	0.05	0.95	0.10

- ▶ REHE offers gain in computational efficiency with little loss in accuracy for fitting **large-scale** linear mixed models.

- ▶ REHE offers gain in computational efficiency with little loss in accuracy for fitting **large-scale** linear mixed models.
- ▶ NetGSA tests for gene set enrichment by incorporating the topology.

- ▶ REHE offers gain in computational efficiency with little loss in accuracy for fitting **large-scale** linear mixed models.
- ▶ NetGSA tests for gene set enrichment by incorporating the topology.
- ▶ NetGSA can leverage existing network information and expression data.

- ▶ REHE offers gain in computational efficiency with little loss in accuracy for fitting **large-scale** linear mixed models.
- ▶ NetGSA tests for gene set enrichment by incorporating the topology.
- ▶ NetGSA can leverage existing network information and expression data.
- ▶ **Caveat in gene set analysis: null hypothesis**

- 1 **Ma J**, Shojaie, A and Michailidis, G. Network-based pathway enrichment analysis with incomplete network information. *Bioinformatics*. 32(20):3165–3174, 2016.
- 2 **Ma J[†]**, Shojaie A and Michailidis G. A comparative study of topology-based pathway enrichment analysis methods. *BMC Bioinformatics*. 20 (546). 2019
- 3 **Ma J**, Karnovsky A, Afshinnia F, Wigginton J, Feldman H, Rader D, Shama K, Porter A, Rahman M, He J, Hamm L, Shafi T, Pennathur S, Michailidis G. Differential network-based enrichment analysis of lipid pathways altered in Chronic Kidney Disease progression. *Bioinformatics*. 35(18):3441–3452, 2019.
- 4 Hellstern M, **Ma J**, Yue K and Shojaie A. netgsa: Fast computation and interactive visualization for topology-based pathway enrichment analysis. *PLoS Computational Biology*. 17(6): e1008979, 2021.
- 5 Yue K, **Ma J**, Thornton T and Shojaie A. REHE: fast variance components estimation for linear mixed models. *Genetic Epidemiology*. 45(8):891–905, 2021.