

2017-04-27

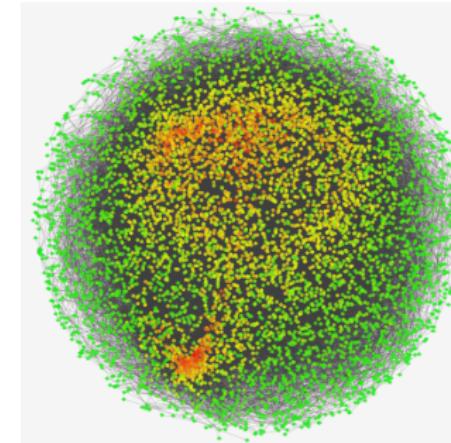
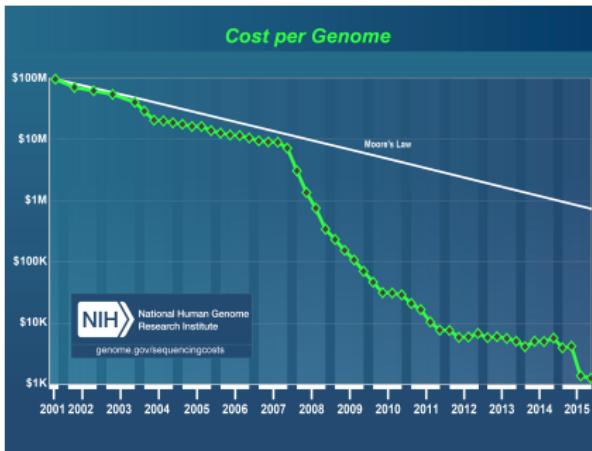
# Methods for Estimating Hidden Structure and Network Transitions in Genomics

Daniel Schlauch, PhD Candidate

Department of Biostatistics, Harvard University

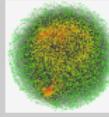
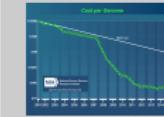
April 27, 2017

# Hidden signals in high dimensions

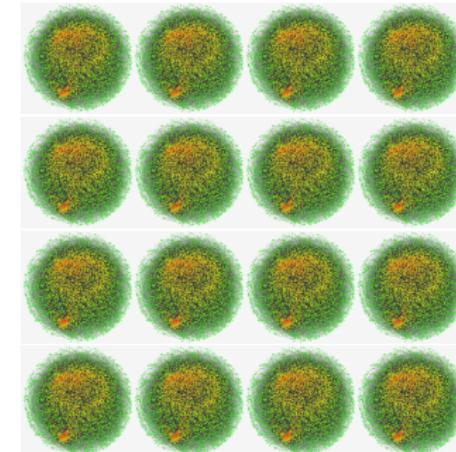
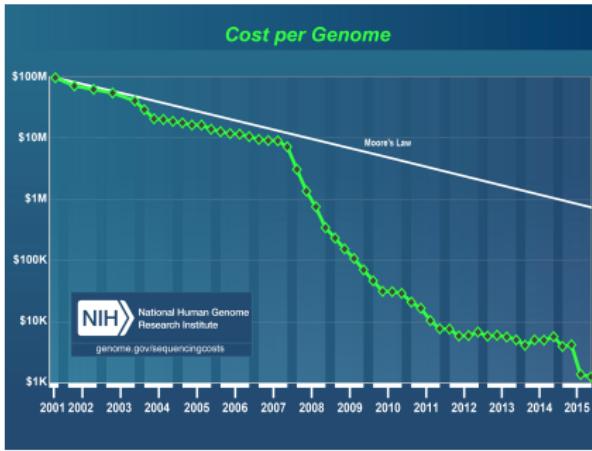


## └ Hidden signals in high dimensions

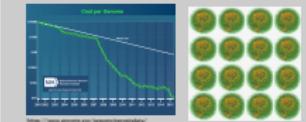
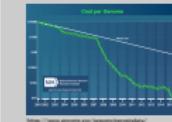
- Economics: New technologies have allowed us access to unprecedented data
- Microarrays, then seq, allowed us access to numerous fields of 'omics.
- Biology isn't simple
- But... snapshots of hairballs, which are noisy and convoluted



# Hidden signals in high dimensions



## └ Hidden signals in high dimensions



- Economics: New technologies have allowed us access to unprecedented data
- Microarrays, then seq, allowed us access to numerous fields of 'omics.'
- Biology isn't simple
- But... snapshots of hairballs, which are noisy and convoluted



# Hidden signals in high dimensions

- Lots of samples across LOTS of features
- Single snapshots of complex, fluid, heterogeneous landscapes
- Noisy, structured data



How do we separate the structure to reveal higher-level components?

## Methods for Estimating Hidden Structure and Network Transitions in Genomics

### Introduction

#### Hidden signals in high dimensions

2017-04-27

##### Hidden signals in high dimensions

- Lots of samples across LOTS of features
- Single snapshots of complex, fluid, heterogeneous landscapes
- Noisy, structured data



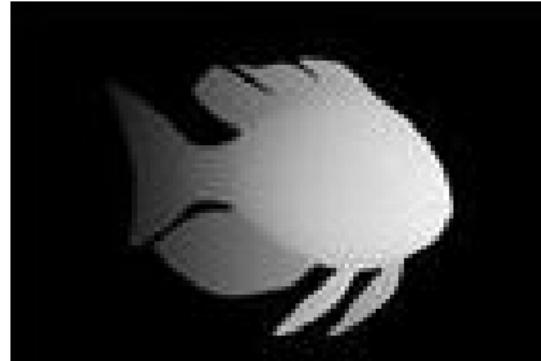
How do we separate the structure to reveal higher-level components?

- High Dimensional data
- This data has high amounts of technical and biological noise
- Single snapshots don't tell the whole picture.
- Collection of often heterogenous cells
- We often want to understand how features interact and lead to higher level phenotypes
- Challenging because structure comes in many forms.



# Hidden signals in high dimensions

- Lots of samples across LOTS of features
- Single snapshots of complex, fluid, heterogeneous landscapes
- Noisy, structured data



How do we separate the structure to reveal higher-level components?

## └ Hidden signals in high dimensions

- What's important is understanding that structure can be many things... pathways to cancer, or
- Microarrays, then seq, allowed us access to numerous fields of 'omics.
- But... snapshots of hairballs

- Lots of samples across LOTS of features
- Single snapshots of complex, fluid, heterogeneous landscapes
- Noisy, structured data



How do we separate the structure to reveal higher-level components?



# Table of Contents

- 1 Introduction
- 2 Identification of Genetic Outliers
- 3 Batch Effect on Covariance Structure
- 4 State Transitions Using Gene Regulatory Network Models
- 5 Future Work
- 6 Appendix

## └ Table of Contents

2017-04-27

Table of Contents

- 1 Introduction
- 2 Identification of Genetic Outliers
- 3 Batch Effect on Covariance Structure
- 4 State Transitions Using Gene Regulatory Network Models
- 5 Future Work
- 6 Appendix



# Identifying Genetic Outliers

## Identification of genetic outliers due to sub-structure and cryptic relationships

Daniel Schlauch<sup>1,2,4</sup>, Heide Fier<sup>1,3</sup> and Christoph Lange<sup>1,4</sup>

<sup>1</sup>Department of Biostatistics, Harvard TH Chan School of Public Health, Boston, MA 02115

<sup>2</sup>Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute, Boston, MA 02115

<sup>3</sup>Institute of Genomic Mathematics, University of Bonn, Bonn, Germany

<sup>4</sup>Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA 02115

## Methods for Estimating Hidden Structure and Network Transitions in Genomics

- └ Identification of Genetic Outliers

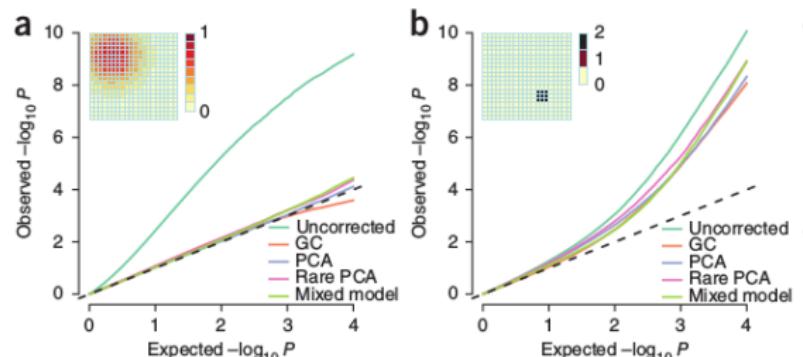
- └ Identifying Genetic Outliers

2017-04-27



## Background

- Individuals may be too similar (due to cryptic relatedness) or too different (due to population structure).
- Many methods exist for addressing some of these concerns (e.g. PCA, LMM).
- Localized effects and rare variants increase confounding



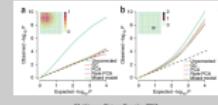
Mathieson, Nature Genetics 2012

2017-04-27

### Identification of Genetic Outliers

- Background
- Background

- Inflation of type I error rate



Mathieson, Nature Genetics 2012

# Motivation

We want to create a similarity measure that...

- is more sensitive to fine scale population stratification
- estimates relatedness
- can be used as a formal test for cryptic relatedness
- can be used as a formal test for population structure

2017-04-27

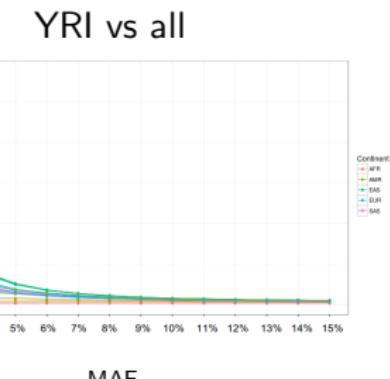
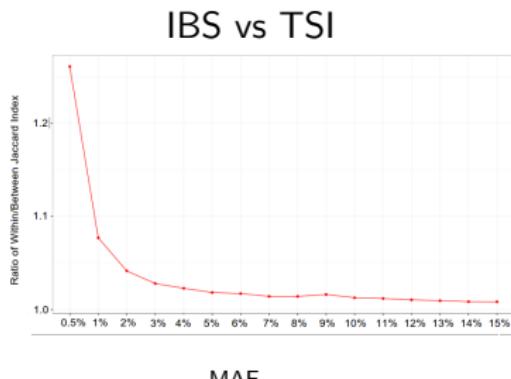
We want to create a similarity measure that...

- is more sensitive to fine scale population stratification
- estimates relatedness
- can be used as a formal test for cryptic relatedness
- can be used as a formal test for population structure



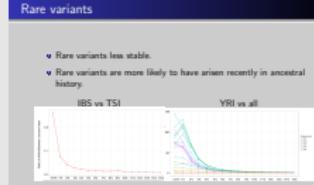
## Rare variants

- Rare variants less stable.
- Rare variants are more likely to have arisen recently in ancestral history.



2017-04-27

Background  
Rare variants



- RVs become fixed at 0% with high probability over a relatively short timeframe.

# Test Statistic

## STEGO: Similarity Test for Estimating Genetic Outliers

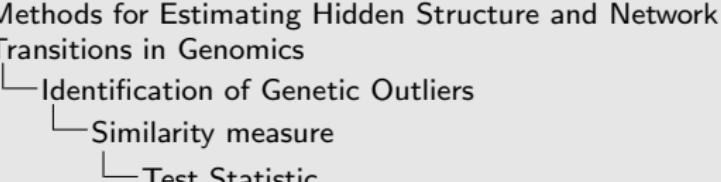
$$s_{i,j} = \frac{\sum_{k=1}^N w_k \mathbf{G}_{i,k} \mathbf{G}_{j,k}}{\sum_{k=1}^N I \left[ \sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1 \right]}$$

where

$$w_k = \begin{cases} \frac{\binom{2n}{2}}{\left( \sum_{l=1}^{2n} \mathbf{G}_{l,k} \right)} & \sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1 \\ 0 & \sum_{l=1}^{2n} \mathbf{G}_{l,k} \leq 1 \end{cases}$$

Under  $H_0$ , homogeneity,

$$s_{i,j} \sim N(1, \sigma_{i,j}^2)$$



2017-04-27

Test Statistic

STEGO: Similarity Test for Estimating Genetic Outliers

$$s_{i,j} = \frac{\sum_{k=1}^N w_k \mathbf{G}_{i,k} \mathbf{G}_{j,k}}{\sum_{k=1}^N I \left[ \sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1 \right]}$$

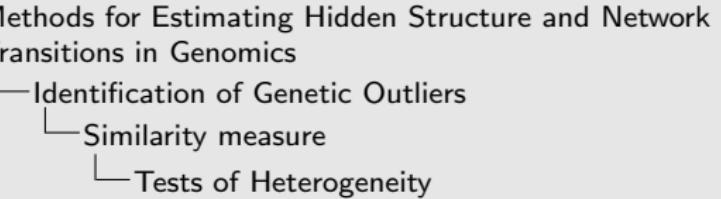
where

$$w_k = \begin{cases} \frac{\binom{2n}{2}}{\left( \sum_{l=1}^{2n} \mathbf{G}_{l,k} \right)} & \sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1 \\ 0 & \sum_{l=1}^{2n} \mathbf{G}_{l,k} \leq 1 \end{cases}$$

Under  $H_0$ , homogeneity,

$$s_{i,j} \sim N(1, \sigma_{i,j}^2)$$

# Tests of Heterogeneity



2017-04-27

$$H_0 : E[s_{i,j}] = 1 \forall i, j \in 1 \dots n$$

$$H_A : \exists i, j \in 1 \dots n \mid E[s_{i,j}] \neq 1$$

Test for cryptic relatedness:

$$R : \max(s_{i,j}) > 1 - \text{probit} \left( \frac{\alpha}{\binom{n}{2}} \right)$$

Test for population structure:

$$K = \sup_x |F_s(x) - \Phi(x)|$$



$$H_0 : E[s_{i,j}] = 1 \forall i, j \in 1 \dots n$$
$$H_A : \exists i, j \in 1 \dots n \mid E[s_{i,j}] \neq 1$$

Test for cryptic relatedness:

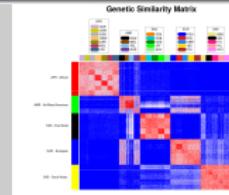
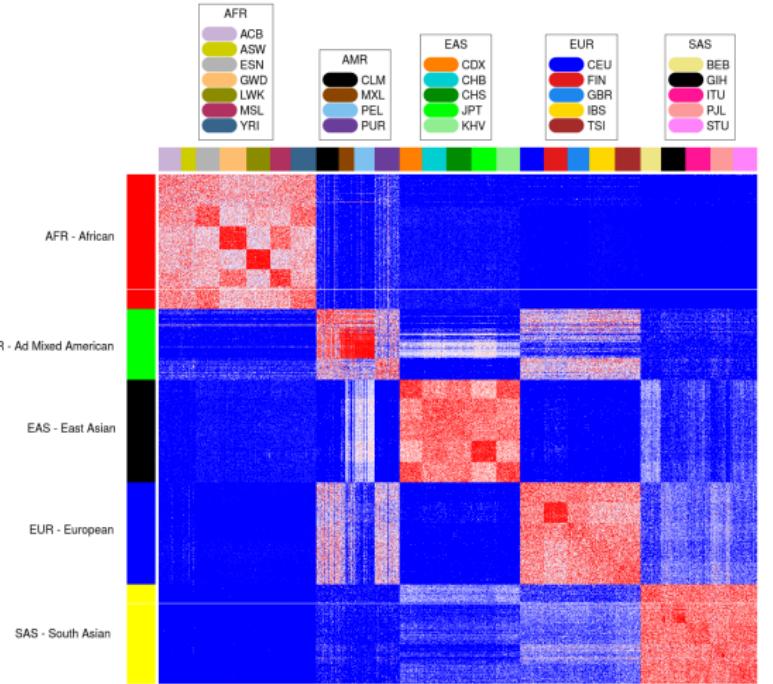
$$R : \max(s_{i,j}) > 1 - \text{probit} \left( \frac{\alpha}{\binom{n}{2}} \right)$$

Test for population structure:

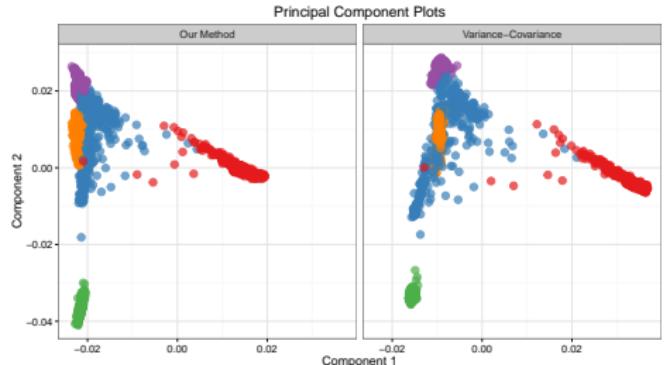
$$K = \sup_x |F_s(x) - \Phi(x)|$$

# Application to 1000 Genomes Project

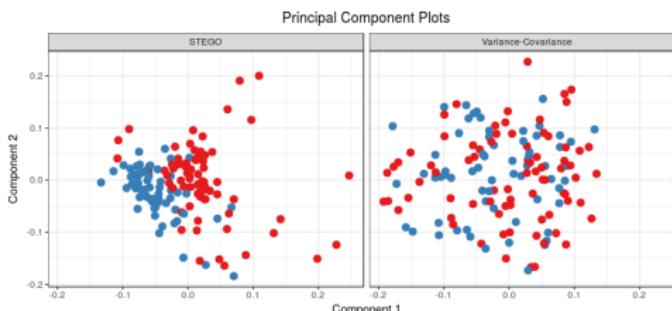
## Genetic Similarity Matrix



# Application to 1000 Genomes Project

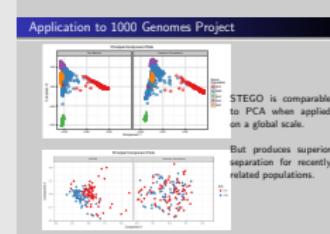


STEGO is comparable  
to PCA when applied  
on a global scale.



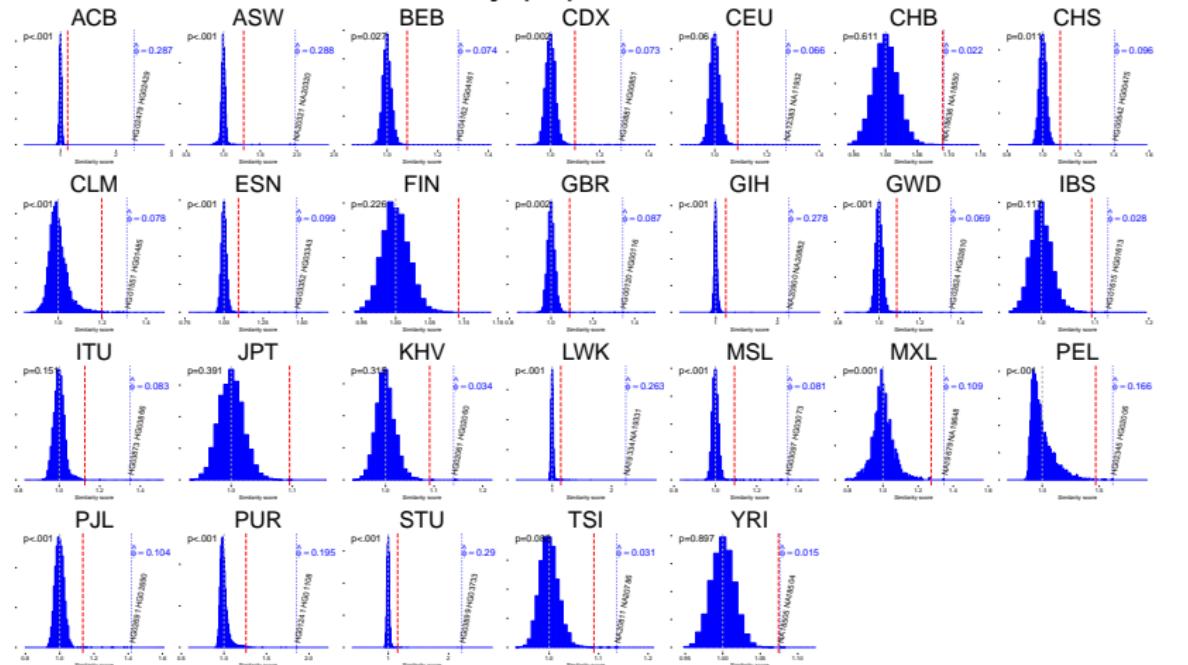
But produces superior  
separation for recently  
related populations.

2017-04-27



# Application to 1000 Genomes Project

## Distribution of $s$ statistics by population

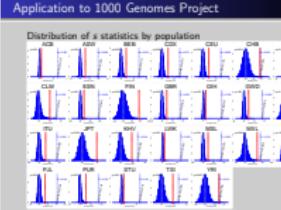


## Methods for Estimating Hidden Structure and Network Transitions in Genomics

### Identification of Genetic Outliers

#### Results

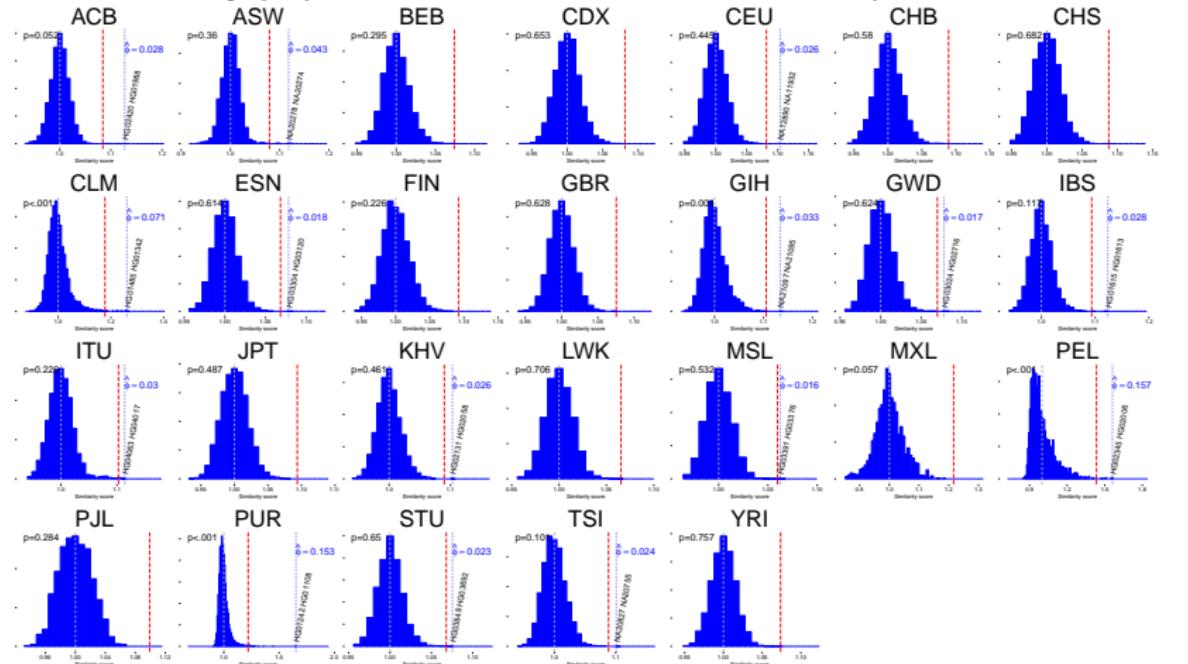
##### Application to 1000 Genomes Project



- Describe red-line, black line, and labeled pair.

# Application to 1000 Genomes Project

s statistics by population after removal of related pairs



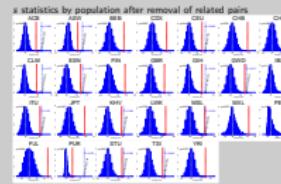
## Methods for Estimating Hidden Structure and Network Transitions in Genomics

### Identification of Genetic Outliers

#### Results

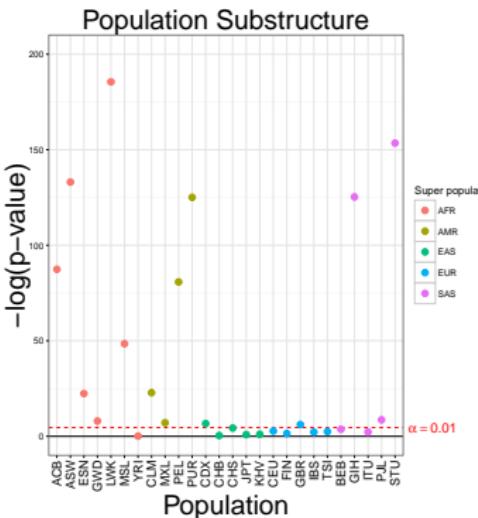
##### Application to 1000 Genomes Project

## Application to 1000 Genomes Project

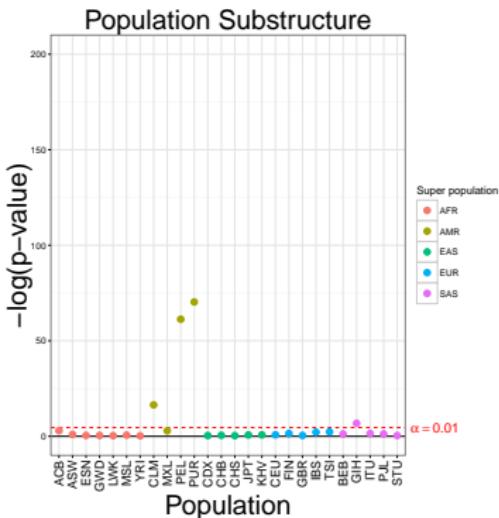


# Application to 1000 Genomes Project

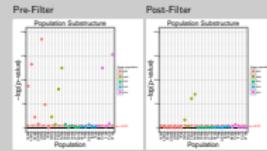
## Pre-Filter



## Post-Filter



2017-04-27



# Batch Effect on Covariance Structure



## Batch effect on covariance structure confounds gene coexpression

Daniel Schlauch<sup>1,2</sup>, Joseph N. Paulson<sup>2</sup>, Kimberly Glass<sup>2,3</sup>, and  
John Quackenbush<sup>1,3</sup>

<sup>1</sup>Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute and Department of Biostatistics, Harvard TH Chan School of Public Health, Boston, MA

<sup>2</sup>Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA

<sup>3</sup>Department of Medicine, Harvard Medical School, Boston, MA

## Methods for Estimating Hidden Structure and Network Transitions in Genomics

- └ Batch Effect on Covariance Structure

- └ Batch Effect on Covariance Structure

Batch Effect on Covariance Structure

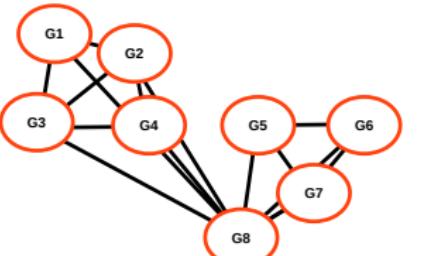
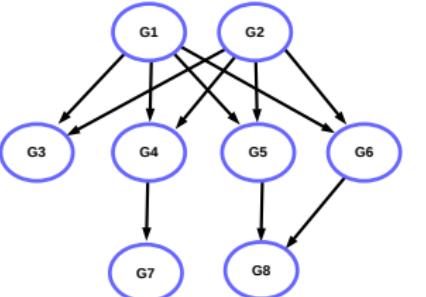
Batch effect on covariance structure  
confounds gene coexpression

Daniel Schlauch<sup>1,2</sup>, Joseph N. Paulson<sup>2</sup>, Kimberly Glass<sup>2,3</sup>, and  
John Quackenbush<sup>1,3</sup>  
<sup>1</sup>Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute and Department of Biostatistics, Harvard TH Chan School of Public Health, Boston, MA  
<sup>2</sup>Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA  
<sup>3</sup>Department of Medicine, Harvard Medical School, Boston, MA

## Background: Differential Network Inference

How do we model functional interactions?

- Gene Regulatory/Coexpression Networks (GRN/GCN)
- Directed/undirected graph
- May imply a sort of physical interaction
- Guilt by association

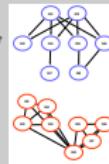


Methods for Estimating Hidden Structure and Network Transitions in Genomics

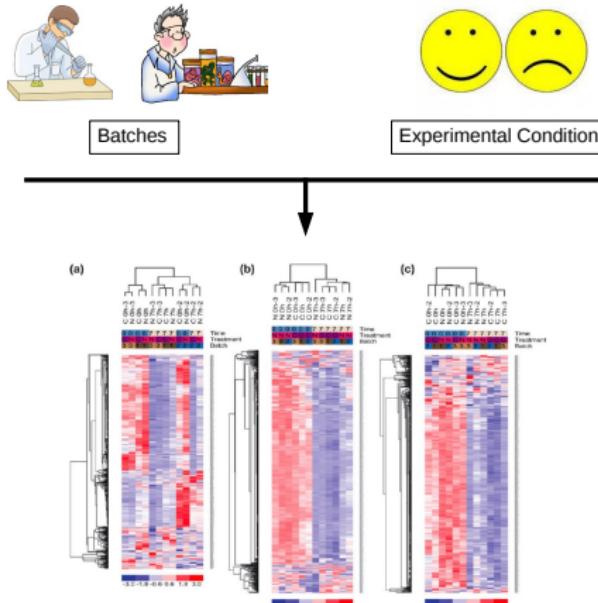
- Batch Effect on Covariance Structure
- Gene Networks and Batch Effect
- Background: Differential Network Inference

2017-04-27

Background: Differential Network Inference



## Background: Batch Effect

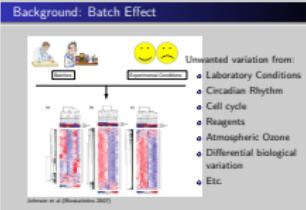


Johnson et al.(Biostatistics 2007)

Methods for Estimating Hidden Structure and Network Transitions in Genomics

- Batch Effect on Covariance Structure
- Gene Networks and Batch Effect
- Background: Batch Effect

2017-04-27



# Methods for Controlling Batch Effect

2017-04-27

Location scale model:

$$Y_{ijg} = \alpha_g + X\beta_g + \gamma_{ig} + \delta_{ig}\epsilon_{ijg}$$

Limitations:

- Gene-specific location/scale assumptions
- Independent effects
- Differential coexpression*

Batch effect removal methods typically return a corrected gene expression matrix (e.g. ComBat) or a correction vector (e.g. SVA).

- point 1

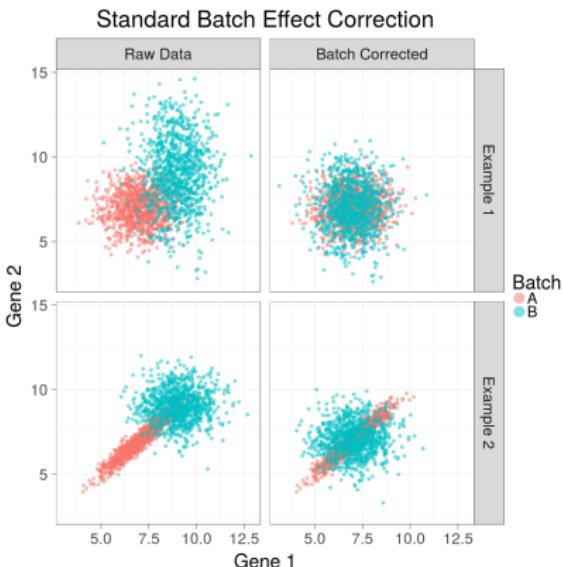


# Limitations to common batch effect correction methods

## Batch Effect on Covariance Structure

### Gene Networks and Batch Effect

#### Limitations to common batch effect correction methods



### Standard corrections:

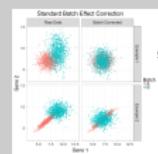
$$f[Gene1|BatchA] = f[Gene1|BatchB]$$

$$f[Gene2|BatchA] = f[Gene2|BatchB]$$

$$f[Gene1, Gene2|BatchA] \neq f[Gene1, Gene2|BatchB]$$

2017-04-27

Limitations to common batch effect correction methods



Standard corrections:

1)  $Gene1|BatchA = f(Gene1|BatchB)$

2)  $Gene2|BatchA = f(Gene2|BatchB)$

3)  $Gene1, Gene2|BatchA = f(Gene1, Gene2|BatchB)$



# Limitations to existing batch effect correction methods

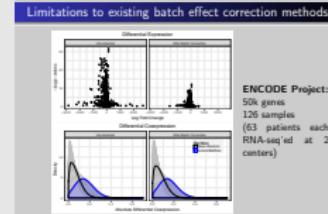
2017-04-27

## Methods for Estimating Hidden Structure and Network Transitions in Genomics

### Batch Effect on Covariance Structure

#### Gene Networks and Batch Effect

#### Limitations to existing batch effect correction methods



## ENCODE Project:

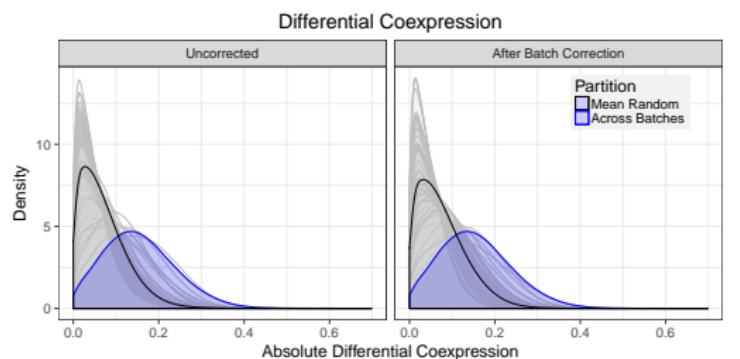
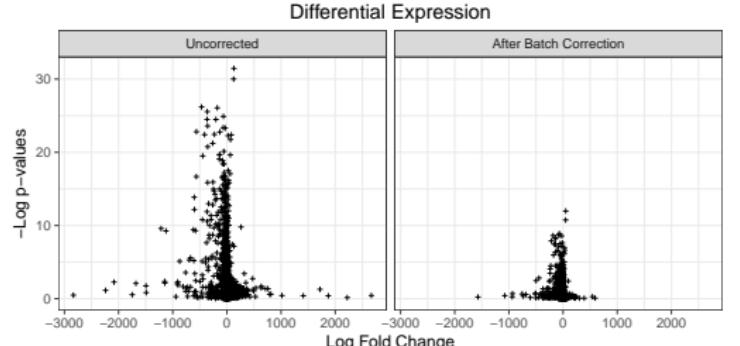
50k genes

126 samples

(63 patients each

RNA-seq'ed at 2

centers)



# Estimating the conditional coexpression matrix

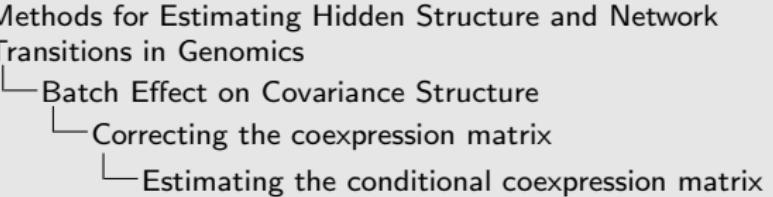
Proposed method: **CMA**, Coexpression model adjustment

Motivating concepts:

- Provide a regression framework for the coexpression matrix.
- Estimate a reduced number of parameters.

Our approach:

- Exploit modular nature of gene expression patterns.
- Define our parameters as functions of components of variation.
- Estimate the eigenvalue contribution of each eigenvector.



2017-04-27

Proposed method: **CMA**, Coexpression model adjustment  
Motivating concepts:

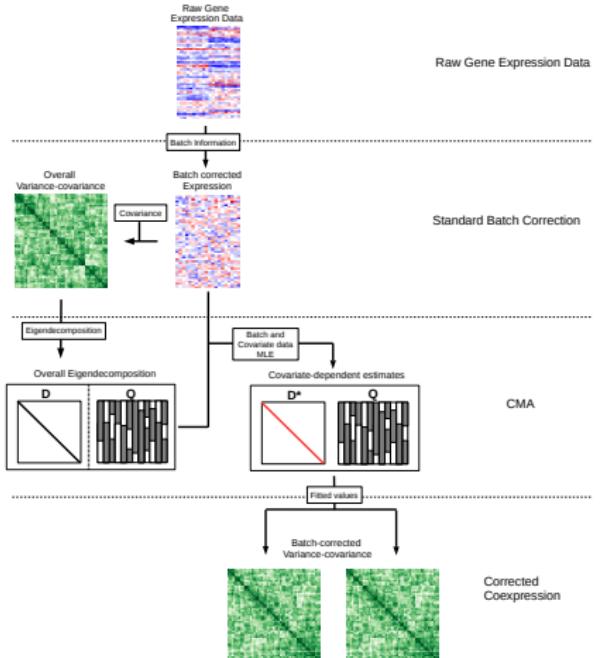
- Provide a regression framework for the coexpression matrix.
- Estimate a reduced number of parameters.

Our approach:

- Exploit modular nature of gene expression patterns.
- Define our parameters as functions of components of variation.
- Estimate the eigenvalue contribution of each eigenvector.



# Example Workflow

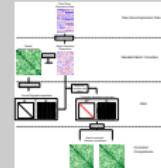


Methods for Estimating Hidden Structure and Network Transitions in Genomics

- └ Batch Effect on Covariance Structure
  - └ Correcting the coexpression matrix
    - └ Example Workflow

2017-04-27

Example Workflow



# Estimator

$$\hat{\Psi}_{\cdot,h} = (\mathbf{X}^T \mathbf{X})^{-1} \sum_{i=1}^N [\mathbf{X}_i^T \mathbf{Q}_h^T \mathbf{G}_i^* \mathbf{G}_i^{*\top} \mathbf{Q}_h]$$

where

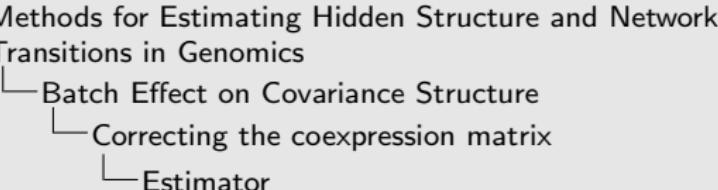
$\mathbf{G}_i^*$  is the residual gene expression for sample  $i$

$\mathbf{X}$  is the design matrix

$\mathbf{Q}_h$  is the  $h^{th}$  eigenvector

$\hat{\Psi}$  is a  $q \times p$  matrix ( $q = \#covariates + 1$ ,  $p = \#genes$ )

(See Appendix)



2017-04-27

Estimator

$$\hat{\Psi}_{\cdot,h} = (\mathbf{X}^T \mathbf{X})^{-1} \sum_{i=1}^N [\mathbf{X}_i^T \mathbf{Q}_h^T \mathbf{G}_i^* \mathbf{G}_i^{*\top} \mathbf{Q}_h]$$

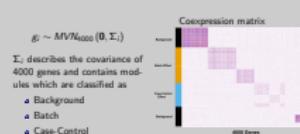
where  
 $\mathbf{G}_i^*$  is the residual gene expression for sample  $i$   
 $\mathbf{X}$  is the design matrix  
 $\mathbf{Q}_h$  is the  $h^{th}$  eigenvector  
 $\hat{\Psi}$  is a  $q \times p$  matrix ( $q = \#covariates + 1$ ,  $p = \#genes$ )  
(See Appendix)

- We can use  $\hat{\Psi}$  in the standard regression way... differential coexpression, fitted coexpression...



## Simulations

2017-04-27

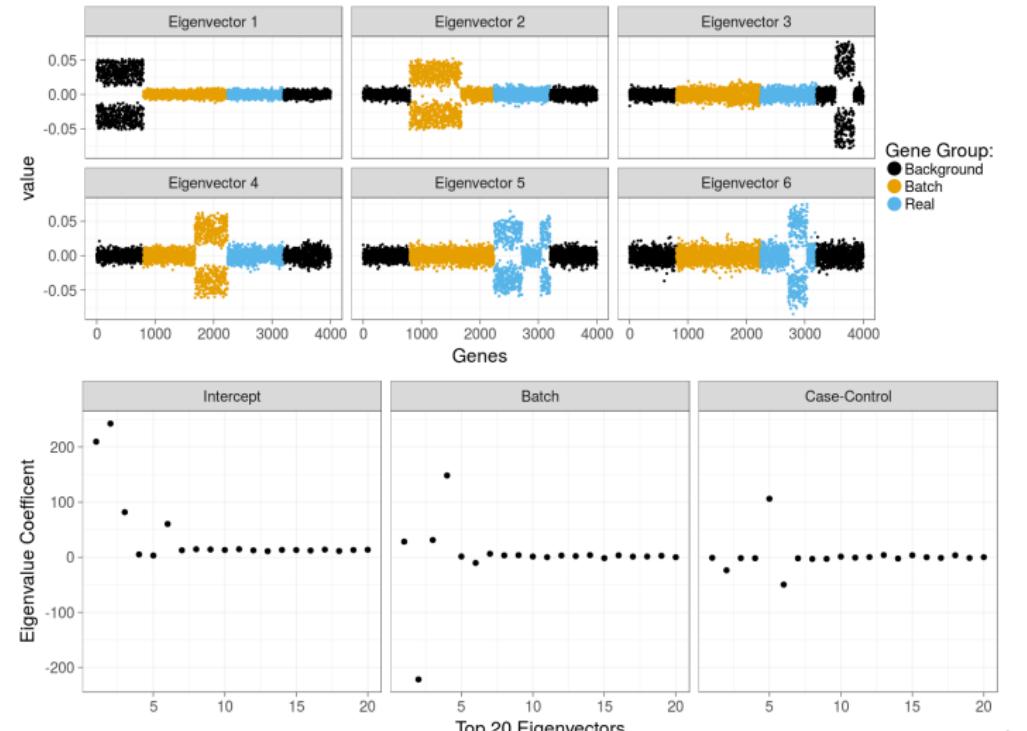


- Describe naive approach

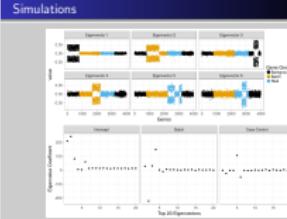
$\Sigma_i$  describes the covariance of 4000 genes and contains modules which are classified as

- Background
- Batch
- Case-Control

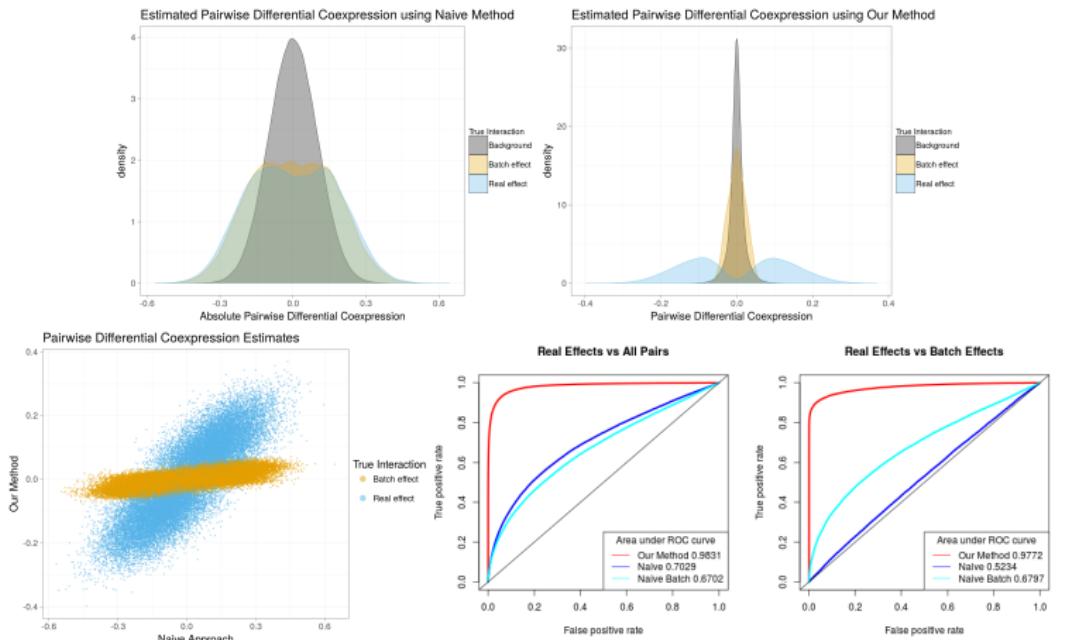
# Simulations



2017-04-27



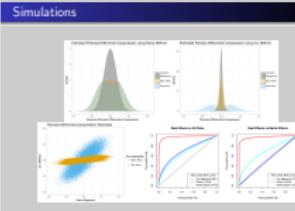
# Simulations



## Methods for Estimating Hidden Structure and Network Transitions in Genomics

- Batch Effect on Covariance Structure
  - Results
  - Simulations

2017-04-27



# Application to data from COPDGene Study

The COPDGene Study (GSE42057):

- 136 Individuals
- Affymetrix Human Genome U133 Plus 2.0 microarrays
- 18,960 genes
- 42 Smoker Controls, 94 COPD subjects

$$\mathbf{S} \sim COPD + Gender + Age + Packyears$$

2017-04-27

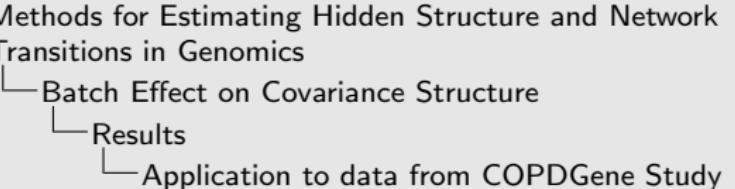
# Application to data from COPDGene Study

## Functional Enrichment in COPD-associated module (Morrow et al. 2015)

GO Term	Genes	p-value	FDR
immune response-regulating cell surface receptor	12	3.33E-010	7.89E-008
B cell proliferation	8	1.33E-010	6.30E-008
antigen receptor-mediated signaling pathway	8	9.25E-009	1.10E-006
regulation of B cell activation	7	1.44E-008	1.37E-006
B cell receptor signaling pathway	6	1.64E-009	2.59E-007
regulation of lymphocyte proliferation	6	1.20E-005	7.34E-004
regulation of mononuclear cell proliferation	6	1.24E-005	7.34E-004
T cell aggregation	6	1.19E-003	2.89E-002

## Functional Enrichment in top differentially coexpressed genes from CMA

GO Term	%	Enrichment	FDR
anatomical structure development	0.26	1.29	2.58E-05
single-organism developmental process	0.26	1.29	2.73E-05
anatomical structure morphogenesis	0.14	1.46	1.60E-04
single-multicellular organism process	0.27	1.25	4.01E-04
system process	0.11	1.50	1.46E-03
regulation of cellular process	0.43	1.12	5.86E-03
single organism signaling	0.28	1.21	7.89E-03
regulation of localization	0.13	1.40	1.15E-02



2017-04-27

Application to data from COPDGene Study

GO Term	Genes	p-values	FDR
immune response-regulating cell surface receptor	12	3.33E-002	7.89E-008
B cell proliferation	8	1.33E-002	6.30E-008
antigen receptor-mediated signaling pathway	8	9.25E-002	1.10E-006
regulation of B cell activation	7	1.44E-002	1.37E-006
B cell receptor signaling pathway	6	1.64E-002	2.59E-007
regulation of lymphocyte proliferation	6	1.20E-002	7.34E-004
regulation of mononuclear cell proliferation	6	1.24E-002	7.34E-004
T cell aggregation	6	1.19E-002	2.89E-002

Functional Enrichment in top differentially coexpressed genes from CMA

GO Term	Enrichment	FDR
anatomical structure development	1.29	2.58E-05
single-organism developmental process	1.29	2.73E-05
anatomical structure morphogenesis	1.46	1.60E-04
single-multicellular organism process	1.25	4.01E-04
system process	1.50	1.46E-03
regulation of cellular process	1.12	5.86E-03
single organism signaling	1.21	7.89E-03
regulation of localization	1.40	1.15E-02



# State Transitions Using Gene Regulatory Network Models

Methods for Estimating Hidden Structure and Network  
Transitions in Genomics  
└ State Transitions Using Gene Regulatory Network Models  
└ State Transitions Using Gene Regulatory Network  
Models

2017-04-27

State Transitions Using Gene Regulatory Network Models

Estimating Drivers of Cell State Transitions  
Using Gene Regulatory Network Models

Daniel Schlauch<sup>1,2</sup>, Kimberly Glass<sup>2,3</sup>, Craig P. Hersh<sup>2,3,4</sup>, Edwin  
K. Silverman<sup>2,3,4</sup> and John Quackenbush<sup>1,2,3</sup>

## Estimating Drivers of Cell State Transitions Using Gene Regulatory Network Models

Daniel Schlauch<sup>1,2</sup>, Kimberly Glass<sup>2,3</sup>, Craig P. Hersh<sup>2,3,4</sup>, Edwin  
K. Silverman<sup>2,3,4</sup> and John Quackenbush<sup>1,2,3</sup>

<sup>1</sup>Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute and Department of  
Biostatistics, Harvard TH Chan School of Public Health, Boston, MA

<sup>2</sup>Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA

<sup>3</sup>Department of Medicine, Harvard Medical School, Boston, MA

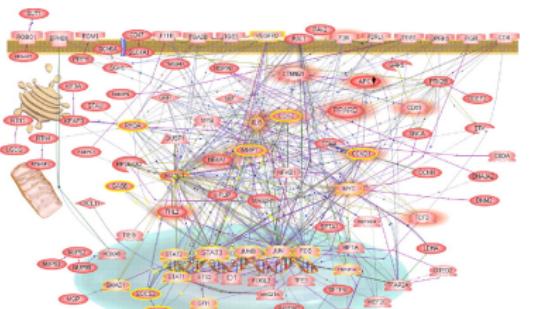
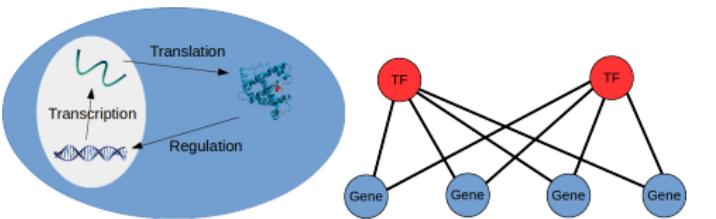
<sup>4</sup>Pulmonary and Critical Care Division, Brigham and Women's Hospital and Harvard Medical School, Boston, MA



## Background

### Why Study Gene Regulatory Networks?

- Genes are not independent objects.
- Regulation of higher level pathways and processes.

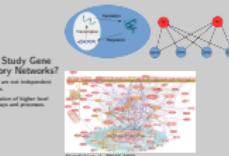


Abdollahi et al. PNAS 2007

### Methods for Estimating Hidden Structure and Network Transitions in Genomics

#### State Transitions Using Gene Regulatory Network Models

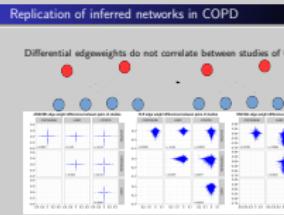
- Background
- Background



- Measurements of gene expression are at the mRNA level.
- Measurements only consist of mRNA abundance.
- Experimental data is collected as static snapshots.
- Biological variability can be difficult to induce
- 
- Gene expression measurements are noisy.
- Model complexity may require the estimate of too many model parameters.
- May be computationally intractable.
- May be statistically undetermined. "The curse of dimensionality"

# Replication of inferred networks in COPD

2017-04-27



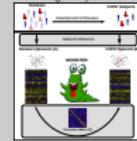
# Algorithm Overview

## Methods for Estimating Hidden Structure and Network Transitions in Genomics

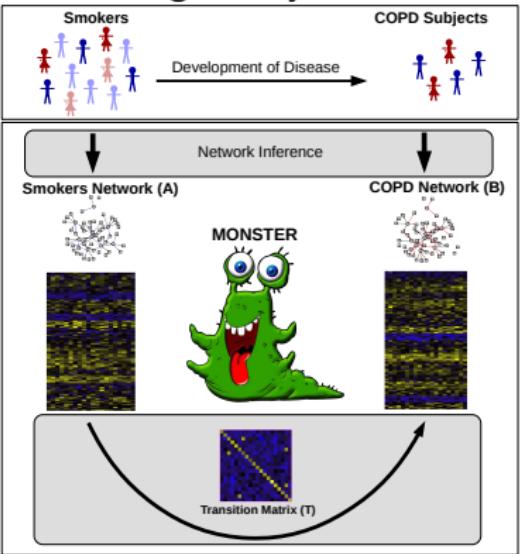
- └ State Transitions Using Gene Regulatory Network Models
  - └ Network Inference
    - └ Algorithm Overview

### Algorithm Overview

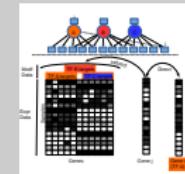
Modeling Network State Transitions from Expression and  
Regulatory data



## Modeling Network State Transitions from Expression and Regulatory data



# Network Inference



# Network Inference

Direct Evidence:

$$d_{i,j} = \text{cor}(g_i, g_j | \{g_{k,-i} : k \neq i, k \in \text{TF}\})^2$$

Indirect Evidence:

$$\text{logit}(E[M_i]) = \beta_0 + \beta_1 g_{(1)} + \cdots + \beta_N g_{(n)}$$

$$e_{i,j} = \frac{1}{1 + e^{\beta_0 + \beta_1 g_{j,(1)} + \cdots + \beta_k g_{j,(k)}}}$$

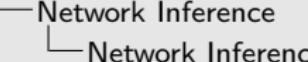
Edgeweight:

$$w_{i,j} = (1 - \alpha) [d_{i,j}] + \alpha [e_{i,j}]$$

## Methods for Estimating Hidden Structure and Network

### Transitions in Genomics

#### State Transitions Using Gene Regulatory Network Models



2017-04-27

## Network Inference

### Direct Evidence:

$$d_{i,j} = \text{cor}(g_i, g_j | \{g_{k,-i} : k \neq i, k \in \text{TF}\})^2$$

### Indirect Evidence:

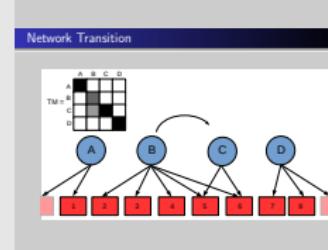
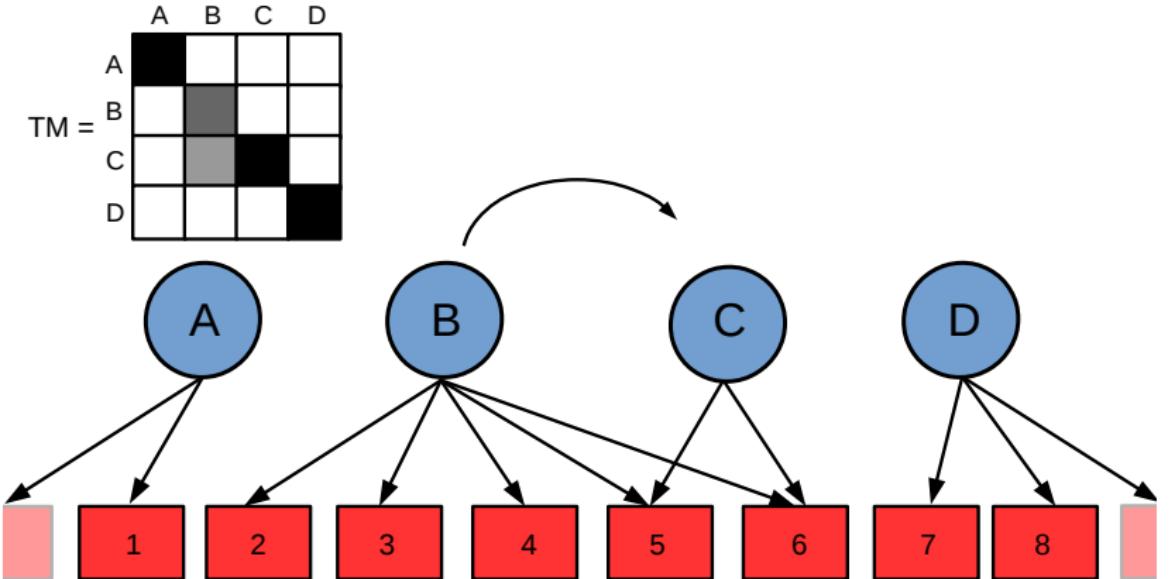
$$\text{logit}(E[M_i]) = \beta_0 + \beta_1 g_{(1)} + \cdots + \beta_N g_{(n)}$$

### Edgeweight:

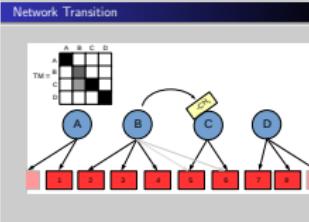
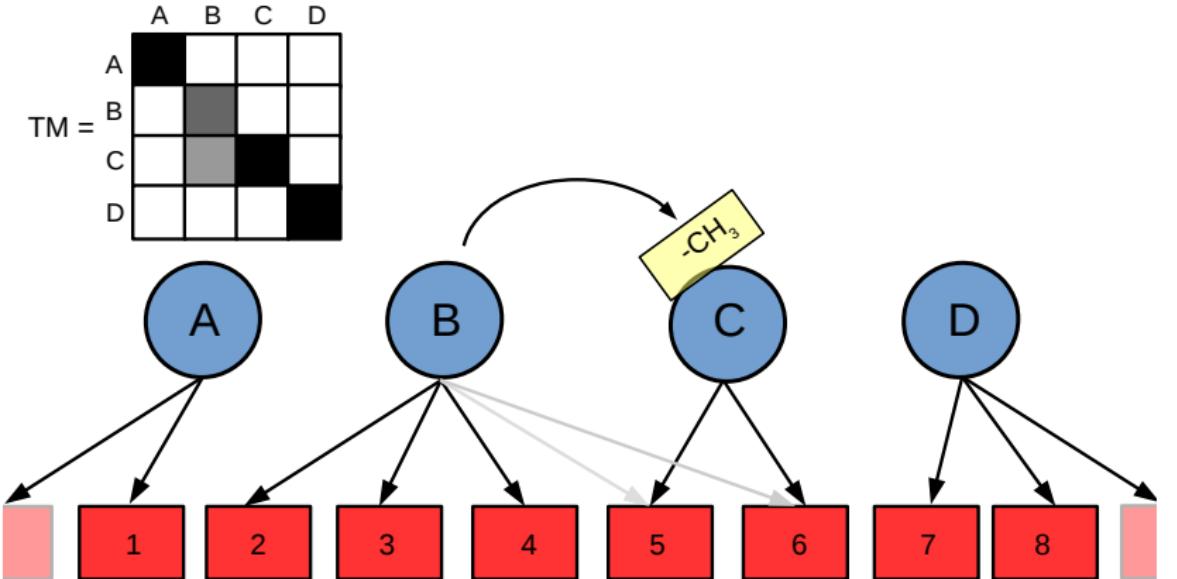
$$w_{i,j} = \frac{1}{1 + e^{\beta_0 + \beta_1 g_{j,(1)} + \cdots + \beta_k g_{j,(k)}}}$$

$$w_{i,j} = (1 - \alpha) [d_{i,j}] + \alpha [e_{i,j}]$$

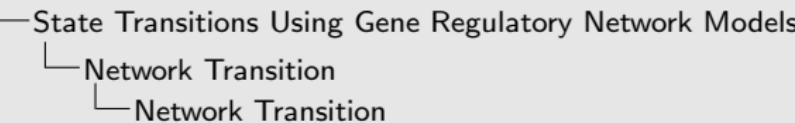
# Network Transition



# Network Transition



# Network Transition



2017-04-27

$$E[b_i - a_i] = \tau_{1,i}a_1 + \cdots + \tau_{m,i}a_m$$

where  $b_i$  and  $a_i$  are column-vectors in  $\mathbf{B}$  and  $\mathbf{A}$ .

In the simplest case, this can be solved with normal equations,

$$\hat{\tau}_i = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (b_i - a_i)$$

to generate each of the columns of the transition matrix  $\mathbf{T}$  such that

$$\hat{\mathbf{T}} = [\hat{\tau}_1, \hat{\tau}_2, \dots, \hat{\tau}_m]$$



$$E[b_i - a_i] = \tau_{1,i}a_1 + \cdots + \tau_{m,i}a_m$$

where  $b_i$  and  $a_i$  are column-vectors in  $\mathbf{B}$  and  $\mathbf{A}$ .

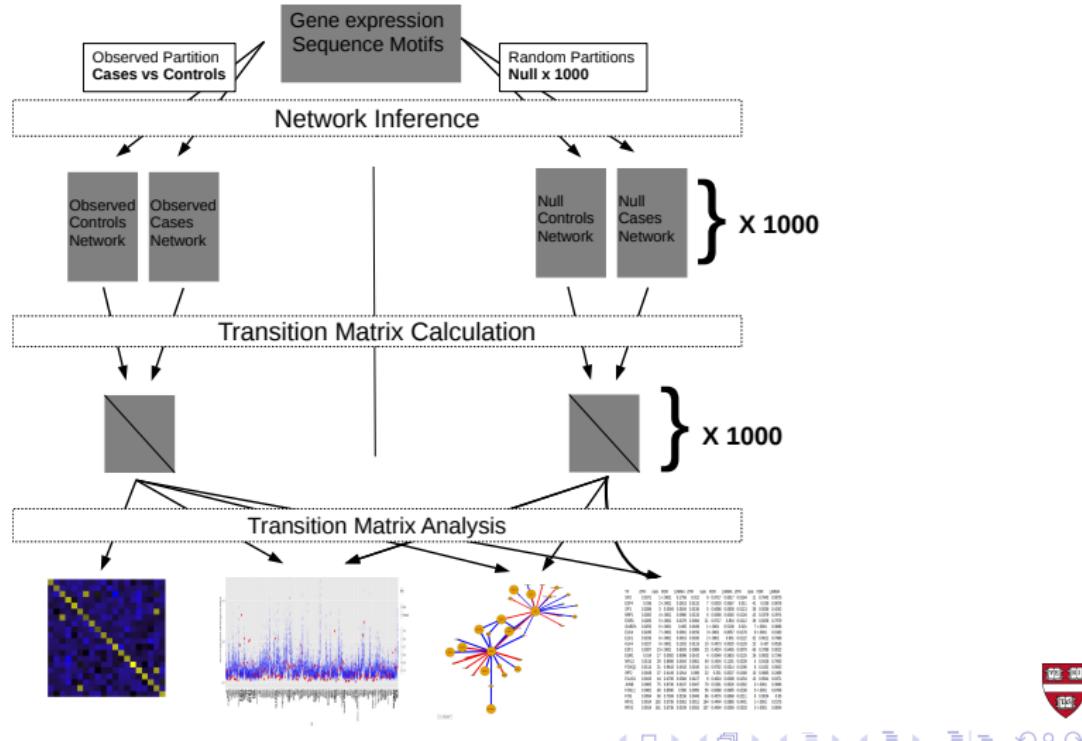
In the simplest case, this can be solved with normal equations,

$$\hat{\tau}_i = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (b_i - a_i)$$

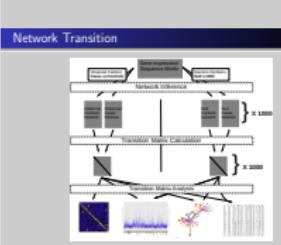
to generate each of the columns of the transition matrix  $\mathbf{T}$  such that

$$\mathbf{T} = [\hat{\tau}_1, \hat{\tau}_2, \dots, \hat{\tau}_m]$$

# Network Transition

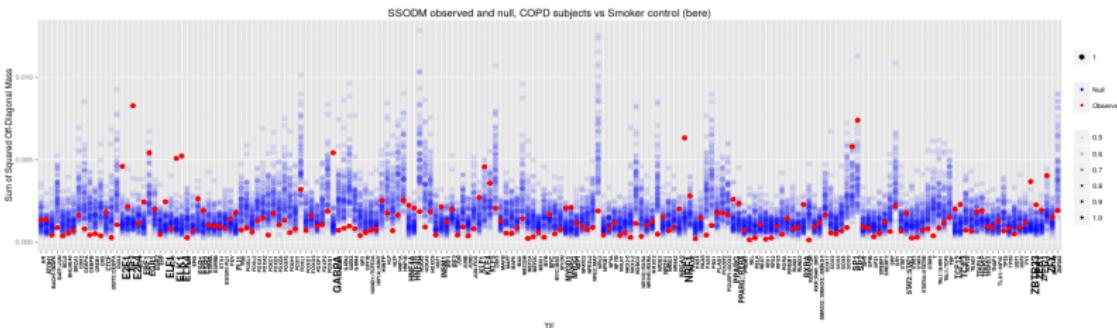
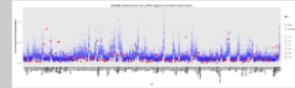


2017-04-27



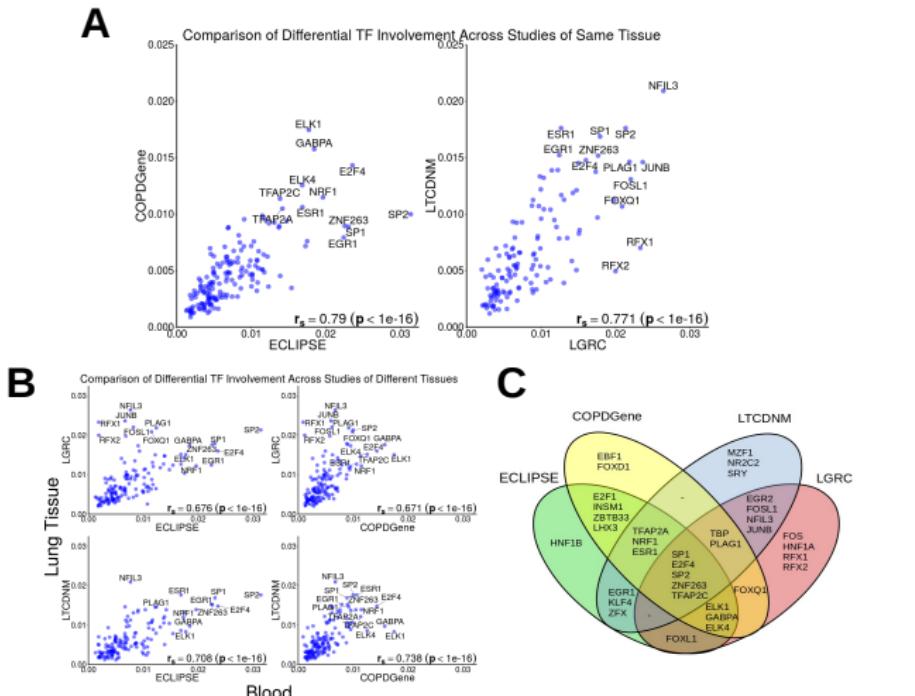
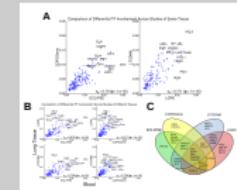
# Evaluating Transition Matrix

$$d\hat{TFI}_j = \frac{\sum_{i=1}^m I(i \neq j) \hat{\tau}_{i,j}^2}{\sum_{i=1}^m \hat{\tau}_{i,j}^2}$$

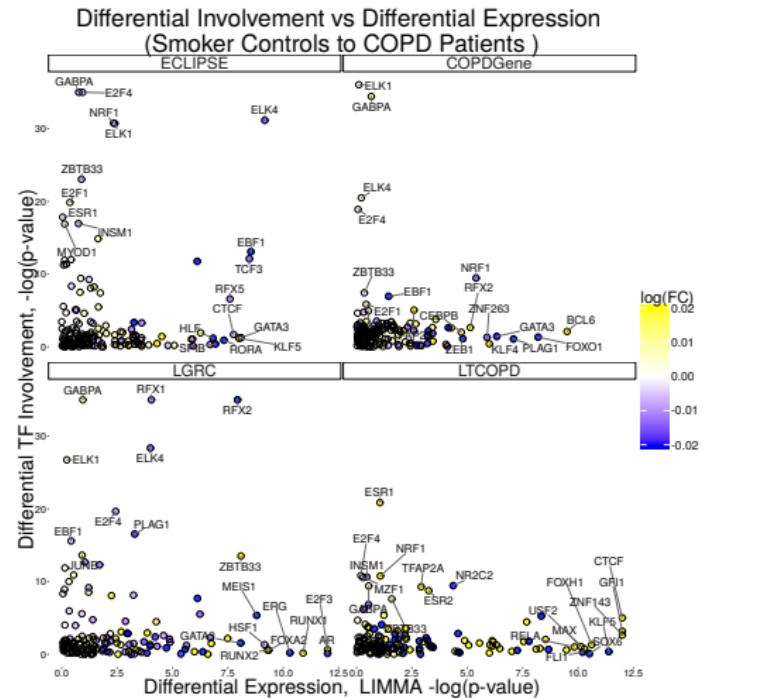


# Reproducibility and novel results

- └ State Transitions Using Gene Regulatory Network Models
  - └ Results
    - └ Reproducibility and novel results



# Reproducibility and novel results

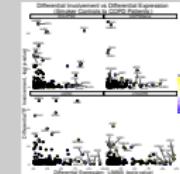


Methods for Estimating Hidden Structure and Network  
Transitions in Genomics

- State Transitions Using Gene Regulatory Network Models
- Results
- Reproducibility and novel results

2017-04-27

Reproducibility and novel results



## Future Work

### Identification of Genetic Outliers:

- Using STEGO in association studies using Linear Mixed Models. Does it reduce false positives and retain power?
- Can STEGO be used to reduce excess Type I error in rare variants (Mathieson & McVean)?

### Coexpression Batch Effect:

- Regularization approaches for estimated eigenvalues.
- Alternative sources for eigenvectors (functional categories, external validation sets).
- Significance analysis for coefficients.

2017-04-27

### Future Work

#### Future Work

##### Identification of Genetic Outliers:

- Using STEGO in association studies using Linear Mixed Models. Does it reduce false positives and retain power?
- Can STEGO be used to reduce excess Type I error in rare variants (Mathieson & McVean)?

##### Coexpression Batch Effect:

- Regularization approaches for estimated eigenvalues.
- Alternative sources for eigenvectors (functional categories, external validation sets).
- Significance analysis for coefficients.



## Future Work

### Estimating Network State Transitions:

- Integrating complimentary data sources into network inference and transition step.
- Use off-diagonal elements to predict protein-protein interactions for validation.
- Experimental validation of predictions.

### Methods for Estimating Hidden Structure and Network Transitions in Genomics

#### Future Work

#### Future Work

2017-04-27

Future Work

#### Estimating Network State Transitions:

- Integrating complimentary data sources into network inference and transition step.
- Use off-diagonal elements to predict protein-protein interactions for validation.
- Experimental validation of predictions.



# Acknowledgements

## Dissertation Committee

- John Quackenbush
- Christoph Lange
- Kimberly Glass

## Channing Division of Network Medicine

- Ed Silverman
- Craig Hersh

## JQ Lab

- Joe Barry
- Joey Chen
- Maud Fagny
- Marieke Kuijjer
- Camila Lopes-Ramos
- Megha Padi
- Joe Paulson
- John Platig
- Heather Selby
- Abhijeet Sonawane
- Nicole Trotman

## Methods for Estimating Hidden Structure and Network Transitions in Genomics

### Future Work

#### Acknowledgements

2017-04-27

## Acknowledgements

All Lab
Joe Barry
Joey Chen
Maud Fagny
Marieke Kuijjer
Camila Lopes-Ramos
Megha Padi
Joe Paulson
John Platig
Heather Selby
Abhijeet Sonawane
Nicole Trotman

Channing Division of Network Medicine
Ed Silverman
Craig Hersh



└ Appendix

2017-04-27

# Appendix



# Batch covariance model

Consider a set of  $N$  samples with  $q$  covariates measuring gene expression across  $p$  genes. Let  $\mathbf{x}_i = (x_{i1}, \dots, x_{iq})$  denote the covariates for sample  $i$  and let  $\mathbf{g}_i = (g_{i1}, \dots, g_{ip})^T$  denote the gene expression values for sample  $i$  for the  $p$  genes.

We can express a model for the gene expression as

$$\mathbf{g}_i = \beta^T \mathbf{x}_i + \epsilon_i \text{ for } i = 1, \dots, N$$

where  $\epsilon_i \sim MVN_p(\mathbf{0}, \Sigma_i)$ . Notably, the covariance of  $\epsilon_i$  differ according to  $i$ .

$$\Sigma_i = \mathbf{Q} \mathbf{D}_i \mathbf{Q}^T$$

where  $\mathbf{D}_i$  is a diagonal matrix with diagonal defined as  $\mathbf{X}_i \Psi_{q \times p}$ .

$$\mathbf{S}_i = \sum_{j=1}^p \mathbf{Q}_j \mathbf{X}_i \Psi_{\cdot j} \mathbf{Q}_j^T + \mathbf{E}_i$$



## Batch covariance model

### Batch covariance model

Consider a set of  $N$  samples with  $q$  covariates measuring gene expression across  $p$  genes. Let  $\mathbf{x}_i = (x_{i1}, \dots, x_{iq})$  denote the covariates for sample  $i$  and let  $\mathbf{g}_i = (g_{i1}, \dots, g_{ip})^T$  denote the gene expression values for sample  $i$  for the  $p$  genes.

We can express a model for the gene expression as

$$\mathbf{g}_i = \beta^T \mathbf{x}_i + \epsilon_i \text{ for } i = 1, \dots, N$$

where  $\epsilon_i \sim MVN_p(\mathbf{0}, \Sigma_i)$ . Notably, the covariance of  $\epsilon_i$  differ according to  $i$ .

$$\Sigma_i = \mathbf{Q}_i \mathbf{D}_i \mathbf{Q}_i^T$$

$$\mathbf{S}_i = \sum_{j=1}^p \mathbf{Q}_j \mathbf{X}_i \Psi_{\cdot j} \mathbf{Q}_j^T + \mathbf{E}_i$$

# Least Squares Estimator For Batch Correction

To calculate least squares solution,  $\hat{\Psi}$ , we solve separately for each column of  $\mathbf{Q}$ .

and note that the residual matrices should be orthogonal to the hyperplane spanned by  $X^T$ .

$$\text{Recall } 0 = \mathbf{X}^T (\mathbf{Y} - \mathbf{X}\hat{\beta})$$

$$\mathbf{0}_q = \sum_{i=1}^N \mathbf{X}_i^T \left[ \mathbf{Q}_h^T \left[ \mathbf{G}_i^* \mathbf{G}_i^{*T} - \mathbf{Q}_h \mathbf{X}_i \hat{\Psi}_h \mathbf{Q}_h^T \right] \mathbf{Q}_h \right]$$

$$\hat{\Psi}_h = (\mathbf{X}^T \mathbf{X})^{-1} \sum_{i=1}^N [\mathbf{X}_i^T \mathbf{Q}_h^T \mathbf{G}_i^* \mathbf{G}_i^{*T} \mathbf{Q}_h]$$

## Least Squares Estimator For Batch Correction

- point 1
- point 2

To calculate least squares solution,  $\hat{\Psi}$ , we solve separately for each column of  $\mathbf{Q}$  and note that the residual matrices should be orthogonal to the hyperplane spanned by  $X^T$ .  
Recall  $0 = \mathbf{X}^T (\mathbf{Y} - \mathbf{X}\hat{\beta})$

$$\mathbf{0}_q = \sum_{i=1}^N \mathbf{X}^T \left[ \mathbf{Q}_h^T \left[ \mathbf{G}_i^* \mathbf{G}_i^{*T} - \mathbf{Q}_h \mathbf{X}_i \hat{\Psi}_h \mathbf{Q}_h^T \right] \mathbf{Q}_h \right]$$

$$\hat{\Psi}_h = (\mathbf{X}^T \mathbf{X})^{-1} \sum_{i=1}^N [\mathbf{X}_i^T \mathbf{Q}_h^T \mathbf{G}_i^* \mathbf{G}_i^{*T} \mathbf{Q}_h]$$



# The Corrected Coexpression Matrix

With the estimates obtained with our method, it is straightforward to see how fitted values for the coexpression matrix for each sample or experimental condition can be obtained. Given an estimate for  $\Psi$ ,  $\hat{\Psi}$ , we can now estimate the batch-independent coexpression structure as

$$\hat{\mathbf{S}} = \mathbf{Q} \text{diag}(\bar{\mathbf{X}}\hat{\Psi}) \mathbf{Q}^T \text{ or } \hat{\mathbf{S}} = \sum_{i=1}^p \bar{\mathbf{X}}\hat{\Psi}_i \mathbf{Q}_i \mathbf{Q}_i^T$$

The differential coexpression matrix between two conditions, defined in binary as column 2 of  $\mathbf{X}$ , is computed

$$\hat{\mathbf{W}} = \mathbf{Q} \text{diag}(\hat{\Psi}_{2,.}) \mathbf{Q}^T$$

## The Corrected Coexpression Matrix

- point 1
- point 2

### The Corrected Coexpression Matrix

With the estimates obtained with our method, it is straightforward to see how fitted values for the coexpression matrix for each sample or experimental condition can be obtained. Given an estimate for  $\Psi$ ,  $\hat{\Psi}$ , we can now estimate the batch-independent coexpression structure as

$$\hat{\mathbf{S}} = \mathbf{Q} \text{diag}(\bar{\mathbf{X}}\hat{\Psi}) \mathbf{Q}^T \text{ or } \hat{\mathbf{S}} = \sum_{i=1}^p \bar{\mathbf{X}}\hat{\Psi}_i \mathbf{Q}_i \mathbf{Q}_i^T$$

The differential coexpression matrix between two conditions, defined in binary as column 2 of  $\mathbf{X}$ , is computed

$$\hat{\mathbf{W}} = \mathbf{Q} \text{diag}(\hat{\Psi}_{2,.}) \mathbf{Q}^T$$



# Test Statistic

In the absence of population structure, cryptic relatedness and dependence between loci the distribution of the similarity index,  $s_{i,j}$

$$s_{i,j} \sim N(1, \sigma_{i,j}^2)$$

Where the variance of  $s_{ij}$  can be estimated by

$$\hat{\sigma}_{i,j}^2 = \hat{Var}(s_{i,j}) = \frac{\sum_{k=1}^N (w_k - 1)}{\left(\sum_{k=1}^N I\left[\sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1\right]\right)^2}$$

$$s_{i,j}^{(diploid)} = \frac{\sum_{k=1}^N [w_k \mathbf{H}_{i,k} \mathbf{H}_{j,k}] / 4}{\sum_{k=1}^N I[(\sum_{l=1}^n \mathbf{H}_{l,k}) > 1]}$$

## Test Statistic

### Test Statistic

In the absence of population structure, cryptic relatedness and dependence between loci the distribution of the similarity index,  $s_{i,j}$

$$s_{i,j} \sim N(1, \sigma_{i,j}^2)$$

Where the variance of  $s_{ij}$  can be estimated by

$$\hat{\sigma}_{i,j}^2 = \hat{Var}(s_{i,j}) = \frac{\sum_{k=1}^N (w_k - 1)}{\left(\sum_{k=1}^N I\left[\sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1\right]\right)^2}$$

$$s_{i,j}^{(diploid)} = \frac{\sum_{k=1}^N [w_k \mathbf{H}_{i,k} \mathbf{H}_{j,k}] / 4}{\sum_{k=1}^N I[(\sum_{l=1}^n \mathbf{H}_{l,k}) > 1]}$$



# Estimating relatedness

2017-04-27

## └ Estimating relatedness

$$\hat{\phi}_{i,j} = \frac{s_{i,j} - 1}{\left[ \frac{\sum_{k=1}^N \hat{p}_k w_k}{\sum_{k=1}^N I[\sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1]} - 1 \right]}$$

$$R : \max(s_{i,j}) > 1 - \text{probit} \left( \frac{\alpha}{\binom{n}{2}} \right)$$

$$P(\text{Reject } H_0 | \phi_{i,j} = \gamma) = \alpha + (1 - \alpha) \left( 1 - \Phi \left( \frac{\mu_{i,j} - 1}{\sqrt{\hat{\sigma}_{i,j}^2}} \right) \right)$$

Estimating relatedness

$$\hat{\phi}_{i,j} = \frac{s_{i,j} - 1}{\left[ \frac{\sum_{k=1}^N \hat{p}_k w_k}{\sum_{k=1}^N I[\sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1]} - 1 \right]}$$

$$R : \max(s_{i,j}) > 1 - \text{probit} \left( \frac{\alpha}{\binom{n}{2}} \right)$$

$$P(\text{Reject } H_0 | \phi_{i,j} = \gamma) = \alpha + (1 - \alpha) \left( 1 - \Phi \left( \frac{\mu_{i,j} - 1}{\sqrt{\hat{\sigma}_{i,j}^2}} \right) \right)$$



# Network Transition Regularization

2017-04-27

## Network Transition Regularization

$$\mathbf{Q}_{i,j} = \begin{cases} 1 & \text{for } i = j \neq k \\ 0 & \text{elsewhere} \end{cases},$$

which results in the minimization of the penalized residual sum of squares

$$PRSS(\mathbf{T}_{\cdot,k}) = \sum_{i=1}^p \left( \mathbf{B}_{i,k} - \sum_{j=1}^m A_{i,j} \mathbf{T}_{j,k} \right)^2 + \lambda \sqrt{\mathbf{T}'_{\cdot,k} \mathbf{Q} \mathbf{T}_{\cdot,k}}$$

An implementation of this extension is available in the R package MONSTER.



$$\mathbf{Q}_{i,j} = \begin{cases} 1 & \text{for } i = j \neq k \\ 0 & \text{elsewhere} \end{cases},$$

which results in the minimization of the penalized residual sum of squares

$$PRSS(\mathbf{T}_{\cdot,k}) = \sum_{i=1}^p \left( \mathbf{B}_{i,k} - \sum_{j=1}^m A_{i,j} \mathbf{T}_{j,k} \right)^2 + \lambda \sqrt{\mathbf{T}'_{\cdot,k} \mathbf{Q} \mathbf{T}_{\cdot,k}}$$

An implementation of this extension is available in the R package MONSTER.