

Methods for Estimating Hidden Structure and Network Transitions in Genomics

Daniel Schlauch, PhD Candidate

Department of Biostatistics
Harvard School of Public Health

April 24, 2017

Table of Contents

- 1 Batch Effect on Covariance Structure
- 2 Identification of genetic outliers
- 3 State Transitions Using Gene Regulatory Network Models



Batch Effect on Covariance Structure

Batch effect on covariance structure confounds gene coexpression

Daniel Schlauch^{1,2}, Joseph N. Paulson², Kimberly Glass^{2,3}, and
John Quackenbush^{1,3}

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

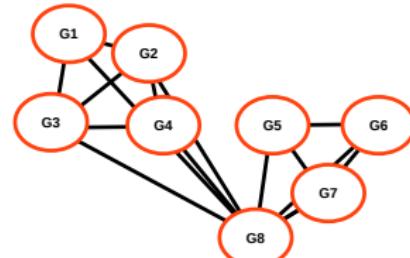
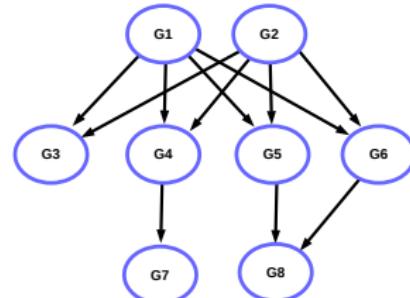
²Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA

³Department of Medicine, Harvard Medical School, Boston, MA

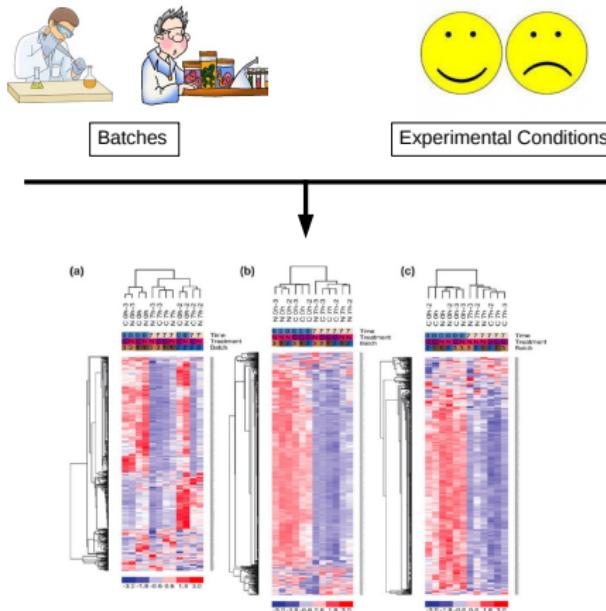
Differential Network Inference

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



Batch Effect



Unwanted variation from:

- Laboratory Conditions
- Circadian Rhythm
- Cell cycle
- Reagents
- Atmospheric Ozone
- Differential biological variation
- Etc.

Johnson et al. (Biostatistics 2007)



Methods for Controlling Batch Effect

Location scale model:

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

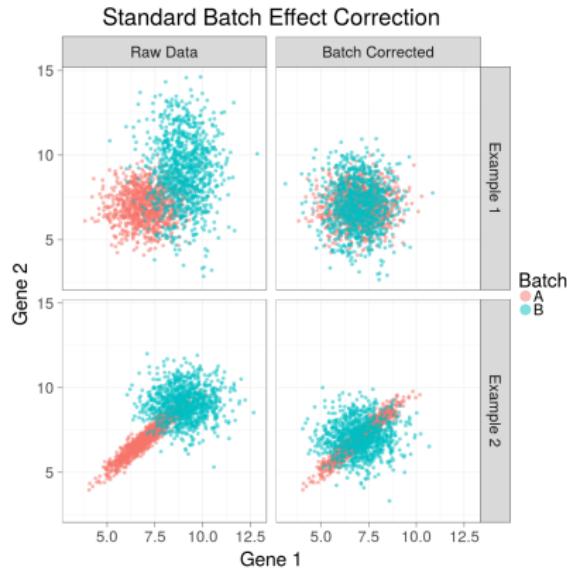
Limitations:

- 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).



Limitations to common batch effect correction methods



Standard corrections:

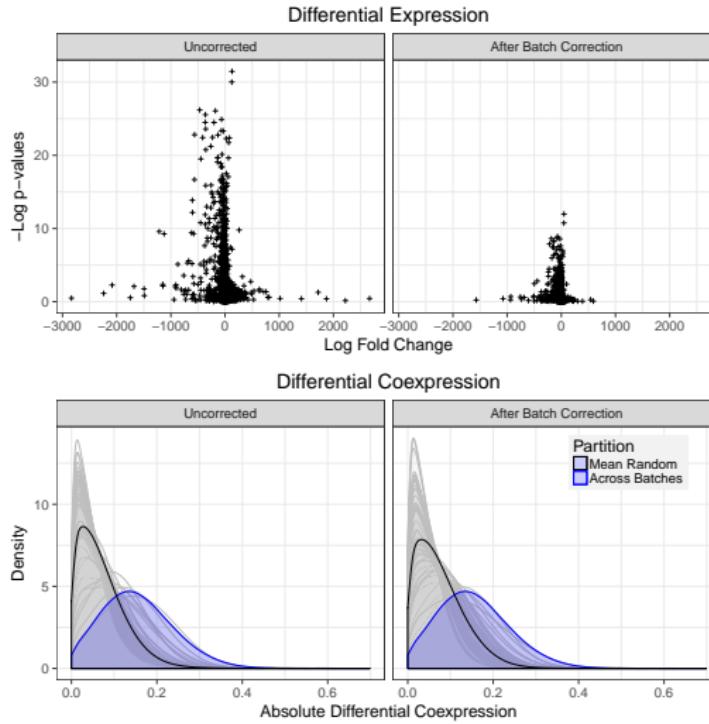
$$f[Gene1|BatchA] = f[Gene1|BatchB]$$

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

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



Limitations to existing batch effect correction methods

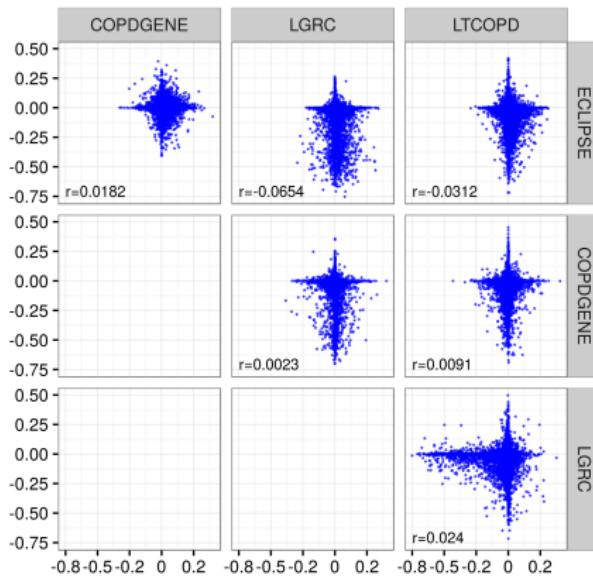


ENCODE Project:
50k genes
126 samples (63 patients RNA-seq'd at 2 centers)



Challenges with batch effect on differential coexpression

WGCNA edge weight differences between pairs of studies



- Ultra-high dimensionality
- Differential coexpression
- Modularity



Estimating the conditional coexpression matrix

Motivating concepts:

- Provide a regression framework for the coexpression matrix
- Estimate a reduced number of parameters
- Exploit modular nature of gene expression patterns

Our proposal:

- Define our parameters as functions of components of variation.
- Estimate the eigenvalue contribution of each eigenvector.



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

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

where $\boldsymbol{\epsilon}_i \sim MVN_p(\mathbf{0}, \Sigma_i)$. Notably, the covariance of $\boldsymbol{\epsilon}_i$ 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 $\mathbf{X}_i \Psi_{q \times p}$

Likelihood Function

$$\mathcal{L}(\mu, \Sigma) = \prod_{i=1}^N \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma_i|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{G}_i - \mu)^T \Sigma_i^{-1} (\mathbf{G}_i - \mu)}$$

Where we define Σ_i ,

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

Where \mathbf{Q} is a matrix with columns defined as the eigenvectors of the estimated coexpression matrix, $\mathbf{G}^* \mathbf{G}^{*T} / N$.



Least Squares Estimator

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 .

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

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

$$\hat{\Psi}_h = (\mathbf{X}^T \mathbf{X})^{-1} \sum_{i=1}^N [\mathbf{X}_i^T \mathbf{Q}_h^T \mathbf{G}_i^* \mathbf{G}_i^{*T} \mathbf{Q}_h] \quad (2)$$



The Corrected Coexpression Matrix

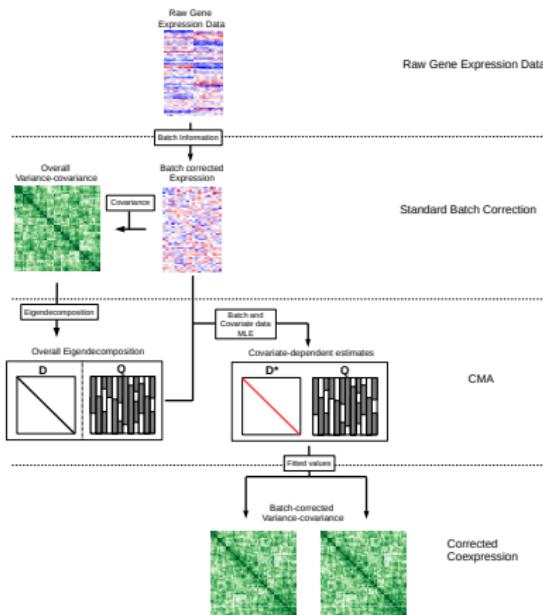
With the estimates obtained with our method, it is straightforward to see how fitted values for the coexpression matrix for each sample or experimental condition can be obtained. Given an estimate for Ψ , $\hat{\Psi}$, we can now estimate the batch-independent coexpression structure as

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

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

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

Example Workflow



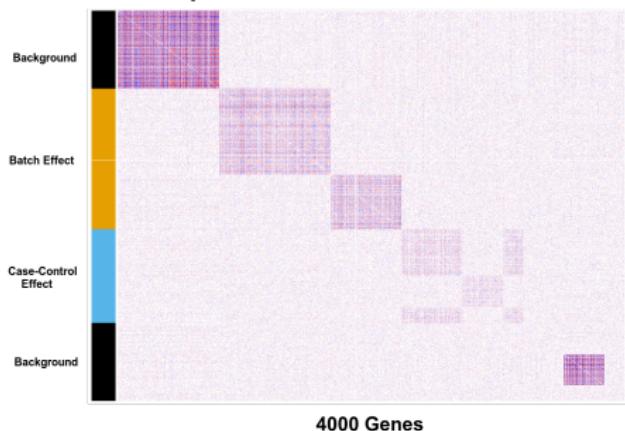
Simulations

$$g_i \sim MVN_{4000} (\mathbf{0}, \Sigma_i)$$

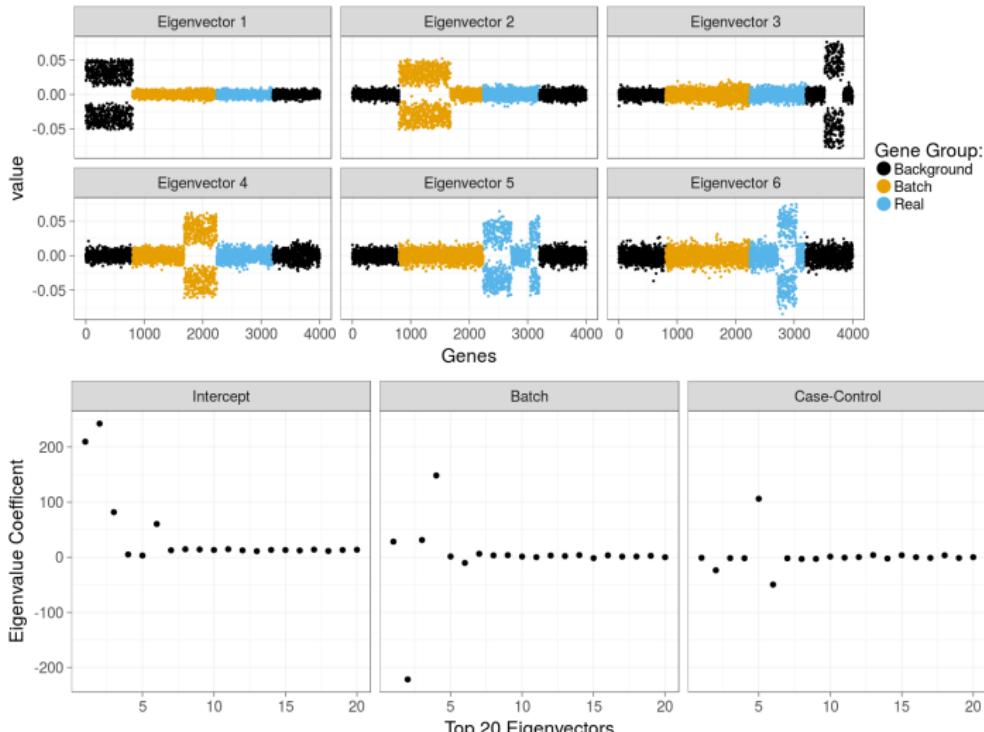
Σ_i describes the covariance of 4000 genes and contains modules which are classified as

- Background
- Batch
- Case-Control

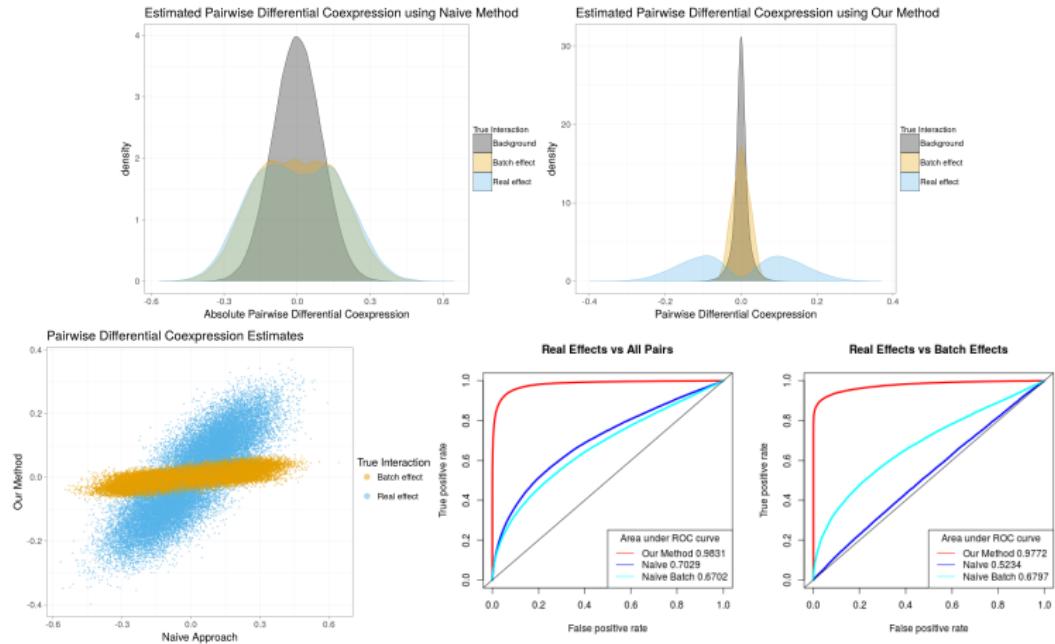
Coexpression matrix



Simulations



Simulations



Application to data from COPDGene Study

GO Term	Count	%	Enrichment	FDR
anatomical structure development	309	0.26	1.29	2.58E-05
single-organism developmental process	309	0.26	1.29	2.73E-05
anatomical structure morphogenesis	168	0.14	1.46	1.60E-04
single-multicellular organism process	324	0.27	1.25	4.01E-04
system process	132	0.11	1.50	1.46E-03
regulation of cellular process	514	0.43	1.12	5.86E-03
single organism signaling	328	0.28	1.21	7.89E-03
regulation of localization	151	0.13	1.40	1.15E-02
regulation of multicellular organismal process	156	0.13	1.35	6.56E-02

Table : GO categories for differential coexpression in COPDGene identified with CPBA found with FDR<0.1.

Identifying Genetic Outliers

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

Daniel Schlauch^{1,2,4}, Heide Fier^{1,3} and Christoph Lange^{1,4}

¹Department of Biostatistics, Harvard TH Chan School of Public Health, Boston, MA 02115

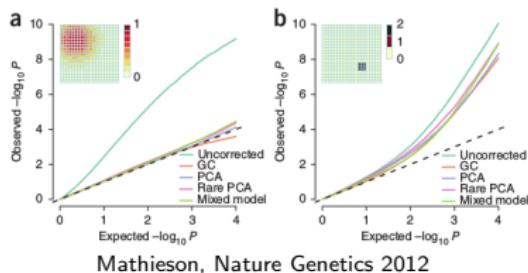
²Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute, Boston, MA 02115

³Institute of Genomic Mathematics, University of Bonn, Bonn, Germany

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

Background

- Individuals may be too similar (due to cryptic relatedness) or too different (due to population structure).
- Both features may lead to spurious results, inflation of type I error.
- Many methods exist for addressing some of these concerns (e.g. PCA, LMM).
- Limitations exist, such as with rare alleles and sharp localized effects, or with the assumption of linear or discrete population structure.



Mathieson, Nature Genetics 2012

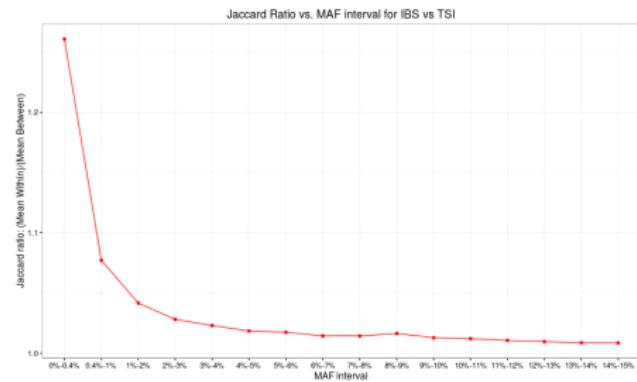
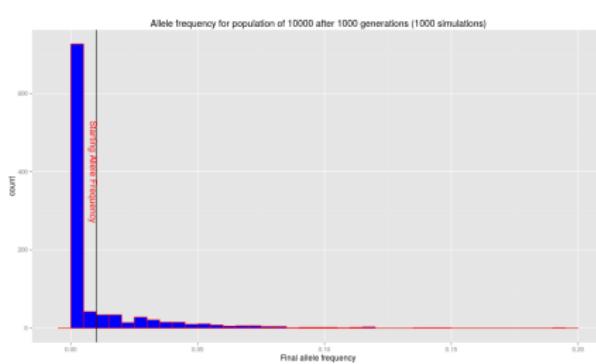
We want to create a similarity measure that...

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



Basis for measure

- Rare variants are recent variants.
- In the absence of selection, rare variants become fixed at 0% with high probability over a relatively short timeframe.



$$P[\text{Fixation} | n=10000, g=1000, \text{maf}=.01] = .678$$

- Key Idea: **Less frequent variants are more informative of ancestry.**

Test Statistic

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

where

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

$$E [s_{i,j}] = 1$$

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

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

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

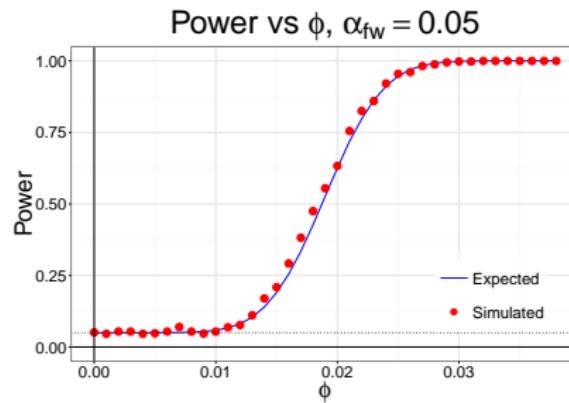


Tests of Heterogeneity

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

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

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



Tests of Heterogeneity

$$H_0 : \mu_{i,j} = 1 \forall i, j \in 1 \dots n$$

$$H_A : \exists i, j \in 1 \dots n | \mu_{i,j} \neq 1$$

Test for population structure:

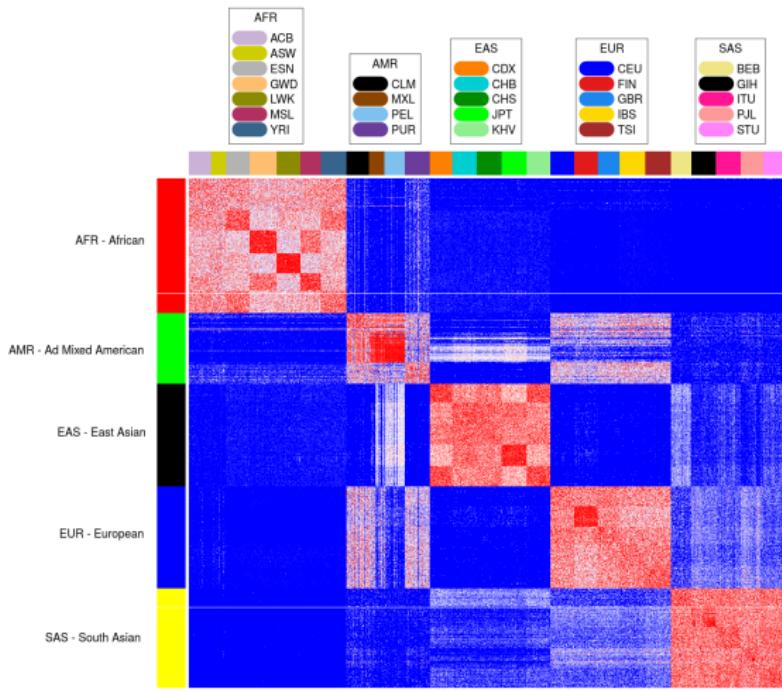
$$K = \sup_x |F_s(x) - \Phi(x)|$$

Test for cryptic relatedness:

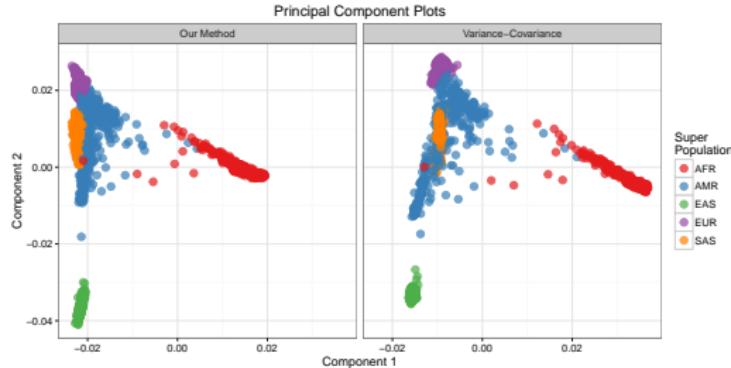
$$R : \max(s_{i,j}) > 1 - \text{probit} \left(\frac{\alpha}{\binom{n}{2}} \right)$$

Application to 1000 Genomes Project

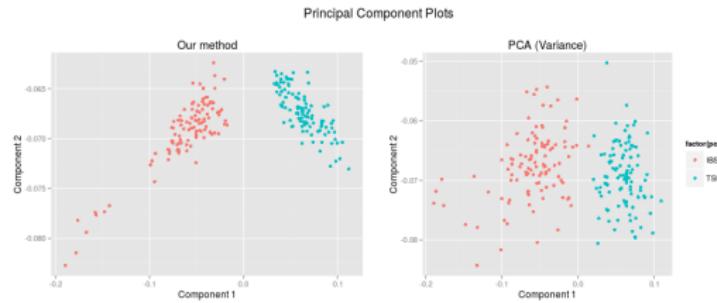
Genetic Similarity Matrix



Application to 1000 Genomes Project



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

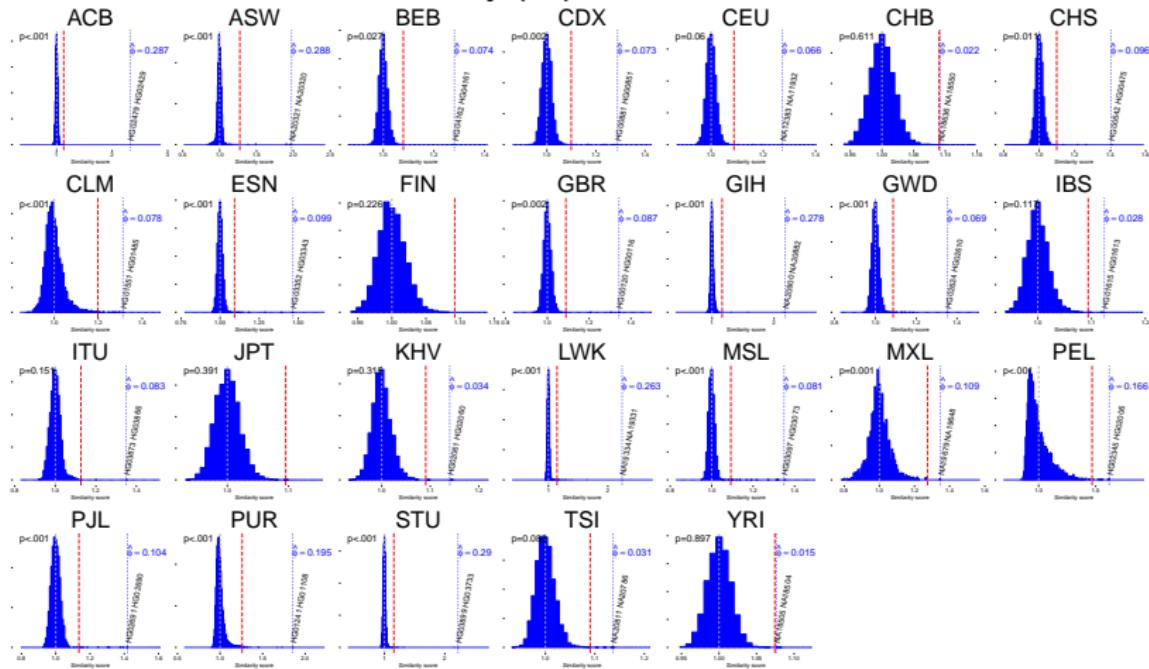


But produces superior separation for recently related populations.



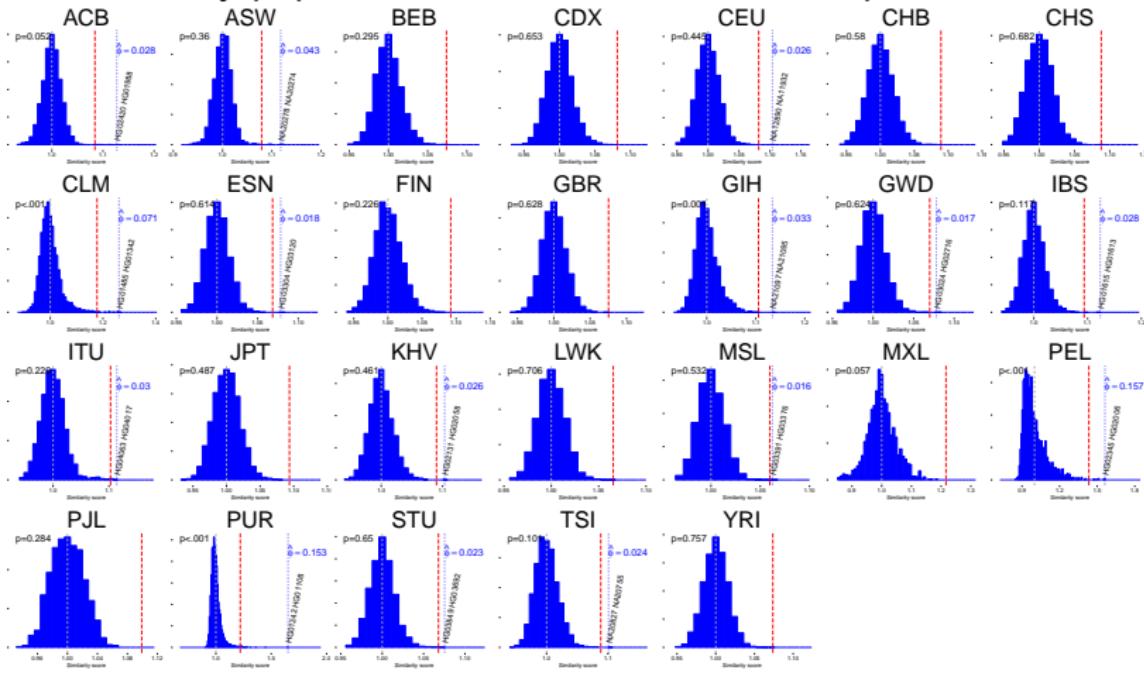
Application to 1000 Genomes Project

Distribution of s statistics by population



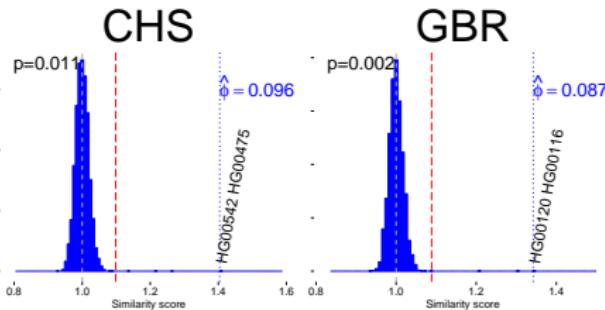
Application to 1000 Genomes Project

s statistics by population after removal of related pairs

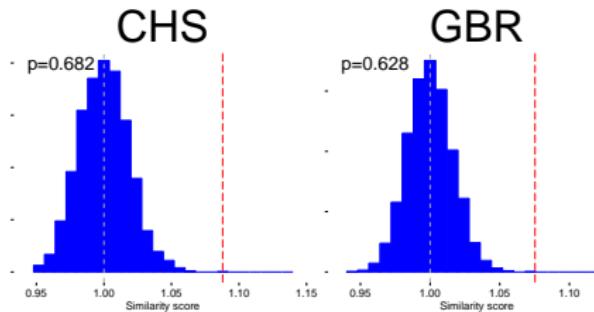


Application to 1000 Genomes Project

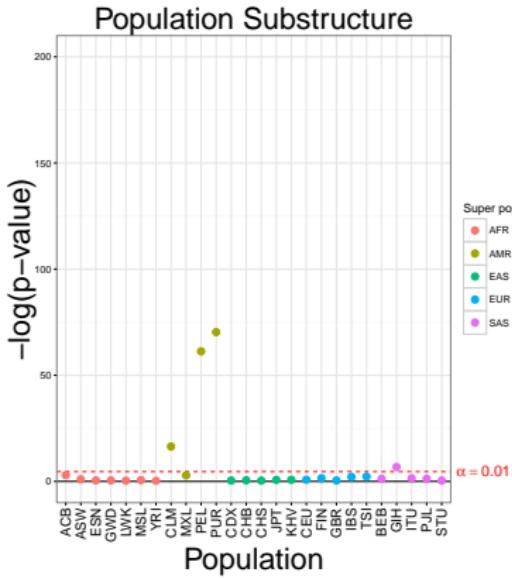
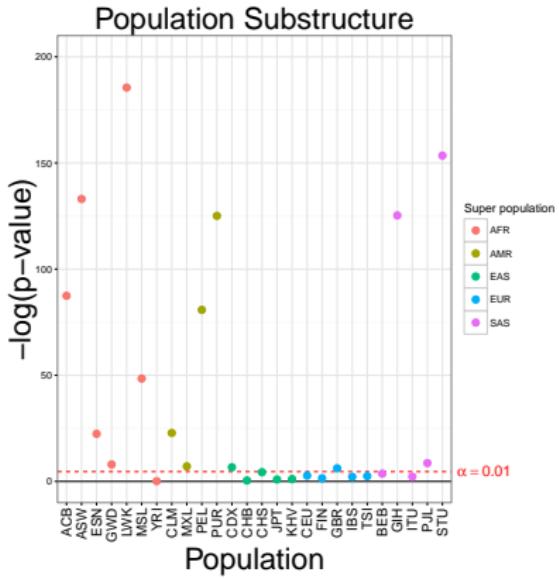
Original Data



After removal of
related pairs



Application to 1000 Genomes Project



State Transitions Using Gene Regulatory Network Models

Estimating Drivers of Cell State Transitions Using Gene Regulatory Network Models

Daniel Schlauch^{1,2}, Kimberly Glass^{2,3}, Craig P. Hersh^{2,3,4}, Edwin K. Silverman^{2,3,4} and John Quackenbush^{1,2,3}

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

²Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA

³Department of Medicine, Harvard Medical School, Boston, MA

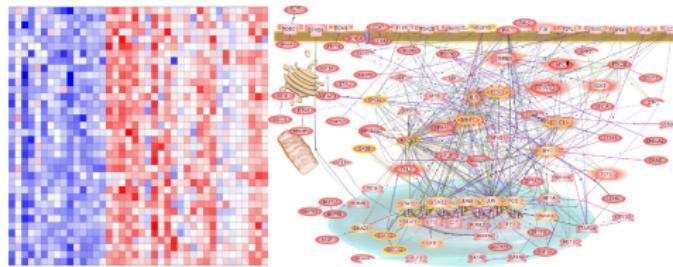
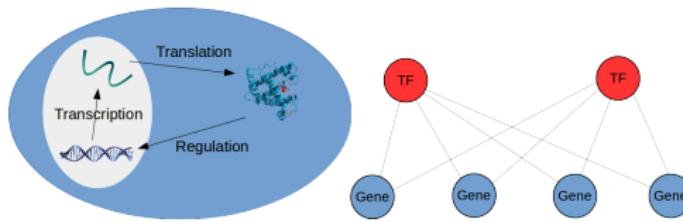
⁴Pulmonary and Critical Care Division, Brigham and Women's Hospital and Harvard Medical School, Boston, MA



Background

Why Study Gene Regulatory Networks?

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



Abdollahi et al. PNAS 2007



Background

Biological Challenges

- 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

Statistical Challenges

- 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”



Background

Biological Challenges

- 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

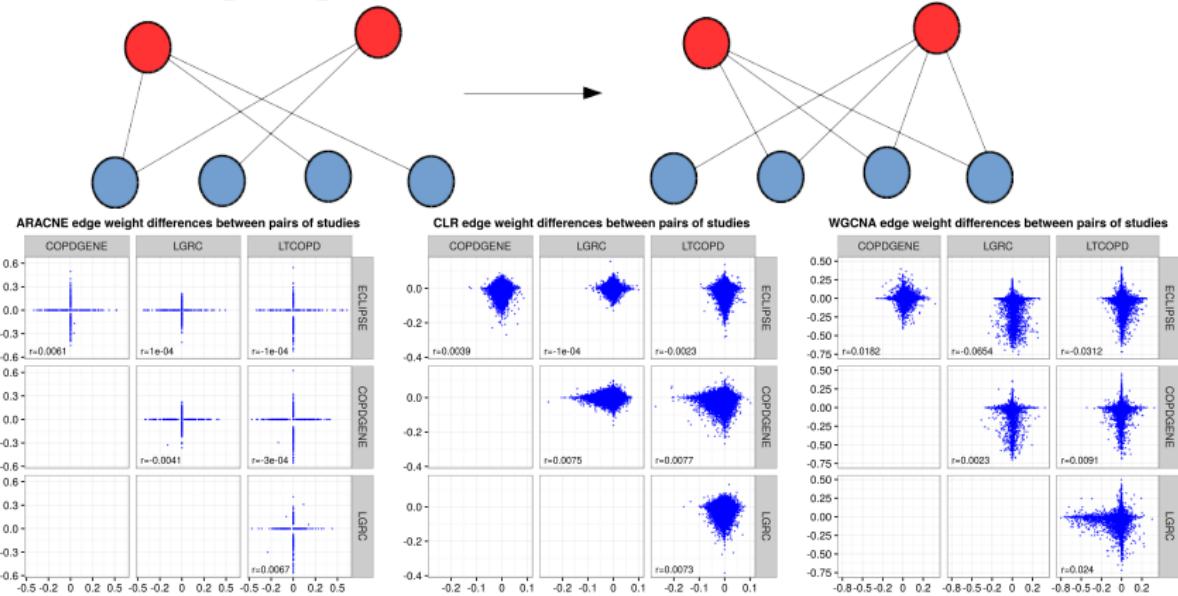
Statistical Challenges

- 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”

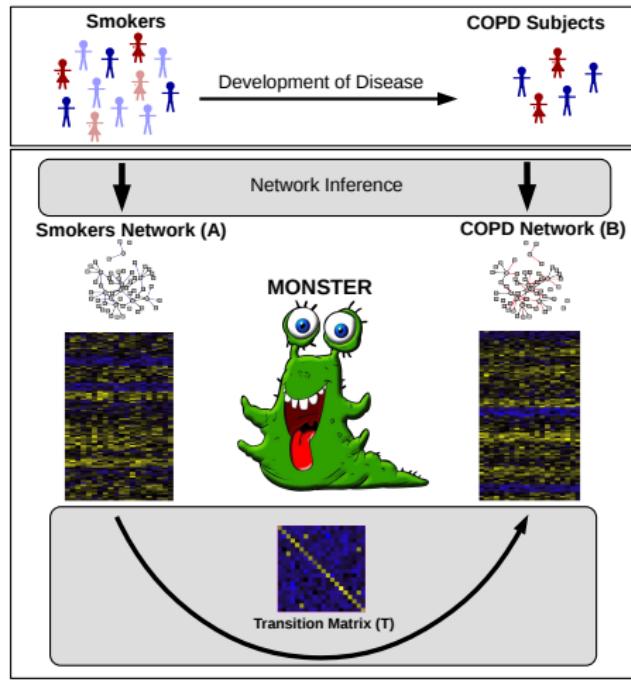


Replication of inferred networks in COPD

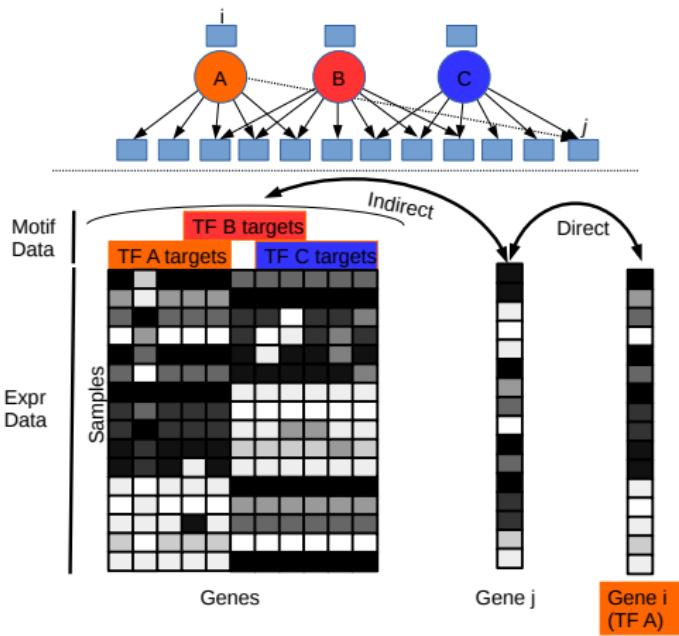
Differential edgeweights do not correlate between studies of COPD



Algorithm Overview



Network Inference



Network Inference

Direct Evidence:

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

Indirect Evidence:

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

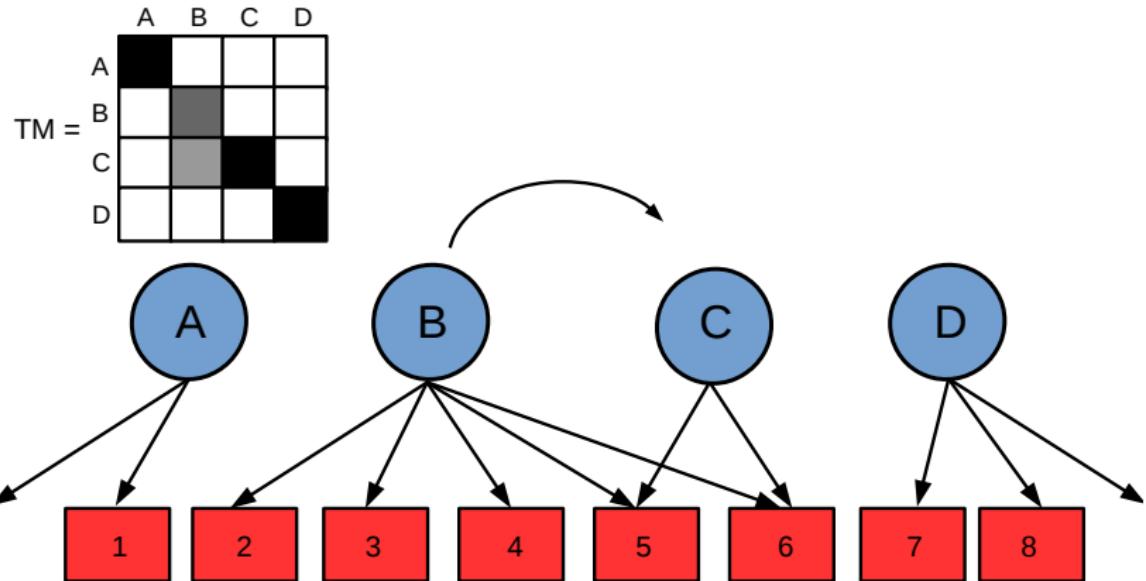
$$e_{i,j} = \frac{1}{1 + e^{\beta_0 + \beta_1 g_{j,(1)} + \cdots + \beta_k g_{j,(k)}}}$$

Edgeweight:

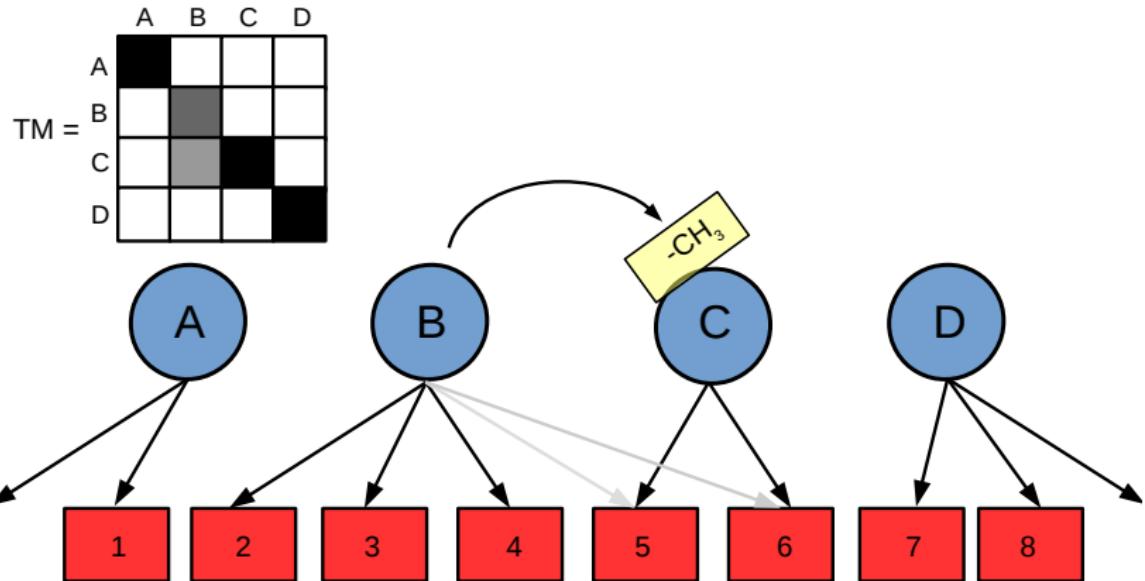
$$w_{i,j} = (1 - \alpha) [d_{i,j}] + \alpha [e_{i,j}]$$



Network Transition



Network Transition



Network Transition

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

where b_i and a_i are column-vectors in \mathbf{B} and \mathbf{A} that describe the regulatory targeting of transcription factor i in the final and initial networks, respectively.

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

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

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

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

Network Transition

Regularization:

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

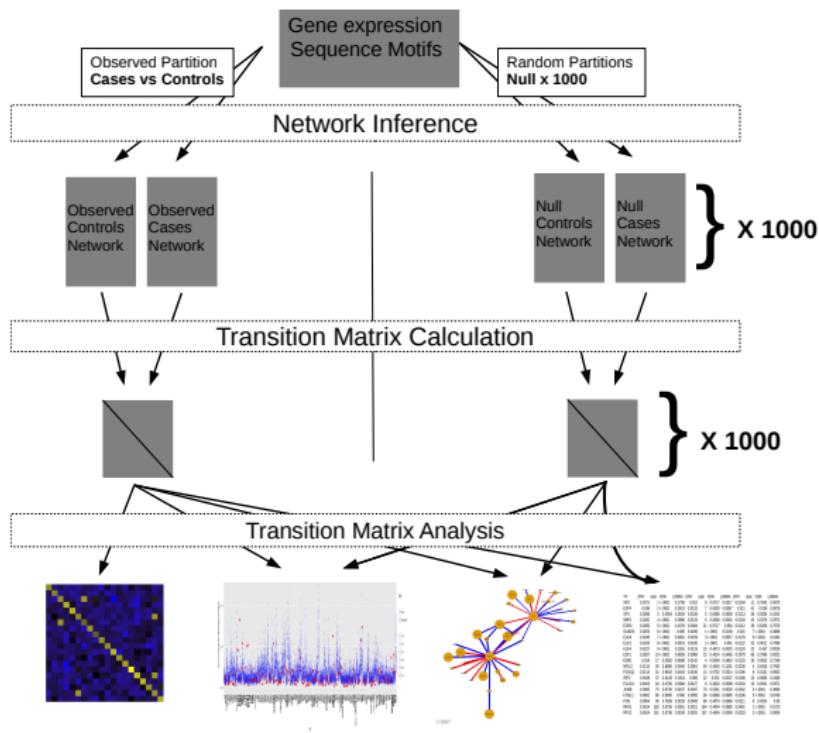
which results in the minimization of the penalized residual sum of squares

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

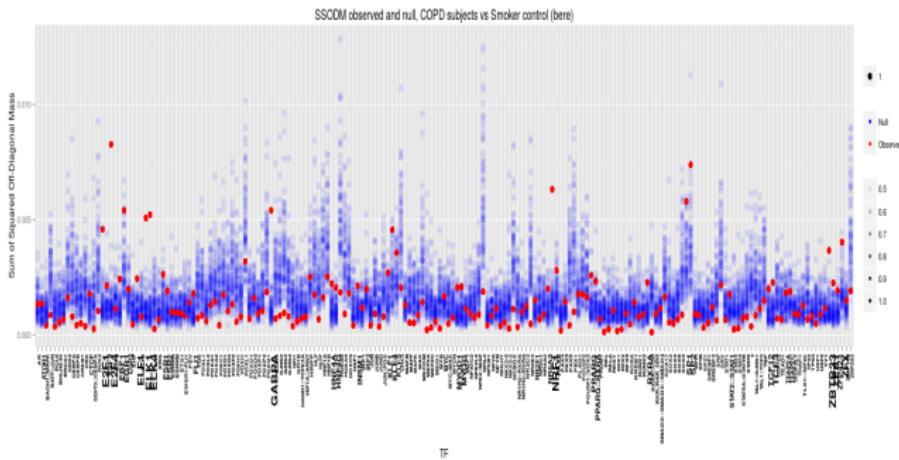
An implementation of this extension is available in the R package MONSTER.



Network Transition



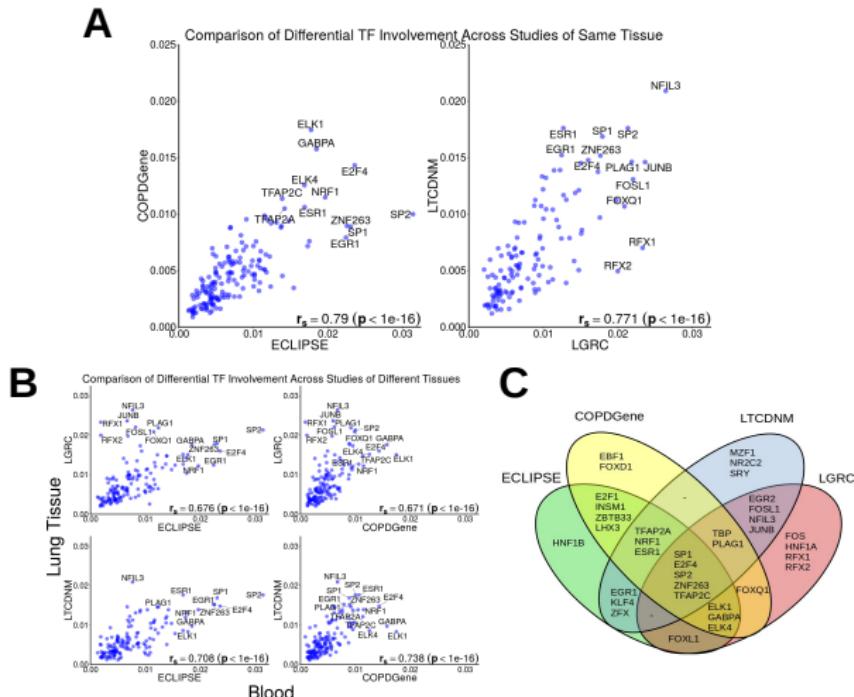
Evaluating Transition Matrix



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

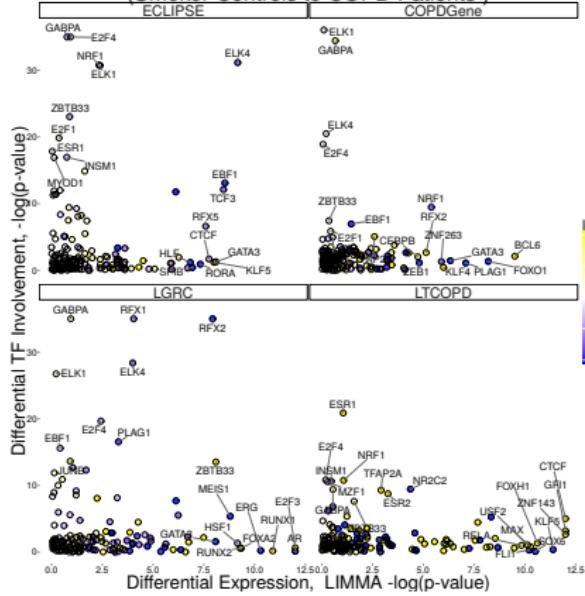


Reproducibility and novel results



Reproducibility and novel results

Differential Involvement vs Differential Expression (Smoker Controls to COPD Patients)



Differential TF Involvement

transcription factor	ECLIPSE			COPDGene			LGRC			LTCOPD		
	dTFI	rank	FlR	dTFI	rank	FlR	dTFI	rank	FlR	dTFI	rank	FlR
SMP2	0.014	1	0.007	0.007	9	0.013	0.01	2	0.016	2	0.005	
E2F4	0.036	2	< .0001	0.043	14	0.037	0.016	8	0.048	1	0.005	
ELK1	0.024	3	0.007	0.019	19	0.026	0.012	10	0.027	1	0.005	
ZBTB33	0.025	4	0.008	0.011	18	0.027	0.017	11	0.018	6	0.007	
ESR1	0.024	5	0.007	0.016	23	0.024	0.016	26	0.020	5	0.005	
NR1	0.024	6	0.007	0.016	24	0.024	0.016	25	0.020	4	0.005	
GABPA	0.025	7	< .0001	0.017	2	< .0001	0.016	12	< .0001	32	0.001	
ELK4	0.017	8	< .0001	0.014	17	< .0001	0.011	40	< .0001	0.003	0.004	
ZFX	0.018	9	0.007	0.016	24	0.018	0.016	27	0.018	26	0.005	
RFX4	0.017	10	0.005	0.007	28	0.012	0.003	21	0.019	55	0.008	
ESR2	0.009	11	0.007	0.016	7	0.001	0.017	3	0.008	0.016	0.004	
ELK5	0.012	12	0.007	0.016	14	0.010	0.016	19	0.017	0.021	0.009	
TEAD2C	0.020	13	0.005	0.014	6	0.005	0.016	19	0.017	0.021	0.009	
PLAGL1	0.024	21	0.01	0.002	15	0.019	5	0.001	0.046	8	0.004	
FOSQ2	0.015	26	0.007	0.006	10	0.005	0.009	20	0.013	4	0.005	
ELK3	0.014	27	0.007	0.006	11	0.005	0.009	21	0.013	17	0.005	
NRF1	0.007	62	0.005	0.007	33	0.004	0.004	1	0.006	7121	0.005	
FOS	0.005	73	0.01	0.007	26	0.005	0.009	9	0.012	26	0.009	
ELK2	0.001	93	0.007	0.006	43	0.003	0.009	14	< .0001	0.001	0.005	
RFX1	0.003	195	0.012	0.008	164	< .0001	0.023	3	< .0001	0.005	48	0.005
RFX2	0.009	258	0.001	0.012	183	0.002	0.000	8	< .0001	0.009	81	0.005

Differential Expression

transcription factor	ECLIPSE			COPDGene			LGRC			LTCOPD		
	dTFI	rank	LIMMA p	dTFI	rank	LIMMA p	dTFI	rank	LIMMA p	dTFI	rank	LIMMA p
SMP2	1	1758	0.001	0.011	6	0.017	0	0.005	2	0.005		
E2F4	2	3013	3	0.007	14	0.019	7	0.018	7	0.012		
ELK1	3	214	0.001	0.019	19	0.019	8	0.018	8	0.012		
ZNF263	4	9814	16	0.008	31	0.021	6	0.019	15	0.019		
ESR1	5	4379	23	0.004	28	0.019	26	0.019	5	0.018		
ELK5	6	2668	5	0.008	30	0.019	30	0.019	11	0.019		
GABPA	7	1000	13	0.001	33	0.019	33	0.019	17	0.019		
ELK4	8	3813	1	0.015	37	0.028	40	0.028	40	0.005		
ZFX	9	8253	24	0.015	40	0.024	38	0.024	58	0.017		
ELK2	10	3411	20	0.003	25	0.019	25	0.019	25	0.019		
ESR2	11	3696	7	0.003	33	0.019	33	0.019	33	0.017		
ELK4	12	3061	4	0.007	38	0.018	38	0.018	38	0.014		
TEAD2C	17	2118		0.004	39	0.014		0.014	39	0.014		
ELK3	18	2100	15	0.008	39	0.014	39	0.014	39	0.014		
FOSQ2	28	4543	10	0.014	53	0.023	57	0.023	57	0.014		
FOSL1	57	5850	41	0.005	64	0.018	41	0.018	17	0.005		
NRF1	73	3556	48	0.003	66	0.023	51	0.023	51	0.013		
JUND	77	3597	43	0.005	52	0.016	39	0.016	9	0.007		
RPX1	109	3381	164	0.005	5	0.015	48	0.015	3	0.005		
RPX2	138	3309	163	0.004	8	0.004	81	0.004	81	0.005		

Acknowledgements

Dissertation Committee

- John Quackenbush
- Christoph Lange
- Kimberly Glass

Channing Division of Network Medicine

- Ed Silverman
- Craig Hersh

JQ Lab

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

