

## Network Analysis and Applications in Biology: Introduction

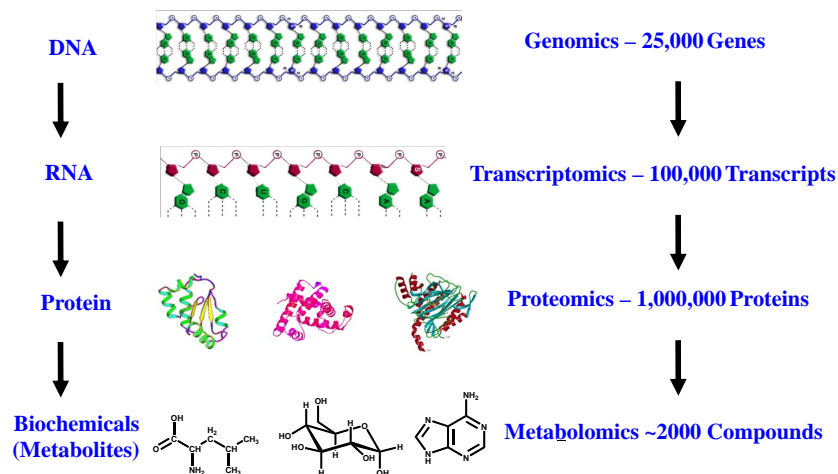
Ali Shojaie & George Michailidis

ENAR 2020

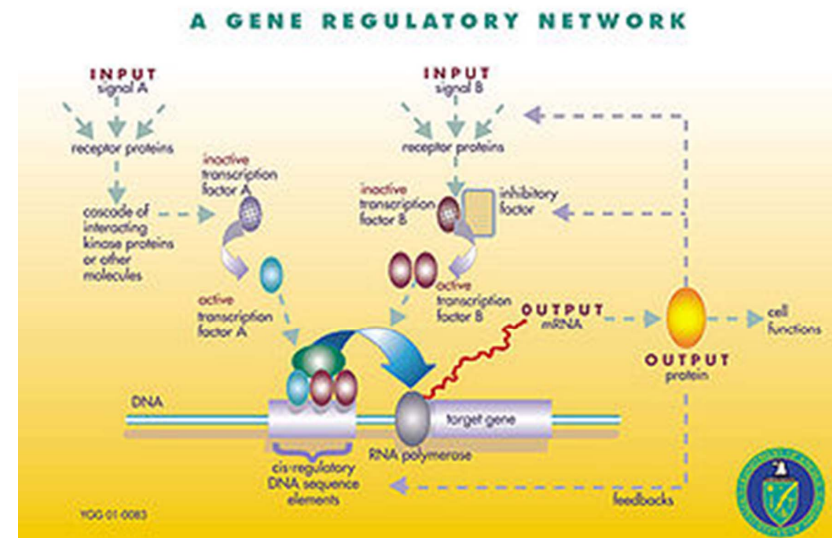
## Why Study Networks?

- ▶ Components of biological systems (genes, proteins etc) interact with each other to carry out cell functions.
- ▶ Examples of such interactions include signaling, regulation and interactions between proteins.
- ▶ We cannot understand the function and behavior of biological systems by studying individual components ( $2 + 2 \neq 4!$ ).
- ▶ Networks provide an efficient representation of complex interactions in cells, and a basis for mathematical/statistical models to study these systems.

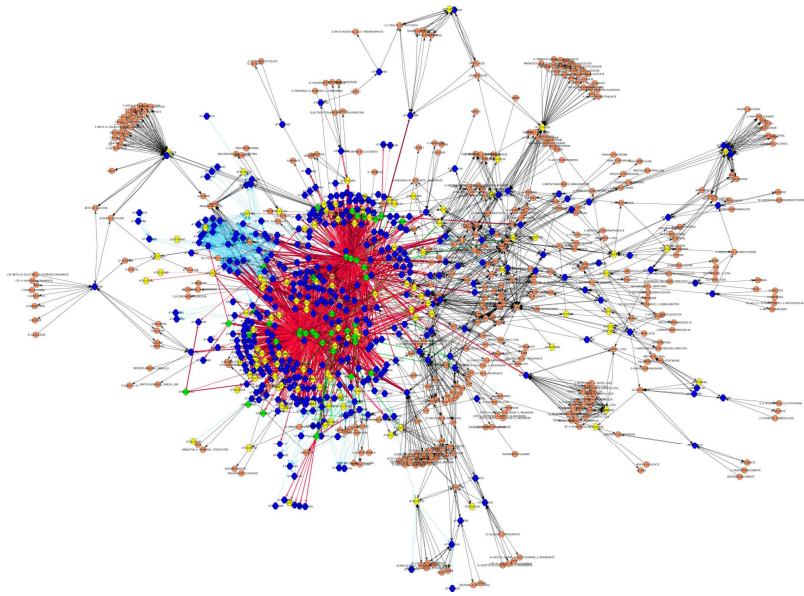
## Central Dogma of Molecular Biology (Extended)



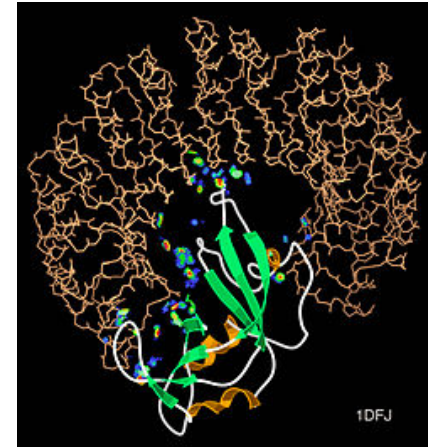
## Networks in Biology: Gene Regulatory Interactions



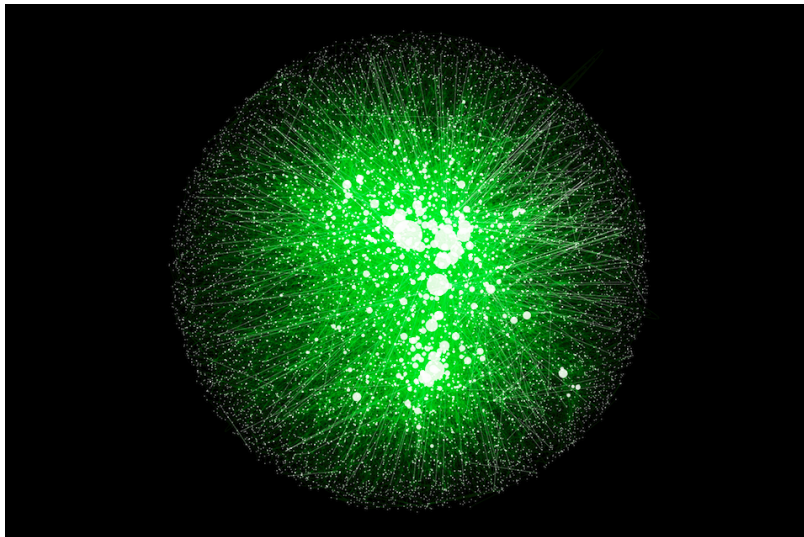
## Networks in Biology: Gene Regulatory Networks



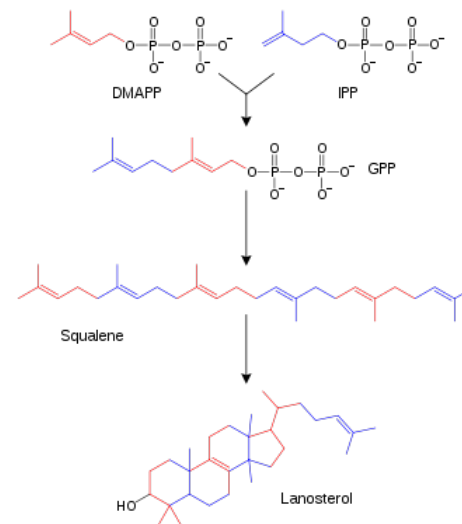
## Networks in Biology: Protein-Protein Interaction



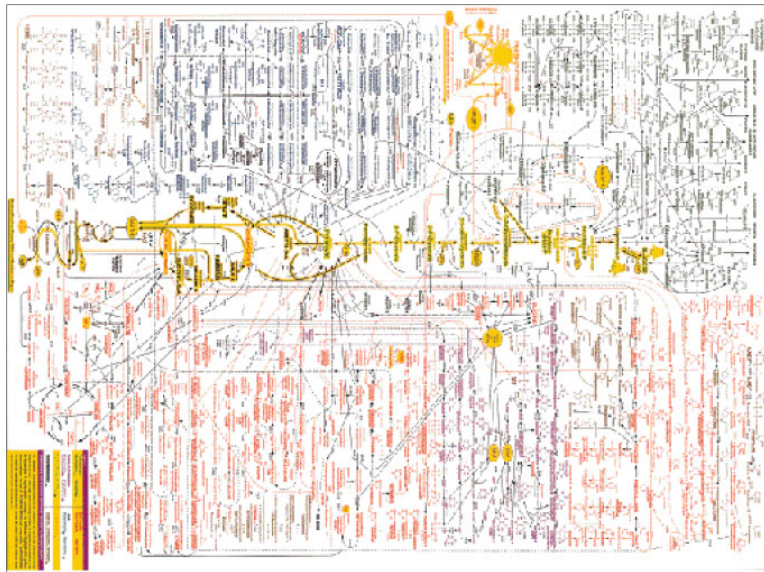
## Networks in Biology: Protein-Protein Interactions (PPI)



## Networks in Biology: Metabolic Reactions



## Networks in Biology: Metabolic Pathways



## But Do Networks Matter?

- ▶ They Do!
- ▶ Recent studies have linked changes in gene/protein networks with many human diseases.

### Systems Biology and Emerging Technologies

#### Gene Networks and microRNAs Implicated in Aggressive Prostate Cancer

Liang Wang,<sup>1</sup> Hui Tang,<sup>2</sup> Venugopal Thayanithy,<sup>3</sup> Subbaya Subramanian,<sup>3</sup> Ann L. Oberg,<sup>2</sup> Julie M. Cunningham,<sup>1</sup> James R. Cerhan,<sup>2</sup> Clifford J. Steer,<sup>4</sup> and Stephen N. Thibodeau<sup>1</sup>

<sup>1</sup>Departments of Laboratory Medicine and Pathology and <sup>2</sup>Health Sciences Research, Mayo Clinic, Rochester, Minnesota; and Departments of <sup>3</sup>Laboratory Medicine and Pathology, <sup>4</sup>Medicine, and Genetics, Cell Biology, and Development, University of Minnesota, Minneapolis, Minnesota

## But Do Networks Matter?

0888-8809/07/\$15.00/0  
Printed in U.S.A.

Molecular Endocrinology 21(9):2112-2123  
Copyright © 2007 by The Endocrine Society  
doi: 10.1210/me.2006-0474

### Estrogen-Regulated Gene Networks in Human Breast Cancer Cells: Involvement of E2F1 in the Regulation of Cell Proliferation

Joshua D. Stender, Jonna Frasor, Barry Komm, Ken C. N. Chang, W. Lee Kraus, and Benita S. Katzenellenbogen

Departments of Biochemistry (J.D.S.) and Molecular and Integrative Physiology (J.F., B.S.K.), University of Illinois at Urbana-Champaign, Urbana, Illinois 61801-3704; Women's Health and Musculoskeletal Biology (B.K., K.C.N.C.), Wyeth Research, Collegeville, Pennsylvania 19426; and Department of Molecular Biology and Genetics (W.L.K.), Cornell University, Ithaca, New York 14853-4203

## But Do Networks Matter?



Cancer Cell  
Article

### A Transcriptional Signature and Common Gene Networks Link Cancer with Lipid Metabolism and Diverse Human Diseases

Heather A. Hirsch,<sup>1,7</sup> Dimitrios Iliopoulos,<sup>1,7</sup> Amita Joshi,<sup>1,7</sup> Yong Zhang,<sup>2</sup> Savina A. Jaeger,<sup>3</sup> Martha Bulyk,<sup>3,4,5</sup> Philip N. Tschlis,<sup>6</sup> X. Shirley Liu,<sup>2</sup> and Kevin Struhl<sup>1,\*</sup>

<sup>1</sup>Department of Biological Chemistry and Molecular Pharmacology, Harvard Medical School, Boston, MA 02115, USA

<sup>2</sup>Department of Biostatistics and Computational Biology, Dana Farber Cancer Institute, Harvard School of Public Health, Boston, MA 02115, USA

<sup>3</sup>Division of Genetics, Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA 02115, USA

<sup>4</sup>Department of Pathology, Brigham and Women's Hospital and Harvard Medical School, Boston, MA 02115, USA

<sup>5</sup>Harvard/MIT Division of Health Sciences and Technology (HST), Harvard Medical School, Boston, MA 02115, USA

<sup>6</sup>Molecular Oncology Research Institute, Tufts Medical Center, Boston, MA 02111, USA

<sup>7</sup>These authors contributed equally to this work

\*Correspondence: kevin@hms.harvard.edu

DOI 10.1016/j.ccr.2010.01.022

## But Do Networks Matter?

And, incorporating the knowledge of networks **improves our ability to find causes of complex diseases.**

Molecular Systems Biology 3; Article number 140; doi:10.1038/msb4100180  
Citation: *Molecular Systems Biology* 3:140  
© 2007 EMBO and Nature Publishing Group. All rights reserved 1744-4292/07  
www.molecularsystemsbiology.com



### REPORT

## Network-based classification of breast cancer metastasis

Han-Yu Chuang<sup>1,5</sup>, Eunjung Lee<sup>2,3,5</sup>, Yu-Tsueng Liu<sup>4</sup>, Doheon Lee<sup>3</sup> and Trey Ideker<sup>1,2,4,\*</sup>

<sup>1</sup> Bioinformatics Program, University of California San Diego, La Jolla, CA, USA, <sup>2</sup> Department of Bioengineering, University of California San Diego, La Jolla, CA, USA, <sup>3</sup> Department of Bio and Brain Engineering, Korea Advanced Institute of Science and Technology, Daejeon, Korea and <sup>4</sup> Cancer Genetics Program, Moores Cancer Center, University of California San Diego, La Jolla, CA, USA

<sup>5</sup> These authors contributed equally to this work

\* Corresponding author. Department of Bioengineering, University of California San Diego, La Jolla, CA 92093, USA. Tel.: +1 858 822 4558; Fax: +1 858 534 5722; E-mail: trey@bioeng.ucsd.edu

## Networks: A Short Primer

- ▶ A network is a collection of **nodes**  $V$  and **edges**  $E$ .
- ▶ We assume there are  $p$  nodes in the network, and that the **nodes correspond to random variables**  $X_1, \dots, X_p$ .
- ▶ Edges can be **undirected**  $X - Y$  or **directed**  $X \rightarrow Y$ .

- 
- ▶ Consider the **node set**  $V = \{1, 2, 3\}$ .
  - ▶ Then **edges** can be:

$$\text{undirected: } E_1 = \{1 - 2, 2 - 3\}$$

$$\text{directed: } E_2 = \{1 \rightarrow 3, 3 \rightarrow 2\}$$

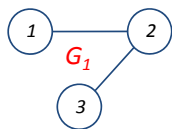
- ★ We focus primarily on **undirected** networks.

## Networks: A Short Primer

- ▶ A convenient way to represent the **edges** of the network is to use an **adjacency matrix**  $A$
- ▶ Adjacency matrix is a **square** matrix, with a **nonzero entry in  $(i, j)$  and  $(j, i)$  if there is an edge between nodes  $i$  and  $j$**

$$A = \begin{bmatrix} \cdot & \mathbf{x} & \cdot \\ \mathbf{x} & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \rightarrow \mathbf{x} \text{ shows an edge between 1 and 2}$$

Example:



$$A = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

## What Do Edges in Biological Networks Mean?

- ▶ In **gene regulatory networks**, an edge from gene  $i$  to gene  $j$  often means that  $i$  **controls the expression of  $j$** : as  $i$ 's expression changes,  $j$ 's expression also increases/decreases.
- ▶ In **protein-protein interaction networks**, an edge between proteins  $i$  and  $j$  often means that *the two proteins bind together and form a protein complex*. Therefore, we expect that these proteins are generated at similar rates.
- ▶ In **metabolic networks**, an edge between compound  $i$  and  $j$  often means that *the two compounds are involved in the same reaction*, meaning that they are generated at relative rates.
- ▶ Thus, edges represent some type of **association among genes, proteins or metabolites**, defined generally to include *linear or nonlinear* associations; more later....

## Statistical Models for Biological Networks

- ▶ We use the framework of **graphical models**
- ▶ In this setting, **nodes correspond to “random variables”**
- ▶ In other words, each node of the network represents one of the variables in the study
  - ▶ In gene regulatory networks, **nodes  $\equiv$  genes**
  - ▶ In PPI networks, **nodes  $\equiv$  proteins**
  - ▶ In metabolic networks, **nodes  $\equiv$  metabolites**
- ▶ In practice, we observe  $n$  measurements of each of the variables (genes/proteins/ metabolites) for say different individuals, and want to determine which variables are connected, or use their connection for statistical analysis

## Our Plan

We will cover the following topics

- ▶ Methods for **detecting signal on known networks**
  - ▶ Network analysis based on **centrality and clustering**
  - ▶ **Topology-based pathway enrichment analysis**
- ▶ Methods for **learning undirected networks**
  - ▶ Co-expression networks
  - ▶ ARACNE
  - ▶ Conditional independence graphs
    - ▶ Gaussian observations (glasso, etc)
    - ▶ Non-Gaussian and non-linear data (nonparanormal, etc)
- ▶ [Will not discuss methods for **learning directed networks**]

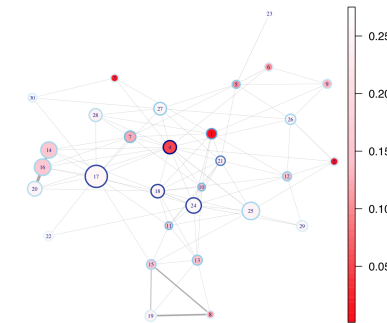
## Network Analysis and Applications in Biology: Analysis of Network-Structured Data

Ali Shojaie & George Michailidis

ENAR 2020

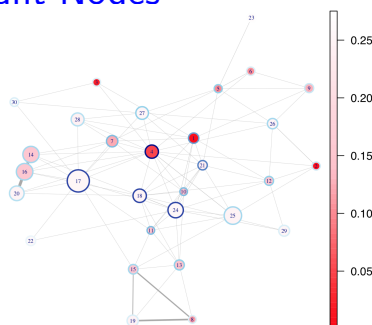
## Introduction

Suppose we observe **activities of individual nodes** (genes, proteins, brain regions, etc) on a **network** (gene regulatory network, structural connectivity network, etc)



How can we identify the **important nodes**?  
*and what does this even mean?*

## Identifying Important Nodes



How can we identify the **important nodes**?

- ▶ We can select the **significant nodes** based on p-values, after adjusting for multiple comparisons (FDR, etc)
- ▶ But the signal is often weak for lots of tests
- ▶ If we believe the network is informative, it may make sense to **use the network to guide our selection**

## Identifying Important Nodes

Possible strategies:

- ▶ Identify **individual nodes** associated with the outcome by incorporating the network (signal detection on network)
- ▶ Test if (pre-specified) **subnetworks** are associated with the outcome (**topology-based pathway enrichment analysis**)
- ▶ Identify **collections of (connected) nodes** that are associated with the outcome (*de-novo identification of enriched modules*)

## Signal Detection on Networks

## Signal Detection on Networks

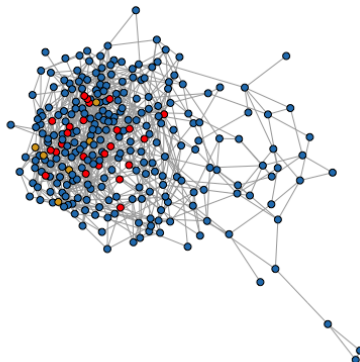
How can we identify the **important nodes in a network**?

The simplest option is to limit our search/testing to the **central nodes** in the network:

- ▶ Nodes connected to many other nodes, aka **hub nodes**
- ▶ Nodes that are **close to many other nodes** (**closeness**)
- ▶ Nodes that are **on many network paths** (**betweenness**)

## Example: Functional Relevance of Hub Nodes

- ▶ Inferred genetic interaction network of cancer-related pathway in prostate cancer (data from TCGA)
- ▶ Hubs defined as nodes whose degrees are at the 75th percentile of the degree distribution



## Other Measures of Centrality

- ▶ **Closeness**: Total distance of each node to other nodes:

$$cl_j = \left( \sum_{k \in V} d(j, k) \right)^{-1}$$

where  $d(j, k)$  is the (shortest path) distance between  $j$  and  $k$ .

- ▶ **Betweenness**: The number of *paths* that go through a node:

$$bw_j = \sum_{i \neq j \neq k} \frac{\pi_{ik}(j)}{\pi_{ik}}$$

where  $\pi_{ik}(j)$  is the number of paths between  $i$  and  $k$  that go through  $j$ , and  $\pi_{ik}$  is the total number of paths between them.

## Identifying “Central” Nodes

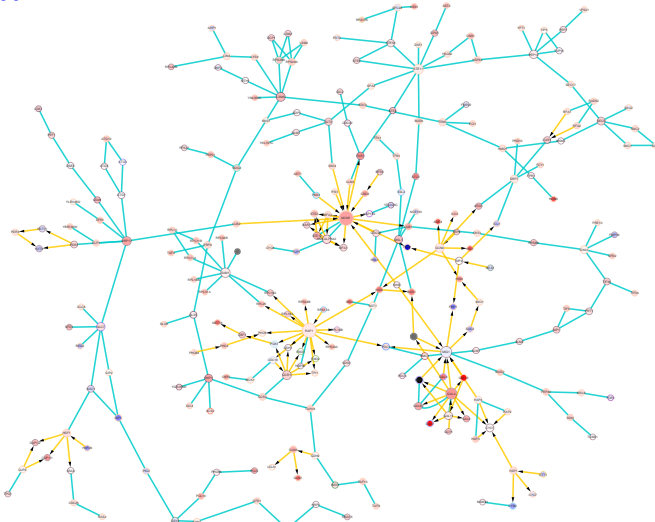
Calculating centrality measures using igraph:

- ▶ Hub nodes: `hub_score(graph)`
- ▶ Closeness: `closeness(graph, vids)`
  - ▶ use `estimate_closeness()` for larger networks)
- ▶ Betweenness: `betweenness(graph, vids)`
  - ▶ use `estimate_betweenness()` for larger networks)

## Topology-Based Pathway Enrichment Analysis

## Yeast GAL Pathway

Ideker et al, 2001



## Topology-Based Pathway Enrichment Analysis

Test for **changes in activities of node** (genes, brain ROIs, etc) in **pre-specified subnetworks**, while **incorporating network information**

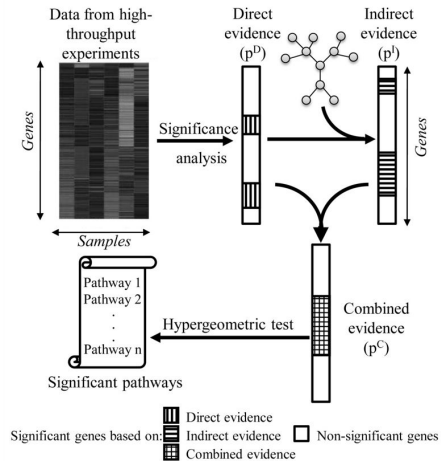
Two possible null hypotheses:

- ▶ **Competitive** null hypothesis: activity of each pathway is compared with other pathways, often using a **permutation test**
  - ▶ Assume few genes are differentially connected, and may be sensitive to the choice of gene sets
- ▶ **Self-contained** null hypothesis: activity of each pathway is compared **against the null** distribution
  - ▶ More rigorous, but may be sensitive to modeling assumptions (*Goemen & Buhlmann (07)*, *Ackermann & Strimmer (09)*)



## PathNet<sup>1</sup>

A simple topology-based pathway enrichment method:



## PathNet: Details

- ▶ Each gene's  $p$ -value from differential expression is combined with  $p$ -values of its neighbors using **Fisher's methods**

$$SI_j = \sum_{k \in \text{ne}(j)} \left\{ -\log_{10} \left( p_k^D \right) \right\}.$$

- ▶ The indirect  $p$ -value,  $p^I$  is calculated from  $SI_j$  by permutation
- ▶ Direct ( $p_j^D$ ) and indirect ( $p_j^I$ )  $p$ -values are then combined ( $p_j^C$ )
- ▶ The significance of  $p_j^C$  for genes in each pathway is assessed using a **hypergeometric test**
- ▶ Implemented in Bioconductor package PathNet

## topologyGSA<sup>2</sup>

- ▶ topologyGSA (Gene Set Analysis Exploiting Pathway Topology) assumes that data are normally distributed:

$$X^1 \sim N(\mu^1, \Sigma^1), \quad X^2 \sim N(\mu^2, \Sigma^2)$$

- ▶ It obtains estimates of  $\Sigma^1$  and  $\Sigma^2$  based on the networks (think graphical lasso, but with **known nonzero entries**)
- ▶ It then **performs two tests**:
  - ▶ equality of covariance matrices:  $H_0^c : \Sigma^1 = \Sigma^2$
  - ▶ equality of means  $H_0^m : \mu^1 = \mu^2$  — it uses different methods depending on the result of  $H_0^c$
- ▶ Implemented in R-package topologyGSA (also in graphite)

<sup>2</sup>Massa et al (2010)

## Signaling Pathway Impact Analysis (SPIA)<sup>3</sup>

- ▶ Combines overrepresentation analysis (ORA) with measure of perturbation of a given pathway under a given condition
- ▶ A bootstrap procedure is used to assess the significance of the observed pathway perturbation (difficult to extend to comparison of  $> 2$  conditions)
- ▶ Currently not applicable to all pathways (more later)
- ▶ Analyzes each pathway separately (ignores connections between pathways)
- ▶ Implemented in the Bioconductor package SPIA

<sup>3</sup>Tarca et al (2009)

## The SPIA Methodology

SPIA combines two types of evidence

- (i) the **overrepresentation** of DE genes in a given pathway
    - ▶ measured by the p-value for the given number of DE genes
- $$P_{NDE} = P(X \geq N_{DE} | H_0)$$

## The SPIA Methodology

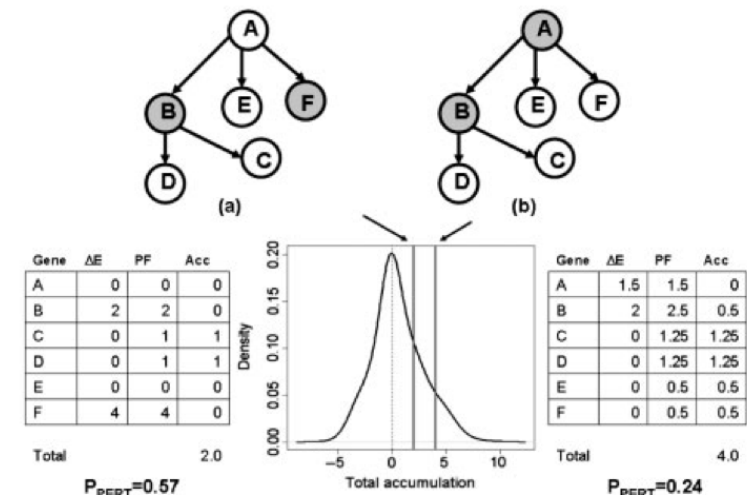
SPIA combines two types of evidence

- (ii) the **abnormal perturbation of the pathway**
  - ▶ the **perturbation for each gene** in the pathway is defined as
 
$$PF(g_i) = \Delta E(g_i) + \sum_{j=1}^p \beta_{ij} \frac{PF(g_j)}{N_{DS}(g_j)}$$
    - ▶  $PF(g_i)$  is the **perturbation factor** of gene  $i$  (not known)
    - ▶  $\beta_{ij}$  is the **magnitude of effect** of gene  $j$  on gene  $i$ ; currently,  $\beta_{ij} = 1$  if  $j \rightarrow i$
    - ▶  $\Delta E(g_i)$  is the **fold change** in expression of gene  $i$
    - ▶  $N_{DS}(g_j)$  is the **number of downstream** genes from gene  $j$

## The SPIA Methodology

- ▶ The **accumulated activity of each gene** can then be calculated as  $ACC(g_i) = B \cdot (I - B)^{-1} \Delta E$ 
  - ▶  $B$  is the **normalized matrix of  $\beta$ 's**:  $B_{ij} = \beta_{ij} / N_{DS}(g_j)$
  - ▶  $\Delta E$  is the **vector of fold changes**
  - ▶ **Requires  $B$  to be invertible**; would not work otherwise
- ▶ The **total accumulated perturbation of the pathway** is then given by  $t_A = \sum_i ACC(g_i)$
- ▶ The **p-value** for pathway perturbation is given by  $P_{PERT} = P(T_A \geq t_A | H_0)$ , which is calculated using a bootstrap approach

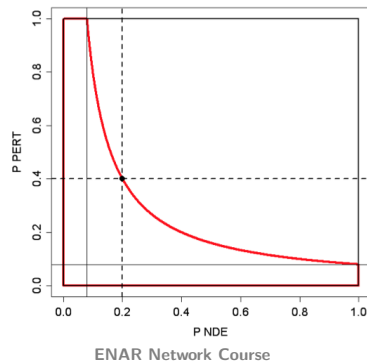
## The SPIA Methodology



## The SPIA Methodology

SPIA combines two types of evidence

- ▶ The **final p-value for each pathway** is calculated based on the p-values from parts (i) and (ii):
  - ▶  $P_G(i) = c_i - c_i \ln(c_i)$
  - ▶  $c_i = P_{NDE}(i)P_{PERT}(i)$



## An Example in R: Data on Colorectal Cancer

```
data(colorectalcancer)

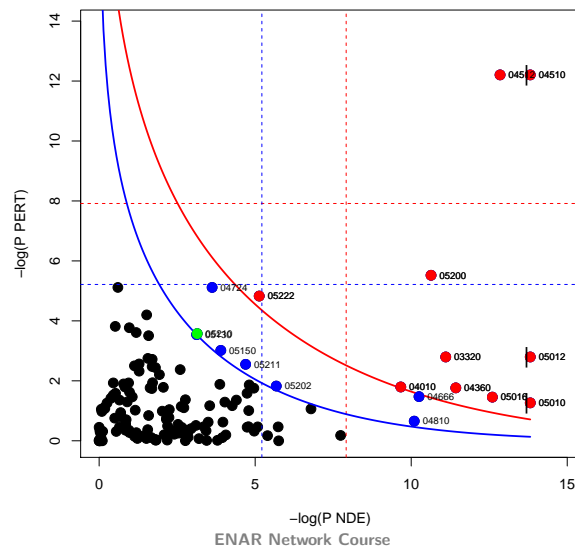
#pathway analysis using SPIA
#use nB=2000 or higher for more accurate results
#uses older version of KEGG signaling pathways graphs
res <- spia(de=DE_Colorectal, all=ALL_Colorectal, organism="hsa", beta=NULL,
            nB=2000, plots=FALSE, verbose=TRUE, combine="fisher")

#now combine pNDE and pPERT using the normal inversion method without
#running spia function again
res$pG=combfunc(res$pNDE,res$pPERT,combine="norminv")
res$pGFdr=p.adjust(res$pG,"fdr")
res$pGFWER=p.adjust(res$pG,"bonferroni")
plotP(res,threshold=0.05)

#highlight the colorectal cancer pathway in green
points(I(-log(pPERT))~I(-log(pNDE)),data=res[res$ID=="05210",],col="green",
       pch=19,cex=1.5)
```

## The SPIA Methodology

SPIA two-way evidence plot



## Network-Based Gene Set Analysis (NetGSA)<sup>4</sup>

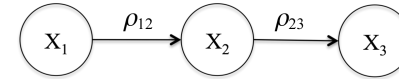
- ▶ Generalizes SPIA, to allow for more complex experiments & incorporate interactions among pathways
- ▶ Assesses the overall behavior of arbitrary subnetworks (pathways): **changes in gene expression & network structure**
- ▶ Uses **latent variables** to model the interaction between genes defined by the network
- ▶ Uses **mixed linear models** for inference in complex data
- ▶ Computationally challenging for large networks, unless pathways separately analyzed (similar to SPIA)

<sup>4</sup>S & M (2009, 2010); Ma, S & M (2016)

## Problem Setup

- ▶ Gene (protein/metabolite) expression data for  $K$  experimental conditions and  $J_k$  time points
- ▶ Network information (partially) available in the form of a **directed weighted graph**  $G = (V, E)$ , with vertex set  $V$  corresponding to the genes/proteins/metabolites and edge set  $E$  capturing their associations
- ▶ Network edges can be **directed**  $j \rightarrow k$  or **undirected**  $j \leftrightarrow k$
- ▶ Edges defines the **effect** of nodes on their immediate neighbors; the weight associated with each edge corresponds to the value of **partial correlation**
- ▶ Represent the network by its **adjacency matrix**  $A$ :  $A_{jk} \neq 0$  iff  $k \rightarrow j$  & for undirected edges,  $A_{jk} = A_{kj}$

## The Latent Variable Model: Main Idea



$$X_1 = \gamma_1$$

$$X_2 = \rho_{12}X_1 + \gamma_2 = \rho_{12}\gamma_1 + \gamma_2$$

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

Thus  $X = \Lambda\gamma$  where

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

## The Latent Variable Model

- ▶ Let  $Y$  be the  $i$ th sample in the expression data
- ▶ Let  $Y = X + \varepsilon$ , with **signal**  $X$  and **noise**  $\varepsilon \sim N_p(0, \sigma_\varepsilon^2 I_p)$
- ▶ The **influence matrix**  $\Lambda$  measures the **propagated effect of genes on each other** through the network, and can be calculated based on the adjacency matrix  $A$
- ▶ Using  $X = \Lambda\gamma$ , we get

$$Y = \Lambda\gamma + \varepsilon, \Rightarrow Y \sim N_p(\Lambda\mu, \sigma_\gamma^2 \Lambda\Lambda' + \sigma_\varepsilon^2 I_p)$$

where  $\gamma \sim N_p(\mu, \sigma_\gamma^2 I_p)$  are **latent variables**

## Mixed Linear Model Representation

Rearranging the expression matrix into  $np$ -vector  $\mathbf{Y}$ , we can write

$$\mathbf{Y} = \Psi\beta + \Pi\gamma + \varepsilon$$

where  $\beta$  and  $\gamma$  are fixed and random effect parameters and

$$\varepsilon \sim N_{np}(\mathbf{0}, R(\theta_\varepsilon)), \quad \gamma \sim N_{np}(\mathbf{0}, \sigma_\gamma^2 \mathbf{I}_{np})$$

- **Temporal Correlation** incorporated through  $R$

In general, the **design matrices**,  $\Psi$  and  $\Pi$  depend on the experimental settings (similar to ANOVA), and are **functions of  $\Lambda$**

## Estimation of MLM Parameters

MLE for  $\beta$ :

$$\hat{\beta} = (\Psi' \hat{W}^{-1} \Psi)^{-1} \Psi' \hat{W}^{-1} \mathbf{Y}$$

where  $W = \sigma_\gamma^2 \Pi \Pi' + R$ .

$\hat{\beta}$  depends on estimates of  $\sigma_\gamma^2$  and  $\theta_\varepsilon^2$  (estimated using **restricted maximum likelihood (REML)**).

## Inference using MLM

- ▶ Let  $\ell$  be a **contrast vector** (a linear combination of fixed effects), and consider the test:

$$H_0 : \ell\beta = 0 \quad \text{vs.} \quad H_1 : \ell\beta \neq 0$$

- ▶ Use t-test to test the significance of each hypothesis separately

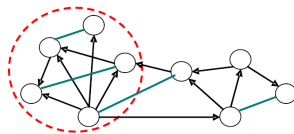
$$T = \frac{\ell\hat{\beta}}{\sqrt{\ell\hat{C}\ell'}}$$

where  $C = (\Psi' W^{-1} \Psi)^{-1}$

- ▶ Under the null hypothesis,  $T$  is approximately  $t$ -distributed with degrees of freedom that needs to be estimated

## “Optimal” Choice of Contrast Vector

- ▶ An intuitive choice is the **indicator (membership) vector** for the pathway,  $\mathbf{b}$ , but this only captures changes in mean
- ▶ Need to **de-couple the effect of subnetwork** from other nodes



## “Optimal” Choice of Contrast Vector

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

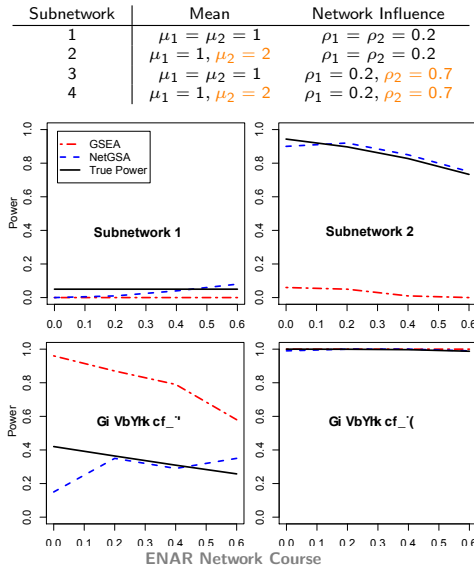
Consider the set,  $\mathbf{b} = (0, 1, 1)$ ; then

$$(\mathbf{b}\Lambda) = (\rho_{12} + \rho_{12}\rho_{23}, 1 + \rho_{23}, 1)$$

On the other hand,

$$(\mathbf{b}\Lambda \cdot \mathbf{b}) = (0, 1 + \rho_{23}, 1)$$

## Comparison in Simulated Data



33

## Yeast Galactose Utilization Pathway

Ideker et al (2001) data on yeast Galactose Utilization Pathway

- ▶ Gene expression data for 2 experimental conditions: (gal+) and (gal-)
- ▶ Gene-gene and protein-gene interactions as well as association weights found from previous studies
- ▶ Q: which pathways respond to the change in growth medium?

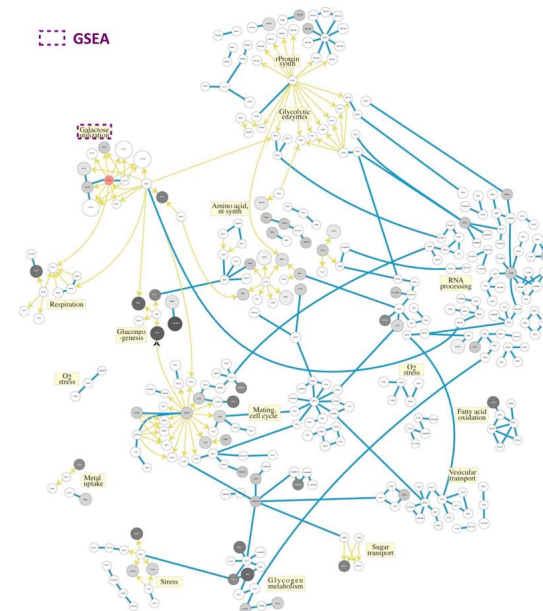
©Ali Shojaie

ENAR Network Course

34

## Analysis of Yeast GAL Data

- ▶ Data:
  - ▶ gene expression data for 343 genes
  - ▶ 419 interactions found from previous studies and integration with protein expression (association among genes also available)
- ▶ Results:
  - ▶ GSEA finds Galactose Utilization Pathway significant
  - ▶ NetGSA finds several other pathways with biologically meaningful functions related to survival of yeast cells in gal-

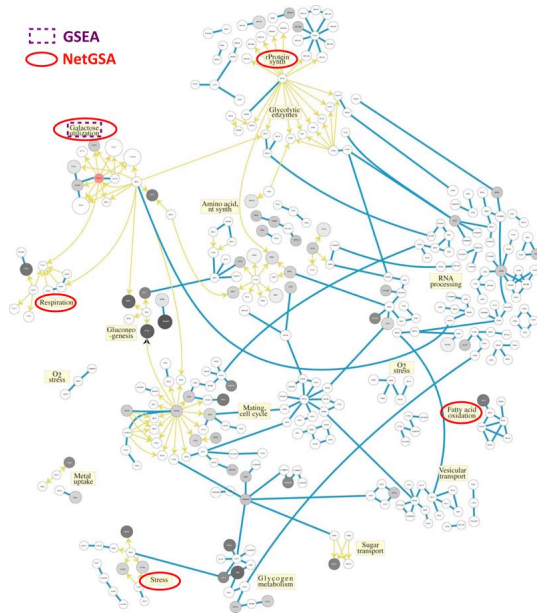


35

©Ali Shojaie

ENAR Network Course

36



## Environmental Stress Response in Yeast

Gene expression data on Yeast Environmental Stress Response (ESR) (*Gasch et al., 2000*)

- ▶ 3 combinations of experimental factor, heat shock and osmotic changes (sorbitol), over 3 time points
- ▶ **Temporal correlation**
- ▶ **Network correlation**
- ▶ **Q**: Which **pathways** indicate response to environmental stress
  - ▶ in different **experimental conditions**
  - ▶ over **time**

## Yeast ESR Data

*Gasch et al (2000)*

### ▶ Gene Expression Data

Experiment	Obs. Time (after 33C)
Mild heat shock (29C to 33C), no sorbitol	5, 15, 30 min
Mild Heat Shock, 1M sorbitol at 29C & 33C	5, 15, 30 min
Mild Heat Shock, 1M sorbitol at 29C	5, 15, 30 min

### ▶ Network Data

- ▶ Use **YeastNet** (*Lee et al., 2007*) for gene-gene interactions (102,000 interactions among 5,900 yeast genes)
- ▶ Use independent experiments of *Gasch et al.* to **estimate weights**
- ▶ Pathways are defined using **GO** functions

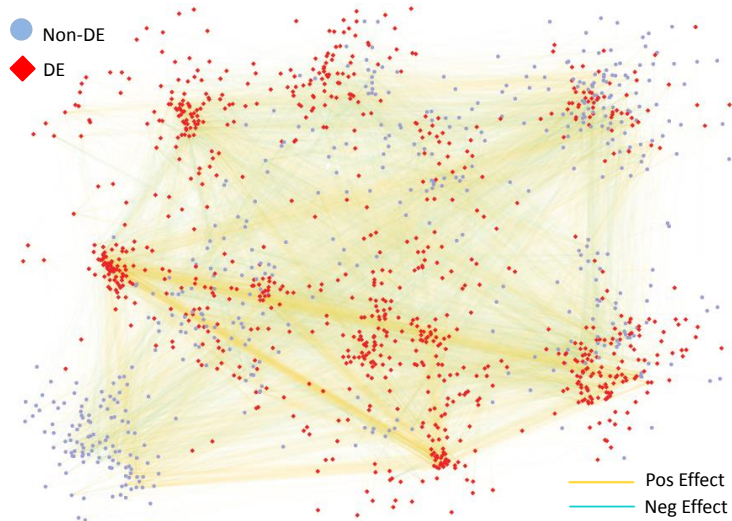
## Model and Results

- ▶ Model: Let  $j$  and  $k$  be indices for **time** and levels of **sorbitol**

$$\mathbb{E}Y_{11} = \Lambda\mu, \quad \mathbb{E}Y_{jk} = \Lambda(\mu + \alpha_j + \delta_k) \quad j, k = 2, 3$$

- ▶ **Temporal correlation** is modeled directly via  $R$  (as  $AR(1)$  process)
- ▶ Results:
  - ▶  $\sim 3000$  genes,
  - ▶ 47 pathways showed significant changes of expression
  - ▶ 24 pathways showed changes over **time**
  - ▶ 29 pathways showed changes in response to different **sorbitol** levels
  - ▶ 12 pathways showed **both** types of changes
  - ▶ Significant pathways overlap with the gene functions recognized by *Gasch et al.*

## Yeast ESR Network

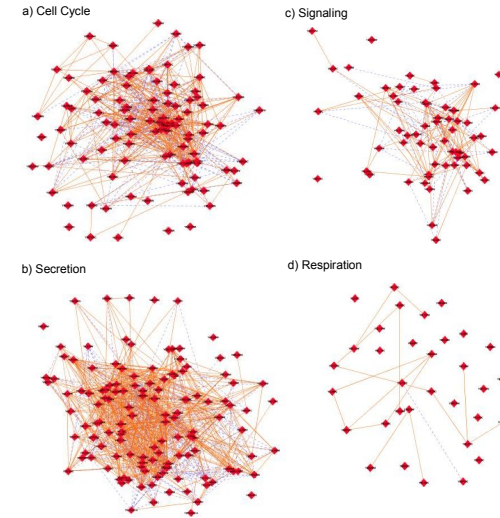


©Ali Shojaie

ENAR Network Course

41

## Significant subnetworks



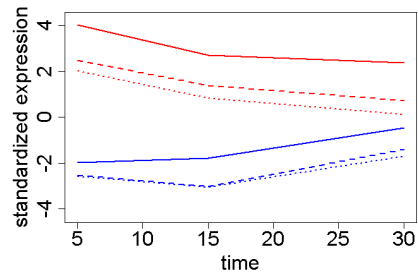
©Ali Shojaie

ENAR Network Course

42

## Expression Profiles

Average Standardized Expression Levels of Pathways



- ▶ Induced and Suppressed Pathways
- ▶ Can observe the transient patterns of expressions as predicted by Gasch et al.

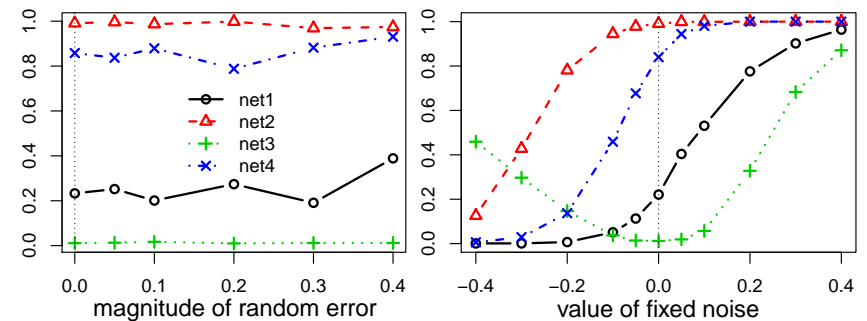
©Ali Shojaie

ENAR Network Course

43

## Effect of Noise In Network Information

- ▶ Let  $\tilde{A}$  be observed network information, and  $A$  be the truth.
- ▶ It can be shown that, if  $\|\tilde{A} - A\|$  is small then, NetGSA still works (is asymptotically most powerful unbiased test)



©Ali Shojaie

ENAR Network Course

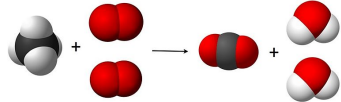
44



## Metabolic Profiling in Bladder Cancer

Targeted metabolic profiling of bladder cancer (BCa) (*Putluri et al.*, 2012)

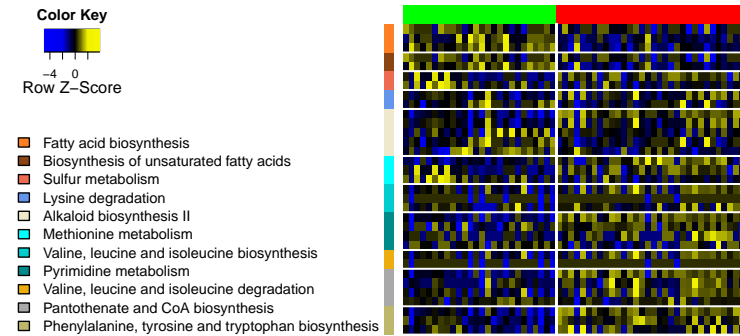
- ▶ 58 bladder cancer and adjacent benign samples
- ▶ Pathways information obtained from **KEGG**



- ▶ Varying number of identified metabolites per pathway (3-15)
- ▶ **Q**: Which **pathways** show differential activity in BCa?

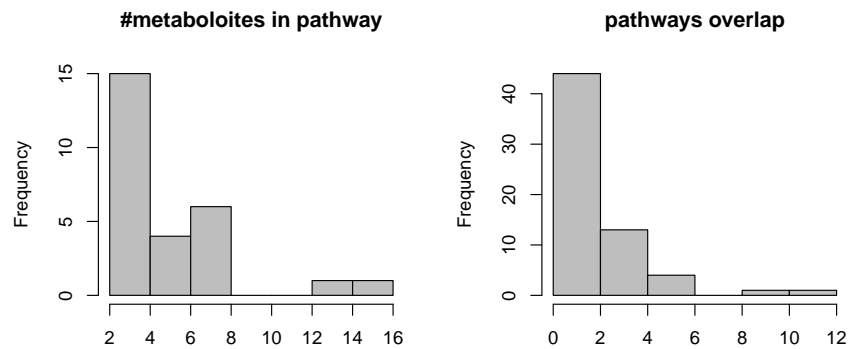
## Metabolic Profiling in BCa

- ▶ 63 metabolites identified, mapped to 70 pathways
- ▶ 27 pathways with at least 3 members



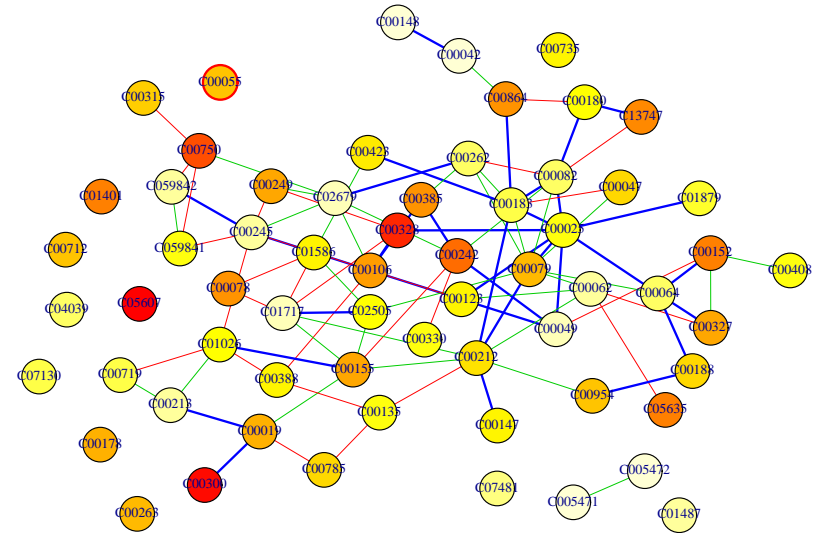
## Metabolic Profiling in BCa

- ▶ Small pathway sizes & significant overlap among pathways



- ▶ Existing methods may not work well

## Metabolic Interaction Network

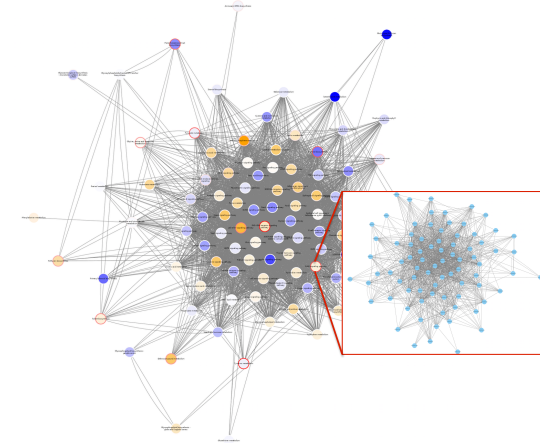


## Significant Pathways

- ▶ **GSEA** does not identify any pathway as differential
- ▶ **GSA** identifies **Fatty Acid Biosynthesis** as differential
- ▶ **NetGSA** identifies another 7 pathways corresponding to role of **Amino Acid Metabolism** in BCa, similar to *Putluri et al* (2012)

## R-Package netgsa

```
adjmats <- prepareAdjMat(data, groups, edges, TRUE)
res <- NetGSA(adjmats$Adj, data, groups, pathways, "REHE")
plot(res) #interactive plotting in Cytoscape
```



## Selected Topology Based Pathway Enrichment Methods

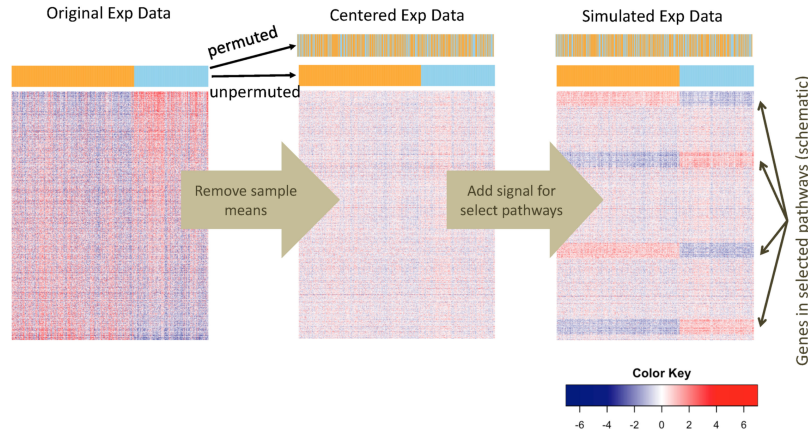
Method	Null hypothesis	Input
Pathway-Express	competitive	DE genes & $p$ -values; sample labels; pathway topology
NetGSA	self-contained	expression matrix; sample labels; pathway membership; network information
SPIA	competitive	DE genes with $p$ -values; sample labels; pathway topology
topologyGSA	self-contained	Gene expression matrix; sample labels; pathway topology
CAMERA	competitiv	Gene expression matrix; sample labels; pathway membership
DEGraph	self-contained	Gene expression matrix; sample labels; pathway topology
PathNet	competitive	DE genes with $p$ -values; sample labels; pathway topology

Overview of tested pathway enrichment methods. All methods return the  $p$ -values before and/or after correcting for multiple comparisons.

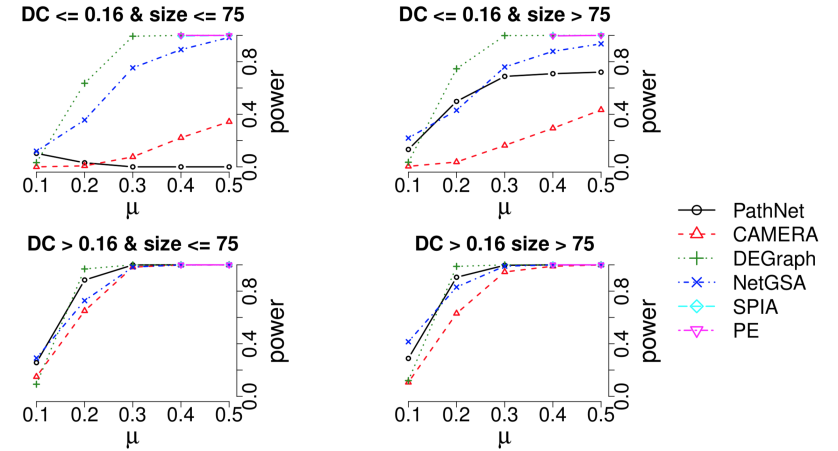
## Comparison of these Methods Using Synthetic Data (Ma, Shojaie, Michailidis, 2019)

- ▶ Comparison of topology-based pathway enrichment methods using two synthetic data sets
  - ▶ Gene expression data  $p \approx 3000$
  - ▶ Metabolomics data  $p \approx 100$
- ▶ *In silico* data sets with known signal:
  1. Remove the original signal, but **keep the correlation structure**
  2. **Perturb means in one condition** (differential expression) for nodes in selected pathways
  3. Also use sample permutation to **create data with equal correlation structure**

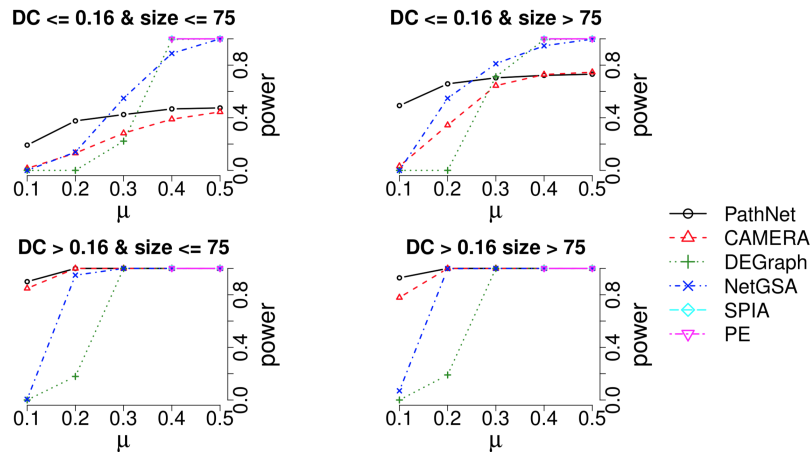
## Comparison Using Synthetic Data



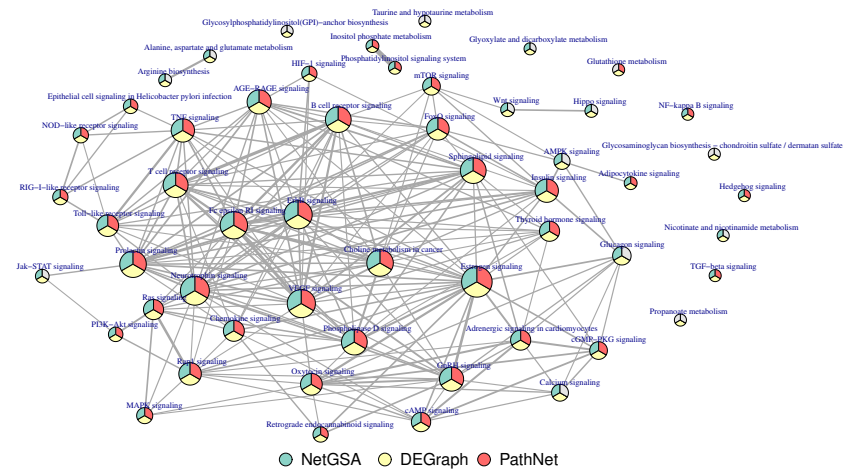
## Results for Gene Expression Data — Equal Covariance



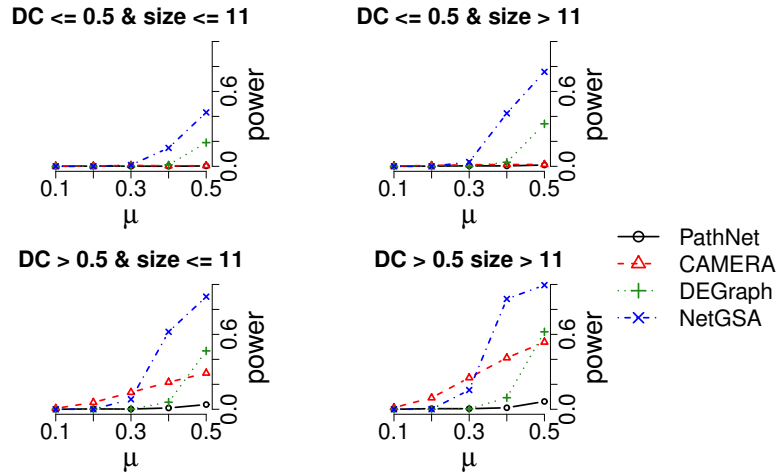
## Results for Gene Expression Data — Diff Covariance



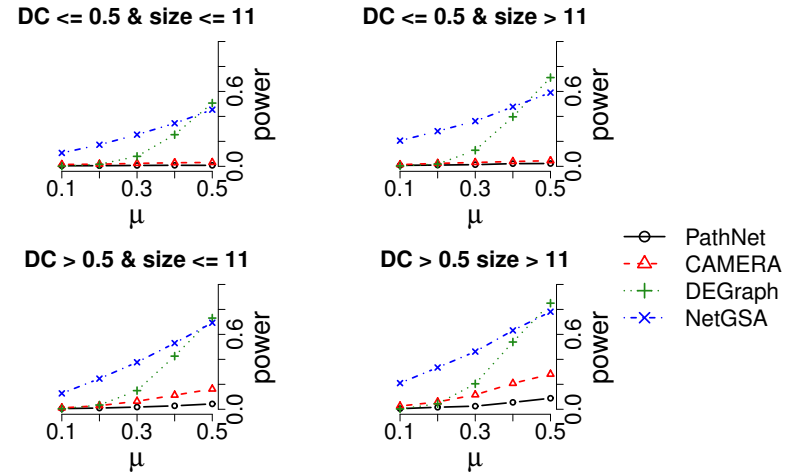
## Results for Gene Expression Data



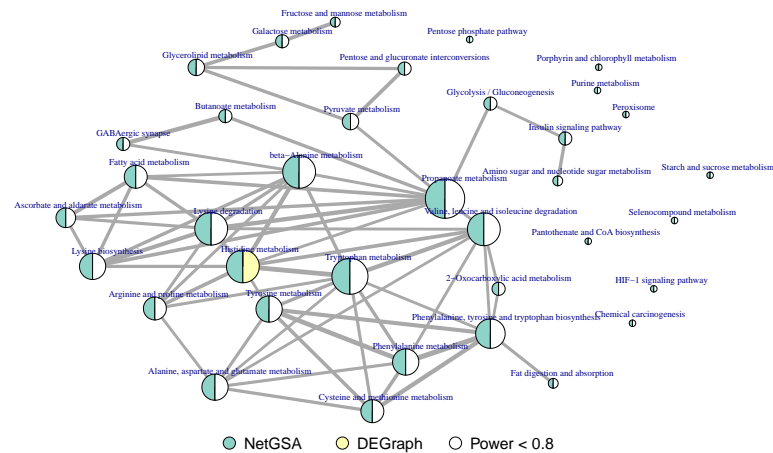
## Results for Metabolomics Data — Equal Covariance



## Results for Metabolomics Data — Diff Covariance

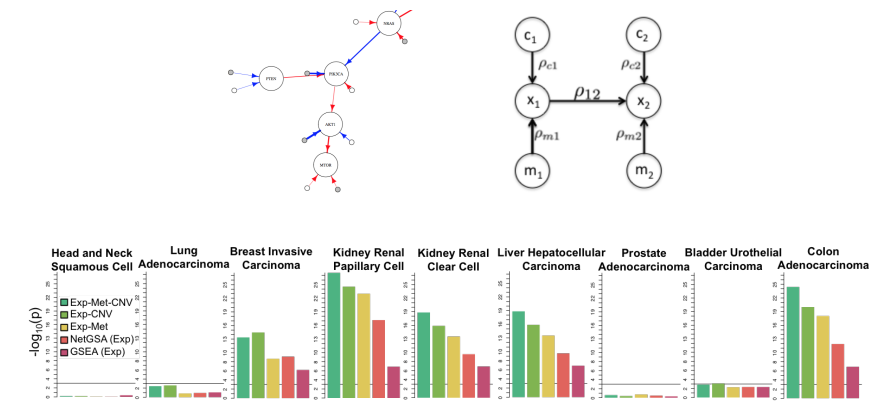


## Results for Metabolomics Data



## Multi-Omics NetGSA

Pan-cancer integration of expression, methylation and CNV in BRAF (TCGA data)<sup>5</sup>



<sup>5</sup>Zhang et al (2018)

## Identifying Enriched Modules in Networks

## Identifying Enriched Modules in Networks

Two general strategies:

- ▶ Assess the significance of **data-driven modules** (WGCNA):
  1. Identify modules (network clustering, etc)
  2. Assess the significance of modules
- ▶ **Search** for enriched (connected) subnetworks (often using greedy search methods)
- ▶ Advantage: No need to rely on known pathways — especially useful when known pathways are not complete, etc
- ▶ Disadvantage: Interpretation may become challenging...

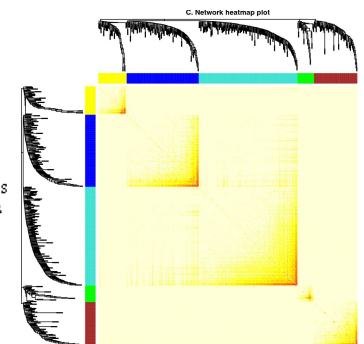
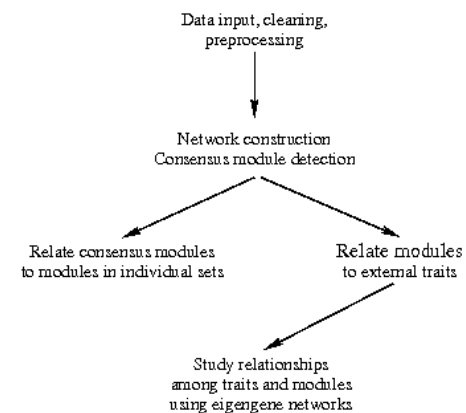
## WGCNA<sup>6</sup>

- ▶ We previously talked about weighted gene co-expression (WGCNA), but for **estimating** networks
- ▶ However, WGCNA is also used for topology-based enrichment analysis, although in a different way than many other topology-based methods
- ▶ Here's how it works:
  1. Estimate the **co-expression network** (more in the next lecture)
  2. Find **modules** by **clustering** the nodes in the estimated network
  3. Summarize the expressions of genes in each module using PCA (**eigen-genes**)
  4. Test if the eigen-genes are associated with the outcome

<sup>6</sup>Horvath & Zhang (2005); Langfelder et al (2008)

## WGCNA

- ▶ Here's how it works:



Let's look at an example in R...

## Walktrap<sup>7</sup>

- ▶ Searches for connected modules containing significant genes
  - ▶ Weights each edges based on the **significance of its corresponding nodes**

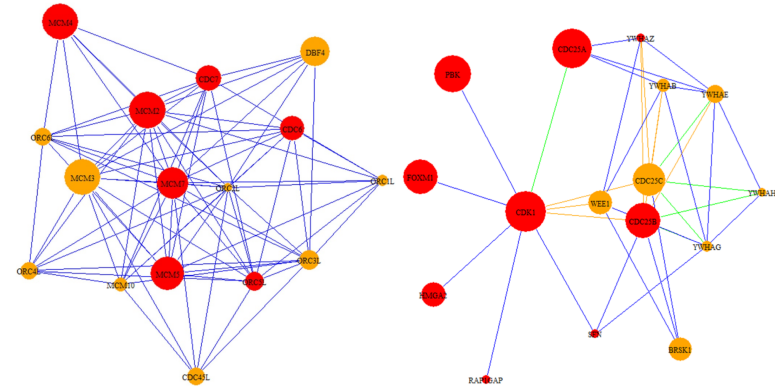
$$w_{ij} = (|FC_i| + |FC_j|)/2$$

- ▶ Connected significant modules are found through **community detection** using a **random walk** with transition probability

$$P_{ij} = \frac{w_{ij}}{\sum_j w_{ij}}$$

<sup>7</sup>Petrochilos et al (2013)

## Identifying Cancer-Related Modules



## Summary

- ▶ Network-based methods (centrality-based, pathway topology, etc) rely on network information — helpful if correct network information avail
- ▶ What if network information is not available?
- ▶ What about differences in network structures — **differential network biology**<sup>8</sup>?

<sup>8</sup>Ideker & Krogan (2012)

## Network Analysis and Applications in Biology: Learning Undirected Networks

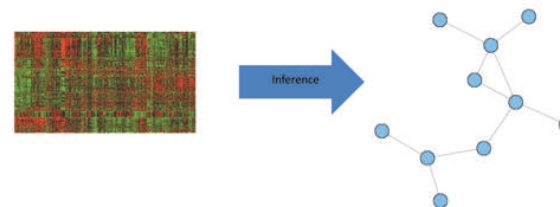
Ali Shojaie & George Michailidis

ENAR 2020

## Learning Undirected Networks

Learn network from data (**structure learning**):

- ▶ Data matrix:  $X_{n \times p}$ .
- ▶ Features correspond to the  $p$  nodes in the network.
- ▶ Goal: Learn edges between nodes  $\equiv$  learn the **statistical relationships** between features.



## Why Do We Need Network Inference?

- ▶ Despite progress, our knowledge of interactions is limited.
- ▶ The entire genome is a vast landscape, and **experiments for discovering networks are very expensive**.
- ▶ From a statistical point of view, **network estimation is related to estimation of covariance matrices**, which has many independent applications in statistical inference and prediction (*more about this later*).
- ▶ Finally, and perhaps most importantly, **gene and protein networks are dynamic** and changes in these networks have been attributed to complex diseases.

## Network Inference — An Overview

Two general classes of network inference methods:

- ▶ Methods based on **marginal measures of association**:
  - ▶ Co-expression Networks (based on linear measures of association)
  - ▶ Methods based on **mutual information** (can accommodate non-linear associations)
- ▶ Methods based on **conditional measures of association**:
  - ▶ Methods assuming (multivariate) normality (glasso, etc)
  - ▶ Generalizations to allow for nonlinear dependencies (nonparanormal, etc)

## Graphical Models

### Probabilistic Graphical Models<sup>1</sup>

Joint multivariate probability distribution where dependencies can be represented as a network.

Advantages:

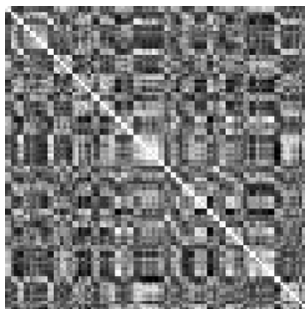
- ▶ Graphical models offer efficient factorized forms for joint distributions with easily interpretable dependencies.
  - ▶ **Conditional dependencies** denoted via an edge in network.
- ▶ Convenient visual representation.

<sup>1</sup>For a detailed technical introduction, see *Graphical Models, Exponential Families, and Variational Inference* by Wainwright & Jordan (2008)

## Marginal Association Networks

## Correlation Networks (Association Networks)

- ▶ Simplest (and most-widely used!) method for estimating networks — key assumption: large correlation  $\equiv$  presence of an edge
- ▶ Let  $r(i, j)$  be correlation between  $X_i$  and  $X_j$ ; we claim an **edge between  $i$  and  $j$  if  $|r(i, j)| > \tau$** .
  - ▶  $\tau$ : a user-specified threshold (**tuning parameter**).



Correlation matrix



Thresholded correlation matrix

## Limitations of Correlation Networks

1. The estimation is highly dependent on the **choice of  $\tau$** .
2. Correlations capture **linear associations**, but **many real-world relationships are nonlinear**.
3. Large correlations can occur due to **confounding**.



## Limitations of Correlation Networks

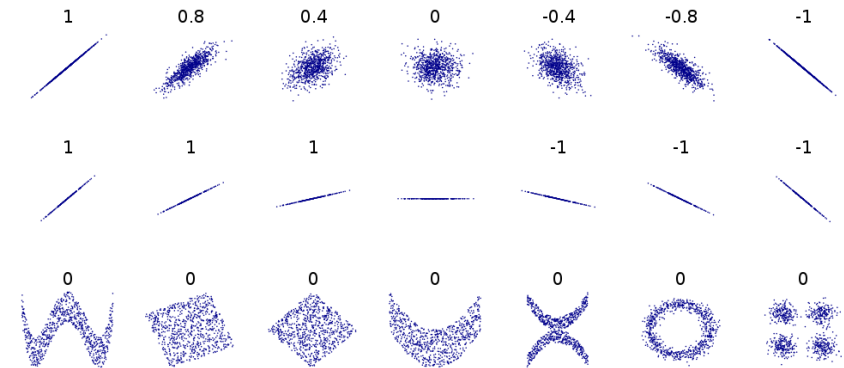
The estimation is highly dependent on the choice of  $\tau$ .

- ▶ We can work with **weighted co-expression networks** (WGCNA)
- ▶ We can instead test  $H_0 : r_{xy} = 0$ 
  - ▶ A commonly used test is based on the **Fisher transformation**

$$Z = \frac{1}{2} \ln \left( \frac{1+r}{1-r} \right) = \text{artanh}(r) \sim_{H_0} N \left( 0, \frac{1}{\sqrt{n-3}} \right)$$

## Limitations of Correlation Networks

Correlations capture **linear** associations, but **many real-world relationships are nonlinear**.



## Limitations of Correlation Networks

Correlations capture **linear** associations, but **many real-world relationships are nonlinear**.

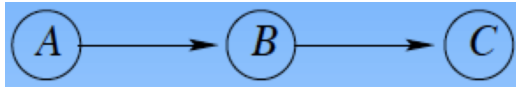
- ▶ We can use other measures of association, for instance, **Spearman correlation** or **Kendal's  $\tau$** .
  - ▶ These methods define the correlation between two variables, based on the **ranking** of observations, and not their exact values.
  - ▶ They can better capture non-linear associations.
- ▶ We can instead use **mutual information**; this has been used in many algorithms, e.g. ARACNE.

## ARACNE: Algorithm for the Reconstruction of Accurate Cellular NEtworks<sup>2</sup>

1. Identifies statistically significant gene-gene co-regulation based on mutual information
2. It then eliminates indirect relationships in which two genes are co-regulated through one or more intermediates

<sup>2</sup>Margolin et al (2006)

## Key Idea: Data Processing Inequality (DPI)



$$I(A, C) \leq \min[I(A, B), I(B, C)]$$

where

$$I(g_i, g_j) = \log P(g_i, g_j) / P(g_i)P(g_j)$$

- ▶ Look at every triplet and remove the weakest link
- ▶ Need to estimate marginal and joint (pairwise) probabilities (using Gaussian Kernel)

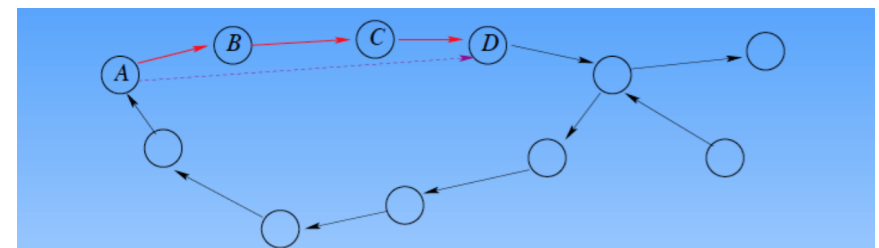
## Algorithm Details

- ▶ The algorithm examines each gene triplet for which all pairwise MIs are greater than a cut-off and removes the edge with the smallest value based on DPI.
  - ▶ Each triplet is analyzed even if its edges have been selected for removal by prior DPI applications to other triplets.
  - ▶ The least of the three MIs can come from indirect interactions only, and checking against the DPI may identify **gene pairs that are not independent, but still do not interact**.

## Rationale and Guarantees

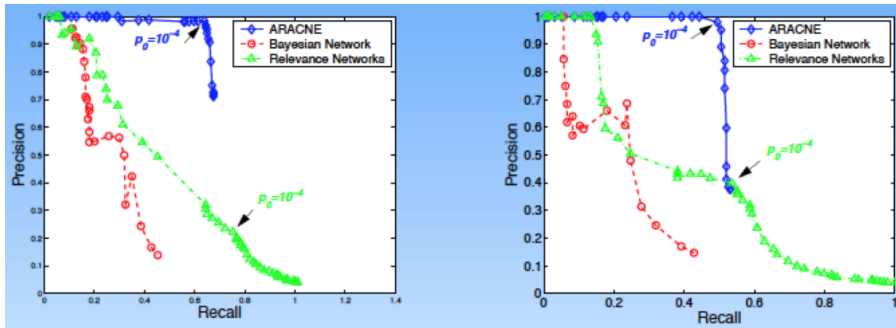
- ▶ If MIs are estimated with no errors, then ARACNE reconstructs the underlying interaction network exactly, if the network is a tree and has only pairwise interactions.
- ▶ The maximum MI spanning tree is a subnetwork of the network built by ARACNE.

## Rationale and Guarantees

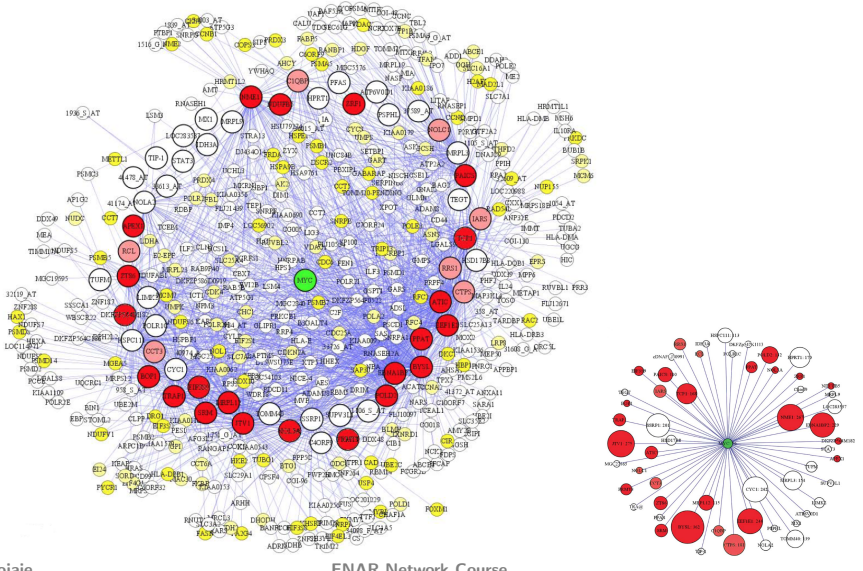


**Theorem.** Let  $\pi_{ik}$  be the set of nodes forming the shortest path in the network between nodes  $i$  and  $k$ . Then, if MIs can be estimated without errors, ARACNE reconstructs an interaction network without false positives edges, provided: (a) the network consists only of pairwise interactions, (b) for each  $j \in \pi_{ik}$ ,  $I_{ij} \geq I_{ik}$ .  
 Further, ARACNE does not produce any false negatives, and the network reconstruction is exact iff (c) for each directly connected pair  $ij$  and for any other node  $k$ , we have  $I_{ij} > \min[I_{ik}, I_{jk}]$ .

## Performance on Synthetic Data



## Application: B-lymphocytes Expression Data



## Application: B-lymphocytes Expression Data

- ▶ MYC (proto-oncogene) subnetwork (2063 genes)
- ▶ 29 of the 56 (51.8%) predicted first neighbors biochemically validated as targets of the MYC transcription factor.
- ▶ New candidate targets were identified, 12 experimentally validated.
  - ▶ 11 proved to be true targets.
- ▶ The candidate targets that have not been validated are possibly also correct.

## Software

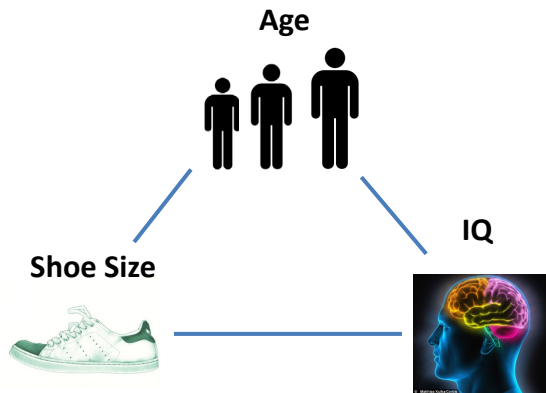
- ▶ Implemented in the R-package minet:
 

```
source("http://bioconductor.org/biocLite.R")
biocLite("minet")
```
- ▶ Main estimation function `aracne(mim, eps=0)`
  - ▶ `mim`: mutual information matrix
 

```
mim <- build.mim(syn.data, estimator="spearman")
```
  - ▶ `eps`: threshold for setting an edge to zero, prior to searching over triplets

## Limitations of Correlation Networks

Large correlations can occur due to **confounding**.

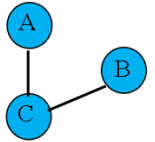
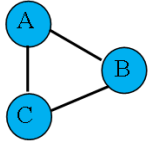


## Markov Networks

Markov Network

An *undirected graphical model* that characterizes **conditional dependence** ( $\equiv$  direct relationships).

- ▶ *Edge*: Two nodes are **conditionally dependent**.
- ▶ *No edge*: Two nodes are **conditionally independent**.
- ▶ Conditions on all other nodes.



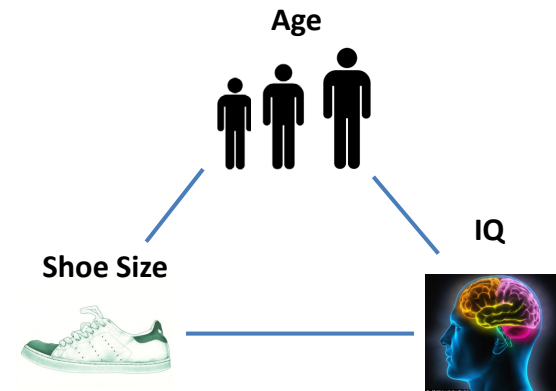
$$A \perp B \mid C$$

## Markov Networks — Conditional Dependence

Regression Interpretation:

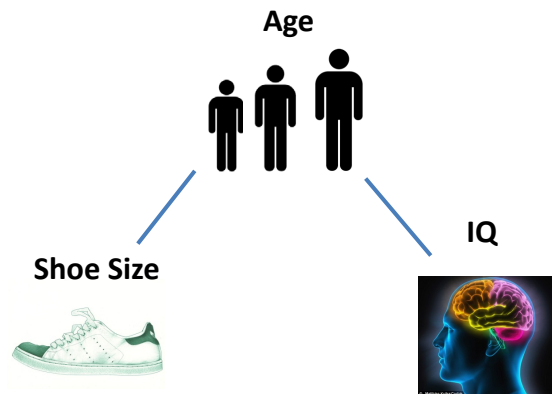
- ▶ Imagine trying to predict the observations in **Node A** (response) by the observations of all other nodes (predictors).
- ▶ **Node B** predictive of **Node A** (with all other nodes in model).
  - ▶ **A** is conditionally dependent on **B**.
  - ▶ Edge.
- ▶ Because of other nodes in model, **Node B** does not add any predictive value for **Node A**.
  - ▶ **A** is conditionally independent of **B**.
  - ▶ No Edge.

## Markov Networks — Conditional Dependence



Correlation.

## Markov Networks — Conditional Dependence



Conditional Dependence.

## Markov Networks — Conditional Dependence

How can we learn conditional dependencies?

- ▶  $A$  and  $B$  are conditionally independent given  $C$  if

$$P(A, B | C) = P(A | C)P(B | C)$$

- ▶ Generally difficult (need to estimate multivariate densities).
- ▶ Alternatively, can use nonparametric approaches, e.g. **conditional mutual information**, but not easy in high dimensions.
- ▶ Often resort to models, or simple measures, such as **partial correlations**...

## Partial Correlation

- ▶ Partial correlation measures the **correlation between  $A$  and  $B$  after the effect of the other variables are removed**.
  - ▶ In our example, this means correlation between shoe size and IQ, **after adjusting for age**.
- ▶ The partial correlation between  $A$  and  $B$  **given  $C$**  is given by:

$$\rho_{AB \cdot C} \equiv \rho(A, B | C) = \frac{\rho_{AB} - \rho_{AC}\rho_{BC}}{\sqrt{1 - \rho_{AC}^2}\sqrt{1 - \rho_{BC}^2}}$$

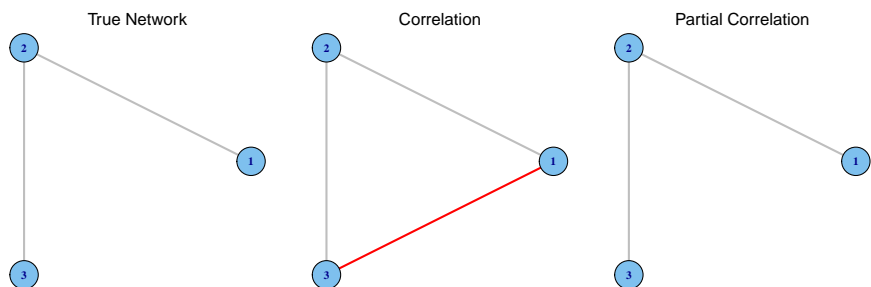
- ▶ Alternatively, **regress**  $A$  on  $C$  and get the residual,  $r_A$ ; do the same for  $B$  to get  $r_B$ . The partial correlation between  $A$  and  $B$  **give  $C$**  is  $\text{Cor}(r_A, r_B)$ .

## Partial Correlation

- ▶ Partial correlation is **symmetric**  $\Rightarrow$  **undirected network**
- ▶ Partial correlation takes values **between -1 and 1**
- ▶ In partial correlation networks, we **draw an edge** between  $A$  and  $B$ , **if the partial correlation between them is large**
- ▶ Calculation of partial correlation is more involved

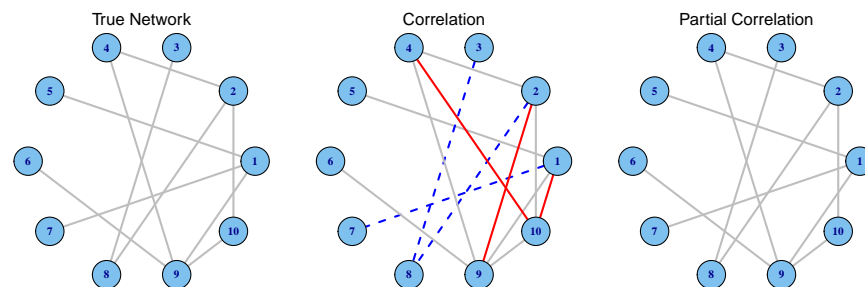
## A Simple Example

$$\text{Correlation} = \begin{bmatrix} 1 & .8 & .7 \\ .8 & 1 & .8 \\ .7 & .8 & 1 \end{bmatrix} \quad \text{PartialCorr} = \begin{bmatrix} 1 & .6 & 0 \\ .6 & 1 & .6 \\ 0 & .6 & 1 \end{bmatrix}$$



## A Larger Example

- ▶ A network with 10 nodes and 20 edges
- ▶  $n = 100$  observations
- ▶ Estimation using correlation & partial correlation (20 edges)

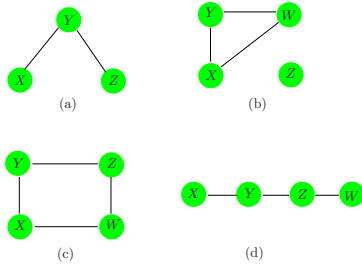


## Gaussian Graphical Models (GGMs)

## Partial Correlation for Gaussian Random Variables

- ▶ For Gaussian (multivariate normal) random variables, partial correlation between  $X_i$  and  $X_j$  given all other variables is given by the inverse of the (standardized) covariance matrix  $\Sigma$ .
  - ▶ The  $(i, j)$  entry in  $\Sigma^{-1}$  gives the partial correlation between  $X_i$  and  $X_j$  given all other variables  $X_{\setminus i, j}$ .
  - ▶ Multivariate normal:  $X \sim N(0, \Sigma)$
  - ▶  $\Theta \equiv \Sigma^{-1}$  = inverse covariance/precision/concentration matrix.
  - ▶ Zeros in  $\Theta \implies$  conditional independence!
  - ▶ Edges correspond to non-zeros in  $\Theta$ .

## Partial Correlation for Gaussian Random Variables



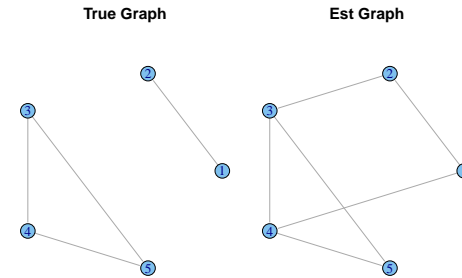
$$\begin{pmatrix} - & \times & 0 \\ \times & - & \times \\ 0 & \times & - \end{pmatrix} \quad \begin{pmatrix} - & \times & \times & 0 \\ \times & - & \times & 0 \\ \times & \times & - & 0 \\ 0 & 0 & 0 & - \end{pmatrix}$$

$$\begin{pmatrix} - & \times & 0 & \times \\ \times & - & \times & 0 \\ 0 & \times & - & \times \\ \times & 0 & \times & - \end{pmatrix} \quad \begin{pmatrix} - & 0 & 0 & \times \\ 0 & - & \times & 0 \\ 0 & \times & - & \times \\ \times & 0 & \times & - \end{pmatrix}$$

## Estimating GGMs

- From the discussion so far, to estimate the network, we can
1. Calculate the **empirical covariance matrix**: for (centered)  $n \times p$  data matrix  $X$ ,  $S = (n - 1)^{-1} X^T X$ .
  2. **Get the inverse of  $S$** . Non-zero values of  $S^{-1}$  give the edges.

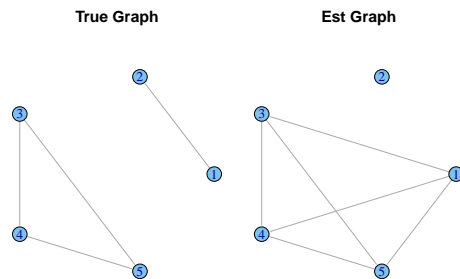
While simple, this may not work well in practice, even with large samples!



## Estimating GGMs in High Dimensions

Many problems arise in high-dimensional settings, when  $p \gg n$ .

- ▶ First,  $S$  is **not invertible** if  $p > n$ !
- ▶ Even if  $p < n$ , but  $n$  is not very large, we may still get poor estimates, and many false positives/negatives.



## Estimating GGMs in High Dimensions

- ▶ A number of methods have been recently proposed for estimating GGMs in high dimensions.
- ▶ The main idea in most of these methods is to **use a regularization penalty**, like the **lasso**.
- ▶ We discuss two approaches:
  - ▶ neighborhood selection
  - ▶ graphical lasso

## Estimating GGMs in High Dimensions – Method 1

The idea behind **neighborhood selection**, is to estimate the graph by fitting a **penalized regression of each variable on all other variables**.

- ▶ Find **neighbors** of each node  $X_j$  by  $l_1$ -penalized regression or lasso:

$$\text{minimize}_{\beta^j} \|X_j - X_{\neq j}\beta^j\|_2^2 + \lambda \sum_{k \neq j} |\beta_k^j|$$

- ▶ The final estimate is found by combining all of the edges from these individual regression problems.
  - ▶ Symmetry —  $\beta_k^j$  not always same as  $\beta_j^k$ .
  - ▶ Use min or max rule.

## Estimating GGMs in High Dimensions – Method 2

Estimate a sparse  $\Theta$  via penalized maximum likelihood estimation (MLE).

Graphical Lasso (glasso)

$$\text{maximize}_{\Theta} \log\det(\Theta) - \text{tr}(S\Theta) - \lambda \|\Theta\|_1$$

- ▶ **Blue**: Log-likelihood;  $\log\det$  denotes the logarithm of the determinant of  $\Theta$  and  $\text{tr}$  the trace (sum of diagonal elements)  $S\Theta$ .
- ▶ **Red**: Penalty term encourages zeros on the off-diagonal elements of  $\Theta$ .

## Comparing the Two Approaches

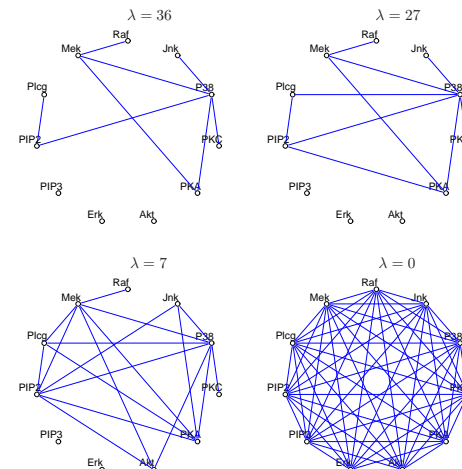
- ▶ Neighborhood selection is an **approximation for graphical lasso**:
  - ▶ Consider regression of  $X_j$  on  $X_k, j \neq k$
  - ▶ Then, the regression coefficient for neighborhood selection is related to the  $j, k$  element of  $\Theta$ :

$$\beta_k^j = -\frac{\Theta_{jk}}{\Theta_{jj}}$$

- ▶ Neighborhood selection is computationally more efficient, and may give better estimates, but doesn't give an estimate of  $\Theta$ !

## A Real Example

- ▶ **Flow cytometry** proteomics in single cells (Sachs et al, 2003).
- ▶  $p = 11$  proteins measured in  $n = 7466$  cells





## How to Choose $\lambda$ ?

- ▶  $\lambda$  modulates trade-off between **model fit** and **network sparsity**:
  - ▶  $\lambda = 0$  gives a dense network (no sparsity).
  - ▶ As  $\lambda$  increases, network becomes more sparse.
- ▶ A number of approaches proposed in the literature and used in practice
  1. **Cross-Validation** — tends to yield overly dense networks.
  2. **Extended BIC** — adjusted BIC for high dimensions.
  3. **Controlling the probability of falsely connecting disconnected components** at level  $\alpha$  (Banerjee et al, 2008):

$$\lambda(\alpha) = \frac{t_{n-2}(\alpha/2p^2)}{\sqrt{n-2 + t_{n-2}(\alpha/2p^2)}},$$

( $t_{n-2}(\alpha)$  is the  $(100 - \alpha)\%$  quantile of  $t$ -dist with  $n - 2$  d.f.)

4. **Stability selection** — Choose  $\lambda$  that gives the most **stable network** (De la Fuente et al, 2008)

## Other Types of Graphical Models

## Nonparanormal (Gaussian Copula) Models

- ▶ Suppose  $X \approx N(0, \Sigma)$ , but there **exist monotone functions**  $f_j, j = 1, \dots, p$  such that  $[f_1(X_1), \dots, f_p(X_p)] \sim N(0, \Sigma)$ 
  - ▶  $X$  has a nonparanormal distribution  $X \sim NPN_p(f, \Sigma)$ .
  - ▶  $f$  and  $\Sigma$  are parameters of the distribution, and estimated from data.
  - ▶ For continuous distributions, the nonparanormal family is **the same as the Gaussian copula family**
- ▶ To estimate the nonparanormal network:
  - i) **transform the data**:  $[f_1(X_1), \dots, f_p(X_p)]$
  - ii) **estimate the network of the transformed data** (e.g. calculate the empirical covariance matrix of the transformed data, and apply glasso or neighborhood selection)

## A Related Procedure

- ▶ Liu et al (2012) and Xue & Zou (2012) proposed a closely related idea using **rank-based correlation**
  - ▶ Let  $r_j^i$  be the **rank of  $x_j^i$**  among  $x_j^1, \dots, x_j^n$  and  $\bar{r}_j = (n + 1)/2$  be the average rank
  - ▶ Calculate **Spearman's  $\rho$**  or **Kendall's  $\tau$**

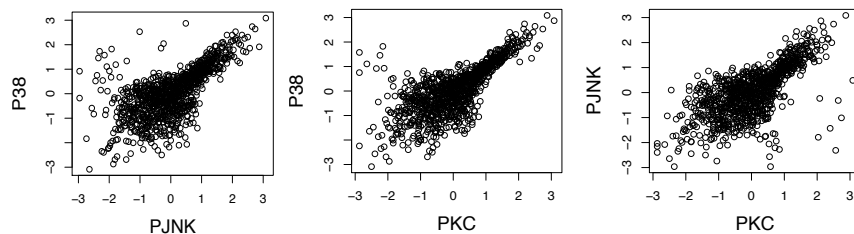
$$\hat{\rho}_{jk} = \frac{\sum_{i=1}^n (r_j^i - \bar{r}_j)(r_k^i - \bar{r}_k)}{\sqrt{\sum_{i=1}^n (r_j^i - \bar{r}_j)^2 \sum_{i=1}^n (r_k^i - \bar{r}_k)^2}}$$

$$\hat{\tau}_{jk} = \frac{2}{n(n-1)} \sum_{1 \leq i < i' \leq n} \text{sign}((x_j^i - x_j^{i'})(x_k^i - x_k^{i'}))$$

- ▶ If  $X \sim NPN_p(f, \Sigma)$ , then  $\Sigma_{jk} = 2 \sin(\rho_{jk}\pi/6) = \sin(\tau_{jk}\pi/2)$
- ▶ Therefore, we can estimate  $\Sigma^{-1}$  by **plugging in rank-based correlations into graphical lasso** (R-package huge)

## A Real Data Example

- ▶ Protein cytometry data for cell signaling (Sachs et al, 2005)
- ▶ Transform the data using a **Gaussian copula** (Liu et al, 2009), giving marginal normality
- ▶ Pairwise relationships still seem **non-linear**



- ▶ Shapiro-Wilk test rejects multivariate normality:  
 $p < 2 \times 10^{-16}$

## Graphical Models for Discrete Random Variables

- ▶ In many cases, biological data are not Gaussian: SNPs, RNAseq, etc
- ▶ Need to estimate CIG for other distributions: **binomial**, **poisson**, etc
- ▶ In this case, the estimators do not have a closed-form!
- ▶ A special case, which is computationally more tractable, is the class of **pairwise MRFs**

## Pairwise Markov Random Fields

- ▶ The idea of **pairwise MRFs** is to “assume” that **only two-way interactions among variables** exist
  - ▶ The pairwise MRF associated with graph  $G$  over the random vector  $X$  is the family of probability distributions  $P(X)$  that can be written as

$$P(X) \propto \exp \sum_{(j,k) \in E} \phi_{jk}(x_j, x_k)$$

- ▶ For each edge  $(j, k) \in E$ ,  $\phi_{jk}$  is called the **edge potential function**
- ▶ For discrete random variables, any MRF can be transformed to an MRF with pairwise interactions by introducing additional variables<sup>3</sup>

<sup>3</sup>Wainwright & Jordan, 2008

## Graphical Models for Binary Random Variables

- ▶ Suppose  $X_1, \dots, X_p$  are binary random variables, corresponding to, e.g. SNPs, or DNA methylation
- ▶ A special case of discrete graphical models is the **Ising model for binary random variables**

$$P_{\theta}(x) = \frac{1}{Z(\theta)} \exp \left\{ \sum_{(j,k) \in E} \theta_{jk} x_j x_k \right\}$$

- ▶ A **pairwise MRF** for binary data, with  $\phi_{jk}(x_j, x_k) = \theta_{jk} x_j x_k$
- ▶  $x^i \in \{-1, +1\}^p$
- ▶ The **partition function**  $Z(\theta)$  ensures that the distribution sums to 1
- ▶  $(j, k) \in E$  iff  $\theta_{jk} \neq 0$ !

## Graphical Models for Binary Random Variables

- ▶ We can consider a **neighborhood selection**<sup>4</sup> approach with an  $\ell_1$  (lasso) penalty to find the neighborhood of each node  $N(j) = \{k \in V : (j, k) \in E\}$
- ▶ For  $j = 1, \dots, p$ , need to solve (after some algebra)

$$\min_{\theta} \left\{ n^{-1} \sum_{i=1}^n \left[ f(\theta; x^i) - \sum_{k \neq j} \theta_{jk} x_j^i x_k^i + \lambda \|\theta_{-j}\|_1 \right] \right\}$$

- ▶  $f(\theta; x) = \log \left\{ \exp \left( \sum_{k \neq j} \theta_{jk} x_k \right) + \exp \left( - \sum_{k \in -j} \theta_{jk} x_k \right) \right\}$
- ▶ This is equivalent to **solving  $p$  penalized logistic regression** problems, which is straightforward (R-package `glmnet`)

<sup>4</sup>Ravikumar et al (2010)

## Other Non-Gaussian Distributions

- ▶ Assume a **pairwise graphical model**

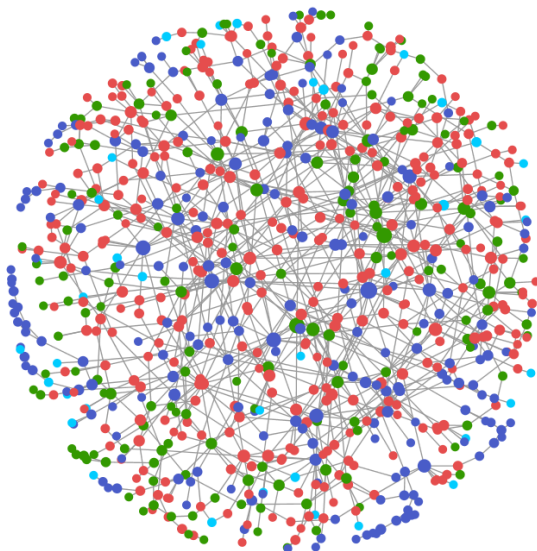
$$P(X) \propto \exp \left\{ \sum_{j \in V} \theta_j \phi_j(X_j) + \sum_{(j,k) \in E} \theta_{jk} \phi_{jk}(X_j, X_k) \right\}$$

- ▶ Then, similar to the Ising model, graphical models can be learned for other members of the **exponential family**
  - ▶ Poisson graphical models (for e.g. RNAseq), Multinomial graphical models, etc
  - ▶ All of these can be learned using a **neighborhood selection approach**, using the `glmnet` package<sup>5</sup>
  - ▶ We can even learn networks with multiple types of nodes (gene expression, SNPs, and CNVs)<sup>6</sup>

<sup>5</sup>Yang et al (2012)

<sup>6</sup>Yang et al (2014), Chen et al (2015)

## Mixed Graphical Models



## A General Approach for Estimation of Graphical Models

- ▶ Consider  $n$  iid observations from a  $p$ -dimensional random vector  $x = (X_1, \dots, X_p) \sim \mathcal{P}$
- ▶ Consider the (undirected) graph  $G = (V, E)$  with vertices  $V = \{1, \dots, p\}$
- ▶ Want to estimate edges  $E \subset V \times V$  that satisfy  $\forall j \in V, \exists N(j)$  such that:

$$p_j(X_j | \{X_k, k \neq j\}) = p_j(X_j | \{X_k : k \in N(j)\}) = p_j(X_j | \{X_k : (k, j) \in E\})$$

- ▶  $N(j)$  is the minimal set of variables on which the conditional densities depend

## Estimating Conditional Independencies

Question: how to condition?

- ▶ **Approach 1:** Estimate the joint density  $f(X_1, \dots, X_p)$ ; then get the conditionals  $f_j(X_j | X_{-j})$ 
  - ▶ Efficient, coherent
  - ▶ Computationally challenging
  - ▶ Restrictive: how many joint distributions do you know?
  - ▶ Hard to check if assumptions hold!
- ▶ **Approach 2:** Estimate the conditionals directly  $f_j(X_j | X_{-j})$ 
  - ▶ Computationally easy
  - ▶ Leads to easy & flexible models (regression)!
  - ▶ May not be efficient or coherent

## A Semi-parametric Approach

- ▶ Consider additive non-linear relationships (additive model):

$$X_j | X_{-j} = \sum_{k \neq j} f_{jk}(X_k) + \varepsilon$$

- ▶ Then if  $f_{jk}(X_k) = f_{kj}(X_j) = 0$ , we conclude that  $X_j$  and  $X_k$  are **conditionally independent**, given the other variables
- ▶ In other words, we **assume that conditional distributions and conditional means depend on the same set of variables**
- ▶ We then use a semi-parametric approach for estimating the conditional dependencies

## SpaCE JAM<sup>7</sup>

- ▶ Sparse Conditional Estimation with Jointly Additive Models (SpaCE JAM)

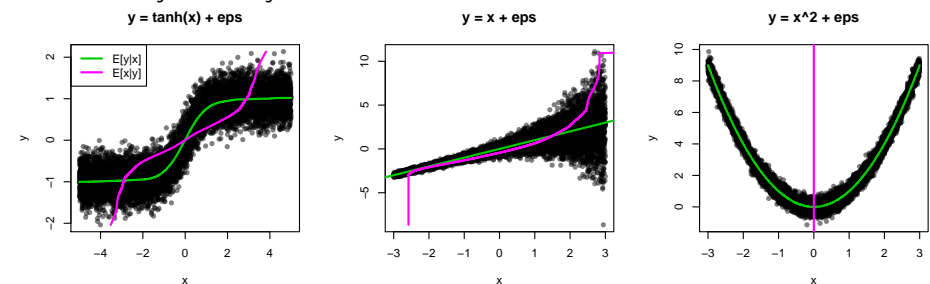
$$\text{minimize}_{f_{jk} \in \mathcal{F}} \frac{1}{2n} \sum_{j=1}^p \|x_j - \sum_{k \neq j} f_{jk}(x_k)\|_2^2 + \lambda \sum_{k > j} (\|f_{jk}(x_k)\|_2^2 + \|f_{kj}(x_j)\|_2^2)^{1/2}$$

- ▶  $f_{jk}(x_k) = \Psi_{jk} \beta_{jk}$
- ▶  $\Psi_{jk}$  is a  $n \times r$  matrix of basis functions for  $f_{jk}$
- ▶  $\beta_{jk}$  is an  $r$ -vector of coefficients
- ▶ The **standardized group lasso** penalty for functions  $\|f_{jk}\|_2$
- ▶ This is a **convex** problem, and **block coordinate descent** converges to the global minimum

<sup>7</sup>Voorman et al (2014), R-package spacejam

## SpaCE JAM

Estimating  $f_{jk}$  and  $f_{kj}$  seems redundant...



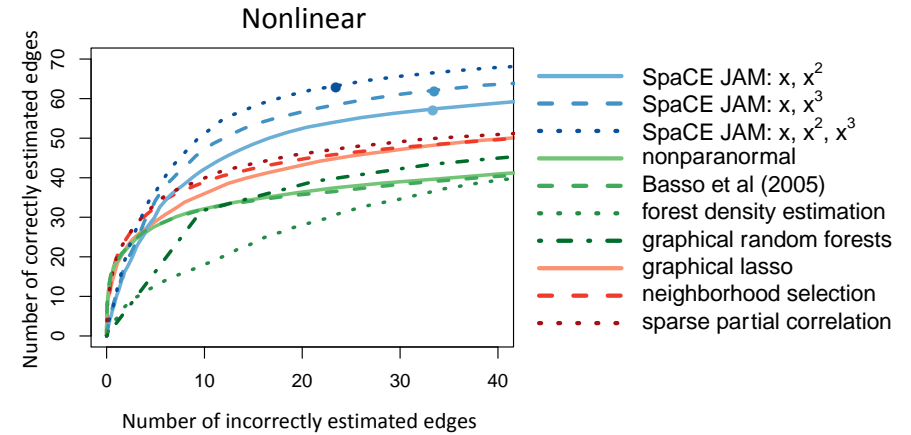
but **necessary for non-linear functions**

## Other Flexible Procedures

- ▶ **Forest density estimation** (Liu et al, 2011) assumes that underlying graph is a forest, and estimates the bivariate densities non-parametrically.
- ▶ **Graphical random forests** (Fellinghauer et al, 2013) uses random forests to flexibly model conditional means
  - ▶ They consider conditional dependencies through conditional mean
  - ▶ They allow for general random variables, discrete or continuous
  - ▶ Use a **random forest** to estimate  $E[X_j | X_{\setminus j}]$  non-parametrically
  - ▶ Theoretical properties have not yet been justified

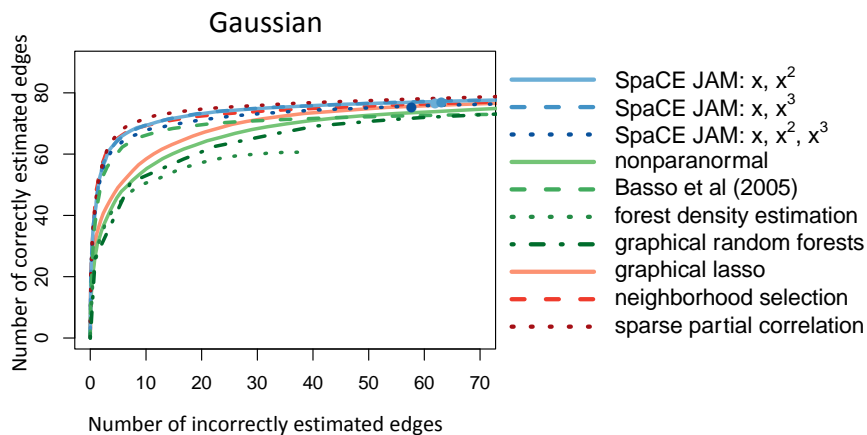
## Comparison on Simulated Data

non-linear relationships ( $p = 100, n = 50$ )

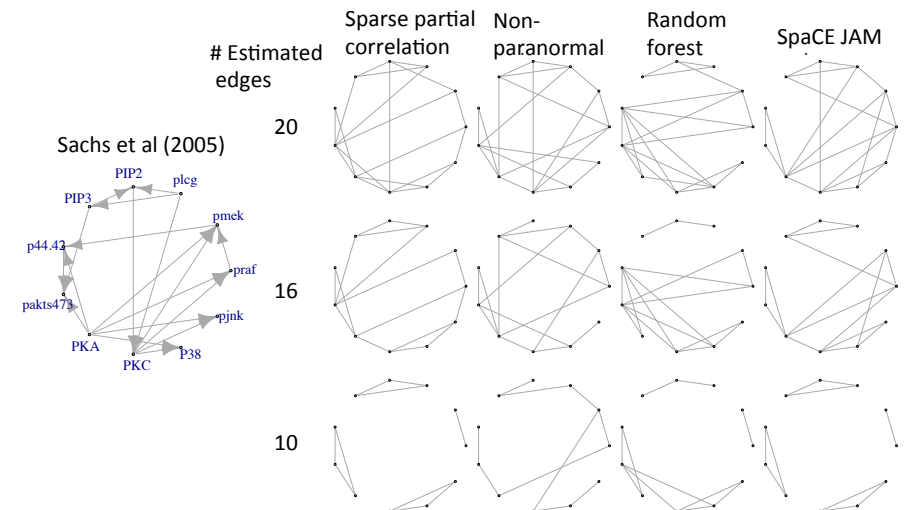


## Comparison on Simulated Data

linear relationships ( $p = 100, n = 50$ )



## Estimation of Cell Signaling Network



## Other Extensions of GGMs

- ▶ Multiple Graphical Models
  - ▶ For groups of observations, estimate graphical models with shared structure across groups and individual structure within groups.
- ▶ Time Varying Graphical Models
  - ▶ Smoothly varying graph over time estimated via local kernel smoothers.
  - ▶ Change points in graph structure over time estimated via fusion penalties.
- ▶ Latent Variable Graphical Models
  - ▶ Assume observed features are dependent on latent variables which exhibit a low-rank effect. Estimate a sparse (graph structure) plus low-rank inverse covariance matrix.

## Joint Estimation of Multiple Graphical Models

### A Brief Introduction

Key idea:

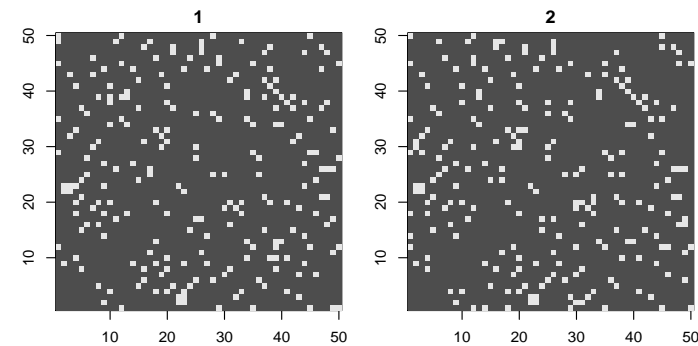
- ▶ We observe data from different **a priori known sub-populations**
  - ▶ Sub-populations may correspond to sub-types of a disease (e.g. neural, proneural, mesenchymal and classical in glioblastoma)
- ▶ For each sub-population, there exists an undirected graphical model
- ▶ The underlying graphical models **share common structure**

Main issue: **What type of common structure is assumed, and how to enforce it?**

## Illustrative framework: Gaussian case

- ▶  $X_k \sim N(0, \Sigma_k), k = 1, \dots, K.$
- ▶  $\Theta_k \equiv \Sigma_k^{-1}, k = 1, \dots, K$  so that  $\Theta_k \sim \Theta_\ell$ , for all  $k \neq \ell$
- ▶ Different approaches build either on maximum likelihood estimation or on neighborhood selection

## Pictorial Motivation - I



## Selected Approaches - I

Rich literature on the topic with many variants appearing in statistics, machine learning and bioinformatics literature

- ▶ **Hierarchical penalty**<sup>8</sup>:
  - ▶ Let  $\Theta_k(i, j) = \alpha(i, j)\gamma_k(i, j)$ ,  $k = 1, \dots, K$
  - ▶ For identifiability, assume  $\alpha(i, j) \geq 0$  for all variable pairs  $(i, j)$
  - ▶  $P(\Theta) = \lambda_\alpha |\alpha|_1 + \lambda_\gamma \sum_{k=1}^K |\gamma_k|_1$  — combine two lasso penalties

<sup>8</sup>Guo et al (2011)

## Selected Approaches - II

- ▶ **Fusing penalties**<sup>9</sup>:

- ▶ Variant 1:

$$P(\Theta) = \sum_{k=1}^K \lambda_k |\Theta_k|_1 + \sum_{k \neq \ell} \lambda_{k,\ell} |\Theta_k - \Theta_\ell|_1$$

element-wise fused lasso penalty, encourages **similarities between all elements** of the  $K$  graphical models

- ▶ Variant 2:

$$P(\Theta) = \sum_{k=1}^K \lambda_k |\Theta_k|_1 + \sum_{i,j} \sqrt{\sum_{k=1}^K (\Theta_k(i, j))^2}$$

encourages **strong fusing towards a common graphical model**

<sup>9</sup>Danaher et al (2014)

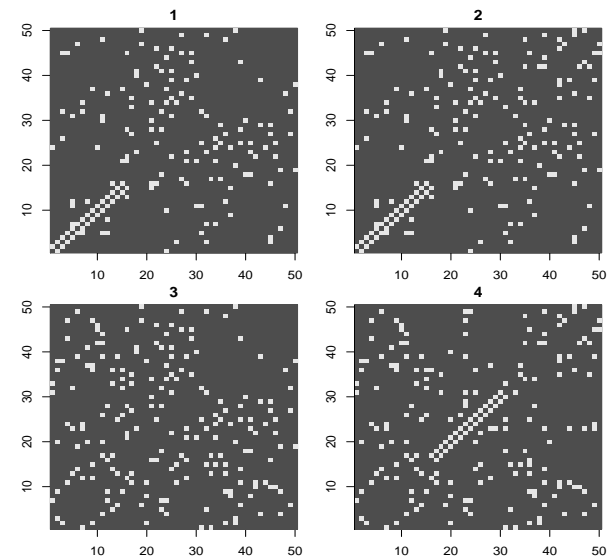
## Selected Approaches - III

- ▶ Cai et al (2016) propose a mixed  $\ell_\infty/\ell_1$  norm:

$$\begin{aligned} & \min_{\{\Theta\}_{k=1}^K} (\max_{1 \leq k \leq K} |\Theta_k|_1) \\ \text{s.t. } & \max_{i,j} \left( \sum_{k=1}^K \frac{n_k}{n} |S_k \Theta_k - I|_{(i,j)}^2 \right)^{1/2} \leq t_n \end{aligned}$$

- ▶ The objective function encourages sparsity across all  $K$  models. The constraint is imposed on the maximum of the element-wise group  $\ell_2$  norm to encourage the groups to share a common graphical structure.

## Pictorial Motivation - II



## Selected Approaches - IV

- ▶ Saegusa & S (2016) encode **similarity** between different sets of edges for pairs of models  $(k, \ell)$  through a Laplacian penalty
- ▶ Ma & M (2016) use group lasso penalties across different subsets of the edges

Both approaches require external information through prior knowledge (e.g. functional pathways, literature, etc.)

## Application: Lipid Interaction Networks in CKD

- ▶ Chronic Kidney Disease (CKD) is strongly linked to cardiovascular morbidity and mortality.
- ▶ Despite the diversity of human plasma lipidome, studies of CKD have been traditionally limited to measuring total cholesterol, triglycerides, and lipoproteins
- ▶ New technologies allow researchers to profile a large number of lipid species ( $p \sim 450$ ) from various lipid classes

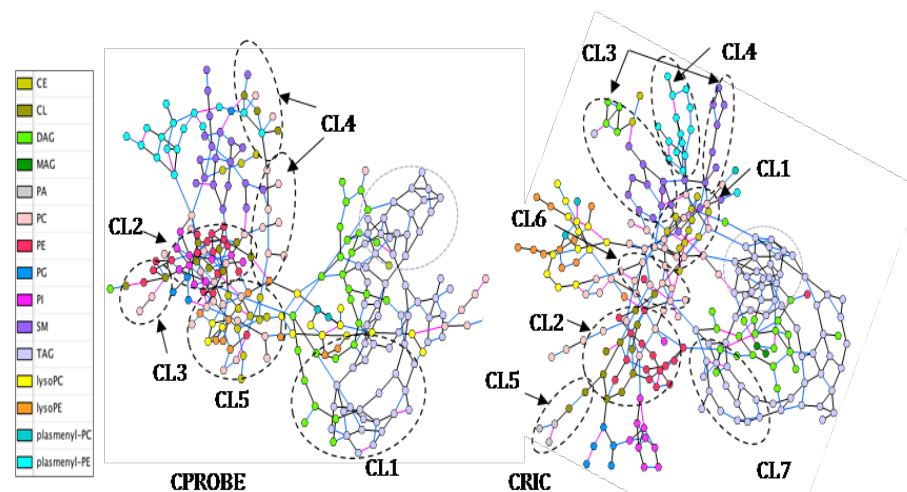
## Lipid Interaction Networks in CKD

Objective: Understand lipid interactions from two related study cohorts<sup>10</sup>

- ▶ Clinical Phenotyping Resource and Biobank Core (CPROBE) — progressors vs non-progressors patients
- ▶ Chronic Renal Insufficiency Cohort (CRIC) — early stage CKD vs late stage CKD patients

<sup>10</sup>Analysis pipeline and results in Ma et al (2019)

## CPROBE/CRIC Modules





## CPROBE/CRIC Differential Sub-Networks<sup>11</sup>

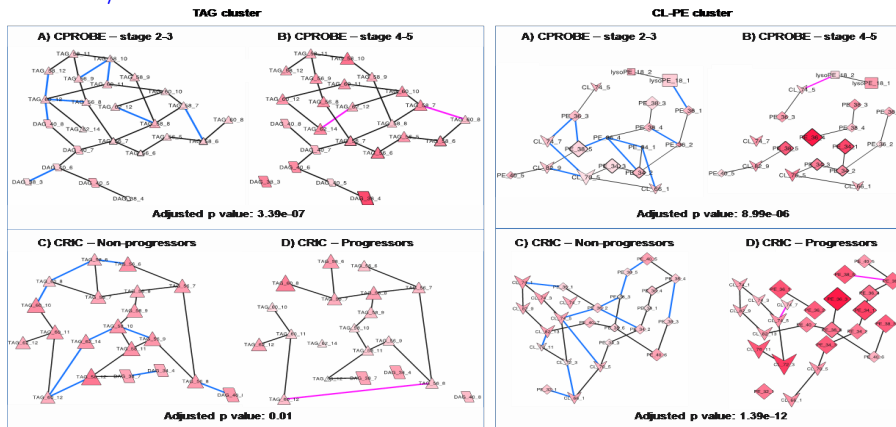


Figure: Black edges: common backbone; blue edges: present in early stage/NP; red edges: present in late stage/P

<sup>11</sup>Based on NetGSA enrichment analysis — see Lecture 2

## Biological Relevance of Discovered Modules

- ▶ Of interest is the “disregulation” of a module comprising of triacylglycerols (TAGs) and another one of cardiolipins with phosphatidylethanolamines (CL-PE) in CRIC/CPROBE
- ▶ Of particular interest is the second (CL-PE) module, that points to role of cellular lipid metabolism and specifically the activity of the mitochondrial respiratory chain after checking the module for enrichment; thus, the loss of lipids may lead to decreased mitochondrial fusion and fragmented mitochondria
- ▶ Concordant with recent findings in the literature that mitochondrial damage and dysfunction might be a highly prevalent abnormality in early CKD (eGFR>60)