

Covariance Properties and Graph Selection for High-Dimensional Compositional Data

Camden Lopez

WNAR 2017, Student Paper Session 3

- Compositional microbiome data
- Graphical model selection using SPIEC-EASI
- Covariance relationships and properties
- Graph selection performance

Compositional microbiome data

16S amplicon sequencing

- Sample \rightarrow DNA \rightarrow 16S sequences \rightarrow OTU counts
- **Relative abundances only**

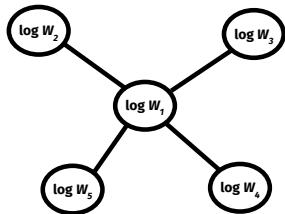
	OTU 1	OTU 2	...	OTU p
Sample 1	136	28	...	10
Sample 2	0	2	...	18
\vdots	\vdots	\vdots	\ddots	\vdots
Sample n	54	25	...	5

OTU = operational taxonomic unit

Graphical model inference using SPIEC-EASI

SPIEC-EASI: SParse **I**nverse **E** Covariance Estimation for **E**cological **A**sociation **I**nference (Kurtz et al. 2015)

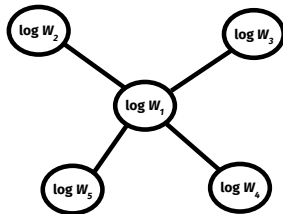
- Relationships among OTU abundances $W = (W_1, \dots, W_p)$?
- Suppose $\log W \sim \text{Normal}(\cdot, \Omega)$
- **Non-zero entries of $\Omega^{-1} \Leftrightarrow$ conditional dependence, graphical model**



Graphical model inference using SPIEC-EASI

SPIEC-EASI: SParse **I**nverse **C**ovariance Estimation for **E**cological **AS**sociation **I**nference (Kurtz et al. 2015)

- **Observe** $W \times ? \rightarrow \log W + ?$
- **Centering** $\log W + ? \rightarrow \text{clr } W$
- **Assumption:**
 $\text{cov}(\text{clr } W) = \Gamma \approx \Omega = \text{cov}(\log W)$
- **Graphical model inference:**
 $\hat{\Gamma} \rightarrow \hat{\Omega}^{-1}$ (graphical lasso, e.g.)



Properties of Γ

- $\gamma_{ij} = \omega_{ij} - \bar{\omega}_{i.} - \bar{\omega}_{.j} + \bar{\omega}_{..}$
- Rows/col's sum to zero
- p fewer free parameters than Ω

Covariance relationships and properties

Properties of Γ

- $\gamma_{ij} = \omega_{ij} - \bar{\omega}_{i.} - \bar{\omega}_{.j} + \bar{\omega}_{..}$
- Rows/col's sum to zero
- p fewer free parameters than Ω

$\Gamma \approx \Omega$

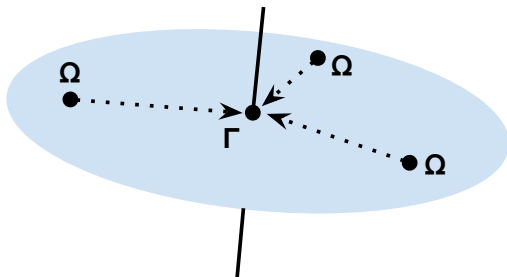
- Small or mostly negative correlations
- Approx. equal variances
- **Small**
“compositional effect”

$\Gamma \not\approx \Omega$

- Mostly positive correlations
- Unequal variances
- **Large**
“compositional effect”

Covariance relationships and properties

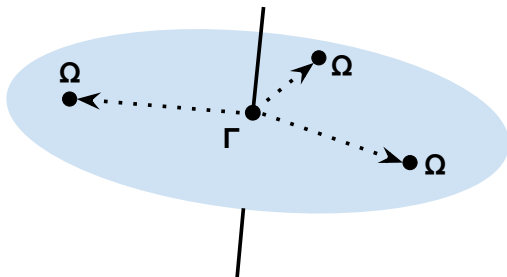
One $\Gamma \leftrightarrow$ many Ω



- For each Γ , p -dimensional space of **potential** Ω s

Covariance relationships and properties

One $\Gamma \leftrightarrow$ many Ω



- Can solve for **potential** Ω s (must check $\Omega \succ 0$)

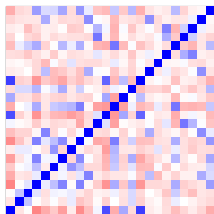
Covariance relationships and properties

Relationships can vary among potential Ω s

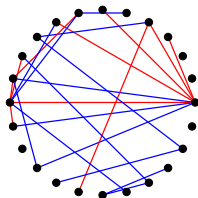
- Somewhat constrained, but not entirely

Example ($p = 24$, red = negative, blue = positive)

$\hat{\Gamma}$ (correlations)



Graph (24 edges)



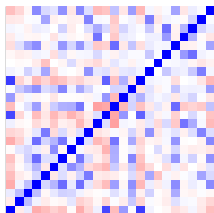
Covariance relationships and properties

Relationships can vary among potential Ω s

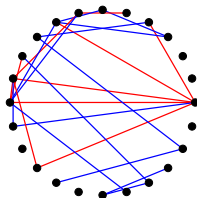
- Somewhat constrained, but not entirely

Example ($p = 24$, red = negative, blue = positive)

Potential $\hat{\Omega}$ #1



Graph (24 edges)



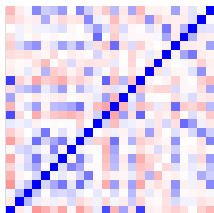
Covariance relationships and properties

Relationships can vary among potential Ω s

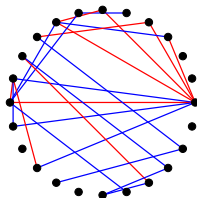
- Somewhat constrained, but not entirely

Example ($p = 24$, red = negative, blue = positive)

Potential $\hat{\Omega}$ #2



Graph (24 edges)



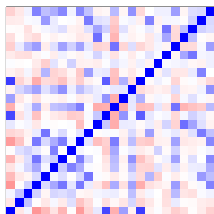
Covariance relationships and properties

Relationships can vary among potential Ω s

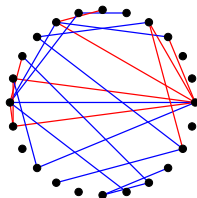
- Somewhat constrained, but not entirely

Example ($p = 24$, red = negative, blue = positive)

Potential $\hat{\Omega}$ #3



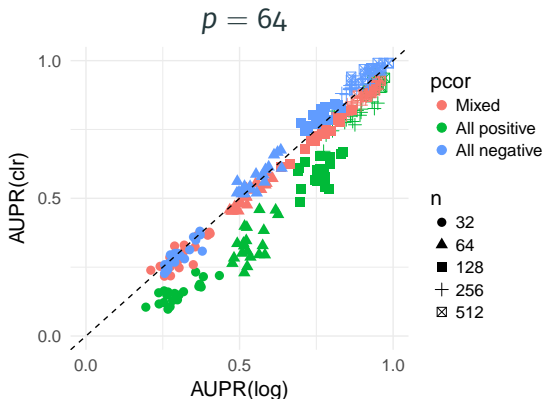
Graph (24 edges)



Graph selection performance

Performance with **small compositional effect**

- Comparable to graph selection from log data

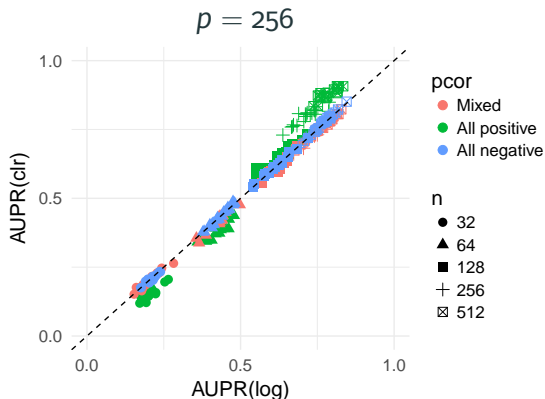


- Sparse graph
- Approx. equal variances
- Partial correlations (pcor) ± 0.25
- AUPR = area under precision-recall curve

Graph selection performance

Performance with **small compositional effect**

- Comparable to graph selection from log data

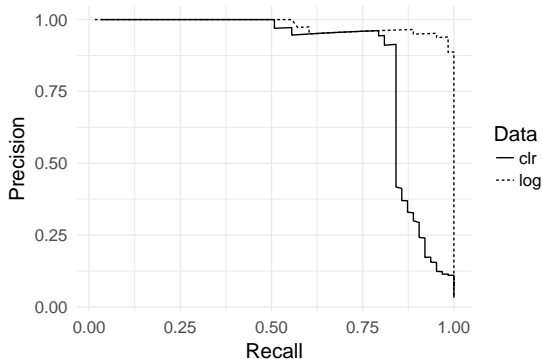


- Sparse graph
- Approx. equal variances
- Partial correlations (pcor) ± 0.25
- AUPR = area under precision-recall curve

Graph selection performance

Performance with **large compositional effect**

- Affected by distortion of covariances



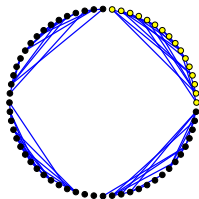
- Cluster graph,
 $p = 64$
- $25\times$ larger
variances in
one cluster
- $n = 1024$

Graph selection performance

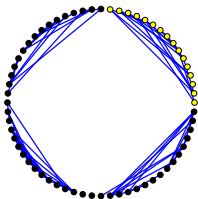
Performance with **large compositional effect**

- Affected by distortion of covariances

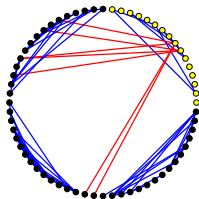
Graph



Data = log



Data = clr



- **Limitation of compositional data:** One $\Gamma \leftrightarrow$ many Ω , uncertainty about $\log W$ relationships
- **SPIEC-EASI graph selection** for $\log W$ based on $\text{clr } W$ data **performs well ...** provided the compositional effect is not too large
- **Large compositional effect** distorts covariances and causes erroneous edges in graph

SPIEC-EASI paper:

- Kurtz., Z. D., Müller, C. L., Miraldi, E. R., Littman, D. R., Blaser, M. J., and Bonneau, R. A. (2015). Sparse and compositionally robust inference of microbial ecological networks. *PLoS Computational Biology* 11, e1004226.

Thank you Oregon State University faculty:

- Yuan Jiang, Duo Jiang, Sarah Emerson (Statistics)
- Thomas Sharpton (Microbiology and Statistics)

Centered log-ratio transformation

Given sample $x = (x_1, \dots, x_p)$ with geometric mean $g(x) = (\prod_{i=1}^p x_i)^{\frac{1}{p}}$,

$$\begin{aligned}\text{clr}(x) &= \left(\log \frac{x_1}{g(x)}, \dots, \log \frac{x_p}{g(x)} \right) \\ &= \left(\log x_1 - \frac{1}{p} \sum_{i=1}^p \log x_i, \dots, \log x_p - \frac{1}{p} \sum_{i=1}^p \log x_i \right)\end{aligned}$$

and for $c > 0$,

$$\text{clr}(cx) = \text{clr}(x)$$

Covariance relationships

Matrix form:

$$\text{clr } W = G \log W$$

$$G = \begin{pmatrix} 1 - \frac{1}{p} & \cdots & -\frac{1}{p} \\ \vdots & \ddots & \vdots \\ -\frac{1}{p} & \cdots & 1 - \frac{1}{p} \end{pmatrix}$$

$$\Rightarrow \text{cov}(\text{clr } W) = G \text{cov}(\log W)G$$

$$\Rightarrow \Gamma = G\Omega G$$

Graphical lasso estimation of Ω^{-1}

$$\widehat{\Omega^{-1}}_{\text{glasso}} = \arg \max_{\Omega^{-1} \succeq 0} \left[\log \det(\Omega^{-1}) - \text{tr}(\widehat{\Omega} \Omega^{-1}) - \lambda \|\Omega^{-1}\|_1 \right]$$

$$\widehat{\Omega^{-1}}_{\text{SPIEC-EASI}} = \arg \max_{\Omega^{-1} \succeq 0} \left[\log \det(\Omega^{-1}) - \text{tr}(\widehat{\Gamma} \Omega^{-1}) - \lambda \|\Omega^{-1}\|_1 \right]$$

Solving for potential Ω

Option 1: Choose $\bar{\omega}_{1.}, \dots, \bar{\omega}_{p.}$

$$\begin{aligned}\gamma_{ij} &= \omega_{ij} - \bar{\omega}_{i.} - \bar{\omega}_{.j} + \bar{\omega}_{..} \\ \Rightarrow \omega_{ij} &= \gamma_{ij} + \bar{\omega}_{i.} + \bar{\omega}_{.j} - \bar{\omega}_{..}\end{aligned}$$

Option 2: Choose $\omega_{11}, \dots, \omega_{pp}$

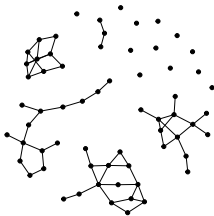
$$\omega_{ij} = \gamma_{ij} + \frac{1}{2} (\omega_{ii} - \gamma_{ii} + \omega_{jj} - \gamma_{jj})$$

Graphs used in simulations

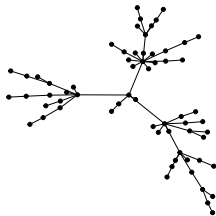
$$p = 64, e = p - 1$$



Band



Cluster



Scale-free