

Network-based Methods for Analysis of Microbiome and Metabolomic Data

Jing Ma

Division of Public Health Sciences
Fred Hutchinson Cancer Research Center

3 Nov 2021

Outline

Regression Analysis of Structured Microbiome Data

Differential Network Enrichment Analysis

Conclusion

Outline

Regression Analysis of Structured Microbiome Data

Differential Network Enrichment Analysis

Conclusion

Microbiome Science

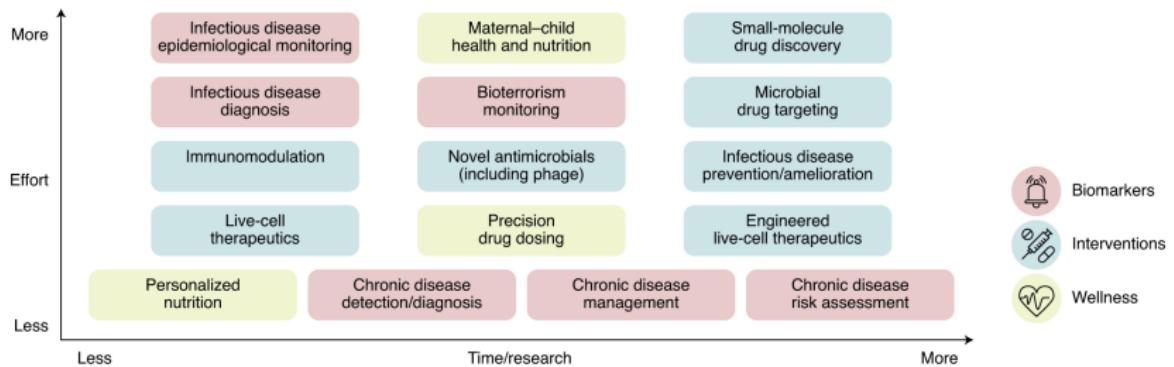
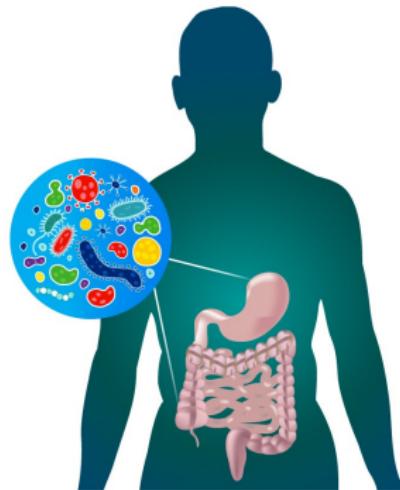


Figure: Plot from Wilkinson et al. (21') Nat Medicine.

The Gut Microbiome



Interest in microorganisms (their genome and products) that live in the digestive tracts of humans and other animals.

Microbiome Data

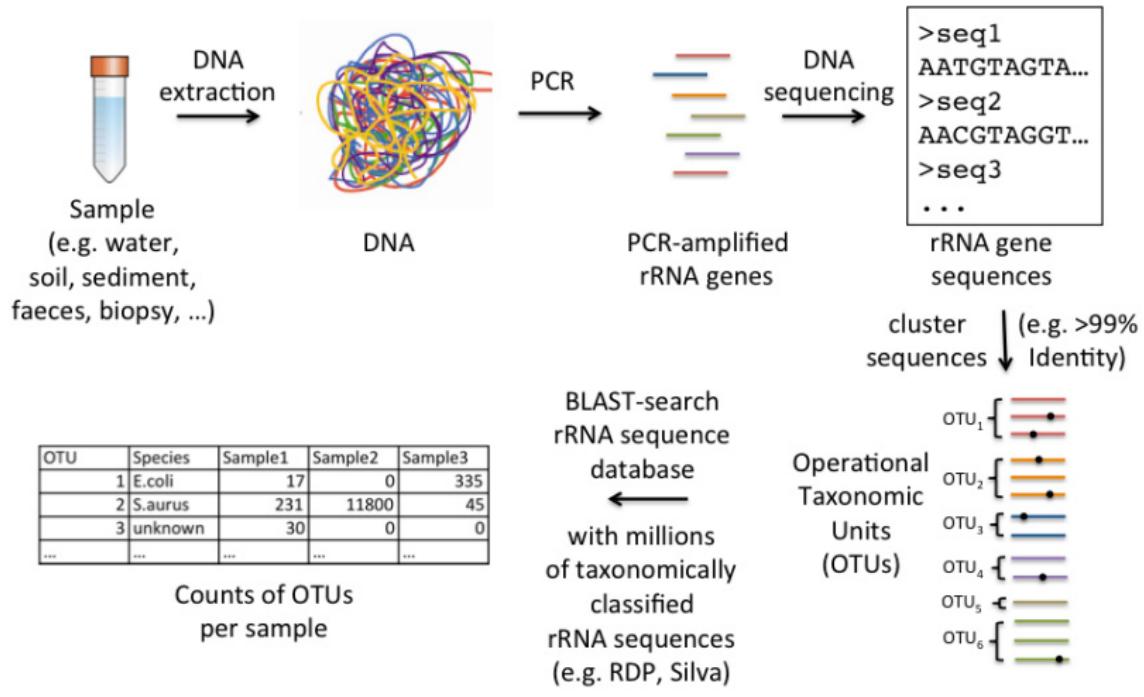
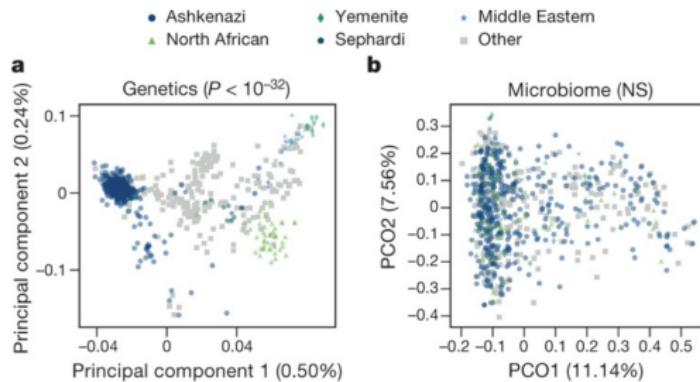


Figure: Plot courtesy to Anders Andersson

Samples are Correlated

- ▶ Microbiome compositions are shaped by environment, diet, and culture
- ▶ Sample similarity is often available (e.g. UniFrac distance)



D Rothschild *et al.* *Nature* **555**, 210–215 (2018) doi:10.1038/nature25973

Variables are Correlated

- ▶ Taxon similarity can be computed by the patristic distance between tips of the tree

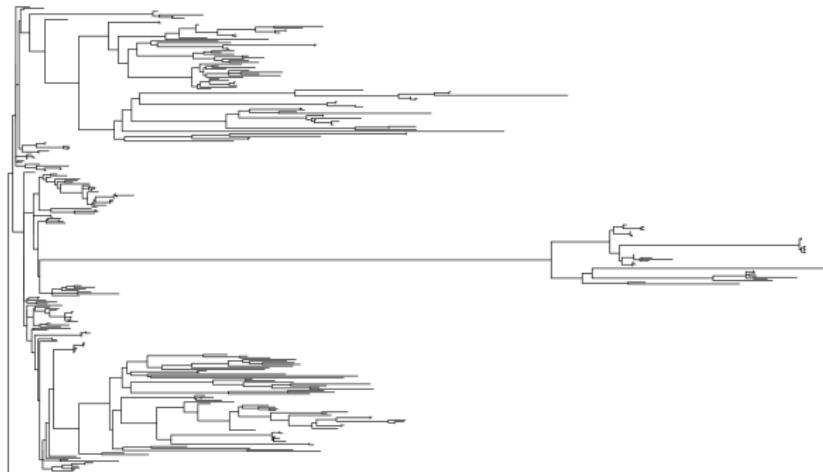


Figure: A phylogenetic tree based on data from Wilmanski et al. (21') Nat Metabolism.

Variables are Correlated

- Taxon similarity can be computed by measuring metabolic functional similarities

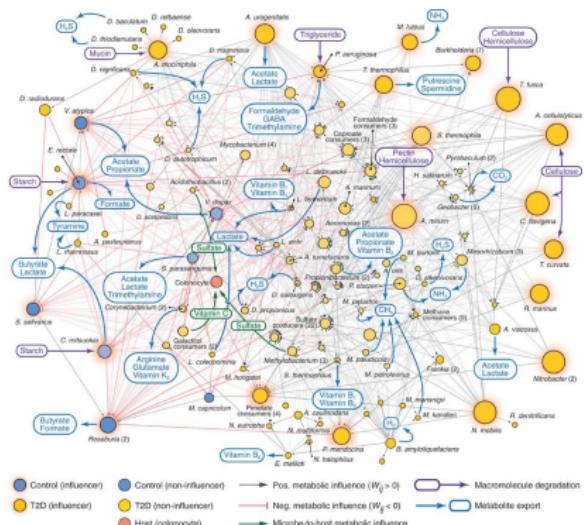
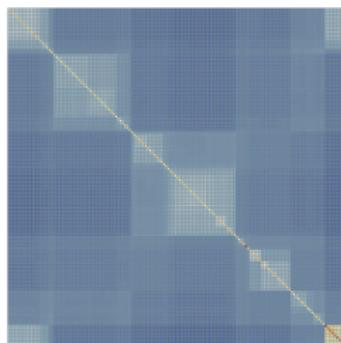
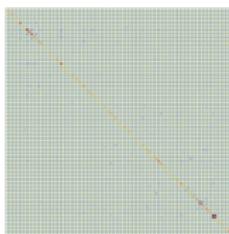


Figure: A metabolic interaction network from Sung et al. (17') Nat Commun.

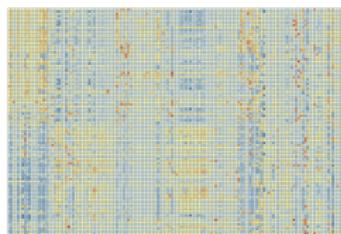
Two-way Structured Data



$Q : p \times p$



$H : n \times n$



$X : n \times p$

Biplot via GMD

- ▶ Generalized matrix decomposition (GMD¹) extends PCA/SVD to two-way structured data

Biplot via GMD

- ▶ Generalized matrix decomposition (GMD¹) extends PCA/SVD to two-way structured data
- ▶ SVD solves

$$\arg \min_{U, S, V} \|X - USV^T\|_F$$

where $\|A\|_F = \text{trace}(A^T A)$.

Biplot via GMD

- ▶ Generalized matrix decomposition (GMD¹) extends PCA/SVD to two-way structured data
- ▶ SVD solves

$$\arg \min_{U, S, V} \|X - USV^T\|_F$$

where $\|A\|_F = \text{trace}(A^T A)$.

- ▶ GMD solves

$$\arg \min_{U, S, V} \|X - USV^T\|_{H, Q}$$

where $\|A\|_{H, Q} = \text{trace}(A^T HAQ)$.

¹Allen et al. (14') JASA

GMD Improves Visualization of Samples

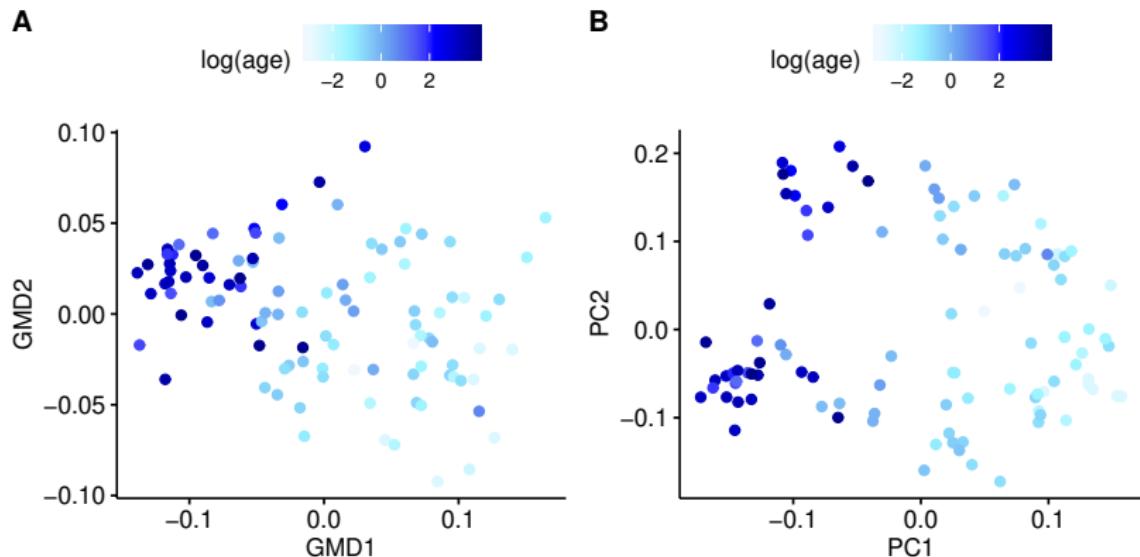


Figure: Plot of samples based on data from Yatsunenko et al. (12') Nature.

GMD Regression

$$\mathbf{H} \underset{n \text{ samples}}{\underset{p \text{ variables}}{\approx}} \mathbf{Q} = \mathbf{U}_1 \mathbf{V}_1^T + \cdots + \mathbf{U}_K \mathbf{V}_K^T$$

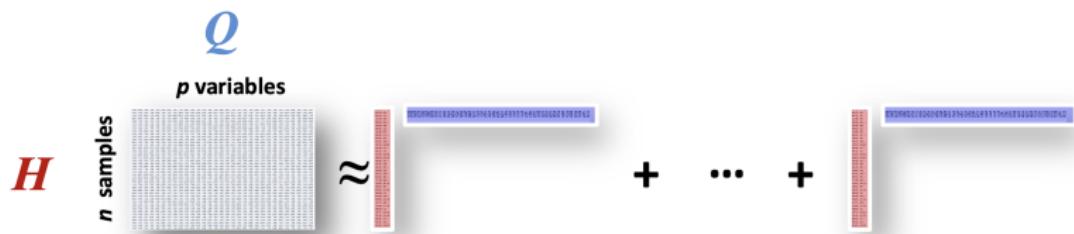
GMD:

$$\mathbf{X} \approx \sigma_1 \mathbf{U}_1 \mathbf{V}_1^T + \cdots + \sigma_K \mathbf{U}_K \mathbf{V}_K^T$$

Weighted LS:

$$\hat{\gamma} = \operatorname{argmin}_{\gamma} \|y - \sigma_1 \mathbf{U}_1 \gamma_1 - \cdots - \sigma_K \mathbf{U}_K \gamma_K\|_{\mathbf{H}}^2$$

GMD Regression



GMD:

$$\mathbf{X} \approx \sigma_1 \mathbf{U}_1 \mathbf{V}_1^\top + \cdots + \sigma_K \mathbf{U}_K \mathbf{V}_K^\top$$

Weighted LS:

$$\hat{\gamma} = \operatorname{argmin}_{\gamma} \|y - \sigma_1 \mathbf{U}_1 \gamma_1 - \cdots - \sigma_K \mathbf{U}_K \gamma_K\|_{\mathbf{H}}^2$$

Remark: rank is selected based on variable importance and generalized cross-validation.

GMD Regression

GMDR estimator:

$$\hat{\beta}_{GMDR} = \mathbf{QV}\hat{\gamma}$$

GMD Regression

GMDR estimator:

$$\hat{\beta}_{GMDR} = \mathbf{QV}\hat{\gamma}$$

KPR estimator²:

$$\hat{\beta}_{KPR} = \arg \min_{\beta} \{ \|y - \mathbf{X}\beta\|_{\mathbf{H}}^2 + \lambda \|\beta\|_{\mathbf{Q}^{-1}}^2 \}$$

GMD Regression

GMDR estimator:

$$\hat{\beta}_{GMDR} = \mathbf{Q}\mathbf{V}\hat{\gamma}$$

KPR estimator²:

$$\hat{\beta}_{KPR} = \arg \min_{\beta} \{ \|y - \mathbf{X}\beta\|_{\mathbf{H}}^2 + \lambda \|\beta\|_{\mathbf{Q}^{-1}}^2 \}$$

The GMD class of estimators

$$\hat{\beta}^W = \mathbf{Q}\mathbf{V}\mathbf{W}\mathbf{S}^{-1}\mathbf{U}^\top \mathbf{H}y,$$

for a diagonal weight matrix $\mathbf{W} = \text{diag}(w_1, \dots, w_K)$, $w_j \geq 0$.

- ▶ $w_j = 1_{j \in \mathcal{J}} \rightarrow \hat{\beta}_{GMDR}$
- ▶ $w_j = \sigma_j^2(\sigma_j^2 + \lambda)^{-1} \rightarrow \hat{\beta}_{KPR}$

²Randolph et al. (18') AOAS

GMD Inference

GMD Inference

- ▶ Find the estimation bias

$$B_j = \mathbb{E}[\hat{\beta}_j^{\mathcal{W}} - \beta_j^*] = (\mathbf{Q}\mathcal{V}\mathcal{W}\mathcal{V}^\top\boldsymbol{\beta}^*)_j - \beta_j^*$$

GMD Inference

- ▶ Find the estimation bias

$$B_j = \mathbb{E}[\hat{\beta}_j^W - \beta_j^*] = (\mathbf{Q}\mathcal{V}\mathcal{W}\mathcal{V}^\top\boldsymbol{\beta}^*)_j - \beta_j^*$$

- ▶ Correct the estimation bias via an initial estimator $\hat{\boldsymbol{\beta}}^{init} = D\tilde{\boldsymbol{\beta}}(\lambda)$ where

$$\mathbf{Q} = D\Delta D^\top, \quad \tilde{\boldsymbol{\beta}}(\lambda) = \arg \min_{\boldsymbol{\beta}} \{ \|y - \mathbf{X}D\boldsymbol{\beta}\|_H^2 + \lambda \|\Delta^{-1/2}\boldsymbol{\beta}\|_1 \}.$$

GMD Inference

- ▶ Find the estimation bias

$$B_j = \mathbb{E}[\hat{\beta}_j^W - \beta_j^*] = (\mathbf{Q}\mathcal{V}\mathcal{W}\mathcal{V}^\top\boldsymbol{\beta}^*)_j - \beta_j^*$$

- ▶ Correct the estimation bias via an initial estimator $\hat{\boldsymbol{\beta}}^{init} = D\tilde{\boldsymbol{\beta}}(\lambda)$ where

$$\mathbf{Q} = D\Delta D^\top, \quad \tilde{\boldsymbol{\beta}}(\lambda) = \arg \min_{\boldsymbol{\beta}} \{ \|y - \mathbf{X}D\boldsymbol{\beta}\|_H^2 + \lambda \|\Delta^{-1/2}\boldsymbol{\beta}\|_1 \}.$$

- ▶ Obtain inference for $H_{0,j} : \beta_j^* = 0$ using de-biased $\tilde{\beta}_j^W = \hat{\beta}_j^W - \hat{B}_j$.

Model Assumptions

Linear model: $y = \mathbf{X}\beta^* + \varepsilon$ with $\text{Cov}(\varepsilon \mid \mathbf{X}) = \Psi = L_\Psi^\top L_\Psi$. There exists sub-Gaussian $\tilde{\varepsilon}$ such that $\varepsilon = L_\Psi^\top \tilde{\varepsilon}$.

- ▶ $\|L_\Psi \mathbf{H} L_\Psi^\top - I_n\|_2 = o(1)$ as $n \rightarrow \infty$.
- ▶ β^* is \mathbf{Q} -smooth:

$$\|\mathbf{Q}^{-1/2} \beta^*\|_0 = o\{(n/\log p)^r\} \text{ for } r \in (0, 1/2).$$

- ▶ A compatibility assumption w.r.t \mathbf{Q} and \mathbf{H}

$$0 < \underline{c} \leq \frac{\|\widetilde{\mathbf{X}}_A v\|^2}{n\|v\|^2} \leq \bar{c} < \infty,$$

where $\widetilde{\mathbf{X}} = \mathbf{H}^{1/2} \mathbf{X} \mathbf{Q}^{1/2}$ and $|A| \geq M_1^* s_0 + 1$.

Microbiome Data from Yatsunenko et al. (12)

Data

- ▶ \mathbf{X} : 100 healthy individuals from Venezuela, Malawi & US \times 149 OTUs
- ▶ \mathbf{Q} : variable similarity from the phylogenetic tree
- ▶ \mathbf{H} : sample similarity from functional profiling
- ▶ y : age (proxy to aging)

Questions

Predict age from microbiome data and identify age-associated bacteria.

Age Prediction

- ▶ Lasso
- ▶ Ridge regression
- ▶ KPR1/GMDR1: only \mathbf{Q}
- ▶ KPR2/GMDR2: only \mathbf{H}
- ▶ KPR3/GMDR3: both \mathbf{Q} and \mathbf{H}

Leave-one-out cross-validation:

Method	Lasso	Ridge	KPR1	GMDR1	KPR3	GMDR3
MSE	0.71	1.13	0.61	0.57	0.56	0.52

Significant Marginal Associations

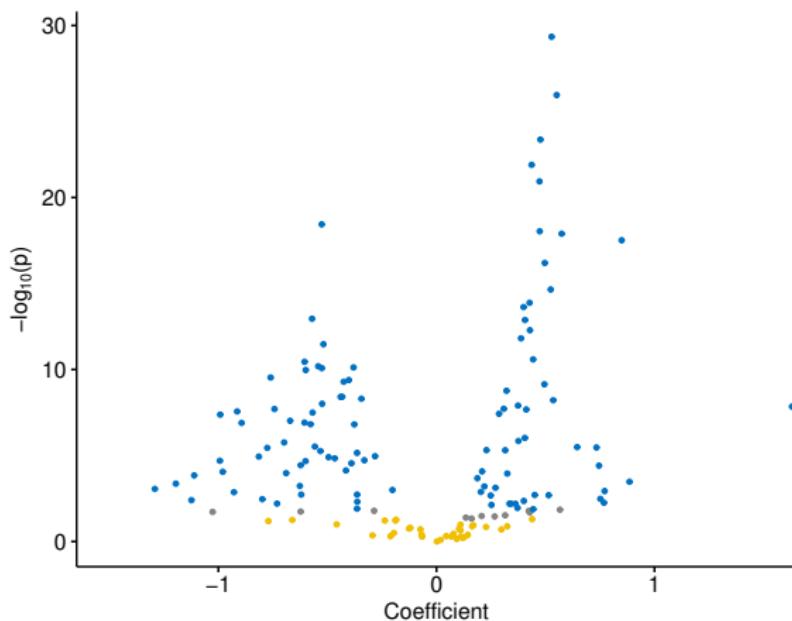


Figure: Univariate association analysis using data in Yatsunenko et al. Nature. (12') (FDR=0.1). Results suggest that either the true signals are not sparse or taxa are highly correlated.

Significant Conditional Associations

Table: Genera found to be associated with age when controlling for FDR=0.1

	Genus	Total
Ridge	(none)	0
LDPE	<i>Atopobium, Campylobacter, Peptoniphilus</i>	3
Grace	<i>Atopobium, Campylobacter</i>	2
GMDI	<i>Acidaminococcus, Aggregatibacter, Akkermansia, Alkalimonas, Anaerofustis, Anaerostipes, Arthrobacter, Atopobium, Bilophila, Brochothrix, Bulleidia, Butyrivibrio, Caloramator, Campylobacter, Citrobacter, Coprobacillus, Desulfitobacterium, Enterobacter, Epulopiscium, Facklamia, Lachnobacterium, Lactobacillus, Leclercia, Megasphaera, Methylophaga, Mitsuokella, Peptococcus, Peptoniphilus, Roseburia, Ruminococcus, Rummeliibacillus, Slackia, Sphingobacterium, Sphingomonas, Streptococcus, Thiomonas, Xenorhabdus</i>	37

Outline

Regression Analysis of Structured Microbiome Data

Differential Network Enrichment Analysis

Conclusion

Chronic Kidney Disease

- ▶ An estimated 37 million adults (15% of the population) in the US have CKD³.
- ▶ CKD leads to marked metabolic changes that can be studied via lipidomics.

Chronic Kidney Disease

- ▶ An estimated 37 million adults (15% of the population) in the US have CKD³.
- ▶ CKD leads to marked metabolic changes that can be studied via lipidomics.

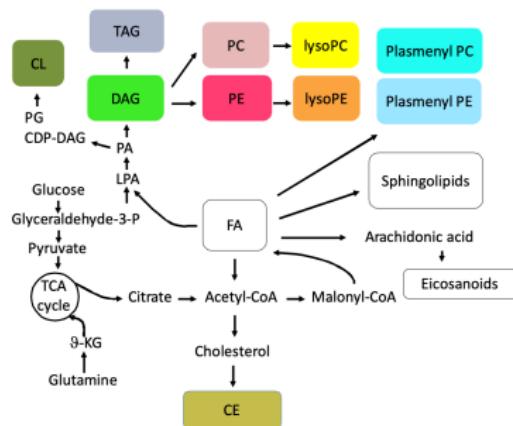


Figure: Overview of lipid biosynthesis from Ma et al. (19') Bioinformatics.

³Johansen et al. (21') American Journal of Kidney Diseases

Data and Study Population

Study cohorts

- ▶ Clinical Phenotyping Resource and Biobank Core (CPROBE⁴): 135 early vs. 79 advanced.
- ▶ Chronic Renal Insufficiency Cohort (CRIC⁵): 121 non-progressors vs. 79 progressors. Baseline only.

Objective

To identify lipidomic signatures of CKD progression.

⁵Afshinnia et al. (18') J. Am. Soc. Nephrol.

⁵Afshinnia et al. (16') Kidney Int. Rep.

Lipids are Correlated!

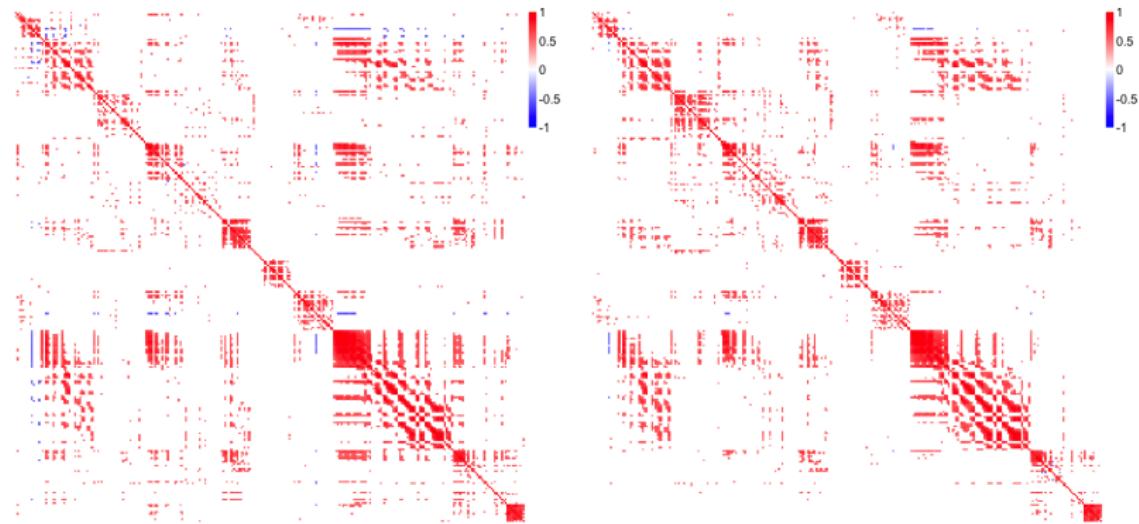


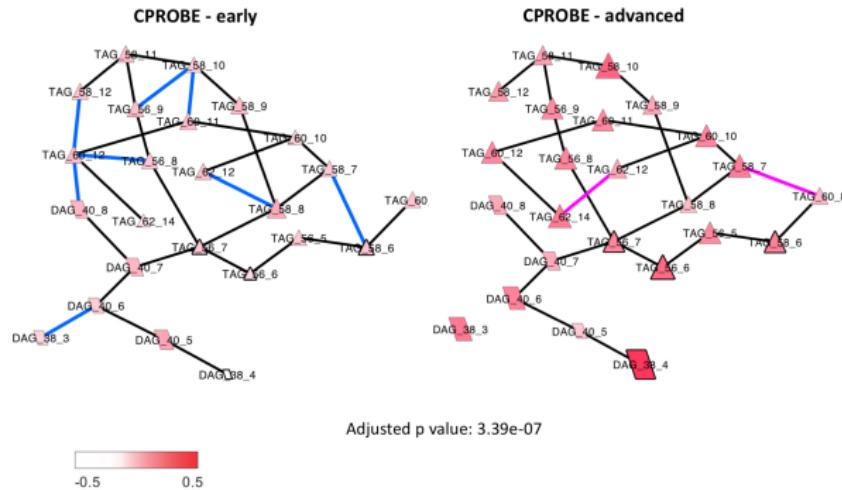
Figure: Heatmap of empirical correlations in non-progressors (left) and progressors (right).
Correlation magnitudes are thresholded at 0.5.

Differential Network Enrichment Analysis

Unlike genetic and genomic data, *no a priori* defined lipid pathways.

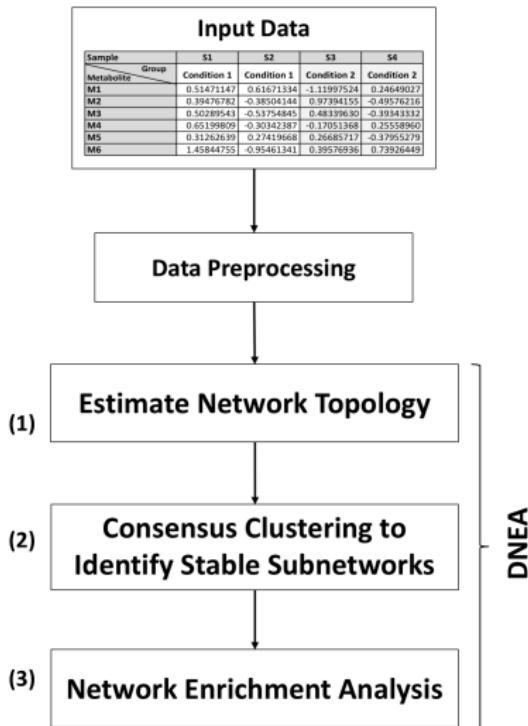
We propose, DNEA⁶, a data-driven approach for subnetwork discovery and enrichment analysis.

Novel Lipid Signatures in CKD



- ▶ Higher abundance of triacylglycerol (TAG) in advanced CKD
- ▶ Fewer edges were associated with advanced CKD

Differential Network Enrichment Analysis



Step I: Estimate Network Topology

We want two lipid partial correlation networks for early and advanced stage, respectively.

Step I: Estimate Network Topology

We want two lipid partial correlation networks for early and advanced stage, respectively.

(:(Sample size is smaller than the number of features (285 lipids).

Step I: Estimate Network Topology

We want two lipid partial correlation networks for early and advanced stage, respectively.

- (:(Sample size is smaller than the number of features (285 lipids).
- (:) Penalized joint estimation strategy to estimate sparse networks⁷:

$$\min_{\Theta_1, \Theta_2} \sum_{k=1}^{K=2} \{ \text{trace}(\hat{\Sigma}_k \Theta_k) - \log \det(\Theta_k) \} + P_\lambda(\Theta_1, \Theta_2).$$

Step I: Estimate Network Topology

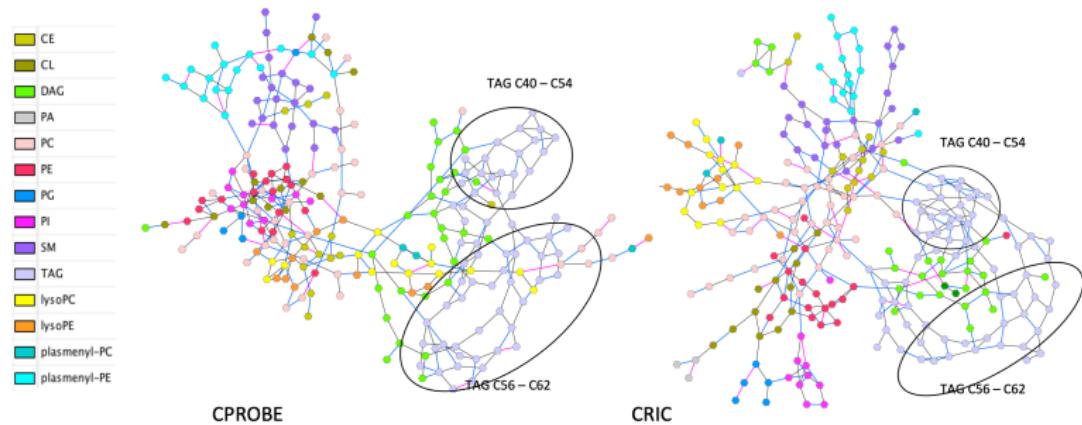
We want two lipid partial correlation networks for early and advanced stage, respectively.

- (:(Sample size is smaller than the number of features (285 lipids).
- (:) Penalized joint estimation strategy to estimate sparse networks⁷:

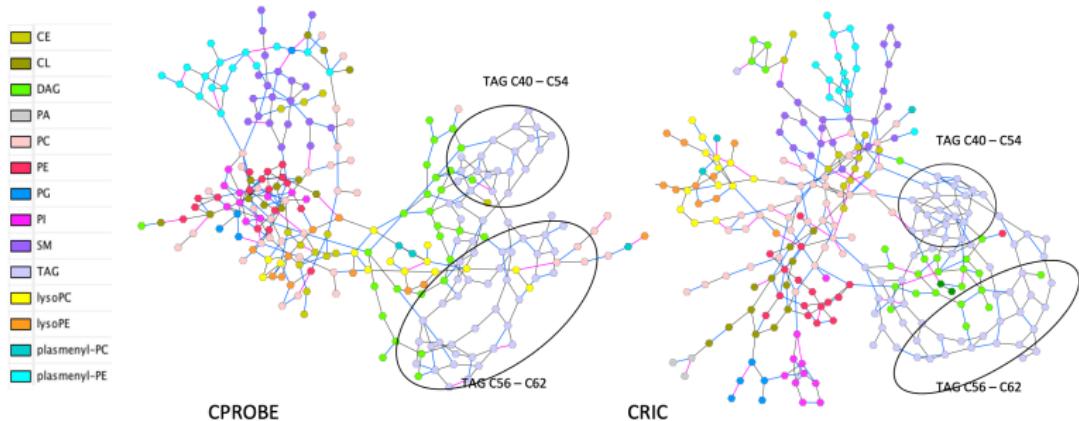
$$\min_{\Theta_1, \Theta_2} \sum_{k=1}^{K=2} \{ \text{trace}(\hat{\Sigma}_k \Theta_k) - \log \det(\Theta_k) \} + P_\lambda(\Theta_1, \Theta_2).$$

Model selection is done via stability selection.

Step II: Consensus Clustering



Step II: Consensus Clustering



Apply consensus clustering to extract stable subnetworks (proxy to pathways)

Step III: Network Enrichment Analysis

Apply NetGSA⁸ to detect enriched subnetworks.

⁹Ma et al. (16') Bioinformatics

⁹Hellstern et al. (21') PLoS Comp Bio; Yue et al. (21') Genet Epidemiol

Step III: Network Enrichment Analysis

Apply NetGSA⁸ to detect enriched subnetworks.

$$Y_j^{(1)} = \Lambda^{(1)}\mu^{(1)} + \Lambda^{(1)}\gamma_j + \epsilon_j, \quad (j = 1, \dots, n_1)$$

$$Y_j^{(2)} = \Lambda^{(2)}\mu^{(1)} + \Lambda^{(2)}\gamma_j + \epsilon_j, \quad (j = n_1 + 1, \dots, n)$$

Obtain significance for each subnetwork $S \subset \{1, \dots, p\}$ by testing the null

$$H_{0,S} : \mathbf{1}^\top \Lambda_{S,S}^{(1)} \mu_S^{(1)} = \mathbf{1}^\top \Lambda_{S,S}^{(2)} \mu_S^{(2)}$$

⁹Ma et al. (16') Bioinformatics

⁹Hellstern et al. (21') PLoS Comp Bio; Yue et al. (21') Genet Epidemiol

Step III: Network Enrichment Analysis

Apply NetGSA⁸ to detect enriched subnetworks.

$$Y_j^{(1)} = \Lambda^{(1)}\mu^{(1)} + \Lambda^{(1)}\gamma_j + \epsilon_j, \quad (j = 1, \dots, n_1)$$

$$Y_j^{(2)} = \Lambda^{(2)}\mu^{(1)} + \Lambda^{(2)}\gamma_j + \epsilon_j, \quad (j = n_1 + 1, \dots, n)$$

Obtain significance for each subnetwork $S \subset \{1, \dots, p\}$ by testing the null

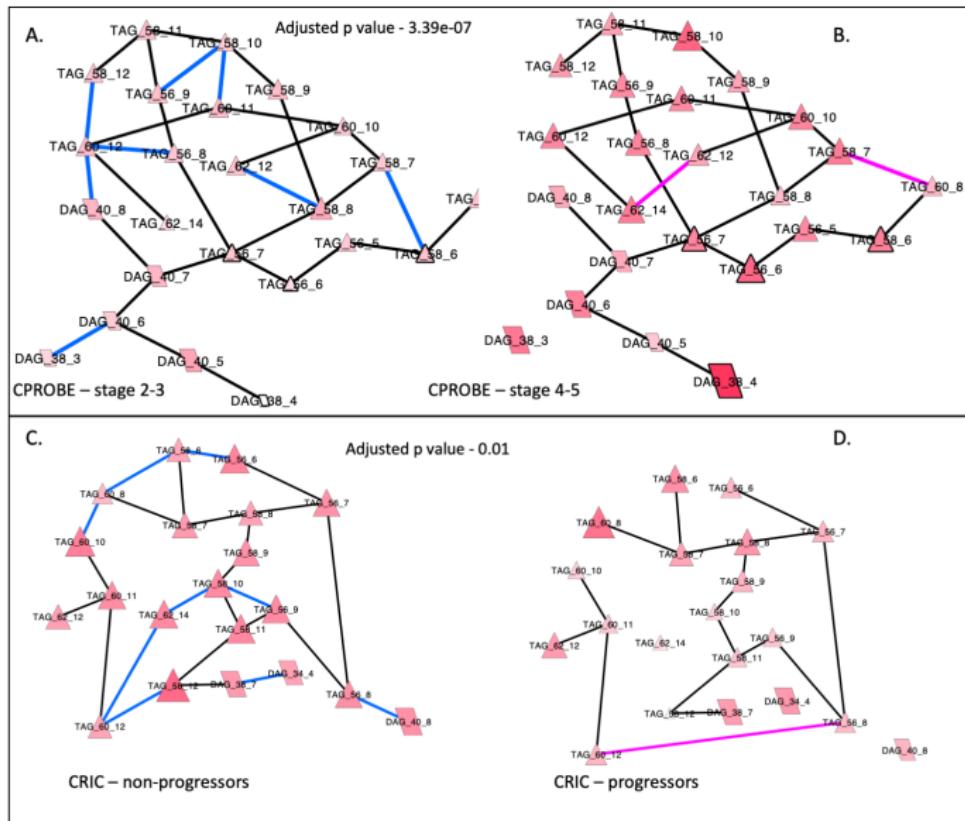
$$H_{0,S} : \mathbf{1}^\top \Lambda_{S,S}^{(1)} \mu_S^{(1)} = \mathbf{1}^\top \Lambda_{S,S}^{(2)} \mu_S^{(2)}$$

NetGSA is now scalable to large datasets⁹.

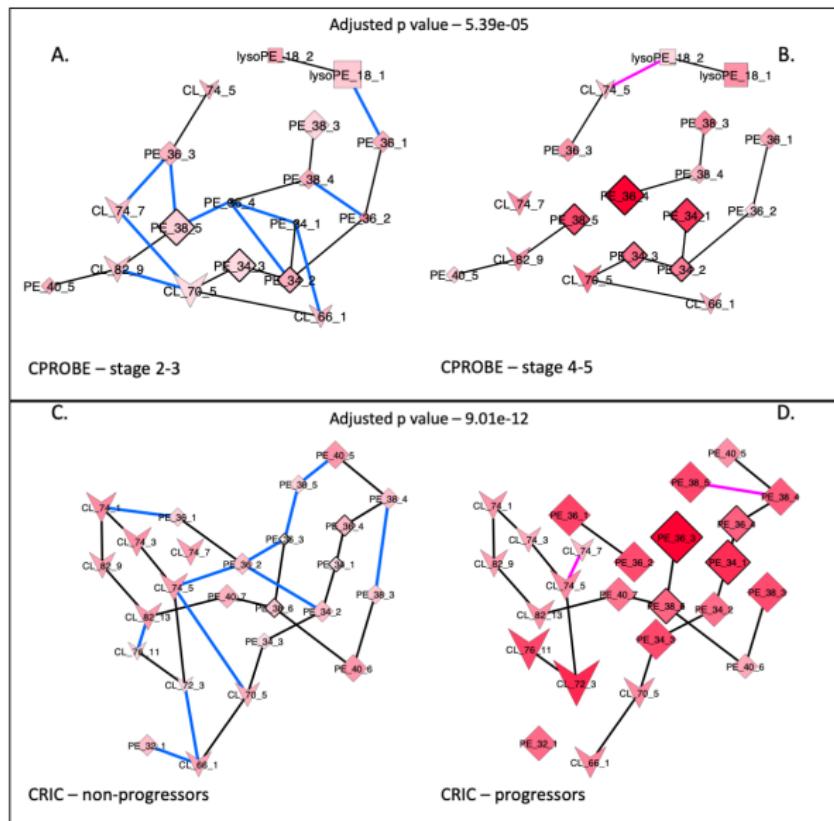
⁹Ma et al. (16') Bioinformatics

⁹Hellstern et al. (21') PLoS Comp Bio; Yue et al. (21') Genet Epidemiol

Lipid Pathways Associated with CKD



Lipid Pathways Associated with CKD



Outline

Regression Analysis of Structured Microbiome Data

Differential Network Enrichment Analysis

Conclusion

Summary

The GMD framework

- ▶ encompasses classical methods, both **unsupervised** (PCA, PCoA/MDS, biplots) and **supervised** (ridge, GLS) methods
- ▶ extends them to non-standard settings (**multi-view data**)

DNEA allows data-driven discovery of pathway signatures of CKD progression.

Acknowledgement

Yue Wang

Kun Yue

Michael Hellstern

Tim Randolph

Ali Shojaie

George Michailidis

Alla Karnovsky

Farsad Afshinnia

- ▶ The GMD-biplot and its application to microbiome data. *mSystems*. 2019
- ▶ Generalized matrix decomposition: estimation and inference for two-way structured data. *arXiv:2104.08408*
- ▶ Differential network enrichment analysis reveals novel lipid pathways in Chronic Kidney Disease. *Bioinformatics*. 2019
- ▶ netgsa: Fast computation and interactive visualization for topology-based pathway enrichment analysis. *PLoS Comp Bio*. 2021
- ▶ REHE fast variance components estimation for linear mixed models. *Genet Epidemiol*. 2021

Thank You!

<http://drjingma.com/>

Simulation Examples

Data generated from linear model with $p = 300$ & $n = 200$

Compare **GMDI-k** and **GMDI-d** with

- ▶ Low-dimension Projection Estimator (LDPE)¹⁰
- ▶ Ridge-based inference¹¹
- ▶ Decorrelated score test (Dscore)¹²
- ▶ Non-sparse high-dimensional inference (Ns-hdi)¹³
- ▶ Grace test¹⁴

⁹van de Geer et al. AoS, 14'; Zhang & Zhang. JRSSB, 14'

¹⁰Bühlmann. Bernoulli, 13'

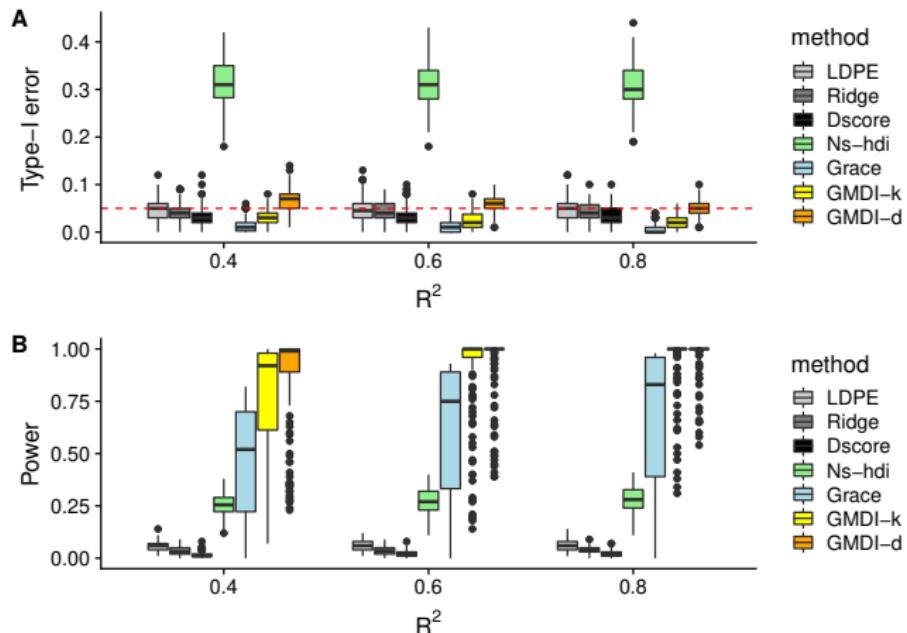
¹¹Ning & Liu. AoS, 17'

¹²Zhu & Bradic. JASA, 18'

¹³Zhao & Shojaie. Biometrics, 15'

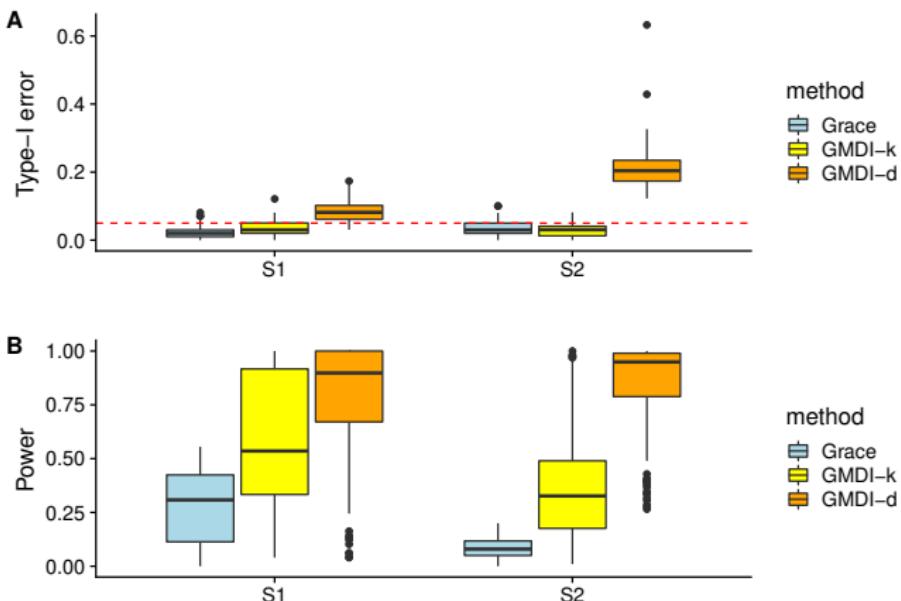
Simulation Example I

► $\mathbf{H} = \mathbf{I}_n$, $\|\mathbf{Q}^{-1/2}\beta^*\|_0 = 10$ but β^* is not sparse.



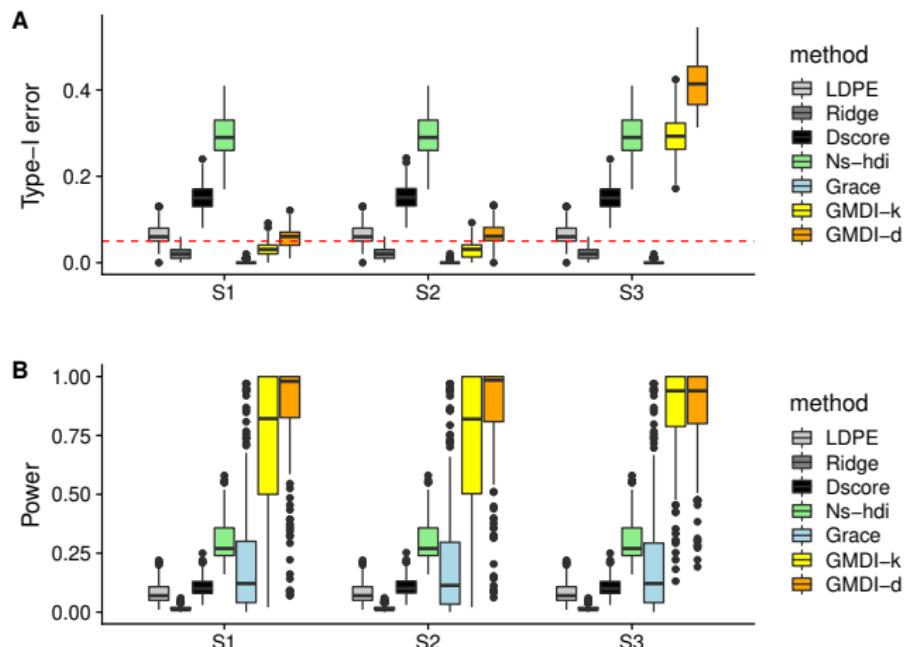
Simulation Example II

- ▶ Similar to Setting 1, but with \mathbf{Q} perturbed. S1: small perturbation; S2: large perturbation.



Simulation Example III

- \mathbf{H} is block diagonal, \mathbf{Q} is the same as in setting 1. S1: correct \mathbf{H} ; S2: perturbed \mathbf{H} ; S3: I_n .



Lipids are Correlated!

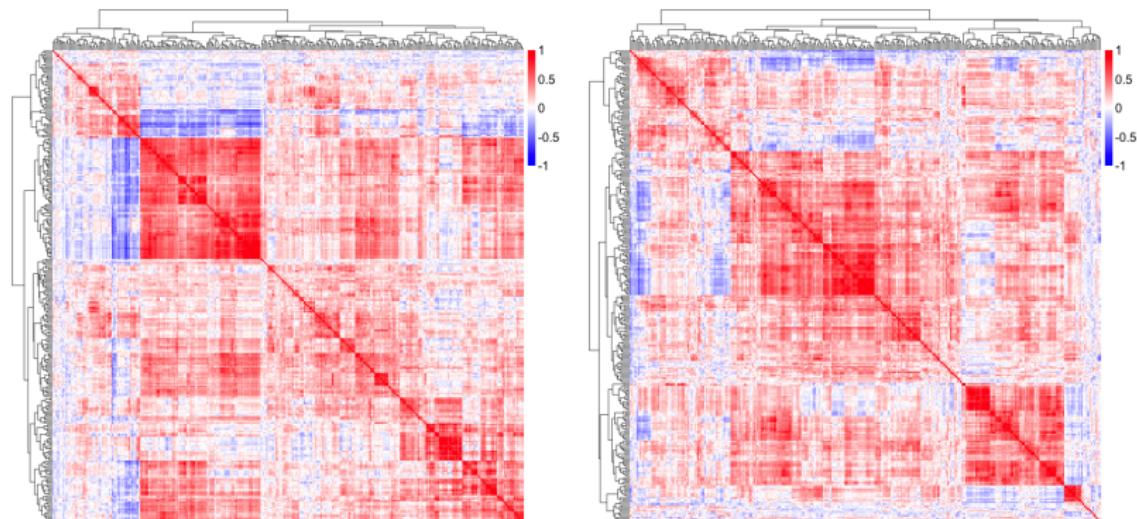


Figure: Heatmap of empirical correlations in non-progressors (left) and progressors (right).

Lipids are Correlated!

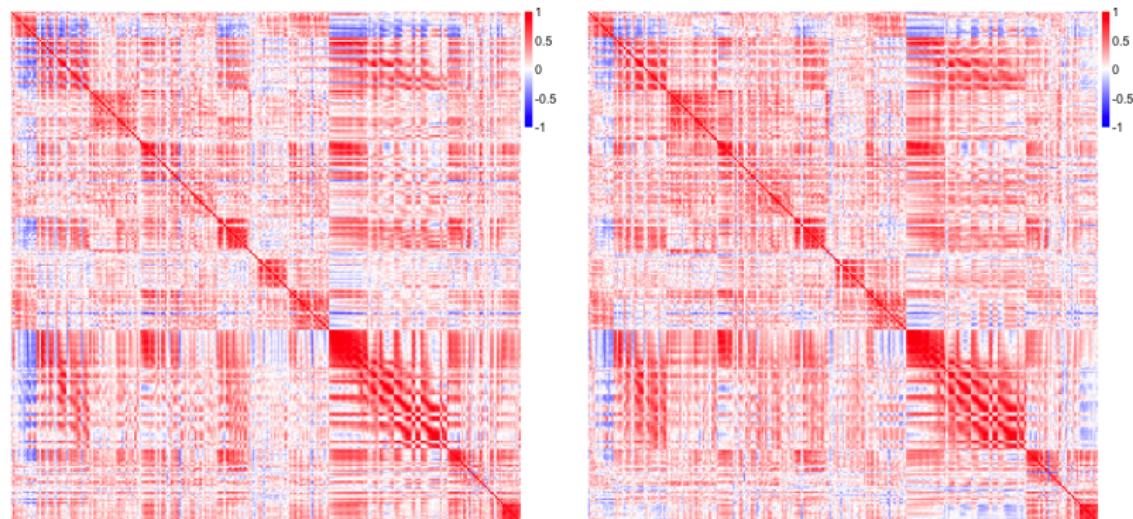


Figure: Heatmap of empirical correlations in non-progressors (left) and progressors (right).