

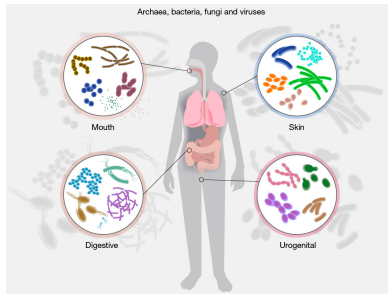
# Network-based Integration of Microbiome and Metabolomic Data

Banff, 17 July 2025

Jing Ma

# The human microbiome

All of the microbes and their genome, mostly bacteria



- More microbial cells than our somatic cells
- More microbial genes than our human genome
- Compositions vary within a person and between individuals
- Highly dynamic yet robust
- Association with many diseases

| OTU | Species | Sample 1 | Sample 2 | Sample 3 |
|-----|---------|----------|----------|----------|
| 1   | E.coli  | 17       | 0        | 335      |
| 2   | S.aurus | 231      | 1180     | 45       |
| 3   | unknown | 30       | 0        | 0        |
| ... | ...     | ...      | ...      | ...      |

Table 1: A typical microbiome contingency table

## Features of microbiome data

- high-dimensional

| OTU | Species | Sample 1 | Sample 2 | Sample 3 |
|-----|---------|----------|----------|----------|
| 1   | E.coli  | 17       | 0        | 335      |
| 2   | S.aurus | 231      | 1180     | 45       |
| 3   | unknown | 30       | 0        | 0        |
| ... | ...     | ...      | ...      | ...      |

Table 1: A typical microbiome contingency table

## Features of microbiome data

- high-dimensional
- sparse: lots of zeros  $\rightarrow$  filtering + pseudocount replacement

| OTU | Species | Sample 1 | Sample 2 | Sample 3 |
|-----|---------|----------|----------|----------|
| 1   | E.coli  | 17       | 0        | 335      |
| 2   | S.aurus | 231      | 1180     | 45       |
| 3   | unknown | 30       | 0        | 0        |
| ... | ...     | ...      | ...      | ...      |

Table 1: A typical microbiome contingency table

## Features of microbiome data

- **high-dimensional**
- **sparse**: lots of zeros → filtering + pseudocount replacement
- **compositional**: only relative abundances are meaningful → normalization

# Metabolomics

Metabolomics is the study of small molecules, known as metabolites, within a biological system.

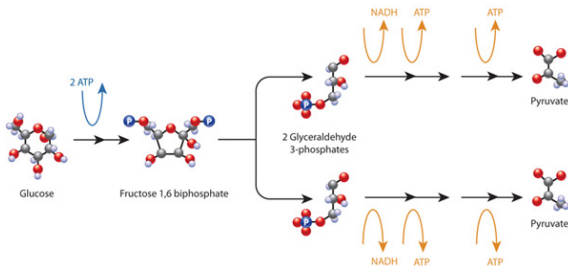
# Metabolomics

Metabolomics is the study of small molecules, known as metabolites, within a biological system.

Metabolites are the building blocks, intermediates, or end products of metabolism.

Metabolomics is the study of small molecules, known as metabolites, within a biological system.

Metabolites are the building blocks, intermediates, or end products of metabolism.



**Figure 1:** Glycolysis: energy is used to convert glucose to a 6 carbon form. Thereafter, energy is generated to create two molecules of pyruvate. Credit to [Nature Education](#).



| Compound      | Sample 1   | Sample 2   | Sample 3   |
|---------------|------------|------------|------------|
| Glucose       | 42,062,493 | 46,507,270 | 48,849,105 |
| Glutamic acid | 1,027,679  | 1,317,161  | 2,527,070  |
| Propionate    | 3,487      | 6,262      | 9,188      |
| ...           | ...        | ...        | ...        |

Table 2: Peak intensities of compounds across samples

## Features of metabolomic data

- high-dimensional

| Compound      | Sample 1   | Sample 2   | Sample 3   |
|---------------|------------|------------|------------|
| Glucose       | 42,062,493 | 46,507,270 | 48,849,105 |
| Glutamic acid | 1,027,679  | 1,317,161  | 2,527,070  |
| Propionate    | 3,487      | 6,262      | 9,188      |
| ...           | ...        | ...        | ...        |

Table 2: Peak intensities of compounds across samples

## Features of metabolomic data

- high-dimensional
- sparse: lots of zeros → filtering + imputation

| Compound      | Sample 1   | Sample 2   | Sample 3   |
|---------------|------------|------------|------------|
| Glucose       | 42,062,493 | 46,507,270 | 48,849,105 |
| Glutamic acid | 1,027,679  | 1,317,161  | 2,527,070  |
| Propionate    | 3,487      | 6,262      | 9,188      |
| ...           | ...        | ...        | ...        |

Table 2: Peak intensities of compounds across samples

## Features of metabolomic data

- high-dimensional
- sparse: lots of zeros → filtering + imputation
- high variance → log transformation

# Microbial metabolites

Microbial metabolites play an important role in host immune system.

Microbial metabolites play an important role in host immune system.

- SCFA: metabolites that are produced by bacteria from dietary components

Microbial metabolites play an important role in host immune system.

- SCFA: metabolites that are produced by bacteria from dietary components
- bile acids: metabolites that are produced by the host and biochemically modified by gut bacteria

Microbial metabolites play an important role in host immune system.

- SCFA: metabolites that are produced by bacteria from dietary components
- bile acids: metabolites that are produced by the host and biochemically modified by gut bacteria
- ATP: metabolites that are synthesized de novo by gut microbes

## Knowledge gap:

- most of the bacterial metabolites remain unidentified.
- many known metabolites have yet to be functionally characterized.

---

<sup>1</sup>Morton et al. 19'. Nat Methods; Reiman et al. 21'. PLOS Comp Bio

<sup>2</sup>Quinn-Bohmann et al. 25'. Nat Microbiol



## Knowledge gap:

- most of the bacterial metabolites remain unidentified.
- many known metabolites have yet to be functionally characterized.

## Current methods for learning microbial–metabolite interactions:

- correlation networks
- machine learning models<sup>1</sup>
- mechanistic models<sup>2</sup>

---

<sup>1</sup>Morton et al. 19'. Nat Methods; Reiman et al. 21'. PLOS Comp Bio

<sup>2</sup>Quinn-Bohmann et al. 25'. Nat Microbiol

Let  $A \in \{0,1\}^{p_1 \times p_2}$  denote the latent network between  $p_1$  microbes and  $p_2$  metabolites.

In this bipartite network, the nodes are microbes/metabolites and the edges represent associations between microbes and metabolites.

Let  $A \in \{0, 1\}^{p_1 \times p_2}$  denote the latent network between  $p_1$  microbes and  $p_2$  metabolites.

In this bipartite network, the nodes are microbes/metabolites and the edges represent associations between microbes and metabolites.

Inference for  $A$  amounts to testing:

$$H_{0,i,j} : A_{i,j} = 0 \quad \text{versus.} \quad H_{1,i,j} : A_{i,j} \neq 0,$$

for all  $1 \leq i \leq p_1, 1 \leq j \leq p_2$ .

# False discovery rate

FDR provides a way of quantifying the statistical significance of multiple hypothesis tests.

|                     | Not significant | Significant | Total |
|---------------------|-----------------|-------------|-------|
| Null is true        | $N_{00}$        | $N_{10}$    | $m_0$ |
| Alternative is true | $N_{01}$        | $N_{11}$    | $m_1$ |
| Total               | $S$             | $R$         | $m$   |

Table 3: Classification of tested hypothesis

$$\text{FDR} = E\left(\frac{N_{10}}{R \vee 1}\right), \quad \text{mFDR} = \frac{E(N_{10})}{E(R)}$$

The Benjamini & Hochberg (BH) procedure<sup>3</sup>

- Choose a desired significance level  $\alpha \in (0, 1)$ .
- Sort the p-values in increasing order:  $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(m)}$ .
- Find the largest index  $k$  such that:

$$p_{(k)} \leq \frac{k}{m} \alpha$$

- Reject all hypotheses with p-values  $p_i \leq p_{(k)}$ .

---

<sup>3</sup>Benjamini & Hochberg. (95') JRSSB

# Local false discovery rate

Let  $X = (x_i) \in \mathbb{R}^{p_1 p_2}$  denote the observations (e.g., z-scores). Efron et al.<sup>4</sup> studied the mixture model

$$f(x) = \pi_0 f_0(x) + \pi_1 f_1(x),$$

for multiple testing, where

- $f_0$ : null distribution (e.g.,  $\mathcal{N}(0, 1)$ )
- $f_1$ : the alternative distribution
- $\pi_0, \pi_1 = 1 - \pi_0$ : prior probabilities

---

<sup>4</sup>Efron et al. (01') JASA

<sup>5</sup>Sun and Cai. (07') JASA

# Local false discovery rate

Let  $X = (x_i) \in \mathbb{R}^{p_1 p_2}$  denote the observations (e.g., z-scores). Efron et al.<sup>4</sup> studied the mixture model

$$f(x) = \pi_0 f_0(x) + \pi_1 f_1(x),$$

for multiple testing, where

- $f_0$ : null distribution (e.g.,  $\mathcal{N}(0, 1)$ )
- $f_1$ : the alternative distribution
- $\pi_0, \pi_1 = 1 - \pi_0$ : prior probabilities

The Empirical Bayes **local false discovery rate**<sup>5</sup> is:

$$\text{IFDR}(x) = \frac{\pi_0 f_0(x)}{f(x)}.$$

Reject hypotheses with  $\text{IFDR}(x) < \alpha$ .

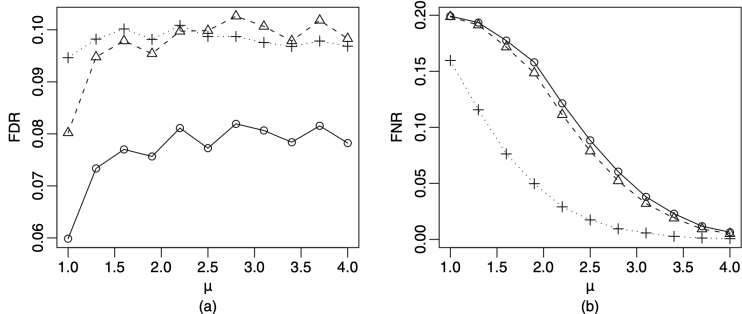
---

<sup>4</sup>Efron et al. (01') JASA

<sup>5</sup>Sun and Cai. (07') JASA

# Incorporating structures

More power can be achieved by exploiting the local dependence structure, e.g., Hidden Markov models<sup>6</sup>.



**Fig. 1.** Comparison of BH (○), AP (△) and OR (+) in an HMM (the FDR level is set at 0.10): (a) FDR versus  $\mu$ ; (b) FNR versus  $\mu$ .

<sup>6</sup>Sun and Cai (09') JRSSB



# Knowledge gap

- Nearly all existing works on large scale multiple testing require vector inputs and are not optimized for **matrix-valued observations**.

---

<sup>7</sup>Cai and Liu (16') JASA

# Knowledge gap

- Nearly all existing works on large scale multiple testing require vector inputs and are not optimized for **matrix-valued observations**.
- Methods for correlation testing exist<sup>7</sup>, but they do not take into account the **structural information in the data**.

---

<sup>7</sup>Cai and Liu (16') JASA

- Nearly all existing works on large scale multiple testing require vector inputs and are not optimized for **matrix-valued observations**.
- Methods for correlation testing exist<sup>7</sup>, but they do not take into account the **structural information in the data**.

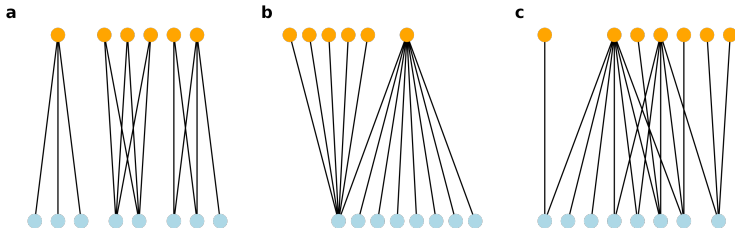


Figure 2: Topology of interest: (a) three biclusters, (b) fully nested graph, and (c) preferential attachment.

---

<sup>7</sup>Cai and Liu (16') JASA

# Latent graph model

Let  $X = (x_{i,j}) \in \mathbb{R}^{p_1 \times p_2}$  denote the matrix-valued observations (e.g., z-scores).

Let  $X = (x_{i,j}) \in \mathbb{R}^{p_1 \times p_2}$  denote the matrix-valued observations (e.g., z-scores).

We model  $X$  using a latent bipartite stochastic block model (biSBM), defined with respect to row clustering  $Z_1 = (Z_{i,1})$  and column clustering  $Z_2 = (Z_{j,2})$ . For  $i = 1, \dots, p_1$  and  $j = 1, \dots, p_2$ :

$$\begin{aligned} Z_{i,1} &\sim \text{Multi}(1, \alpha_1), \\ Z_{j,2} &\sim \text{Multi}(1, \alpha_2), \\ A_{i,j} \mid Z_1, Z_2 &\sim \text{Bern}(\pi_{Z_{i,1}, Z_{j,2}}), \\ x_{i,j} \mid A_{i,j}, Z_1, Z_2 &\sim A_{i,j} g_{\nu_{Z_{i,1}, Z_{j,2}}} + (1 - A_{i,j}) g_{0, \nu_0}. \end{aligned} \tag{1}$$

Model parameters are  $\theta = (\alpha_1, \alpha_2, \pi, \nu, \nu_0)$ .

# Multiple testing procedure

Ideally, we need

$$P(A_{i,j} = 0 \mid X)$$

to control the IFDR. However, they are intractable in our context.

# Multiple testing procedure

Ideally, we need

$$P(A_{i,j} = 0 \mid X)$$

to control the IFDR. However, they are intractable in our context.

Instead, we use the *structured  $\ell$ -value*

$$P(A_{i,j} = 0 \mid X, Z_1, Z_2).$$

Reject the hypotheses if the  $\ell$ -value is small, where the threshold is chosen to control mFDR.

Ideally, we need

$$P(A_{i,j} = 0 \mid X)$$

to control the IFDR. However, they are intractable in our context.

Instead, we use the *structured  $\ell$ -value*

$$P(A_{i,j} = 0 \mid X, Z_1, Z_2).$$

Reject the hypotheses if the  $\ell$ -value is small, where the threshold is chosen to control mFDR.

The *structured  $\ell$  values* provide much more information than a single observation  $x_{i,j}$  and will considerably help to make the final decision.



The *Gaussian* noisy biSBM is identifiable under the constraint that all elements of  $\{(0, \sigma_0), (\mu_{q,l}, \sigma_{q,l}), 1 \leq q \leq B_1, 1 \leq l \leq B_2\}$  are distinct.

The *Gaussian* noisy biSBM is identifiable under the constraint that all elements of  $\{(0, \sigma_0), (\mu_{q,l}, \sigma_{q,l}), 1 \leq q \leq B_1, 1 \leq l \leq B_2\}$  are distinct.

In general, the requirement of all elements being distinct is not necessary, as the model is also identifiable if there is a single alternative distribution such that  $\sigma_{q,l} = \sigma$  and  $\mu_{q,l} = \mu$ .

Need to use the EM algorithm due to the latent variables  $A, Z_1, Z_2$ . Let  $Q$  denote a probability distribution of the latent variables.

$$\log \mathcal{L}(X; \theta) = \underbrace{E_Q[\log \mathcal{L}(X, A, Z_1, Z_2; \theta)] + \mathcal{H}(Q)}_{ELBO} + KL(Q \| P_{A, Z_1, Z_2 | X; \theta}), \quad (2)$$

- ① Initialize  $\theta^{(0)}$ .
- ② E-step: evaluate the expectation in (2) with respect to  $Q = P_{A, Z_1, Z_2 | X; \theta^{(t)}}$ .
  - Mean-field approximation
- ③ M-step: update  $\theta^{(t+1)}$  by maximizing the ELBO.
- ④ Iterate between E- and M-step until convergence.

# Selecting the number of clusters

The numbers of blocks  $B_1, B_2$  are unknown.

# Selecting the number of clusters

The numbers of blocks  $B_1, B_2$  are unknown.

We use the integrated classification likelihood (ICL) criterion to select the optimal  $B_1$  and  $B_2$ , allowing them to be different.

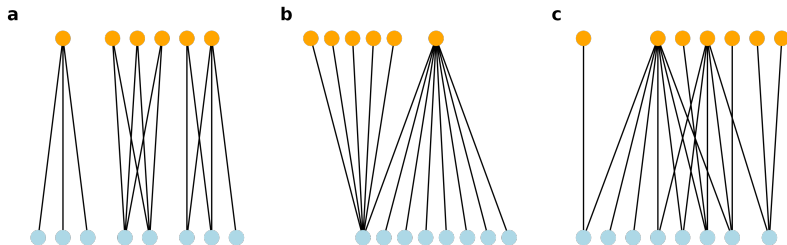


Figure 3: Illustrations of the latent bipartite network used in simulations: (a) three biclusters, (b) fully nested graph, and (c) preferential attachment

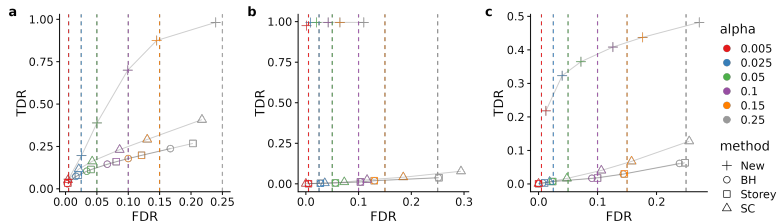
- (a) Modular structures
- (b) Nested graph in ecology: generalist vs specialist
- (c) Preferential attachment: the rich gets richer

- $p_1 = 150, p_2 = 200$
- $\mathcal{N}(0, 1)$  versus  $\mathcal{N}(2, 1)$
- Compare the average performance over 100 simulations
- Both the new and the SC procedures<sup>8</sup> were implemented assuming known null density.

---

<sup>8</sup>Sun and Cai (07') JASA

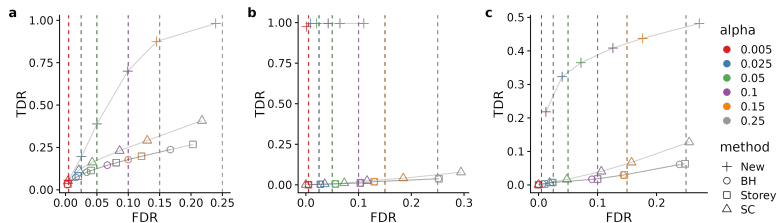
# Simulation results



**Figure 4:** Plot of the empirical (FDR, TDR) as a function of the nominal level  $\alpha$  for the new procedure, BH, Storey's  $q$ -value, and the SC procedure. Dashed lines indicate the nominal level  $\alpha$ .



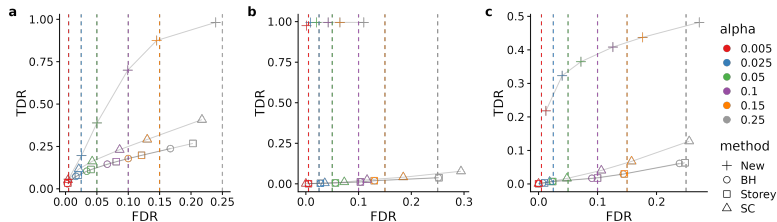
# Simulation results



**Figure 4:** Plot of the empirical (FDR, TDR) as a function of the nominal level  $\alpha$  for the new procedure, BH, Storey's  $q$ -value, and the SC procedure. Dashed lines indicate the nominal level  $\alpha$ .

(a) ICL selected a model with three biclusters correctly in 90% of the simulations.

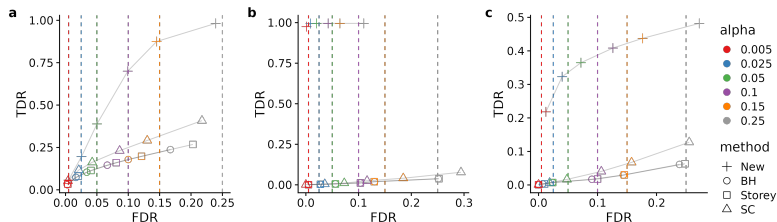
# Simulation results



**Figure 4:** Plot of the empirical (FDR, TDR) as a function of the nominal level  $\alpha$  for the new procedure, BH, Storey's  $q$ -value, and the SC procedure. Dashed lines indicate the nominal level  $\alpha$ .

- (a) ICL selected a model with three biclusters correctly in 90% of the simulations.
- (b) ICL selected two biclusters in 83% of the simulations.

# Simulation results



**Figure 4:** Plot of the empirical (FDR, TDR) as a function of the nominal level  $\alpha$  for the new procedure, BH, Storey's  $q$ -value, and the SC procedure. Dashed lines indicate the nominal level  $\alpha$ .

- (a) ICL selected a model with three biclusters correctly in 90% of the simulations.
- (b) ICL selected two biclusters in 83% of the simulations.
- (c) ICL mostly found three clusters in type I vertex and one cluster in type II vertex.

- BV is the most common vaginal condition, affecting an estimated 30% of women at any given time<sup>9</sup>.
- BV is associated with increased transmission of HIV and STIs as well as increased risk of preterm labour.
- Diagnosis relies on microscopy to identify BV-like bacteria by morphology alone (Nugent Scoring).
- The pathogenesis of BV remains unclear.

---

<sup>9</sup>McMillan et al. (15') Scientific Reports

- $Y_1 \in R^{n \times p_1}$ : relative abundances from 49 genera obtained from 16S; centered log ratio transformed while keep zeros unchanged.

- $Y_1 \in R^{n \times p_1}$ : relative abundances from 49 genera obtained from 16S; centered log ratio transformed while keep zeros unchanged.
- $Y_2 \in R^{n \times p_2}$ : concentrations of 128 metabolites from GC-MS; log transformed

- $Y_1 \in R^{n \times p_1}$ : relative abundances from 49 genera obtained from 16S; centered log ratio transformed while keep zeros unchanged.
- $Y_2 \in R^{n \times p_2}$ : concentrations of 128 metabolites from GC-MS; log transformed
- **Outcome**: normal ( $n = 79$ ) vs BV ( $n = 45$ ) status defined by nugent score

- $Y_1 \in R^{n \times p_1}$ : relative abundances from 49 genera obtained from 16S; centered log ratio transformed while keep zeros unchanged.
- $Y_2 \in R^{n \times p_2}$ : concentrations of 128 metabolites from GC-MS; log transformed
- **Outcome**: normal ( $n = 79$ ) vs BV ( $n = 45$ ) status defined by nugent score
- **Aim** is to understand the microbial functional changes during BV.



Let  $\hat{Y}_1$  and  $\hat{Y}_2$  denote, respectively, the standardized data. The sample correlation is defined by

$$\hat{\rho}_{i,j} = \frac{1}{n} \sum_{k=1}^n \hat{Y}_{1,k,i} \hat{Y}_{2,k,j}.$$

Let  $\hat{Y}_1$  and  $\hat{Y}_2$  denote, respectively, the standardized data. The sample correlation is defined by

$$\hat{\rho}_{i,j} = \frac{1}{n} \sum_{k=1}^n \hat{Y}_{1,k,i} \hat{Y}_{2,k,j}.$$

Let

$$s_{i,j} = \frac{1}{n} \sum_{i=1}^n (2\hat{Y}_{1,k,i} \hat{Y}_{2,k,j} - \hat{\rho}_{k,i} \hat{Y}_{1,k,i} - \hat{\rho}_{k,i} \hat{Y}_{2,k,j})^2.$$

Let  $\hat{Y}_1$  and  $\hat{Y}_2$  denote, respectively, the standardized data. The sample correlation is defined by

$$\hat{\rho}_{i,j} = \frac{1}{n} \sum_{k=1}^n \hat{Y}_{1,k,i} \hat{Y}_{2,k,j}.$$

Let

$$s_{i,j} = \frac{1}{n} \sum_{i=1}^n (2\hat{Y}_{1,k,i} \hat{Y}_{2,k,j} - \hat{\rho}_{k,i} \hat{Y}_{1,k,i} - \hat{\rho}_{k,i} \hat{Y}_{2,k,j})^2.$$

The test statistic

$$x_{i,j} = \frac{2\hat{\rho}_{i,j}}{\sqrt{s_{i,j}/n}} \rightarrow \mathcal{N}(0, 1),$$

under finite fourth moment condition<sup>10</sup>.

---

<sup>10</sup>Cai and Liu (16') JASA

Two-sample inference comparing BV to normal patients:

$$x_{i,j} = \frac{2(\hat{\rho}_{i,j}^{(1)} - \hat{\rho}_{i,j}^{(2)})}{\sqrt{s_{i,j}^{(1)}/n_1 + s_{i,j}^{(2)}/n_2}}$$

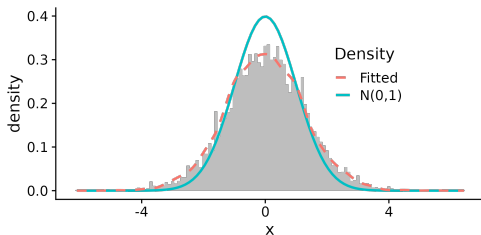
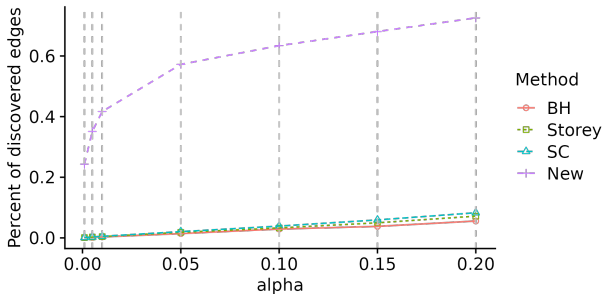
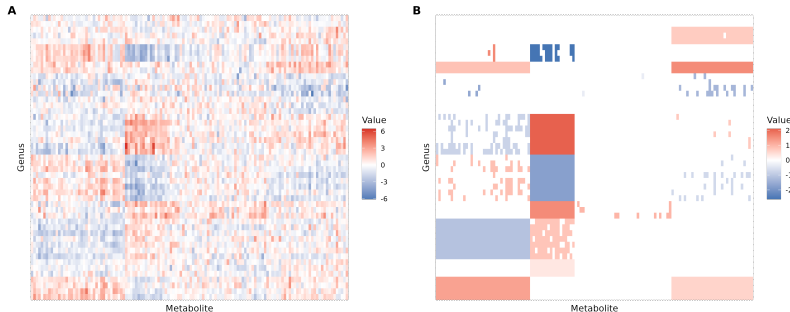


Figure 5: Histogram of observed z-scores compared to the standard normal and the estimated marginal distribution by the proposed approach.



**Figure 6:** Percent of rejected edges as a function of the significance level  $\alpha$  for the different procedures.

The new procedure groups microbes and metabolites with similar association patterns into biclusters.



**Figure 7:** Heat map of the data (A) compared to the estimated graph by the proposed approach at  $\alpha = 0.1\%$  (B). Rows and columns are ordered by the inferred clustering.

# Results

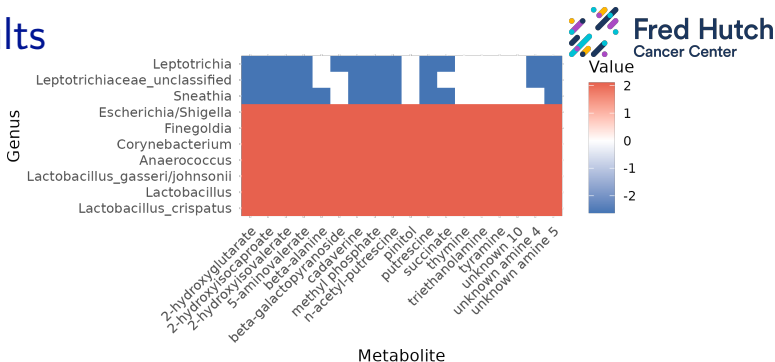


Figure 8: A zoom in view of two biclusters with the largest mean difference.

- Top bicluster consists of *Leptotrichia* and *Sneathia* which are emerging pathogens implicated in BV. Association of these genera with metabolites is higher in BV patients.
- Bottom bicluster consists of *Lactobacillus* species, important for keeping a healthy vaginal microbiome. Association of these genera with metabolites is higher in normal individuals.
- Shed light on uncharacterized metabolites through “Guilt by Association”.

# Summary

## metaMint

- preprint: [arXiv.2506.12275](https://arxiv.org/abs/2506.12275)
- Software: <https://github.com/drjingma/metaMint>

## Future directions

- Degree heterogeneity is not accounted for.
- Computation: the method is slower than existing methods.  
Alternative model estimation and/or selection is helpful.