

### A framework for multi-view analysis of microbiome data

Jing Ma

Fred Hutchinson Cancer Research Center Public Health Sciences Division

2 July 2019

#### Collaborators





Yue Wang



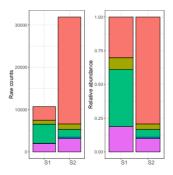
Ali Shojaie



Tim Randolph

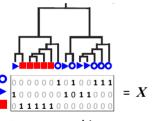


- $\mathbf{X} = (x_{ij})_{n \times p}$  matrix of microbiome data for *n* samples and *p* taxa
- ▶ Due to sample differences, often work with relative abundances





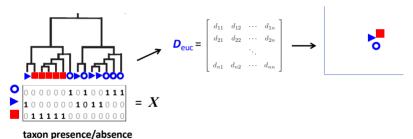
- ► Similarities among samples better captured by phylogenetic tree
- ▶ Many methods for capturing phylogenetic distances, e.g. UniFrac dist.



taxon presence/absence

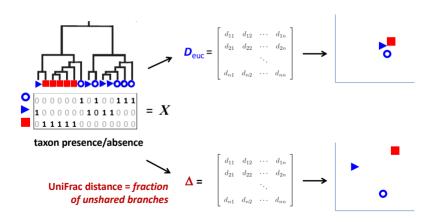


- ► Similarities among samples better captured by phylogenetic tree
- ▶ Many methods for capturing phylogenetic distances, e.g. UniFrac dist.



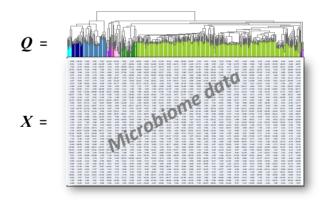


- ► Similarities among samples better captured by phylogenetic tree
- ▶ Many methods for capturing phylogenetic distances, e.g. UniFrac dist.



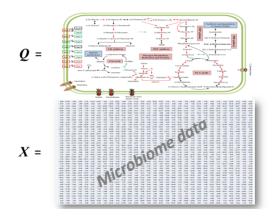


- ► Similarities among samples better captured by phylogenetic tree
- ▶ The phylogenetic tree also captures similarities among taxa.





- ► Similarities among samples better captured by phylogenetic tree
- ▶ Alternatively, can consider information from metabolic pathways.



#### Microbiome Data Analysis



- ▶ Need to capture
  - ► Similarities among samples non-Euclidean (e.g. UniFrac distance)
  - ► Similarities among taxa phylogenetic tree, pathway information, etc.

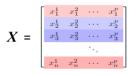
#### Microbiome Data Analysis



- ▶ Need to capture
  - ► Similarities among samples non-Euclidean (e.g. UniFrac distance)
  - ► Similarities among taxa phylogenetic tree, pathway information, etc.
- Often use exploratory data analysis tools
  - PCoA (aka MDS)
  - ► DPCoA (Double PCoA)

#### Exploratory Analysis using PCoA First recall PCA

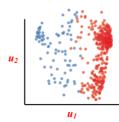




n samples (rows)p variables (cols)

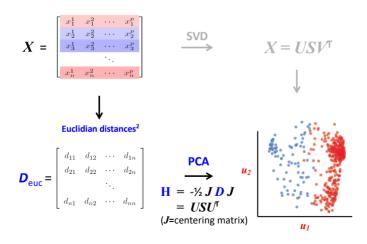
 $\xrightarrow{\mathsf{SVD}}$ 

$$X = USV^{\mathsf{T}}$$



#### Exploratory Analysis using PCoA First recall PCA

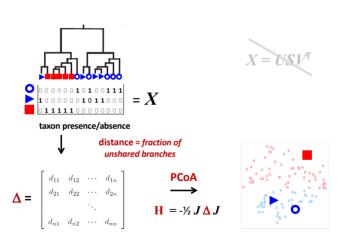




#### Exploratory Analysis using PCoA



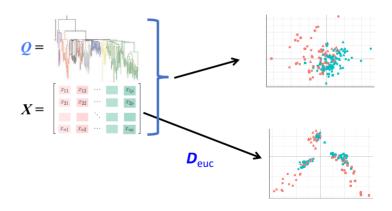
#### UniFrac PCoA



### Exploratory Analysis using DPCoA



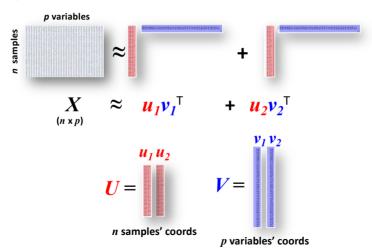
DPCoA<sup>1</sup> based on phylogenetic tree



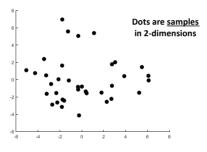
<sup>&</sup>lt;sup>1</sup>Pavoine et al. 2004; Purdom. AOAS, 2011





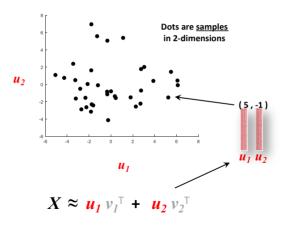




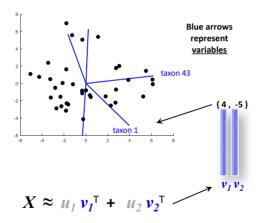


$$X \approx \boldsymbol{u}_1 \boldsymbol{v}_1^{\mathsf{T}} + \boldsymbol{u}_2 \boldsymbol{v}_2^{\mathsf{T}}$$





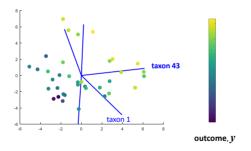






Biplots provide simultaneous visualization of samples and variables

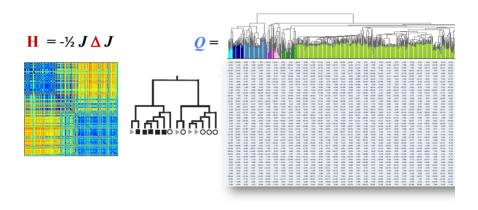
► Can also overlay outcome



$$X \approx u_1 v_1^{\mathsf{T}} + u_2 v_2^{\mathsf{T}}$$

### Microbiome Data are Doubly Structured







Can we draw a biplot that accounts for the external structures?

<sup>&</sup>lt;sup>2</sup>Satten et al. PLOS One, 2017



Can we draw a biplot that accounts for the external structures?

- Unfortunately, there is no biplot with PCoA because
  - ▶ PCA uses X and gives:  $X = USV^{T}$  (SVD)
  - ▶ PCoA uses  $\triangle$  and only gives:  $US^2U^T$  (no V)

<sup>&</sup>lt;sup>2</sup>Satten et al. PLOS One, 2017



Can we draw a biplot that accounts for the external structures?

- Unfortunately, there is no biplot with PCoA because
  - ▶ PCA uses X and gives:  $X = USV^{T}$  (SVD)
  - ▶ PCoA uses  $\triangle$  and only gives:  $US^2U^T$  (no V)
- ► Existing approaches<sup>2</sup> in the field are approximate/add hoc.

<sup>&</sup>lt;sup>2</sup>Satten et al. PLOS One, 2017



Can we draw a biplot that accounts for the external structures?

▶ SVD gives  $X = USV^{T}$  by solving

$$\mathop{\mathsf{arg\,min}}_{oldsymbol{U},oldsymbol{S},oldsymbol{V}} \|oldsymbol{X} - oldsymbol{U}oldsymbol{S}oldsymbol{V}^\intercal\|_F$$

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

<sup>&</sup>lt;sup>3</sup>Allen et al. JASA, 2014



Can we draw a biplot that accounts for the external structures?

▶ SVD gives  $X = USV^{T}$  by solving

$$\mathop{\mathsf{arg\,min}}_{oldsymbol{U},oldsymbol{S},oldsymbol{V}} \|oldsymbol{X} - oldsymbol{U}oldsymbol{S}oldsymbol{V}^\intercal\|_F$$

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

Consider instead a general norm to incorporate H and Q:

$$\|\boldsymbol{X} - \boldsymbol{U}\boldsymbol{S}\boldsymbol{V}^{\intercal}\|_{\boldsymbol{H}, \boldsymbol{Q}}$$

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

<sup>&</sup>lt;sup>3</sup>Allen et al. JASA, 2014



Can we draw a biplot that accounts for the external structures?

▶ SVD gives  $X = USV^{T}$  by solving

$$\mathop{\mathsf{arg\;min}}_{oldsymbol{\mathcal{U}},oldsymbol{\mathcal{S}},oldsymbol{\mathcal{V}}} \|oldsymbol{\mathcal{X}} - oldsymbol{\mathcal{U}}oldsymbol{\mathcal{S}}oldsymbol{\mathcal{V}}^\intercal\|_F$$

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

Consider instead a general norm to incorporate H and Q:

$$\|\boldsymbol{X} - \boldsymbol{U}\boldsymbol{S}\boldsymbol{V}^{\intercal}\|_{\boldsymbol{H}, \boldsymbol{Q}}$$

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

► The <u>GMD</u> (Gen'zd Matrix Decomp<sup>3</sup>) gives  $\mathbf{X} = \mathcal{USV}^{\mathsf{T}}$  such that  $\mathcal{U}^{\mathsf{T}} \mathbf{H} \mathcal{U} = \mathcal{V}^{\mathsf{T}} \mathbf{Q} \mathcal{V} = I_K$ , and  $\mathbf{S}$  is the diagonal matrix of GMD values.

<sup>&</sup>lt;sup>3</sup>Allen et al. JASA, 2014

#### Variance Explained



SVD gives

$$X = USV^{\mathsf{T}}$$

**GMD** gives

$$X = \mathcal{U}\mathcal{S}\mathcal{V}^{\mathsf{T}}$$

S and S are diagonal matrices of weights:

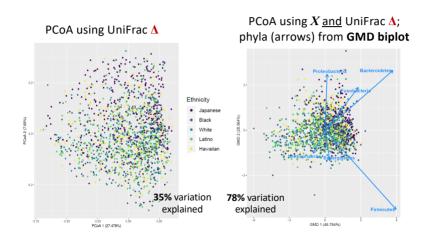
$$m{S} = egin{pmatrix} \sigma_1 & & & & \ & \sigma_2 & & \ & & \sigma_3 \end{pmatrix} \quad ext{(decreasing: } \sigma_1 > \sigma_2 > \sigma_3 ext{)}$$

 $\Rightarrow$  These quantify "variance explained" by components 1, 2, 3, ...

#### **GMD** Biplot



The GMD-biplot displays samples and variables using columns of  $\mathcal U$  and  $\mathcal V$ 



### Supervised Learning with GMD: GMDR



- GMD generalizes SVD for doubly structured data
- ► Can thus use GMD for supervised learning, similar to PCR

## Principal Component Regression



- ▶ Linear model  $y = X\beta + \varepsilon$
- Using two dimensions

$$y \approx \gamma_1 \underline{u_1} + \gamma_2 \underline{u_2} + \varepsilon = \underline{U}_{\scriptscriptstyle (2)} \underline{S}_{\scriptscriptstyle (2)} \underline{V}_{\scriptscriptstyle (2)}^{\mathsf{T}} \beta + \varepsilon$$

Coefficient

$$\hat{eta}_{PCR} = m{V}_{_{(2)}} m{S}_{_{(2)}}^{-1} m{U}_{_{(2)}}^{\mathsf{T}} y$$

### Ridge Regression



- ▶ Linear model  $y = X\beta + \varepsilon$
- ► Using all *p* dimensions

$$y = USV^{\mathsf{T}}\beta + \varepsilon$$

Coefficient

$$\hat{\beta}_{Ridge} = \mathbf{V} \mathbf{W}_{\lambda} \mathbf{S}^{-1} \mathbf{U}^{\mathsf{T}} \mathbf{y},$$

where  $W_{\lambda}$  is a diagonal matrix of weights and  $\lambda$  a tuning parameter.

### **GMD** Regression



- ▶ Linear model  $y = X\beta + \varepsilon$
- Incorporating H and Q

$$y = \mathcal{U}\mathcal{S}\mathcal{V}^{\mathsf{T}}\beta + \varepsilon$$

Coefficient

$$\hat{\beta}_{GMD} = \mathbf{Q} \mathcal{V} \mathcal{W} \mathcal{S}^{-1} \mathcal{U}^{\mathsf{T}} \mathbf{H} y,$$

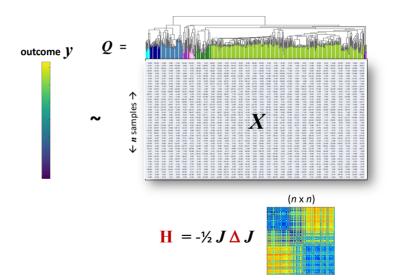
where  $\mathcal{W}$  is a diagonal matrix of weights:

- $\blacktriangleright \ \mathcal{W}_i = \mathbf{1}_{i \in \mathcal{J}} 
  ightarrow \hat{eta}_{GMDR}(\mathcal{J}), \mathcal{J} \subset \{1, \dots, p\}$
- $\mathcal{W} = \mathcal{S}^2 (\mathcal{S}^2 + \lambda I_n)^{-1} \rightarrow \hat{\beta}_{KPR} = \arg\min_{\beta} \{ \|y \boldsymbol{X}\beta\|_{\boldsymbol{H}}^2 + \lambda \|\beta\|_{\boldsymbol{Q}^{-1}}^2 \}^4$

<sup>&</sup>lt;sup>4</sup>Randolph et al. AOAS, 2018

### Summary





#### Summary



- Discussed visualization and regression based on GMD
- ► GMD is closely related to the duality<sup>5</sup> between viewing the data from the perspective of samples and variables

$$\begin{array}{c} \mathbb{R}^{p} \xleftarrow{X^{\top}} \mathbb{R}^{n^{*}} \\ \mathbb{Q} \downarrow & \uparrow \\ \mathbb{R}^{p^{*}} \xrightarrow{X} \mathbb{R}^{n} \end{array}$$

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

$$\begin{array}{cccc} \mathbb{R}^{p} \xleftarrow{X^{\top}} \mathbb{R}^{n^{*}} & \xrightarrow{Z^{\top}} \mathbb{R}^{q} \\ Q \downarrow & \uparrow H & \downarrow R \\ \mathbb{R}^{p^{*}} & \xrightarrow{X} \mathbb{R}^{n} \xleftarrow{Z} \mathbb{R}^{q^{*}} \end{array}$$

<sup>&</sup>lt;sup>5</sup>Escoufier, 1977: de la Cruz and Holmes, AOAS, 2011