

# Distances, Trees, and Types: Methods for Analyzing High-Dimensional Microbiome Data

---

Brendan J. Kelly, MD, MSCE  
Division of Infectious Diseases  
University of Pennsylvania  
2 June 2022

# Outline

- ▷ The problem: too much data.
- ▷ Reducing dimensions:
  - richness, evenness, and diversity
  - ecological distances (UniFrac)
  - PCA & PCoA
  - PERMANOVA (adonis)
- ▷ Cluster analysis:
  - supervised vs unsupervised learning
  - Dirichlet multinomial mixture modeling

# Outline

- ▷ The problem: too much data.
- ▷ Reducing dimensions:
  - richness, evenness, and diversity
  - ecological distances (UniFrac)
  - PCA & PCoA
  - PERMANOVA (adonis)
- ▷ Cluster analysis:
  - supervised vs unsupervised learning
  - Dirichlet multinomial mixture modeling

# Generating Microbiome Data

- ▷ High-density (next-generation, high-throughput) sequencing:
  - “tag” gene with conserved and variable regions (16S, 18S, ITS)
  - “shotgun” metagenomics (pool of randomly amplified nucleic acid)
- ▷ Sequence binning and assignment:
  - operational taxonomic units (OTUs) based on 97% sequence similarity  
→ taxonomic assignment of OTUs
  - assemble contiguous metagenomic sequences → taxonomy
  - unassembled reads → taxonomic assignment

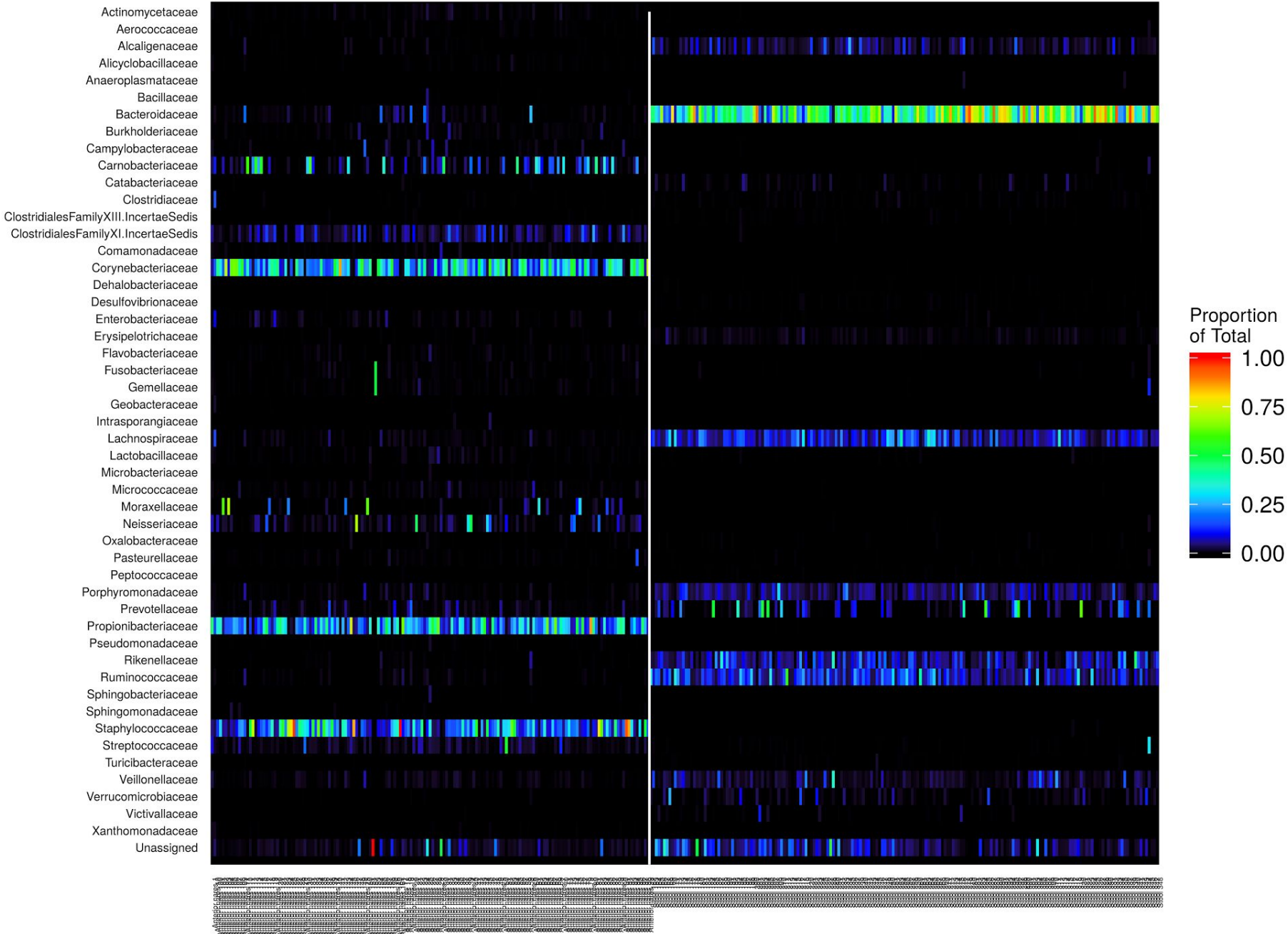
[illegible]

[illegible]

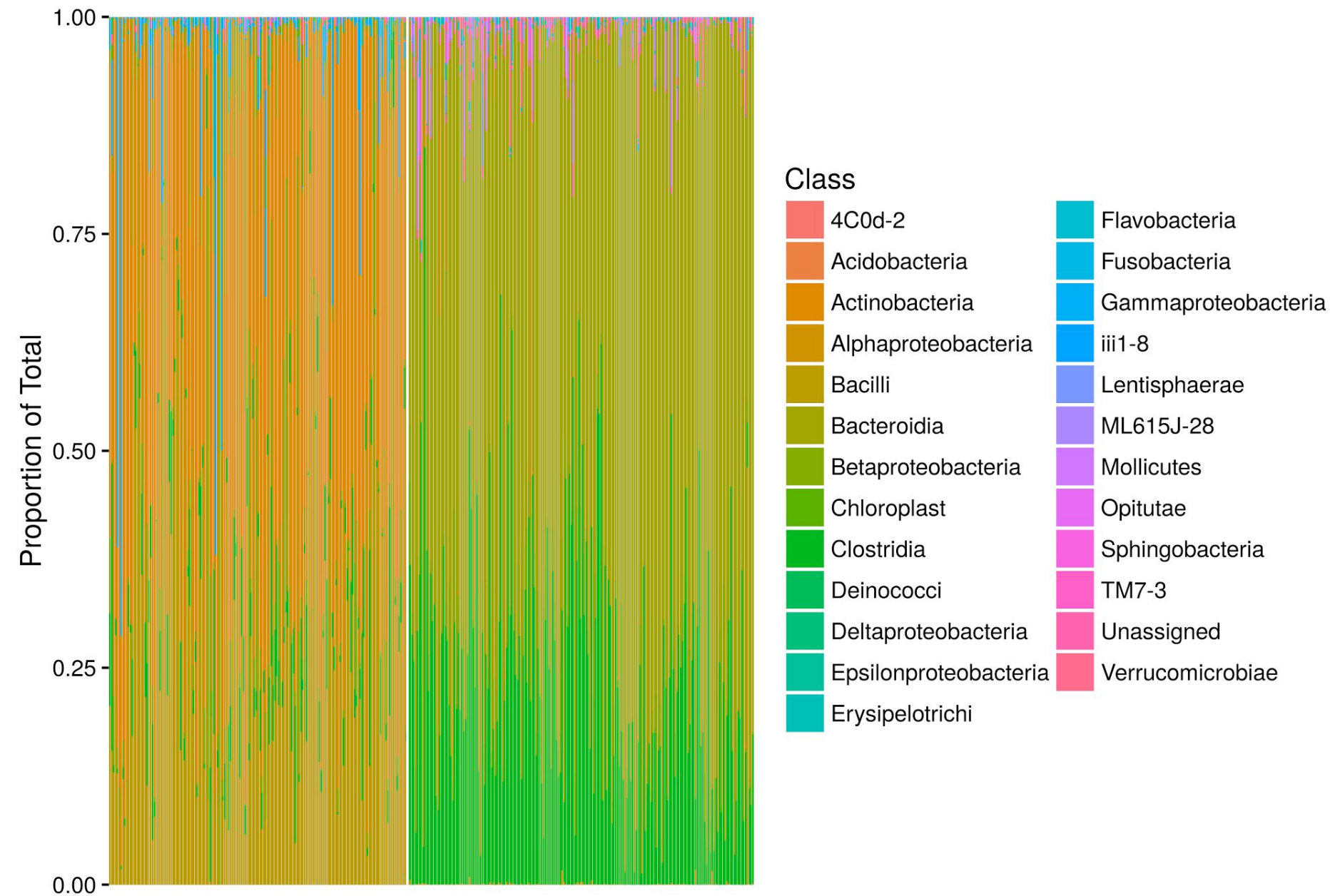
# High-Dimensional Microbiome Data

- ▷ Descriptive:
  - heatmaps
  - stacked barplots
- ▷ Test *a priori* hypotheses regarding specific OTUs/taxa.
- ▷ Reduce dimensions:
  - single summary statistic (alpha diversity)
  - pairwise distances (beta diversity) with PCoA or PERMANOVA
  - community types (mixture modeling)

## Anterior Nares vs Stool







# Descriptive: Heatmaps & Barplots

- ▷ Visualization of OTU table:
  - typically present counts as a proportion of sample total
  - choice of sample order can highlight group differences
- ▷ Limitations:
  - cannot depict full list of OTUs
  - space dictates taxonomic level presented

# Single-Taxon Hypotheses

- ▷ You suspect *Bacteroides* has a relationship with outcome of interest...
  - *Bacteroides* (genus)?
  - *Bacteroidaceae* (family)?
  - *Bacteroidales* (order)?
  - *Bacteroidetes* (class)?
- ▷ Hypotheses focusing on specific taxa often fail to account for possibility of selection bias from culture.

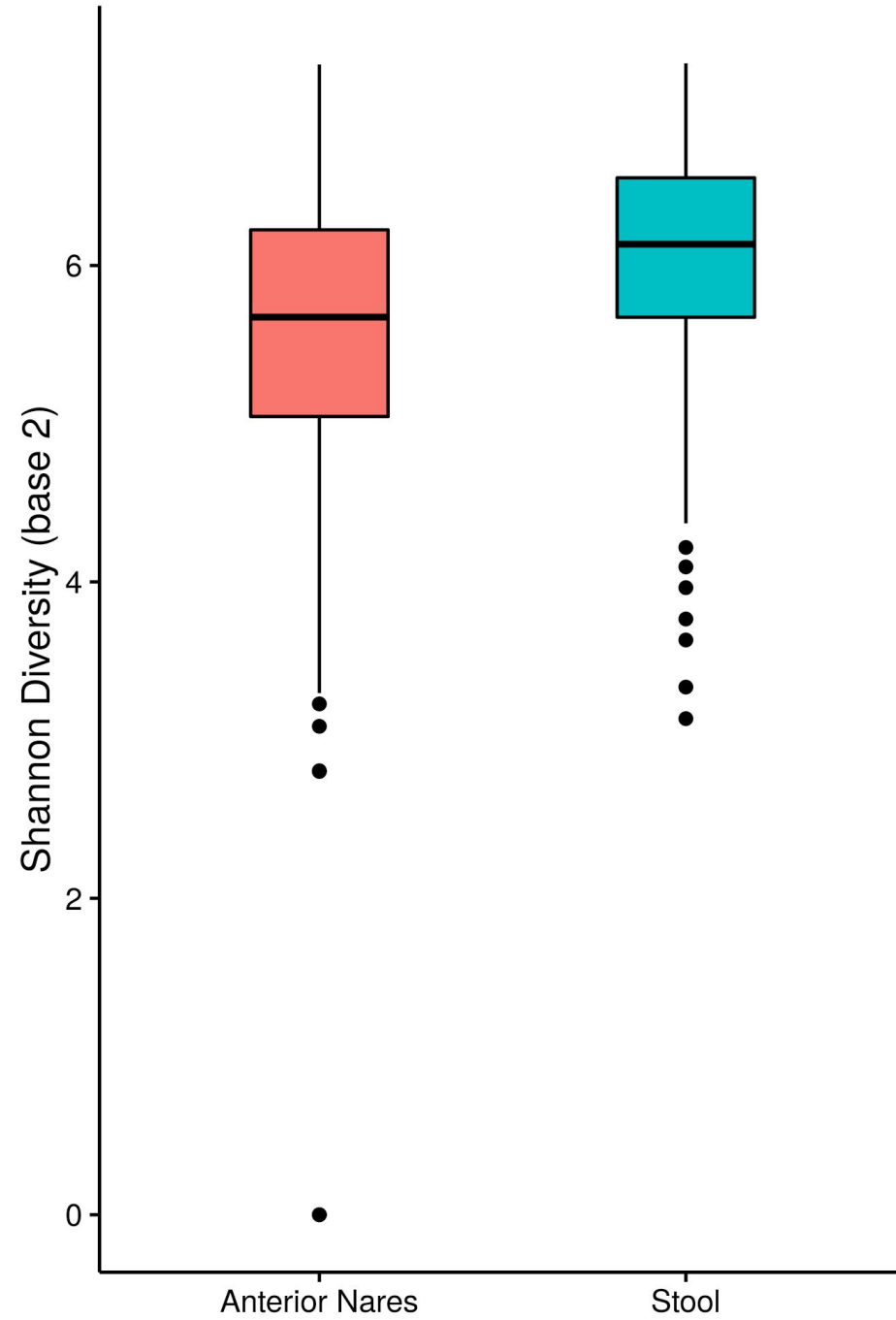
OTU	assignment	phylum	class	order	family	genus
OTU_97.1	Root;p__Firmicutes;c__Bacilli;o__Lactobacillales;f__Lactobacillaceae;g__Lactobacillus	Firmicutes	Bacilli	Lactobacillales	Lactobacillaceae	Lactobacillus
OTU_97.10	Root;p__Firmicutes;c__Clostridia;o__Clostridiales;f__Veillonellaceae;g__Veillonella	Firmicutes	Clostridia	Clostridiales	Veillonellaceae	Veillonella
OTU_97.100	Root;p__Firmicutes;c__Clostridia;o__Clostridiales;f__Veillonellaceae	Firmicutes	Clostridia	Clostridiales	Veillonellaceae	NA
OTU_97.1000	Root;p__Proteobacteria;c__Betaproteobacteria	Proteobacteria	Betaproteobacteria	NA	NA	NA
OTU_97.10000	Root;p__Actinobacteria;c__Actinobacteria;o__Actinomycetales;f__Propionibacteriaceae;g__Propionibacterium	Actinobacteria	Actinobacteria	Actinomycetales	Propionibacteriaceae	Propionibacterium
OTU_97.10001	Root;p__Firmicutes;c__Bacilli;o__Lactobacillales;f__Lactobacillaceae;g__Lactobacillus	Firmicutes	Bacilli	Lactobacillales	Lactobacillaceae	Lactobacillus
OTU_97.10002	Root;p__Firmicutes;c__Clostridia;o__Clostridiales;f__Lachnospiraceae	Firmicutes	Clostridia	Clostridiales	Lachnospiraceae	NA
OTU_97.10003	Root;p__Actinobacteria;c__Actinobacteria;o__Actinomycetales;f__Propionibacteriaceae;g__Propionibacterium	Actinobacteria	Actinobacteria	Actinomycetales	Propionibacteriaceae	Propionibacterium
OTU_97.10004	Root;p__Actinobacteria;c__Actinobacteria;o__Actinomycetales;f__Corynebacteriaceae;g__Corynebacterium	Actinobacteria	Actinobacteria	Actinomycetales	Corynebacteriaceae	Corynebacterium
OTU_97.10005	Root;p__Bacteroidetes;c__Bacteroidia;o__Bacteroidales;f__Bacteroidaceae;g__Bacteroides	Bacteroidetes	Bacteroidia	Bacteroidales	Bacteroidaceae	Bacteroides

# Outline

- ▷ The problem: too much data.
- ▷ Reducing dimensions:
  - richness, evenness, and diversity
  - ecological distances (UniFrac)
  - PCA & PCoA
  - PERMANOVA (adonis)
- ▷ Cluster analysis:
  - supervised vs unsupervised learning
  - Dirichlet multinomial mixture modeling

# Dimension Reduction: Alpha Diversity

- ▷ Summarize each sample's community in a single measure:
  - richness: number of community members
  - evenness: the distribution of member counts
- ▷ Many alpha diversity metrics (weight richness/evenness):
  - species number, Chao1 (singletons & doubletons)
  - Shannon diversity:  $H' = -\sum_i p_i \log_b p_i$
  - (note: may measure similarity or dissimilarity)



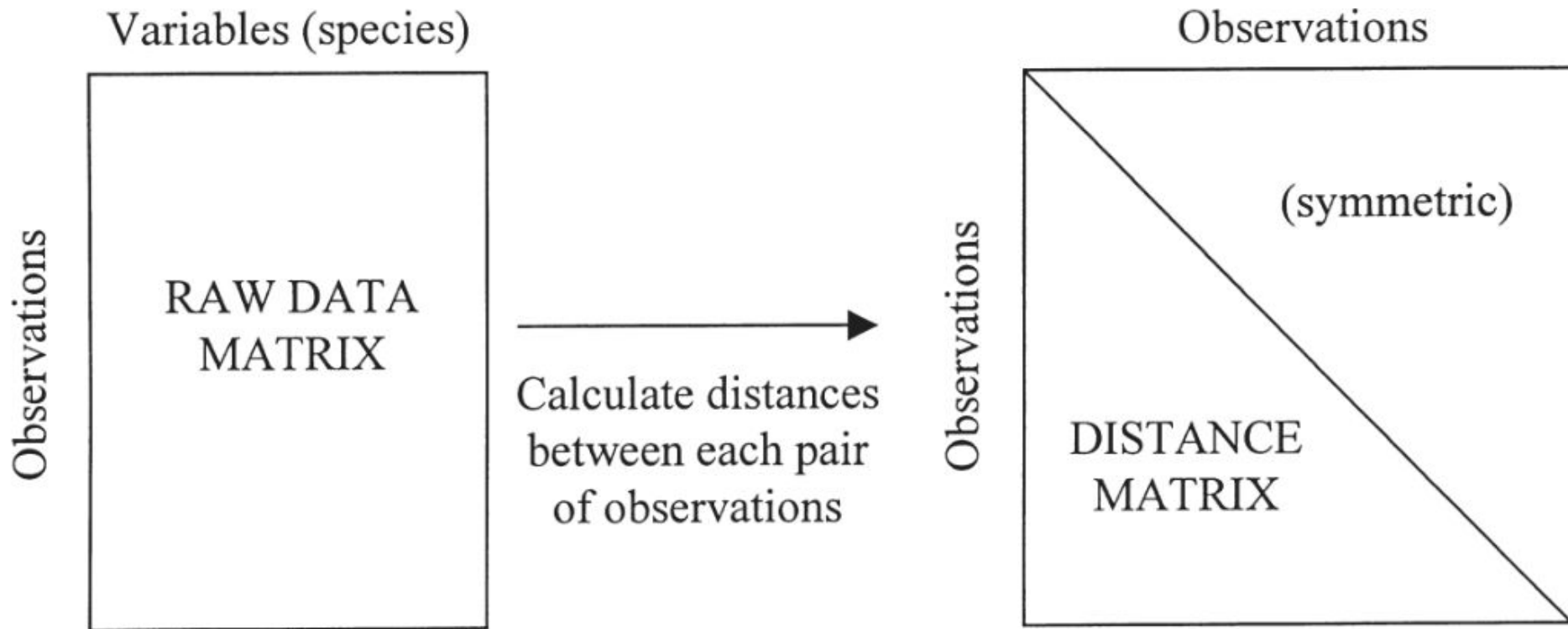
# Outline

- ▷ The problem: too much data.
- ▷ Reducing dimensions:
  - richness, evenness, and diversity
  - ecological distances (UniFrac)
  - PCA & PCoA
  - PERMANOVA (adonis)
- ▷ Cluster analysis:
  - supervised vs unsupervised learning
  - Dirichlet multinomial mixture modeling



# Dimension Reduction: Beta Diversity

- ▷ Summarize each sample's relationship to other samples:
  - pairwise distances
  - OTU table → square matrix
- ▷ Many beta diversity metrics:
  - just counts versus counts + phylogeny
  - weighted versus unweighted
  - (Euclidean versus non-Euclidean)



# Dimension Reduction: Beta Diversity

▷ Just counts versus counts + phylogeny:

- Jaccard:  $J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$

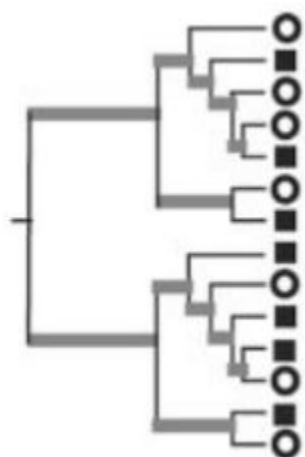
$$d_J(A, B) = 1 - J(A, B) = \frac{|A \cup B| - |A \cap B|}{|A \cup B|}$$

- UniFrac: fraction of unique branch length in tree

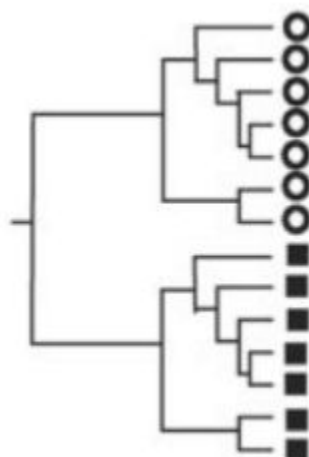
▷ Weighted versus unweighted:

- weighted: counts matter
- unweighted: binary (presence-absence)

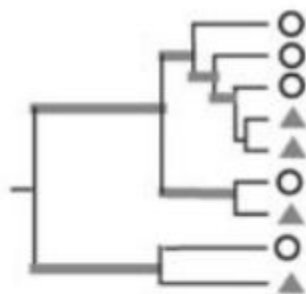
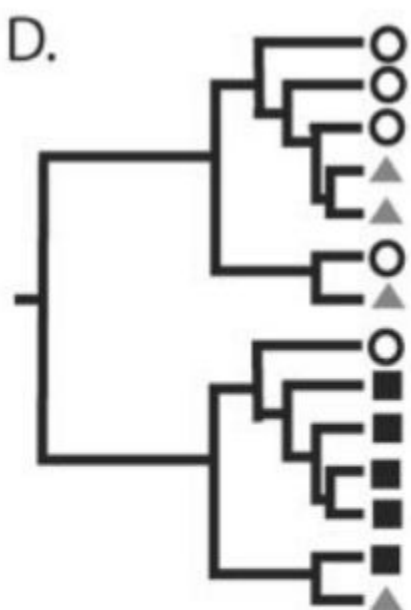
A.



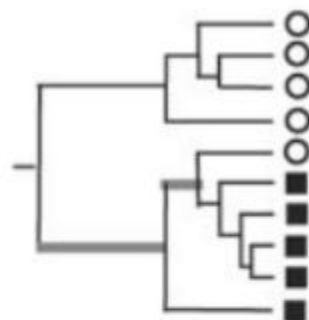
B.



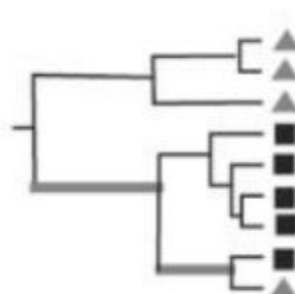
D.



Triangle vs Circle



Square vs Circle



Triangle vs Square

	O	▲	■
O	0	.3	.7
▲	.3	0	.6
■	.7	.6	0

Distance Matrix

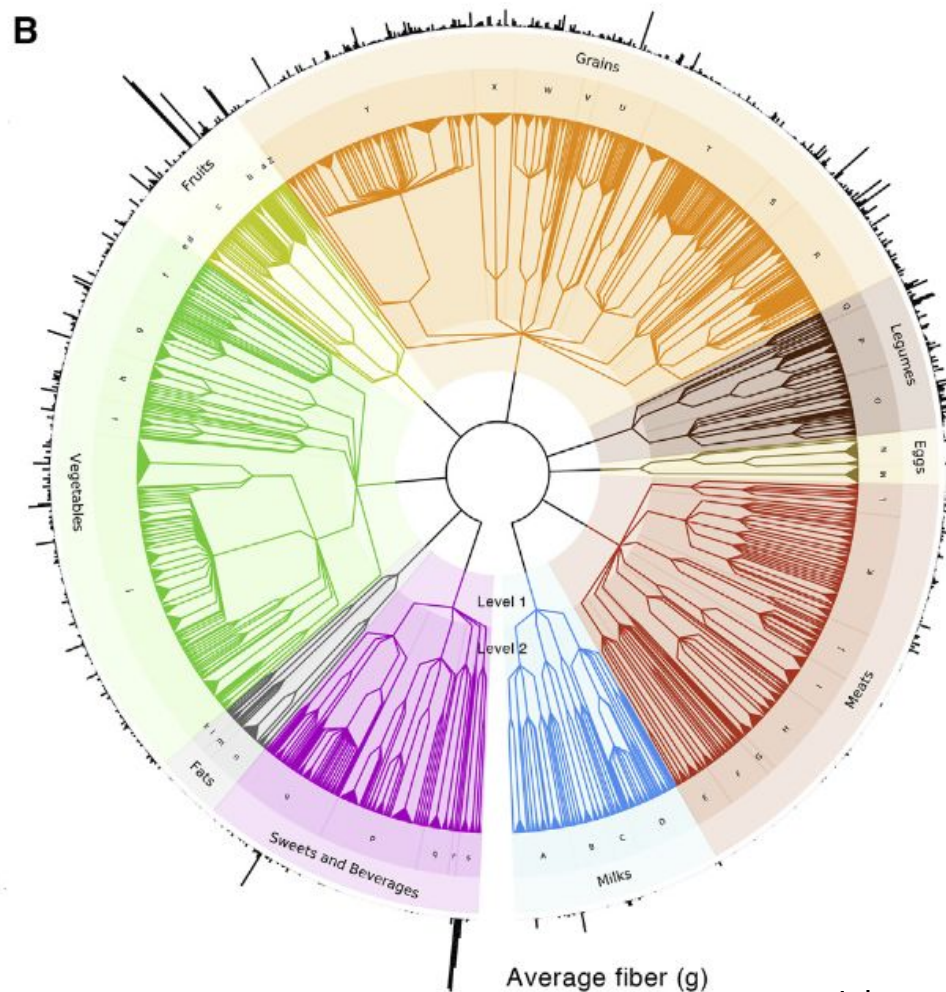


Cluster of environments

# Beta Diversity: How to Choose?

- ▷ Why use Jaccard? UniFrac?
- ▷ Why use weighted? Unweighted?

# Thinking Like a Tree



# Outline

- ▷ The problem: too much data.
- ▷ Reducing dimensions:
  - richness, evenness, and diversity
  - ecological distances (UniFrac)
  - PCA & PCoA
  - PERMANOVA (adonis)
- ▷ Cluster analysis:
  - supervised vs unsupervised learning
  - Dirichlet multinomial mixture modeling

# Original Descriptors: PCA

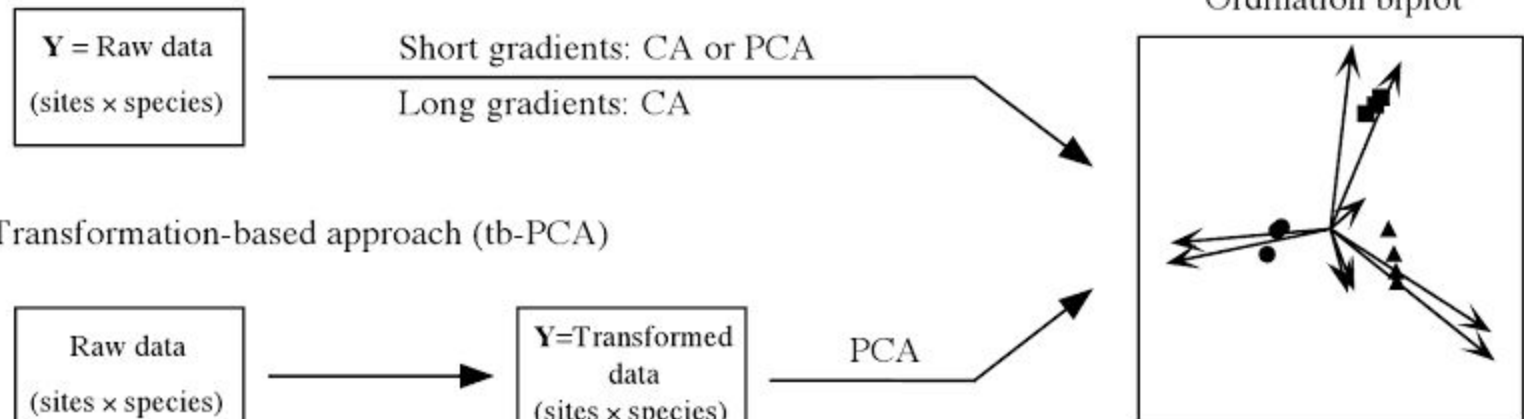
- ▷ PCA: principal component analysis
  - rigid rotation for successive directions of maximum variance
  - lots of restrictions (Euclidean)
  - but allows projection of original descriptors in PCA space



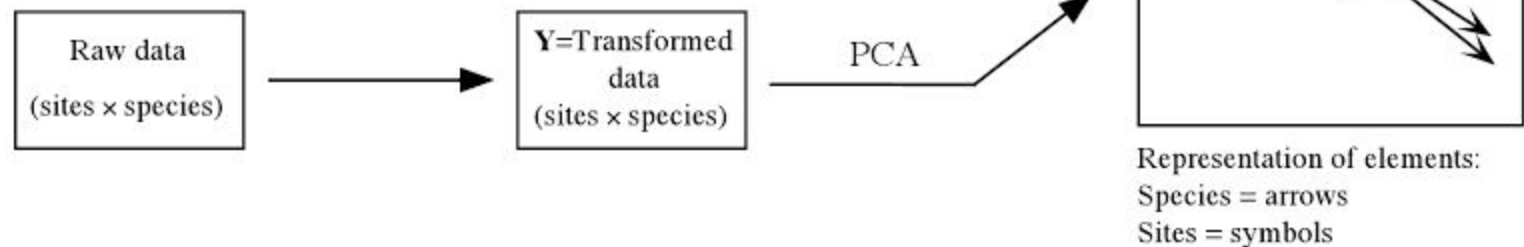
# Pairwise Distances: PCoA

- ▷ PCoA: principal coordinate analysis
  - any metric distance, even if non-Euclidean
  - like PCA, eigenvalue decomposition (maximum variance) but mediated by distance function (no original descriptors)
  - unlike PCA, does not allow projection of original descriptors in reduced-dimension space

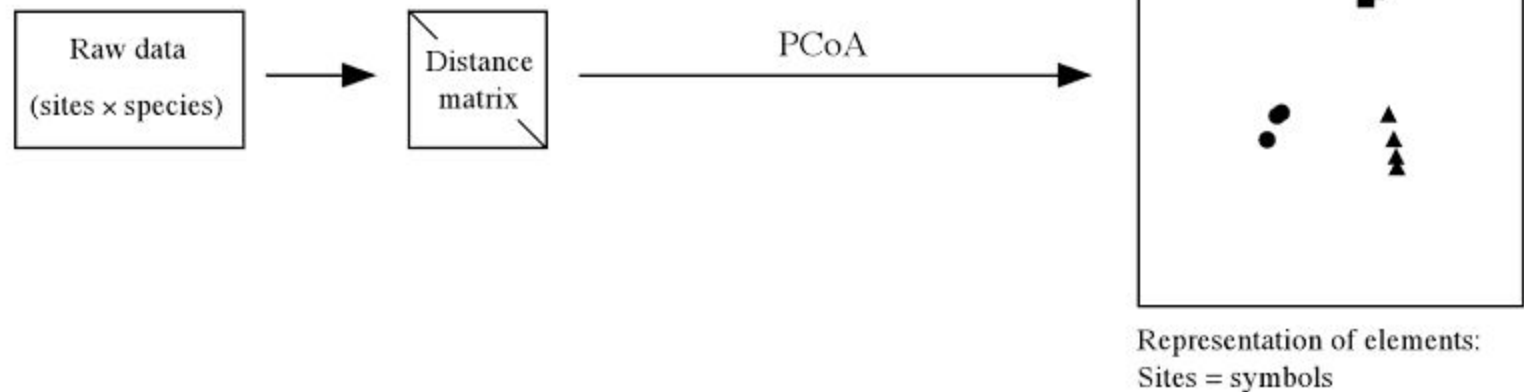
(a) Classical approach



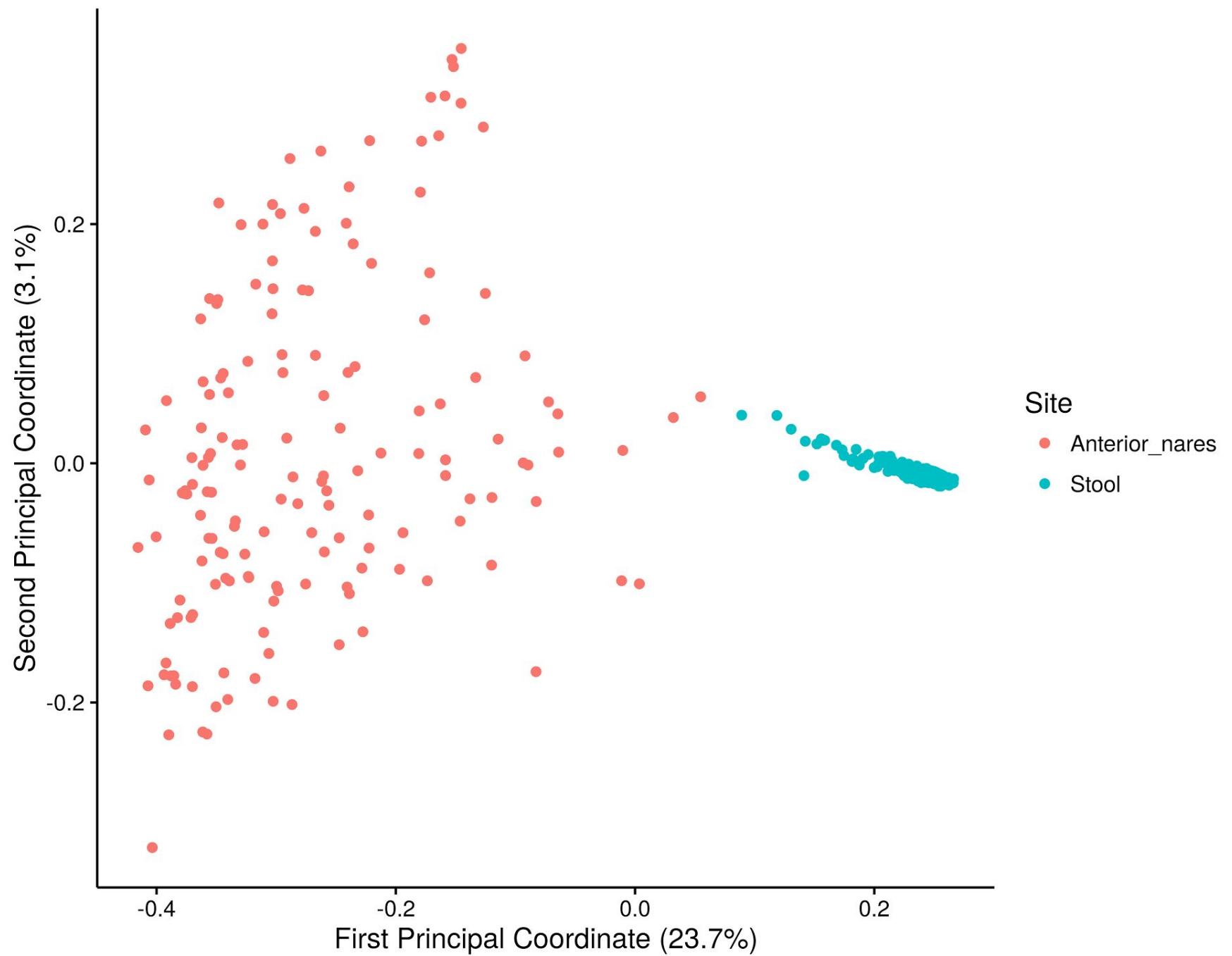
(b) Transformation-based approach (tb-PCA)



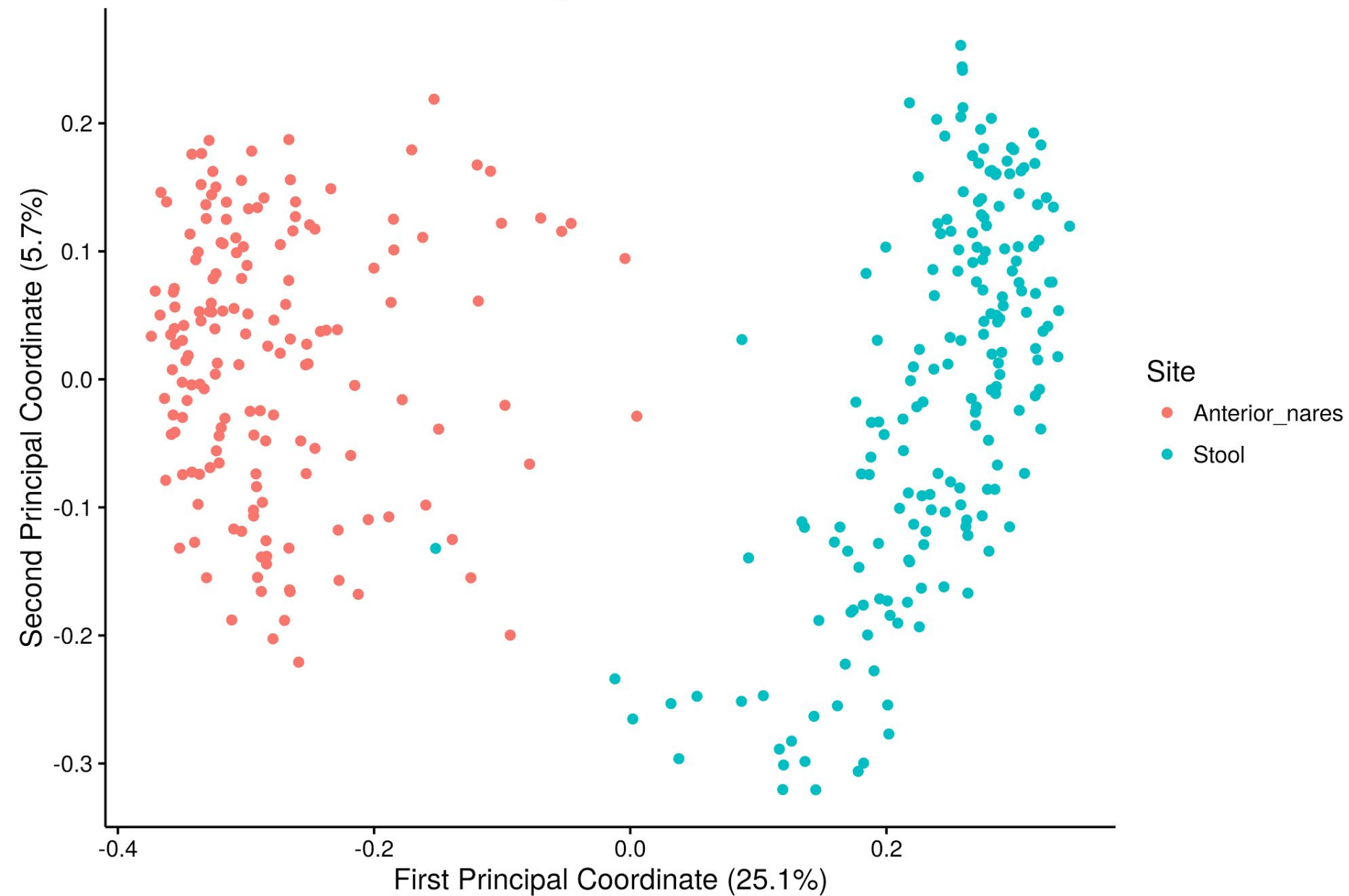
(c) Distance-based approach (PCoA)



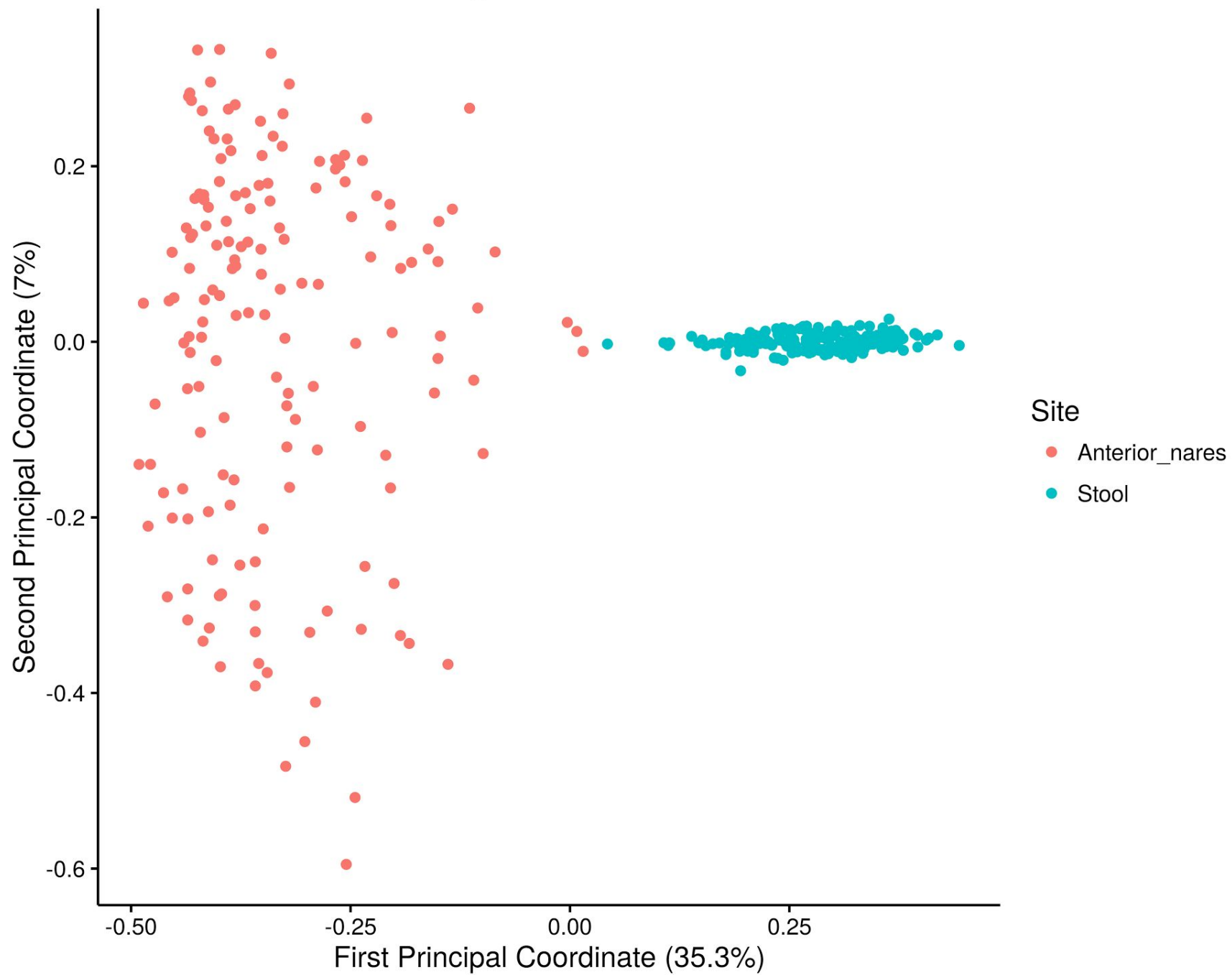
# Weighted UniFrac



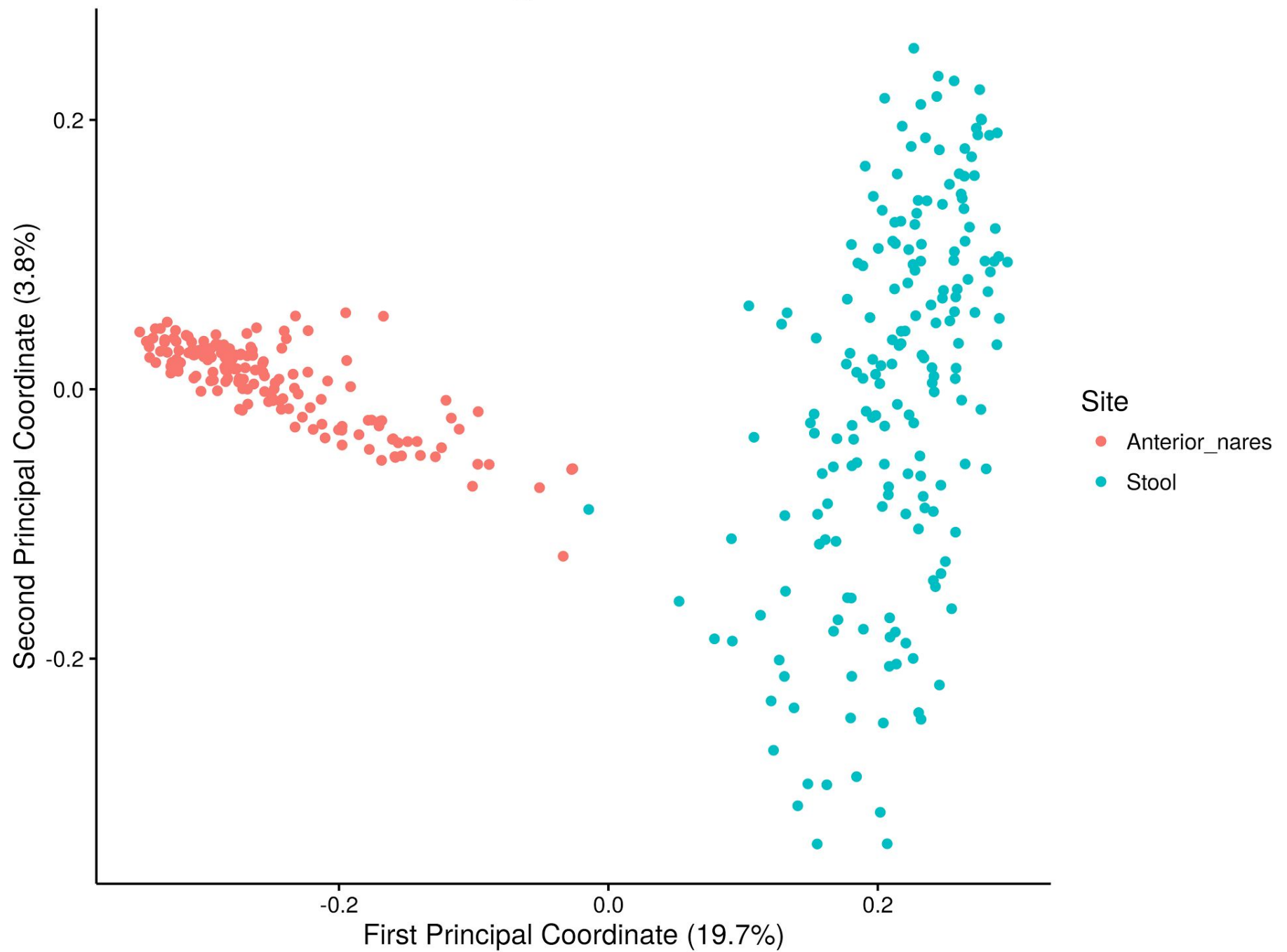
# Unweighted UniFrac



# Weighted Jaccard



# Unweighted Jaccard



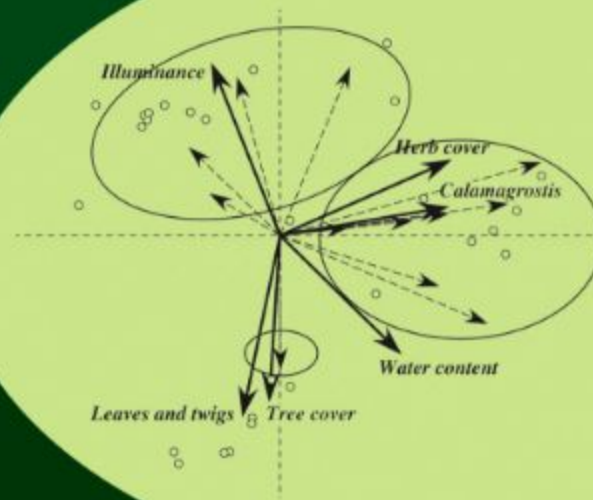


Developments in  
Environmental Modelling  
**Vol. 24**

Third English  
Edition

# Numerical Ecology

Pierre Legendre  
Louis Legendre



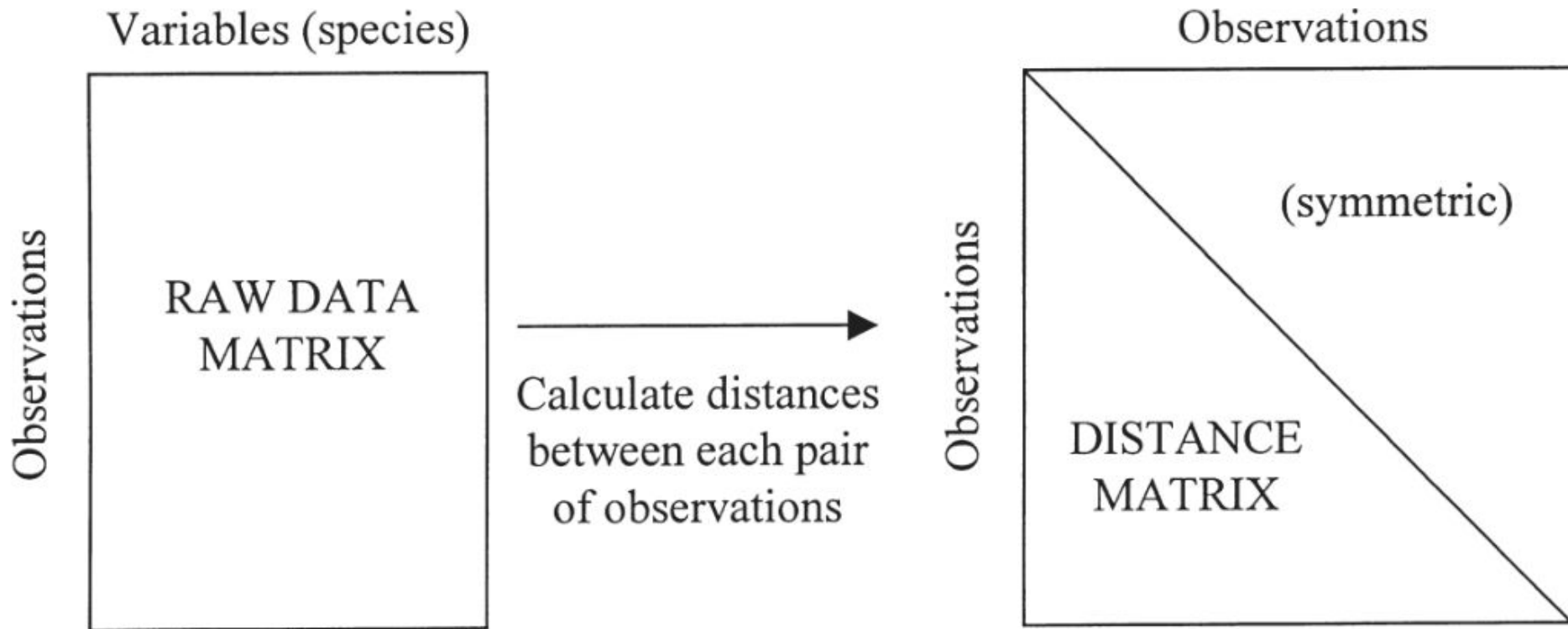
# Outline

- ▷ The problem: too much data.
- ▷ Reducing dimensions:
  - richness, evenness, and diversity
  - ecological distances (UniFrac)
  - PCA & PCoA
  - PERMANOVA (adonis)
- ▷ Cluster analysis:
  - supervised vs unsupervised learning
  - Dirichlet multinomial mixture modeling

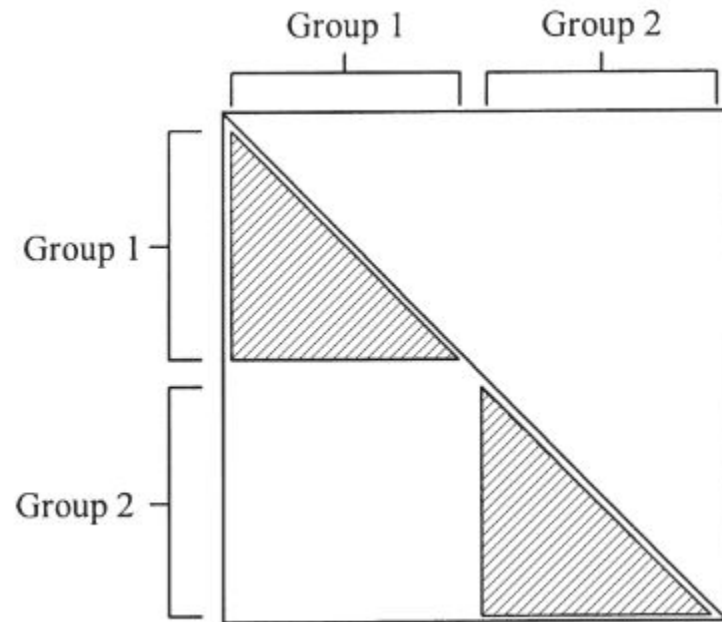
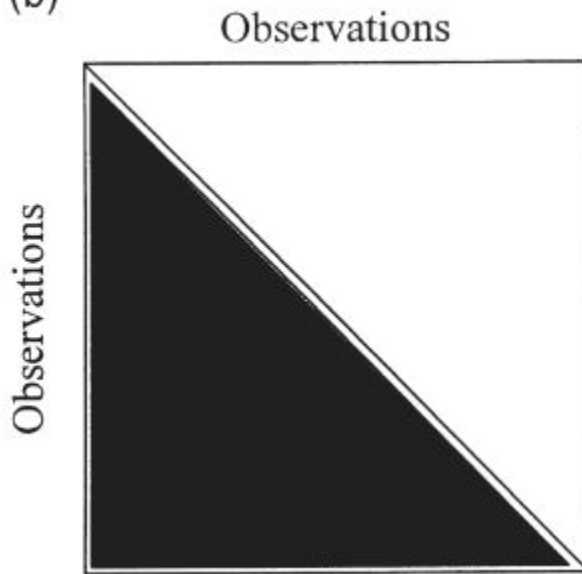


# Pairwise Distances: PERMANOVA

- ▷ Pairwise distance matrix can be partitioned by group assignment and ANOVA-like analysis can be applied to detect difference between groups.
- ▷ PERMANOVA: permutational ANOVA (aka, adonis)
  - pseudo F-ratio: conceptually similar but not F-distributed
  - testing by label permutation
  - quantification of effect size by R-squared or omega-squared (the latter a less biased estimator of true effect)

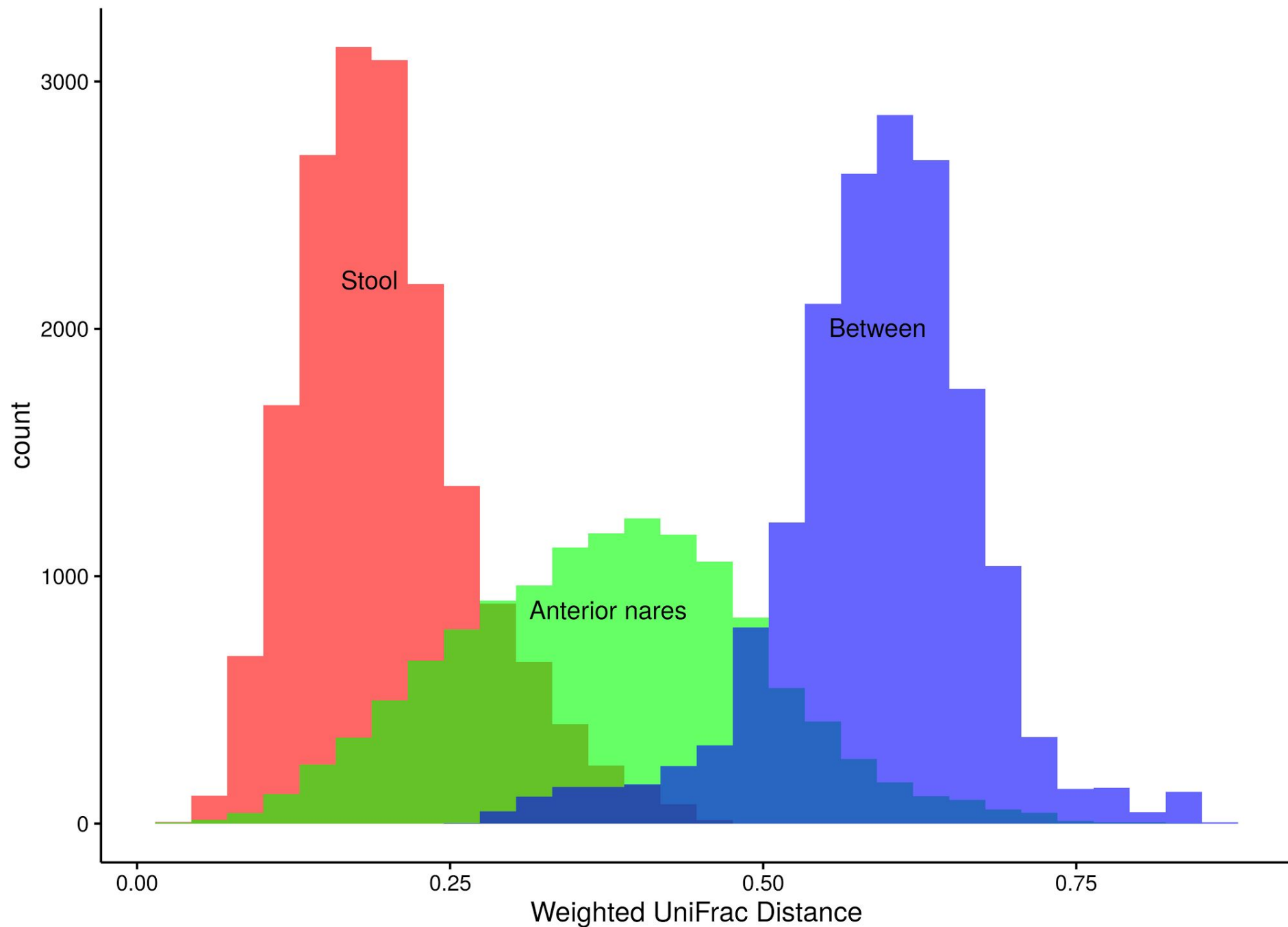


(b)

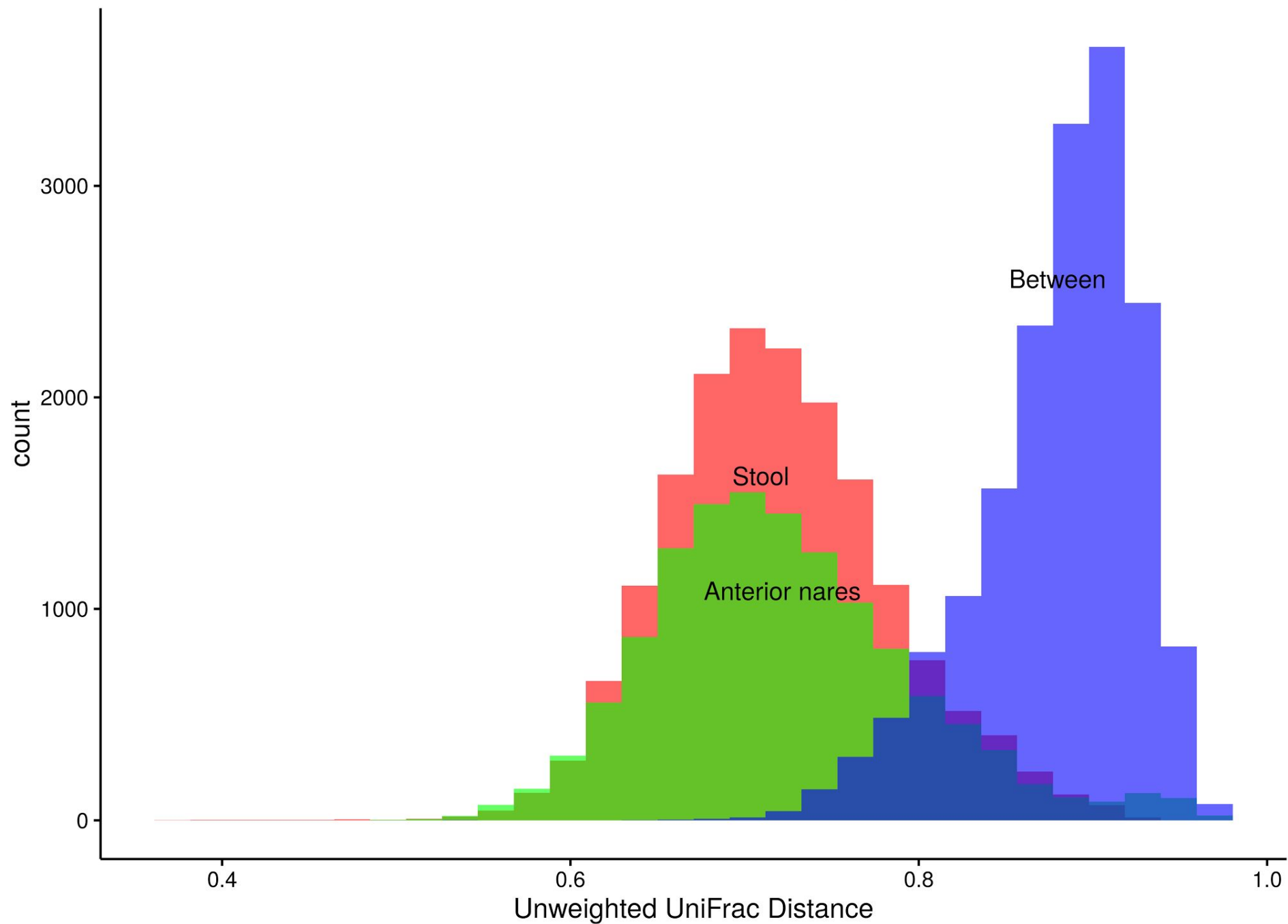


$$F = \frac{SS_A / (a - 1)}{SS_W / (N - a)}$$

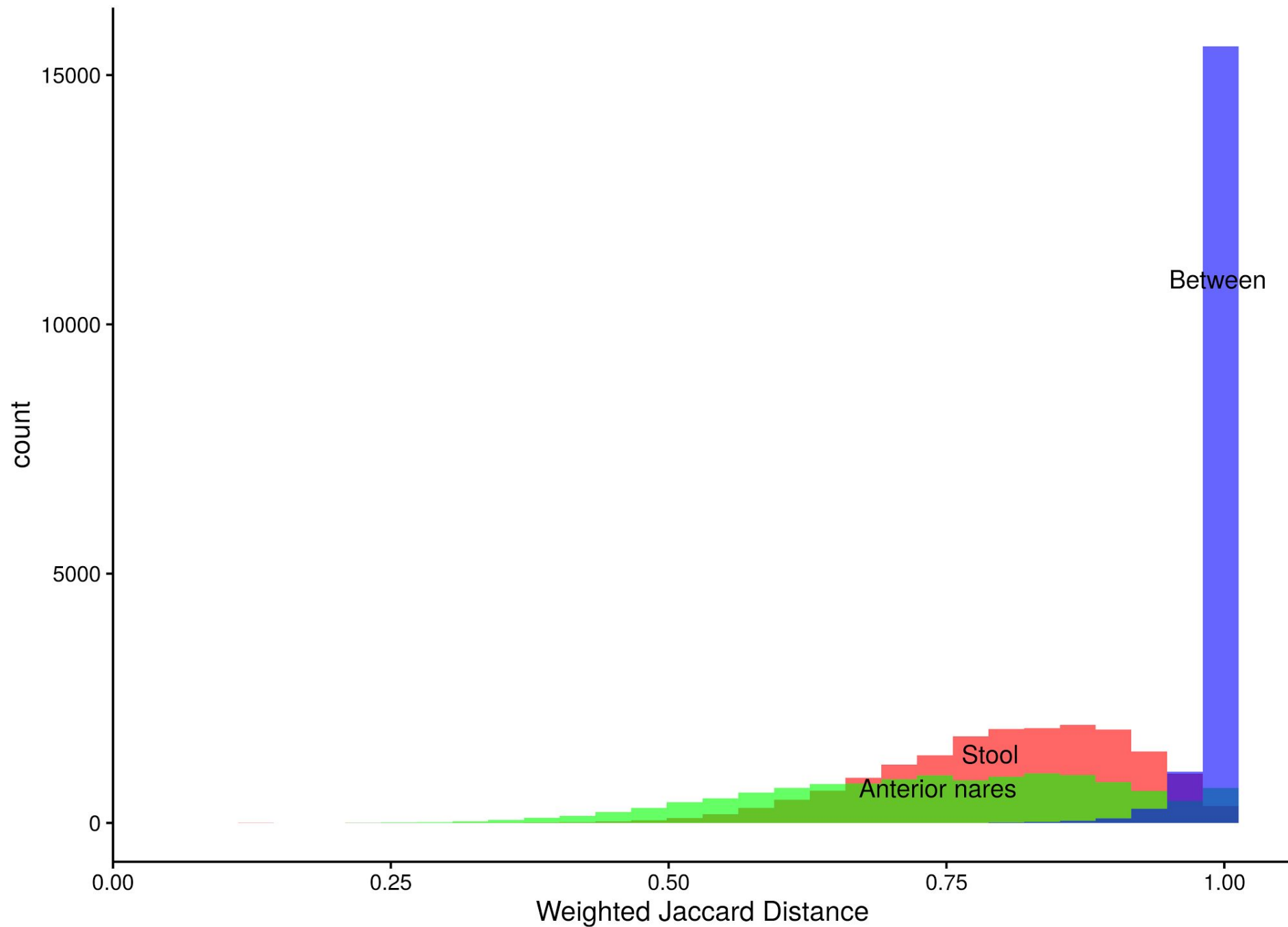
# HMP V1-V3 16S rRNA Amplicon



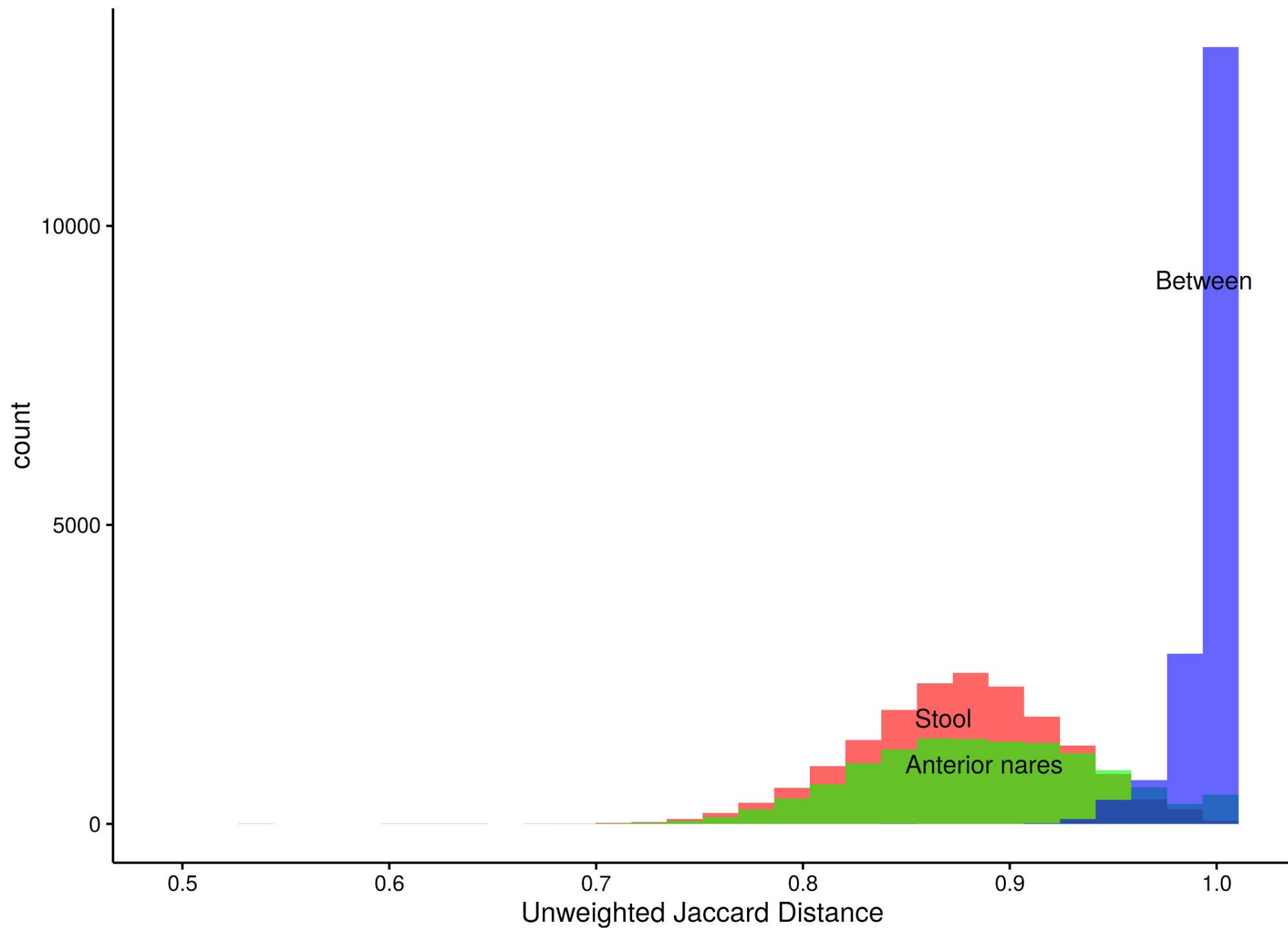
# HMP V1-V3 16S rRNA Amplicon



# HMP V1-V3 16S rRNA Amplicon



# HMP V1-V3 16S rRNA Amplicon



$$R^2 = 1 - \frac{SS_W}{SS_W + SS_A} = \frac{SS_A}{SS_T}$$

$$\omega^2 = \frac{SS_A - (a-1) \frac{SS_W}{N-a}}{SS_T + \frac{SS_W}{N-a}}$$

**Table 1.** Effect sizes observed from various exposures/interventions in studies of various microbiome sampling sites are shown as measured by omega-squared ( $\omega^2$ ) statistics, together with the *P*-values from PERMANOVA test

Site	Comparison groups		$\omega^2$ / <i>P</i> -value				Reference
	Control	Exposure	Weighted UniFrac	Unweighted UniFrac	Weighted Jaccard	Unweighted Jaccard	
Nares	Non-smoker (33)	Smoker (29)	0.042/0.001	0.009/0.001	0.023/0.001	0.007/0.001	<a href="#">Charlson <i>et al.</i> (2010)</a>
Oral	Non-smoker (33)	Smoker (29)	0.032/0.001	0.008/0.001	0.024/0.001	0.007/0.001	<a href="#">Charlson <i>et al.</i> (2010)</a>
Gut	Before feeding (10)	After feeding (10)	0.056/0.138	0.013/0.986	0/0.989	0.014/0.985	<a href="#">Wu <i>et al.</i> (2011)</a>
Oral	No azithromycin (42)	Azithromycin (6)	0.063/0.01	0.039/0.001	0.099/0.004	0.032/0.001	<a href="#">Charlson <i>et al.</i> (2012)</a>
Lung	No azithromycin (34)	Azithromycin (6)	0.065/0.005	0.038/0.001	0.019/0.089	0.033/0.001	<a href="#">Charlson <i>et al.</i> (2012)</a>
Skin	Left retroauricular (186)	Right retroauricular (187)	0.000/0.828	0.0001/0.327	0.000/0.986	0.000/1.000	<a href="#">HMP Consortium (2012b)</a>
Human	Anterior nares (161)	Stool (187)	0.567/0.001	0.201/0.001	0.230/0.001	0.117/0.001	<a href="#">HMP Consortium (2012b)</a>

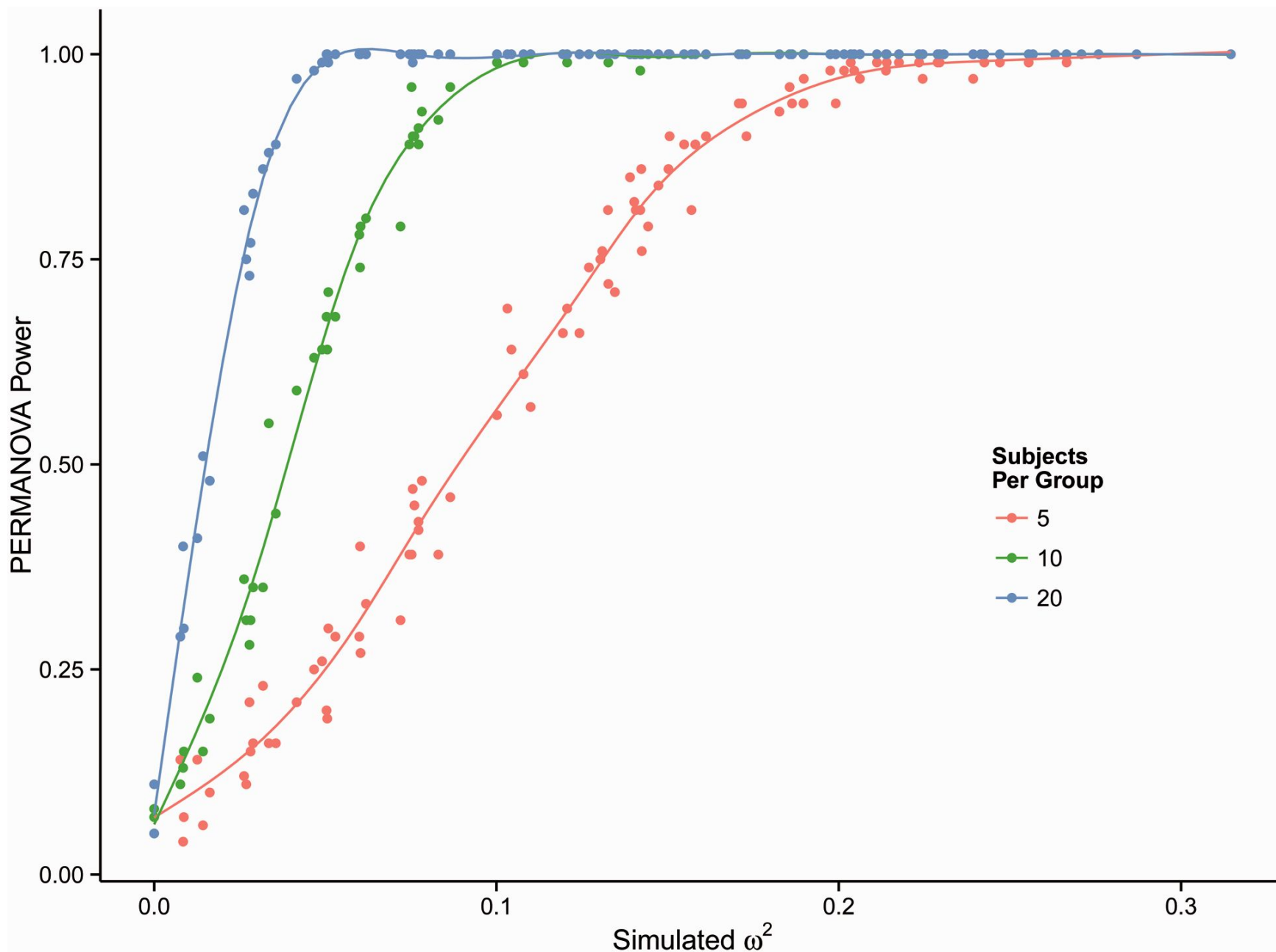


$$R^2 = 1 - \frac{SS_W}{SS_W + SS_A} = \frac{SS_A}{SS_T}$$

$$\omega^2 = \frac{SS_A - (a-1) \frac{SS_W}{N-a}}{SS_T + \frac{SS_W}{N-a}}$$

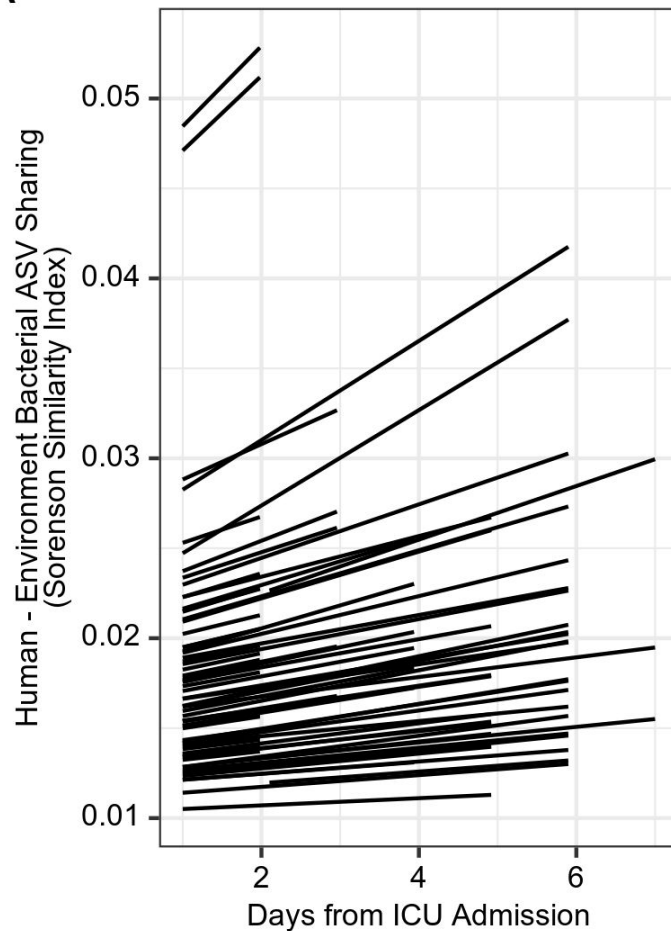
**Table 1.** Effect sizes observed from various exposures/interventions in studies of various microbiome sampling sites are shown as measured by omega-squared ( $\omega^2$ ) statistics, together with the *P*-values from PERMANOVA test

Site	Comparison groups		$\omega^2$ / <i>P</i> -value				Reference
	Control	Exposure	Weighted UniFrac	Unweighted UniFrac	Weighted Jaccard	Unweighted Jaccard	
Nares	Non-smoker (33)	Smoker (29)	0.042/0.001	0.009/0.001	0.023/0.001	0.007/0.001	Charlson <i>et al.</i> (2010)
Oral	Non-smoker (33)	Smoker (29)	0.032/0.001	0.008/0.001	0.024/0.001	0.007/0.001	Charlson <i>et al.</i> (2010)
Gut	Before feeding (10)	After feeding (10)	0.056/0.138	0.013/0.986	0/0.989	0.014/0.985	Wu <i>et al.</i> (2011)
Oral	No azithromycin (42)	Azithromycin (6)	0.063/0.01	0.039/0.001	0.099/0.004	0.032/0.001	Charlson <i>et al.</i> (2012)
Lung	No azithromycin (34)	Azithromycin (6)	0.065/0.005	0.038/0.001	0.019/0.089	0.033/0.001	Charlson <i>et al.</i> (2012)
Skin	Left retroauricular (186)	Right retroauricular (187)	0.000/0.828	0.0001/0.327	0.000/0.986	0.000/1.000	HMP Consortium (2012b)
Human	Anterior nares (161)	Stool (187)	0.567/0.001	0.201/0.001	0.230/0.001	0.117/0.001	HMP Consortium (2012b)

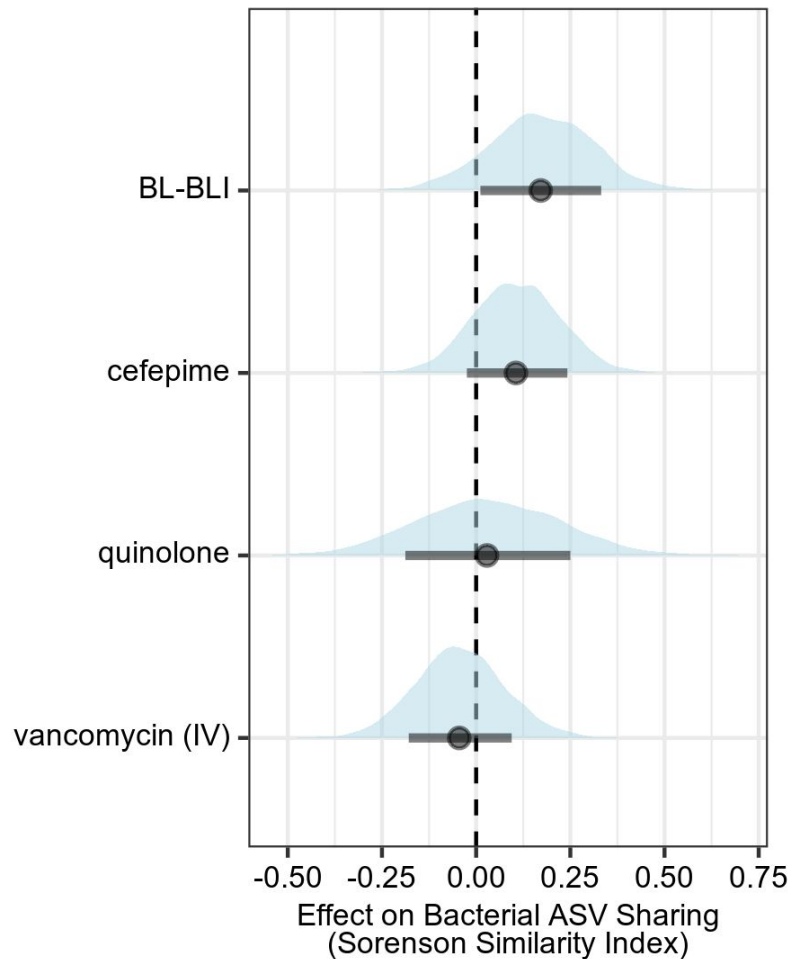


# Other Approaches to Modeling Distance

A



B



# Outline

- ▷ The problem: too much data.
- ▷ Reducing dimensions:
  - richness, evenness, and diversity
  - ecological distances (UniFrac)
  - PCA & PCoA
  - PERMANOVA (adonis)
- ▷ **Cluster analysis:**
  - supervised vs unsupervised learning
  - Dirichlet multinomial mixture modeling

# Statistical / Machine Learning

- ▷ Supervised learning:
  - exposure and outcome
  - regression, linear discriminant analysis, KNN clustering
  - test & training data; cross-validation
- ▷ Unsupervised learning:
  - understand relationships between observations or variables
  - can we reduce the dimensions of microbiome data?

Springer Texts in Statistics

Gareth James  
Daniela Witten  
Trevor Hastie  
Robert Tibshirani

# An Introduction to Statistical Learning

with Applications in R

 Springer

Springer Series in Statistics

Trevor Hastie  
Robert Tibshirani  
Jerome Friedman

# The Elements of Statistical Learning

Data Mining, Inference, and Prediction

Second Edition

 Springer

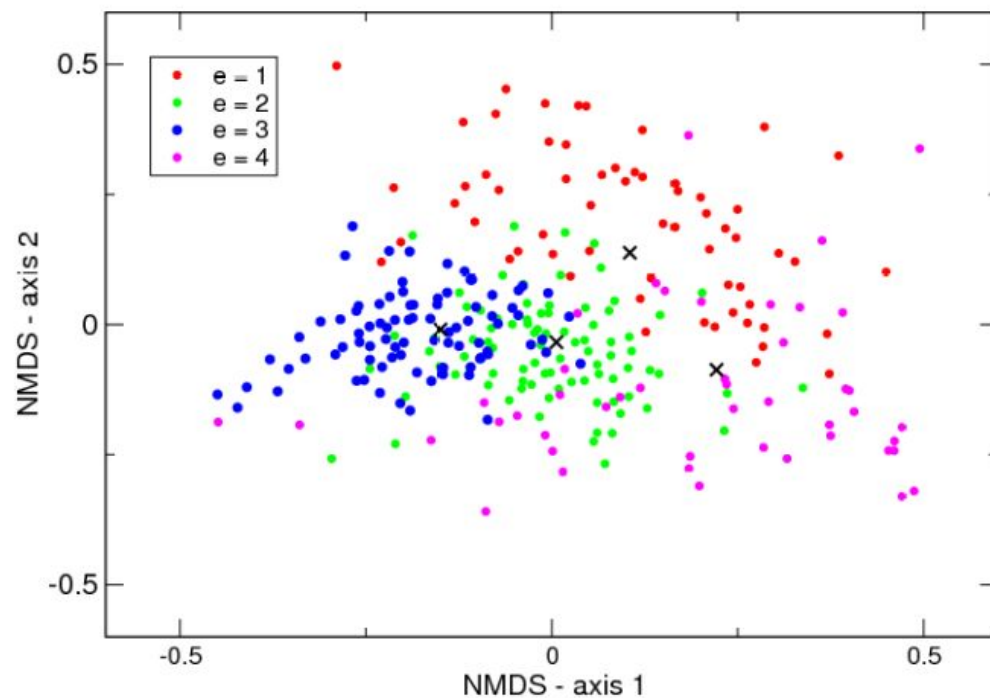
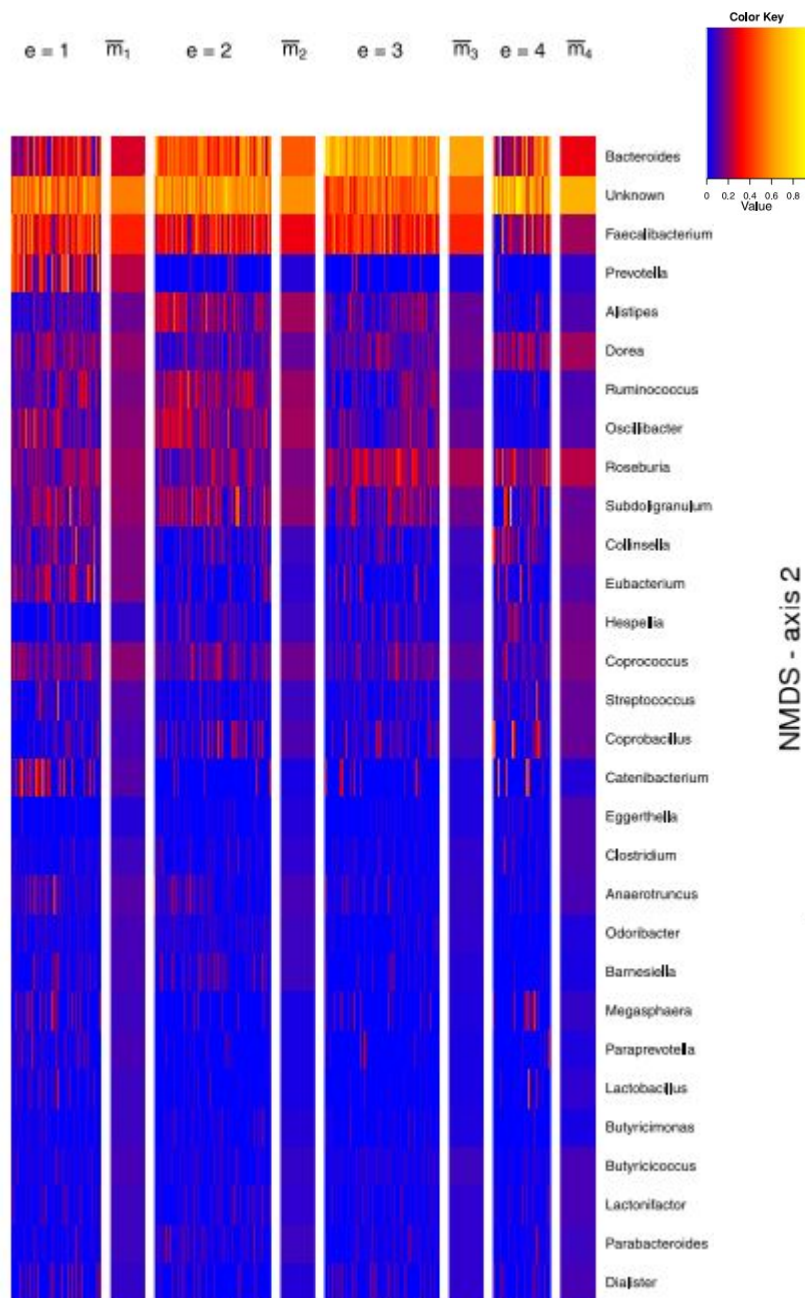
# Outline

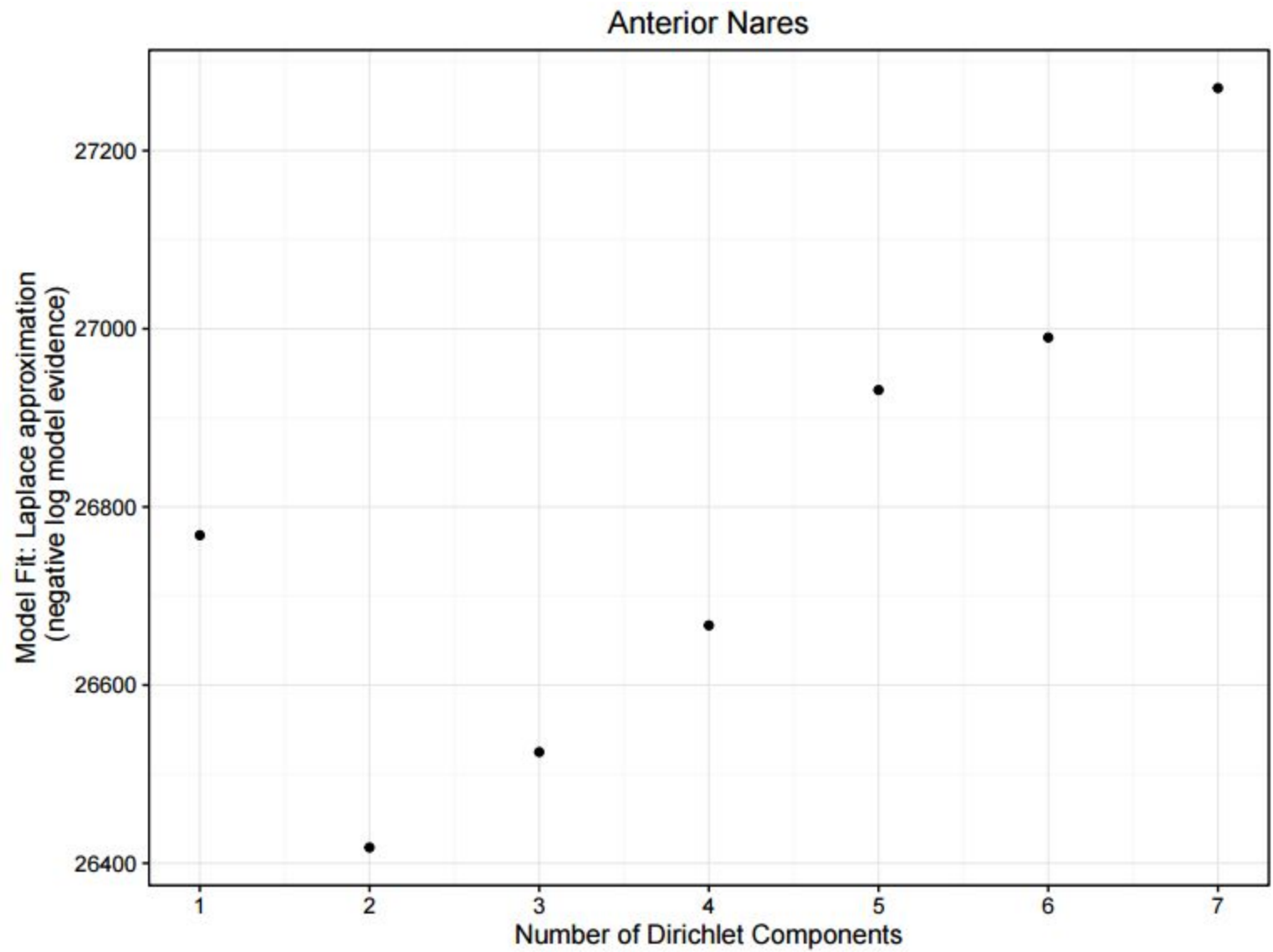
- ▷ The problem: too much data.
- ▷ Reducing dimensions:
  - richness, evenness, and diversity
  - ecological distances (UniFrac)
  - PCA & PCoA
  - PERMANOVA (adonis)
- ▷ **Cluster analysis:**
  - supervised vs unsupervised learning
  - **Dirichlet multinomial mixture modeling**

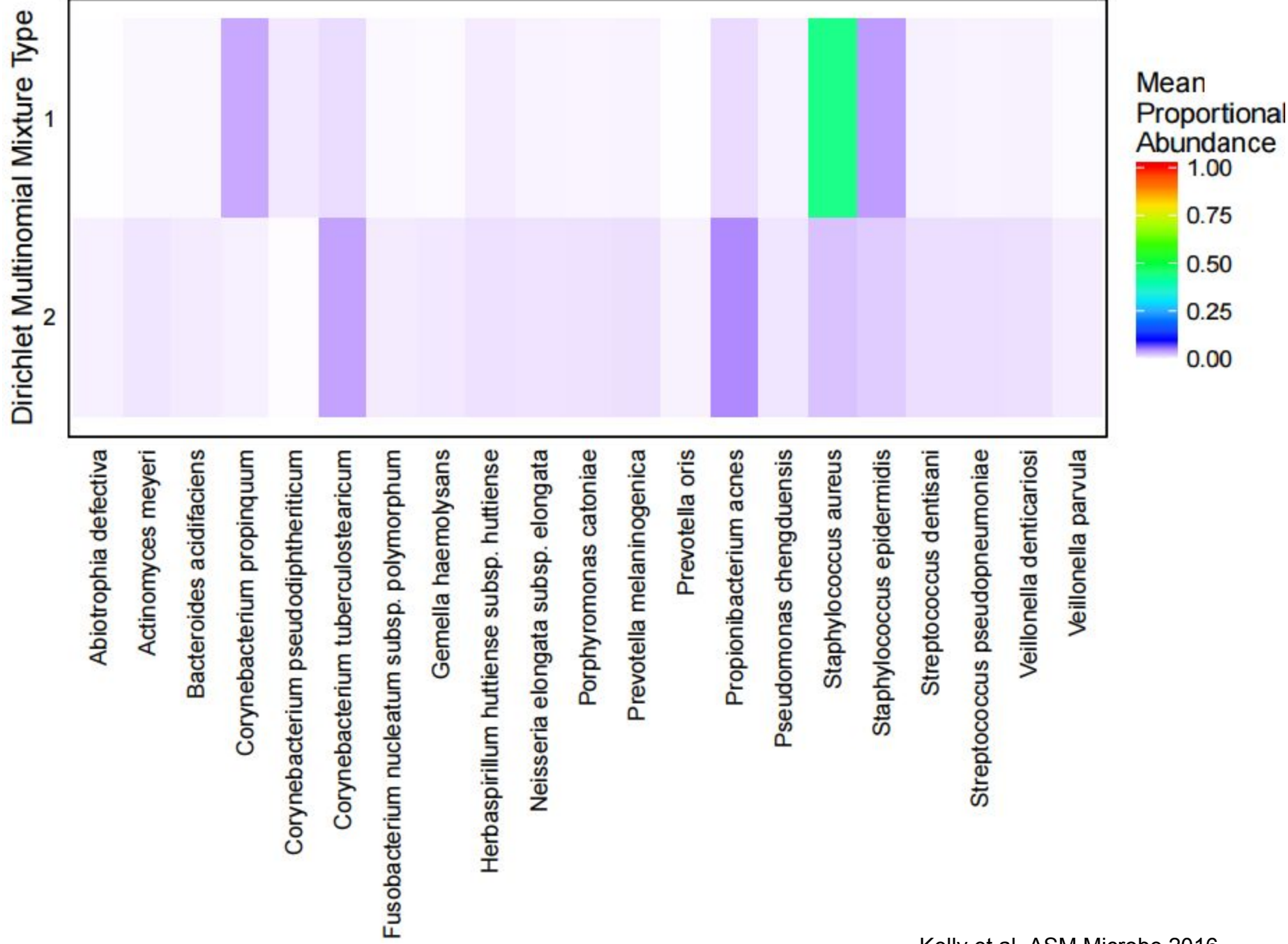
# Dimension Reduction: Mixture Models

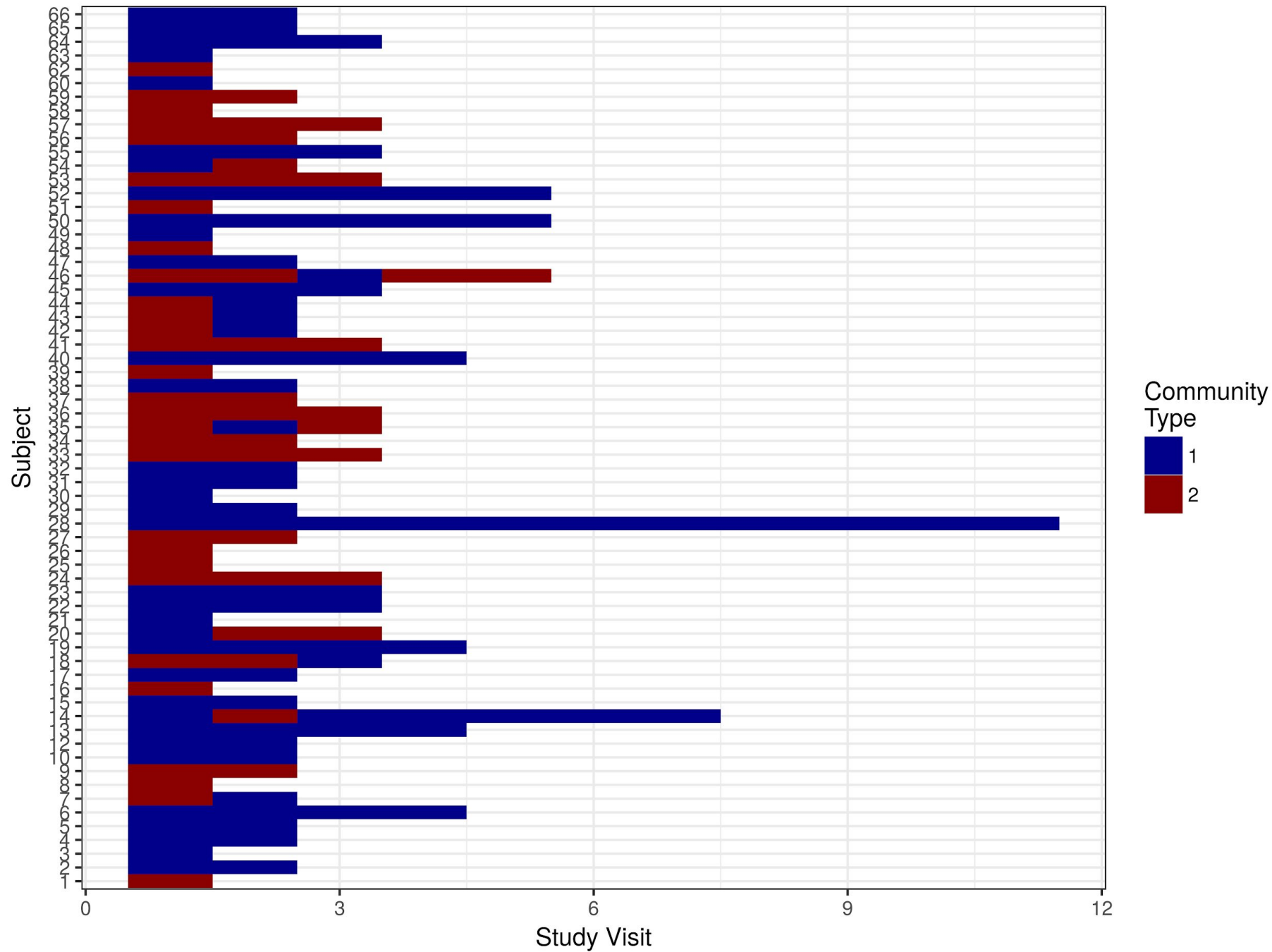
- ▷ Dirichlet-multinomial distribution:
  - compound probability distribution: probability vector drawn from Dirichlet distribution (generalized beta) → observation drawn from multinomial distribution (generalized binomial)
- ▷ D-M mixture modelling:
  - each sample  $\sim$  multinomial from one Dirichlet vector
  - # Dirichlet vectors: minimize  $-\log(\text{model evidence, Laplace approx})$
  - Dirichlet probability vectors = “community types”











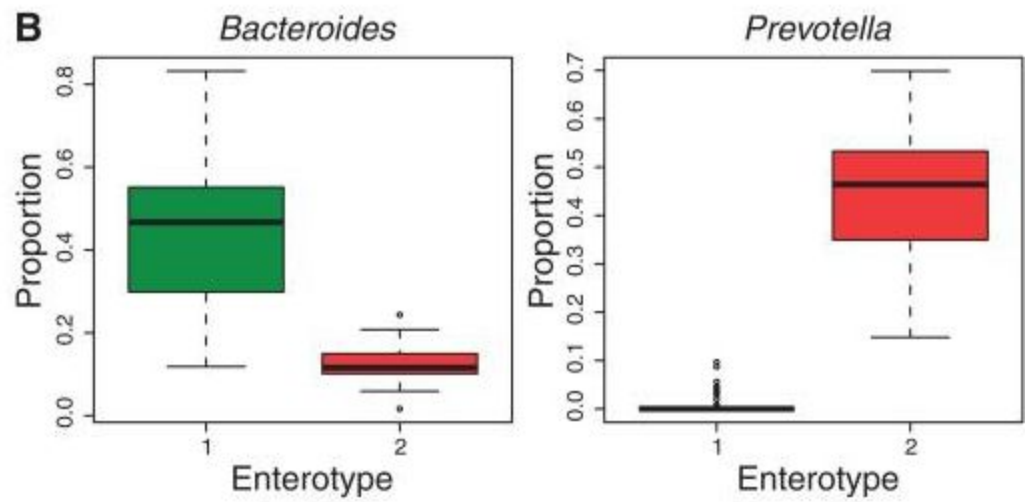
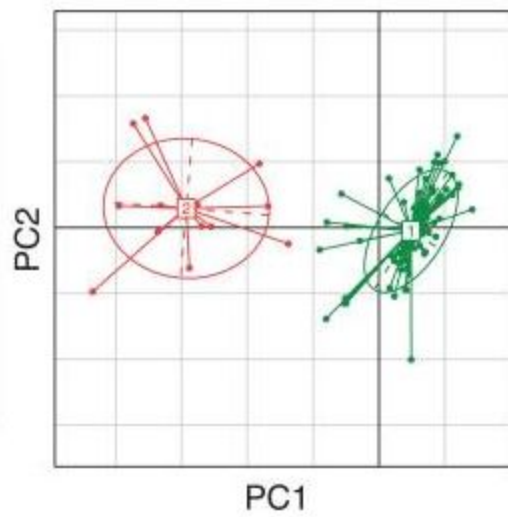
# LETTER

doi:10.1038/nature13178

---

## **Dynamics and associations of microbial community types across the human body**

Tao Ding<sup>1</sup> & Patrick D. Schloss<sup>1</sup>



Wu et al Science 2011.

# Methods for Microbiome Data

- ▷ Visualization: heatmaps and barplots.
- ▷ Single-taxon hypothesis.
- ▷ Alpha diversity: richness and evenness.
- ▷ Pairwise distances: count/phylogeny, weighted/unweighted.
- ▷ Ordination: PCA & PCoA.
- ▷ PERMANOVA: categorical exposure & microbiome outcome  
(allows quantification of effect size)
- ▷ DMM models: unsupervised clustering > “community types”  
(identify relationships among variables/OTUs)

# Conclusions

- ▷ Distance-based analysis and adonis/PERMANOVA testing:
  - microbiome outcome measures
  - *omega2* to define effect size of exposure/intervention
  - power estimation
- ▷ Dirichlet-multinomial mixtures:
  - categorical analysis may correspond with biologic community types
  - identify key species
  - discovery / validation design



