# Distances and Ordination

# Community Distance

Communities are a vector of abundances:
$$\mathbf{x} = \{x_1,\ x_2,\ x_3,\ \ldots\}$$

*E. coli:* ● ● ●
*P. fluorescens:* ●
*B. subtilis:* ●
*P. acnes:*
*D. radiodurans:*
*H. pylori:* ● ● ● ● ● ● ●
*L. crispatus:*

$$\mathbf{x} = \{3,1,1,0,0,7,0\}$$

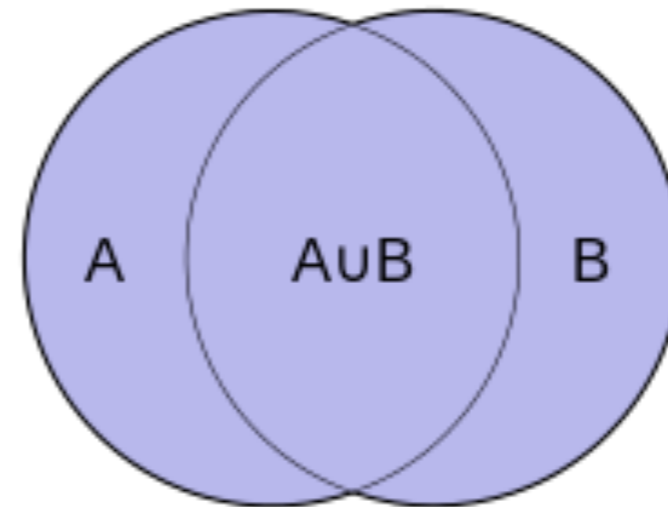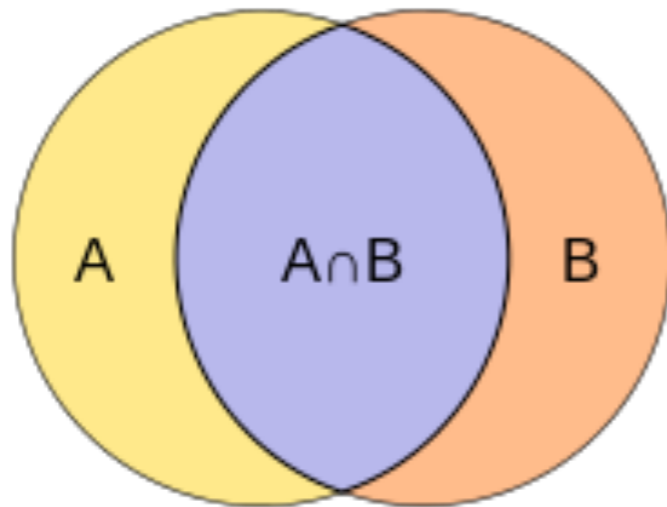# Community Distance Properties

- Range from 0 to 1

- Distance to self is 0

- If no shared taxa, distance is 1

- Triangle inequality (metric)

- Joint absences do not affect distance (biology)

- Independent of absolute counts (metagenomics)

# The Distance Spectrum

|  | Categorical | Phylogenetic |
|---|---|---|
| **Presence/ Absence** | Jaccard | Unifrac |
| **Quantitative Abundance** | Bray-Curtis | Weighted Unifrac |

# Jaccard

$Dist(A, B) = 1 - (A \cap B)/(A \cup B)$
$= ((\mathbf{x_A}>0) \& (\mathbf{x_B}>0))/((\mathbf{x_A}>0) | (\mathbf{x_B}>0))$

# Jaccard

$$\text{Dist}(A, B) = 1 - (A \cap B)/(A \cup B)$$
$$= ((\mathbf{x_A}>0) \,\&\, (\mathbf{x_B}>0))/((\mathbf{x_A}>0) \,|\, (\mathbf{x_B}>0))$$



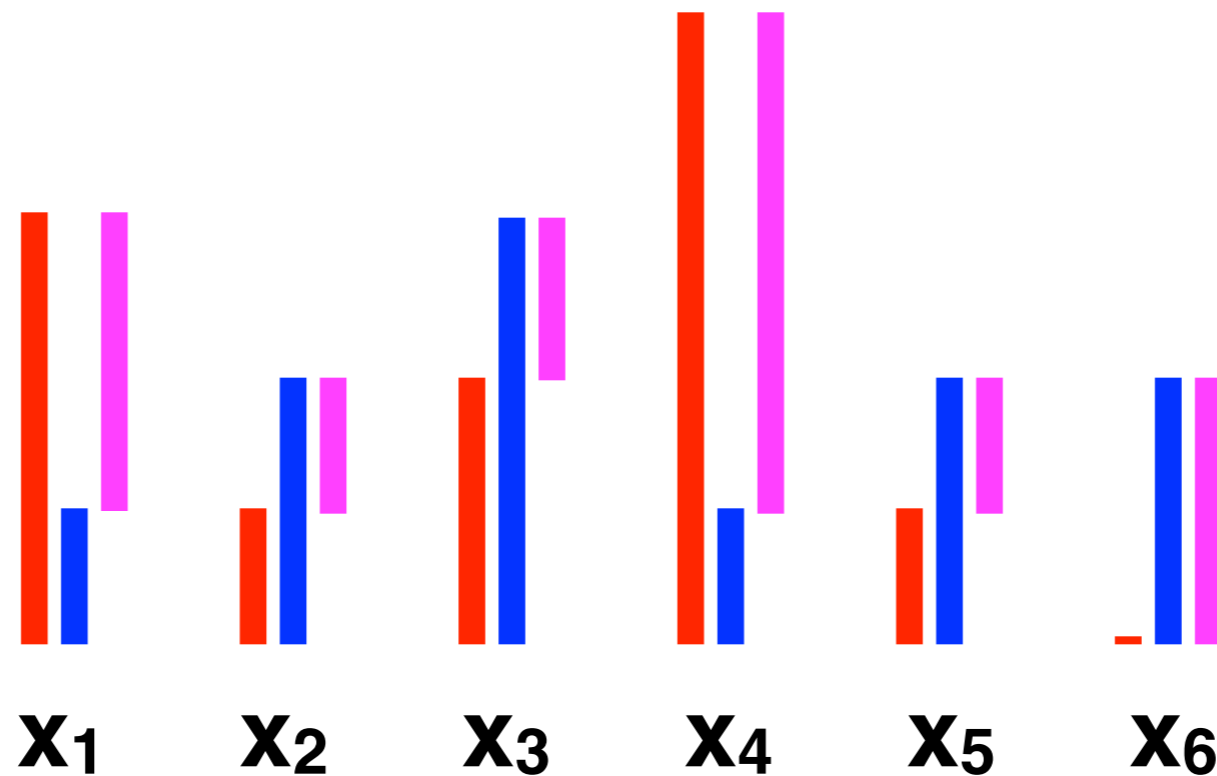**Intuition**: Fraction of shared **types** unique to one of the communities

# Bray-Curtis

$$\text{Dist}(x, y) = \frac{\sum |x_i - y_i|}{\sum x_i + \sum y_i} = \frac{\rule{2em}{0.6em}}{\textcolor{red}{\rule{2em}{0.6em}} + \textcolor{blue}{\rule{2em}{0.6em}}}$$

**x₁**  **x₂**  **x₃**  **x₄**  **x₅**  **x₆**

# Bray-Curtis

$$\text{Dist}(x, y) = \frac{\sum |x_i - y_i|}{\sum x_i + \sum y_i} = \frac{\rule{2em}{0.6em}}{\rule{2em}{0.6em} + \rule{2em}{0.6em}}$$
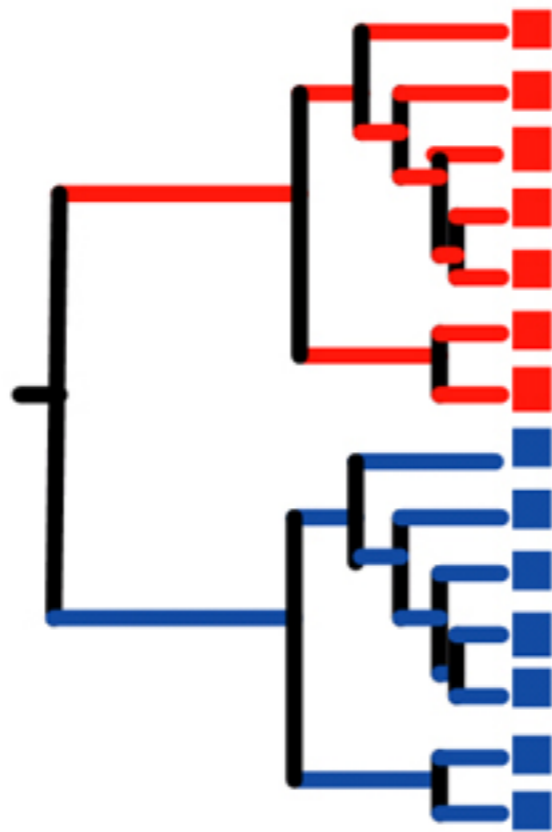


$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6$

**Intuition**: *City block distance*. Sum of absolute differences over total abundance.

# Unifrac



$$\text{Dist}(x, y) = \frac{\textcolor{red}{\rule{1cm}{0.3em}} + \textcolor{blue}{\rule{1cm}{0.3em}}}{\textcolor{red}{\rule{1cm}{0.3em}} + \textcolor{blue}{\rule{1cm}{0.3em}} + \textcolor{purple}{\rule{1cm}{0.3em}}}$$

D = 1

D = ~ 0.5

Lozupone and Knight (2008)

# Unifrac

$$\text{Dist}(x, y) = \frac{\textcolor{red}{\rule{1.2cm}{0.3cm}} + \textcolor{blue}{\rule{1.2cm}{0.3cm}}}{\textcolor{red}{\rule{1.2cm}{0.3cm}} + \textcolor{blue}{\rule{1.2cm}{0.3cm}} + \textcolor{purple}{\rule{1.2cm}{0.3cm}}}$$



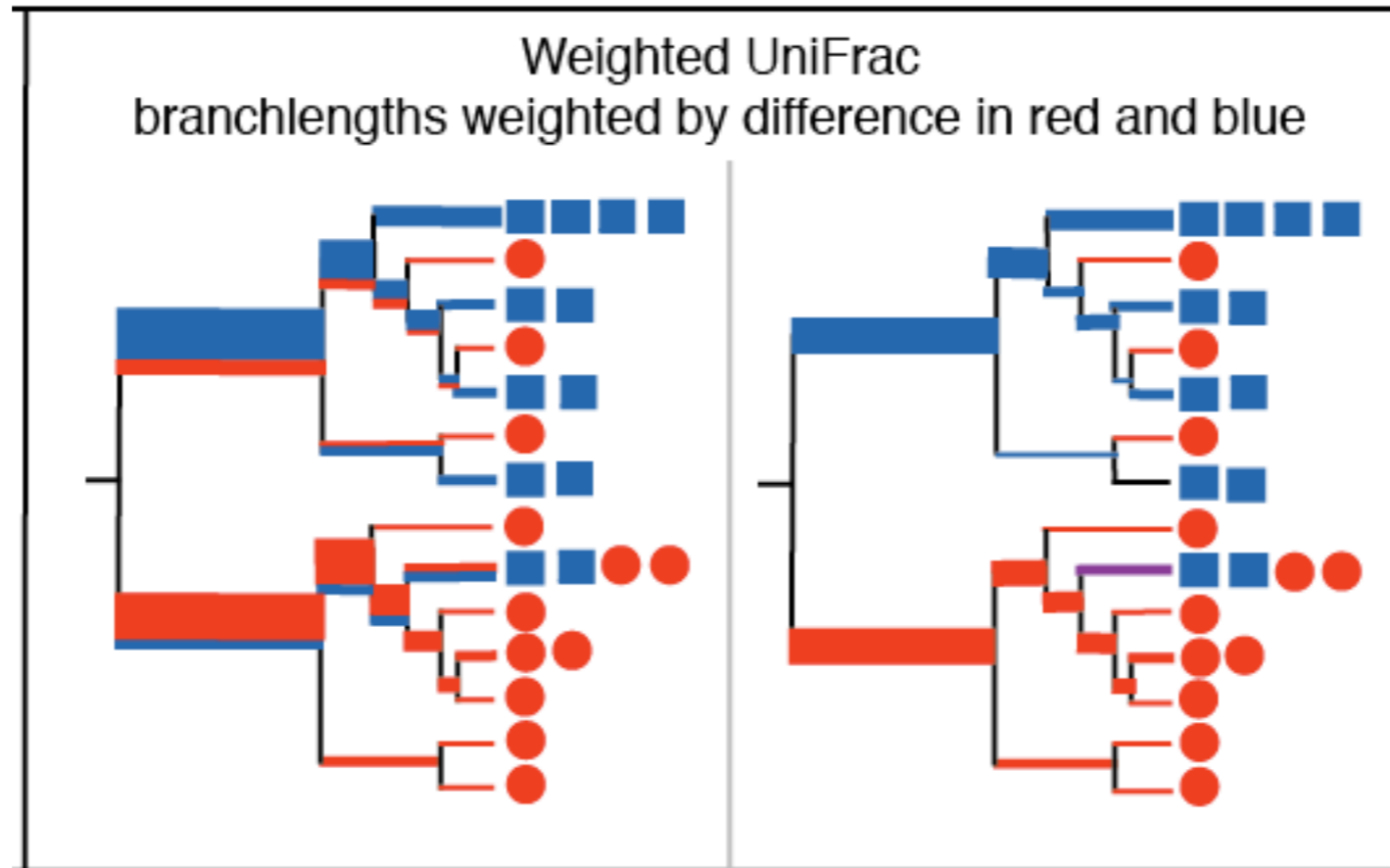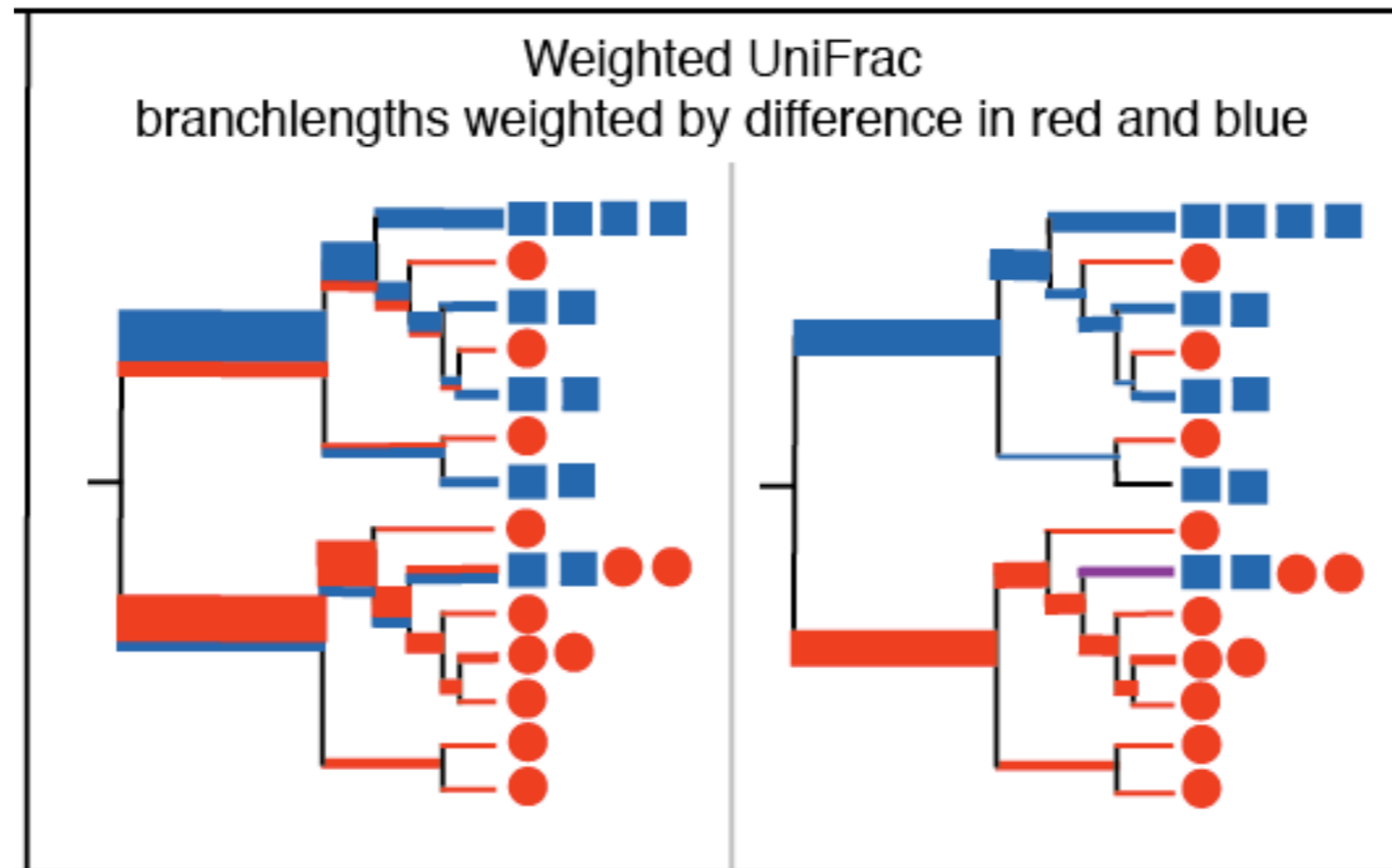D = 1          D = ~ 0.5

**Intuition**: Fraction of shared **tree** unique to one of the communities

Lozupone and Knight (2008)

# Weighted Unifrac



Weighted UniFrac
branchlengths weighted by difference in red and blue

Lozupone et al. (2007)

# Weighted Unifrac
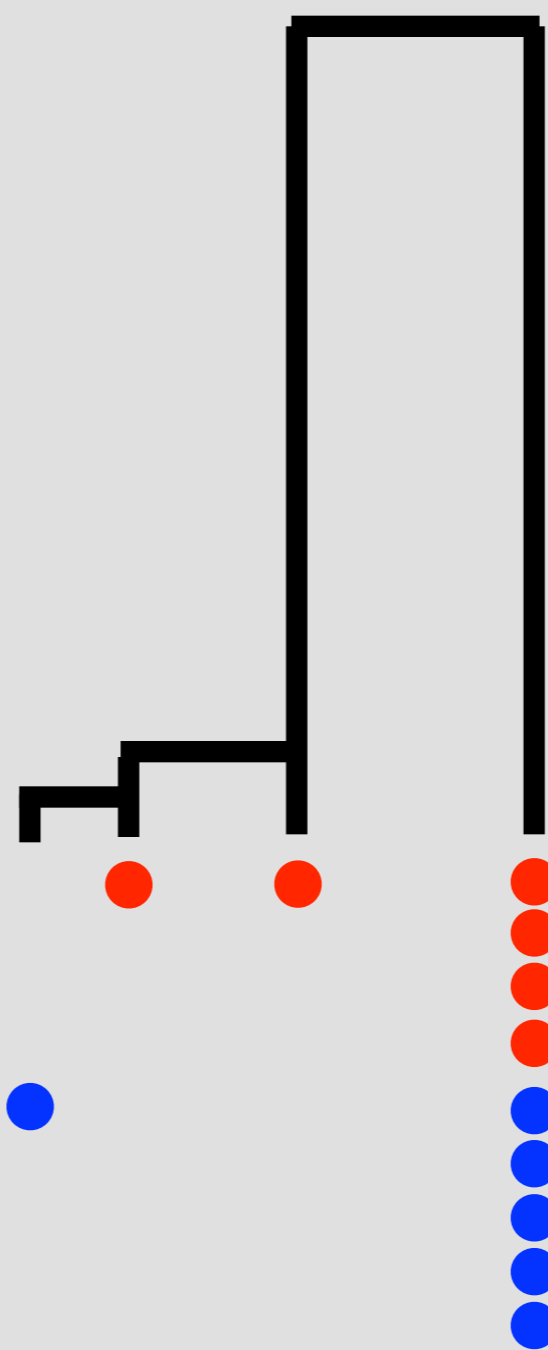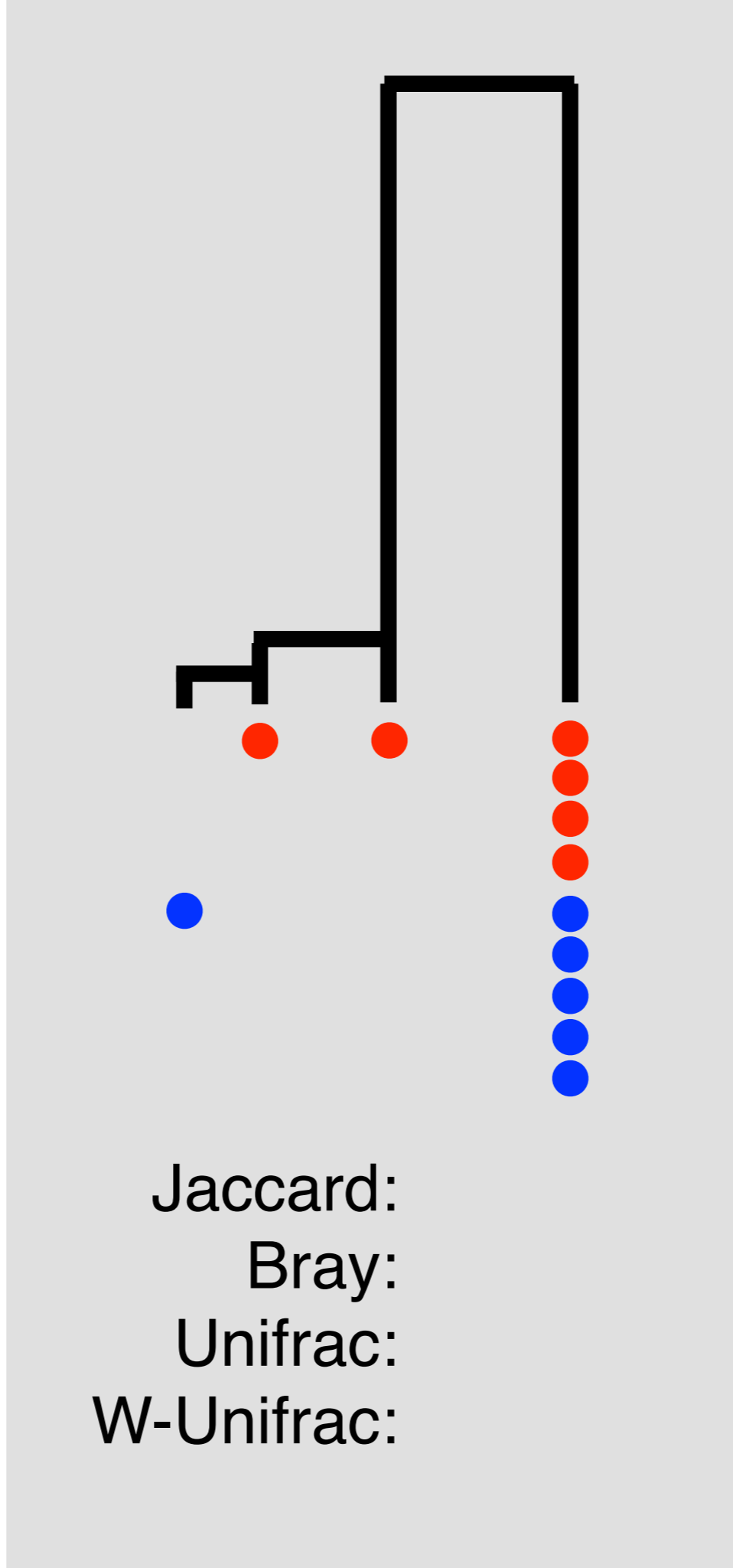


Weighted UniFrac
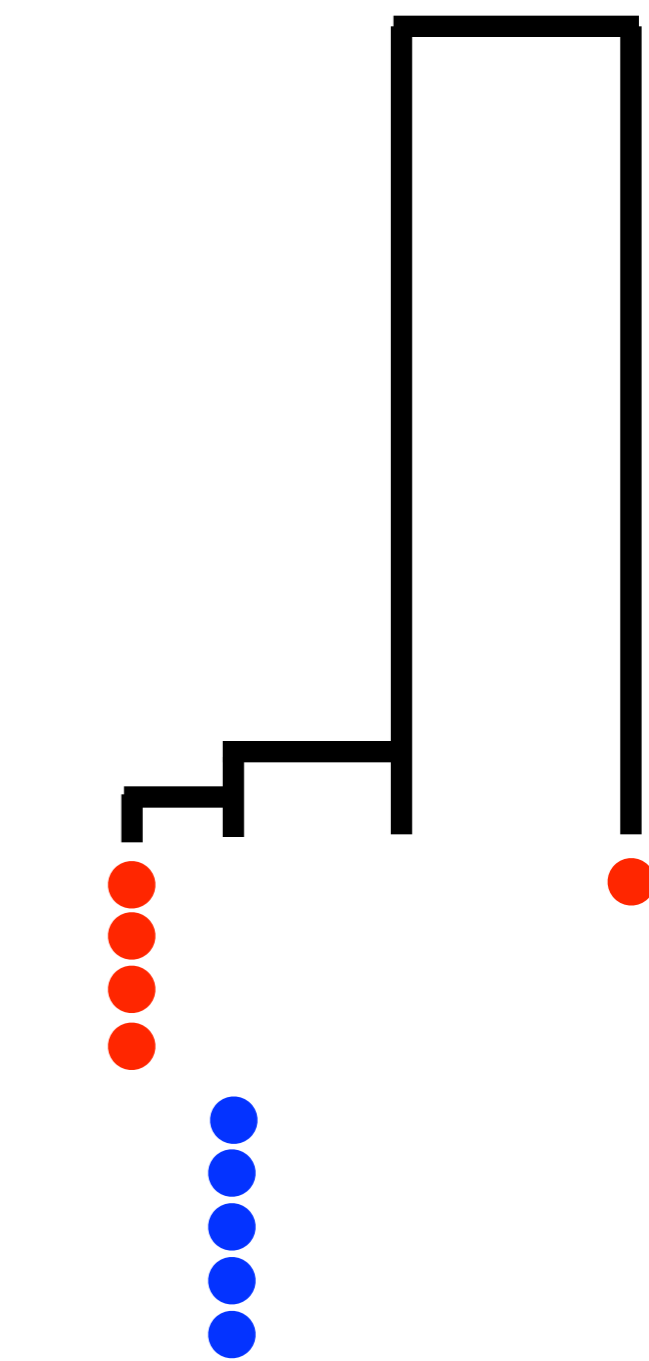branchlengths weighted by difference in red and blue

**Intuition**: The cost of turning one distribution into the other; where the cost is the amount of "dirt" moved times the distance by which it is moved.
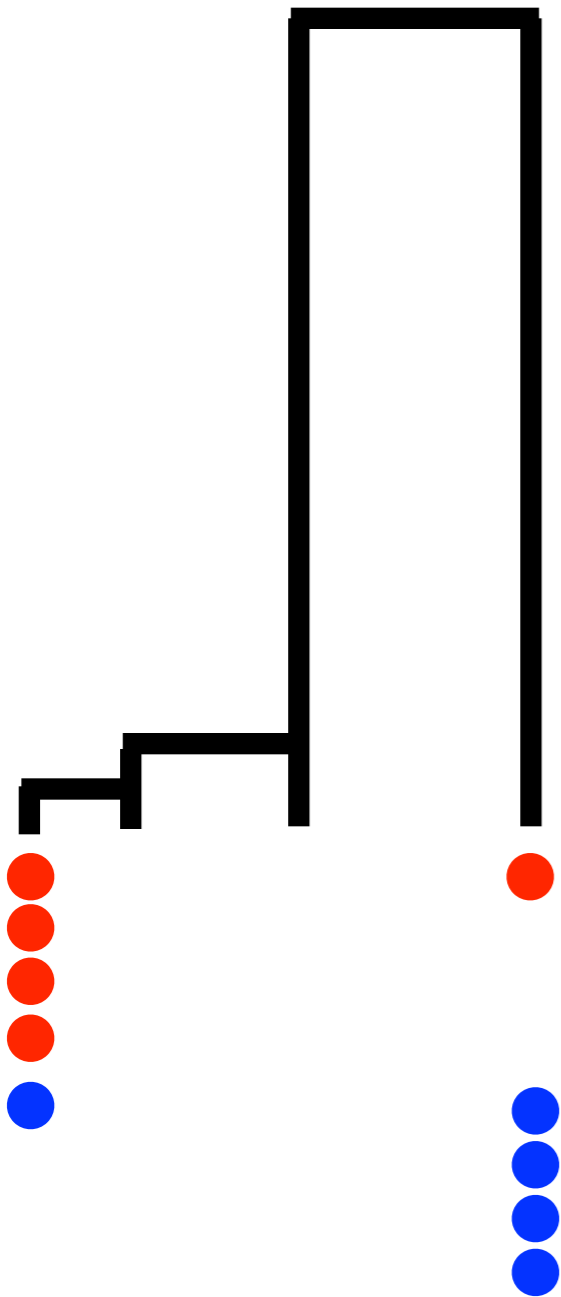
Lozupone et al. (2007)

Jaccard:
Bray:
Unifrac:
W-Unifrac:

Jaccard:
Bray:
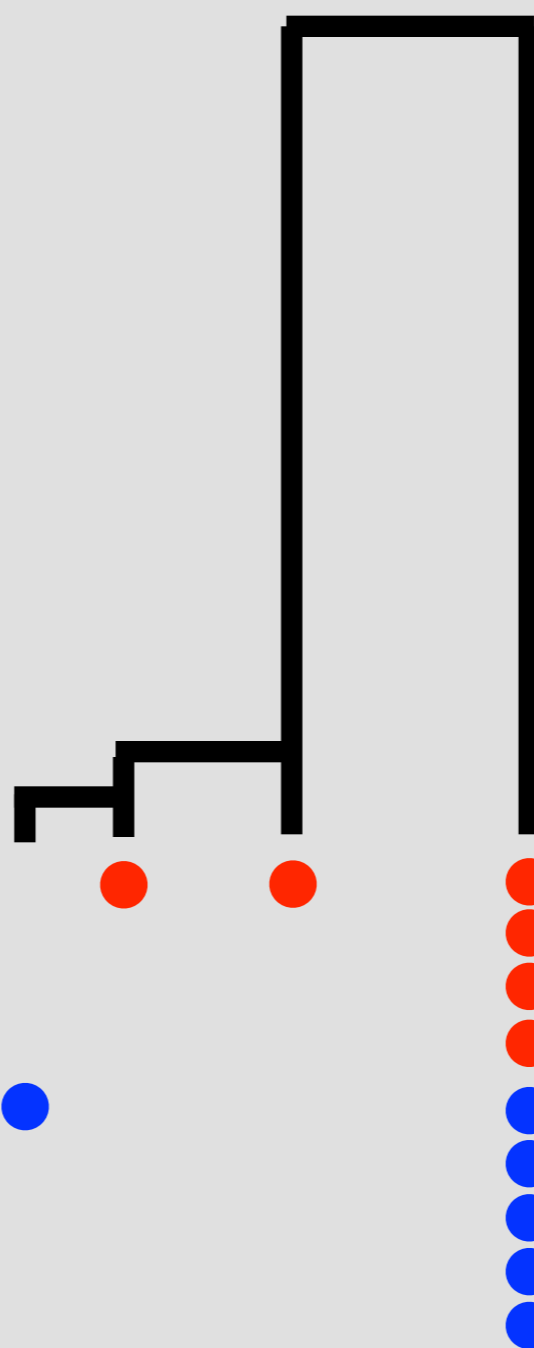Unifrac:
W-Unifrac:

Jaccard:
Bray:
Unifrac:
W-Unifrac:

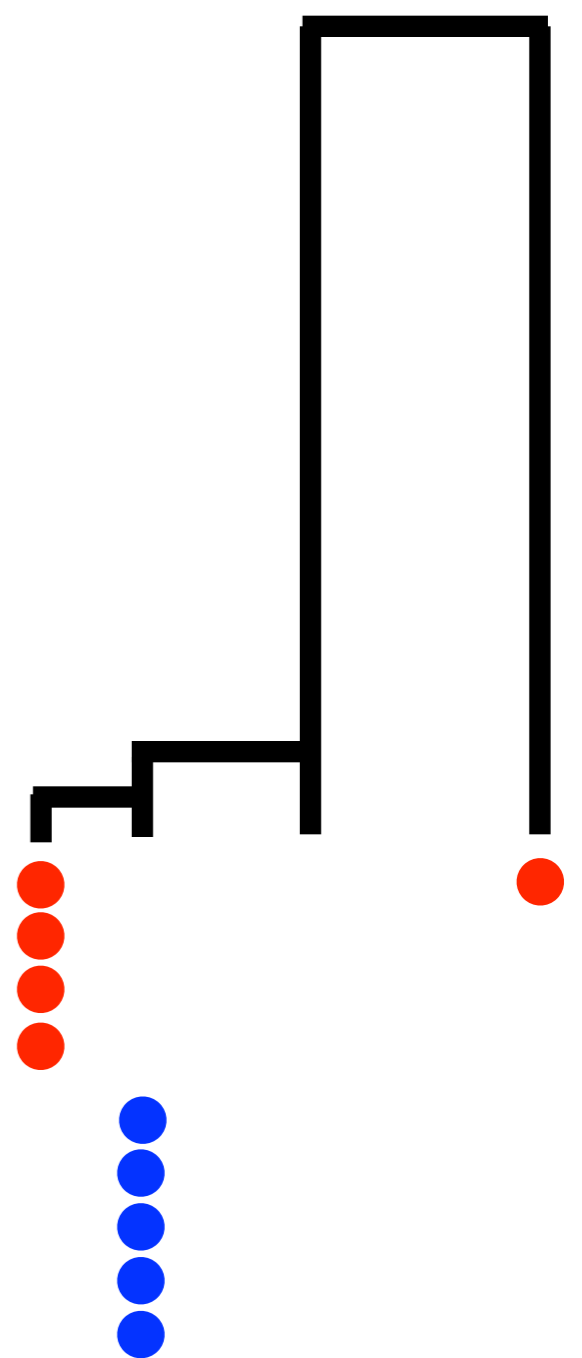Jaccard: d=0
Bray:
Unifrac:
W-Unifrac:

Jaccard: Distant
Bray:
Unifrac:
W-Unifrac:
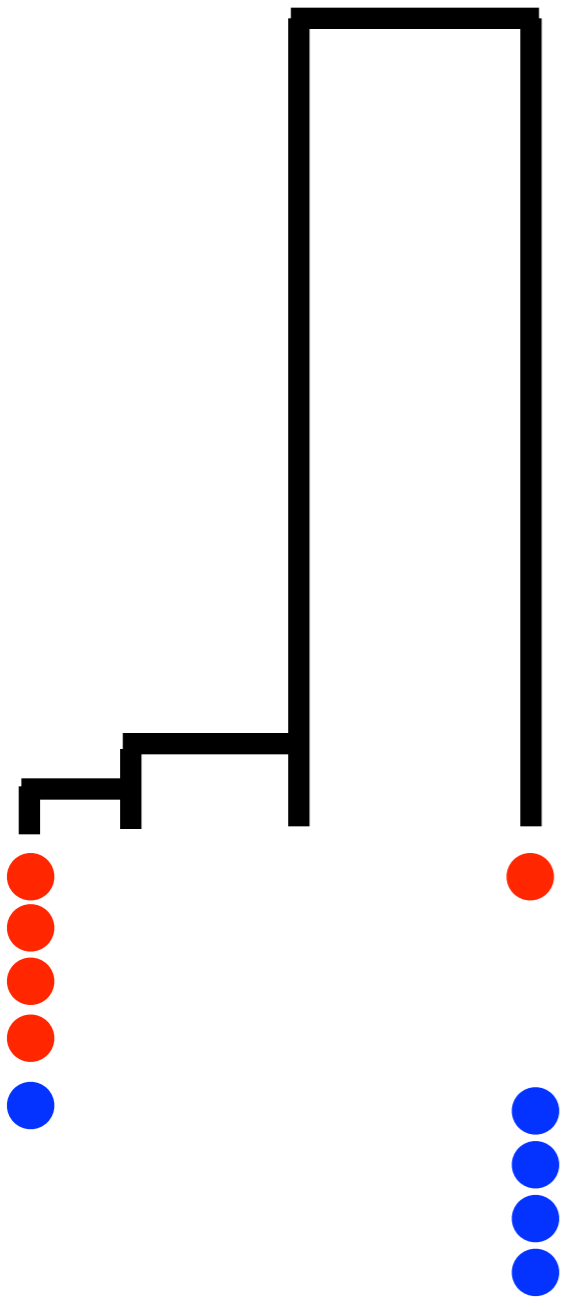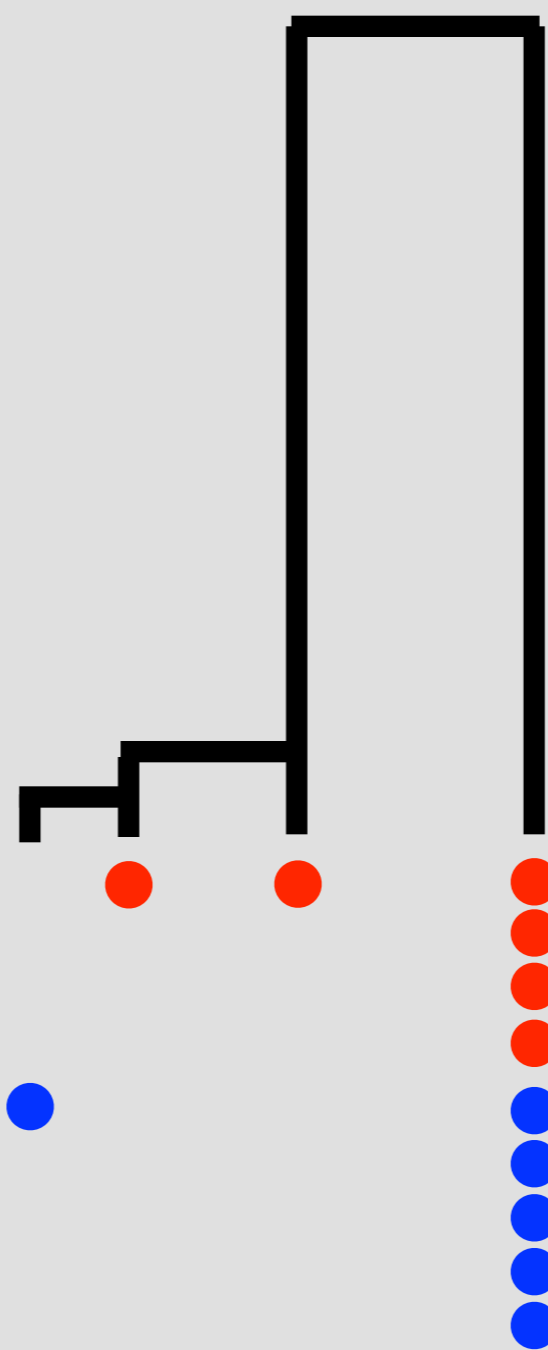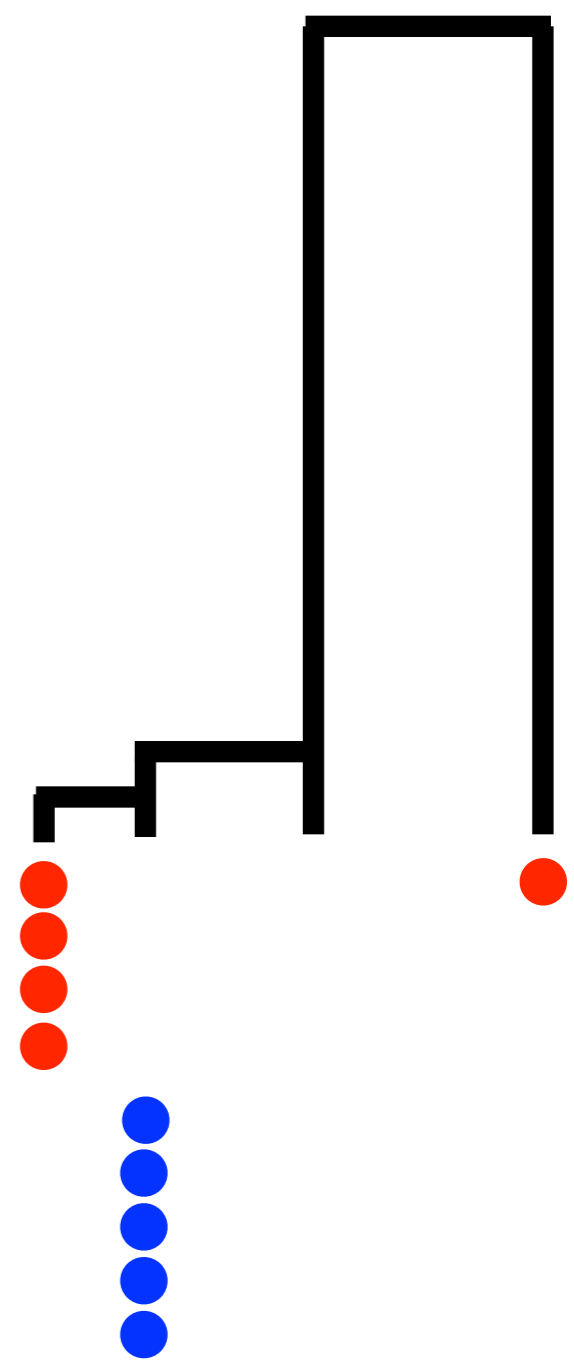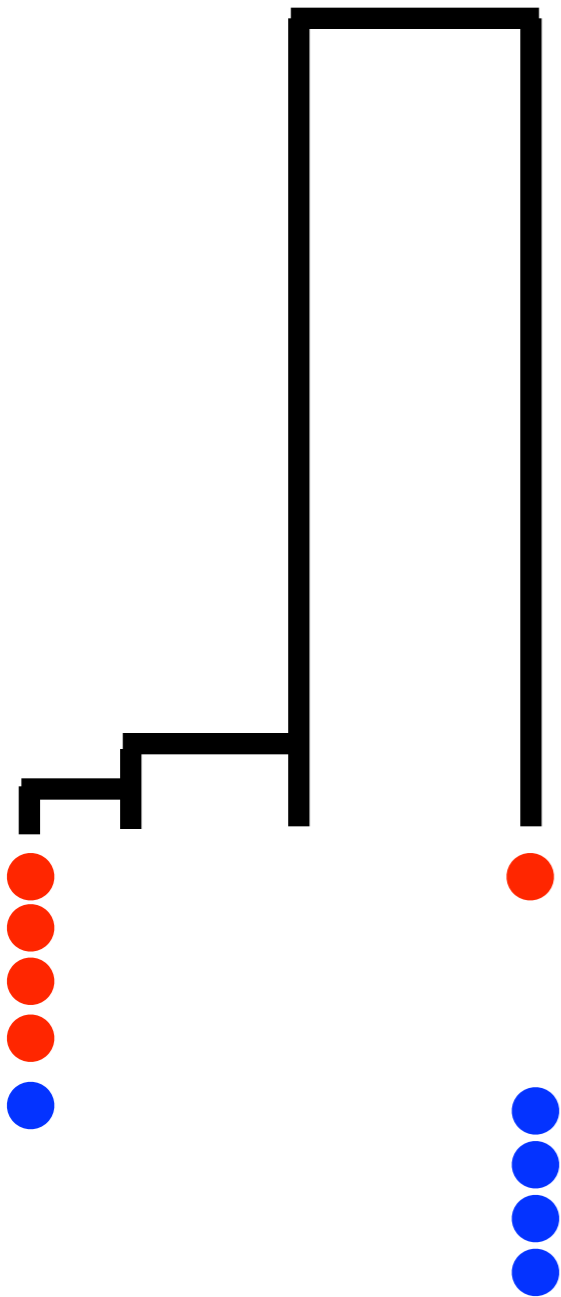
Jaccard: Distant
Bray:
Unifrac:
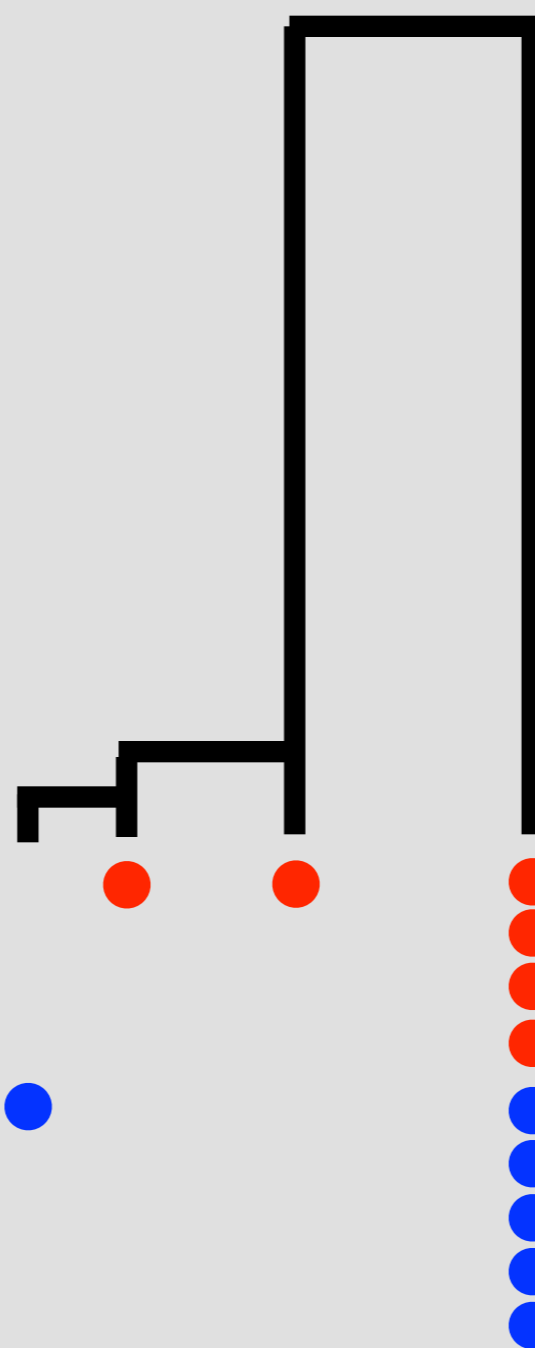W-Unifrac:

Jaccard: d=0
Bray: Distant
Unifrac:
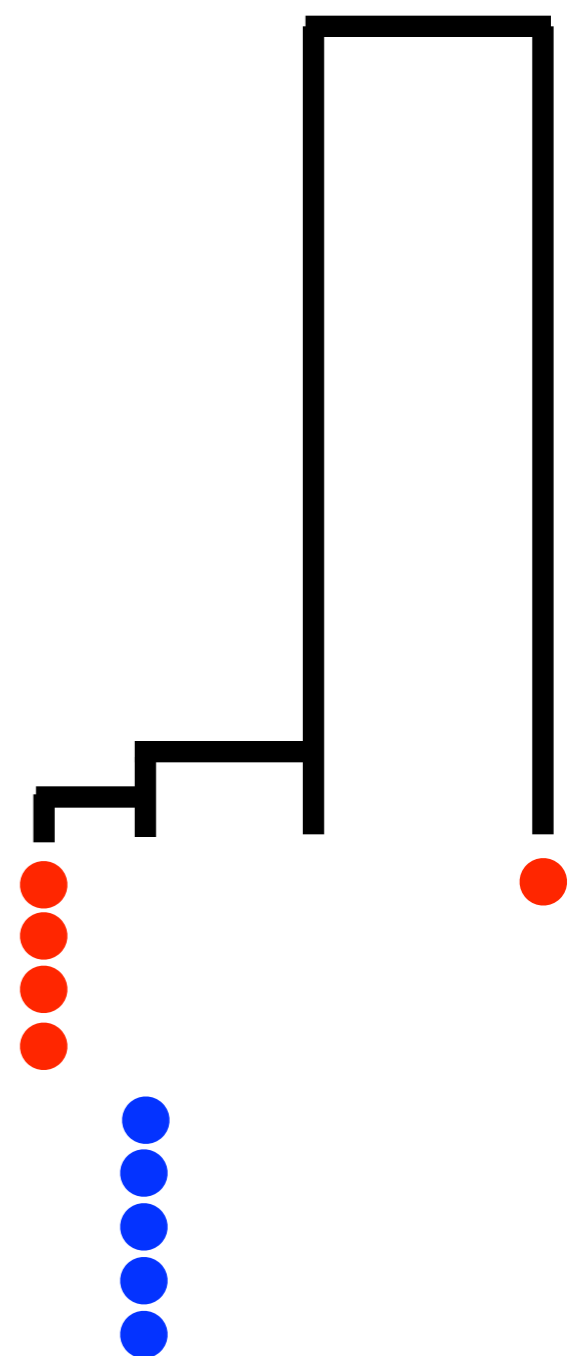W-Unifrac:

Jaccard: Distant
Bray: Similar
Unifrac:
W-Unifrac:

Jaccard: Distant
Bray: Distant
Unifrac:
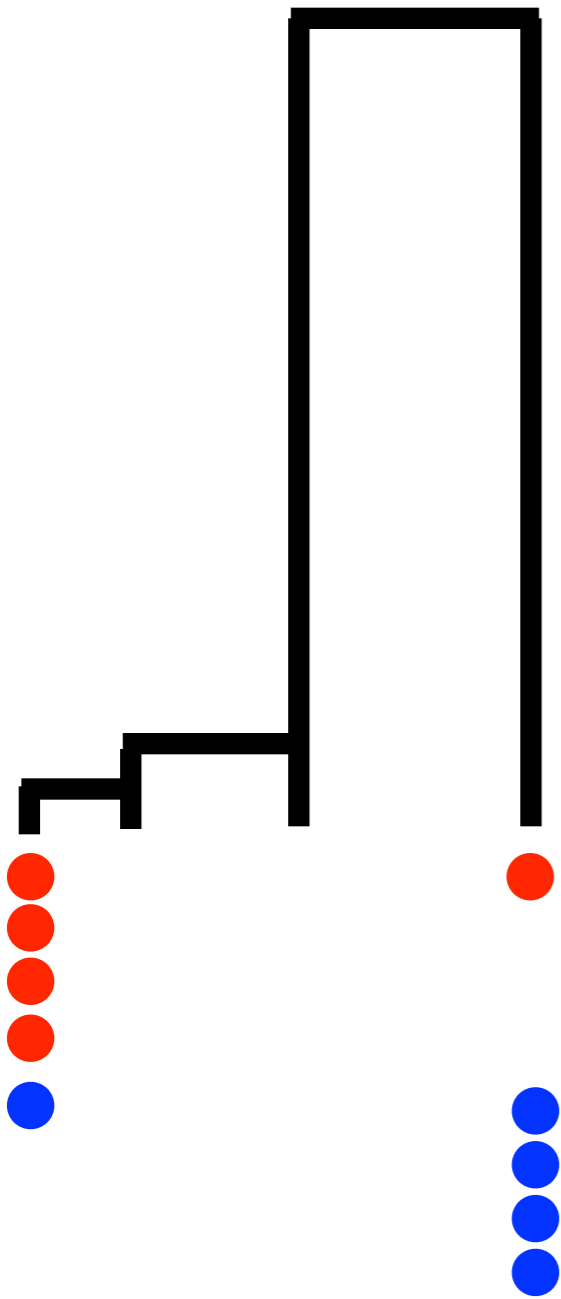W-Unifrac:

Jaccard: d=0
Bray: Distant
Unifrac: d=0
W-Unifrac:
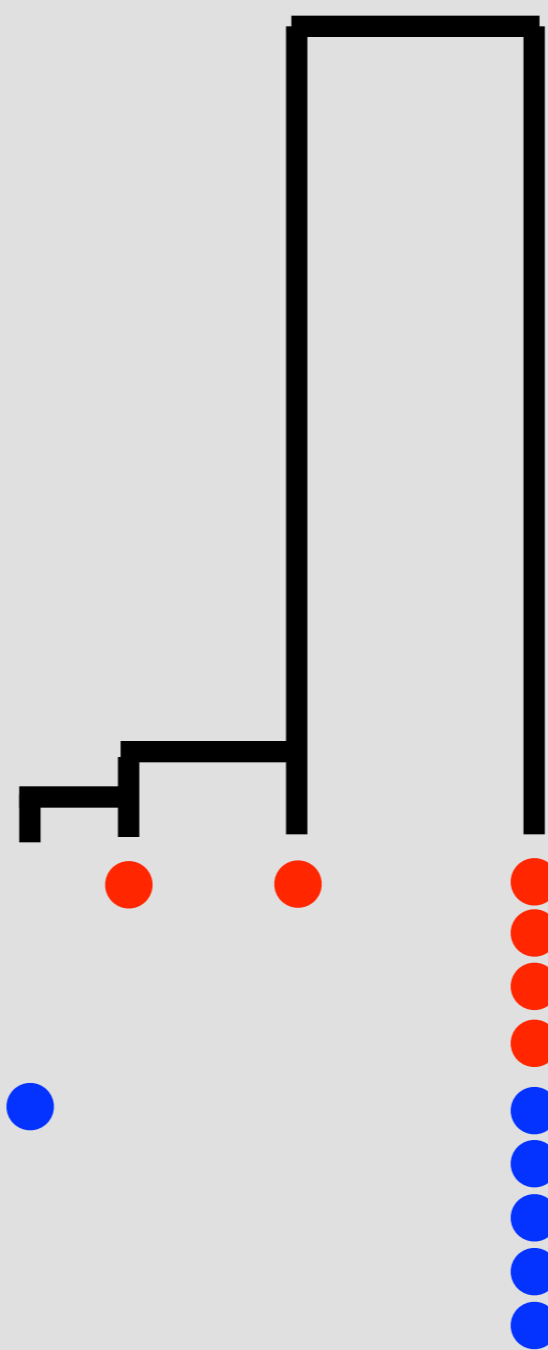
Jaccard: Distant
Bray: Similar
Unifrac: Similar
W-Unifrac:

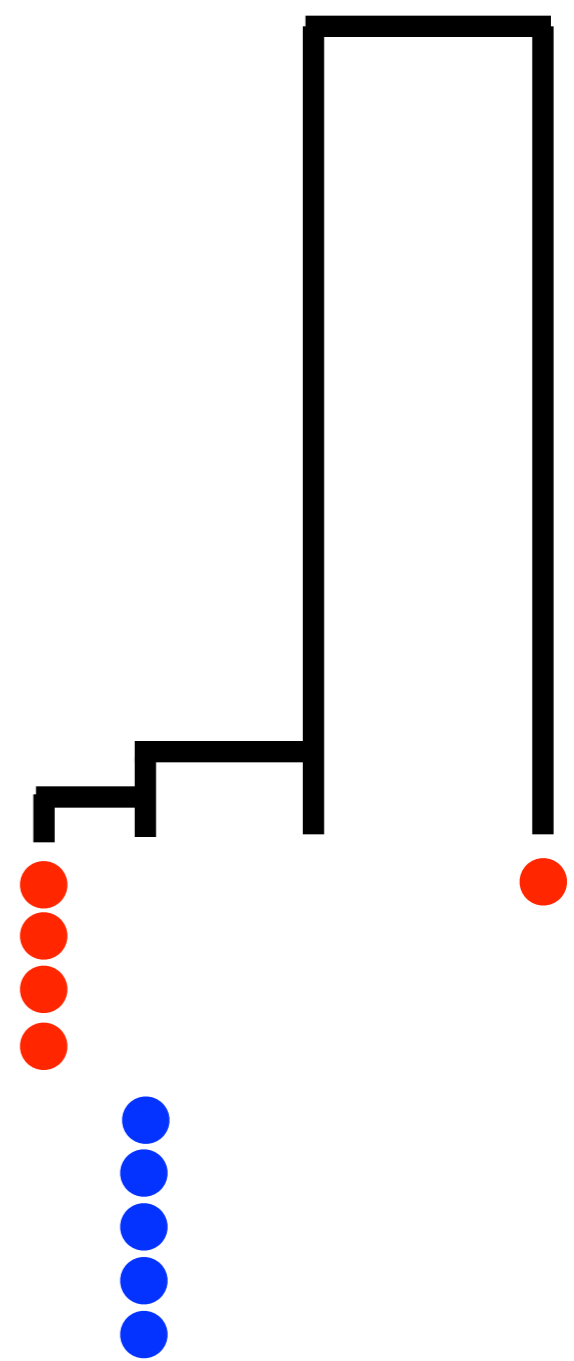Jaccard: Distant
Bray: Distant
Unifrac: Distant
W-Unifrac:

Jaccard: d=0
Bray: Distant
Unifrac: d=0
W-Unifrac: Distant

Jaccard: Distant
Bray: Similar
Unifrac: Similar
W-Unifrac: Similar

Jaccard: Distant
Bray: Distant
Unifrac: Distant
W-Unifrac: Similar

# The Distance Spectrum



|  | Categorical | Phylogenetic |
|---|---|---|
| Presence/Absence | Jaccard | Unifrac |
| Quantitative Abundance | Bray-Curtis | Weighted Unifrac |

**phyloseq distances**
manhattan
euclidean
canberra
bray
kulczynski
jaccard
gower
altGower
morisita-horn
mountford
raup
binomial
chao
cao
jensen-shannon
unifrac
weighted-unifrac
...

# Ordination Methods

Project high-dimensional data onto lower dimensions

**P taxa**

N samples
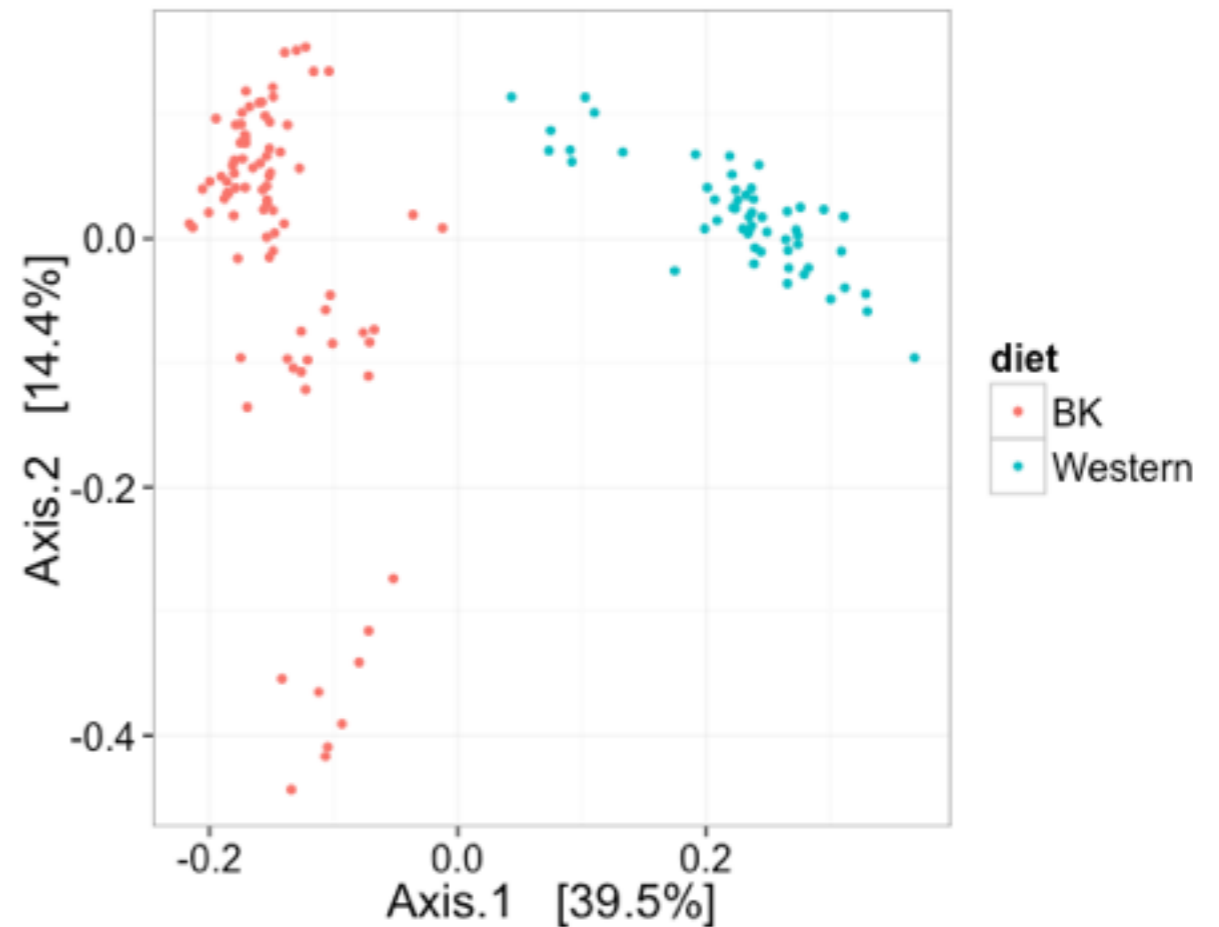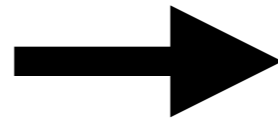
0,1,5,1,0,1,2,1,0,0,9,...
7,2,0,0,0,0,0,0,1,0,0,...
0,0,0,0,0,0,8,0,0,0,1,...
0,0,0,1,0,1,2,0,0,0,5,...
0,1,0,2,0,0,0,1,0,0,4,...
0,0,0,1,9,1,2,5,2,0,1,...
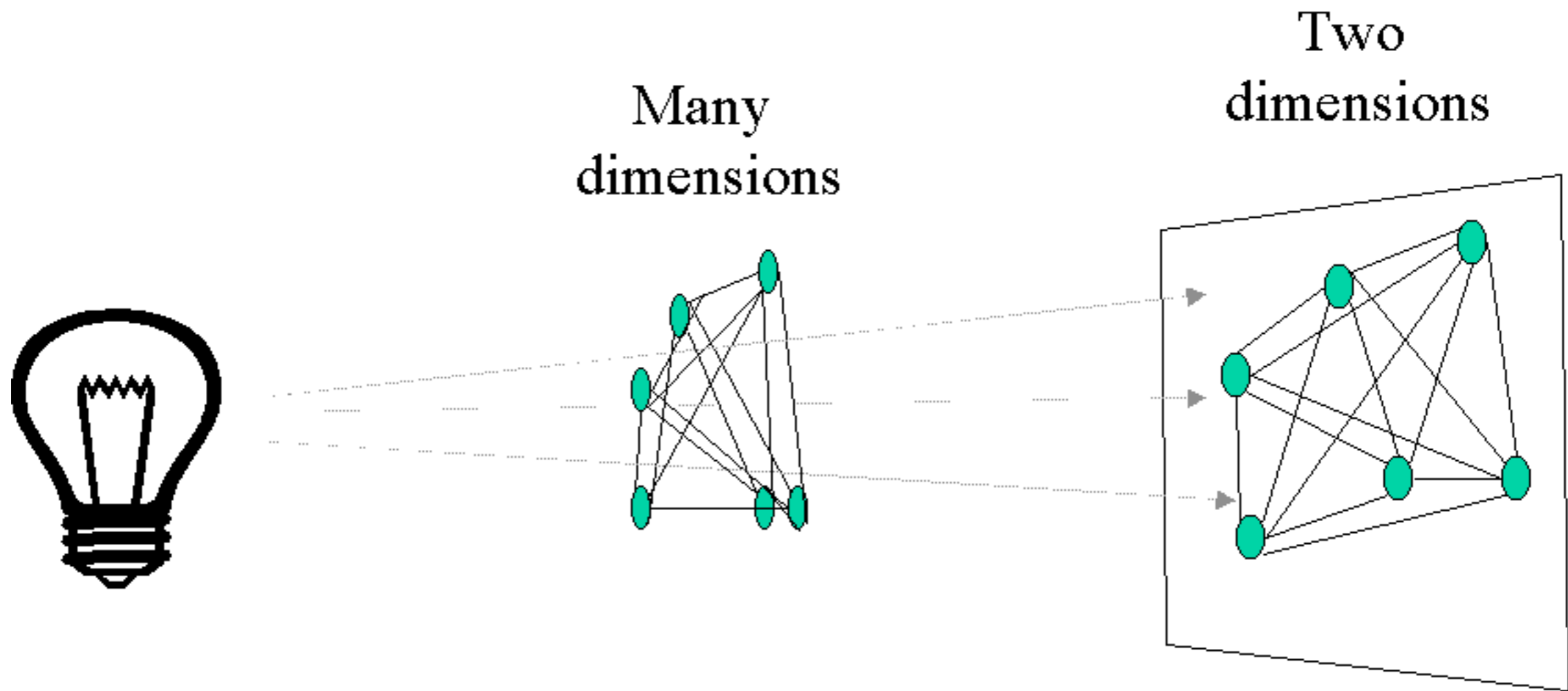0,0,0,0,0,1,2,1,8,0,0,...
0,0,0,0,9,4,0,0,0,0,1,...
.
.

P-dimensions



2-dimensions

# Multi-dimensional Scaling

Why MDS? It works with any distance!



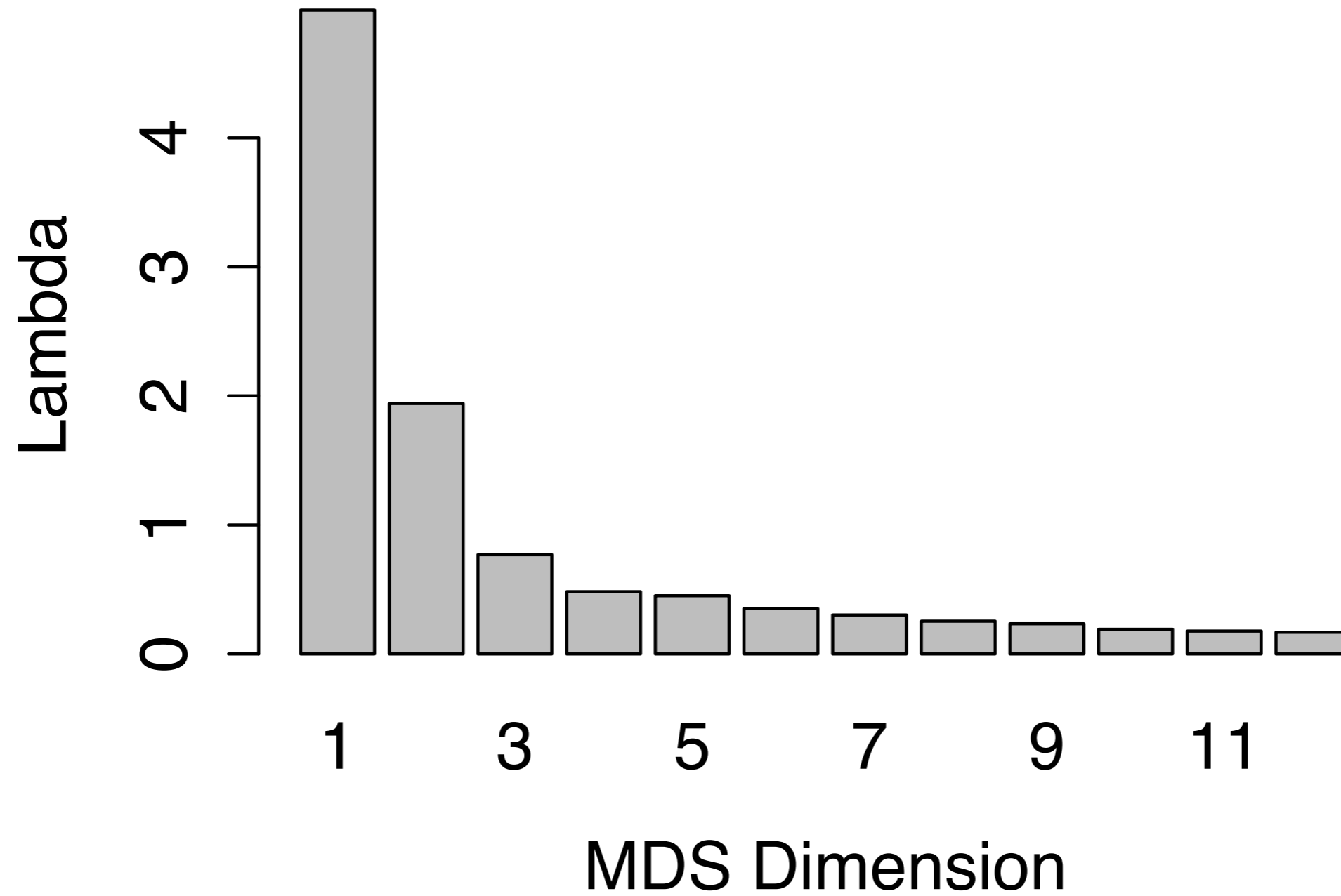Input distance matrix can by Bray-Curtis, Unifrac, …

# MDS Details

Given distances between each observation (sample), MDS finds the closest approximation of that in lower dimensional Euclidean space.

- Algorithm starts from **D** inter-point distances:
  - Center the rows and columns of the distance matrix:
    $$\mathbf{S} = -1/2 \; \mathbf{H} \; \mathbf{D}^{(2)} \; \mathbf{H}$$
  - Compute SVD by diagonalizing S: $\mathbf{S} = \mathbf{U} \; \mathbf{\Lambda} \; \mathbf{U}^{\mathsf{T}}$
  - Extract Euclidean representations: $\mathbf{X} = \mathbf{U} \; \mathbf{\Lambda}^{1/2}$
- The relative values of diagonal elements of $\Lambda$ gives the proportion of variability explained by each of the axes.
- The valued of $\Lambda$ should always be looked at in deciding how many dimensions to retain

---

NMDS is similar, but minimizes a different function (difference in distance ranks)

**MDS Scree Plot**

# Exploratory Analysis

- Looking for patterns (the "I-test")

- Use multiple distances

- phyloseq makes this easy!