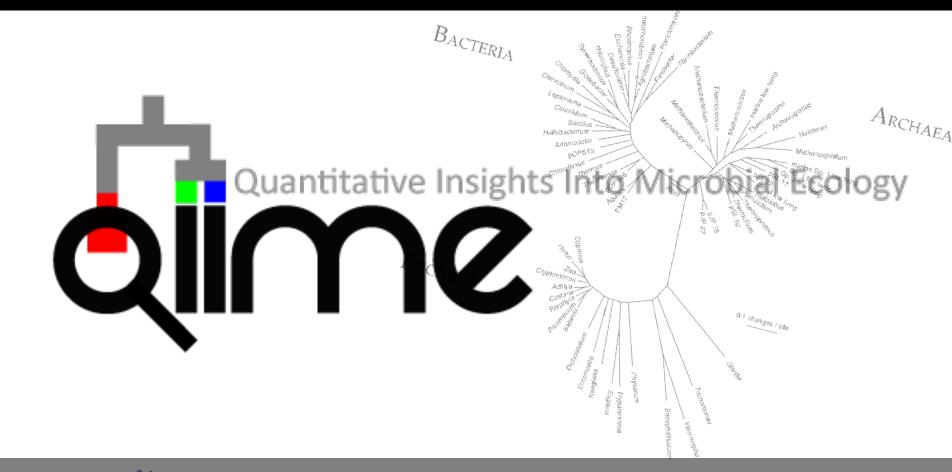
Microbial Communities Profiling via QIIME





Today's schedule

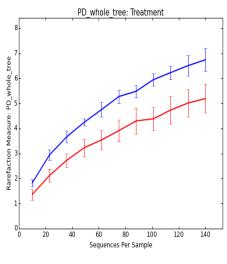
• Day 3.

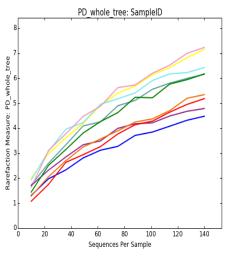
9:00am – 3:00pm: Diversity Analyses and Running Commands

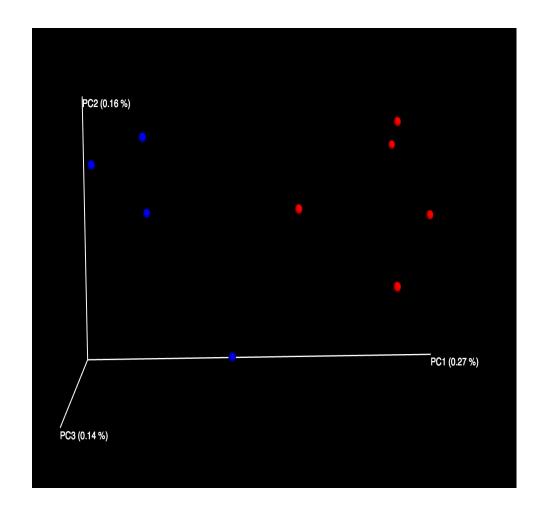
3:30 – 5:00pm: Discuss Group Project

Diversity Analyses

Diversity analysis







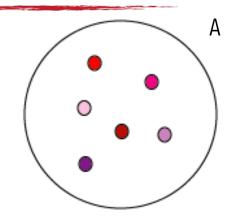
How do we describe and compare diversity?

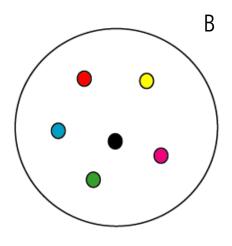
• α Diversity:

- "How many species are in a sample?"
 - (e.g. 6 colors in A and 6 in B)
- e.g.: Are polluted environments less diverse than pristine?

β Diversity:

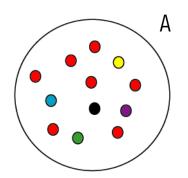
- "How many species are shared between samples?"
 - (e.g. 2 shared colors between A and B)
- e.g.: Does the microbiota differ with different disease states?

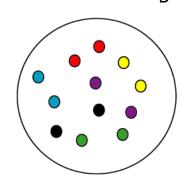




Quantitative versus Qualitative measures

- Qualitative: Considers presence absence only
 - $-\alpha$: How many species are in a sample?
 - e.g.: 6 colors in both A and B.
 - $-\beta$: How many species are shared between samples?
 - e.g.: A and B are identical because the same colors are present in both.
- Quantitative: Also considers relative abundance.
 - $-\alpha$: Accounts for "evenness":
 - e.g. B, where the population is evenly distributed across the 6 species, is more diverse than A, where all species are present but red dominates.
 - β : Samples will be considered more similar if the same species are numerically dominant versus rare.
 - e.g. B and A no longer look identical because of differences in abundance.





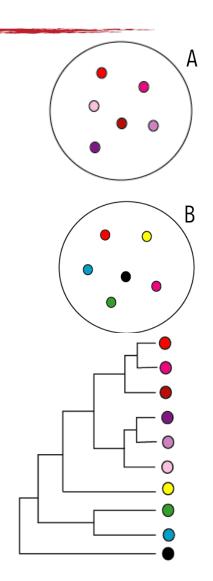
What is a phylogenetic diversity measure?

α Diversity:

- Taxon: "How many species are in a sample?"
- Phylogenetic: "How much phylogenetic divergence is in a sample?"
 - (e.g. B more individually diverse than A more divergent colors)

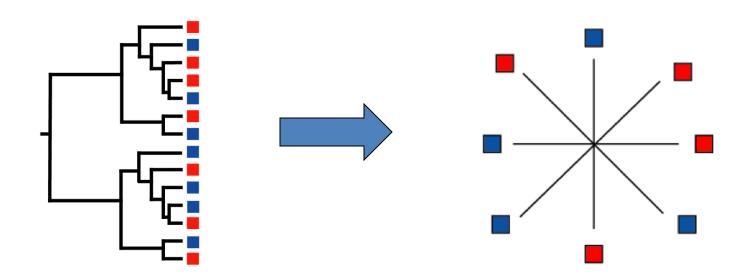
β Diversity:

- Taxon: "How many species are shared between samples?"
- Phylogenetic: "How much phylogenetic distance is shared between samples?"
 - (only related colors from B are in A)



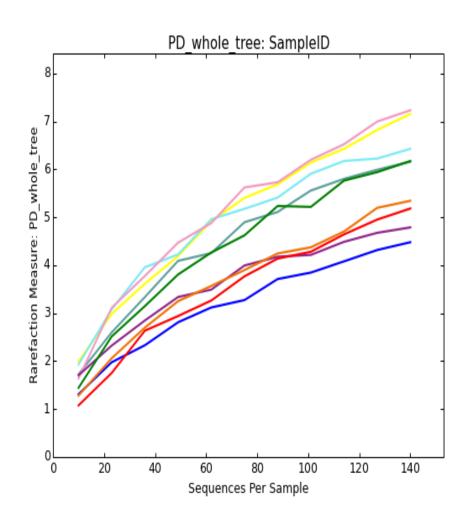
Advantages of phylogenetic techniques.

- Phylogenetically related organisms are more likely to have similar roles in a community.
- Taxon-based methods assume a "star phylogeny"
- Phylogeny and Taxon-based methods can be complementary.



Alpha diversity

Alpha diversity



Basic alpha diversity measure: count number of OTUs.

other measures can be:

- phylogenetic (PD)
- estimators (chao1)
- other statistics (evenness)

PD Rarefaction

- Plot the amount of branch length against the # of observations.
- Shape of curve allows for estimating how far we are from sampling all of the phylogenetic diversity.
- Allows for comparison of phylogenetic diversity between samples.

Comparing diversity

Sample A

Pseudomonas aeruginosa Pseudomonas argentinensis Pseudomonas flavescens

Sample C

Pseudomonas aeruginosa Giardia lamblia Methanobrevibacter smithii

Sample B

Pseudomonas aeruginosa Pseudomonas argentinensis Escherichia coli

Observed species

Sample A

Pseudomonas aeruginosa Pseudomonas argentinensis Pseudomonas flavescens

Sample B

Pseudomonas aeruginosa Pseudomonas argentinensis Escherichia coli

Sample C

Pseudomonas aeruginosa Giardia lamblia Methanobrevibacter smithii

Observed species:

Sample A 3

Sample B 3

Sample C 3

Conclusion:

Samples A, B, and C are equally diverse.

Sample A

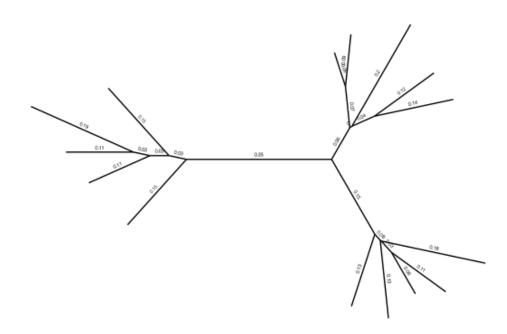
Pseudomonas aeruginosa Pseudomonas argentinensis Pseudomonas flavescens

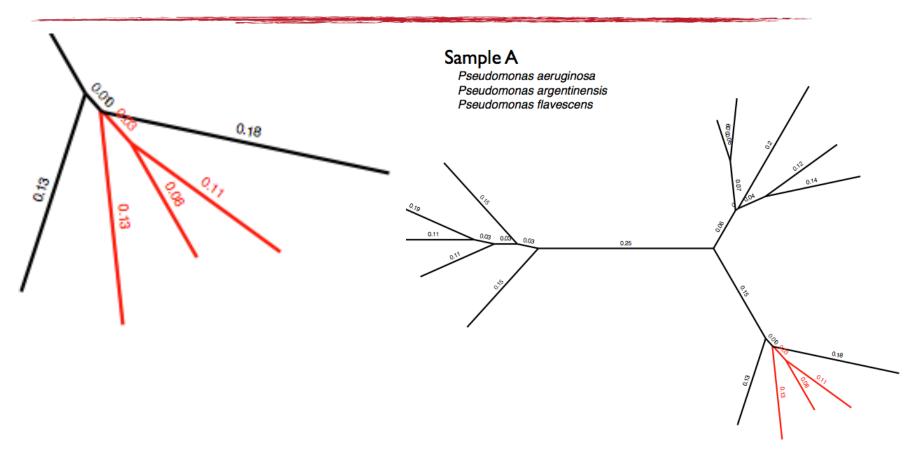
Sample B

Pseudomonas aeruginosa Pseudomonas argentinensis Escherichia coli

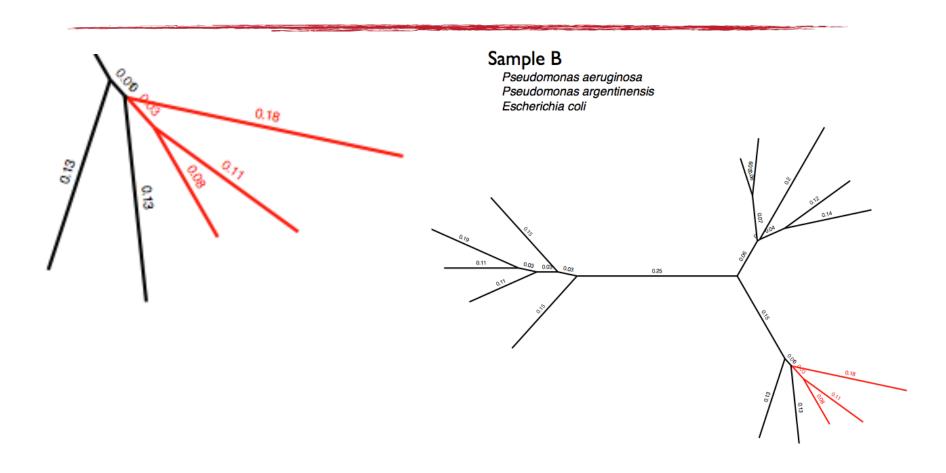
Sample C

Pseudomonas aeruginosa Giardia lamblia Methanobrevibacter smithii

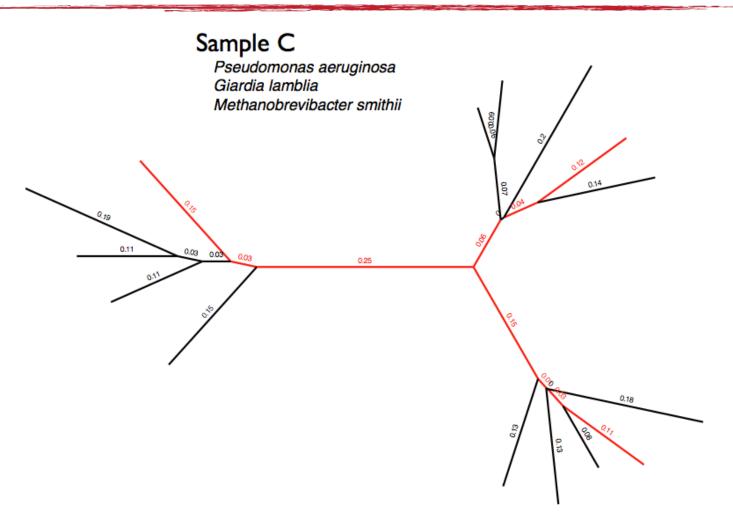




PD = 0.13 + 0.03 + 0.11 + 0.08 = 0.35



PD = 0.18 + 0.03 + 0.11 + 0.08 = 0.40



PD = 0.15+0.03+0.25+0.06+0.04+0.12+0.15+0.01+0.03+0.11=0.95

Sample A

Pseudomonas aeruginosa Pseudomonas argentinensis Pseudomonas flavescens

PD = 0.35

Sample B

Pseudomonas aeruginosa Pseudomonas argentinensis Escherichia coli

PD = 0.40

Sample C

Pseudomonas aeruginosa Giardia lamblia Methanobrevibacter smithii

PD = 0.95

Sample A

Pseudomonas aeruginosa Pseudomonas argentinensis Pseudomonas flavescens

Sample B

Pseudomonas aeruginosa Pseudomonas argentinensis Escherichia coli

Sample C

Pseudomonas aeruginosa Giardia lamblia Methanobrevibacter smithii

$$PD = 0.35$$
 < $PD = 0.40$ < $PD = 0.95$

Conclusion:

Sample C is more diverse than sample B, which is more diverse than sample A.

Beta diversity

Should you rarefy?







RESEARCH ARTICLE

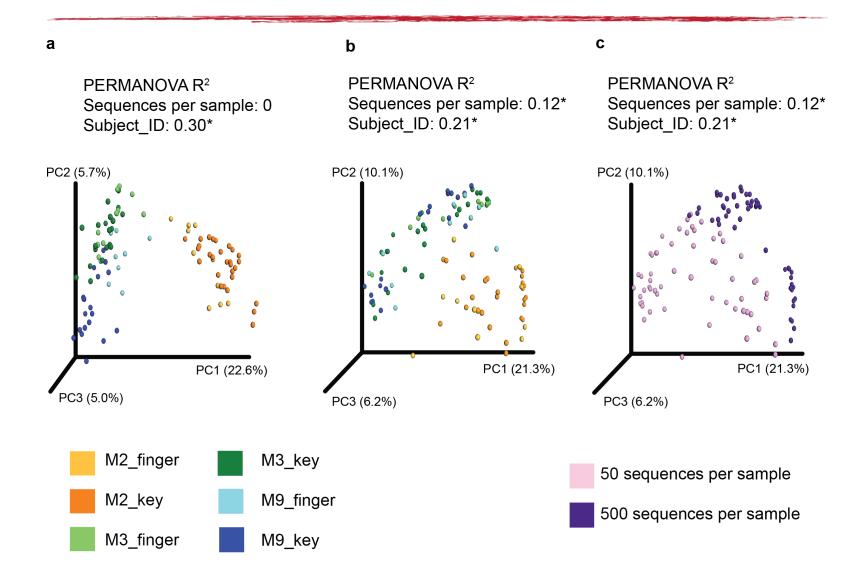
Waste Not, Want Not: Why Rarefying Microbiome Data Is Inadmissible

Paul J. McMurdie, Susan Holmes

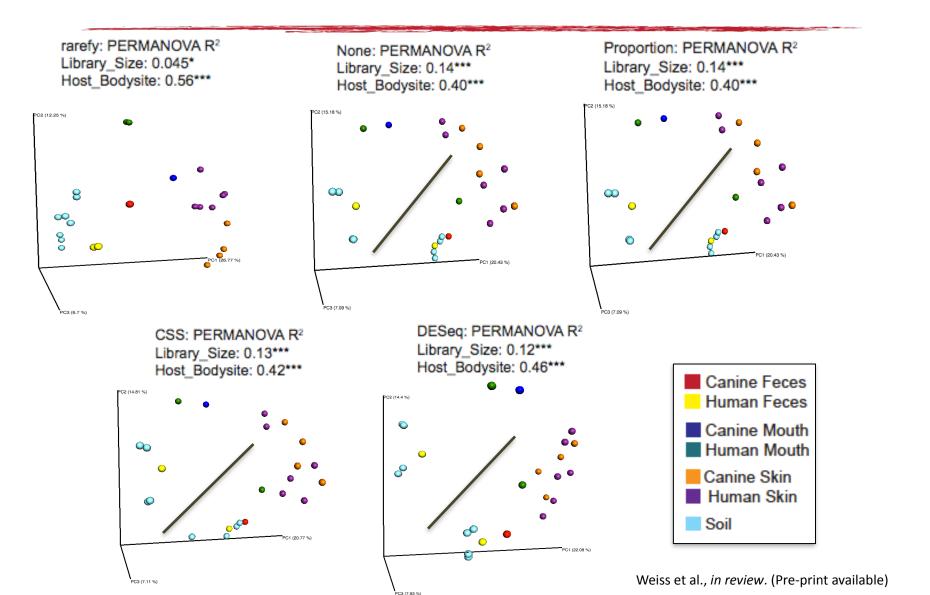
Published: April 3, 2014 • DOI: 10.1371/journal.pcbi.1003531

213	12
Saves	Citations
20,729	84
Views	Shares

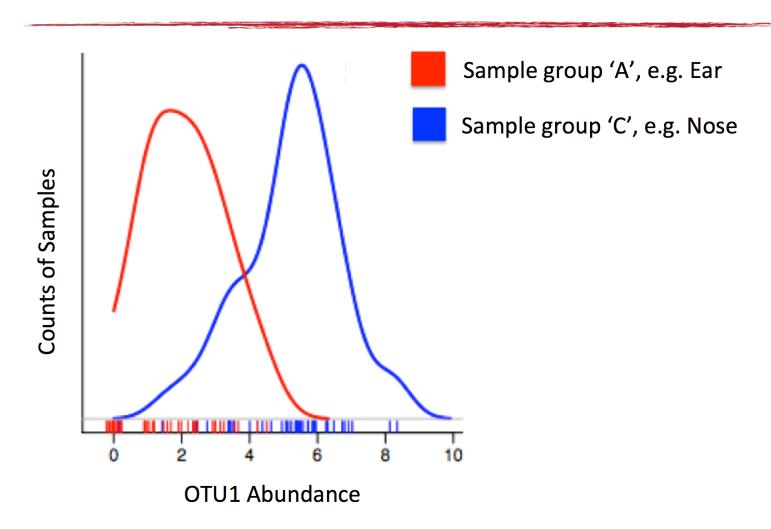
Data normalization



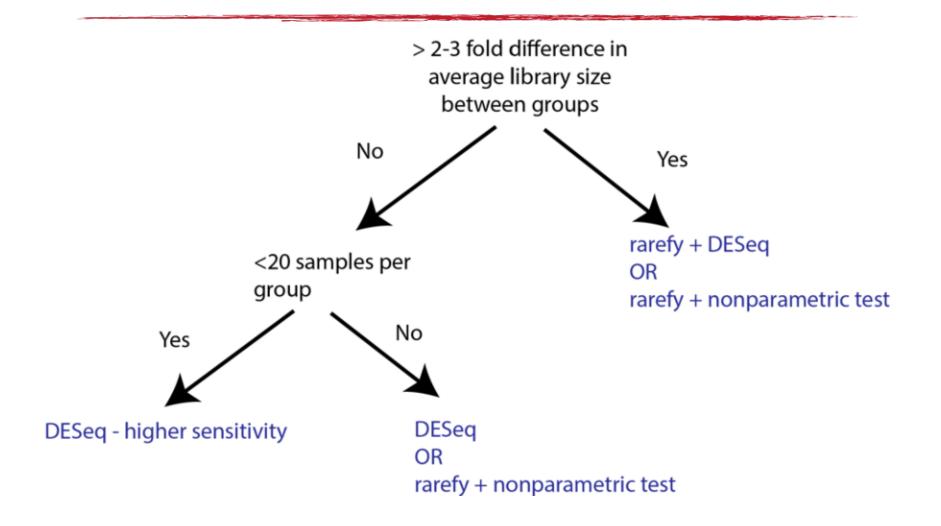
Normalization Methods on Actual Data, Unweighted UniFrac



Differential Abundance Testing: Finding Significant Changes in Abundance



Proposed Differential Abundance Testing Strategy

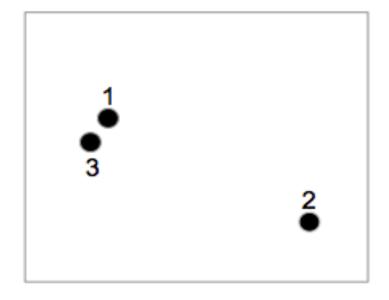


Just to be clear

Normalization ≠ Differential abundance testing

Beta diversity

	1	2	3
OTU1	4	0	4
OTU2	4	0	4
OTU3	0	7	1
OTU4	0	7	0
•••			



Principal Coordinates Analysis (PCoA)

- distance based multivariate approach
- each sample becomes one point, and the closer together points are,
 the more similar the samples they represent are
- dependent on positions of all other points

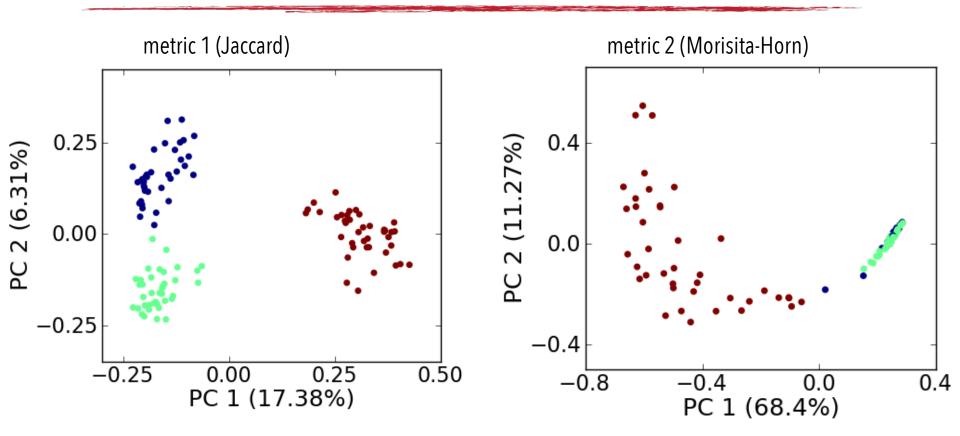
name	formula	re
bray-curtis	$D_{ab} = rac{\sum_{j} A_{aj} - A_{bj} }{\sum_{j} (A_{aj} - A_{bj})}$	[1
ratio (dran	$\sum_{j}(A_{bj}-A_{bj})$	'*
canbeira	$D_{ab} = \frac{1}{Num_{ach}} \sum_{b} \frac{A_{aj} - A_{bj}}{A_{ad} + A_{bd}}$	[1]
(0110).1 (11	$Num_{m'h} \stackrel{\text{\tiny def}}{=} A_{aj} + A_{bj}$	1.
	$\bigcap \left(1 \left(A_{ini} - A_{ini}\right)^2\right)$	
χ^2 (chi-squared)	$D_{ab} = \sqrt{A_{++} \sum_i \left(rac{1}{A_{-j}} \left(rac{A_{nj}}{A_{a+}} - rac{A_{b,j}}{A_{b+}} ight)^2 ight)}$	[1
	/ (11-) (114- 114-)	
	A_{ii} A_{ii} A_{ii}	
chord	$D_{ab} = \sqrt{\sum_{j} \left(rac{A_{aj}}{\sqrt{\sum_{j} A_{aj}^2}} - rac{A_{bj}}{\sqrt{\sum_{j} A_{bj}^2}} ight)^2}$	[2]
cuclidean	$D_{\alpha b} = \sqrt{\sum_{j} (A_{\alpha j} - A_{bj})^2}$	1
460m-Ar	$D_{idi} = \sum_{j} rac{ A_{aj} - A_{bj} }{(max_{oldsymbol{arphi}_{j}} A_{ij} - min_{oldsymbol{arphi}_{i}} A_{ij})}$	[1
gower	$= \frac{D_{igi} - \sum_{i} \overline{(max_{\forall i} A_{ij} - min_{\forall i} A_{ij})}}{(max_{\forall i} A_{ij} - min_{\forall i} A_{ij})}$	'-
hellinger	$D_{nh} = \sqrt{\sum_i \left(\sqrt{rac{A_{ng}}{A_{00}}} - \sqrt{rac{A_{hg}}{A_{00}}} ight)^2}$	[3]
114—π <u>Σ</u> 41	$ = \frac{D_{nb} - \sqrt{\sum_{i} \sqrt{\sqrt{A_{ni}}} - \sqrt{A_{b+1}}}}{\sqrt{A_{b+1}}} $	"'
	$1/\sqrt{\sum \min\{A \in A_{i,j}\}} = \sum \min\{A \in A_{i,j}\}$	
kulczynski (quantitative)	$D_{ab} = 1 - \frac{1}{2} \left(\frac{\sum_{j} min(A_{nj}, A_{nj})}{\sum_{j} A_{aj}} - \frac{\sum_{j} min(A_{nj}, A_{nj})}{\sum_{j} A_{nj}} \right)$	[1]
manhattan	$D_{ab} = \sum A_{ag} - A_{bg} $	[1]
	$2\sum^{l}A_{ml}A_{kl}$	
morisita-horn	$\frac{2\sum_{j}^{T}A_{nj}A_{bj}}{\left(\frac{\sum_{j}A_{nj}^{2}}{\sum_{j}A_{nj}}-\frac{\sum_{j}A_{bj}^{2}}{\sum_{j}A_{bj}}\right)\sum_{j}A_{nj}*\sum_{j}A_{bj}}$	[4
	$\left(rac{\sum_{j}A_{nj}}{\sum_{j}A_{nj}} - rac{\sum_{j}A_{nj}}{\sum_{j}A_{nj}} ight)\sum_{j}A_{nj}*\sum_{j}A_{nj}$	
	$D_{ab} = 1 - \frac{\sum_{j} (A_{aj} - (\frac{A_{aj}}{M})) \sum_{j} (A_{bj} - (\frac{A_{bj}}{M}))}{\sqrt{\sum_{i} (A_{aj} - (\frac{A_{bj}}{M}))^{2}} \sqrt{\sum_{i} (A_{aj} - (\frac{A_{bj}}{M}))^{2}}}$	
bearson	$D_{ab} = 1 - \frac{1}{\sqrt{\sqrt{\chi_2} \chi_4} + (A_{a-1} \chi_2)} \frac{\pi}{\sqrt{\chi_2} \chi_4} \frac{\pi}{\chi_2} \frac{\pi}{\sqrt{\chi_2} \chi_4} \frac{\pi}{\chi_2} \frac{\pi}{\sqrt{\chi_2} \chi_4} \frac{\pi}{\sqrt{\chi_2} \chi_4} \frac{\pi}{\sqrt{\chi_2} \chi_4} \chi$	[5]
soergel	$D_{\alpha\delta} = rac{\sum_{j} A_{\alpha j} - A_{\delta j} }{\sum_{j} max(A_{lpha j}, A_{\delta j})}$	[6
	$\sum_{j} \max(A_{\alpha j}, A_{\delta j})$	
spearrign	$D_{ab} = 1 - \frac{6\sum_{j}\overline{\left(\left(rank\left(A_{aj}\right) - rank\left(A_{bj}\right)\right)^{2}}}{M(M-1)}$ $D_{ab} = \sqrt{\sum_{j}\left(\frac{A_{aj}}{A_{a+}}\right) - \frac{A_{bj}}{A_{b+}}\right)^{2}}$	[7]
	M(M-1)	
. married accountille	$\mathbf{p}_{i,j} = \left[\sum \left(A_{0,j} - A_{0,j}\right)^2\right]$	[2
species profile	$A_{lab} = \sqrt{\frac{2}{A_{n+1}}} \sqrt{\frac{A_{n+1}}{A_{n+1}}}$	[4
	qualitative methods	
χ^{3} (qualitative)	see quantitative y ³	see above
chord	see quantitative enord	see above
cuclidean	see quantitative cuclidean	see above
hamming	$D_{ab} = Num_a + Nvm_{\underline{b}} + 2Nvm_{a \cap b}$, see also manhactan	[8
jaccard	$D_{ab} = 1 (\frac{Nnm_{aCb}}{})$, see also socred	
	$Num_a - Num_b - Num_{a \cap b}$	'
lennon	$D_{ab} = 1 - \frac{N_{acm-1} - min(N_{min} - N_{acm-1}, N_{acm-1}, N_{acm-1})}{N_{acm-1} - min(N_{min} - N_{acm-1}, N_{acm-1}, N_{acm-1})}$	[10]
1 55	$Num_{arb} = \frac{1 \times am_{a + b} - avam_{a + b} + 2 \times am_{b} + 2 \times am_{b$	for a
ochiai	$\begin{aligned} D_{ab} &= Num_a + Nvm_b + 2Nnm_{acb} \text{, see also manhattan} \\ D_{ab} &= 1 - (\frac{Num_{acb}}{Num_a + Num_b + Num_{acb}}) \text{, see also soergel} \\ D_{ab} &= 1 + \frac{Num_{acb} + Num_{acb}}{Num_{acb} + uin(Num_a + Num_{acb}, Num_b + Num_{acb})} \\ &= \frac{D_{ab} - 1 + \frac{Num_{acb}}{\sqrt{Num_aNum_b}}}{\sqrt{Num_aNum_b}} \end{aligned}$	[11]
pearson	see quantitative pearson	see above
	$D_{ab} = 1 - rac{2Num_{acto}}{Num_a + Num_b}$	l rest
dice	$D_{ch} = 1 - \frac{1}{100}$	[12]

A variety of community dissimilarity measures exist

name	formula	ref
bray-curtis	$D_{ab} = rac{\sum_{g} \left[A_{a_3} - A_{b_d} - \sum_{g} \left(A_{a_1} - A_{b_f} \right) ight]}{\sum_{g} \left(A_{a_1} - A_{b_f} \right)}$	[1]
canberra	$D_{ab} = \frac{1}{Nmn_{ab}} \sum_{i} \frac{A_{aj} - A_{bj}}{A_{aj} + A_{bj}}$	D.
χ^2 (chi-squared)	$D_{ab} = \sqrt{A_{++} \sum_{i} \left(rac{1}{A_{-i}} \left(rac{A_{aj}}{A_{a-}} - rac{A_{bj}}{A_{b+}} ight)^2 ight)}$	[1
chord	$D_{ab} = \sqrt{\sum_{j} \left(rac{A_{aj}}{\sqrt{\sum_{j} A_{aj}^{2}}} - rac{A_{bj}}{\sqrt{\sum_{j} A_{bj}^{2}}} ight)^{2}}$	[2
cuclidean	$D_{\alpha b} = \sqrt{\sum_i (A_{\alpha j} - A_{bj})^2}$	ΙŢ
Rower	$D_{cb} = \sum_{j} rac{ A_{aj} - A_{bj} }{(max_{S_s}A_{ij} - min_{oldsymbol{arphi}_i}A_{ij})}$	[1]
he <u>llinge</u> r	$D_{nh} = \sqrt{\sum_{j} \left(\sqrt{rac{A_{ng}}{A_{0}}} - \sqrt{rac{A_{hj}}{A_{0+}}} ight)^2}$	[3]
kulczynski (quantitative)	$D_{ab} = 1 - \frac{1}{2} \left(\frac{\sum_{j} min(A_{aj}, A_{bj})}{\sum_{j} A_{aj}} - \frac{\sum_{j} min(A_{aj}, A_{bj})}{\sum_{j} A_{bj}} \right)$	[1]
manhautan	$D_{ab} = \sum A_{aj} - A_{bj} $	[1]
morisita-horn	$\frac{2\sum_{j}^{j}A_{\alpha j}A_{bj}}{\left(\frac{\sum_{j}A_{\alpha j}^{2}}{\sum_{j}A_{\alpha j}}-\frac{\sum_{j}A_{\alpha j}^{2}}{\sum_{j}A_{\alpha j}}\right)\sum_{j}A_{\alpha j}*\sum_{j}A_{bj}}$	[4]
реаткоп	$D_{ab} = 1 - \frac{\sum_{j} (A_{aj} - (\frac{A_{aj}}{M})) \sum_{j} (A_{bj} - (\frac{A_{bj}}{M}))}{\sqrt{\sum_{i} (A_{aj} - (\frac{A_{bj}}{M}))^{2}} \sqrt{\sum_{i} (A_{aj} - (\frac{A_{bj}}{M}))^{2}}}$	[5.
socrgel	$D_{ab} = rac{\sum_{j} \left[A_{aj} - A_{bj} - \sum_{j} max(A_{aj}, A_{bj}) ight]}{\sum_{j} max(A_{aj}, A_{bj})}$	[6]
spearman	$D_{ab} = 1 - \frac{6\sum_{j}\left(\left(rank(A_{aj}) - rank(A_{bj})\right)^{2}}{M(M-1)}$	17.
species profile	$D_{ab} = \sqrt{\sum_{J} \left(\frac{A_{aJ}}{A_{a+}}\right) - \frac{A_{bJ}}{A_{b+}}}\right)^2}$	[2]
	qualitative methods	
χ^{3} (qualitative)	see quantitative χ^2	see above
chord	see quantitative chord	sec above
cuclidean	see quantitative cuclidean	see above
hamming	$D_{nb} = Num_n + Nvm_b + 2Nnm_{arib}$, see also manhactan Nnm_{arib} ,	[8]
jaccard	$D_{ab} = Num_a + Num_b + 2Num_{arb}$, see also manhattan $D_{ab} = 1 - (\frac{Num_{arb}}{Num_a + Num_b + Num_{arb}})$, see also socrael $D_{ab} = 1 - \frac{Num_b + Num_{arb}}{Num_{arb}}$	9
lennon	$D_{ab} = 1 - \frac{Nem_{a \neg b} - min(Num_a - Nem_{a \neg b}, Nem_b - Num_{a \neg b})}{Nem_{a \neg b} - min(Num_a - Num_{a \neg b}, Nem_b - Num_{a \neg b})}$	[10]
ochiai	$D_{ab} = 1 - \frac{Num_a - Num_b - Num_{a \cap b}}{Num_{a \cap b} - min(Num_a - Num_{a \cap b}, Num_b - Num_{a \cap b})}$ $D_{ab} = 1 - \frac{Num_{a \cap b}}{\sqrt{Num_b Num_b}}$	[11]
pearson	see quantum peresson	
dice	$D_{ab}=1-rac{2Num_{a^{\prime\prime\prime}a}}{Num_{a}-Num_{b}}$	[12]
	1 2 - Dette if 2 - Georgi	1

Is metric choice important?

Absolutely



Fierer et al. PNAS 2010

NATURE METHODS | ANALYSIS

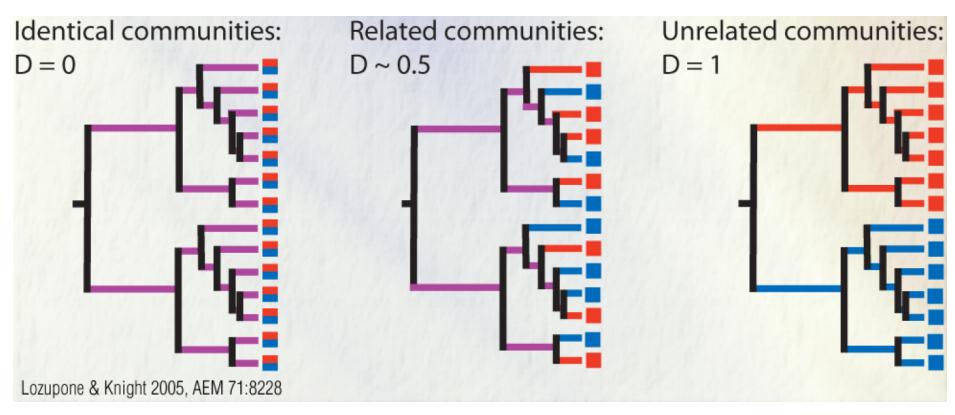




Microbial community resemblance methods differ in their ability to detect biologically relevant patterns

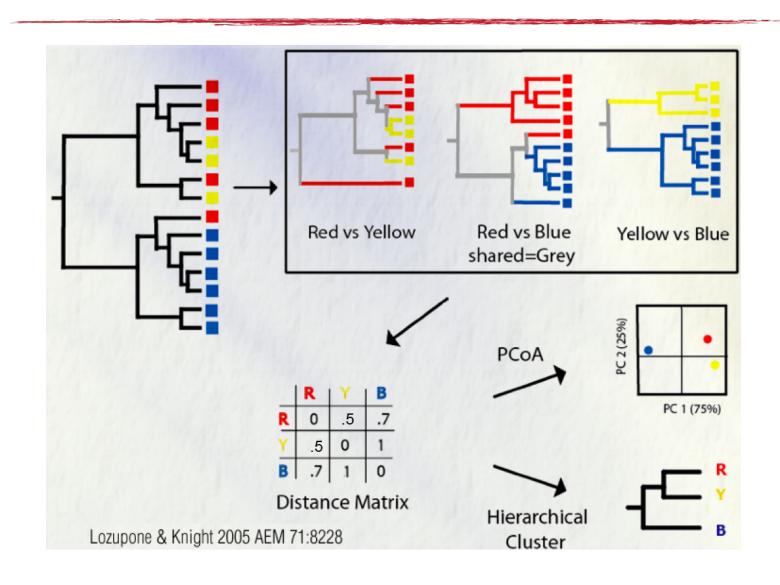
Justin Kuczynski, Zongzhi Liu, Catherine Lozupone, Daniel McDonald, Noah Fierer & Rob Knight

UniFrac metric can be used in PCoA

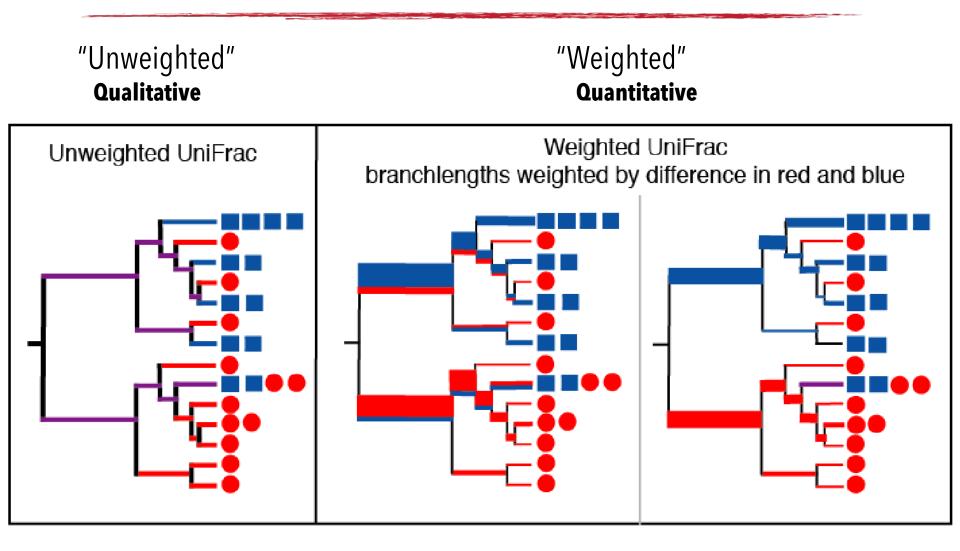


(cited over 1500 times)

UniFrac distance matrix



UniFrac



Lozupone et al., 2007. Appl Environ Microbiol 73:1576

Core diversity analysis

```
core_diversity_analyses.py
-i <biom format otu table>
-m <mapping file>
-t <tree file>
-o <output directory>
-e < rarefaction level>
-c < mapping category/categories >
-aO < number of cores if parallel >
```

What's actually happening:

alpha_rarefaction.py

```
multiple_rarefactions.py
alpha_diversity.py
collate_alpha.py
make_rarefaction_plots.py
```

beta_diversity_through_plots.py

```
beta_diversity.py
principal_coordinates.py
make_emperor.py
```

summarize_taxa_through_plots.py

```
summarize_otu_by_cat.py
summarize_taxa.py
plot_taxa_summary.py
```

- make_distance_boxplots.py
- compare_alpha_diversity.py
- group_significance.py

Running Illumina Tutorial