

CENTER FOR INDIVIDUALIZED MEDICINE

Introductory Microbiome Statistics Jun Chen (chen.jun2@mayo.edu)



Characteristics of microbiome data



Data characteristics

- High-dimensional counts with variable library sizes
- Overdispersion, excessive zeros and outliers
- Phylogeny and hierarchical structure
- Compositional effects





Multi-level summarization of microbiome data

- Raw count data
 - OTU/ASV level
 - Different taxonomic level
- Normalized/Proportion data
 - OTU/ASV level
 - Different taxonomic level
- Diversity (community-level) data
 - α -diversity
 - β -diversity



MAYO CLINIC



Diversity measures quantify communitylevel properties

- α -diversity: a measure of within-sample diversity (richness and evenness)
 - Not dependent on the composition of other samples
 - Sobs, Chao1, Shannon, Simpson, Phylogenetic diversity (PD)
- β-diversity: a measure of between-sample diversity in terms of pairwise ecological distances
 - Measures how different one community is to another
 - Bray-Curtis, Jaccard, UniFrac distance
- Can be split along two major axes:
 - Unweighted vs. weighted measure
 - Phylogenetic vs. Non-phylogenetic measure





Each diversity measure captures a different aspect of the community properties

- Unweighted measure community membership/structure; Weighted measure - community compositional profile;
- Unweighted measure also puts more emphasis on rare/low-abundance species than weighted measure
- Phylogenetic measure captures "clustered" signals
- It is important to study several representative diversity measures to have a full view of the community





Each diversity measure captures a different aspect of the community properties







Different richness, same evenness (Sobs, Chao1)

Same richness, different evenness (Shannon, Simpson)

Same richness, same evenness, different phylogeny (Phylogenetic diversity)



UniFrac distance (Knight group) captures phylogenetically related community difference

Unweighted UniFrac (d^U) : the fraction of the branch length of the tree that is unique to each community/microbiome sample.



- Most efficient in detecting the difference in community membership
- Also difference in rare lineages, differs in probability of being picked up by the sequencing machine

Weighted UniFrac (d^W) : branch length weighted by proportion difference



- Most efficient in detecting difference in abundant lineages
- Absolute difference puts too much weight on abundant lineages

Bray-Curtis/Jaccard distances are UniFrac equivalents without using phylogeny.

Lozupone, C., & Knight, R. (2005). UniFrac: a new phylogenetic method for comparing microbial communities. *Applied and environmental microbiology*, *71*(12), 8228-8235.



Diversity depends on the sampling effort (sequencing depth)

As we sequence more deeply, we will detect more rare species.





Differential sequencing depth could produce false patterns





Diversity measure and rarefaction

- To conduct a fair comparison of diversity measures, esp. for unweighted measures, the sequencing depth should be equal. The process of down-sampling the sequences to the same depth is called rarefaction.
- We usually down-sample the sequences at different depths and create a rarefaction curve. The rarefaction curve can be used to assess the adequacy of the sampling depth.
- Unless rarefied to a very low depth, rarefaction does not lose much statistical power for diversity analysis. It could even increase the power for unweighted analysis.





Rarefaction corrects the false pattern



Center for INDIVIDUALIZED MEDICINE

T



Rarefaction corrects the false pattern





Center for INDIVIDUALIZED MEDICINE

©2012 MFMER | slide-12

Central themes of microbiome studies



Hall, A. B., Tolonen, A. C., & Xavier, R. J. (2017). Human genetic variation and the gut microbiome in disease. *Nature Reviews Genetics*, *18*(11), 690-699. Warnecke, F., & Hugenholtz, P. (2007). Building on basic metagenomics with complementary technologies. *Genome biology*, *8*, 1-5.



Center for INDIVIDUALIZED MEDICINE

©2012 MFMER | slide-13



Center for INDIVIDUALIZED MEDICINE



MAYO

Exploratory data analysis

- Various visualizations through heat map, stacked bar plot, PCoA plot, cladogram
- Exploratory data analysis can
 - Detect outlier samples/features
 - Understand sources of variability, particularly batch effects
 - Detect contaminants by comparing to negative controls
 - Form hypotheses for statistical testing
- Two interactive tools
 - EMPeror (Vzquez-Baeza, et al., 2013, GigaScience): PCoA plot-based.
 - Calour (Xu et al., 2019, mSystems): heat map-based.



Center for INDIVIDUALIZED MEDICINE

Xu, Z. Z., et al. (2019). Calour: an interactive, microbe-centric analysis tool. Msystems, 4(1), 10-1128.

Vázquez-Baeza, Y., et al. (2013). EMPeror: a tool for visualizing high-throughput microbial community data. Gigascience, 2(1), 2047-217X.



Exploratory data analysis (examples)





Center for INDIVIDUALIZED MEDICINE

T

T



Exploratory data analysis (examples)



Time

- 6 wks
- Fecal Slurry

PA Status

- Healthy volunteers
- PI-IBS (High PA)
- PI-IBS (Low PA)





Microbiome association test

- Individual taxon vs community-level test
- Community-level test to establish an overall association
 - α -diversity: close to normal, linear (mixed) model can be applied
 - β -diversity: PERMANOVA and MiRKAT (Zhao et al., 2015, AJHG)
 - Idea: compare microbiome similarity to phenotype similarity
 - PERMANOVA: distance-based regression with the microbiome as the outcome
 - MiRKAT: kernel machine regression with the microbiome as the covariate
 - More powerful methods are under development.
- Individual taxon testing to identify specific taxa (differential abundance analysis)
 - Many methods are available but no consensus exists!
 - Multiple testing correction must be performed!!

	Case samples			Control samples				
Taxon 1	c_{11}	c_{12}		c_{1m}	$c_{1,m+1}$	$c_{1,m+2}$		c_{1n}
Taxon 2	c_{21}	c_{22}	•••	c_{2m}	$c_{2,m+1}$	$c_{2,m+2}$	•••	c_{2n}
•••	•••	•••	•••	• • •	• • •	•••	•••	
Taxon q	c_{q1}	c_{q2}		c_{qm}	$c_{q,m+1}$	$c_{q,m+2}$		c_{qn}
Library size	s_1	s_2		s_m	s_{m+1}	s_{m+2}	•••	s_n







Multiple testing correction

- Assume we are testing 500 taxa, if we use the traditional type I error cutoff 0.05, by definition, we will detect an average of 25 differential taxa while none are truly differential!!
- P-values should be corrected for multiple testing.
 - Family-wise error rate (FWER) control (e.g. Bonferroni correction) controls the probability of making any false discovery in the results. May be too conservative.
 - False discovery rate (FDR) control (e.g. Benjamini-Hochberg procedure) controls the average percentage of false discoveries in the results. Balance of power and false positives.





Challenges of differential abundance analysis

- Microbiome sequencing data are "compositional": the change of the abundance of one taxon will lead to the changes of the relative abundances of other taxa
- The severity of compositional effects depends on the number of differential taxa, their abundance level, and the direction of change



The "relative" abundances (relative to **the total sum**) are changed for all species!



2012 MFMER | slide-20

Assumption for differential abundance analysis based on relative abundance data





Center for INDIVIDUALIZED MEDICINE

Approaches to address compositional effects

• Select a suitable "divisor"

- Use some robust normalizing factor (CSS, GMPR, CLR, etc.) to capture the sequencing effort of the non-differential part (sparse signal)
- Find a single or a set of reference taxa, which are assumed to be relatively invariant with respect to some condition (DACOMP, RAIDA, ZicoSeq)
- Bias correction (ANCOM-BC, LinDA, fastANCOM)

Yang, L., & Chen, J. (2022). A comprehensive evaluation of microbial differential abundance analysis methods: current status and potential solutions. *Microbiome*, *10*(1), 130.
 Zhou, H., He, K., Chen, J., & Zhang, X. (2022). LinDA: linear models for differential abundance analysis of microbiome compositional data. *Genome biology*, *23*(1), 95.
 Lin, H., & Peddada, S. D. (2020). Analysis of compositions of microbiomes with bias correction. *Nature communications*, *11*(1), 3514.



Center for INDIVIDUALIZED MEDICINE

©2012 MFMER | slide-22

A non-exhaustive list of existing differential abundance analysis methods

Method	Model A	ddress compositional effect	s Handling zero	Covariate
DESeq2	Negative binomial model	Scale factor (RLE)	Nodel (Overdispersion)	Yes
edgeR	Negative binomial model	Scale factor (TMM)	Nodel (Overdispersion)	Yes
	Zero-inflated log-normal model	Scale factor		
metagenomeSeq		(CSS/Wrench)	Model (Zero-inflation)	Yes/No
RAIDA	Zero-inflated log-normal model	Reference	Model (Zero-inflation)	No
ANCOM	Wilcoxon	Pairwise log ratio	Pseudo-count	No
ANCOM-BC	Log linear model	Bias correction	Pseudo-count	Yes
fastANCOM	Log linear model	Bias correction	Pseudo-count	Yes
LinDA	Log linear model	Bias correction	Pseudo-count	Yes
MaAsLin2	Log linear model	Scale factor	Pseudo-count	Yes
DACOMP	Wilcoxon	Reference	Not neccesary	No
ZINQ	Logistic+quantile regression	Scale factor	Model (Logistic)	Yes
LDM	Linear model	Scale factor (TSS)	Not necessary	Yes
ZicoSeq	Linear model	Reference	Empirical Bayes	Yes
MicrobiomeDDA		Scale factor (GMPP)		
(mbzinb)	Zero-inflated negative binomial mo	odel	Model (Zero-inflaton)	Yes
Aldex2	Wilcoxon/t test/glm	Scale factor (CLR)	Empirical Bayes	Yes
Beta-binomial (Corncob)	beta-binomial model	Scale factor (TSS)	Model (Overdispersion)	Yes
eBay	Wilcoxon/t test	Scale factor (CLR)	Empirical Bayes	No



LinDA: linear models for differential abundance analysis with bias correction (Zhou et al., 2022, Genome Biology)

https://cran.r-project.org/web/packages/MicrobiomeStat/

True model based on absolute abundance X

$$\log (X_{is}) = u_s \alpha_i + (1, \mathbf{c}_s^{\top}) \boldsymbol{\beta}_i + \epsilon_{is},$$

where u_s is the variable of interest, c_s are covariates, and ϵ_{is} is the error term.

Model based on CLR-transformed relative abundance Y

$$egin{aligned} \mathcal{W}_{is} &:= \log\left\{rac{Y_{is}}{(\prod_{j=1}^m Y_{js})^{1/m}}
ight\} pprox \log\left\{rac{X_{is}}{(\prod_{j=1}^m X_{js})^{1/m}}
ight\} \ &= \log(X_{is}) - rac{1}{m}\sum_{j=1}^m \log(X_{js}) \ &= u_s\left(lpha_i - ar{oldsymbollpha}\right) + (1, \mathbf{c}_s^{ op})\left(oldsymboleta_i - ar{oldsymboleta}
ight) + arepsilon_{is} - ar{arepsilon}_s \end{aligned}$$

where $\bar{\alpha} = m^{-1} \sum_{i=1}^{m} \alpha_i$, $\bar{\beta} = m^{-1} \sum_{i=1}^{m} \beta_i$, and $\bar{\varepsilon}_s = m^{-1} \sum_{i=1}^{m} \varepsilon_{is}$.

The OLS estimator for α based on the CLR-transformed data is **biased** with the bias term being $\bar{\alpha}$. LinDA estimates the bias and performs inference based on the de-biased estimate.



Microbiome-based prediction with Random Forest algorithm

- Algorithm based on decision trees with bootstrap samples and splitting on a random set of features
- Advantages
 - Capture nonlinear effects
 - Capture interaction between taxa
 - Can accommodate a large number of taxa
 - Robust to outliers
 - Provide importance rank of the taxa
 - Can couple with Boruta feature selection to identify biomarker taxa
 - Out-of-bag (OOB) error gives a reasonable error estimate (may overestimate)
- Important parameters
 - Number of trees (default may be small)
 - Number of features to split (default is usually OK)



https://towardsdatascience.com/from-a-singledecision-tree-to-a-random-forest-b9523be65147



Microbiome-based prediction - considerations

- What input to Random Forest?
 - Taxonomic or original OTU/ASV abundance
 - Whether to incorporate phylogenetic information
 - Taxonomic or functional abundance
 - How to normalize the data
 - Whether to perform feature filtering and selection



- Performance evaluation by cross-validation (CV)
 - If all the parameters are set as default, OOB error or traditional CV can be used to evaluate prediction performance.
 - If we tune the parameters and perform feature filtering/selection, twostage nested CV should be used.
 - Feature filtering/selection should be performed in the training data only.



Nested cross-validation for objective assessment of prediction performance





mPower: a power calculation tool for microbiome study design



Please cite the following references if you use mPower:

Yang, L., & Chen, J. (2023). Benchmarking differential abundance analysis methods for correlated microbiome sequencing data. Briefings in bioinformatics, 24(1), bbae607. Yang, L., & Chen, J. (2022). A comprehensive evaluation of microbial differential abundance analysis methods: current status and potential solutions. Microbiome, 10(1), 130. Yang, L., & Chen, J. (2002). A comprehensive evaluation of thor Microbiome Sudy Design (In preparation)

https://microbiomestat.shinyapps.io/mPower/ (ongoing project)



Center for INDIVIDUALIZED MEDICINE

©2012 MFMER | slide-28

MicrobiomeStat: a R package supporting comprehensive microbiome data analysis

SINGLE-POINT ANALYSIS

Introduction

Alpha Diversity Analysis

Beta Diversity Analysis

Feature-level Analysis

One-Click Reports Generation

PAIRED SAMPLES ANALYSIS

Introduction

Alpha Diversity Analysis

Beta Diversity Analysis

Feature-level Analysis

One-Click Reports Generation

LONGITUDINAL ANALYSIS

Introduction

Alpha Diversity Analysis

Beta Diversity Analysis

Feature-level Analysis

One-Click Reports Generation

MicrobiomeStat: Supporting Longitudinal Microbiome Analysis in R



MicrobiomeStat is a specialized tool for analyzing longitudinal and paired microbiome data. It can also be used for longitudinal analysis of other omics datatypes as long as the data are properly normalized and transformed. As a special case, cross-sectional and case-control microbiome data analysis is also supported.

Gitbook LLM Integration

Maximize the potential of MicrobiomeStat by utilizing its Gitbook wiki documentation in conjunction with Gitbook LLM (Large Language Model). In the top right corner of the Gitbook, you'll find the "Ask or Search" feature. Feel free to pose any questions you might have, such as "How to filter microbiome data in MicrobiomeStat?"

Resources

- MicrobiomeStat Repository: <u>GitHub</u>
- MicrobiomeStat Shiny Repository: <u>GitHub</u>
- Shiny Web Application: Link
- MicrobiomeStat Github Pages: Link
- Function Documentation: Link
- Gitbook Wiki: Link

www.microbiomestat.wiki (ongoing development)



Center for INDIVIDUALIZED MEDICINE

©2012 MFMER | slide-29