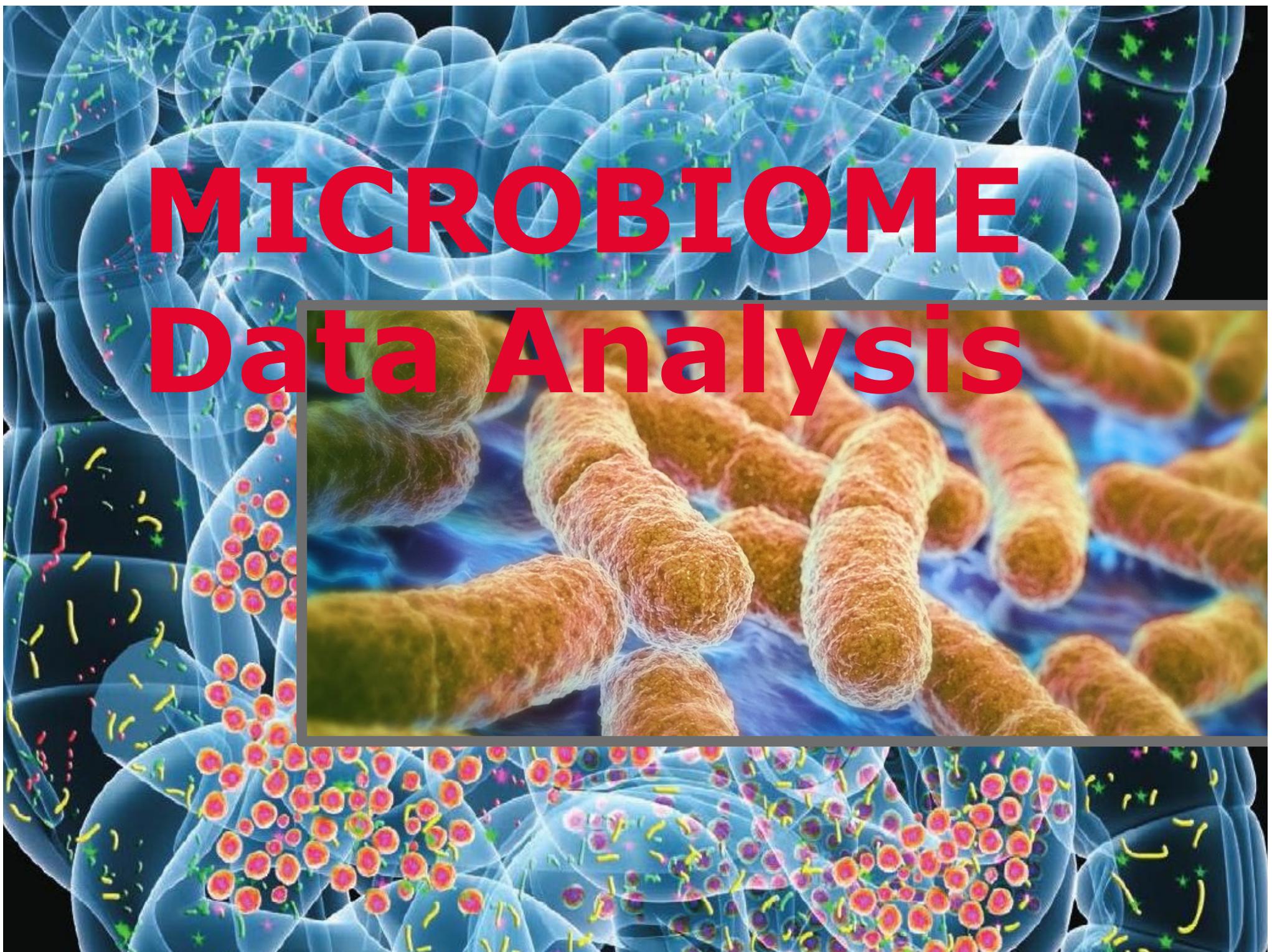


MICROBIOME

Data Analysis



- * A site which points you towards data analysis resources for microbiome data

<https://ucph-foodmicro.github.io/UCPH-FOODMICRO/>

- * The site used for this course including tutorial material and exercises

<https://mortenarendt.github.io/MicrobiomeDataAnalysis/index.html>



Purpose

- To descriptive describe the individual communities.
- To compare with external data
- To integrate with other layers of *-omics* type data.

Problem

- Sparse / Zero inflated
- Compositional
- Over dispersed



Outline

Day1 Morning

Preprocessing

Alpha diversity

Beta diversity

Day2 Morning

Testing design

versus Beta diversity

Day1 Afternoon

Differential

Abundance testing

Day2 Afternoon

Multomics with

heatmaps and CCA



Preprocessing

Problem

- Sparse / Zero inflated
- Compositional
- Over dispersed
- High dimensional

- Normalization/rarefaction
- Filter off rare taxa
- Agglomeration
- Transformation (e.g. log())



Diversity metrics

Alpha diversity

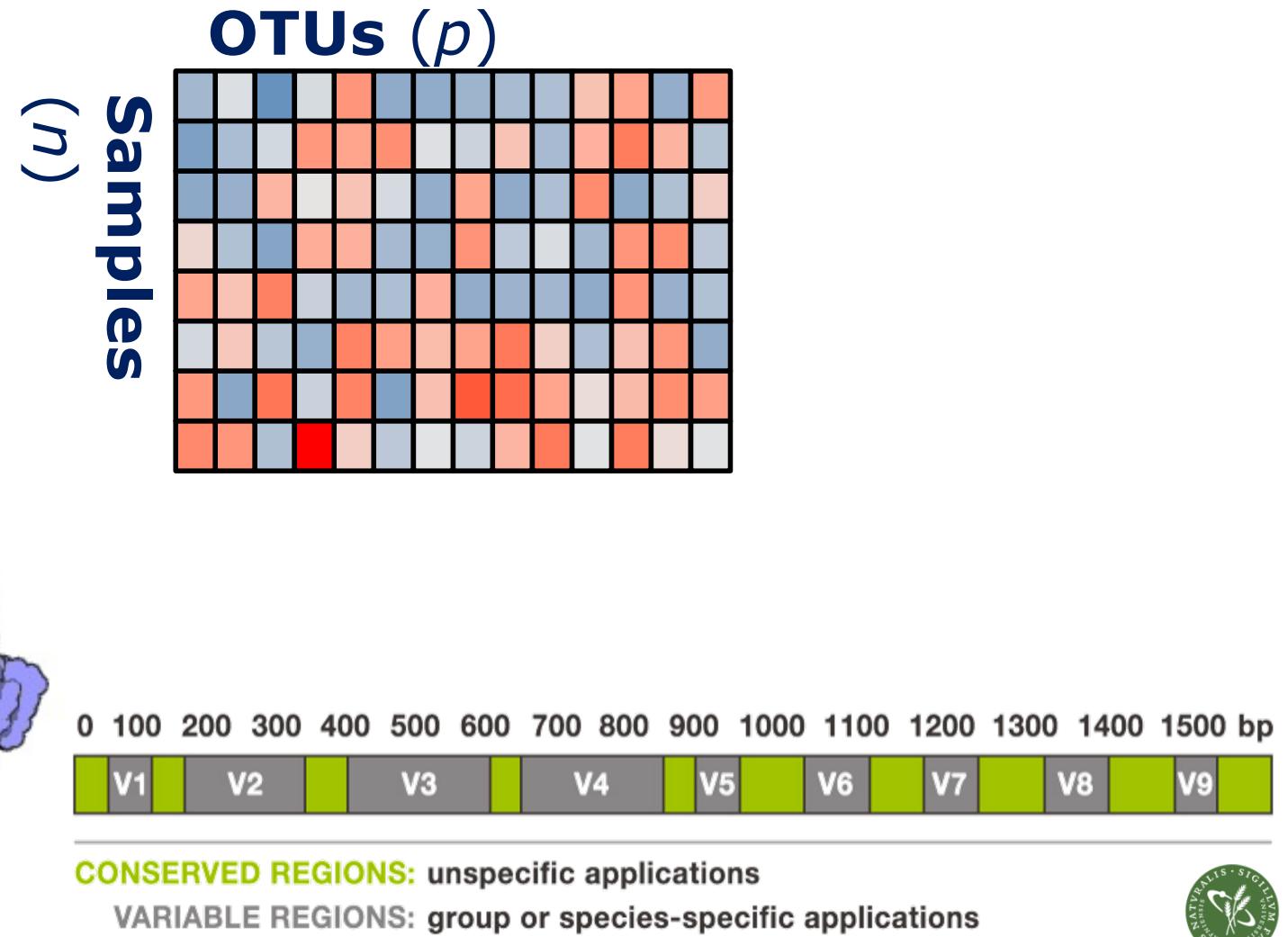
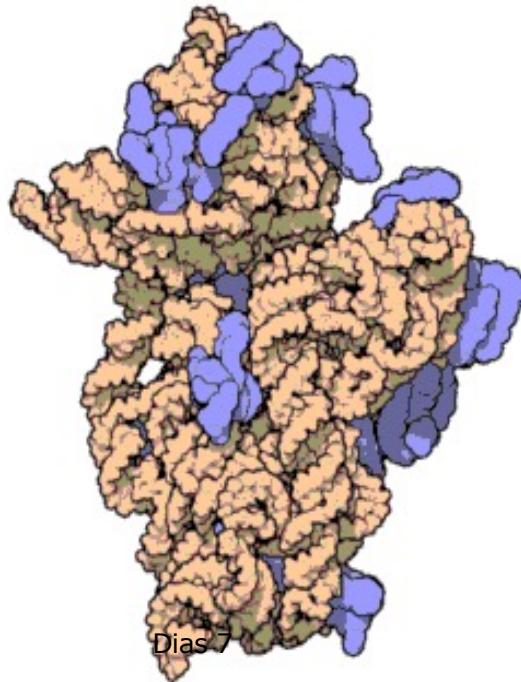
Within sample characteristics

Beta diversity

Between sample characteristics



Amplicon



Alpha diversity

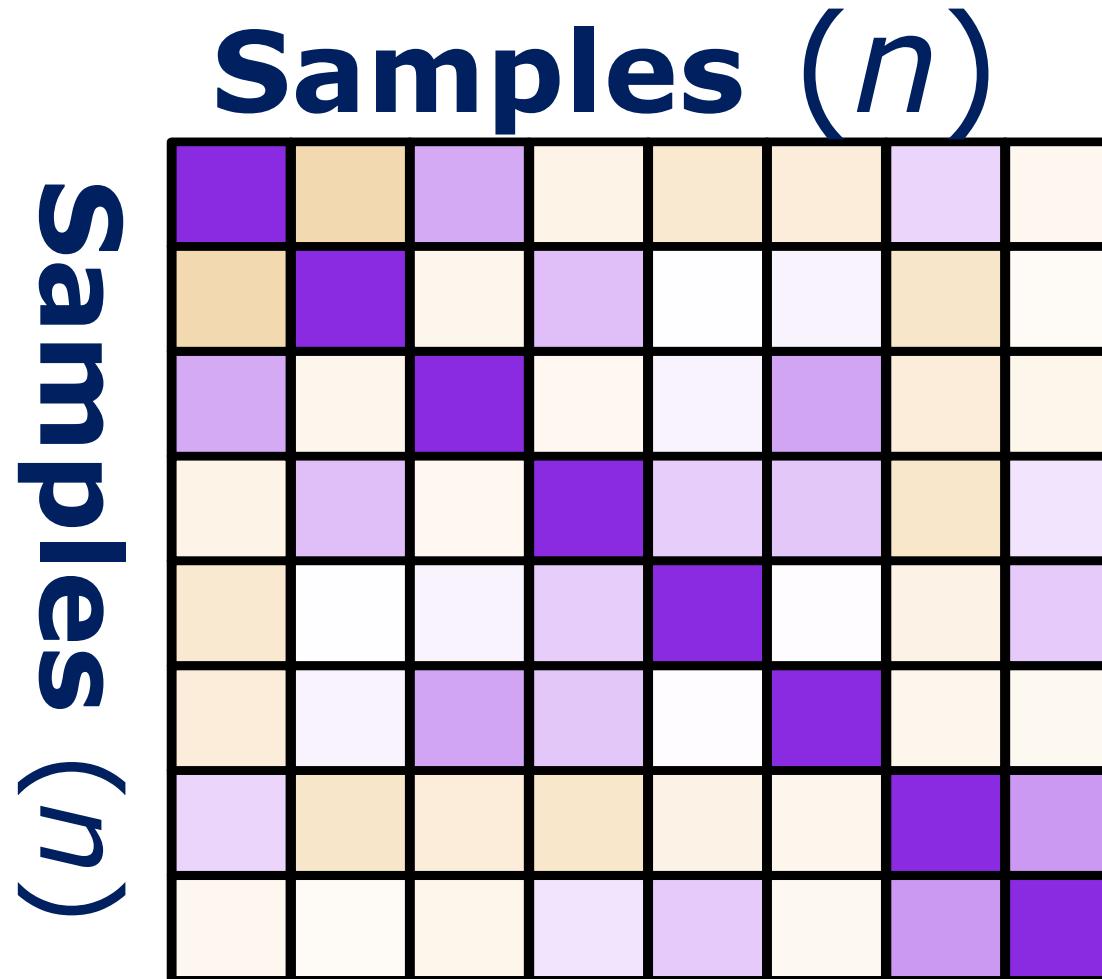
Number of different taxa

Shannon diversity $H = - \sum_{i=1}^p ra_i \cdot \ln(ra_i)$

Simpson $D = \frac{1}{\sum_{i=1}^p ra_i^2}$



β diversity



β diversity

	Presence/absense	Abundance
+Phylo	UNIFRAC PINA	wUNIFRAC wPINA
No-phylo	Jaccard Sørensen ...	Bray-Curtis Euclidian Manhattan ...



Jaccard

		Sample A	
		No. of species present	No. of species absent
Sample B	No. of species present	a	b
	No. of species absent	c	d

$$S_j = \frac{a}{a + b + c}$$



Bray Curtis

$$BC = \frac{\sum_{i=1}^n |X_{ij} - X_{ik}|}{\sum_{i=1}^n (X_{ij} + X_{ik})}$$

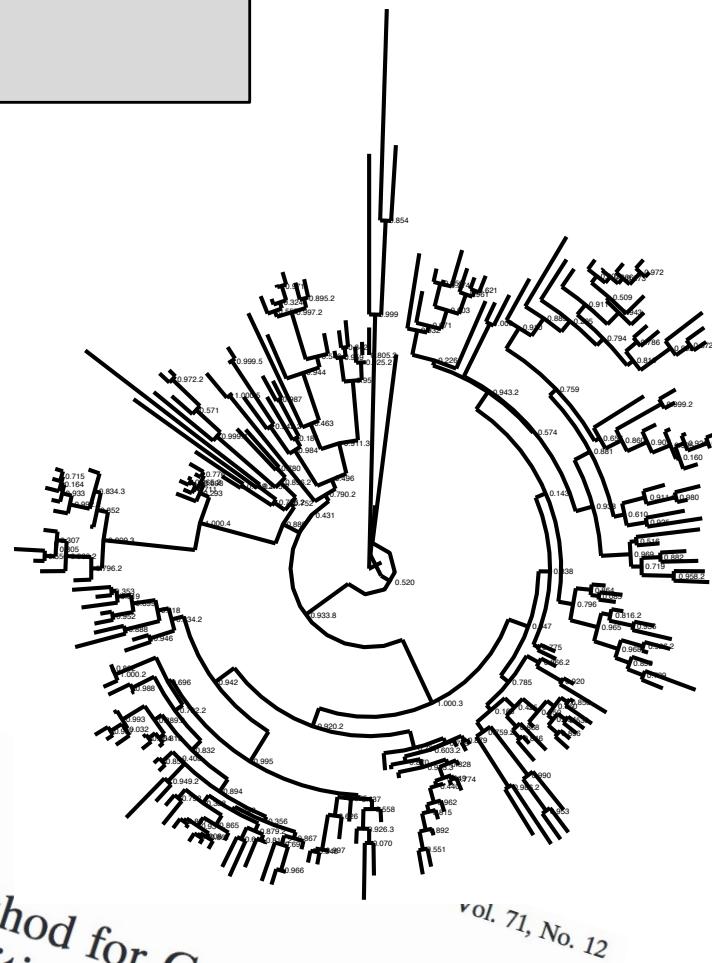
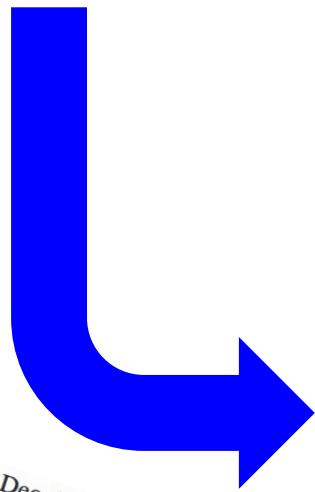
X_{ij}, X_{ik} Number of individuals in species i in each sample (j, k)
 n Total number of species in samples.



UNIFRAC

OTU table

Dist



UniFrac: a New Phylogenetic Method for Comparing
Microbial Communities

Catherine Lozupone¹ and Rob Knight^{2,*}

¹Department of Molecular, Cellular, and Developmental Biology,
Colorado 80309,¹ and Department of Chemistry,
Colorado, P

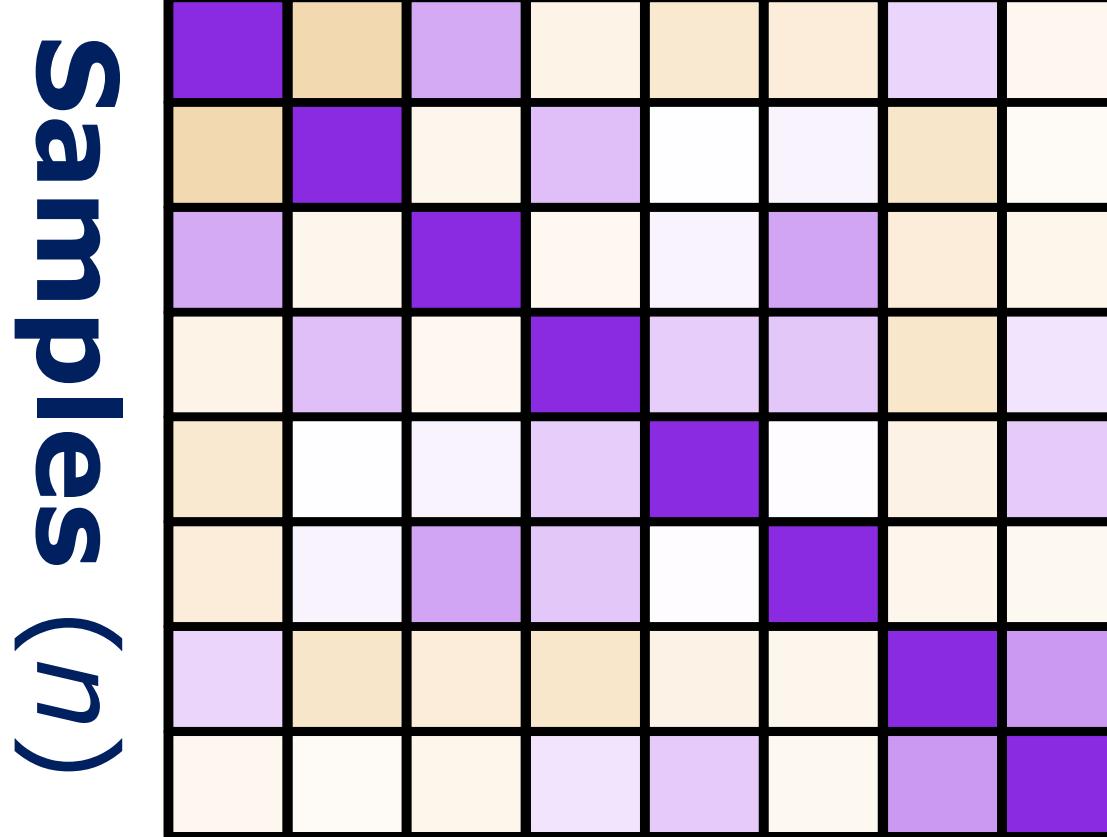


Ordination



β diversity to PCoA

Samples (n)

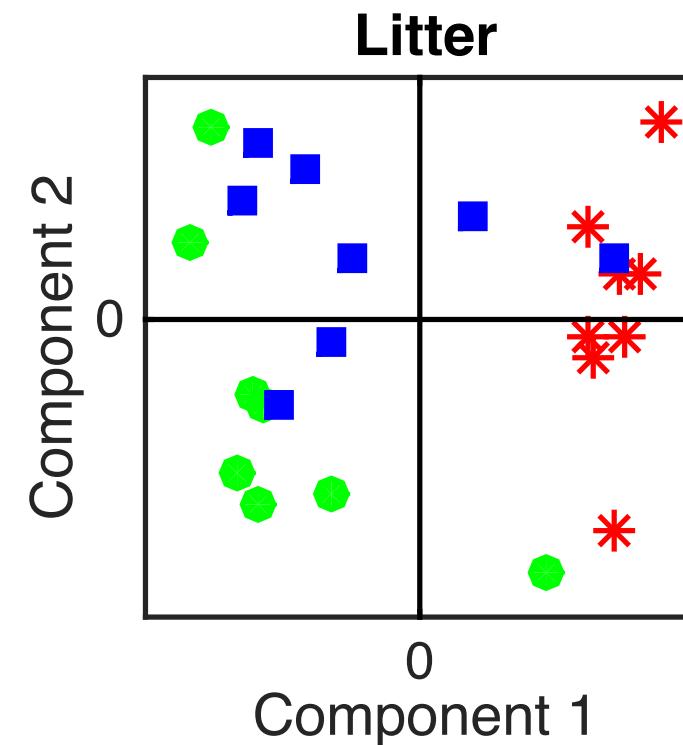
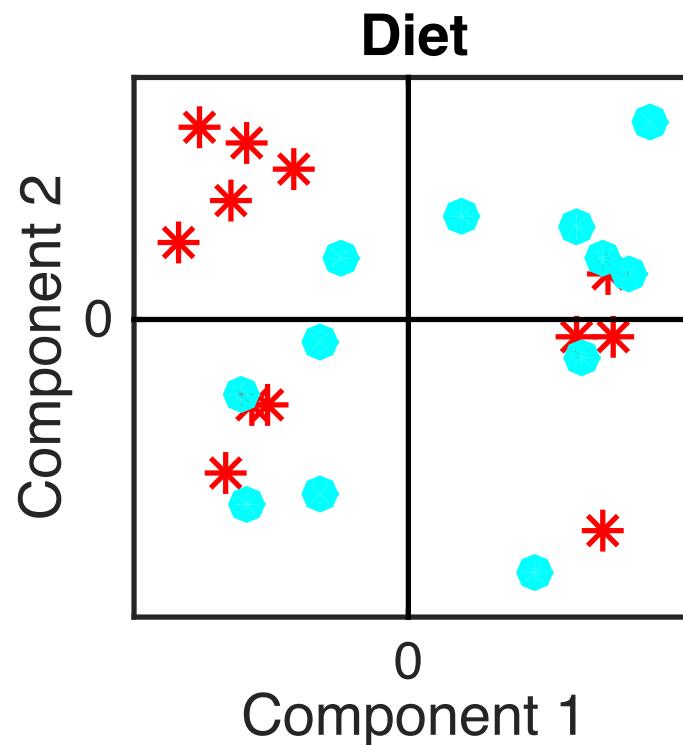


$$= U \Lambda U'$$



Multidimensional scaling

$$\mathbf{U} \Lambda \mathbf{U}^T = \mathbf{M} \exp(-\mathbf{D} \text{ist}) \mathbf{M}^T$$



		Litter		
		1	2	3
Diet	A	4	4	4
	B	4	4	4

$$\mathbf{M} = (\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T)$$



Differential Abundance Testing *or* OTUWAS



Idea

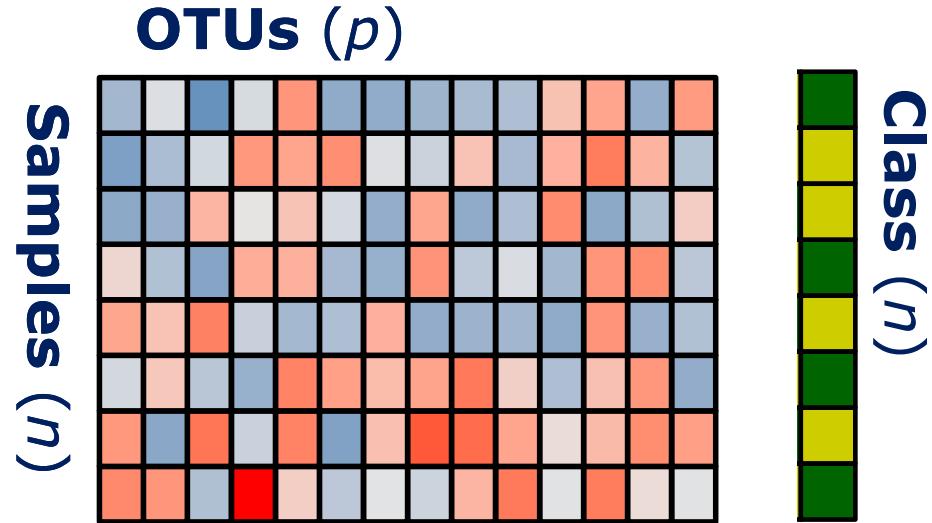
1. Perform p univariate tests recording an inferential statistics (e.g. the p-value)
2. Arrange the p (OTUs) from the most different wrt classes to the least different

$$pv_1 < pv_2 < \dots < pv_p$$

3. Figure out a threshold to separate the p OTUs into discoveries and non-discoveries.

Sted og dato

Dias 18



What to consider?

Choose a powerful statistical method

That is: avoid methods which are wrong in distributional assumptions.

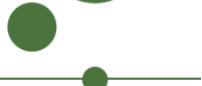
Go parametric if you can!

Utilize actively the multiple estimation to *robustify* the individual estimates.

That is: instead of using maximum likelihood for each of the p variables, shrink these towards a common value.

Problem

- Sparse / Zero inflated
- Compositional
- Over dispersed
- High dimensional



DESeq2

Developed for RNAseq

Based on log2 fold changes between (two) groups

Zeros are handled by regularized logarithm (shrinkage of low abundance towards common value)

Uses **empirical bayes** on

- Dispersion parameter to shrink towards theoretical distribution
- Fold Change (central parameter) to shrink towards zero



MetagenomeSeq

Handles the zero inflation explicitly by a mixture model of

- 1)The zeros and (fitted across OTUs)
- 2)The biological model (fitted for each OTU)

Uses **empirical bayes** on central- and dispersion parameters to shrink towards common value

(uses cumulative sum scaling for prepro)

$$f_{zig}(y_{ij}; s_j, \beta, \mu_i, \sigma_i^2) = \pi_j(s_j) \cdot I_{\{0\}}(y_{ij}) + (1 - \pi_j(s_j)) \cdot f_{count}(y_{ij}; \mu_i, \sigma_i^2)$$



Controlling Type I error

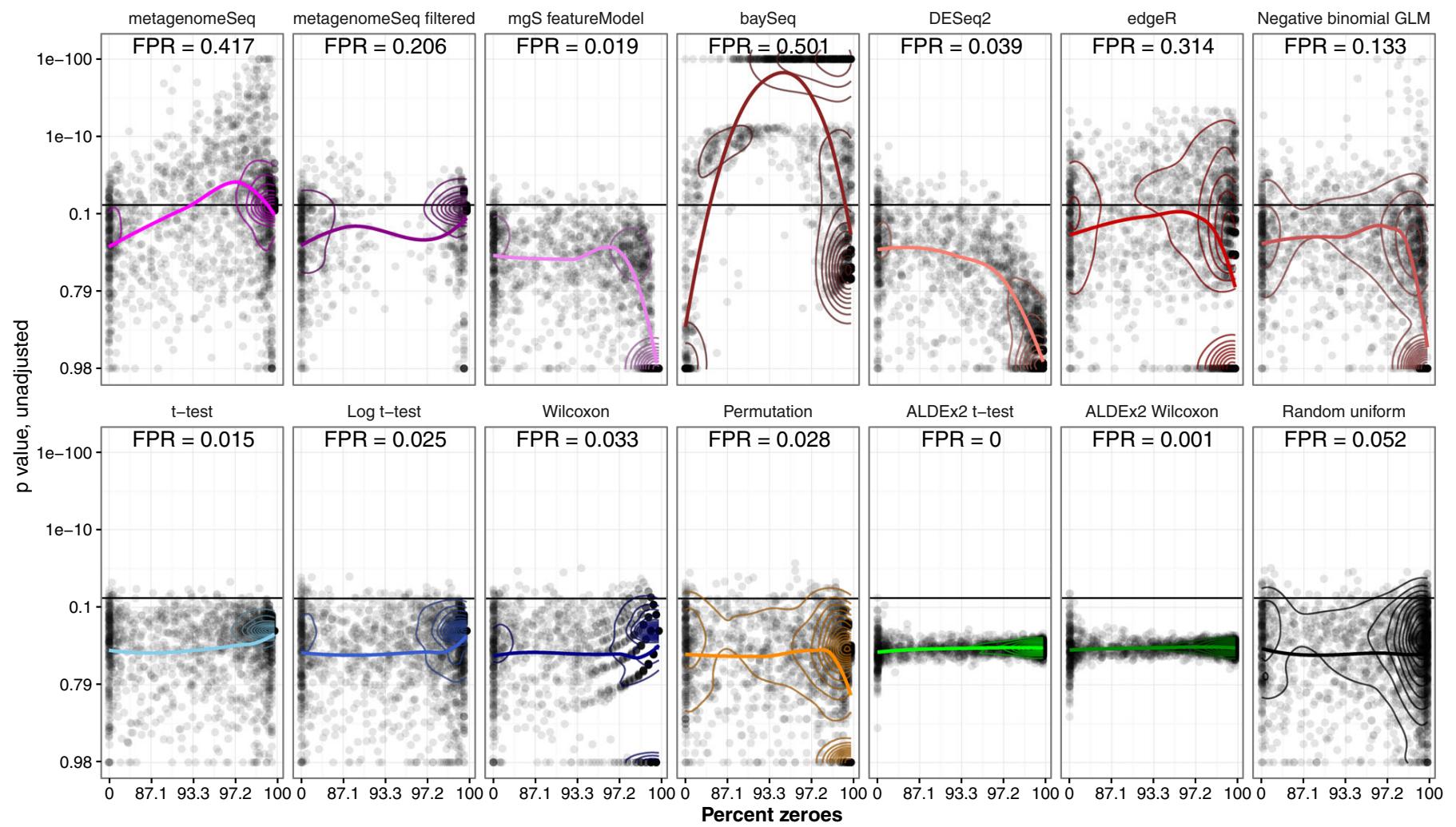


Fig. 2 OTU sparsity vs. p value. Scatterplots of OTU sparsity vs p value with panels representing each differential relative abundance test method in feces dataset A1, with 50% cases. Colored line represents the LOESS regression on data. False positive rate (FPR) is defined as the fraction of OTUs with $p < 0.05$. Each differential relative abundance test represents the median FPR for that method, out of all 150 permutations. Contour lines indicate point density and can be compared to a hypothetical null distribution of p values demonstrated in the final panel ("Random uniform")

Correct ordering

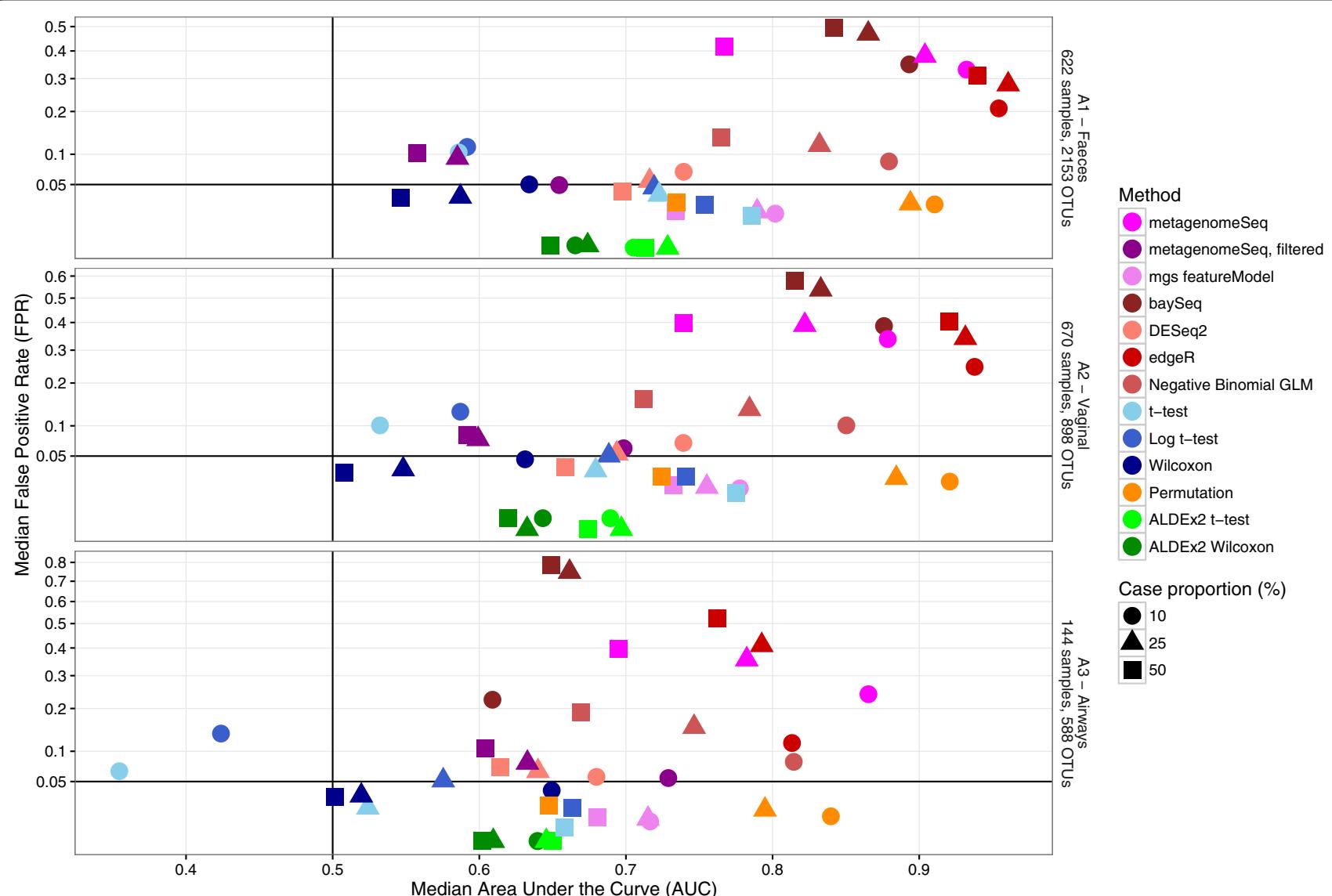


Fig. 3 Median test AUC vs median test FPR. Scatterplots of median test area under the curve (AUC) vs median test false positive rate (FPR), for the datasets A1–A3; three different compartments in the COPSAC₂₀₁₀ cohort. FPR is defined as the fraction of OTUs with $p < 0.05$. Dot color represents differential relative abundance test method, while dot shape represents experiment case/control balance

Where to cut?

Bonferoni's Family Wise Error Rate (FWER)

Any **$p_{v_i} < \alpha / p$** is a discovery

False Discovery Rate (FDR) control

Any **$p_{v_k} < k * \alpha / p$** is a discovery

Where k is the order.

These are under independence assumptions... Which is almost never fulfilled

Alternatives: *permutation testing*



Testing Beta diversity *PermANOVA*



Partitioning when we do not see the original data

The underlying model

$$Y = XB + E = \hat{Y} + E$$

Variance partitioning

$$\text{tr}(Y^T Y) = \text{tr}(\hat{Y}^T \hat{Y}) + \text{tr}(E^T E)$$

$$= \text{tr}(YY^T) = \text{tr}(\hat{Y}\hat{Y}^T) + \text{tr}(EE^T)$$

$$\hat{Y}\hat{Y}^T = H(YY^T)H$$

$$H = X(X^T X)^{-1} X^T$$

$$YY^T \propto \exp(-Dist)$$

McArdle, B. H., & Anderson, M. J. (2001). Fitting multivariate models to community data: a comment on distance-based redundancy analysis. *Ecology*, 82(1), 290-297.

Establishing the *null* distribution

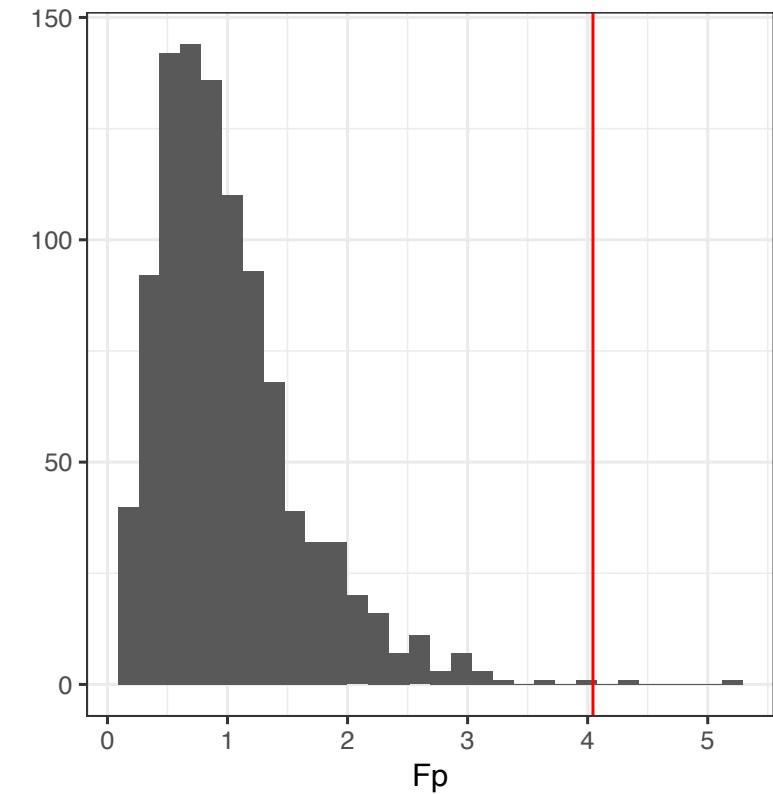
For univariate ANOVA models, we rely on independence and homoschedastic and gaussian *like* residuals.

This make testing easy, as the F-statistic follows the F-distribution under the null hypothesis

For adonis / permanova, we do not know exactly which distribution to use.

What to do?

This one is established by **permuting** the design matrix many times, each time calculating the F-statistics.

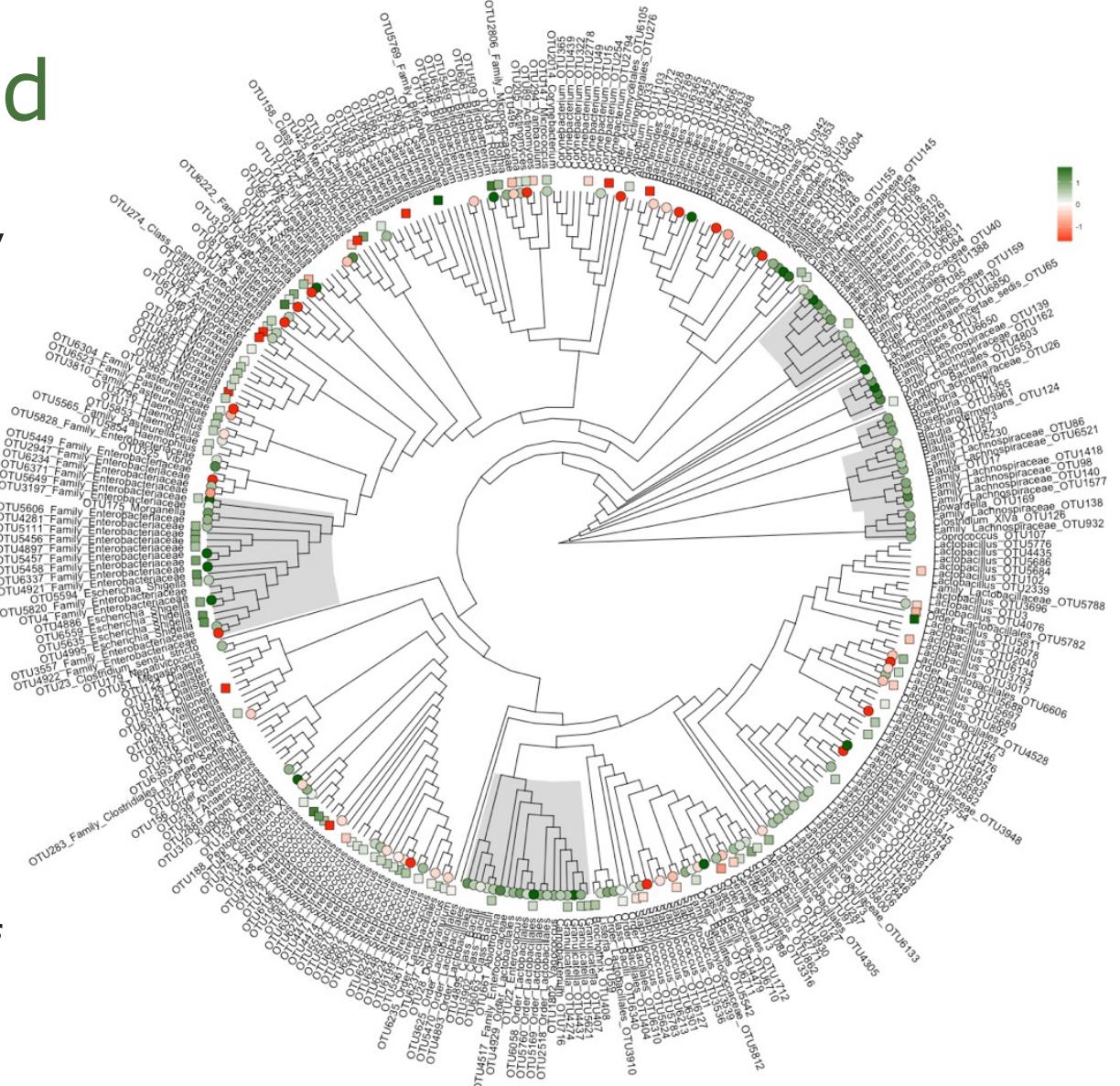


Flip it around

Just as the beta diversity
reflects sample
similarity.

The phylogenetic tree
reflects evolutionary
similarity.

We can use this in a
similar fashion with
adonis to answer
questions related to the
dependency on the
phylogenetic structure of
some relevant univariate
results



ggtree



Multomics



Data integration

-Omics

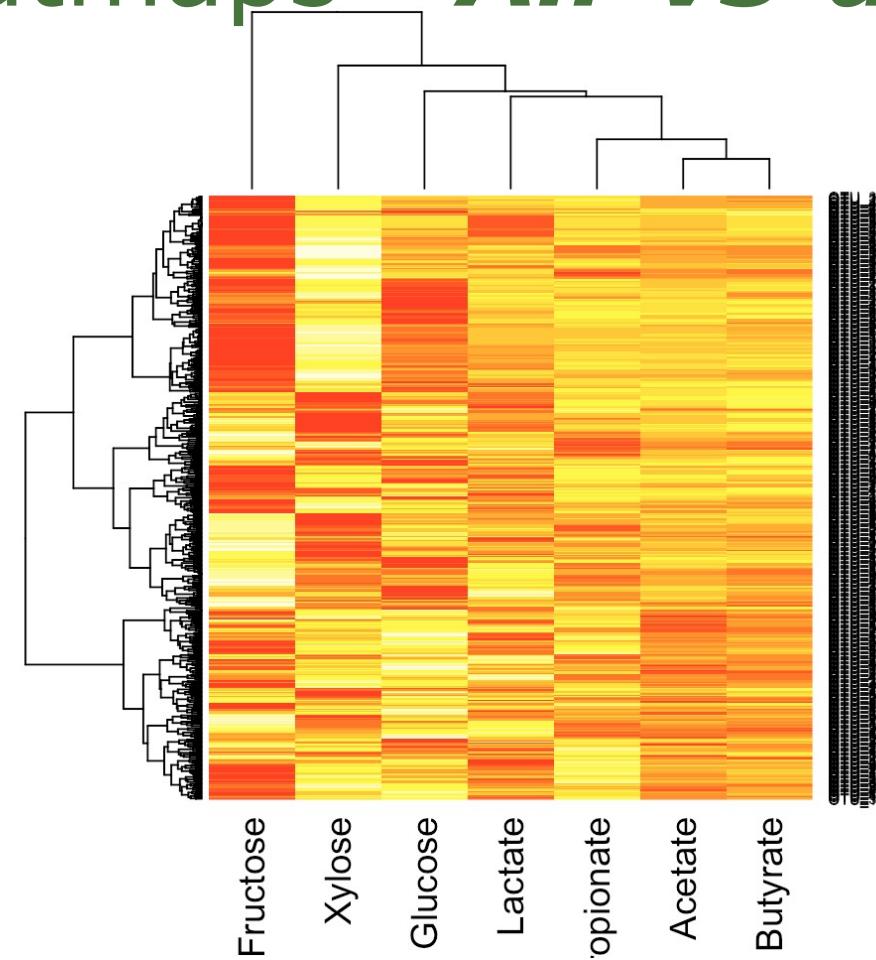
OTU table

$$n \mathbf{X}^{p_X}$$
$$n \mathbf{Y}^{p_Y}$$


Descriptive Heatmaps - *All-vs-all*

Easy and intuitive representation for multi-omics data

univariate as it is basically *univariate* correlations interpreted multivariate visually



Corr matrix of 1500 by 79000



Data integration

-Omics

OTU table

$$n \mathbf{X}^{p_X}$$
$$n \mathbf{Y}^{p_Y}$$


Functionally Supervised

Normal CCA

$$\max_{\mathbf{w}_X, \mathbf{w}_Y} \mathbf{w}_X^T \mathbf{X}^T \mathbf{Y} \mathbf{w}_Y \quad \begin{array}{l} \mathbf{w}_X^T \mathbf{X}^T \mathbf{X} \mathbf{w}_X = 1 \\ \mathbf{w}_Y^T \mathbf{Y}^T \mathbf{Y} \mathbf{w}_Y = 1 \end{array}$$

s.t.

Supervised CCA

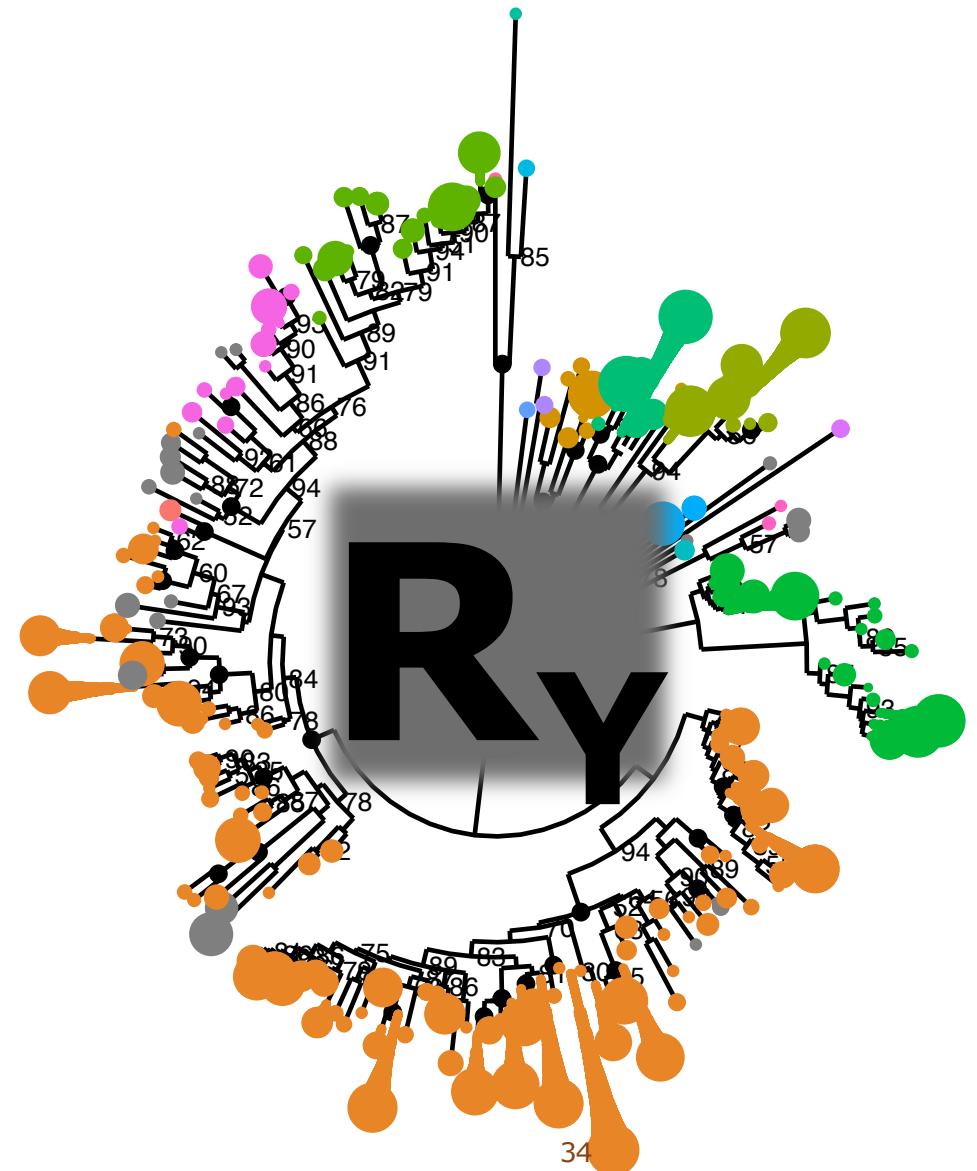
- Let \mathbf{R}_Y be a symmetric matrix reflecting similarity between the features in \mathbf{Y} based on external knowledge



Functionally Supervised

Supervised CCA

- Let \mathbf{R}_Y be a symmetric matrix reflecting similarity between the features in \mathbf{Y} based on **external knowledge**



Functionally Supervised

Supervised CCA

Re-formalize the main objective

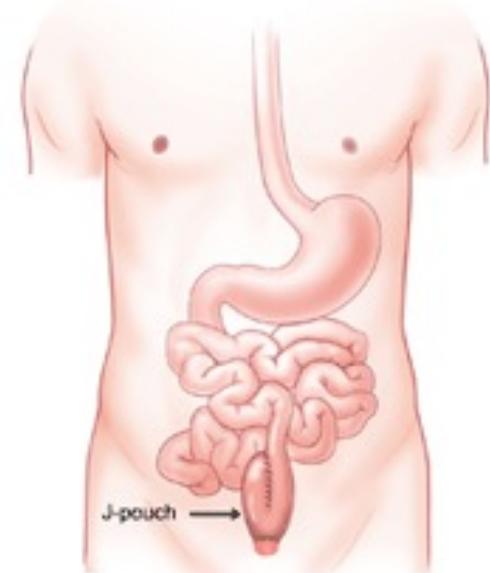
$$\max_{\mathbf{w}_X, \mathbf{w}_Y} \mathbf{w}_X^T \mathbf{X}^T \mathbf{Y} \mathbf{R}_Y \mathbf{w}_Y$$

$$\mathbf{R}_Y = \alpha \mathbf{S} + (1 - \alpha) \mathbf{I}$$

Kernel Smoothing



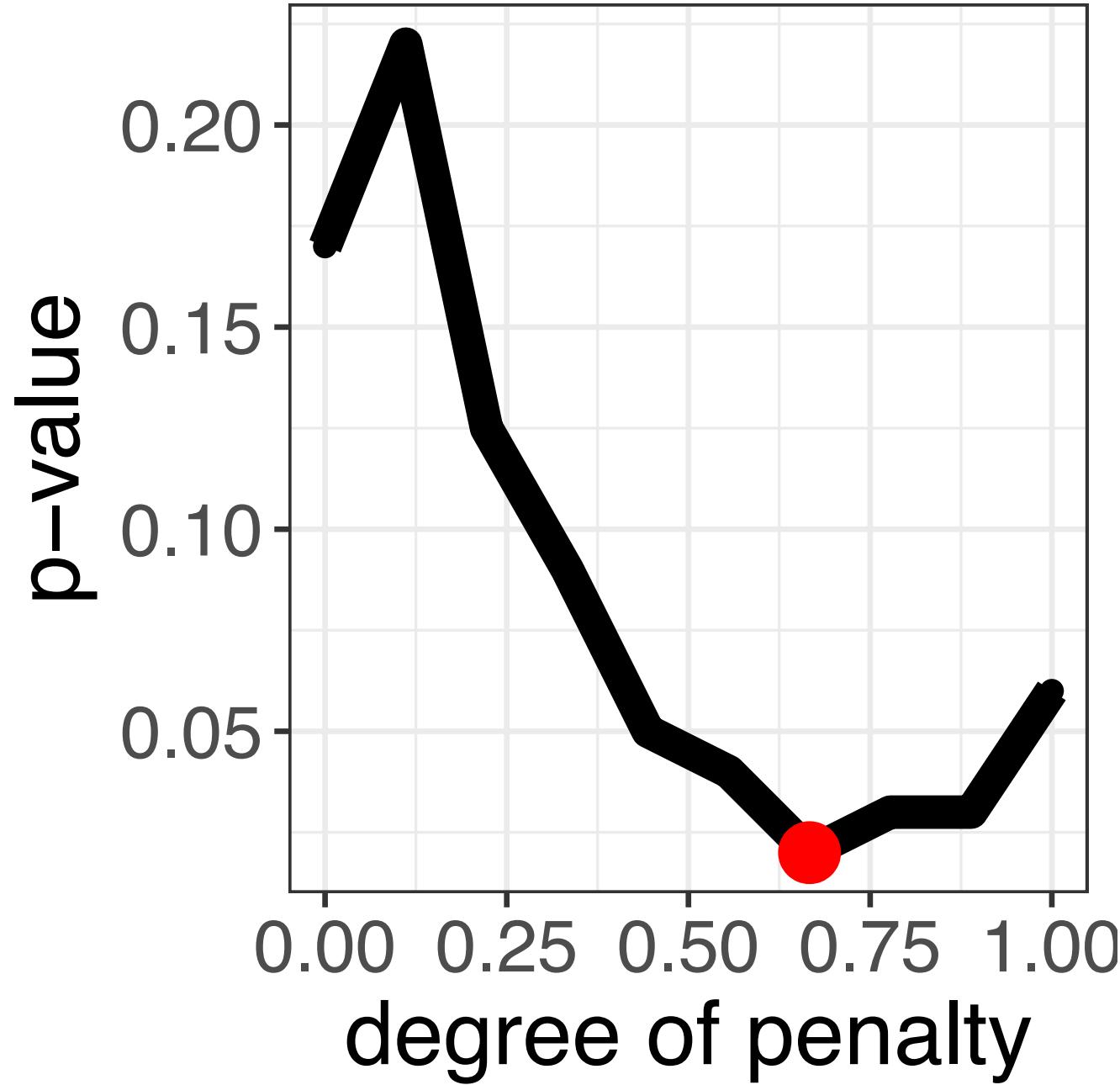
Example Pouchitic cohort

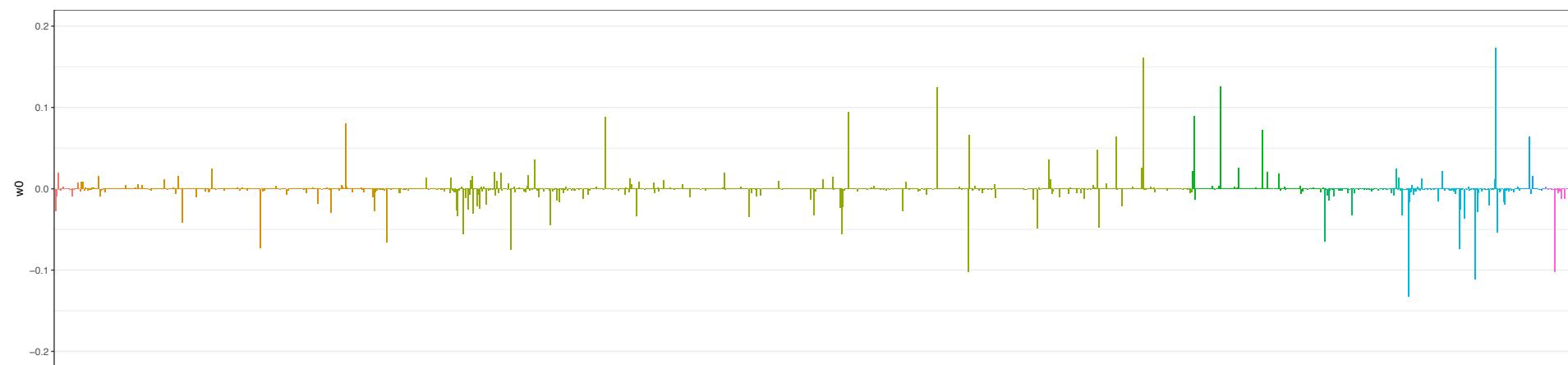
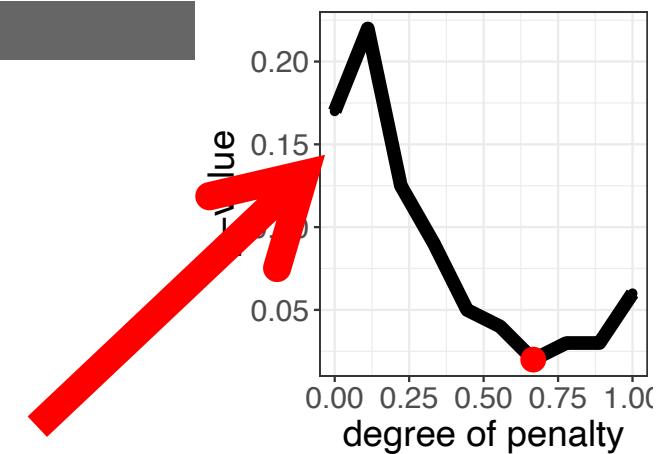


Gene
Expression

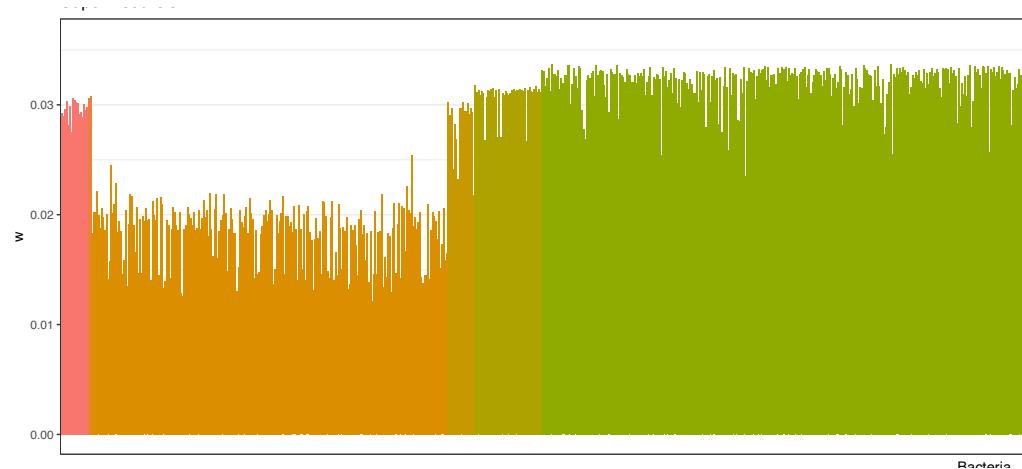
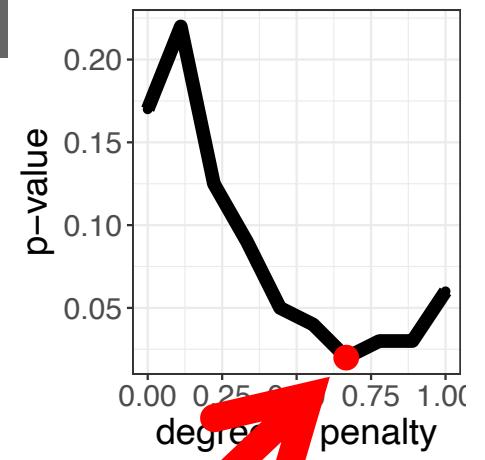
OTU table

248 X¹⁹⁹¹ 248 Y¹³⁸⁰





Order	Actinomycetales	Bifidobacteriales	Coriobacteriales	Flavobacteriales	Lactobacillales	Rhizobiales	Sphingomonadales	Xanthomonadales
	Bacillales	Burkholderiales	Enterobacteriales	Fusobacteriales	Pasteurellales	Rhodocyclales	Turicibacteriales	
	Bacteroidales	Clostridiales	Erysipelotrichales	Gemellales	Pseudomonadales	Sphingobacteriales	Verrucomicrobiales	



Truncates to
common loading
for similar bacteria

Order	Actinomycetales	Bifidobacteriales	Coriobacteriales	Flavobacteriales	Lactobacillales	Rhizobiales	Sphingomonadales	Xanthomonadales
	Bacillales	Burkholderiales	Enterobacteriales	Fusobacteriales	Pasteurellales	Rhodocyclales	Turicibacteriales	
	Bacteroidales	Clostridiales	Erysipelotrichales	Gemellales	Pseudomonadales	Sphingobacteriales	Verrucomicrobiales	

Cook-Book

- Set up data in phyloseq
- Take appropriate preprocessing choices and maybe remove samples with low seq-depth
- Perform alpha diversity analysis versus design
- Perform beta diversity analysis versus design
- Do DA and report overall results as e.g. volcano-plot and maybe reference the results against phylogenetics
 - Maybe tax agglomerate to a higher taxonomic level, and repeat analysis to see at which taxonomic level the associations are pronounced
- Do omics-omics analysis as correlation heatmaps
- Perform a supervised omics-omics using CCA or (s)PLS2 including cross-validation.

