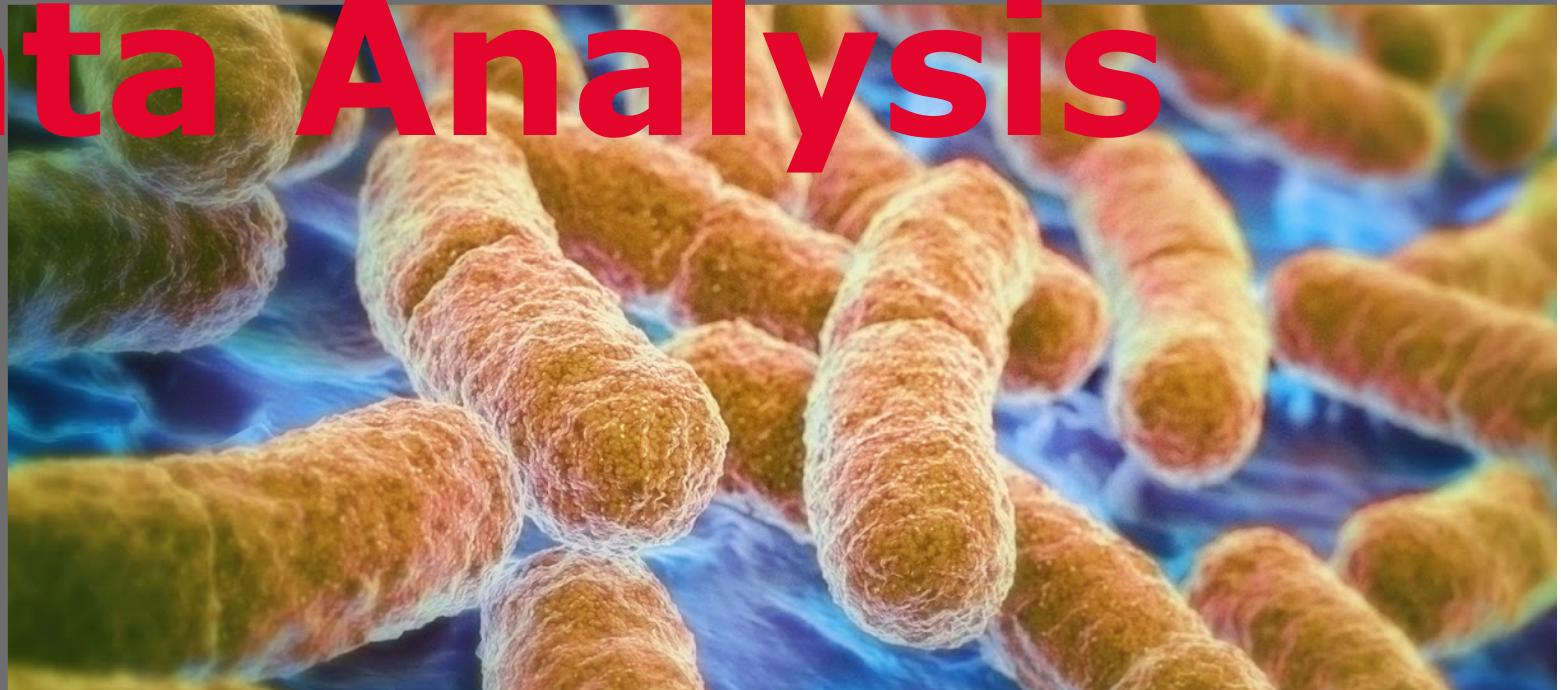


MICROBIOME

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

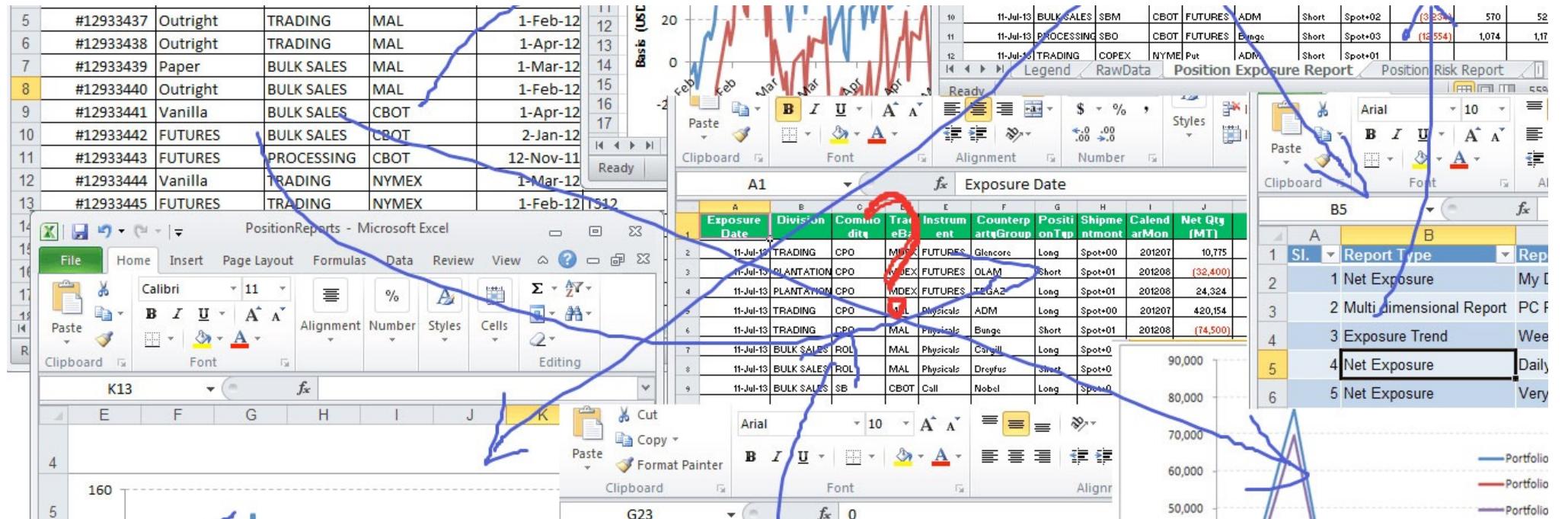
Day1 STAT101 recap Preprocessing	Day2 Morning Alpha diversity Beta diversity Testing design versus Beta diversity Permanova	Day3 Morning More DA Multiomics with heatmaps
	Day2 Afternoon Differential Abundance testing	Day3 Afternoon Multiomics - CCA Mediation and Propensity scoring



Why R??

- Digitalization *is everywhere*
- Reproducibility
- Get quicker insight and More knowledge out of your data





Excel Hell



What you want to do

R version 3.5.1 (2018-07-02) -- "Feather Spray"
 Copyright (C) 2018 The R Foundation for Statistical Computing
 Platform: x86_64-apple-darwin15.6.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
 You are welcome to redistribute it under certain conditions.
 Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
 Type 'contributors()' for more information and
 'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for online help, or
 'help.start()' for an HTML browser interface to help.
 Type 'q()' to quit R.

Do it!

Where it is executed

The output of your effort

Plots, Files, Packages,...

RStudio interface showing:

- Console tab with R startup messages.
- Terminal tab.
- Environment pane showing "Environment is empty".
- Files pane showing a directory tree under "Home".
- Plots pane.
- Packages pane.
- Help pane.
- Viewer pane.
- File menu: RStudio, File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Window, Help.
- Help menu: 100 %, Mon 10.38, Morten Arendt Rasmussen, Search.

Packages

R comes with some functionality, but not everything is covered.

Therefor we need additional functions.

There comes in the form of ***packages*** which is installed directly from within R.

> *install.packages('ggplot2')*
From
CRAN (<http://cran.r-project.org/>)
(13713 available packages)

> *BiocManager::install('metagenomeSeq')*
From
Bioconductor (<https://www.bioconductor.org/>)

> *devtools::install_github('vqv/ggbiplot')*
From
github (<https://github.com/>)



Data import

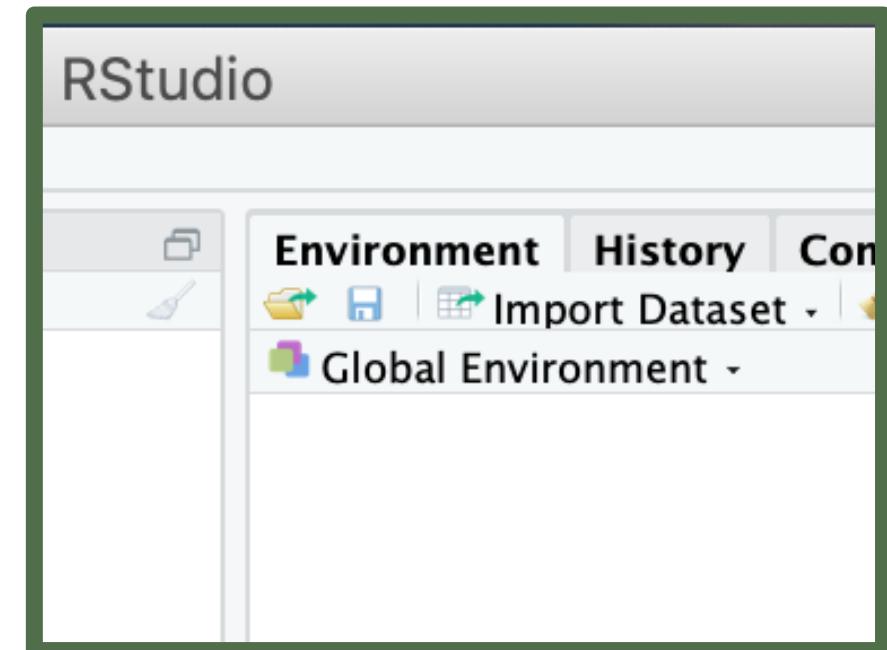
Use the *point-n-click* method in Rstudio

.... Eventually copy paste the commands produced into the script

Or

Use the **rio** package

```
> library(rio)  
> X <-  
import('myExcelFile.xlsx')
```



data.frame()

You should store your data in a data.frame

A data frame is as an excel sheet with first row being variable names.

To be aware of:

Avoid using space, leading numerics and repetitions in names.

Some useful functions to look at the data.frame

> head() / tail()	> View()	> rownames()
> str()	> colnames()	



Descriptive stats

Calculate descriptive stats – *a single number describing a distribution of numbers* – splitted according to grouping information.

Functions

- > mean(), sd(), min(), max(), range(), length()
- > aggregate()

Overlook of a design

How many observations are there for a combination of (design) variables

- > table()



STAT 101 recap

A Linear model

$$y = a(\text{treatment}) + b(\text{time}) + k(\text{batch}) + e$$

The workhorse of understanding effects of interventions, handling covariates and confounders, and getting uncertainties on all this

Test statistics

A metric/scalar measure of observed differences, which we can construct a test for



Plot *everything* with ggplot2



Average age

26.45 40 44 48 52 56 73.6

Dias 15



PLOTS

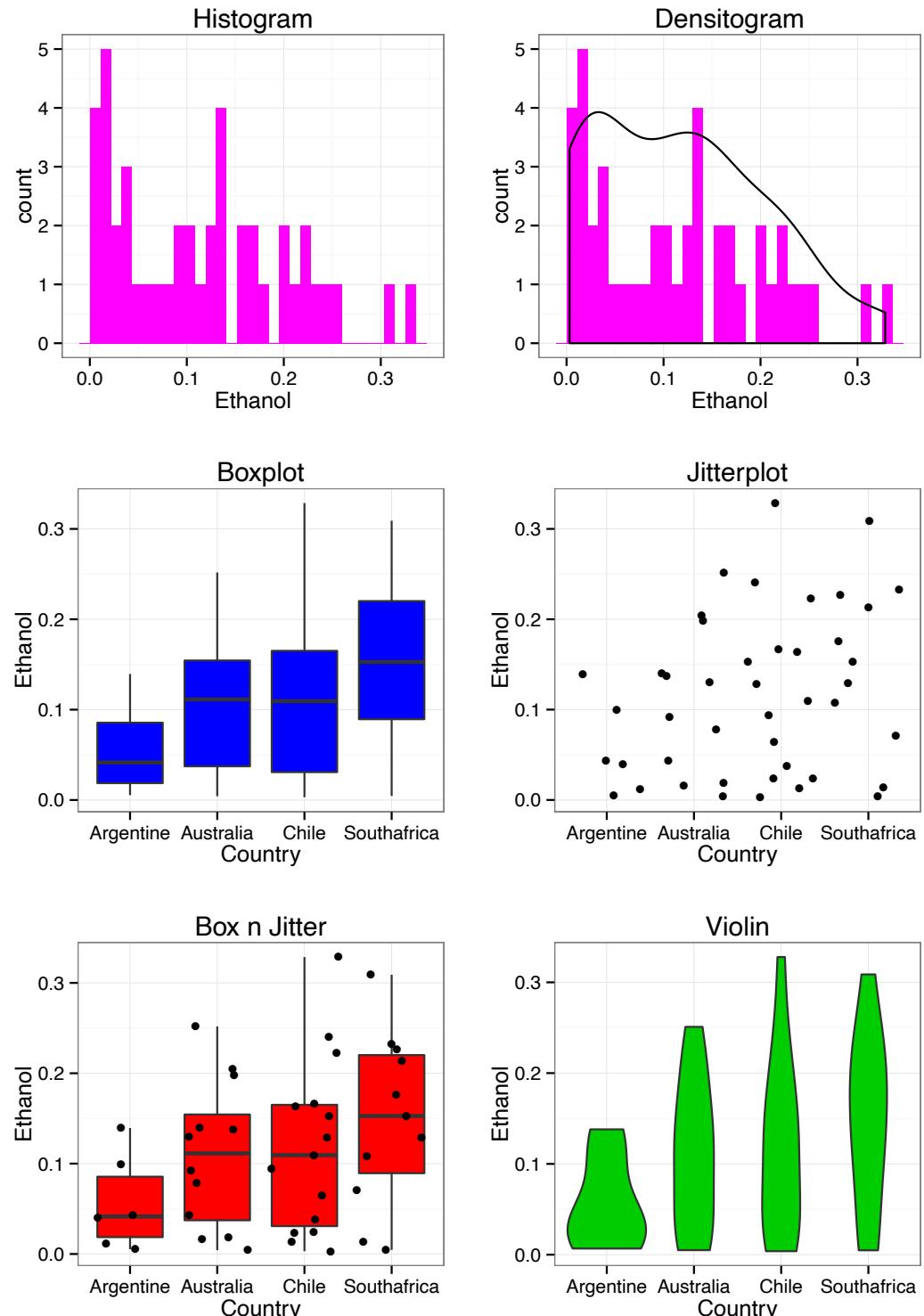
- Histogram, Densitogram, Boxplot, Jitterplot, Violinplot
 - Continuous data
 - Describes distribution
- Lineplot, Stem plot
 - Continuous data
 - Ordinal in e.g. time
- Bar chart, pie chart
 - Categorical data or compositional data
- Scatterplot
 - Bi variate continuous data
- Spline plot
 - Smooth relation between two (or more) variables



Histogram, Densitogram, Boxplot, Jitterplot, Violin plot,

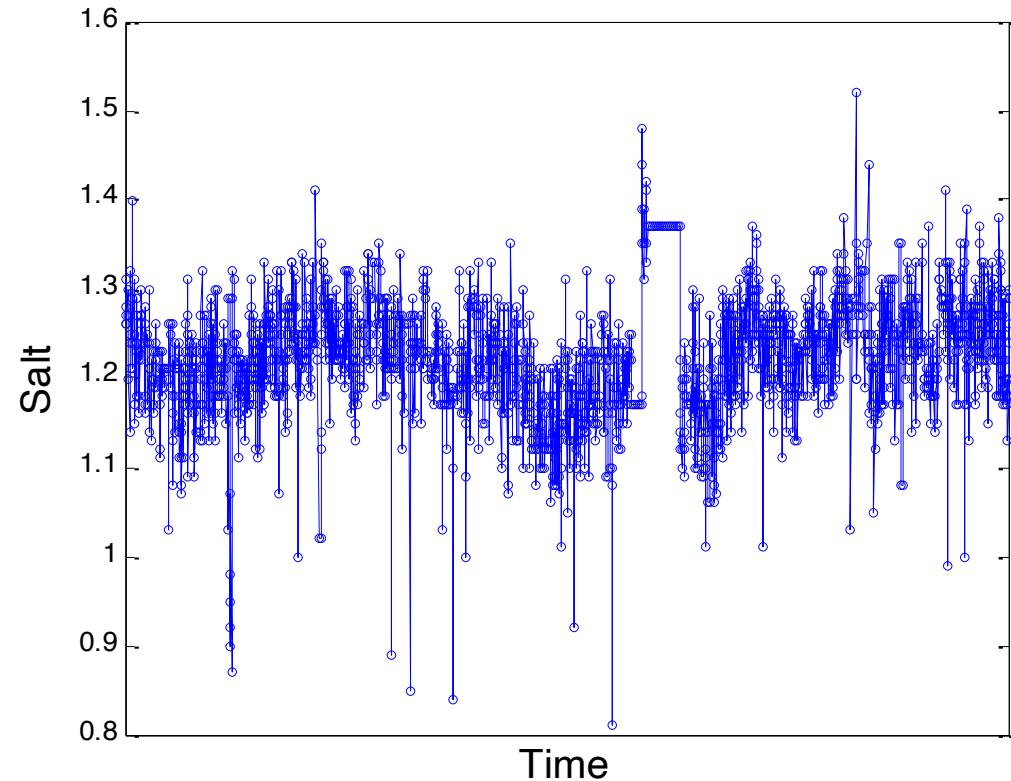
- Jitter shows raw data
- Boxplot and violin can handle many observations
Densitogram and Violin may cheat in the representation of raw data

Continuous data
Describes distribution



Lineplot

- Is used in e.g. Statistical Process Control (SPC)
- Excellent for capturing drift (systematical change over time) in the laboratory or in the production.



Continuous data
Ordinal in e.g. time

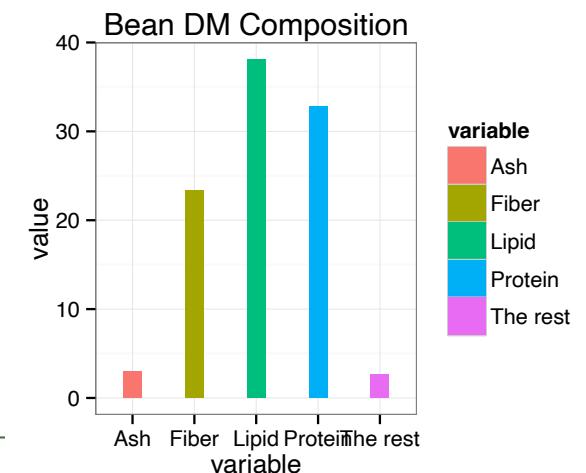
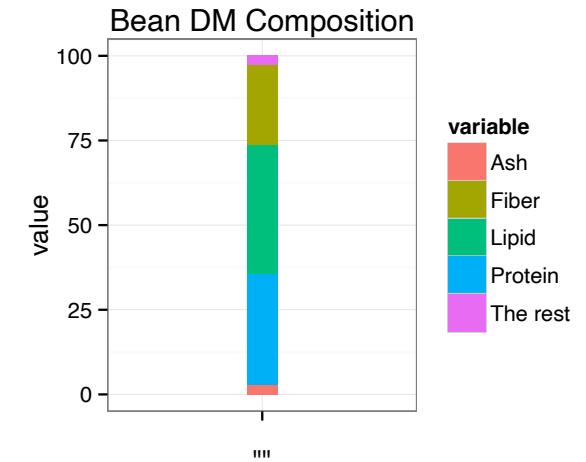
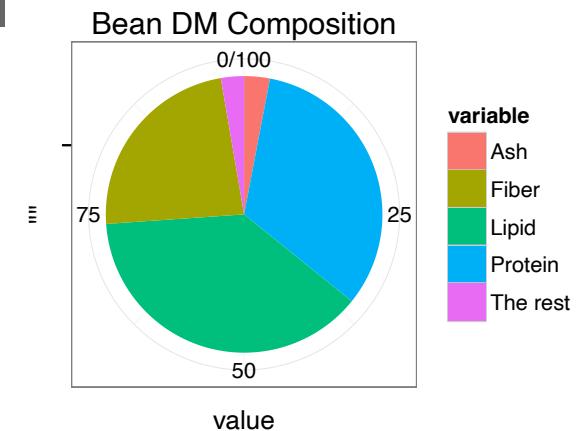


Bar charts

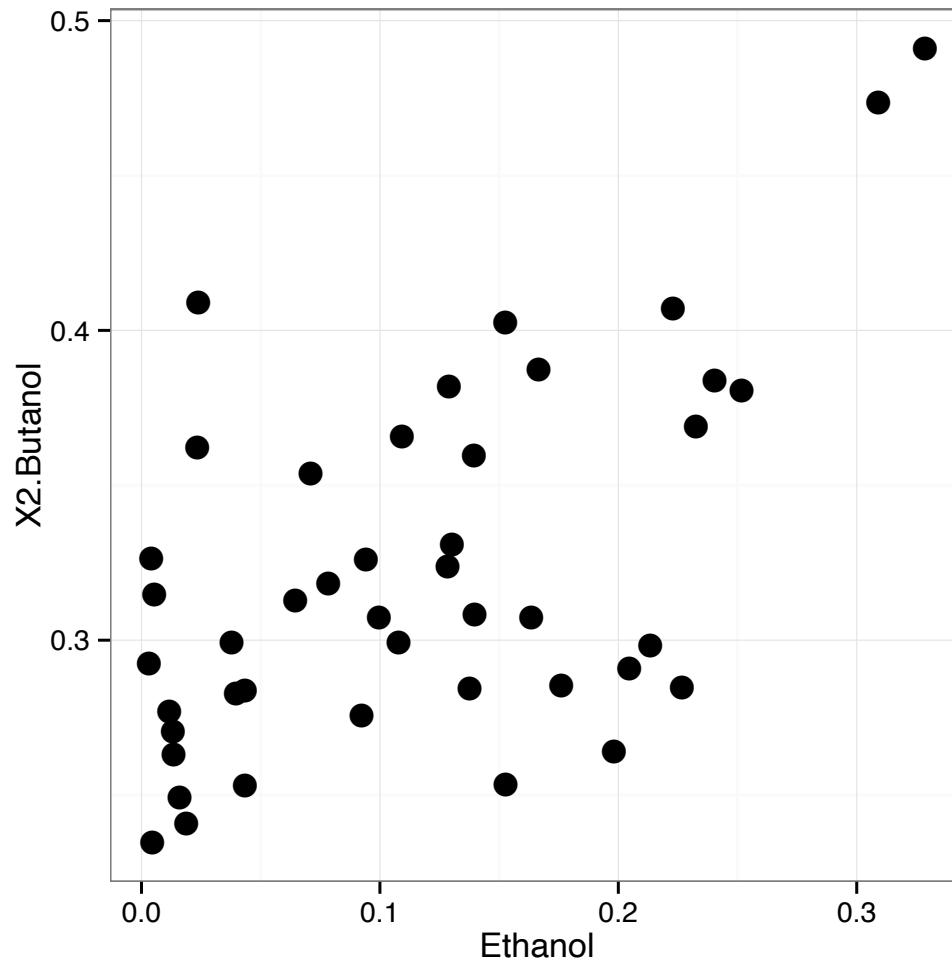
Pie charts

- OBS: Do not replace a boxplot with a bar chart!

Compositional data
Categorical data



Scatterplot

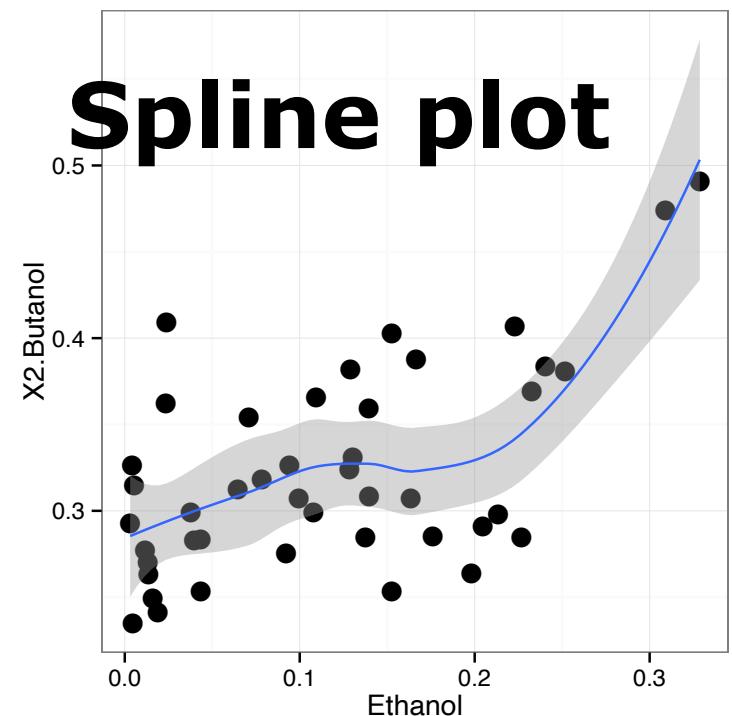
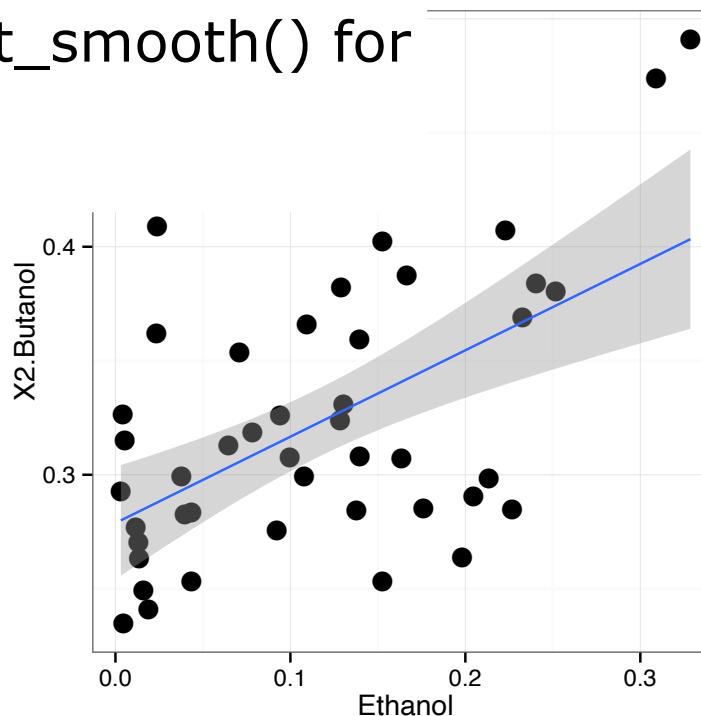
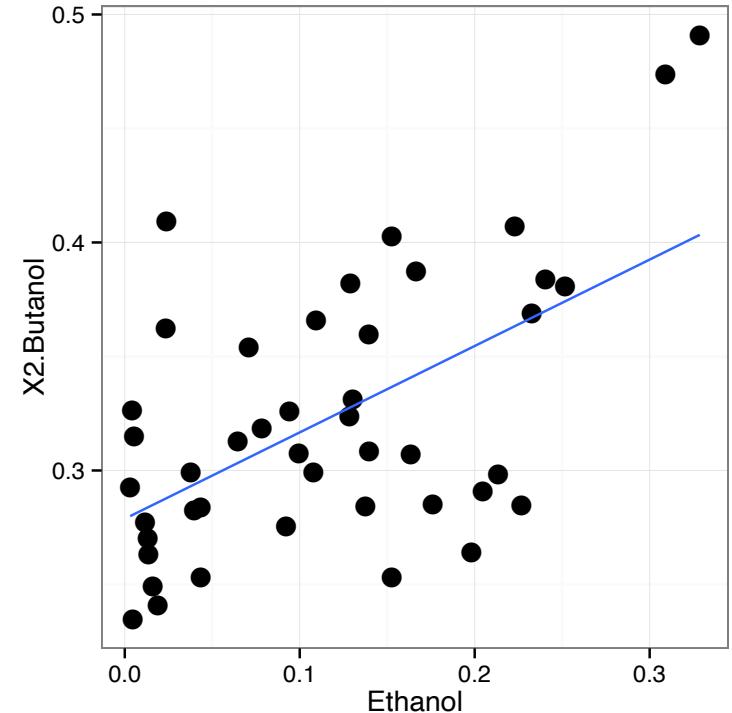
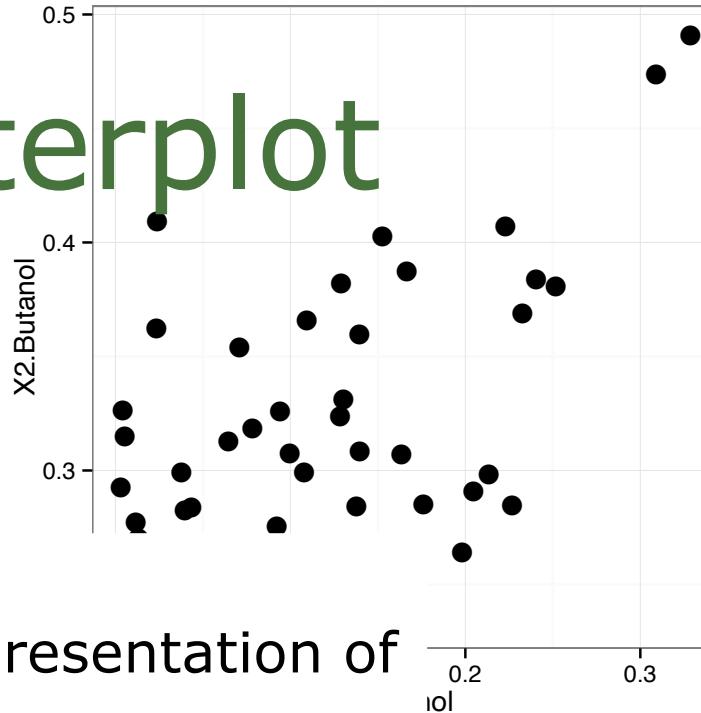


- Two response variables versus each other
- Very usefull for Multivariate data



Scatterplot

- Same data!
- Different representation of the relation
- Check: `+stat_smooth()` for details



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

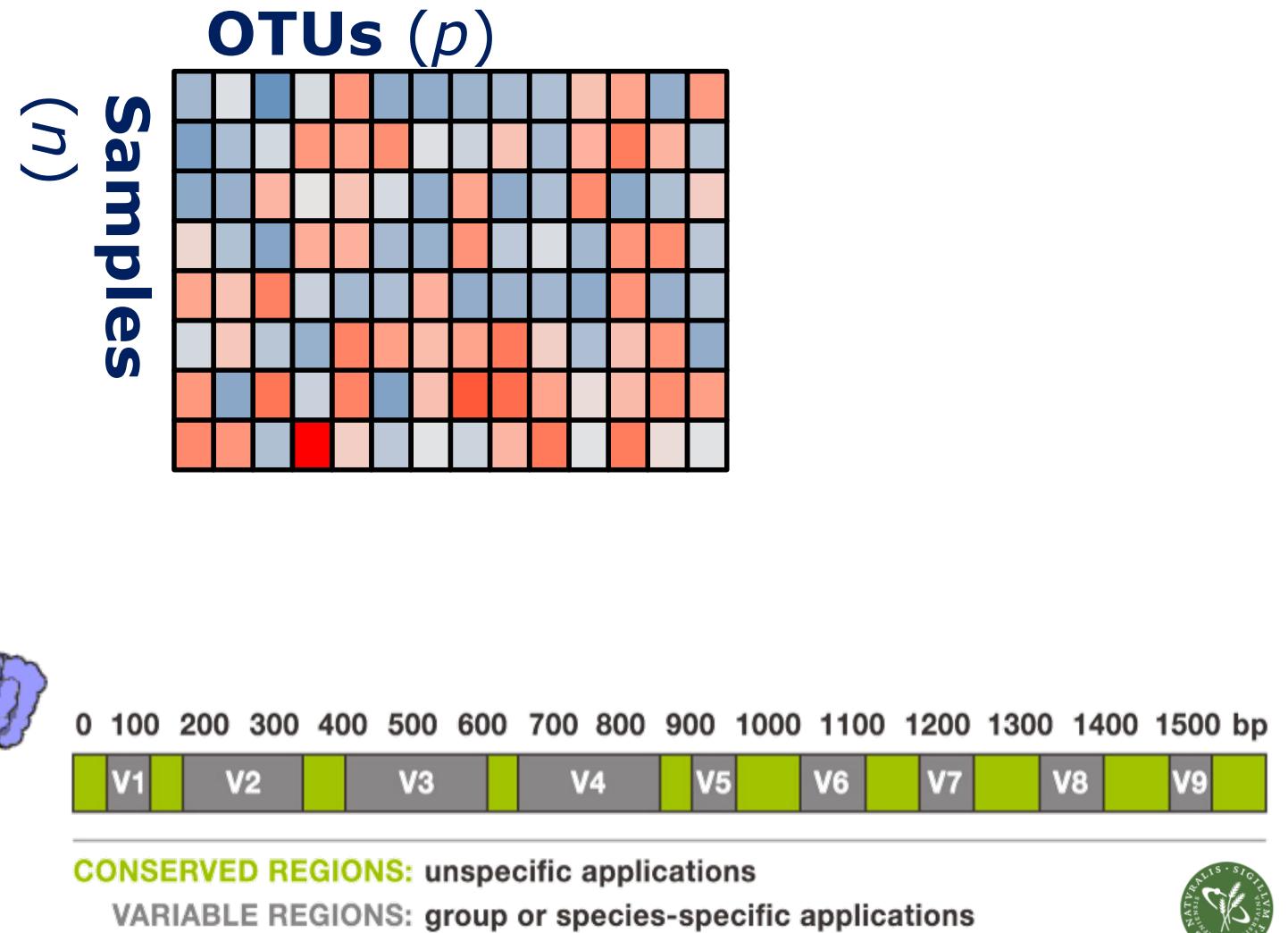
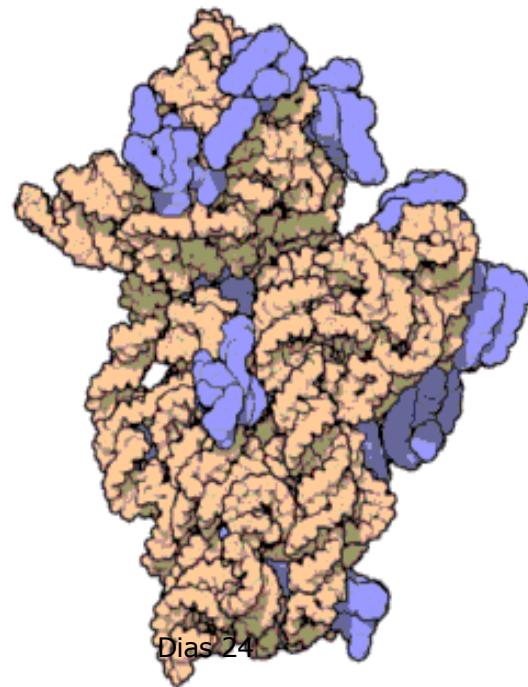
Within sample characteristics

Beta diversity

Between sample characteristics



Amplicon



Alpha diversity

Number of different taxa

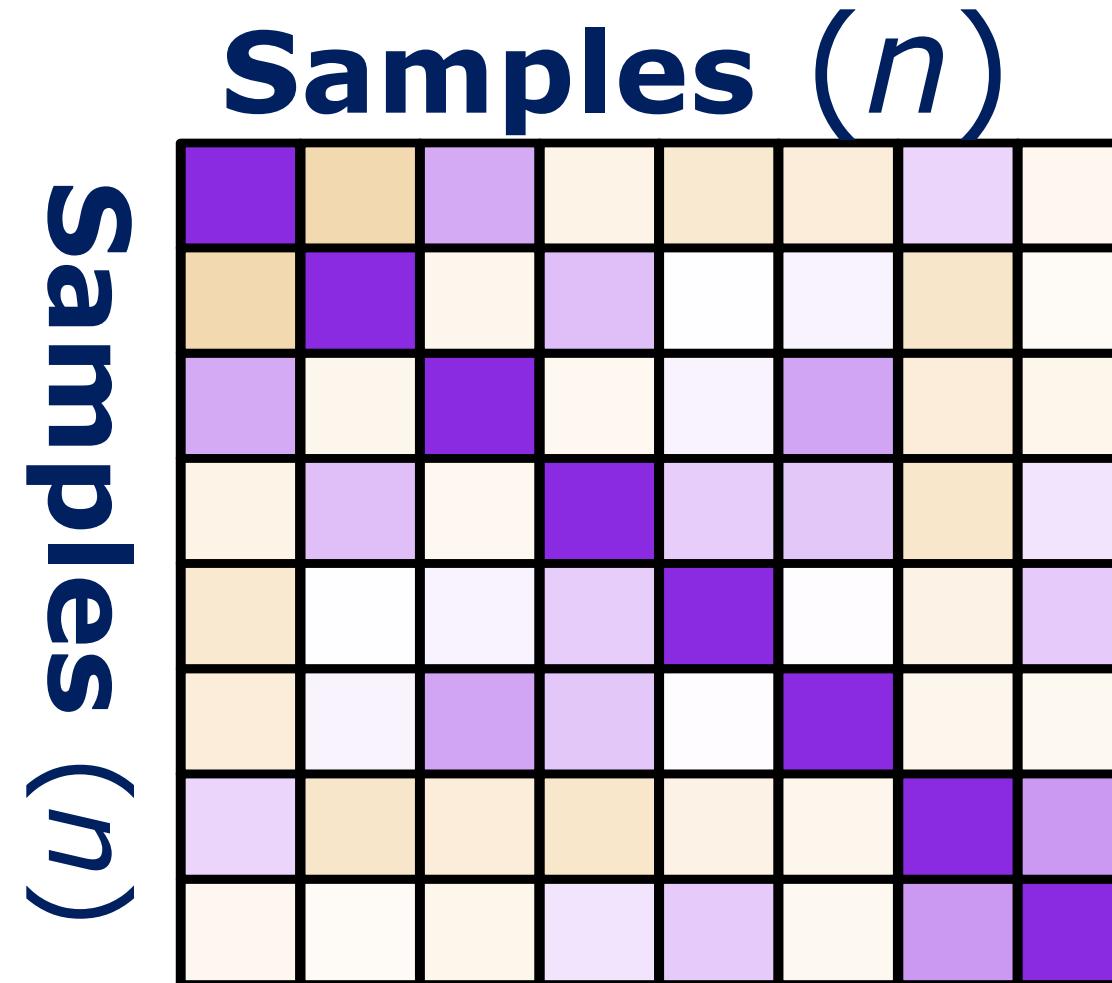
Shannon diversity $H = - \sum_{i=1}^p ra_i \cdot \ln(ra_i)$

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



kinda' ~ $\mathbf{X}\mathbf{X}^T$

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

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

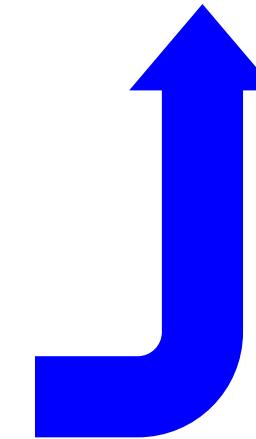
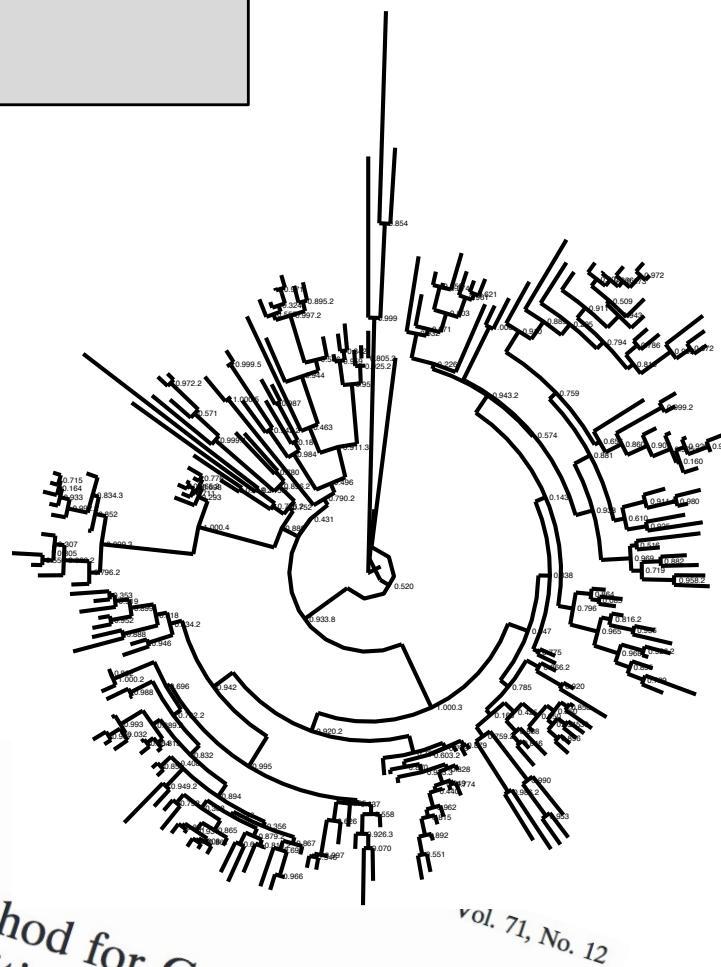
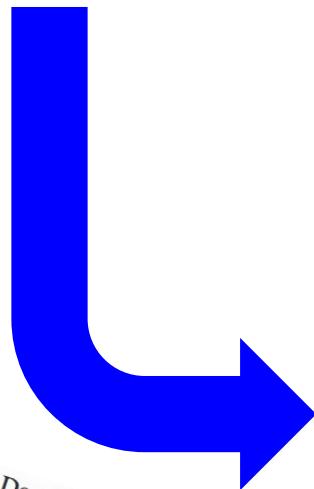
X_{ij}, X_{ik} Number of individuals in species i in each sample (j, k)
 n Total number of species in samples.



UNIFRAC

OTU table

Dist

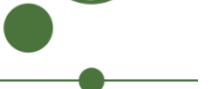


AND ENVIRONMENTAL MICROBIOLOGY, Dec. 2005, p. 8228–8235
/05/\$08.00+0 doi:10.1128/AEM.71.12.8228–8235.2005
© 2005, American Society for Microbiology. All Rights Reserved.

Unifrac: a New Phylogenetic Method for Comparing
Microbial Communities
Catherine Lozupone¹ and Rob Knight^{2,*}
¹Catherine Lozupone¹ and Rob Knight^{2,*}
Department of Molecular, Cellular, and Developmental Biology,
Colorado 80309,¹ and Department of Chemistry,
Colorado, P²

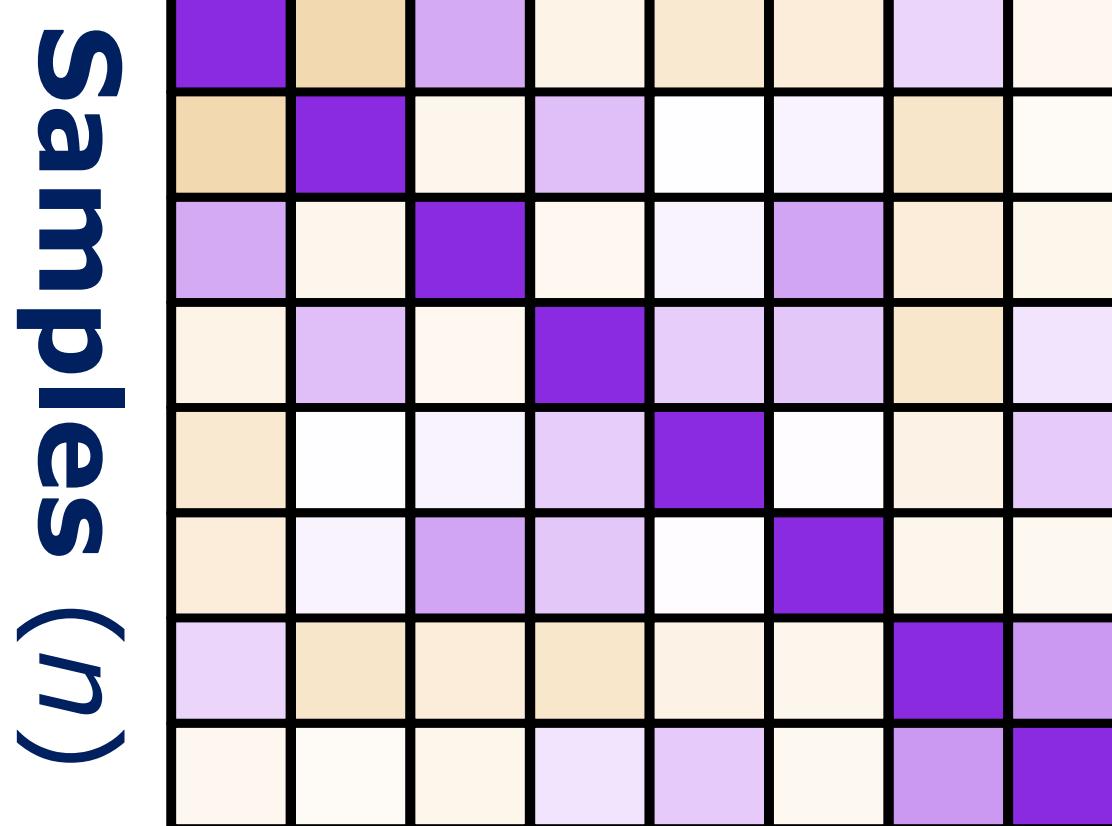


Ordination



β diversity to PCoA

Samples (n)

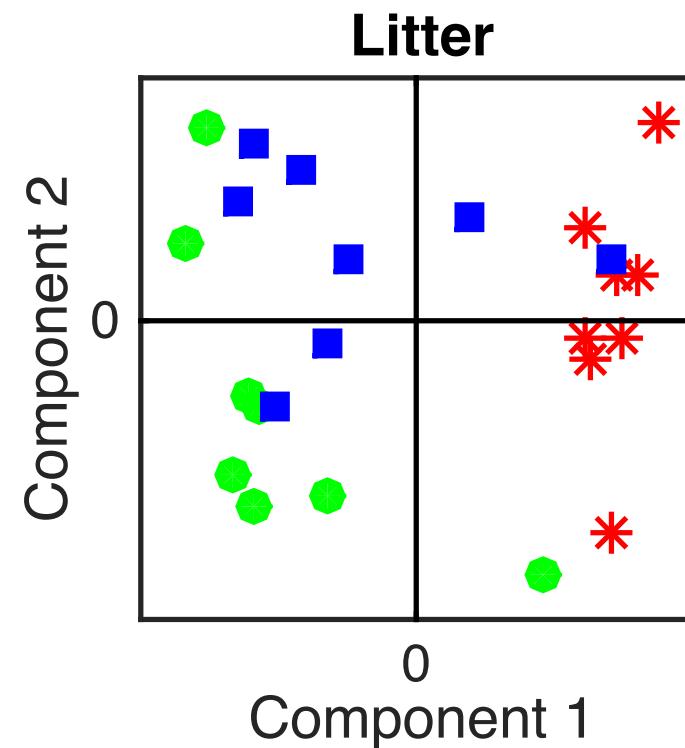
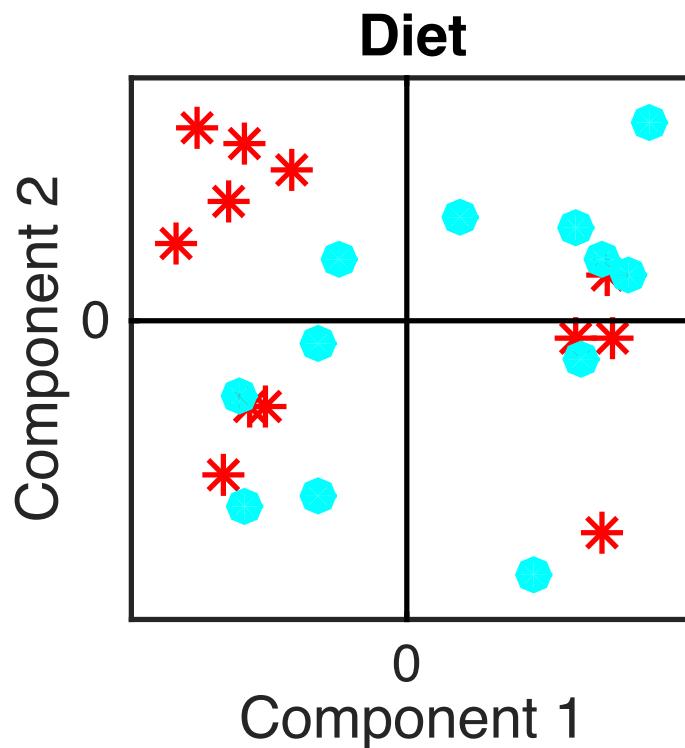


$$= \mathbf{U} \Lambda \mathbf{U}'$$



Multidimensional scaling

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



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

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



Testing Beta diversity *PermANOVA*



Partitioning when we do not see the original data

The underlying model

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

Variance partitioning

$$\begin{aligned} \text{tr}(Y^T Y) &= \text{tr}(\hat{Y}^T \hat{Y}) + \text{tr}(E^T E) \\ &= \text{tr}(YY^T) = \text{tr}(\hat{Y}\hat{Y}^T) + \text{tr}(EE^T) \end{aligned}$$

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

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

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

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

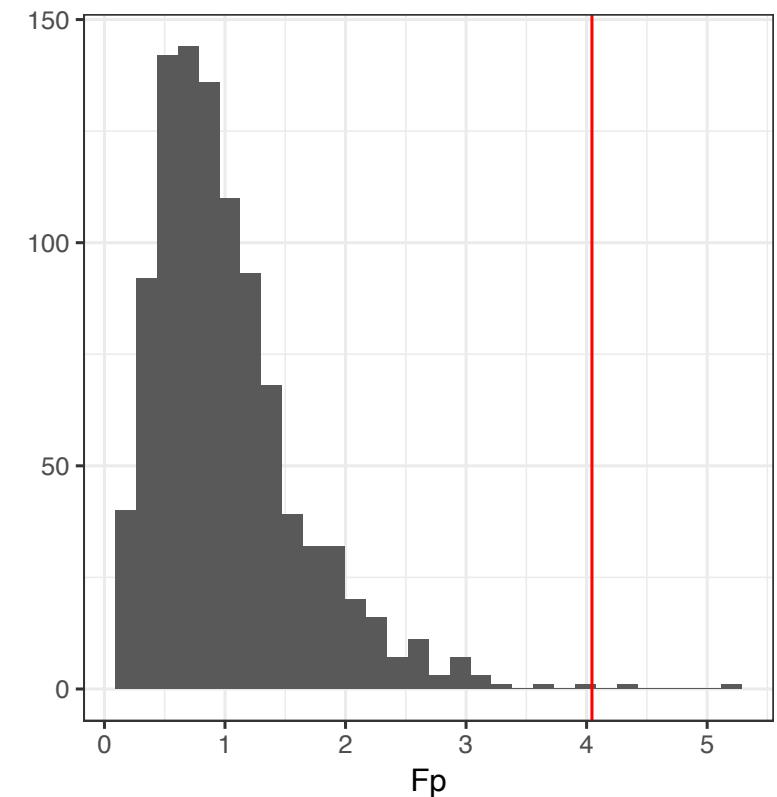
For univariate ANOVA models, we rely on independence and homoschedastic and gaussian *like* residuals.

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.

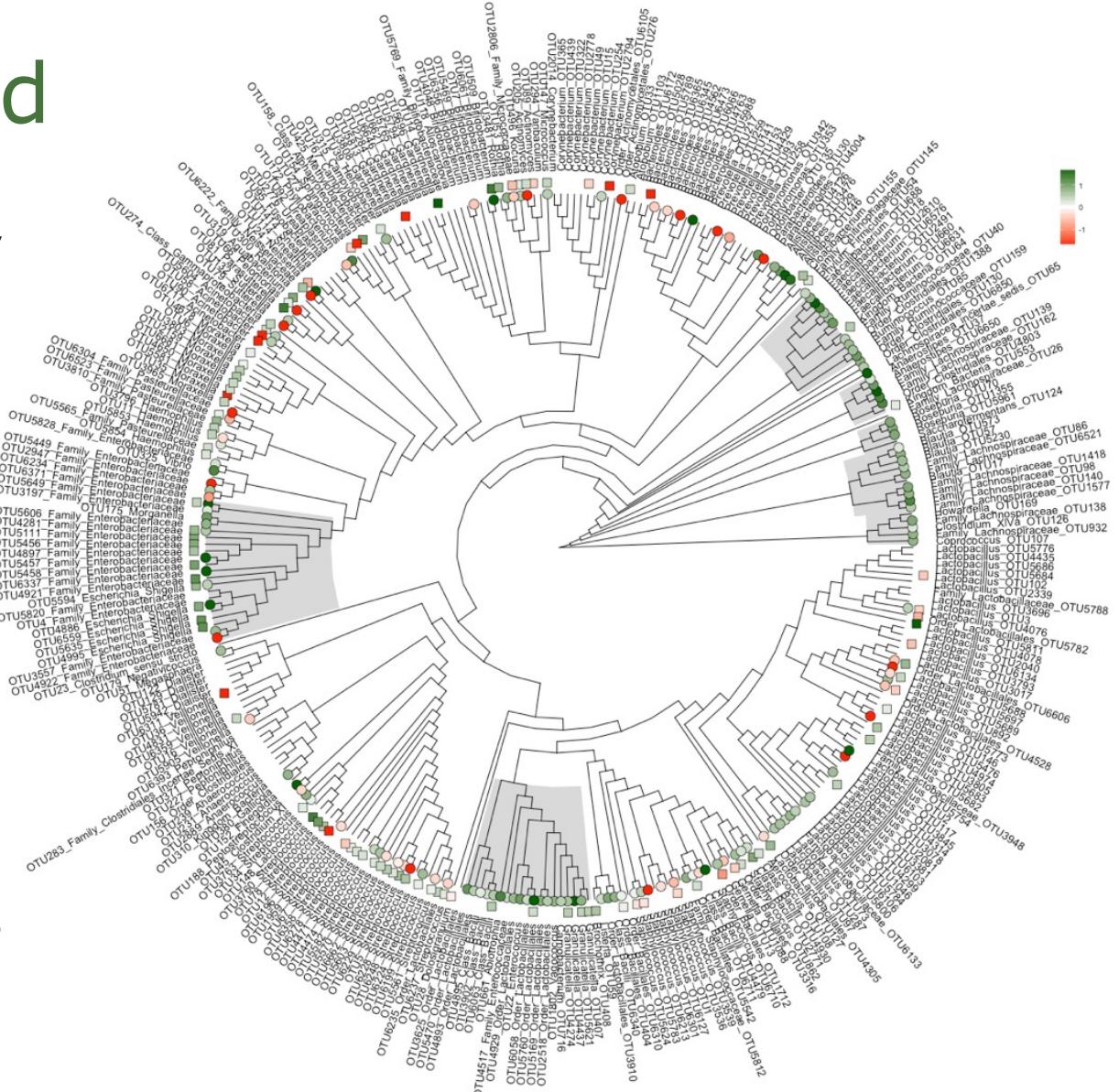


Flip it around

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



Differential Abundance Testing *or* OTUWAS



Idea

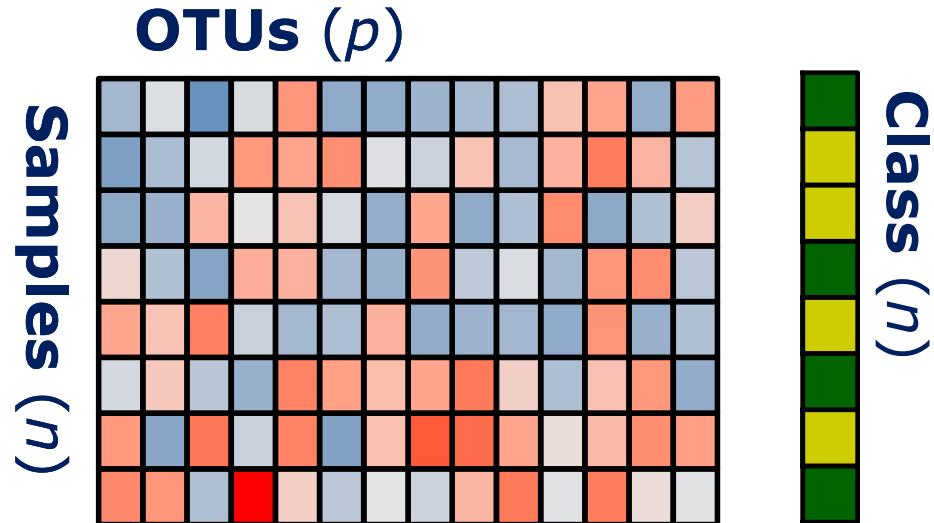
1. Perform p univariate tests recording an inferential statistics (e.g. the p-value)
2. Arrange the p (OTUs) from the most different wrt classes to the least different

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

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

Sted og dato

Dias 39



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

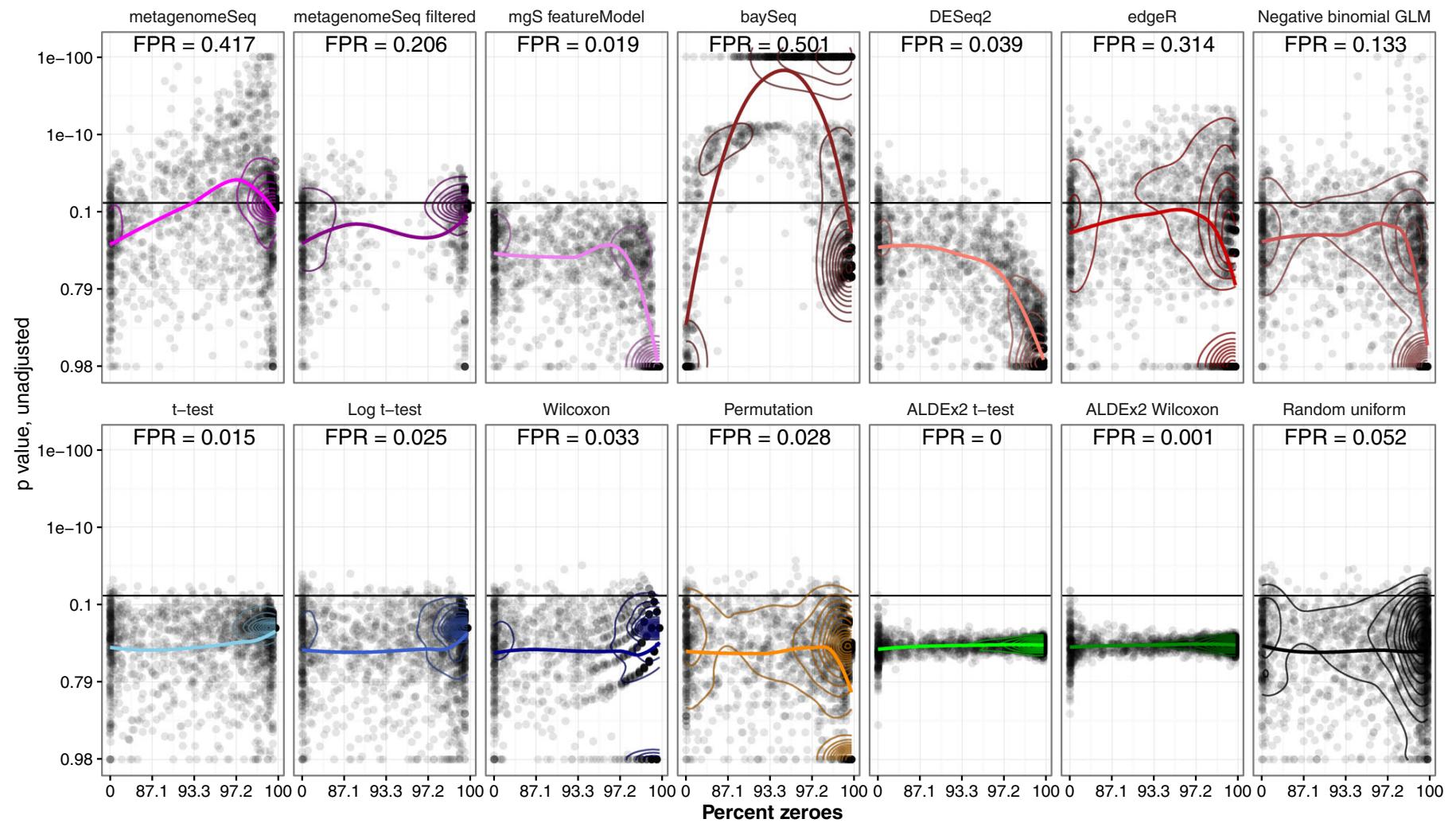


Fig. 2 OTU sparsity vs. p value. Scatterplots of OTU sparsity vs p value with panels representing each differential relative abundance test method in feces dataset A1, with 50% cases. Colored line represents the LOESS regression on data. False positive rate (FPR) is defined as the fraction of OTUs with $p < 0.05$. 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

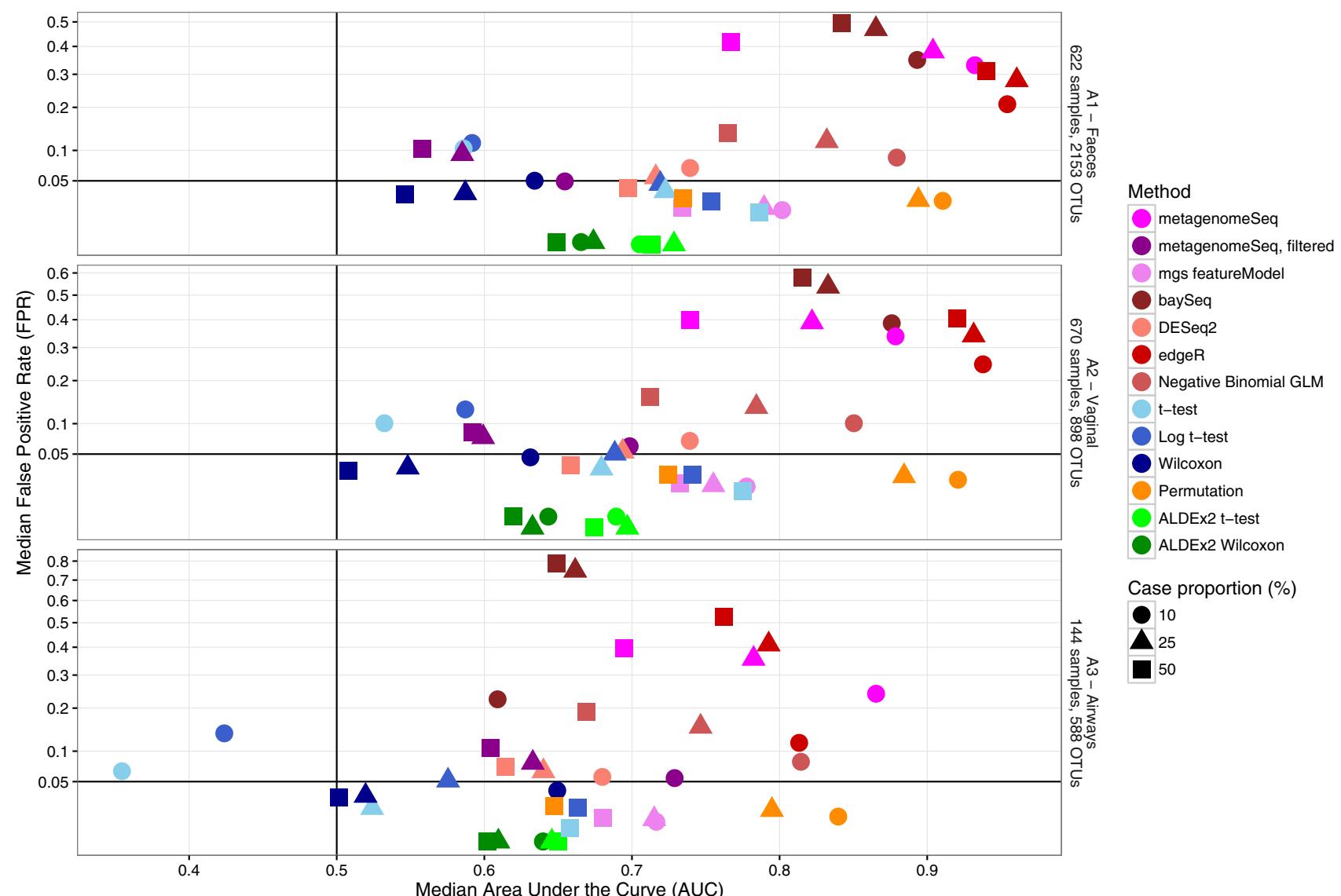


Fig. 3 Median test AUC vs median test FPR. Scatterplots of median test area under the curve (AUC) vs median test false positive rate (FPR), for the datasets A1–A3; three different compartments in the COPSAC₂₀₁₀ cohort. FPR is defined as the fraction of OTUs with $p < 0.05$. Dot color represents differential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

Where to cut?

Bonferroni'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*



Multomics



Data integration

-Omics

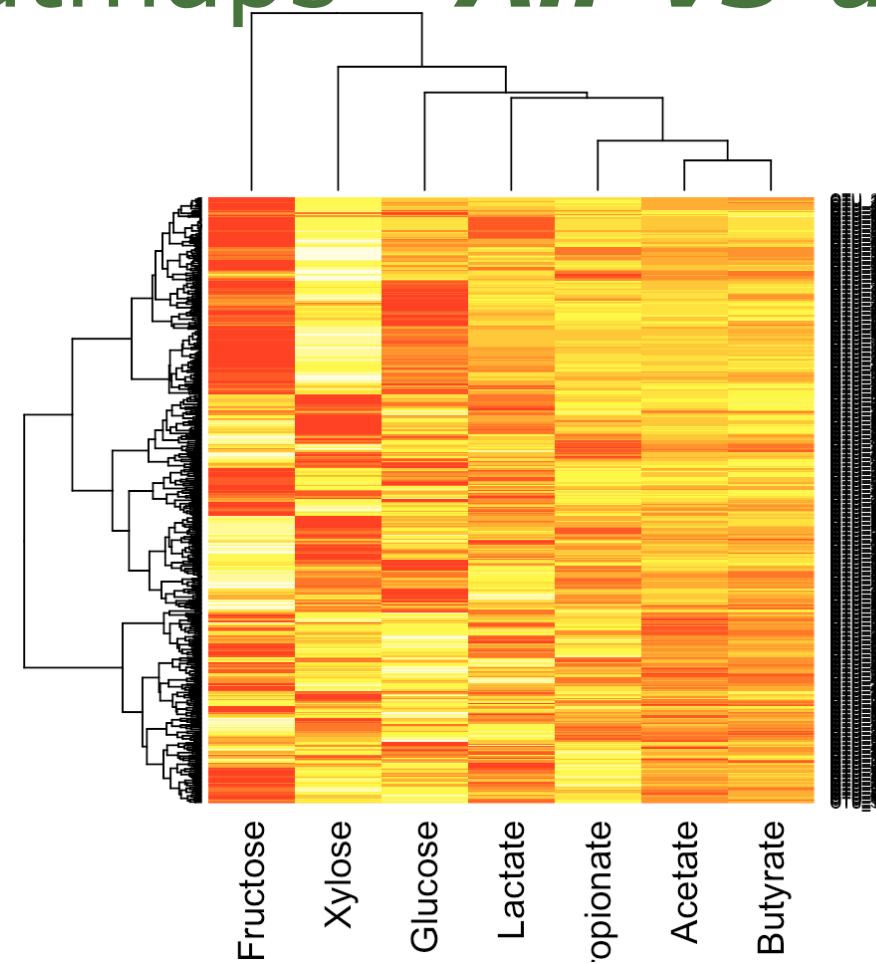
OTU table

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


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

$$\max_{\mathbf{w}_X, \mathbf{w}_Y} \mathbf{w}_X^T \mathbf{X}^T \mathbf{Y} \mathbf{w}_Y \quad \text{s.t.} \quad \begin{aligned} \mathbf{w}_X^T \mathbf{X}^T \mathbf{X} \mathbf{w}_X &= 1 \\ \mathbf{w}_Y^T \mathbf{Y}^T \mathbf{Y} \mathbf{w}_Y &= 1 \end{aligned}$$

Supervised CCA

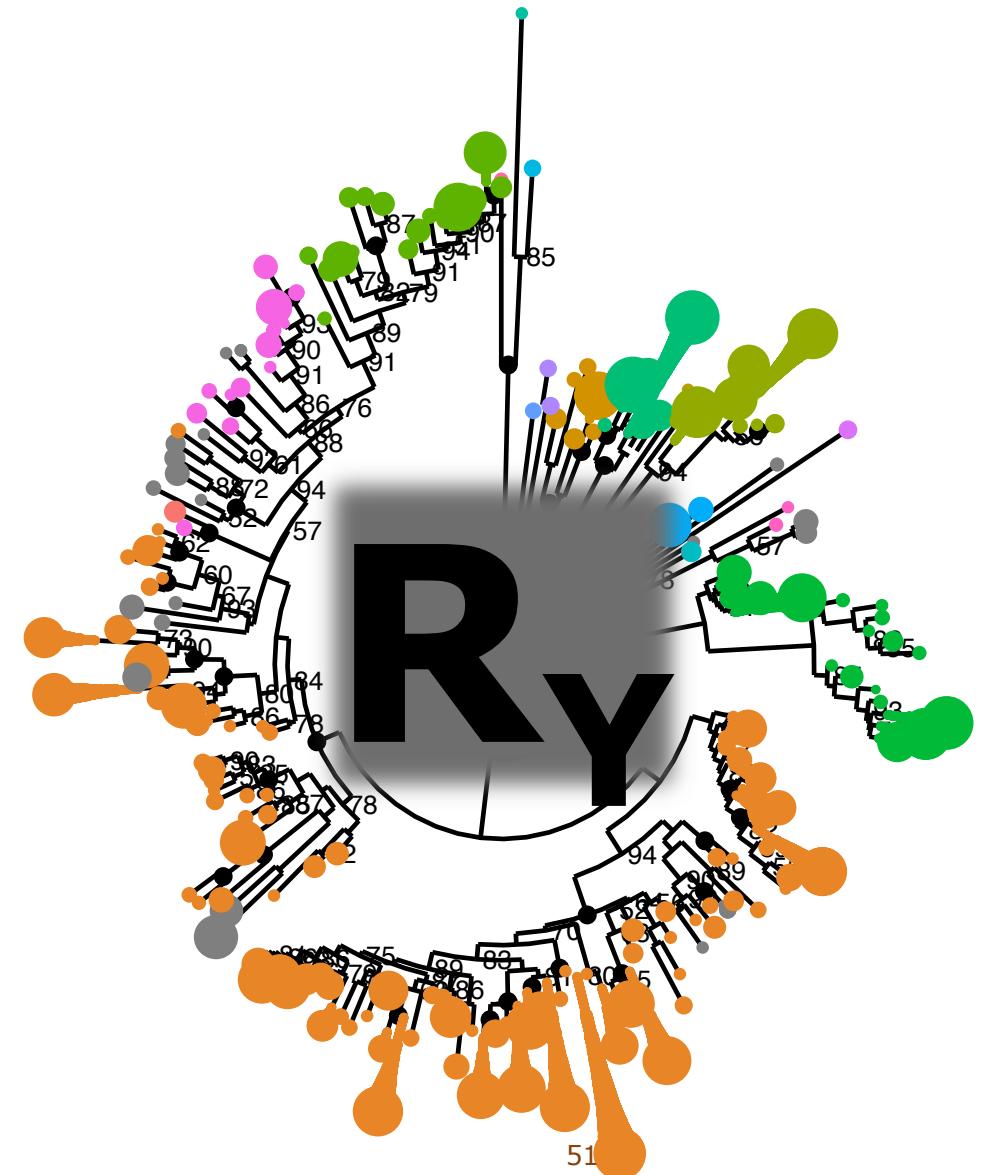
- Let \mathbf{R}_Y be a symmetric matrix reflecting similarity between the features in \mathbf{Y} based on external knowledge



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

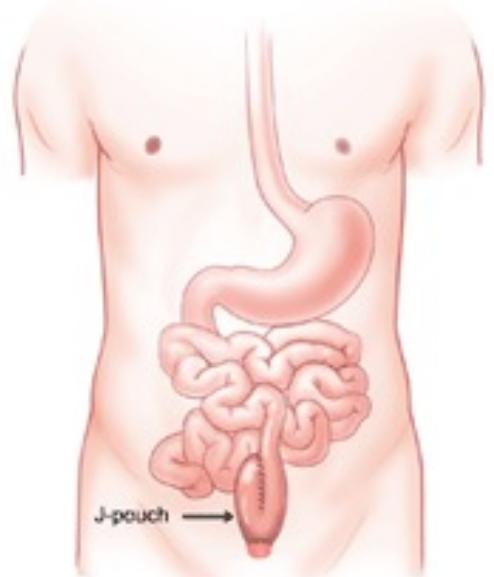
Kernel Smoothing

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

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



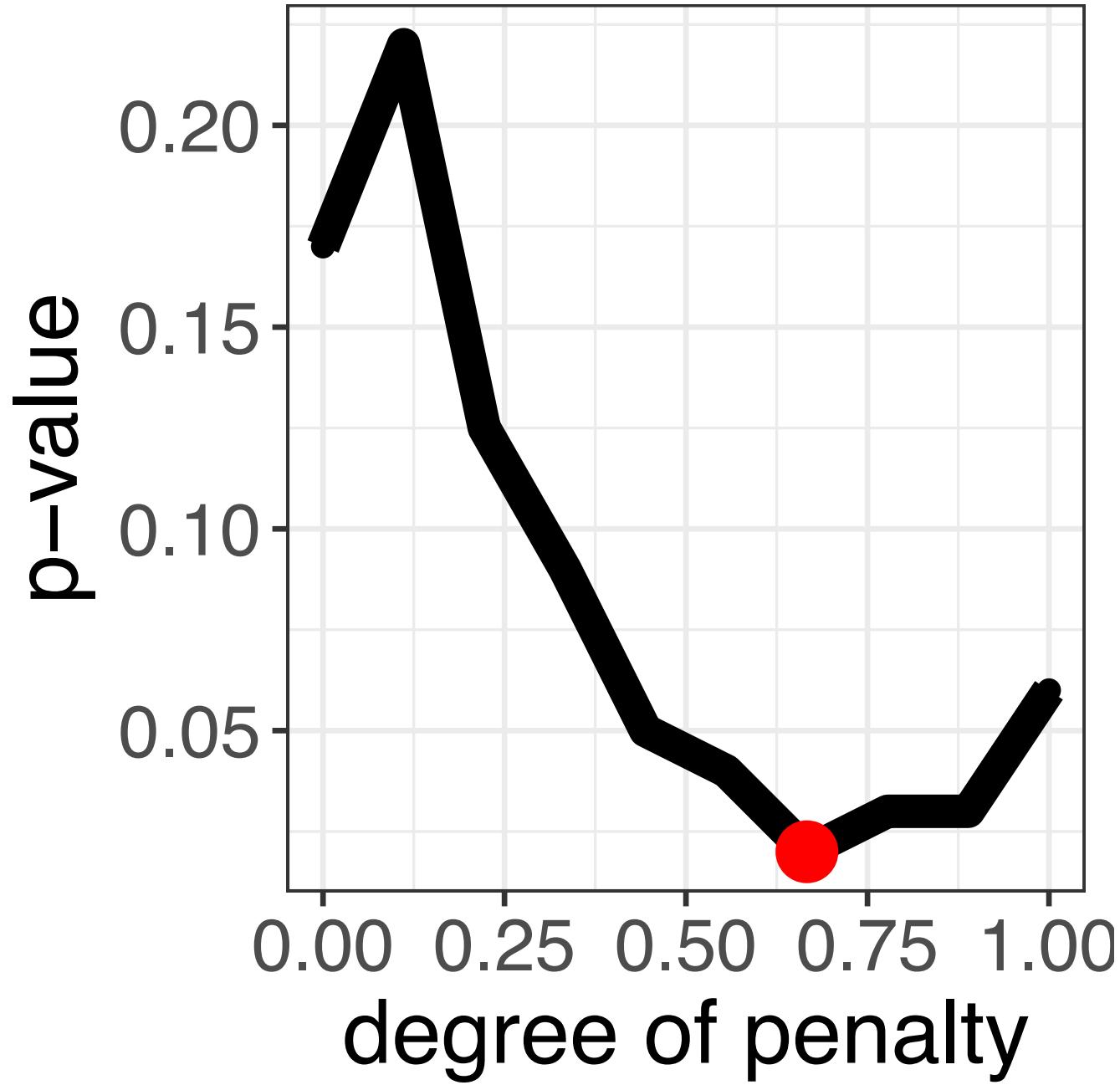
Example Pouchitic cohort

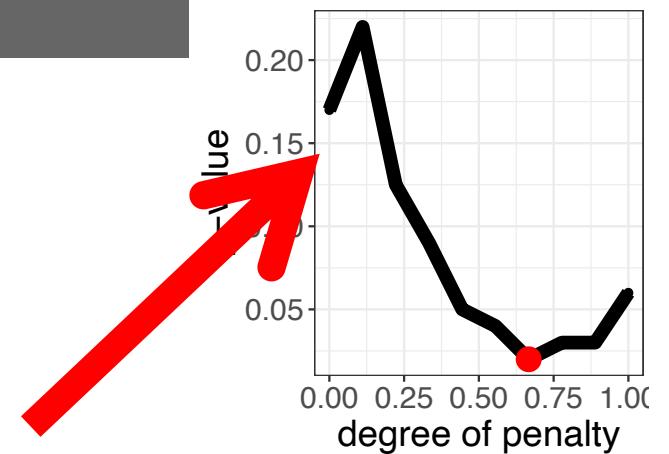


Gene
Expression

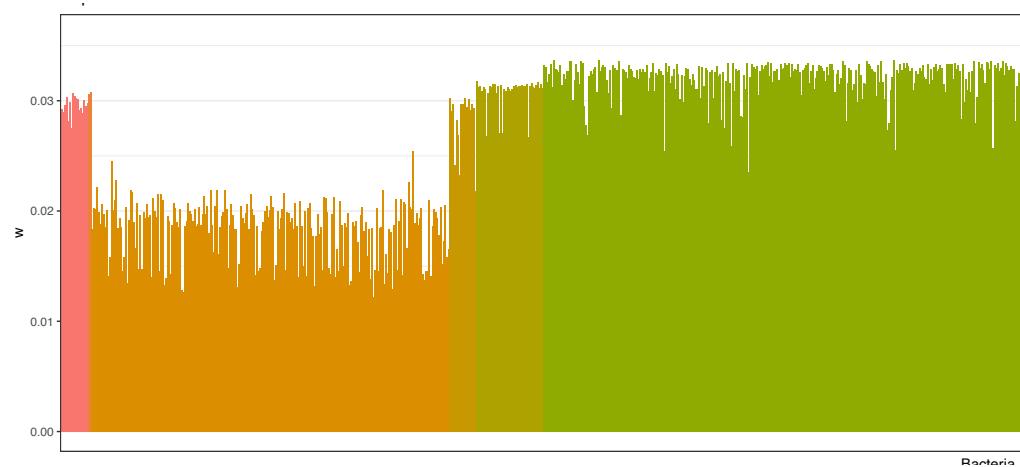
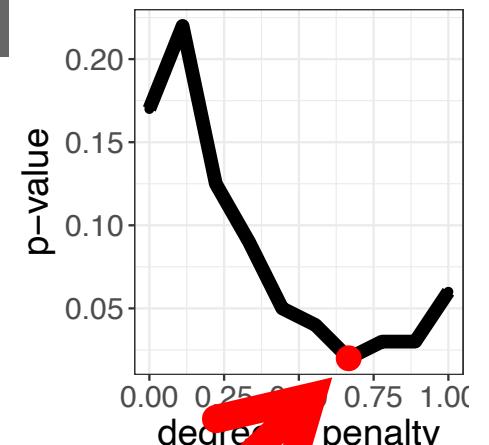
OTU table

248 X¹⁹⁹¹ 248 Y¹³⁸⁰





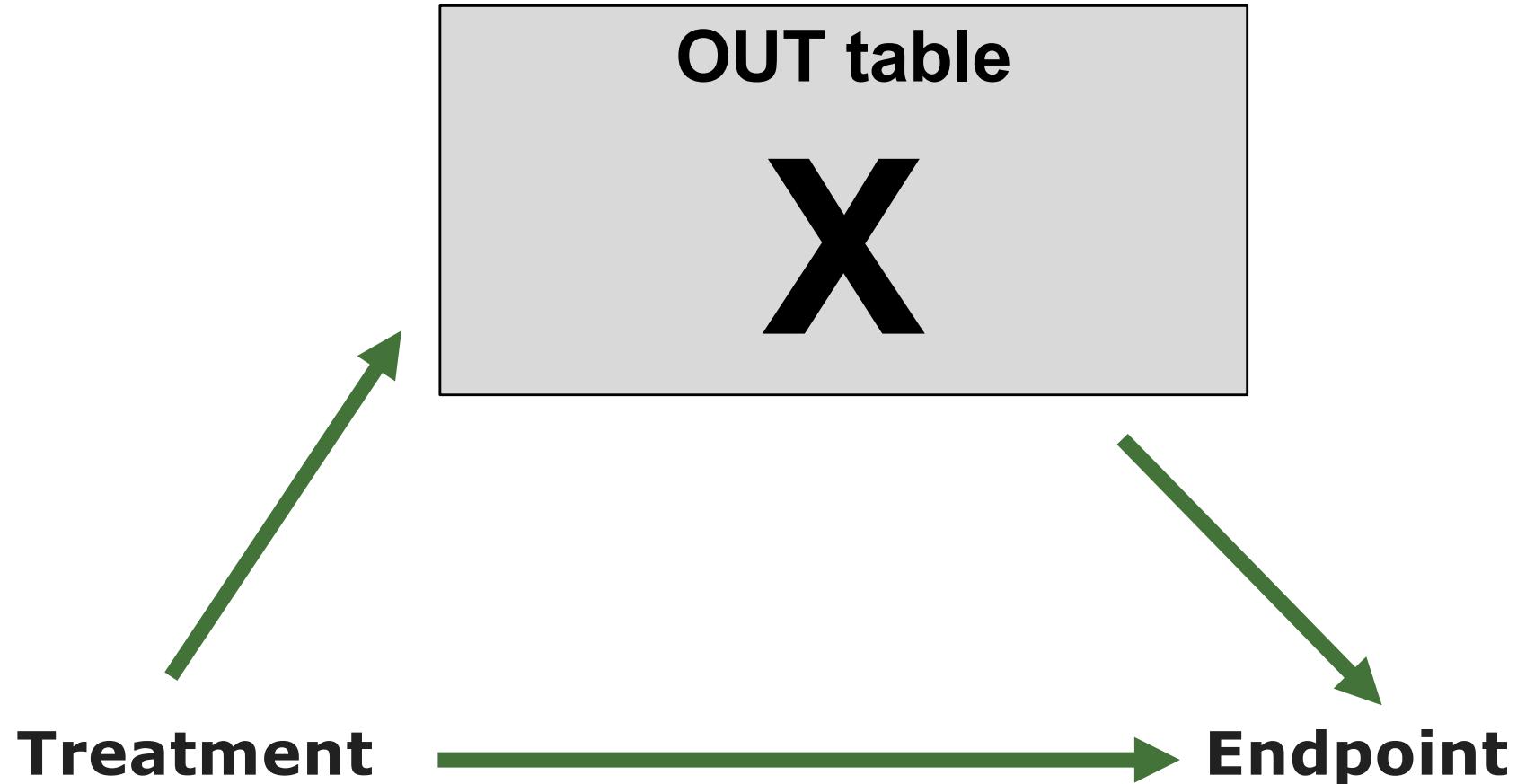
Order	Actinomycetales	Bifidobacteriales	Coriobacteriales	Flavobacteriales	Lactobacillales	Rhizobiales	Sphingomonadales	Xanthomonadales
	Bacillales	Burkholderiales	Enterobacteriales	Fusobacteriales	Pasteurellales	Rhodocyclales	Turicibacteriales	
	Bacteroidales	Clostridiales	Erysipelotrichales	Gemellales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



Truncates to
common loading
for similar bacteria

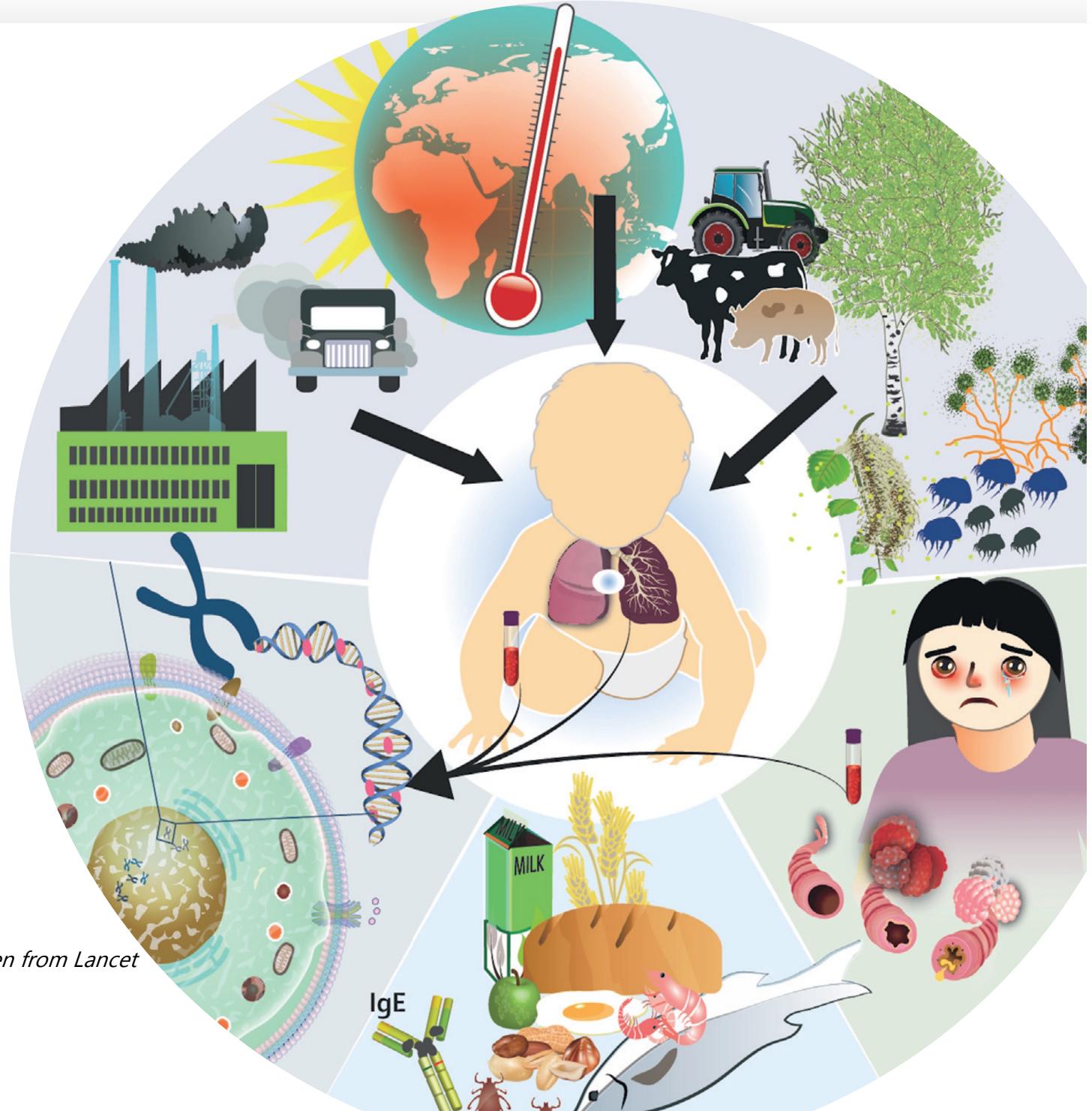
Order	Actinomycetales	Bifidobacteriales	Coriobacteriales	Flavobacteriales	Lactobacillales	Rhizobiales	Sphingomonadales	Xanthomonadales
	Bacillales	Burkholderiales	Enterobacteriales	Fusobacteriales	Pasteurellales	Rhodocyclales	Turicibacteriales	
	Bacteroidales	Clostridiales	Erysipelotrichales	Gemellales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

Mediation analysis

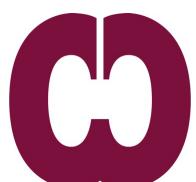
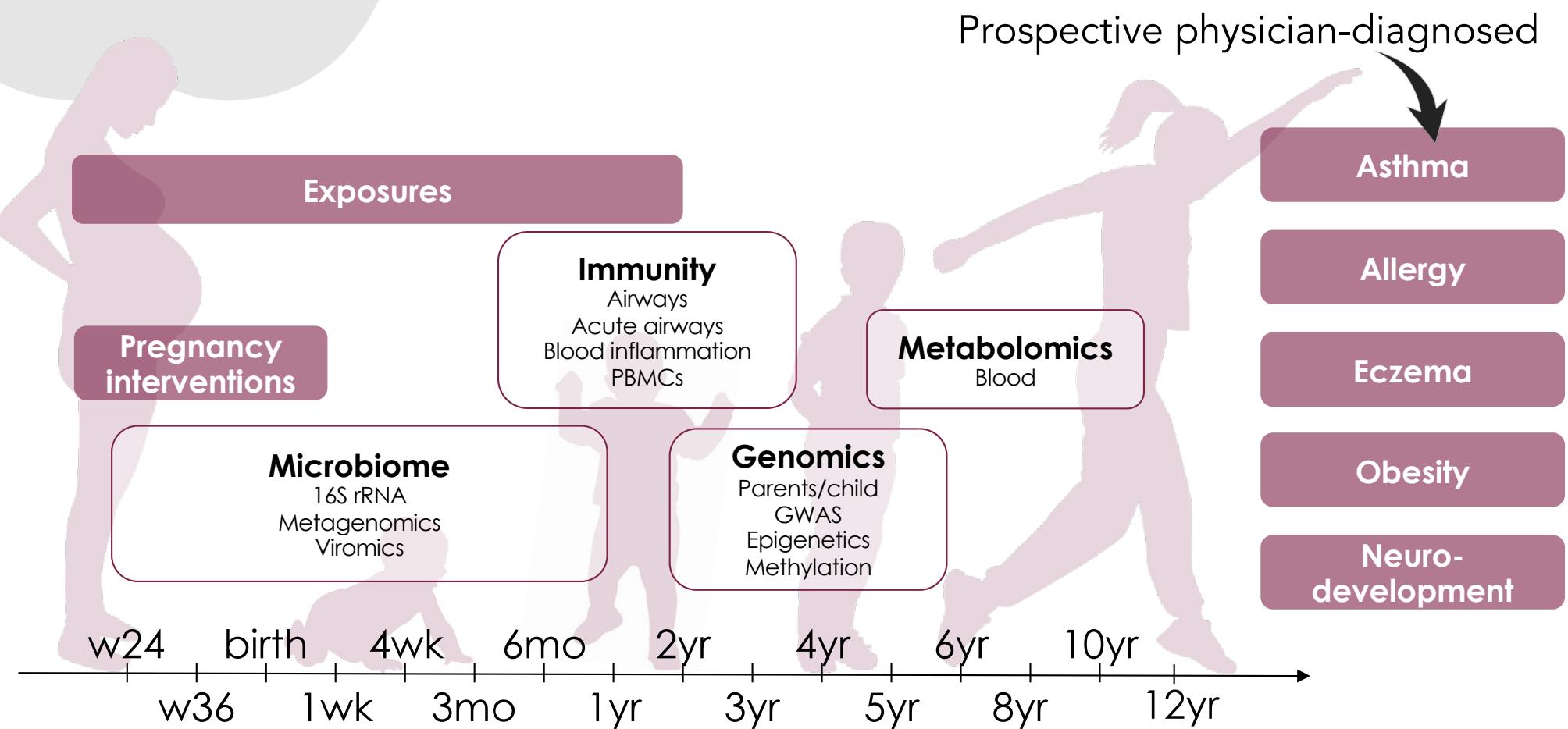


Exposure Propensities

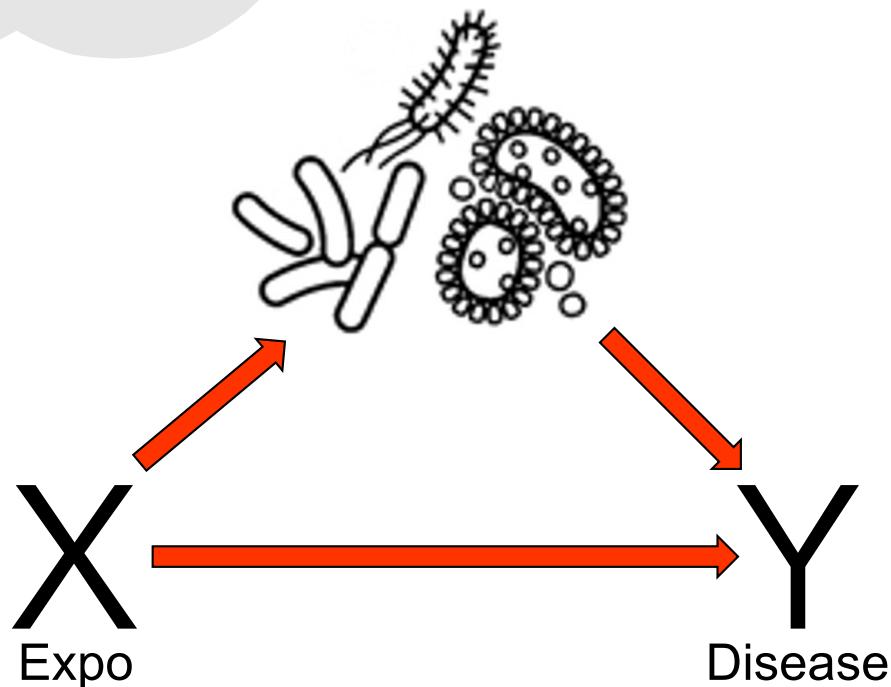
Image stolen from Lancet



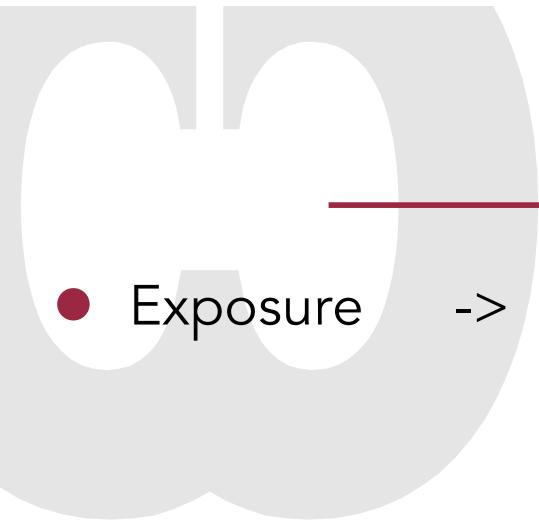
COPSAC₂₀₁₀ cohort design



Triangulation

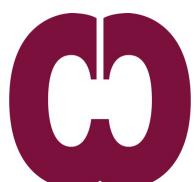


- Uncovers mediation
- Comes with the benefit of statistically controlling the *noisy* omics layer



A common question structure

- Exposure -> Omics -> Disease
- Siblings -> Microbiome -> Allergy
- Rural living -> Blood Metabolome -> Infections
- C-section -> Microbiome -> Asthma
- Diet -> Blood Metabolome -> Neuro





Establishing an: Omics-derived Exposure Propensity Score *Fingerprint*

- Build between exposure models
- Use C model
- Use the prediction

Example

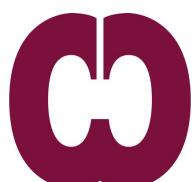
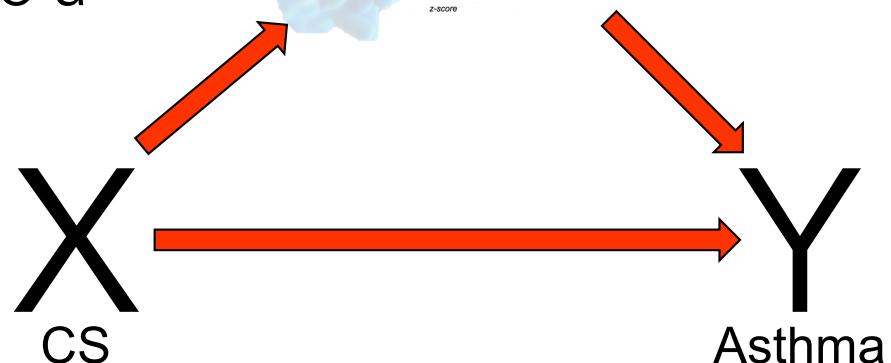
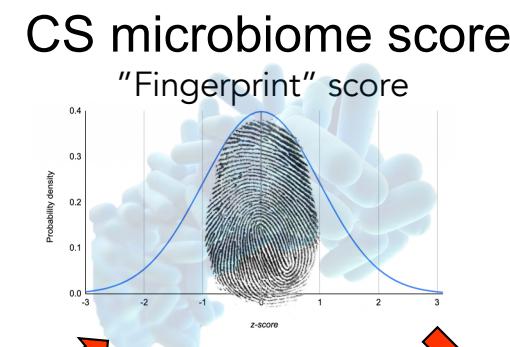
Cesarean section score

High values = your microbiome looks like you were born by cesarean section

Low values = your microbiome looks like you were born vaginally

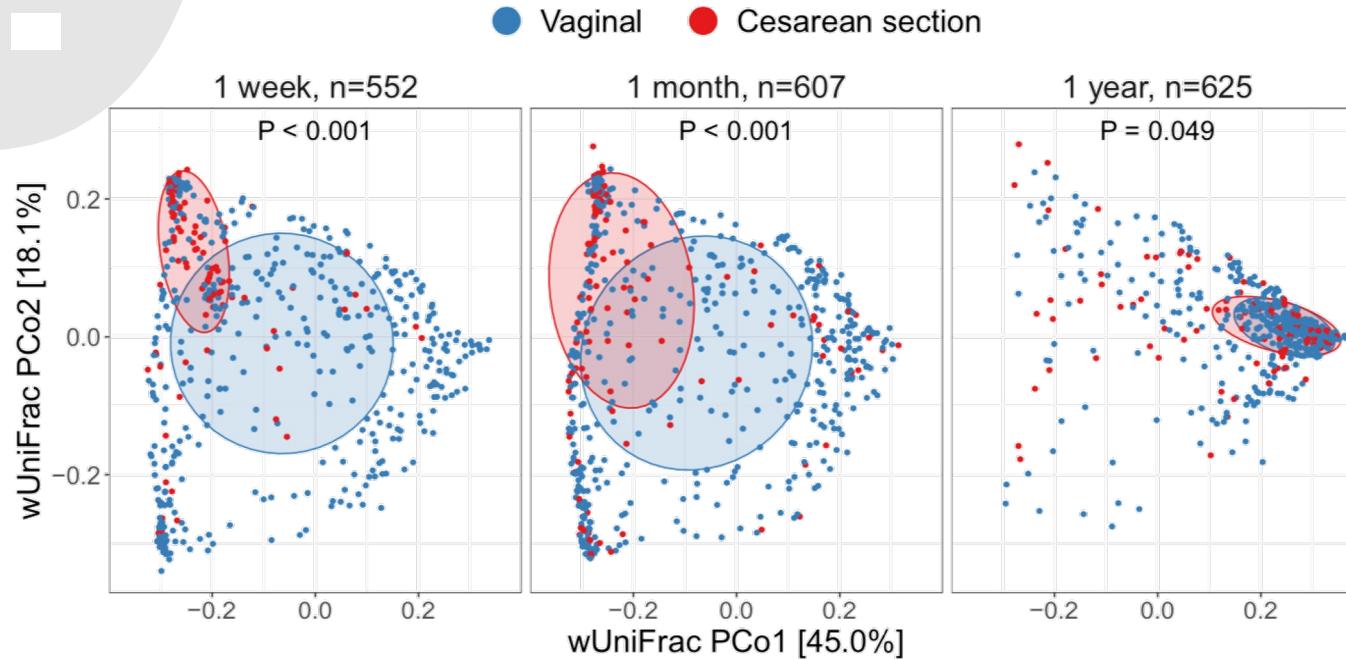
Fingerprint mediation analysis

- 1) Establish clinical association
 - 2) Derive fingerprint
 - 3) Examine mediation
-
- Conclusion: Is the microbiome a mediator of the association between X and Y?



Cesarean section and gut microbiome

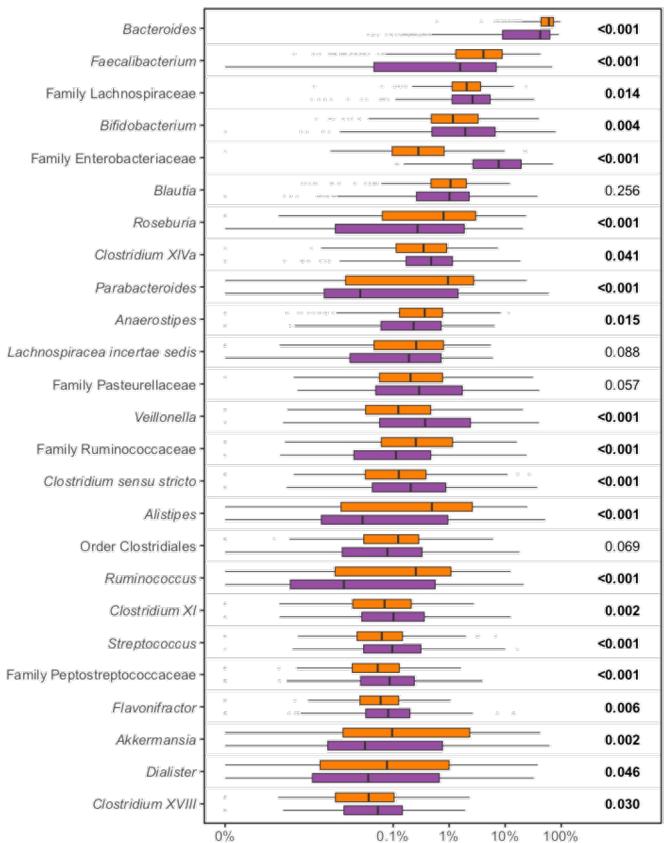
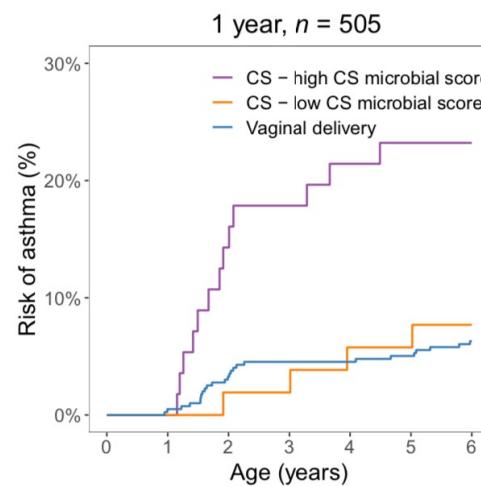
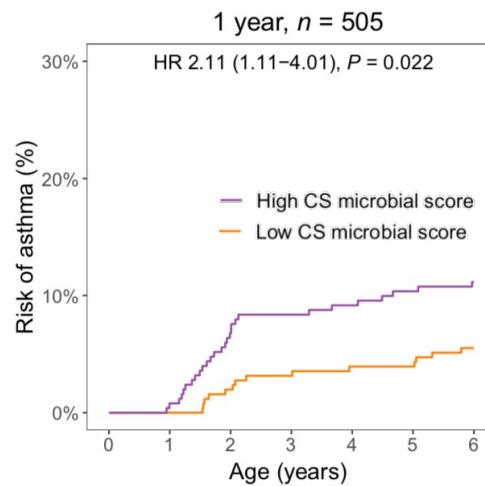
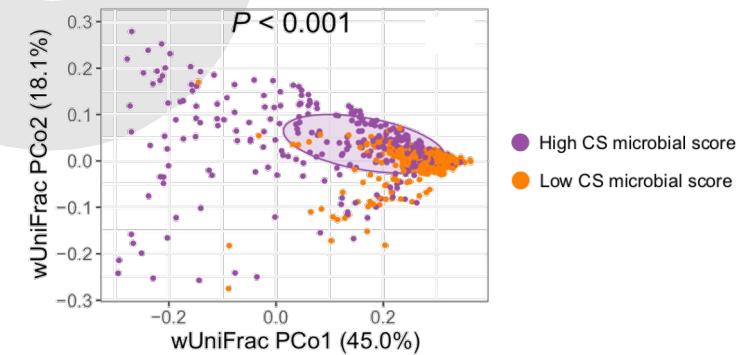
 COPSAC₂₀₁₀ cohort



Is the CS-asthma association mediated by differences in the gut microbiome?

Cesarean section and gut microbiome

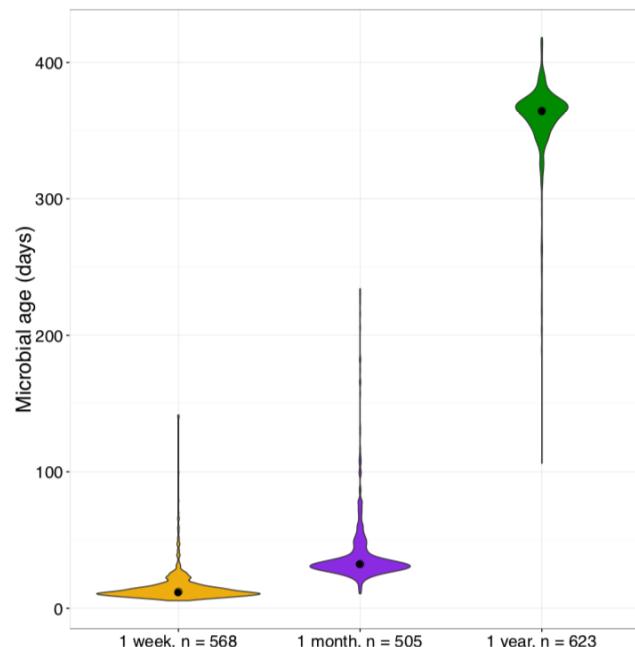
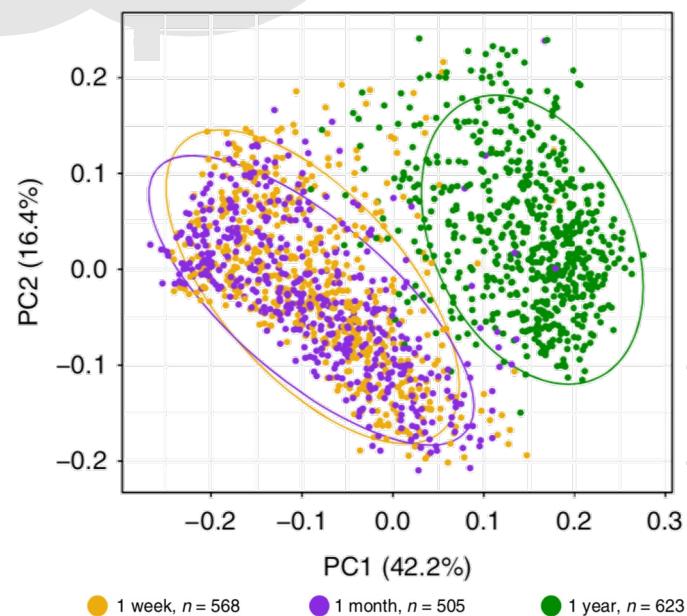
 COPSAC₂₀₁₀ cohort



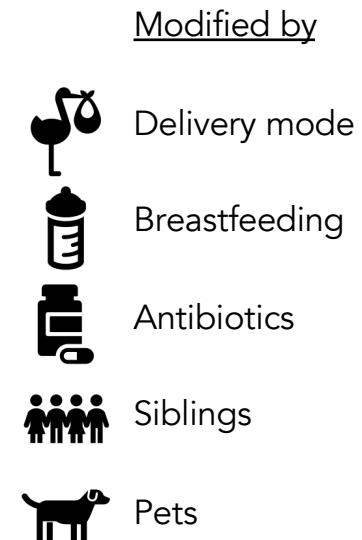
Gut microbiome maturation

- COPSAC₂₀₁₀ cohort, N=700

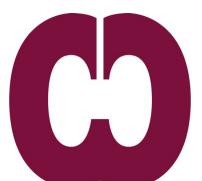
COPSAC₂₀₁₀ cohort



Microbiome maturation
MAZ – Microbiota-by-Age Z-scores
Random Forest model

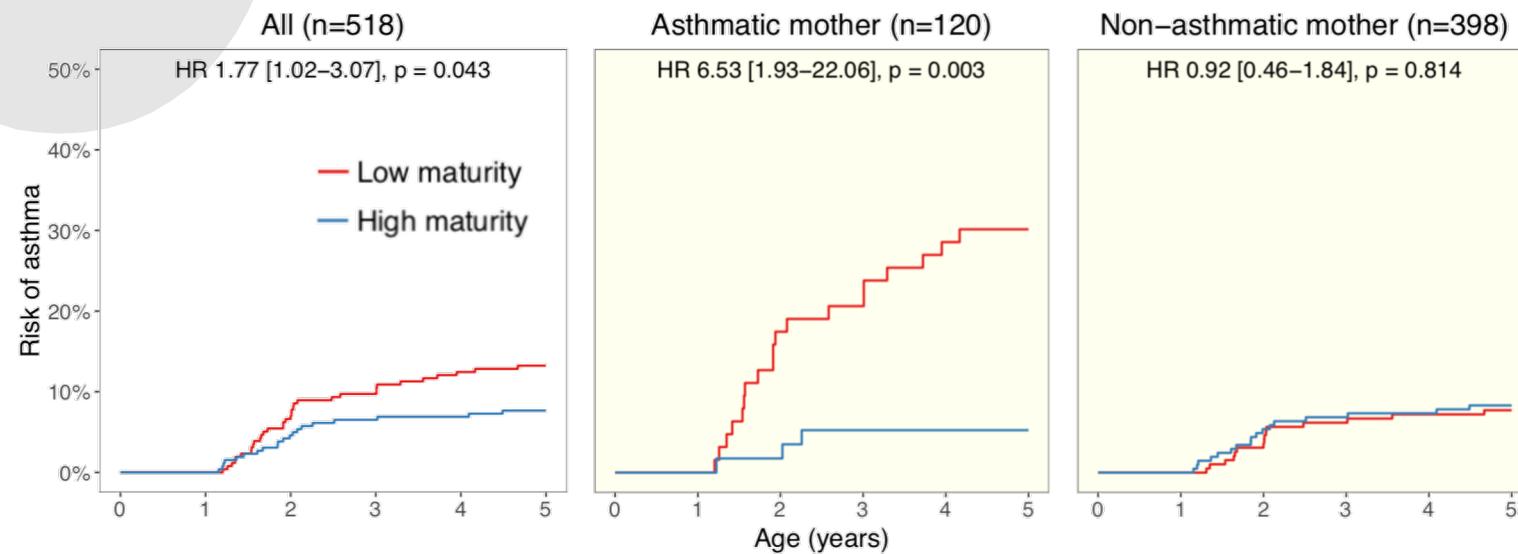


Stokholm et al, *Nat Comm* 2018; Subramanian et al, *Nature* 2014;
Stewart et al, *Nature* 2018; Bokulich et al, *Sci Trans Med* 2016



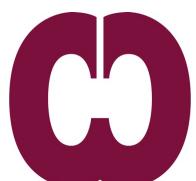
Maturation and asthma

COPSAC₂₀₁₀ cohort

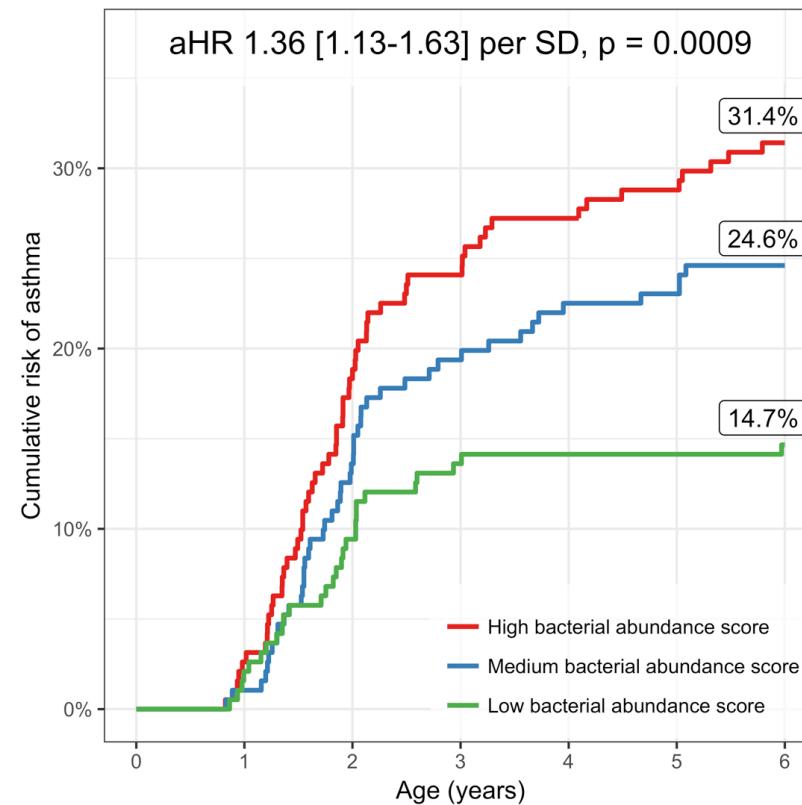
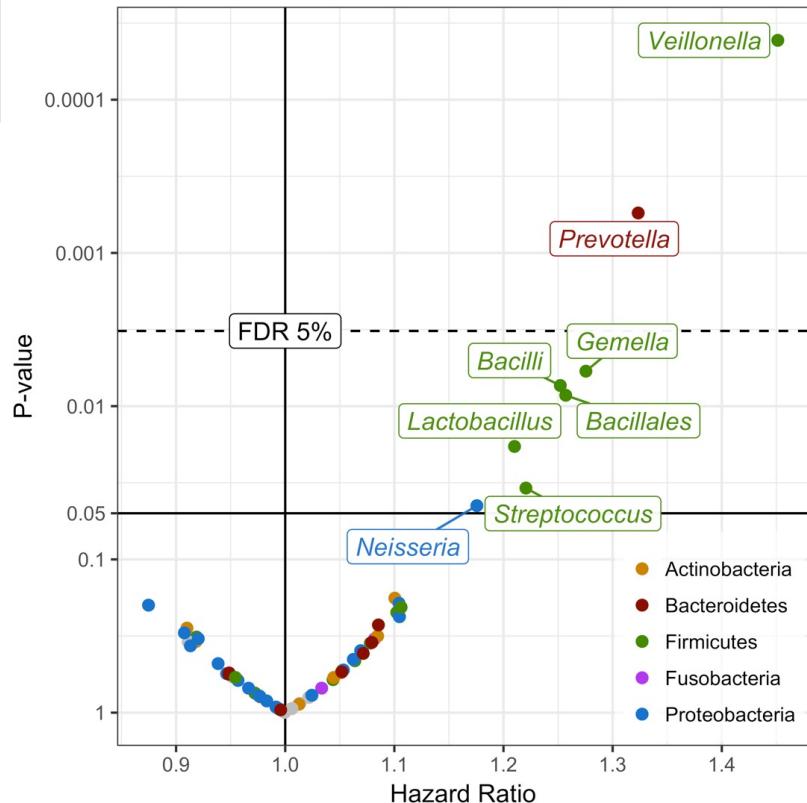


1-year microbiome maturity and asthma by age 5

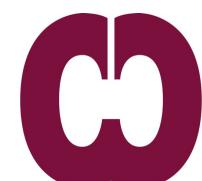
Stokholm et al, *Nat Comm* 2018



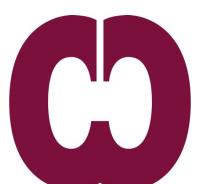
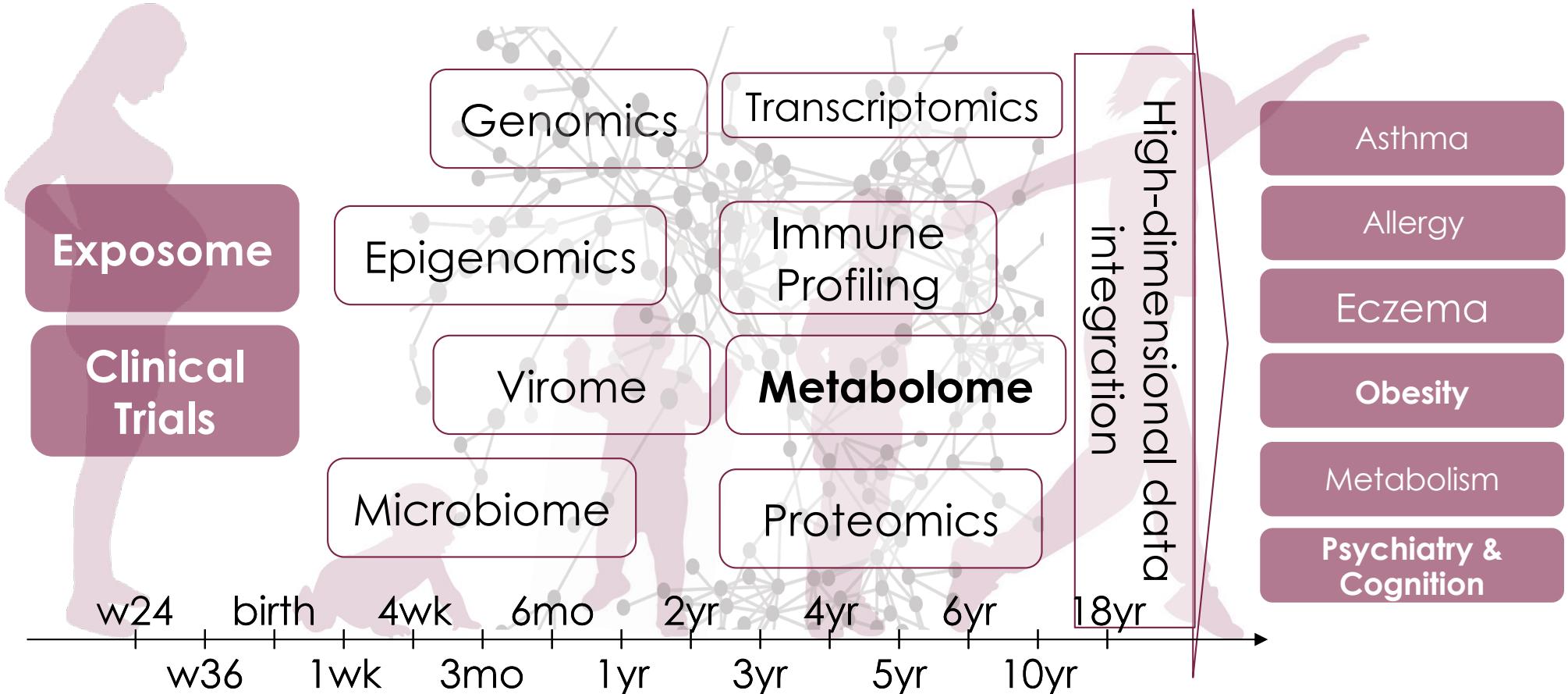
Infant airway microbiome and asthma



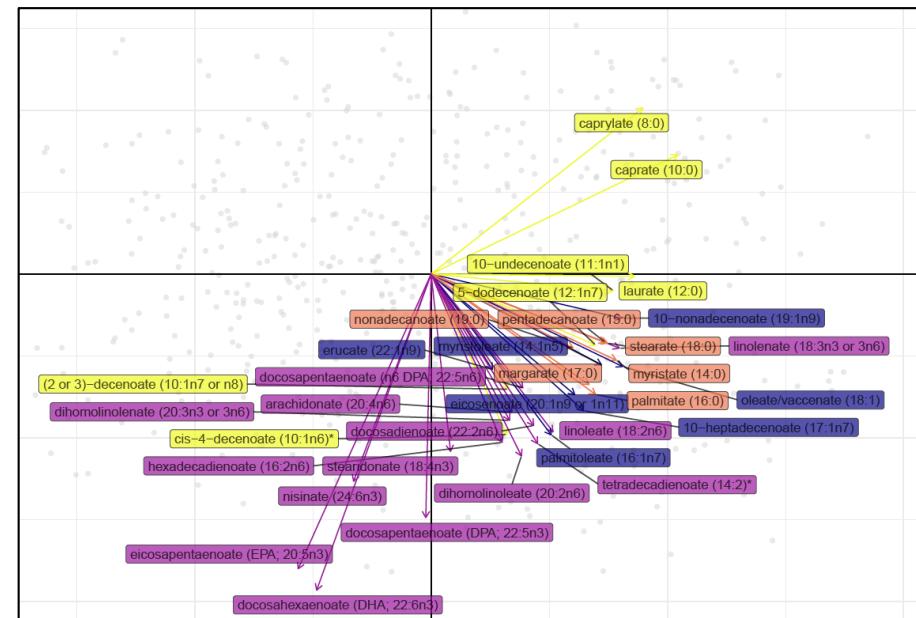
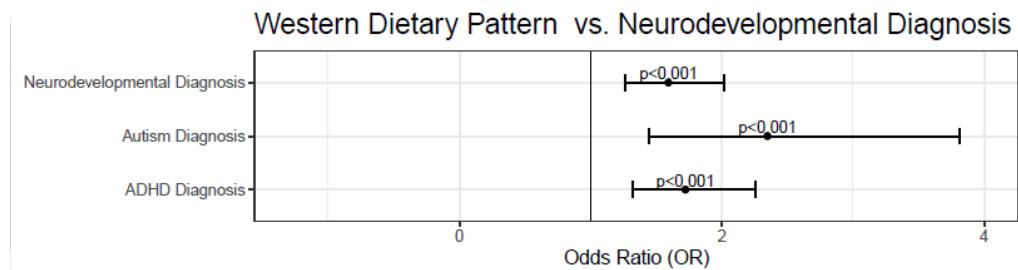
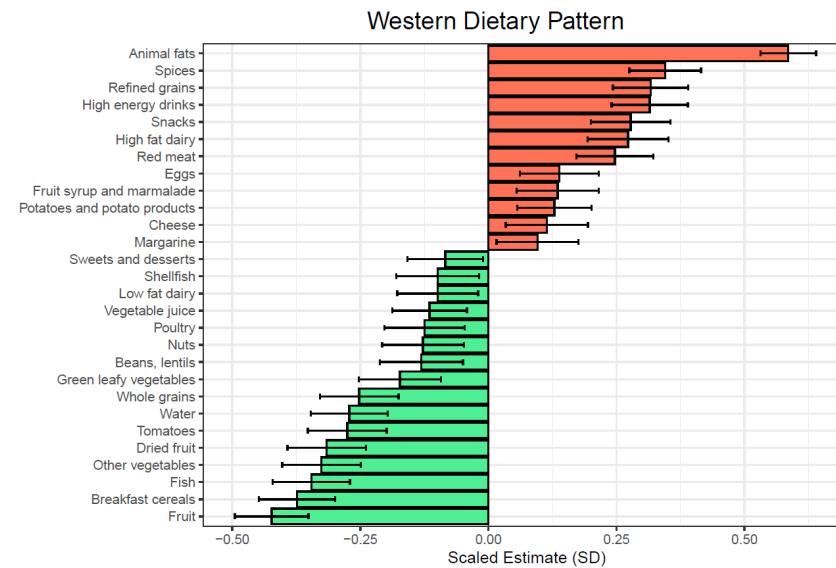
Thorsen et al, Nature Comm. 2019



Another example from COPSAC2010

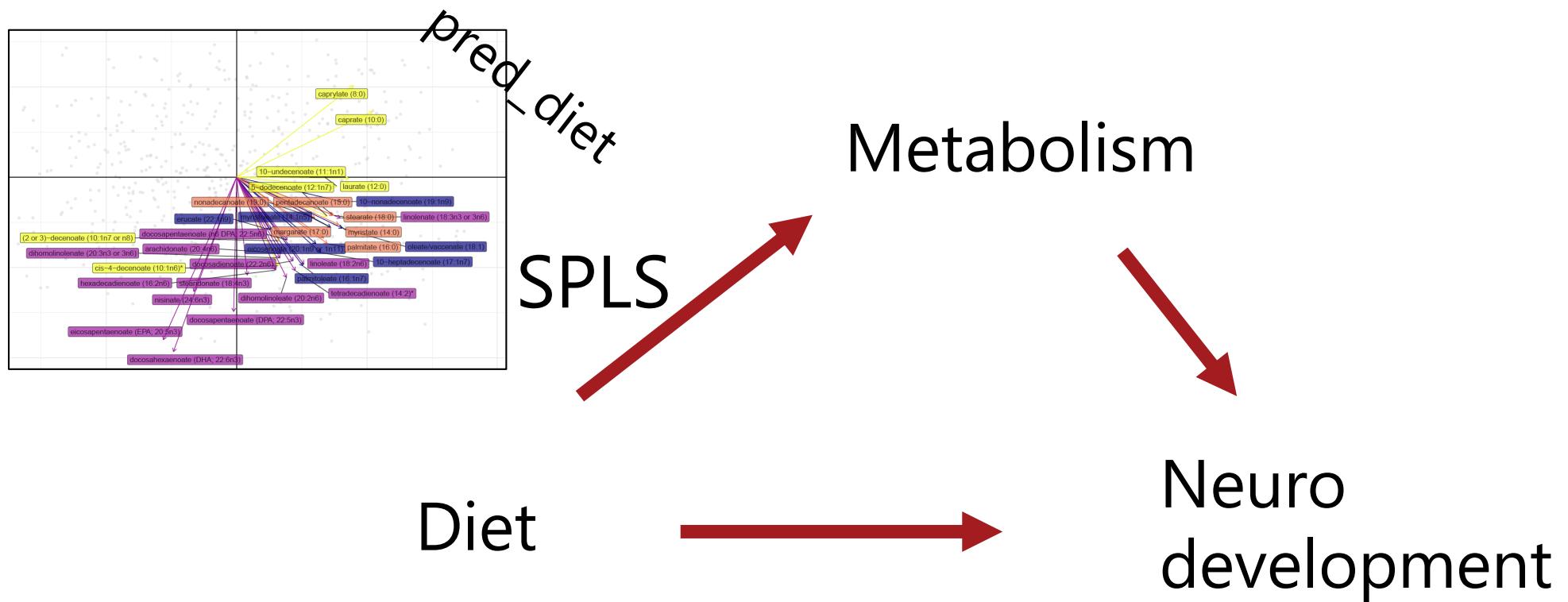


Diet in Pregnancy and offspring health



By the use of **blood metabolomics** in both pregnancy and childhood we can point towards the **critical window** in time is **during pregnancy** and not just a marker of a constant poor family diet. Further, elucidate mechanisms of inflammatory metabolic products

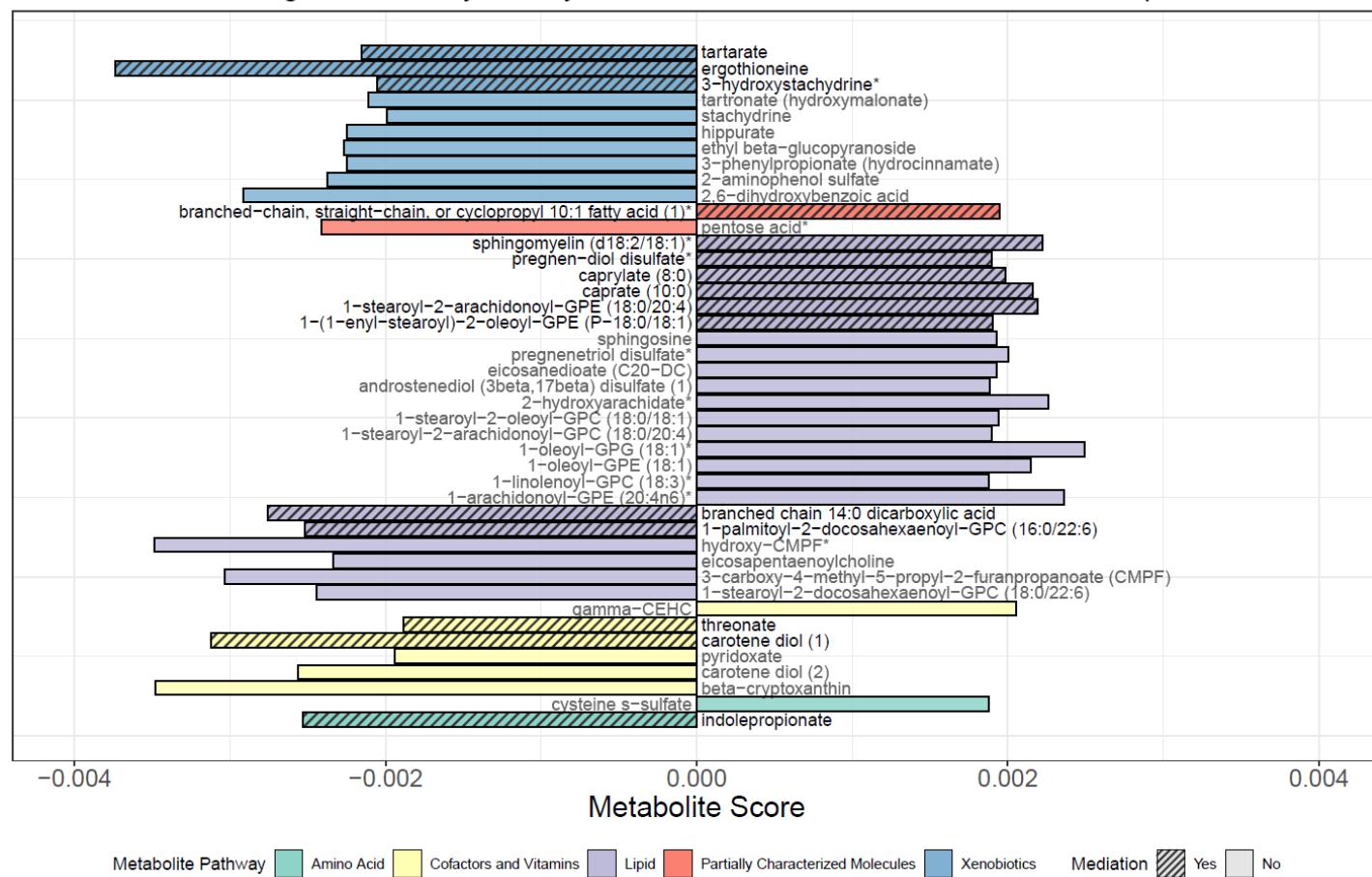
Metabolomics as mediator



Metabolomics as mediator

Boolean
backward
elimination
Keep/kill
without
re-estimation

Metabolites Reflecting an Unhealthy Dietary Pattern and Their Mediation in Neurodevelopmental Disorders



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.
- In case of a triangulation setup, use propensity scoring and mediation to add some “causality”

