

1

UNIVERSITY OF COPENHAGEN Enhedens navn

* 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](https://mortenarendt.github.io/MicrobiomeDataAnalysis/index.html)

Sted og dato
Dias 2



2

1

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

Sted og dato
Dias 3



3

Outline

Day0 STAT101 recap Get R working Preprocessing Alpha diversity	Day1 Morning Beta diversity Testing design versus Beta diversity	Day2 Morning More DA heatmaps
	Day1 Afternoon Differential Abundance testing	Day2 Afternoon Multiomics with heatmaps and CCA

Sted og dato
Dias 5



5

UNIVERSITY OF COPENHAGEN

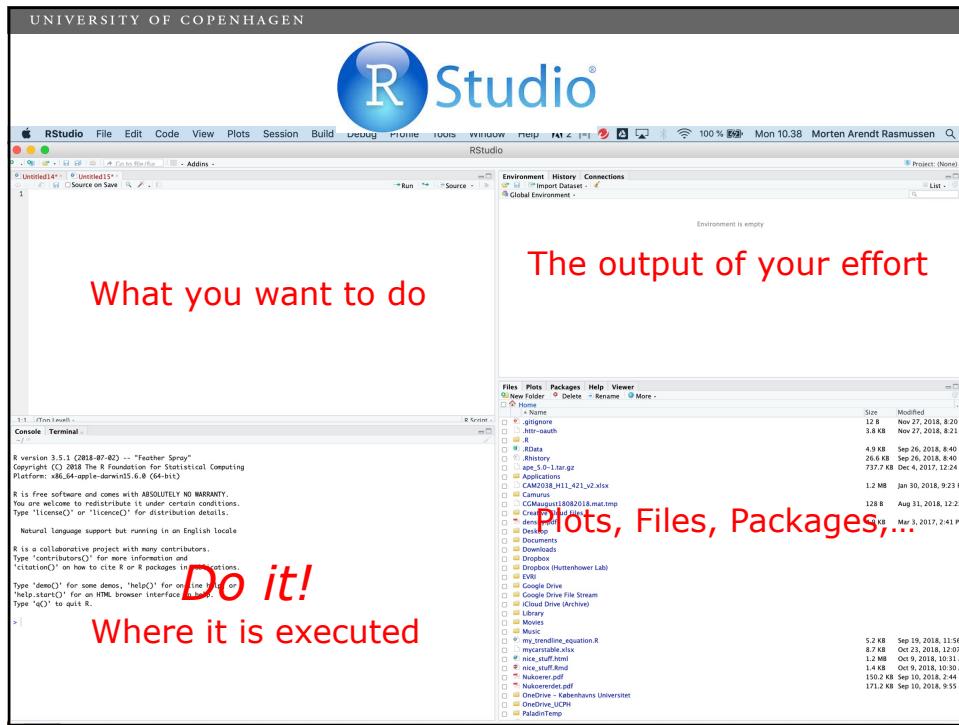
Why R??

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

6

UNIVERSITY OF COPENHAGEN

7



8

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.

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

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

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

Dias 9

9

Data import

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

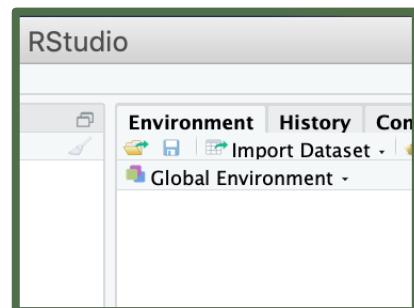
.... Eventually copy paste the commands produced into the script

Or

Use the **rio** package

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

Dias 10



10

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()	

Dias 11



11

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()

Dias 12



12

STAT 101 recap

A Linear model

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

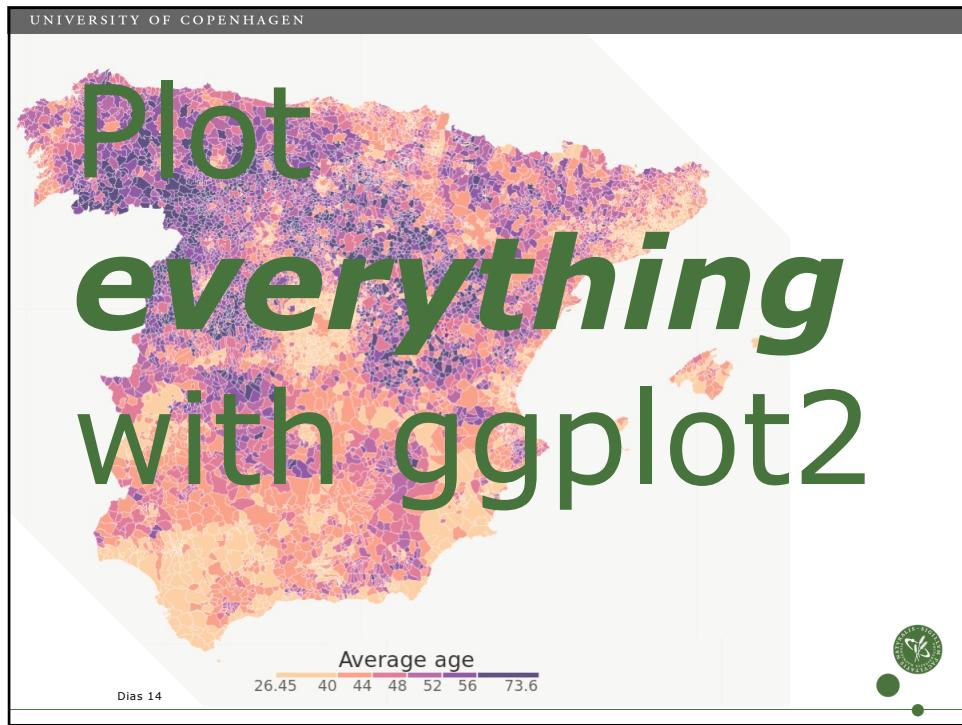
The workhorse of understanding effects of interventions, handling covariates and confounders, and getting uncertainties on all this

Test statistics

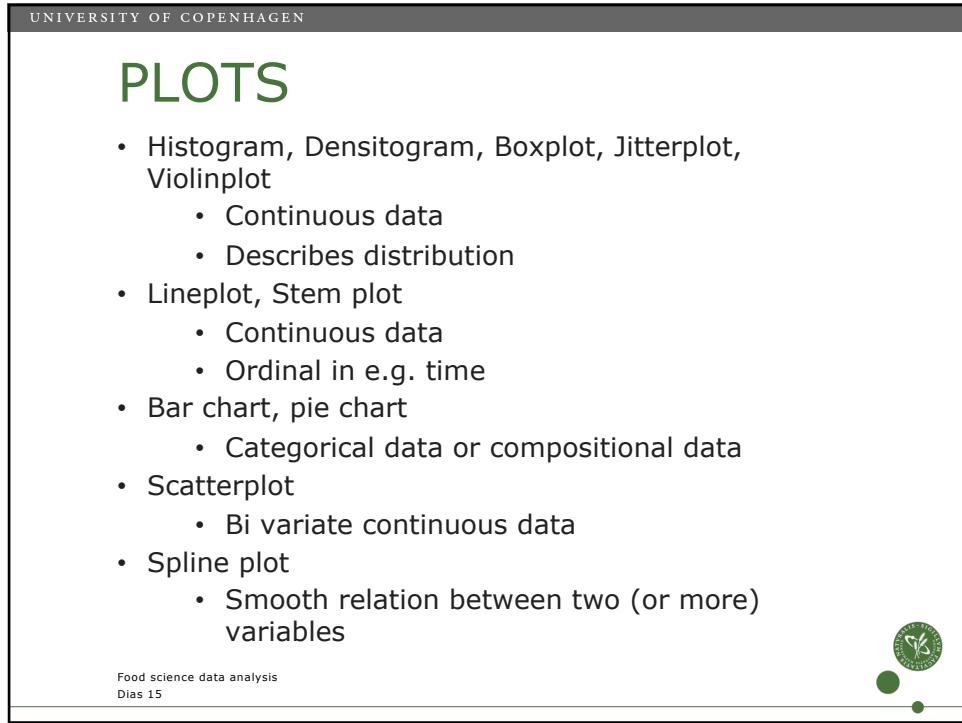
A metric/scalar measure of observed differences, which we can construct a test for

Sted og dato
Dias 13

13



14

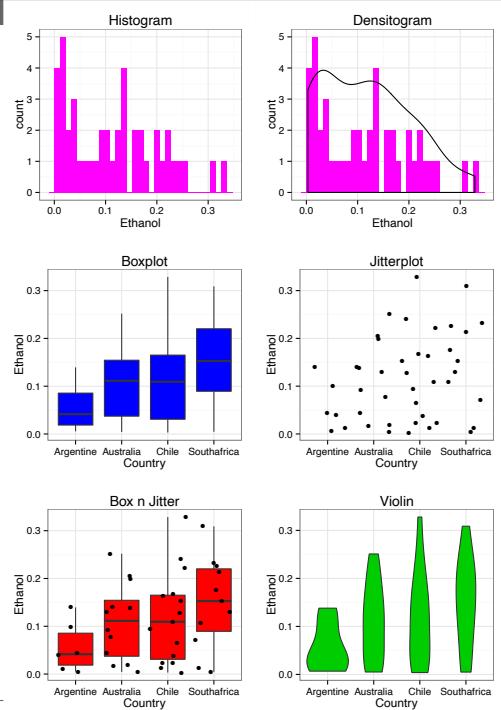


15

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

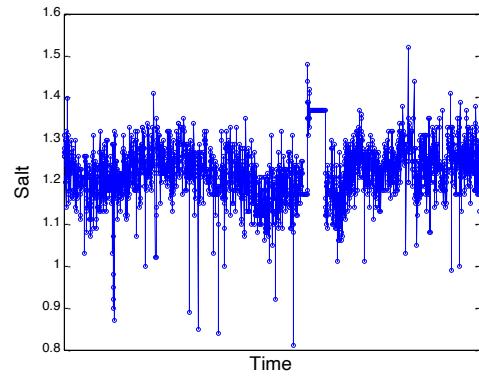
Continuous data
Describes distribution



16

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.



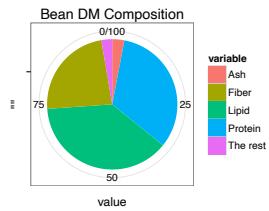
Continuous data
Ordinal in e.g. time



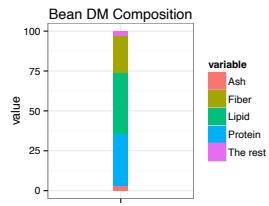
17

Bar charts

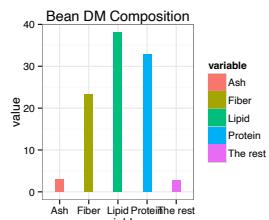
Pie charts



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



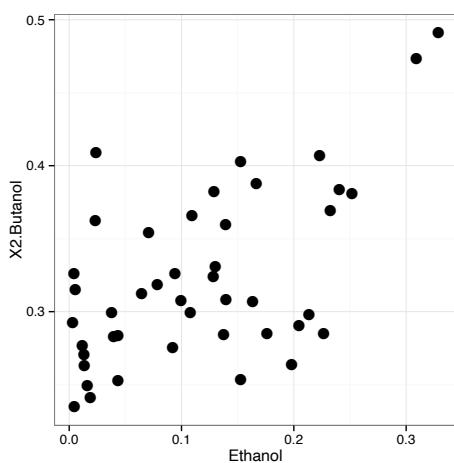
Compositional data
Categorical data



Food science data analysis
Dias 18

18

Scatterplot

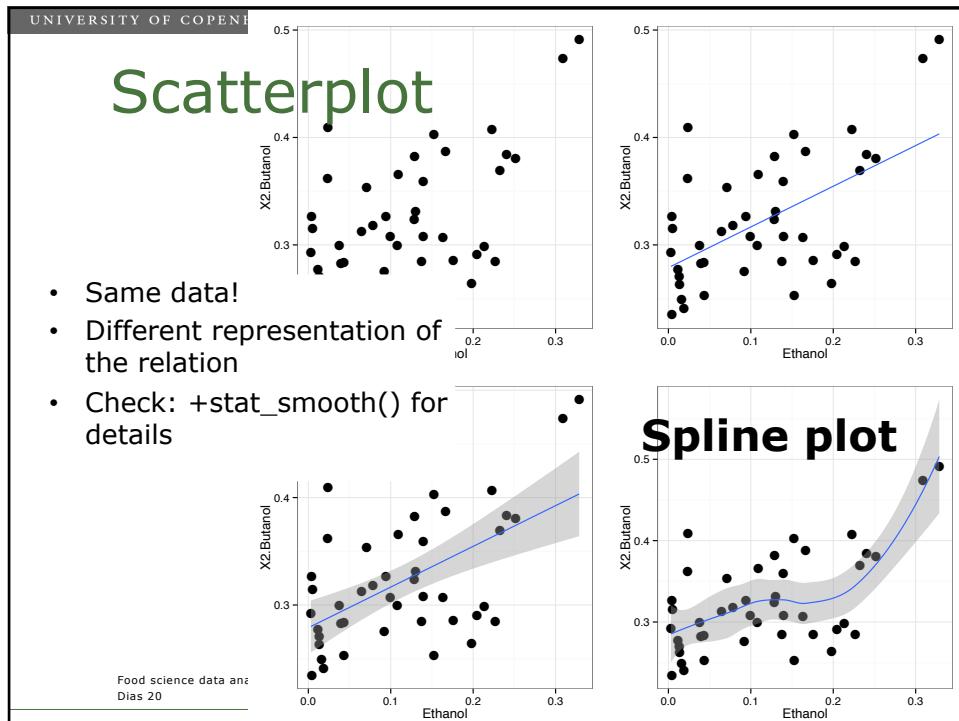


- Two response variables versus each other
- Very useful for Multivariate data

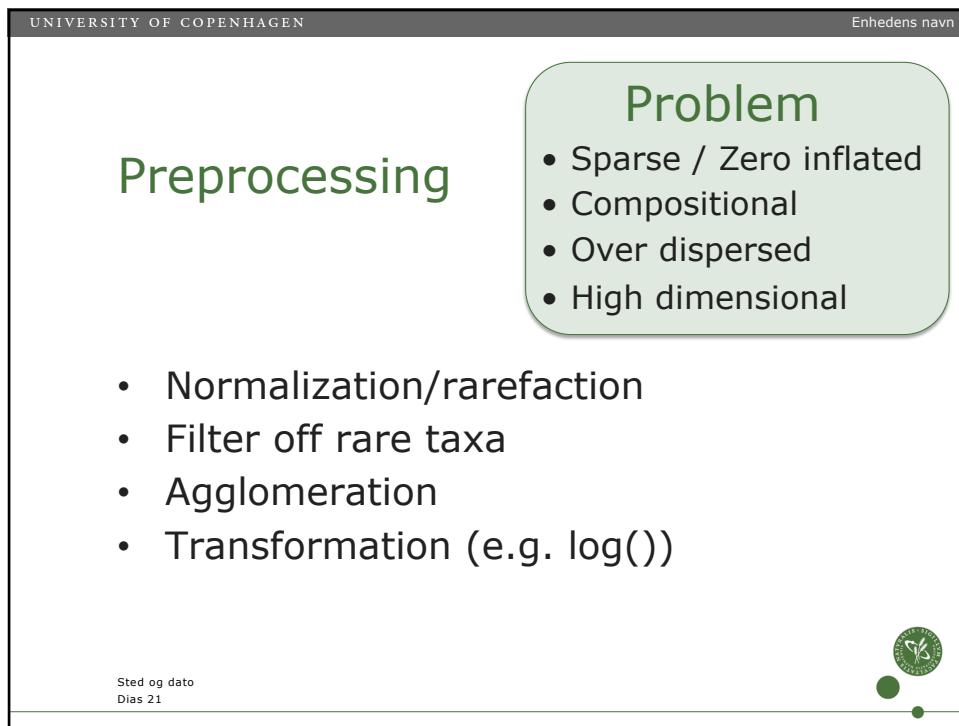
Food science data analysis
Dias 19



19



20



21

UNIVERSITY OF COPENHAGEN

Enhedens navn

Diversity metrics

Alpha diversity

Within sample characteristics

Beta diversity

Between sample characteristics

Sted og dato
Dias 22



22

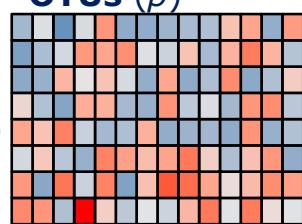
UNIVERSITY OF COPENHAGEN

Quality and Technology

Amplicon

OTUs (p)

Samples (n)



0 100 200 300 400 500 600 700 800 900 1000 1100 1200 1300 1400 1500 bp

V1 V2 V3 V4 V5 V6 V7 V8 V9

CONSERVED REGIONS: unspecific applications
VARIABLE REGIONS: group or species-specific applications



23

Alpha diversity

Number of different taxa

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

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

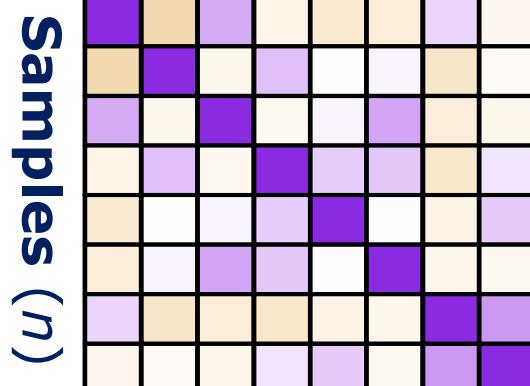
Sted og dato
Dias 24



24

β diversity

Samples (n)



25

25

β diversity

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

Sted og dato
Dias 26



26

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

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

Sted og dato
Dias 27



27

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.

Sted og dato
Dias 28

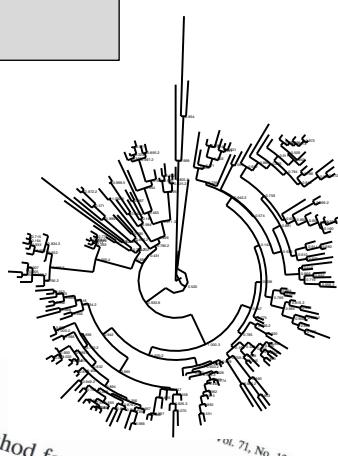


28

UNIFRAC

OTU table

Dist



© 2005, AMERICAN MICROBIOLOGY, Dec. 2005, p. 8228-8235
 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 Knight²
¹Department of Molecular, Cellular, and Developmental Biology, University of Colorado 80309¹ and Department of Computer Science, University of Colorado 80309²



29

Ordination



Dias 30

30

β diversity to PCoA

Samples (n)

Purple	Yellow	Purple	Yellow	Yellow	Yellow	Purple	Yellow
Yellow	Purple	Yellow	White	Purple	Yellow	Yellow	White
Purple	White	White	Purple	White	Purple	Yellow	Yellow
White	Purple	White	Purple	Purple	White	Yellow	Purple
Yellow	White	Purple	White	Purple	White	Yellow	White
Purple	Yellow	White	Yellow	White	Purple	Purple	White
White							
White							

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

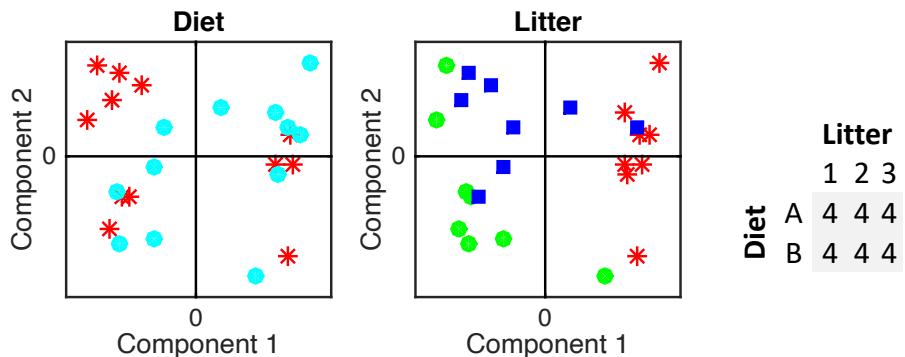


31

31

Multidimensional scaling

$$\mathbf{U}\Lambda\mathbf{U}^T = \mathbf{M}\exp(-\mathbf{Dist})\mathbf{M}^T$$



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



32

Differential Abundance Testing *or* OTUWAS

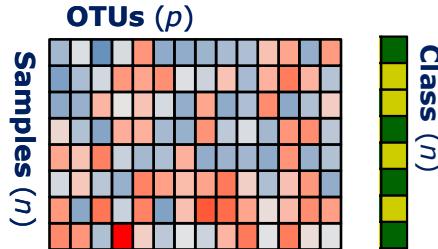


Dias 33

33

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 34



34

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

35

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

Sted og dato
Dias 36



36

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

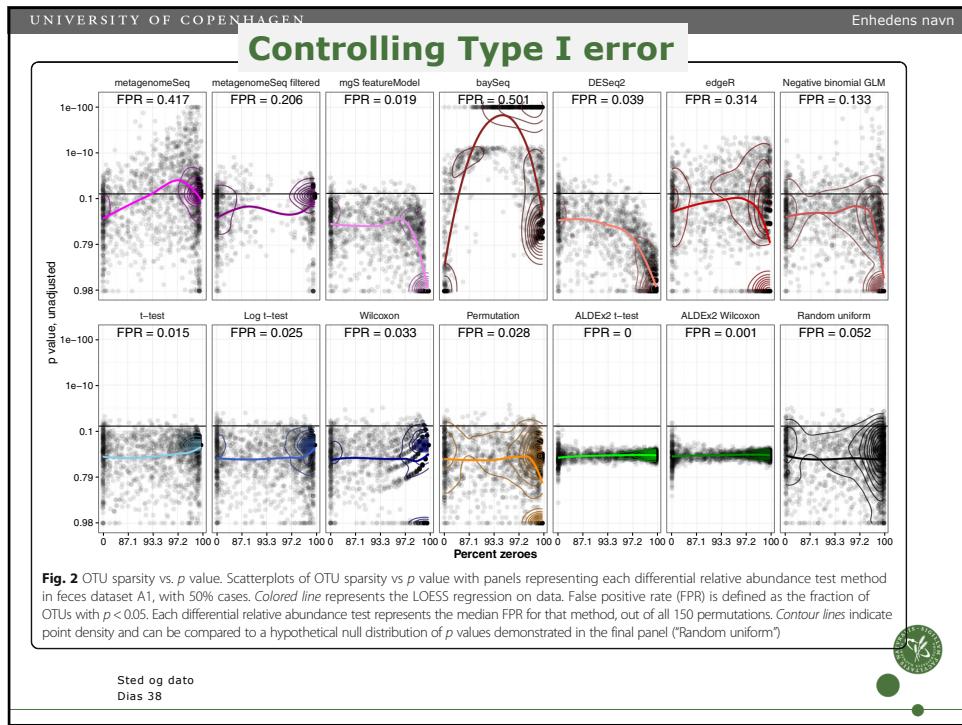
(uses cumulative sum scaling for prepro)

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

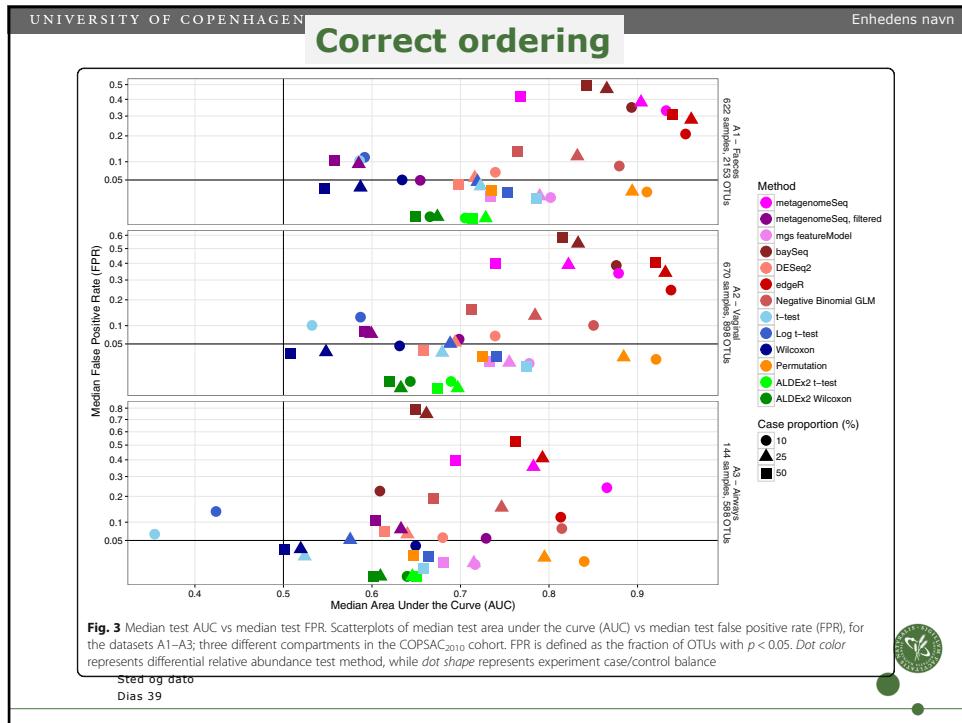
Sted og dato
Dias 37



37



38



39

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*

Sted og dato
Dias 40



40

Testing Beta diversity *PerMANOVA*

Dias 41



41

Partitioning when we do not see the original data

The underlying model

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

Variance partitioning

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

$$\begin{aligned} \hat{Y}\hat{Y}^T &= H(YY^T)H & YY^T \propto \exp(-Dist) \\ H &= X(X^T X)^{-1} X^T \end{aligned}$$

Dias 42

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.

42

Establishing the *null* distribution

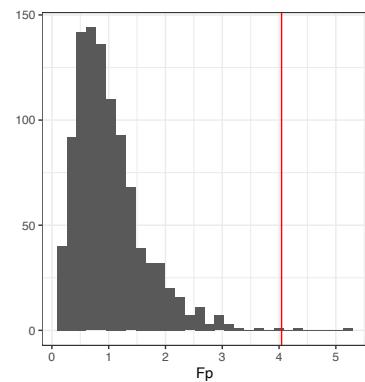
For univariate ANOVA models, we rely on independence and homoschedastic and gaussian *like* residuals.

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?

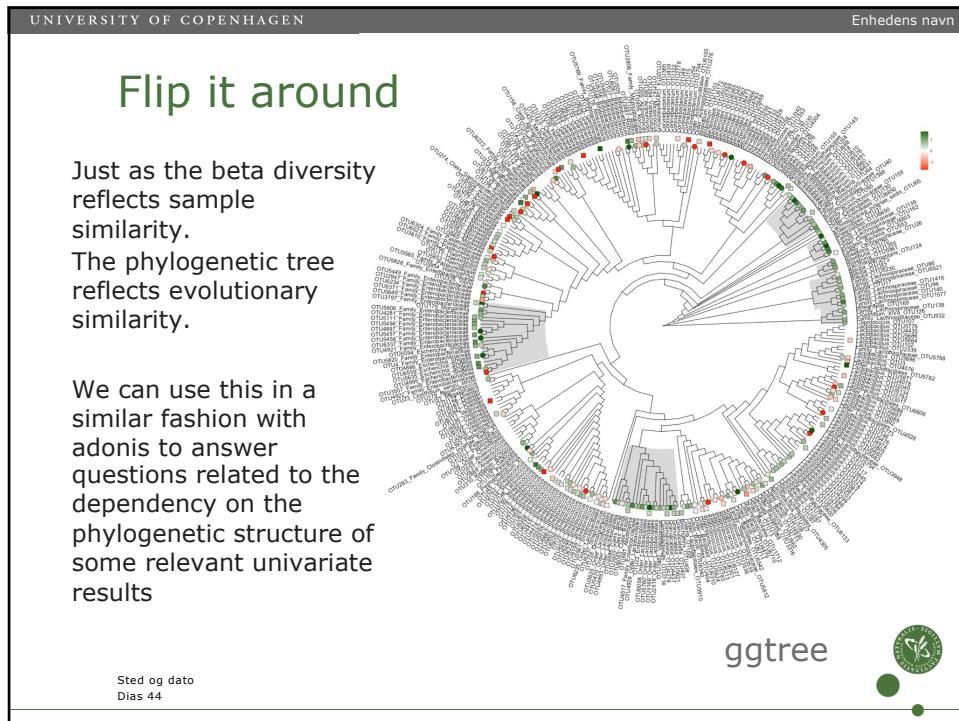
This one is established by **permuting** the design matrix many times, each time calculating the F-statistics.



Sted og dato
Dias 43



43



44



45

Data integration

-Omics

OTU table

$$n^X p_X$$

$$n^Y p_Y$$

46

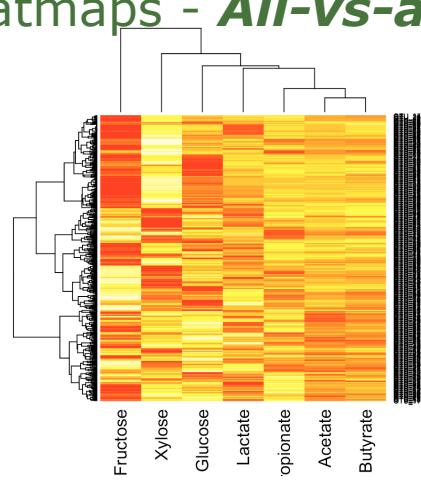


46

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



47

Data integration

-Omics

OTU table

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

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

48



48

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

49



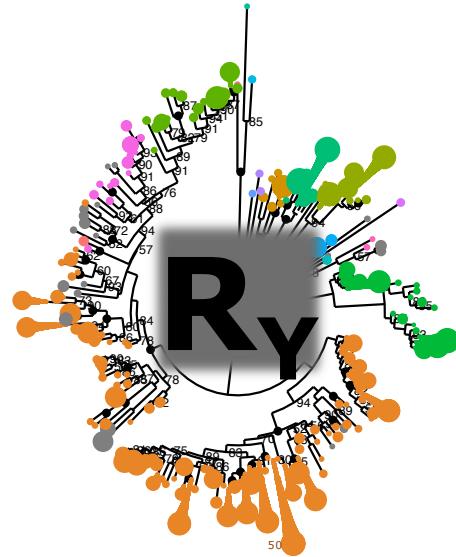
49

24

Functionally Supervised

Supervised CCA

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



50

Functionally Supervised

Supervised CCA

Re-formalize the main objective

Kernel Smoothing

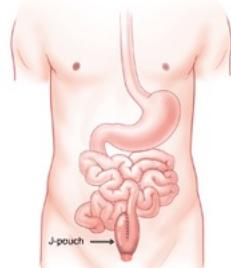
$$\max_{\mathbf{w}_X, \mathbf{w}_Y} \mathbf{w}_X^T \mathbf{X}^T \mathbf{Y} \mathbf{R}_Y \mathbf{w}_Y$$

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



51

Example Pouchitic cohort

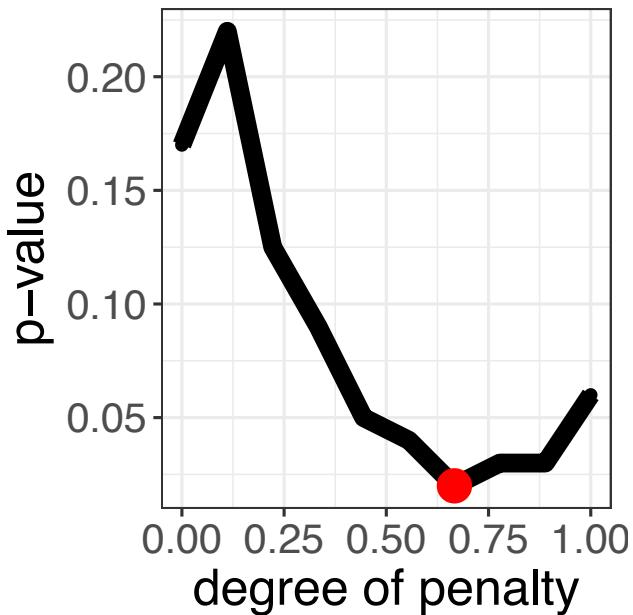


Gene
Expression

OTU table

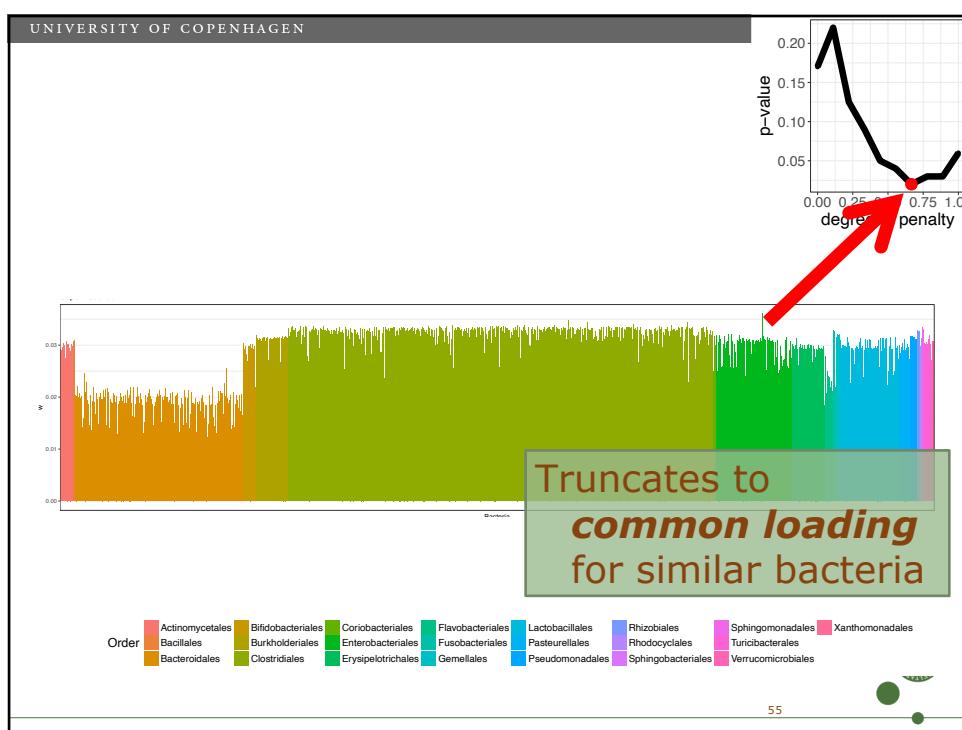
248 X¹⁹⁹¹ 248 Y¹³⁸⁰

52



53

53



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.

Sted og dato
Dias 56

