

1

UNIVERSITY OF COPENHAGEN

Enhedens navn

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

Sted og dato
Dias 2



2

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



Outline

Day1 Morning

Preprocessing
Alpha diversity
Beta diversity

Day2 Morning

Testing design
versus Beta diversity

Day1 Afternoon

Differential
Abundance testing

Day2 Afternoon

Multiomics with
heatmaps and CCA

Sted og dato
Dias 4



Preprocessing

Problem

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

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

Sted og dato
Dias 5



5

Diversity metrics

Alpha diversity

Within sample characteristics

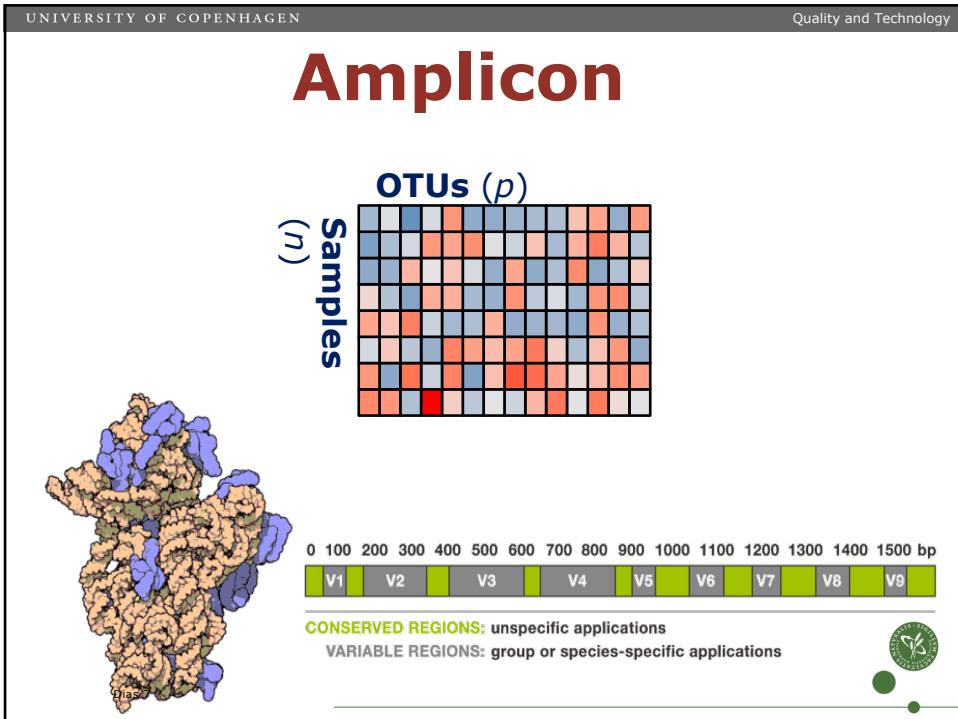
Beta diversity

Between sample characteristics

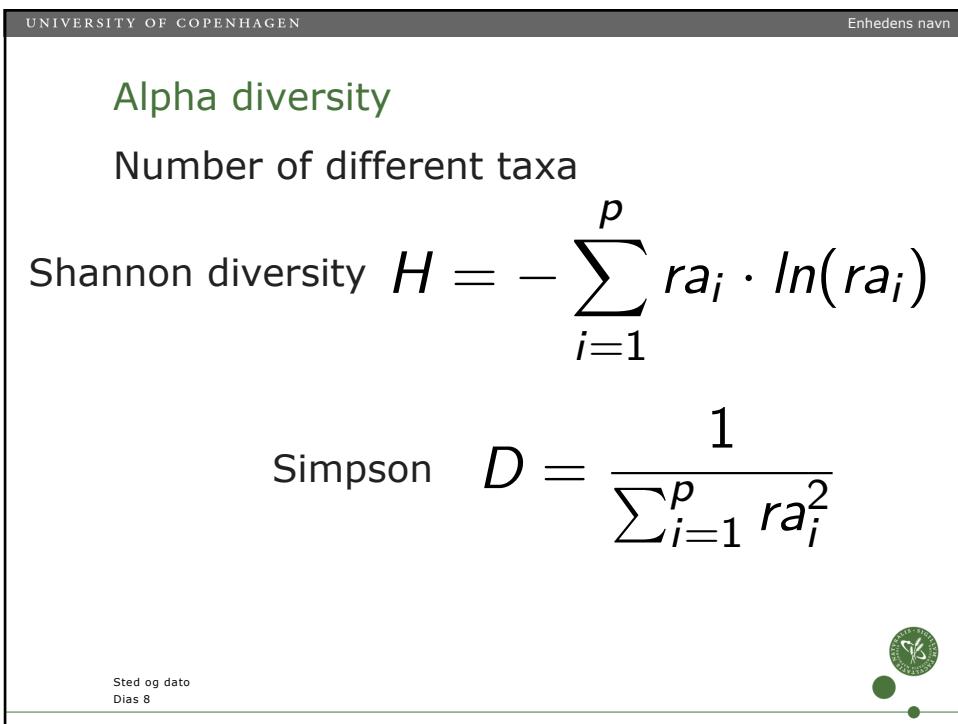
Sted og dato
Dias 6



6



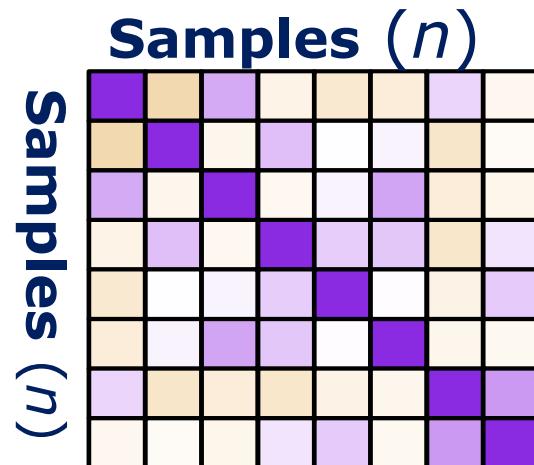
7



8

$\text{kinda}' \sim \mathbf{X}\mathbf{X}^T$

β diversity



9

9

Enhedens navn

β diversity

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

Sted og dato
Dias 10

10

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 11



11

Bray Curtis

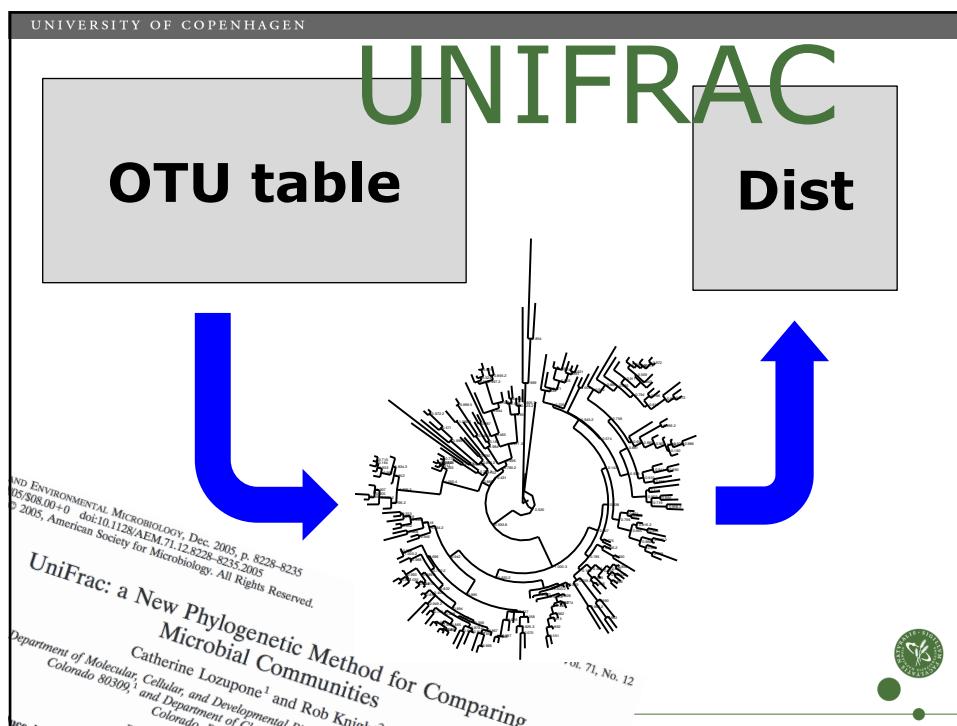
$$BC = \frac{\sum_{i=1}^n |X_{ij} - X_{ik}|}{\sum_{i=1}^n (X_{ij} + X_{ik})}$$

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 12



12



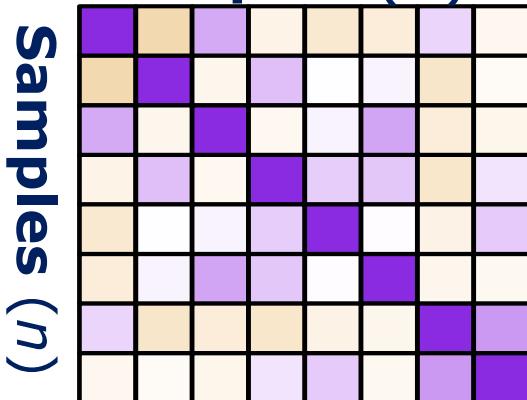
13



14

β diversity to PCoA

Samples (n)



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

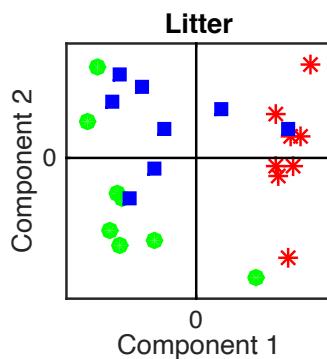
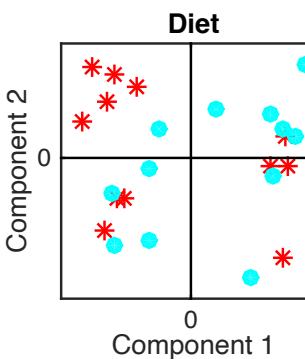


15

15

Multidimensional scaling

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



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

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



16

Differential Abundance Testing or OTUWAS

Dias 17



17

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.

Enhedens navn
Sted og dato
Dias 18

OTUs (p)

Samples (n)

Class (n)

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



UNIVERSITY OF COPENHAGEN Enhedens navn

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)

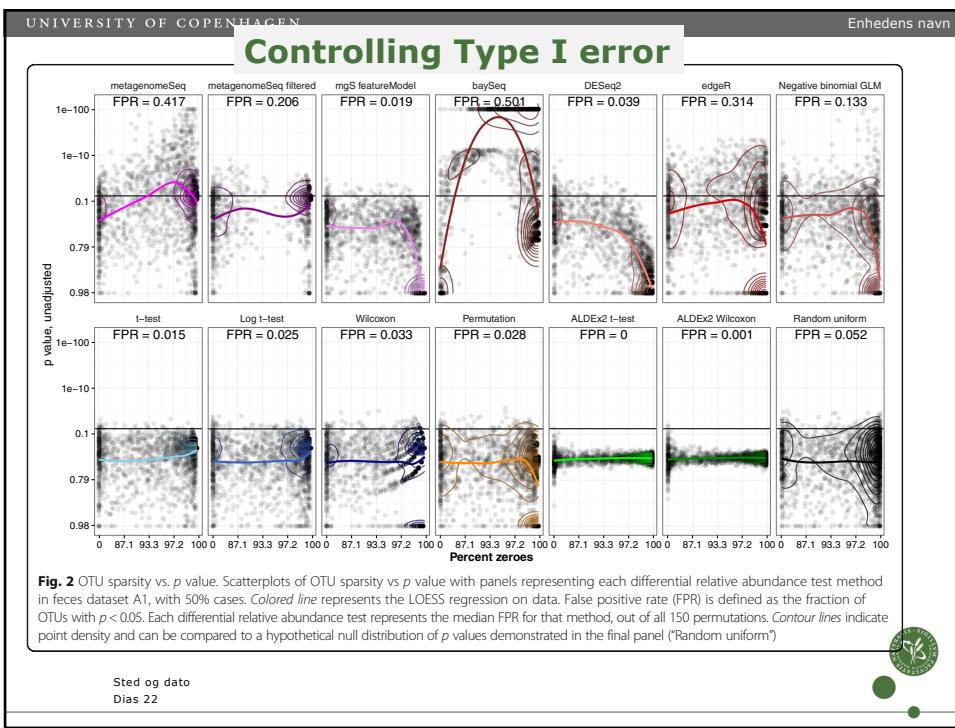
Uses **empirical bayes** on central- and dispersion parameters to shrink towards common value

(uses cumulative sum scaling for prepro)

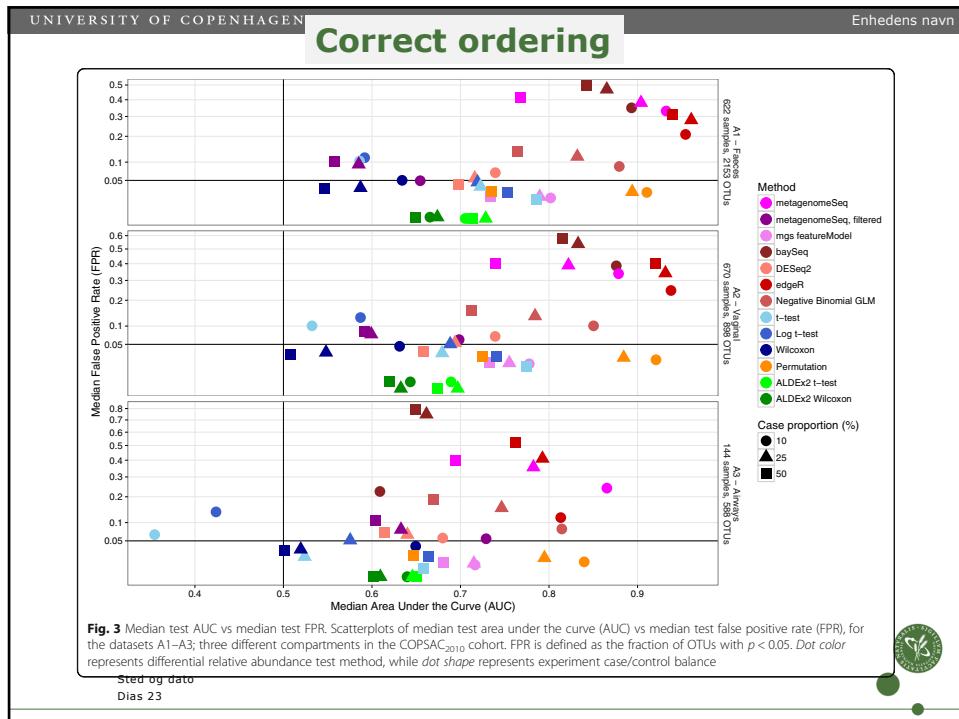
$$f_{sig}(y_{ij}; s_j, \beta, \mu_i, \sigma_i^2) = \pi_j(s_j) \cdot I_{\{0\}}(y_{ij}) + (1 - \pi_j(s_j)) \cdot f_{count}(y_{ij}; \mu_i, \sigma_i^2)$$

Sted og dato
Dias 21

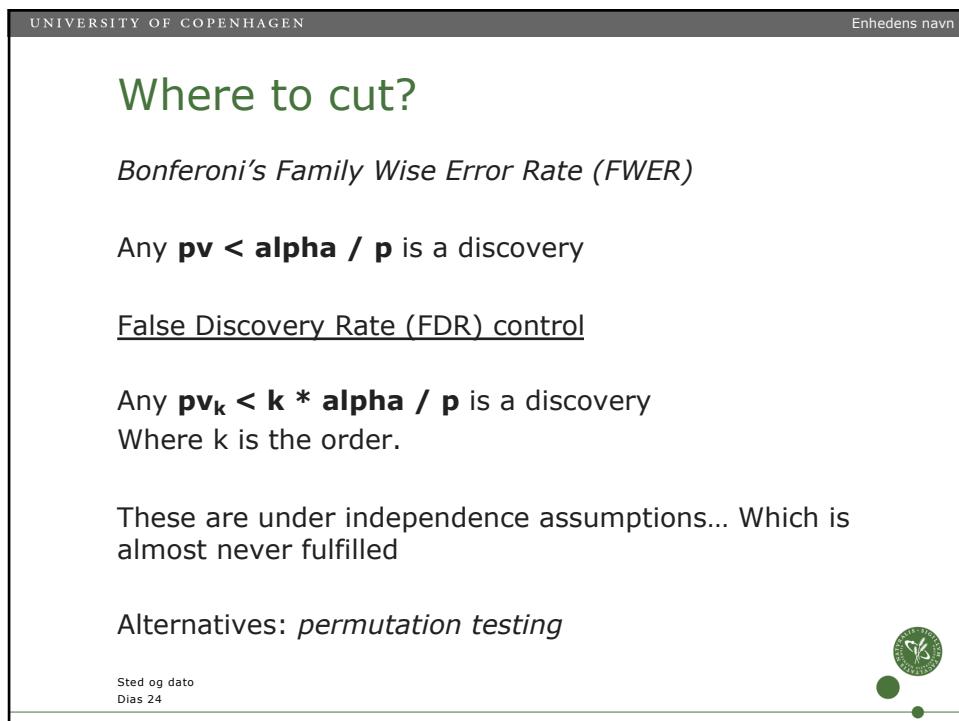
21



22



23



24

Testing Beta diversity *PermANOVA*

Dias 25



25

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.

Dias 26

26

Establishing the *null* distribution

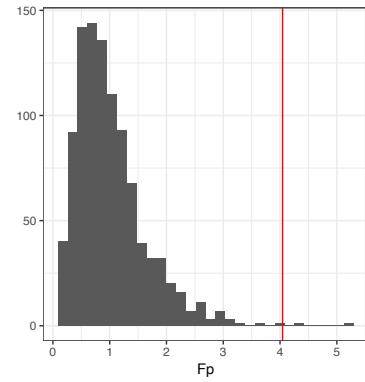
For univariate ANOVA models, we rely on independence and homoschedastic and gaussian *like* residuals.

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 27



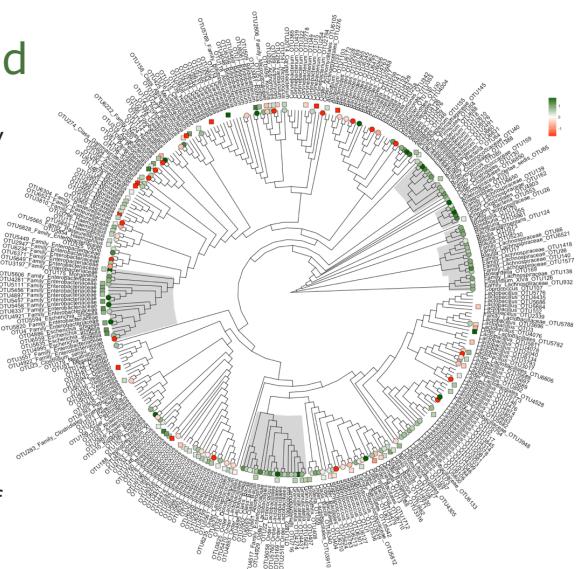
27

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



Sted og dato
Dias 28

28

MultiOmics

Dias 29



29

Data integration

-Omics

OTU table

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

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

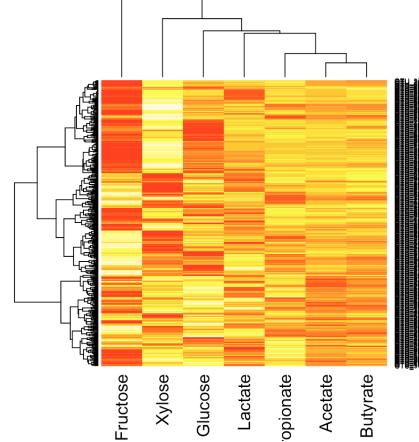
30



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



31

Data integration

-Omics

OTU table

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

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

32



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

33

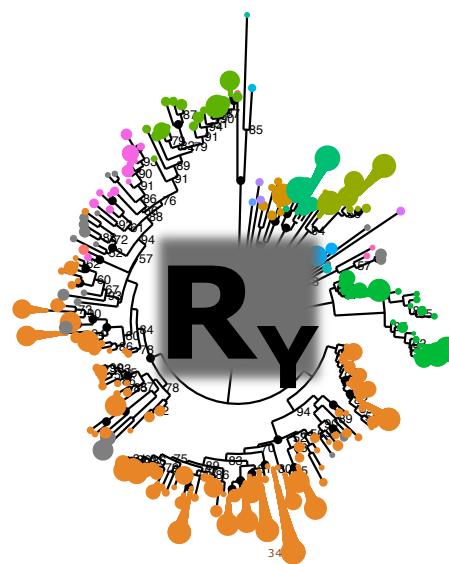


33

Functionally Supervised

Supervised CCA

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



34

Functionally Supervised

Supervised CCA

Re-formalize the main objective

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

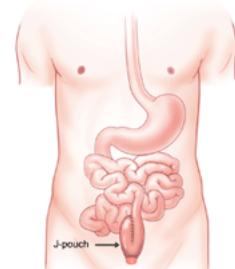
$$\mathbf{R}_Y = \alpha \mathbf{S} + (1 - \alpha) \mathbf{I}$$

Kernel Smoothing



35

Example Pouchitic cohort

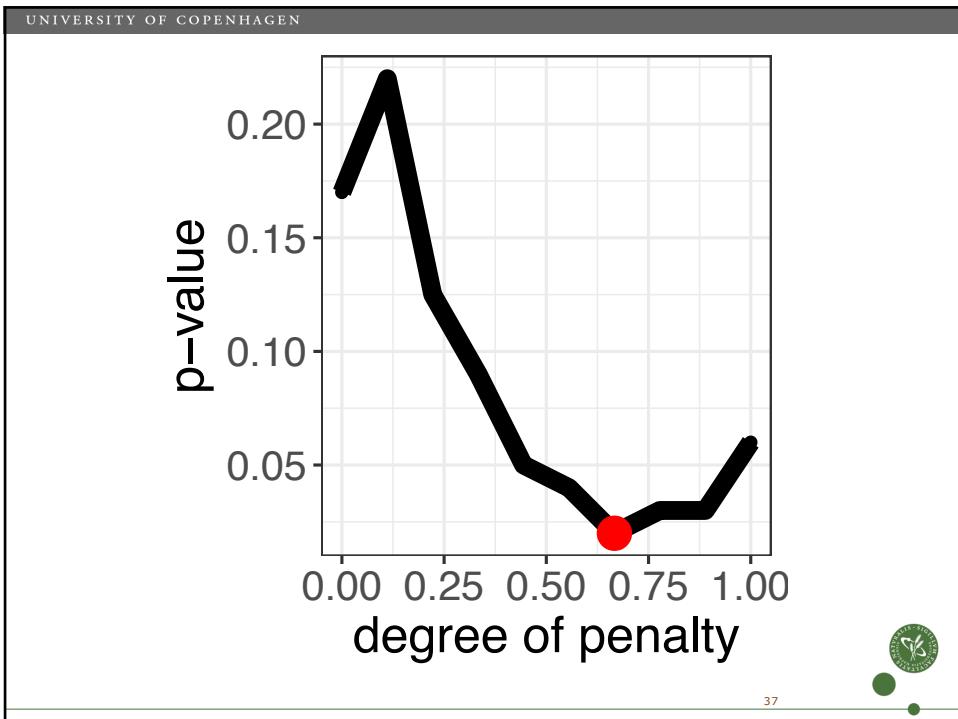


Gene
Expression

OTU table

248 \mathbf{X}^{1991} 248 \mathbf{Y}^{1380}

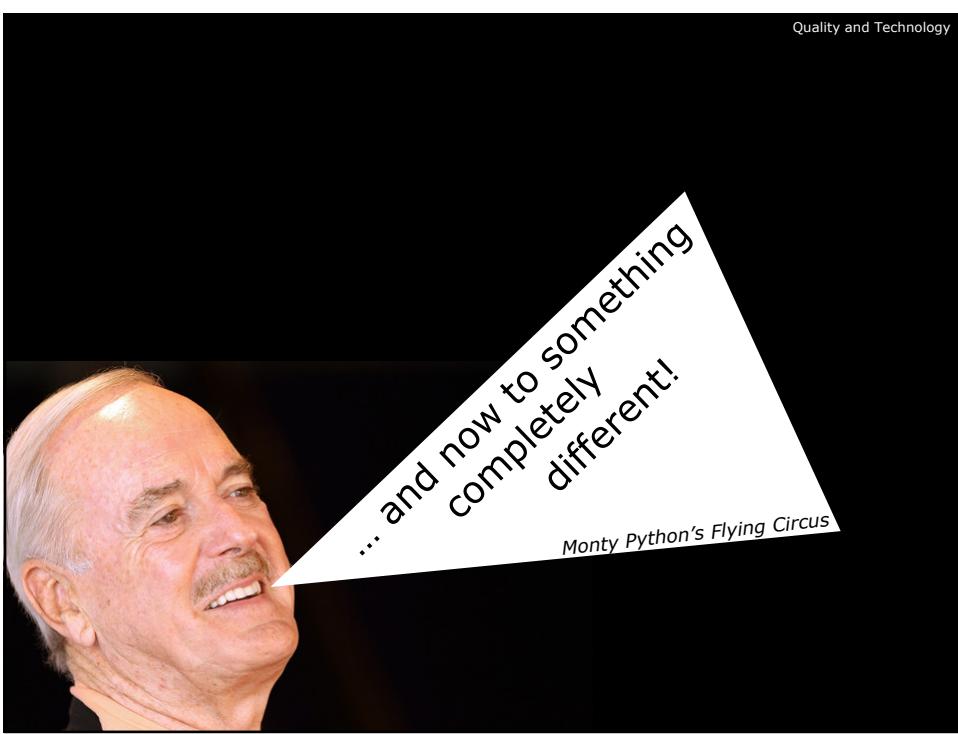
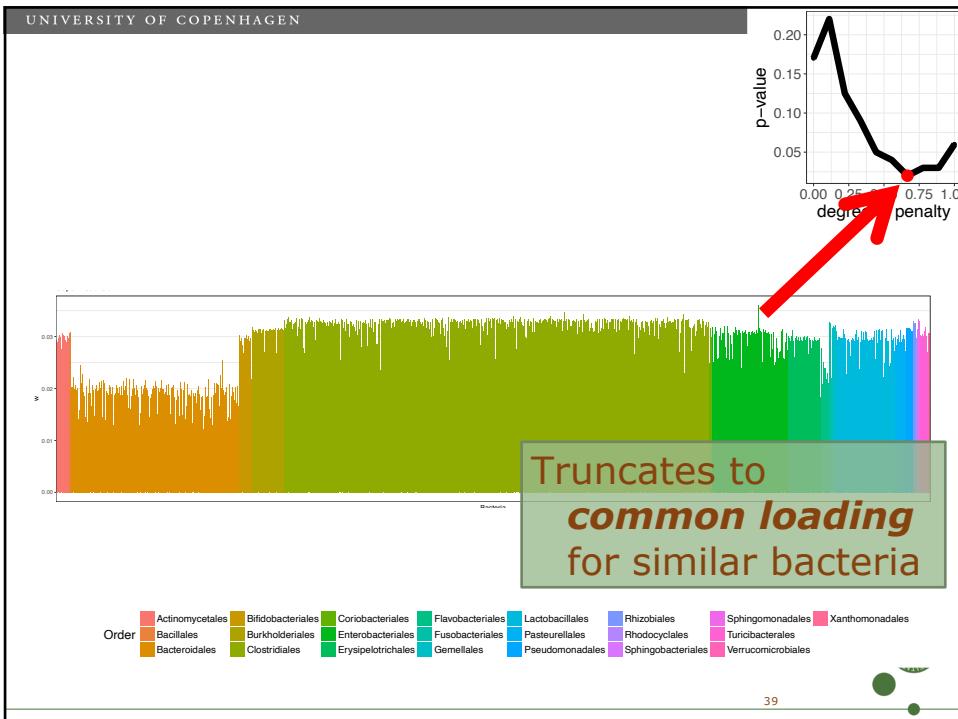
36



37



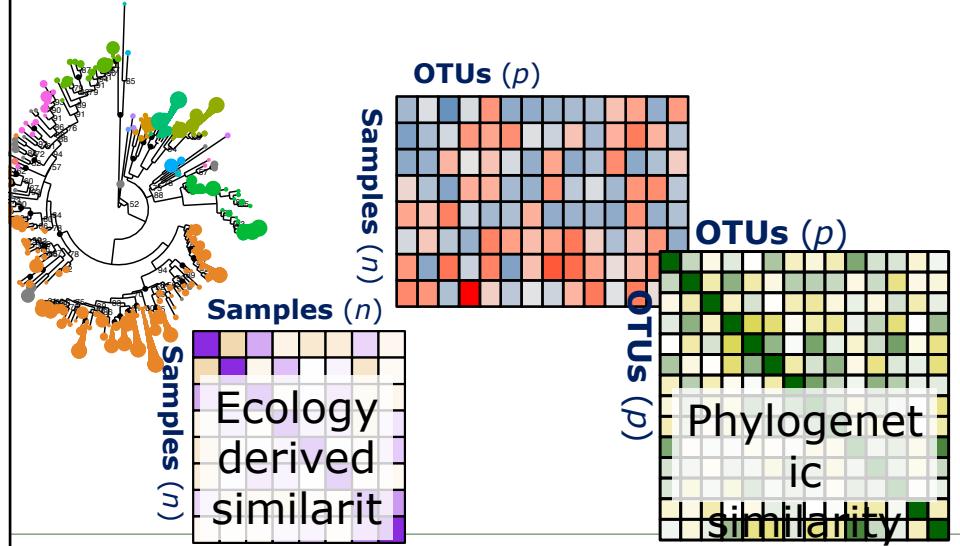
38



40

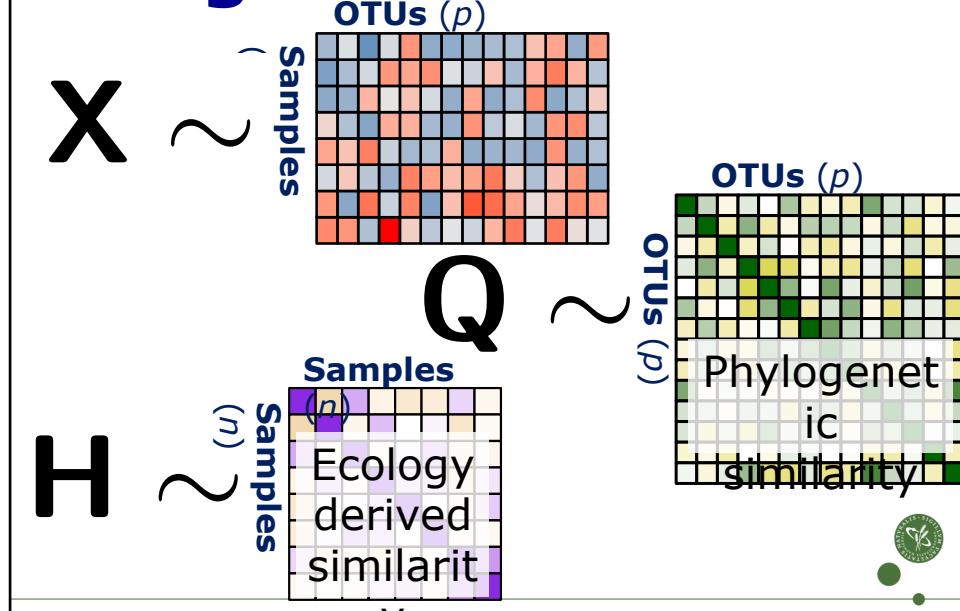
Data analysis

We know *more* about the features than *just* their number/name!



41

Regularized PLS



42

Regularized PLS

Objective

Find a PLS solution such that

1. The covariance between \mathbf{X} and \mathbf{Y} are high
2. The weights resemble known feature2feature structure
3. The scores resemble ecological structure
4. (...including variable selection)



43

Regularized PLS

Objective

1. Covariance btw \mathbf{X} and \mathbf{Y}
2. Resemble feature structure
3. Resemble ecological structure
4. Variable selection

$$\max_{\mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X} \mathbf{w})$$

$$\max_{\mathbf{w}} (\mathbf{w}^T \mathbf{Q} \mathbf{w})$$

$$\max_{\mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{H} \mathbf{X} \mathbf{w})$$

$$\|\mathbf{w}\|_1 \leq c \quad \|\mathbf{w}\|_2 = 1$$

44

Regularized PLS

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} (\mathbf{w}^T (\alpha_1 \mathbf{Q} \mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X} \mathbf{Q} + (1 - \alpha_1) \mathbf{X}^T \mathbf{H} \mathbf{X}) \mathbf{w})$$

$$\hat{\mathbf{w}} \leftarrow \text{softthr}(\hat{\mathbf{w}}, \lambda)$$

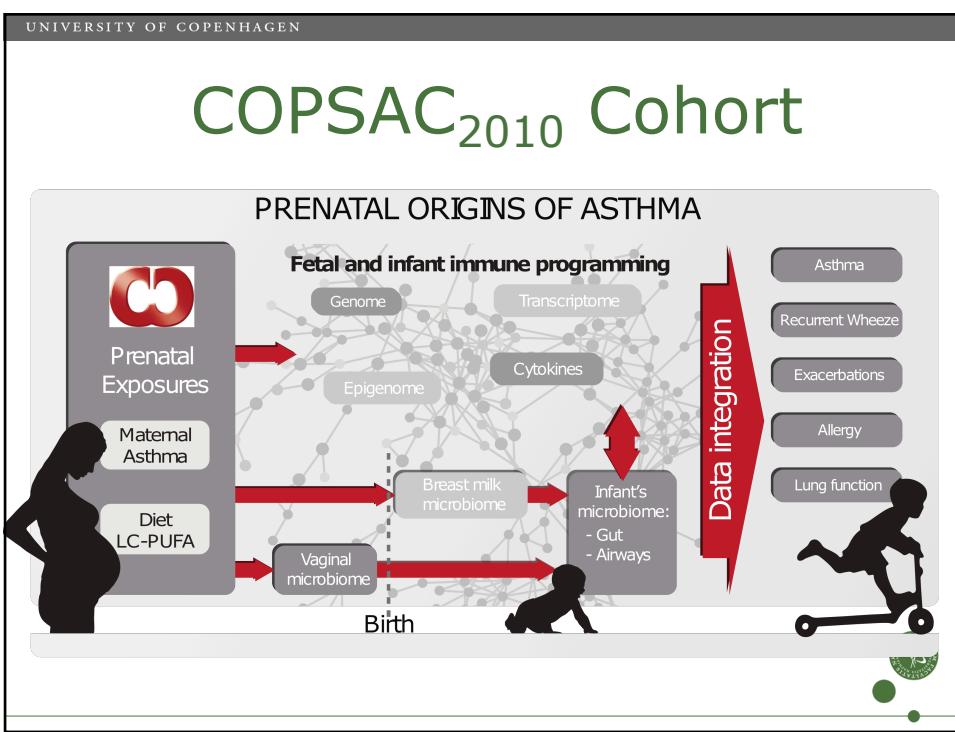
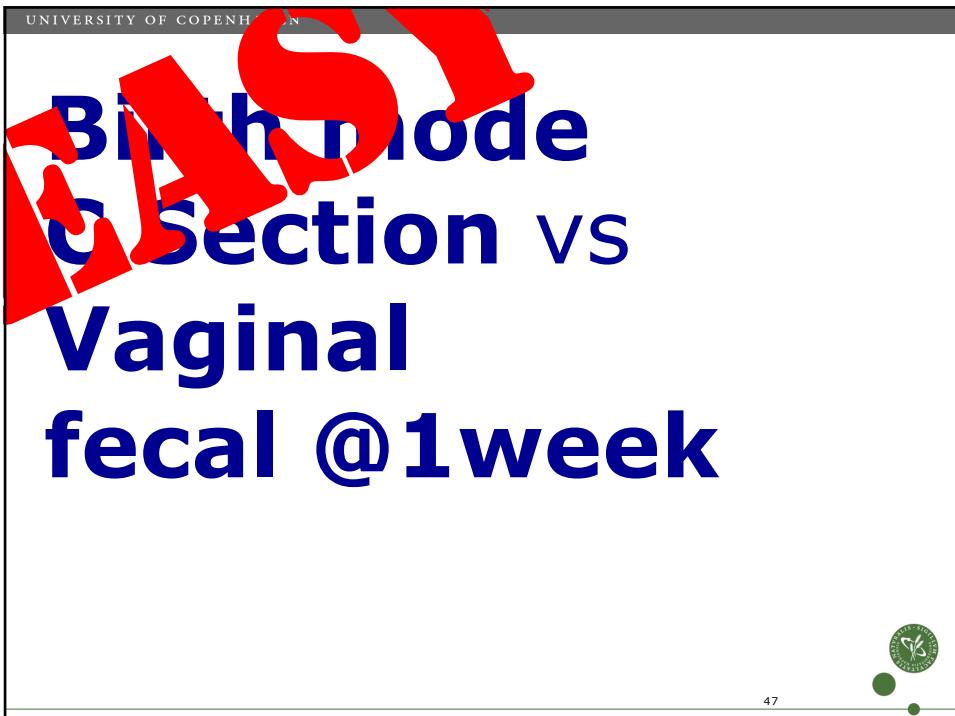
$$\mathbf{Q} = \alpha_2 \mathbf{Q}^* + (1 - \alpha_2) \mathbf{I}$$

$$\|\mathbf{w}\|_2 = 1$$

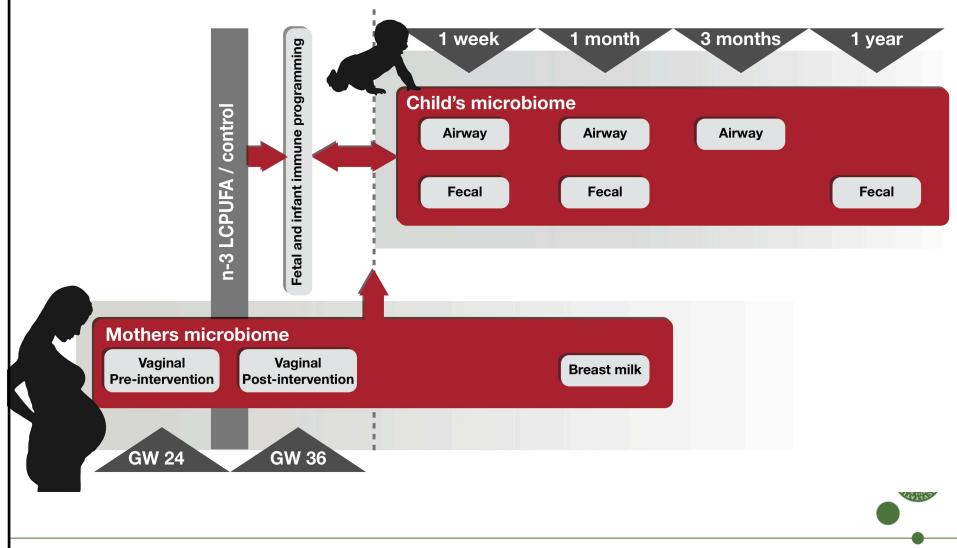
45



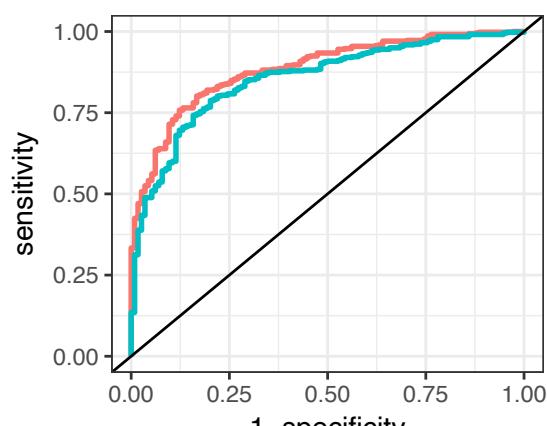
46



Microbiome data



49

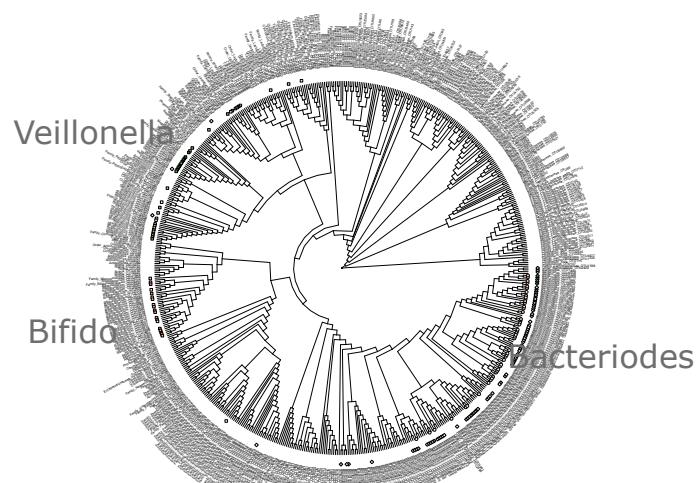


model — Cal — Val

50

50

Loadingplot



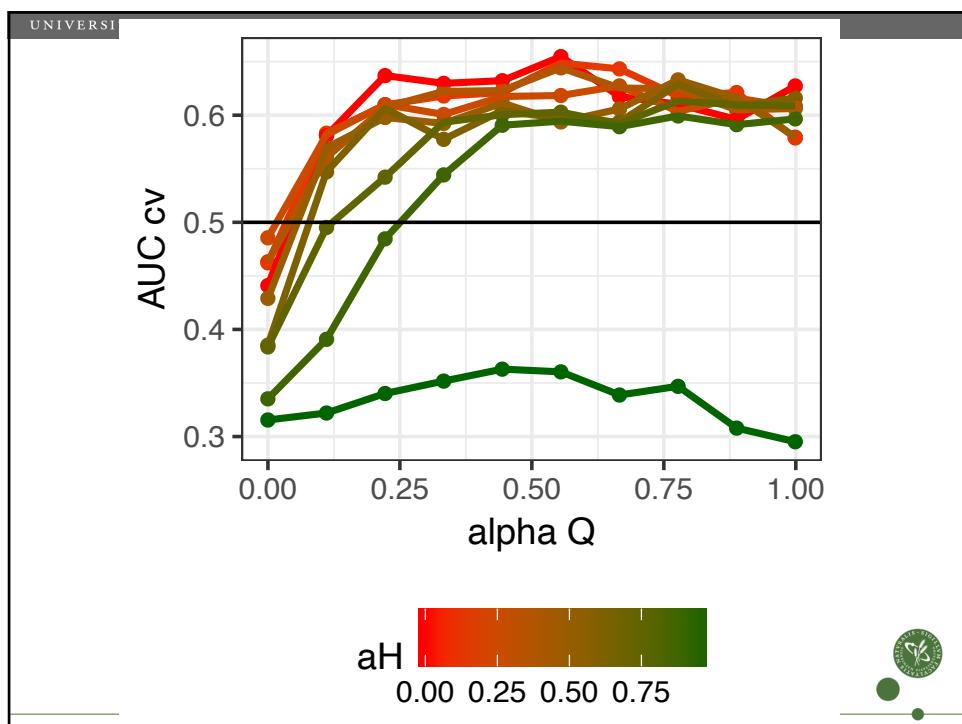
51

51

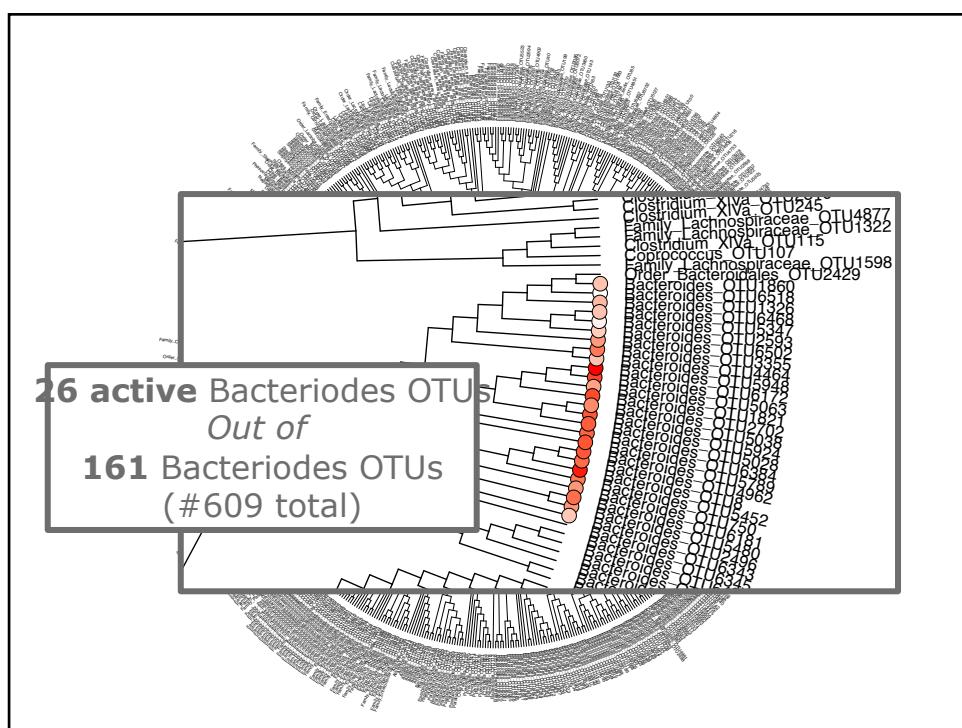
**asthma 6y
fecal @1year**
- in asthmatic mothers-

52

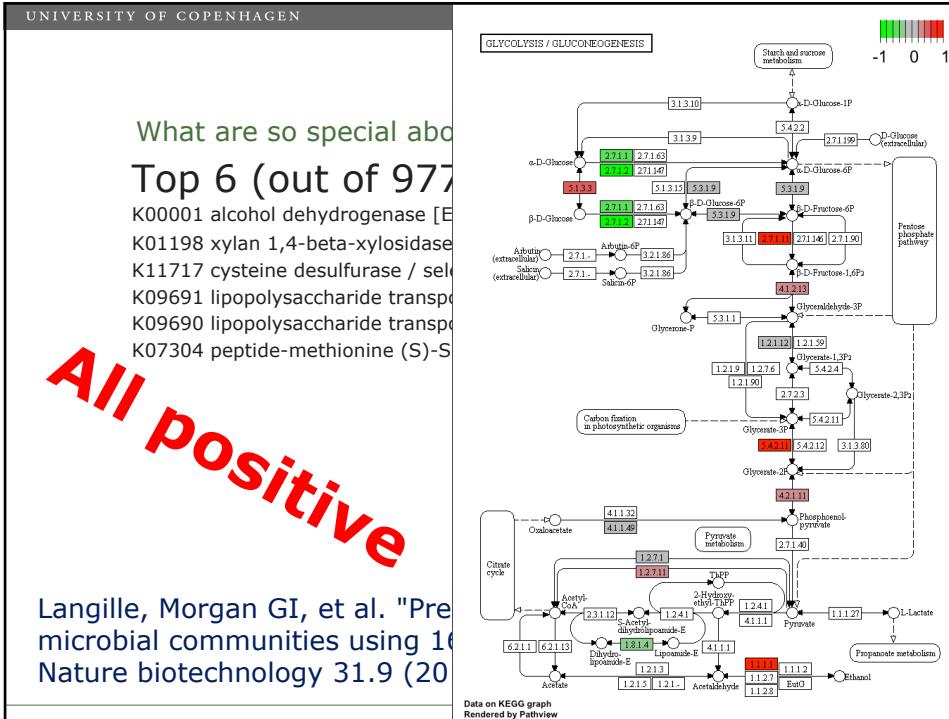
52



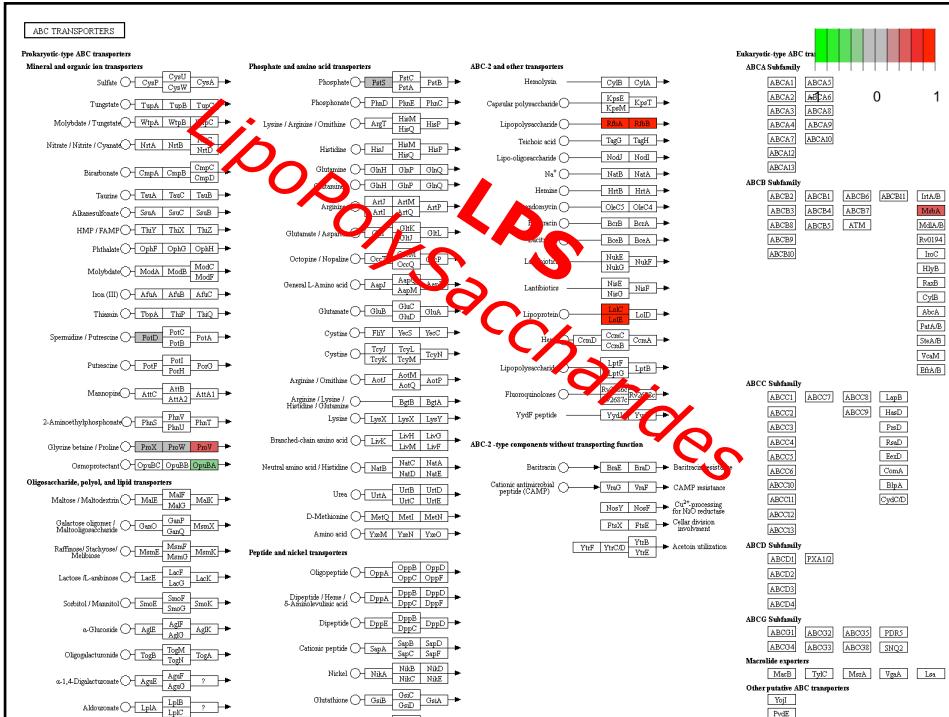
53



54



55



56