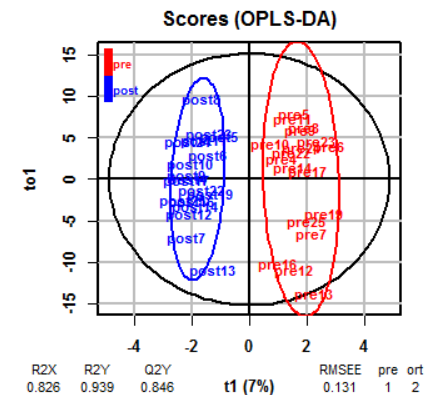
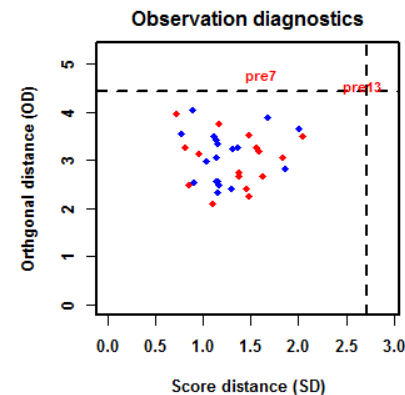
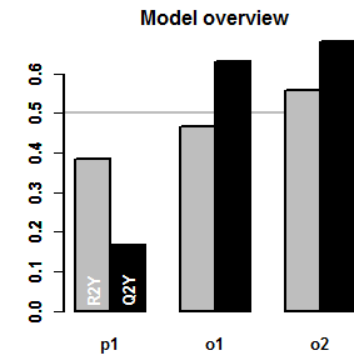
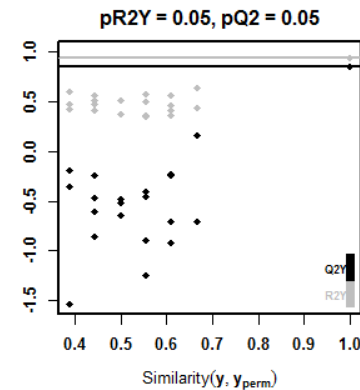
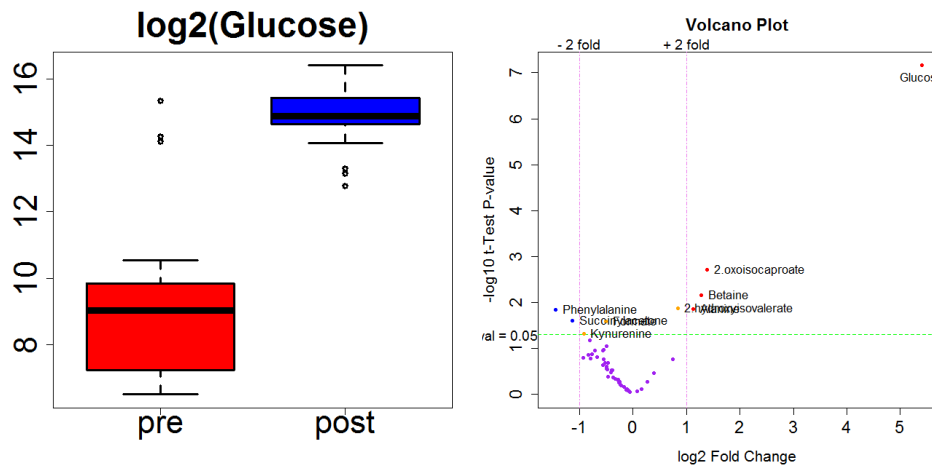




# Vergleich einer univariaten und einer multivariaten «OMICS» - Datenanalyse; zB. für Genomics, Transcriptomics, Proteomics, *Metabolomics*, Foodomics, ...

Dominik Guggisberg, Agroscope

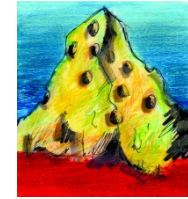


CAS/DAS Angewandte Statistik,

04.05.2018

www.agroscope.ch | gutes Essen, gesunde Umwelt

# Agenda



## 1. Einleitung / Ziele

## 2. Univariat; traditionelle Tests, parametrisch, nicht-parametrisch, FDR (false discovery rate) Biologische Effekte («Volcano Plot» $\cong$ Scatterplot)

R: library(metabolomics) ([cran.r-project.org/web/packages](http://cran.r-project.org/web/packages))

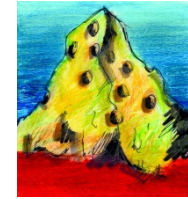
## 3. Multivariat; PCA, PLS-DA, OPLS-DA

«Feature selection», VIP (Variable Importance in Projection)

R: library(ropls) ([bioconductor.org](http://bioconductor.org))

## 4. Zusammenfassung/Ausblick

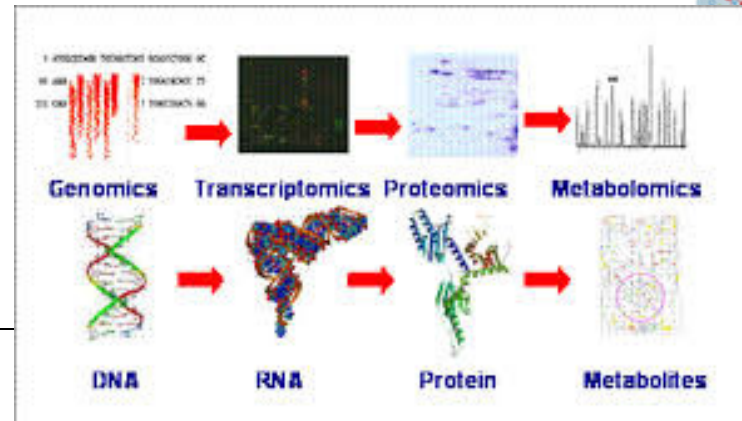
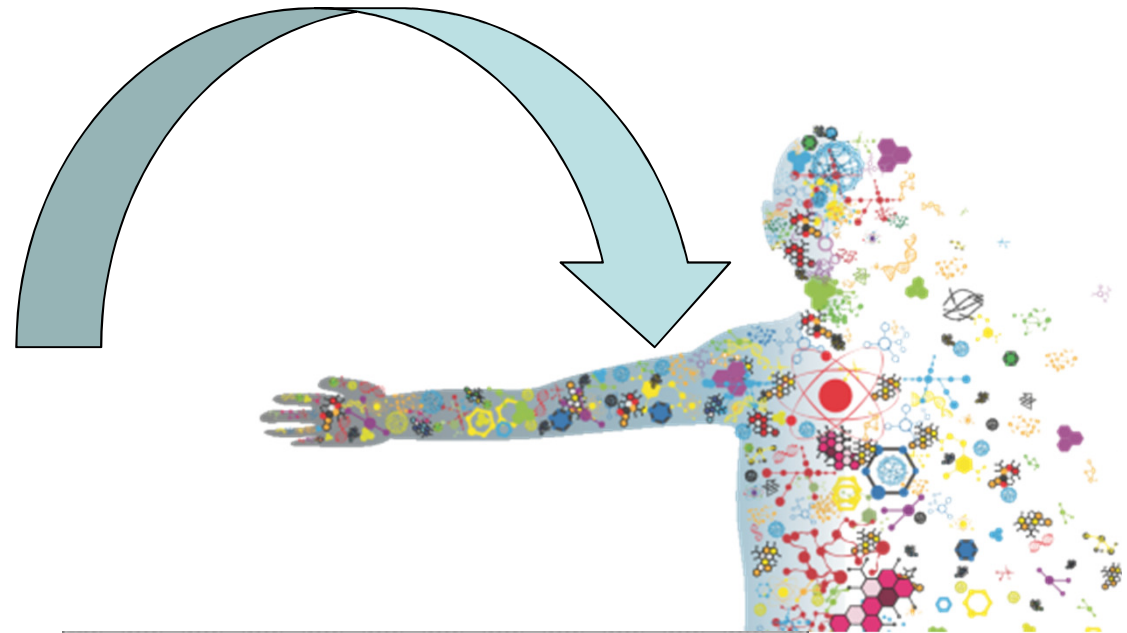
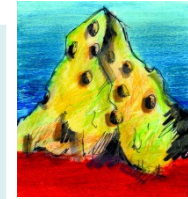
# 1. Einleitung / Ziele:



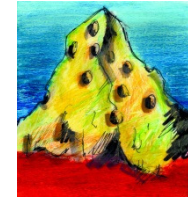
Die Erforschung des Metaboloms wird als **Metabolomik** bezeichnet (im englischen: *metabolomics*). Diese umfasst die Wechselwirkung der darin enthaltenen Metaboliten, deren Identifizierung und Quantifizierung. (=> Einfluss auf den Stoffwechsel verstehen.)



# zB. Metabolomics



# 🇨🇭 Datensatz: library(metabolomics)\*; data(treated)



▪ S Group		met1	met2	met3	...	metN
▪ S1	A	0.6358	0.0851	0.3665	...	1.0024
▪ S2	A	0.5871	0.0935	0.3421	...	1.0329
▪ .....	...	....	....	....	....	....
▪ S19	B	0.6650	1.0705	0.6710	...	0.7319
▪ S20	B	0.6907	1.0341	0.6858	...	0.7376 ...

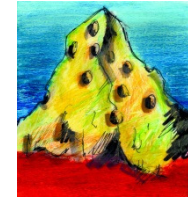
A metabolomics study with **paired** observations.

**Description:** A data frame with data collected from several subjects **before** and **after** a specific treatment. (= > Diabetes-Forschung, 2012)

```
n<- ncol(treated) # n= 54 (53 Metaboliten)
```

```
m<- nrow(treated) # m= 36 (= 2x18 Subjekte)
```

## 2. Univariate Analyse

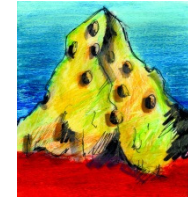


- Parametrisch: T-Test, ANOVA
- Nicht-parametrisch: Wilcoxon-Test (Mann-Whitney), Kruskal-Wallis
- Probleme: “Alpha-Fehler - Kumulierung” (=> multiples Testen)
- => False Discovery Rate (FDR): Korrektur mit Benjamini-Hochberg<sup>1</sup>
- => Filterung der Daten, damit weniger “Features” vorhanden sind: Nachteil; ev. “verlieren” wir wichtige “Features”?
- Grafik: “Volcano Plot”

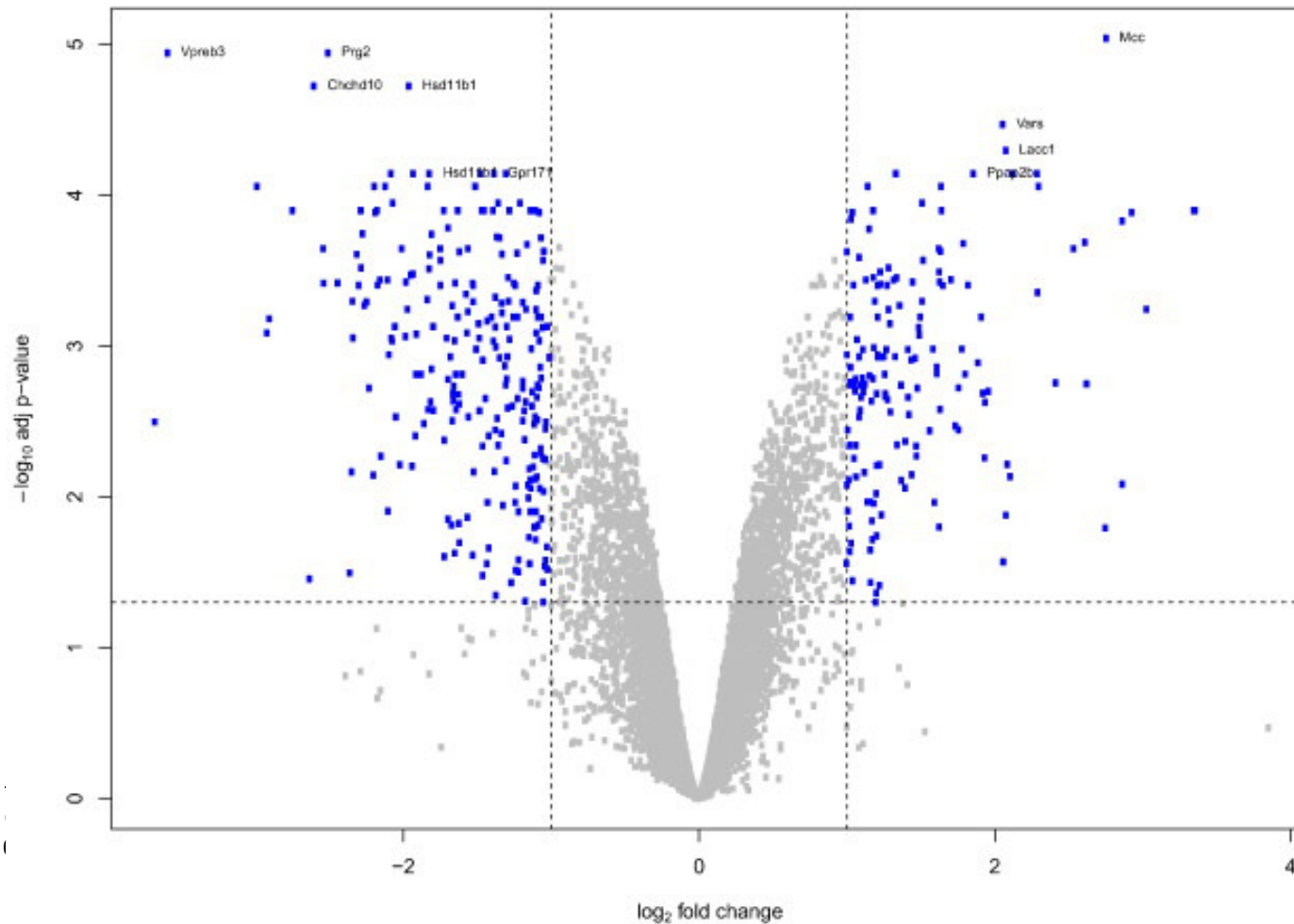
Der Begriff wurde erstmals 1995 von [Yoav Benjamini](#) und [Yosi Hochberg](#) definiert.<sup>[1]</sup>

Benjamini, Yoav; Hochberg, Yosef: ["Controlling the false discovery rate: a practical and powerful approach to multiple testing"](#) In: Journal of the Royal Statistical Society, Series B Nr. 57, 1995, S. 289–300

# Volcano Plot (Typ: Scatterplot)

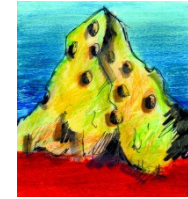


- x-Achse:  $\log_2$  fold change =  $\log_2$  (mean post/mean pre)
- y-Achse:  $-\log_{10}$  (p-values)





# Volcano plot für data(treated) (library(metabolomics))



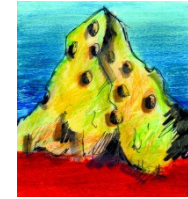
- #### Volcano plot
- data(treated)
- treated.log <- LogTransform(treated, base = 2)\$output
- results <- TwoGroup(treated.log, **paired = TRUE**)\$output
- (Description TwoGroup: This function computes t- statistics, p-values, adjusted p-values, fold changes and standard errors for each metabolite given a series of replicates and two biological conditions).

```
> results
```

	t-statistic	p-value	BH	Adjusted p-value	fold change	standard error
2.hydroxybutyrate	-1.8408367	8.316721e-02		1.377457e-01	-0.28430458	0.1544431
2.hydroxyisobutyrate	0.3673505	7.178907e-01		7.460433e-01	0.07671834	0.2088423
2.hydroxyisovalerate	4.1505989	6.693494e-04		5.912587e-03	0.83338203	0.2007860
2.hydroxyvalerate	-2.5875902	1.916541e-02		4.416378e-02	-0.40917314	0.1581290
2.methylglutarate	-1.7963295	9.023429e-02		1.449217e-01	-0.33423740	0.1860669
2.oxoisocaproate	8.0900778	3.132787e-07		8.301887e-06	1.38511010	0.1712110



# Volcano plot für data(treated) (library(metabolomics))



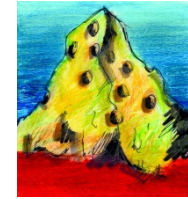
- #### Volcano plot
- data(treated)
- treated.log <- LogTransform(treated, base = 2)\$output
- results <- TwoGroup(treated.log, **paired = TRUE**)\$output
- (Description TwoGroup: This function computes t- statistics, p-values, adjusted p-values, fold changes and standard errors for each metabolite given a series of replicates and two biological conditions).

```
> results
```

	t-statistic	p-value	BH Adjusted p-value	fold change	standard error
2.hydroxybutyrate	-1.8408367	8.316721e-02	1.377457e-01	-0.28430458	0.1544431
2.hydroxyisobutyrate	0.3673505	7.178907e-01	7.460433e-01	0.07671834	0.2088423
2.hydroxyisovalerate	4.1505989	6.693494e-04	5.912587e-03	0.83338203	0.2007860
2.hydroxyvalerate	-2.5875902	1.916541e-02	4.416378e-02	-0.40917314	0.1581290
2.methylglutarate	-1.7963295	9.023429e-02	1.449217e-01	-0.33423740	0.1860669
2.oxoisocaproate	8.0900778	3.132787e-07	8.301887e-06	1.38511010	0.1712110



### 3. Multivariate Analyse:



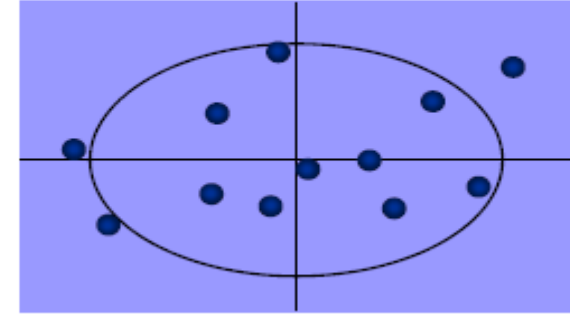
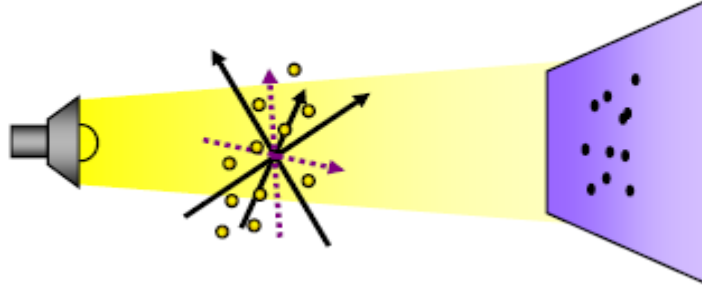
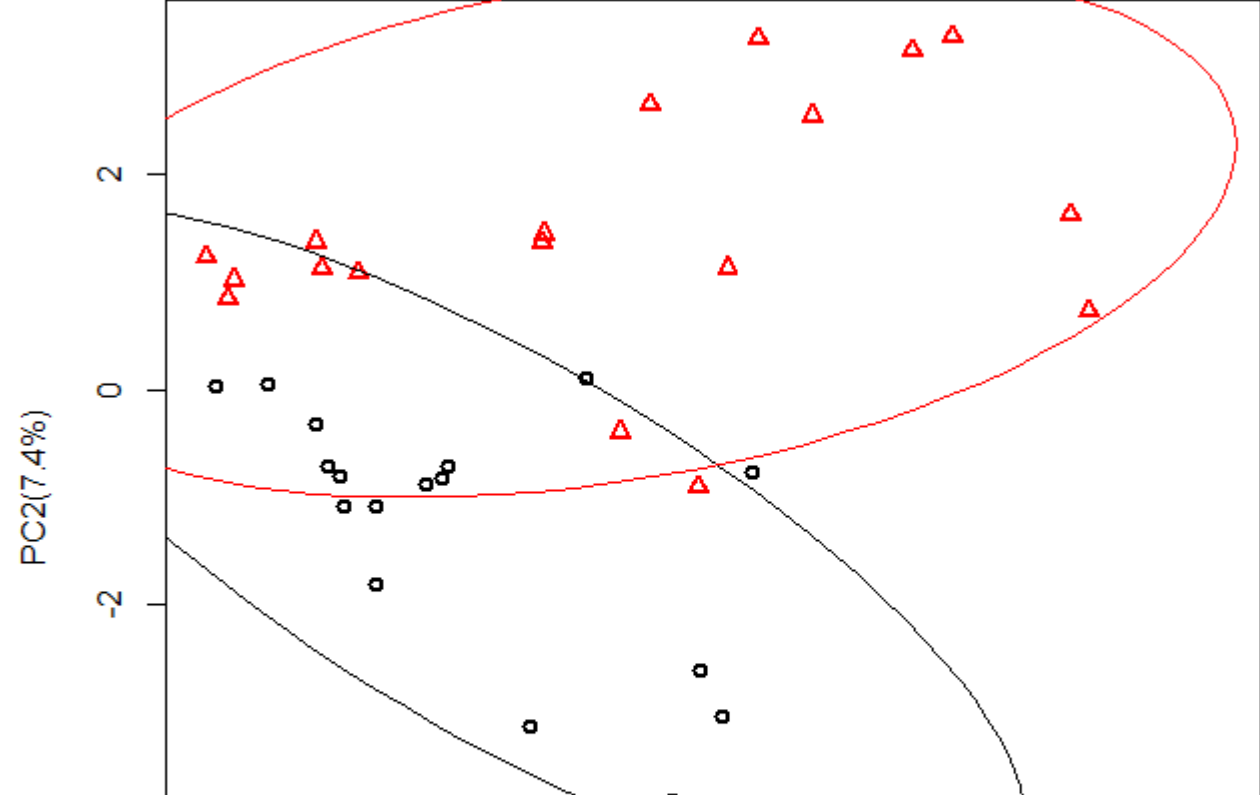
- The ropls R package implements the **PCA**, **PLS(-DA)** and **OPLS(-DA)** approaches with the original, **NIPALS**-based, versions of the algorithms. It includes the **R2** and **Q2** quality metrics, the permutation **diagnostics**, the computation of the **VIP** values, the score and orthogonal distances to detect **outliers**, as well as many **graphics** (scores, loadings, predictions, diagnostics, outliers, etc).



# PCA

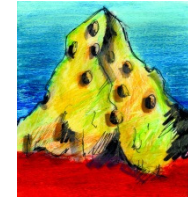
center=T,  
scale=T,  
(nicht-  
log)

## Hauptkomponentenanalyse



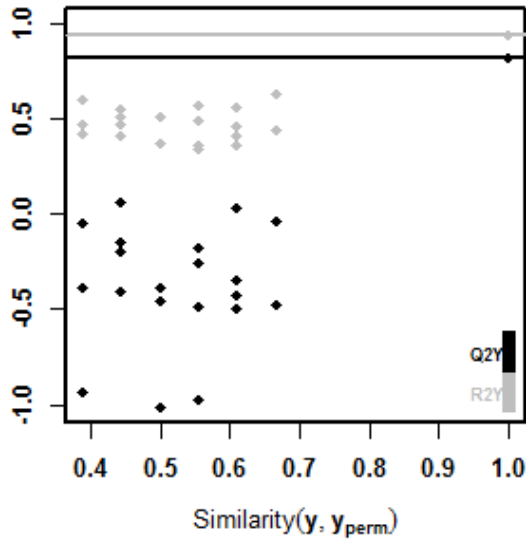


# PLS-DA (log-transformierte Daten)

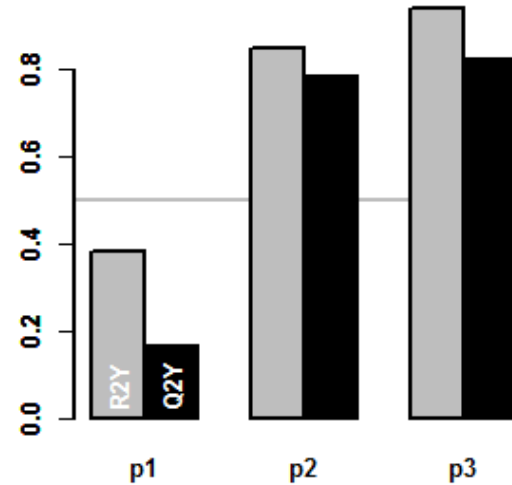


Latent variable regression based on covariance between X (predictors) and Y (response). (ideal for multi-collinear predictors!)

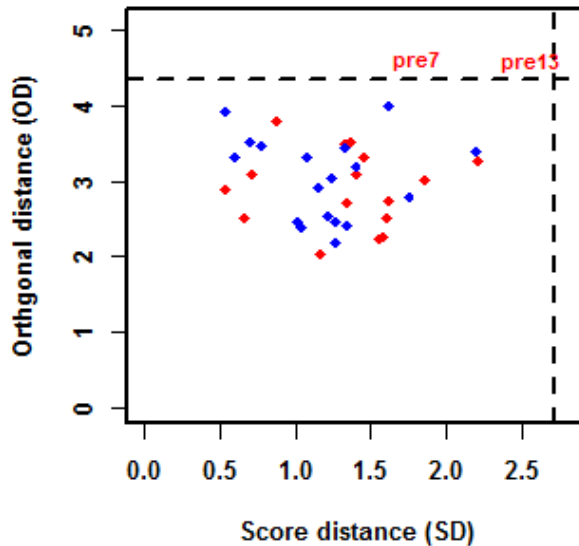
pR2Y = 0.05, pQ2 = 0.05



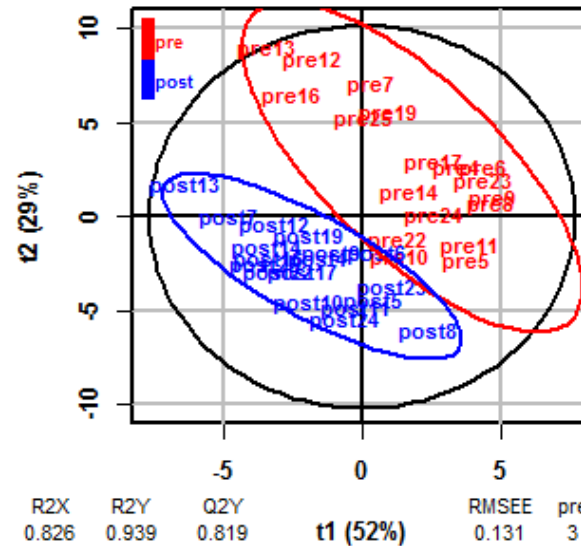
Model overview



Observation diagnostics

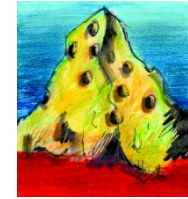


Scores (PLS-DA)

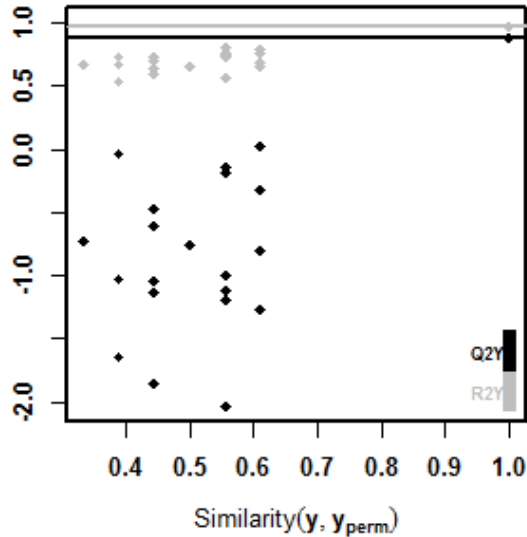




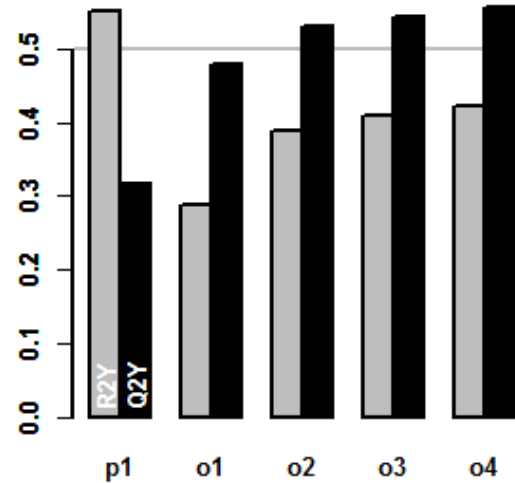
# OPLS-DA (log-transformierte Daten)



pR2Y = 0.05, pQ2 = 0.05

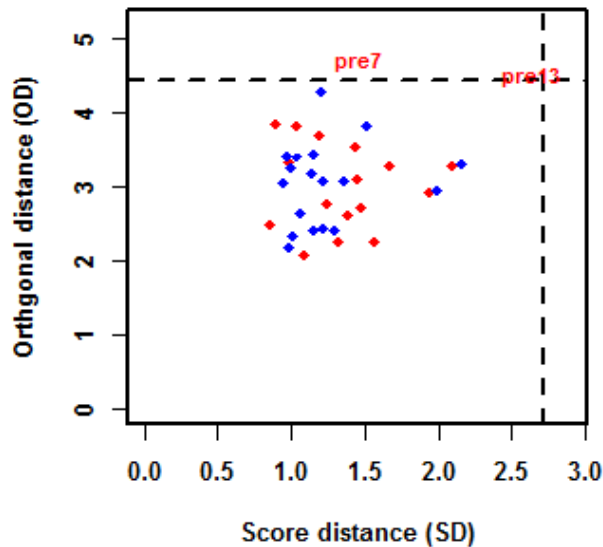


Model overview

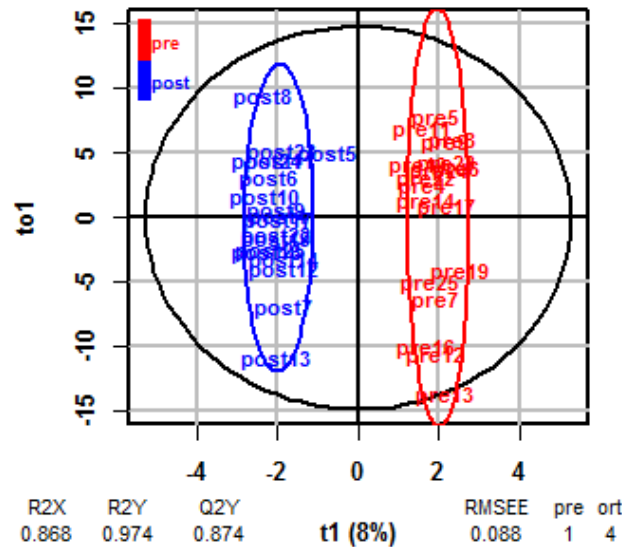


Model separates the variations of X (predictors) correlated and orthogonal to Y (response). Q2Y metrics and permutation testing to avoid overfitting. VIP (loading weights and variability of the response explained by this component: feature selection!)

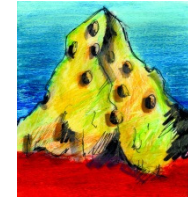
Observation diagnostics



Scores (OPLS-DA)



# VIP (Variable Importance in Projection)



▪ head(newdata1, 10), log-transformierte Daten

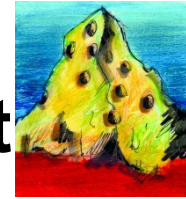
▪	<b>Glucose</b>	<b>2-oxoisocaproate</b>	<b>2-hydroxyisovalerate</b>	<b>Betaine</b>
▪	3.178962	1.906171	1.664372	1.657276
▪	<b>Alanine</b>	<b>Phenylalanine</b>	<b>Formate</b>	<b>Succinylacetone</b>
▪	1.657078	1.603227	1.525911	1.499841
▪	<b>Kynurenine</b>	<b>Tryptophan</b>		
▪	1.321389	1.301598		

(sortiert!)





## 4. Zusammenfassung: Univariat/Multivariat



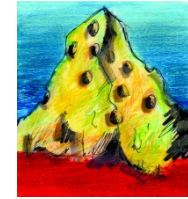
### ■ Univariat (Volcano-Plot)

- Glucose
- Oxoisocaproate
- Alanine
- Betaine
  
- Succinylacetone
- Phenylalanine

### ■ Multivariat (OPLS-DA, log-Daten)

- Glucose
- Oxoisocaproate
- **Hydroxyisovalerate**
- Betaine
- Alanine
- Phenylalanine
- **Formate**
- Succinylacetone
- ...

# Skalierung: Unit variance oder Pareto



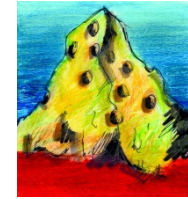
- Unit variance scaling: jede Variable wird durch die Standardabweichung dividiert.
- Pareto scaling: jede Variable wird durch die Wurzel der Standardabweichung dividiert (häufig in Metabolomics-Studien verwendet).

- In der Literatur wird zB. auch empfohlen:

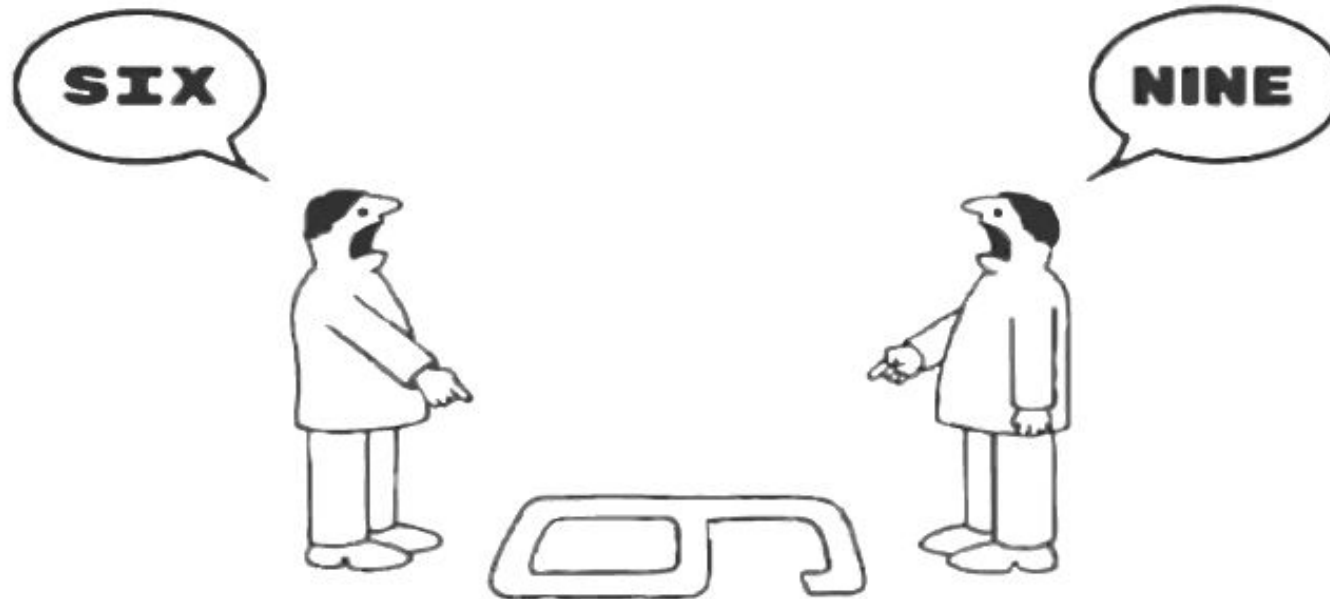
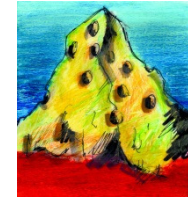
$$\text{Range scaling: } \tilde{x}_{ij} = \frac{x_{ij} - \bar{x}_i}{(x_{i_{max}} - x_{i_{min}})}$$

- [Robert A van den Berg](#),<sup>1</sup> [Huub CJ Hoefsloot](#),<sup>2</sup> [Johan A Westerhuis](#),<sup>2</sup> [Age K Smilde](#),<sup>1,2</sup> and [Mariët J van der Werf](#)<sup>1</sup>, «Centering, scaling, and transformations: improving the biological information content of metabolomics data», BMC Genomics, 2006.

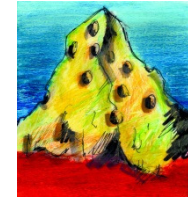
## 4. Ausblick



- Univariat oder multivariat?
- Unit variance, Pareto oder Range scaling?
- Log-Transformation?
- Plausibilitätsprüfung , „Know-how“ aus der Praxis!



Just because you are right,  
does not mean, I am wrong.  
You just haven't seen life  
from my side.

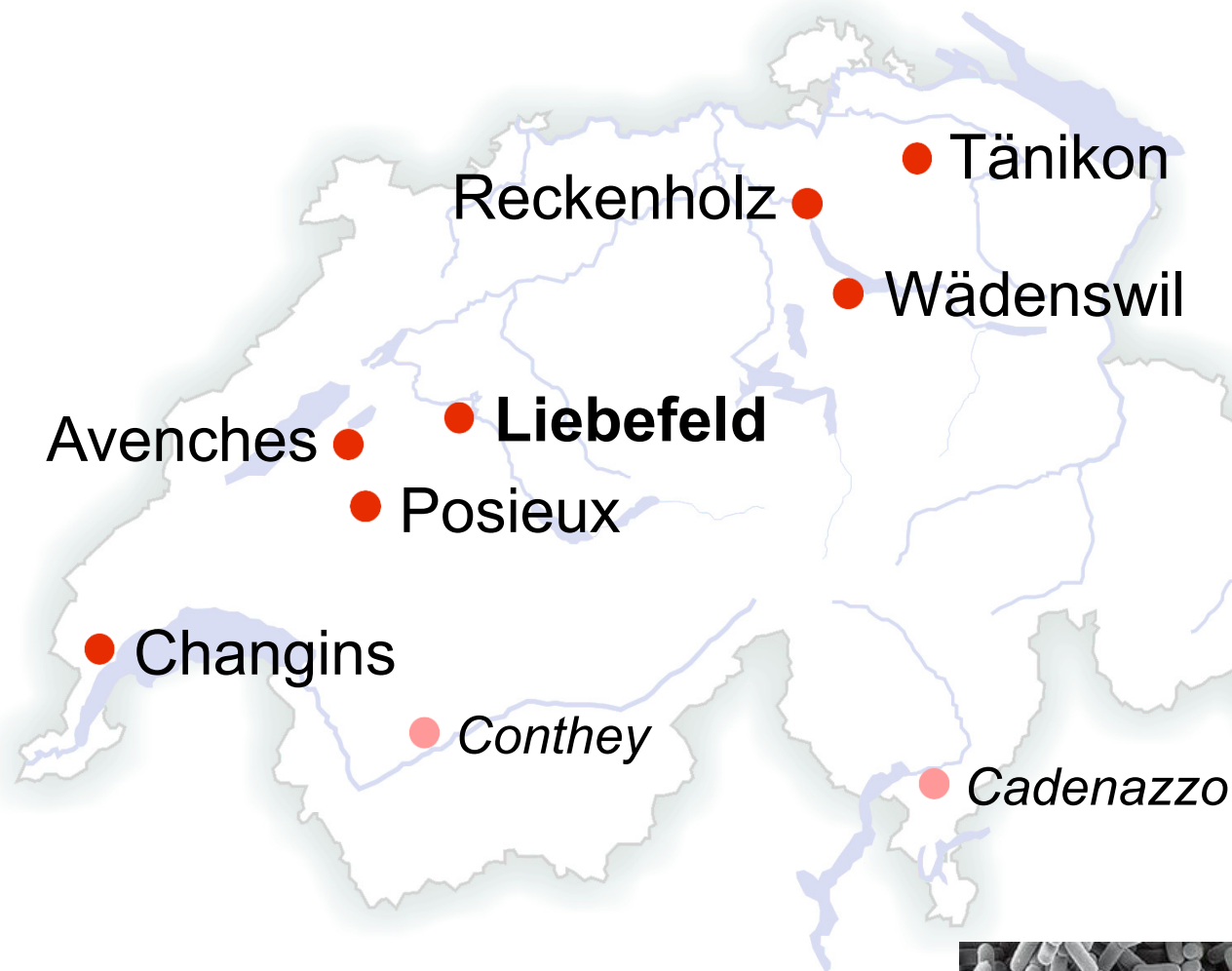
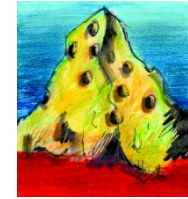


# Fragen?

## Besten Dank für eure Aufmerksamkeit

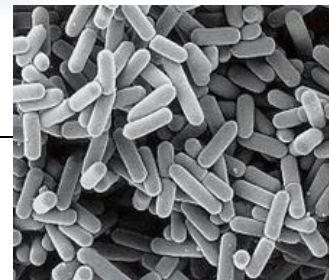


# Agroscope



[www.agroscope.ch](http://www.agroscope.ch)

CAS/DAS Angewandte Statistik  
04. Mai 2018





# Agroscope: Liebefeld Milchsäurebakterien

