► **Data Sciences for Molecular Phenotyping and Precision Medicine team**

► **CEA, INRAE, Paris Saclay University, MetaboHUB, 91191 Gif-sur-Yvette, France**

► https://scidophenia.github.io

**etienne.thevenot@cea.fr**

DE LA RECHERCHE À L'INDUSTRIE

# *ProMetIS*: Proteomics and metabolomics data integration

**Alyssa Imbert and Etienne Thévenot (*ProMetIS* consortium)**

**with the help from Camilo Broc and Olivier Sand**
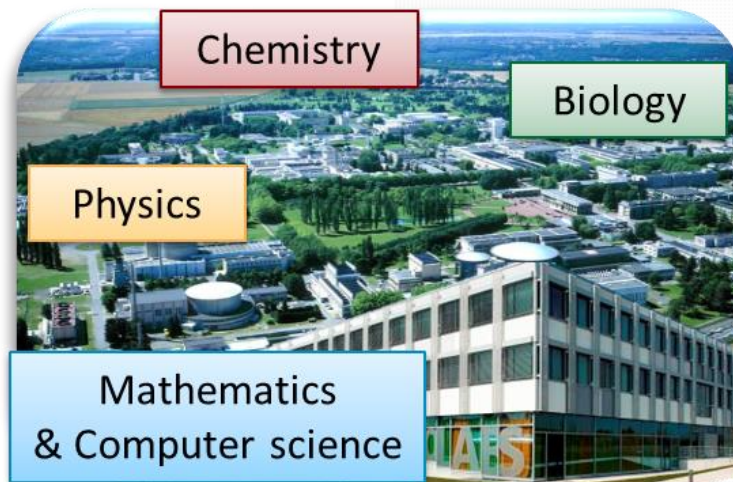
DU Bioinformatique intégrative (DUBii)

# Who we are

▶ **Cluster for data sciences**
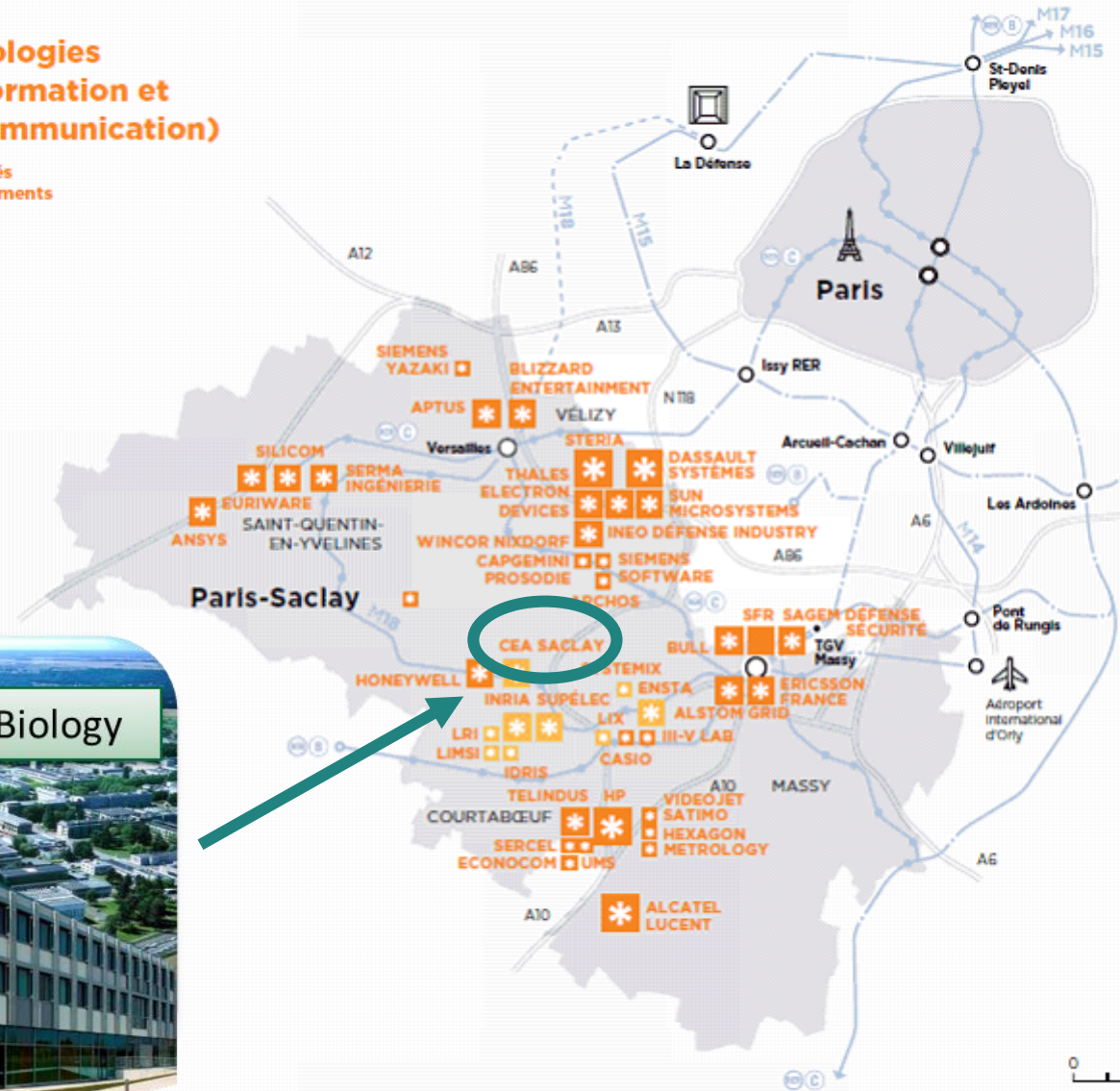▶ **Interdisciplinarity**

- ▶ **omics science**
- ▶ **dedicated to small molecules (< 1kDa)**
- ▶ **involved in metabolic chemical reactions**



Wu RQ, J Dent Res. 2011

Wishart, 2019. Metabolomics for investigating physiological and pathophysiological processes. Physiological Reviews.
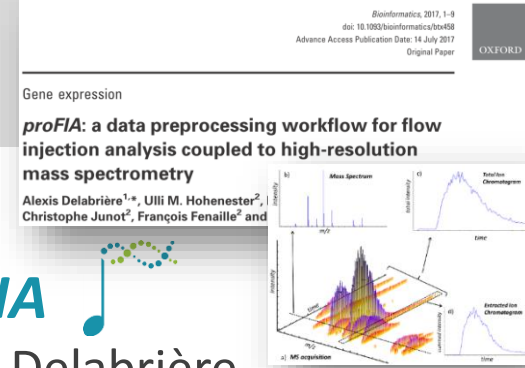
phenomis

ropls

biosigner

biodb

Importing

```
eSet <- phenomis::reading(dirC)
```

Post-processing

Quality control
```
eSet <- phenomis::inspecting(eSet)
```

Signal drift and batch effect correction
```
eSet <- phenomis::correcting(eSet)
```

Transformation
```
eSet <- phenomis::transforming(eSet, methodC = 'log2')
```

Statistics

Univariate hypothesis testing
```
eSet <- phenomis::hypotesting(eSet,  testC = 'limma',
                              factorNamesVc = 'gender', adjustC = 'BH')
```

PCA
```
setPca <- ropls::opls(eSet)
eSet <- ropls::getEset(setPca)
```

Clustering
```
eSet <- phenomis::clustering(eSet, clustersVi = c(2, 2))
```

(O)PLS(-DA)
```
setPlsda <- ropls::opls(eSet, 'gender')
eSet <- ropls::getEset(setPlsda)
```

Feature selection
```
setBiosign <- biosigner::biosign(eSet, 'gender', seedI = 123)
eSet <- biosigner::getEset(setBiosign)
```

Annotation

Chemical annotation
```
eSet <- phenomis::annotating(eSet, databaseC = c('chebi', 'local.ms'))
```

Exporting
```
phenomis::writing(eSet, dirC = getwd())
```

## 1. Clinical questions

Q1

Food allergy in childhood

Q2

Q1: Characterize **early breast milk**
Q2: Predict the development of **food allergy**

## 3. Materials and Methods

EDEN cohort
300 mother/child



## 3. Preliminary results

- Development of a **multi-omics statistical framework** (incl. differential analysis, classification and feature selection)
- Next: **data integration** (multi-blocs, correlation analysis)



## Statistician Post-doctorate
Camilo Broc

## 4. Leverage effects

- ProMetIS (PIA: France Genomique, MetaboHUB, ProFI, IFB),
- MICROB-PREDICT (H2020)
- Master 2 « Systems Immunology » (Sorbonne)



## 5. Multidisciplinary consortium
CEA (DRF/DRT) - INRA

# Introduction: proteomics and metabolomics integration

# Proteomics and metabolomics

## Proteomics

▶ large-scale study of proteins

▶ post-translational modifications

## Metabolomics

▶ small molecule substrates, intermediates, and products of metabolism

▶ peptides, carbohydrates, lipids, nucleosides

▶ "functional readout of the physiological state"

## Proteins – Metabolites

▶ **Interactions**
- Building blocks of proteins
- Substrates, cofactors, products of enzymatic reactions
- Allosteric regulators (enzymes, receptors, transcription factors)
- Post-translational modifications by covalent link to metabolites

Piazza *et al.* (2018). A map of protein-metabolite interactions reveals principles of chemical communication. *Cell*, **172**:358–372.

Figure 5. Key Proteins and Metabolites Characterized in Severe COVID-19 Patients in a Working Model
SARS-CoV-2 may target alveolar macrophages via ACE2 receptor, leading to an increase of secretion of cytokines including IL-6 and TNF-α, which subsequently induce the elevation of various APPs such as SAP, CRP, SAA1, SAA2, and C6, which are significantly upregulated in the severe group. Proteins involved in macrophage, lipid metabolism, and platelet degranulation were indicated with their corresponding expression levels in four patient groups.

Shen *et al.* (2020). Proteomic and metabolomic characterization of COVID-19 patient sera. *Cell*, **9**:59–72.

Parker *et al.* (2019). An integrative systems genetic analysis of mammalian lipid metabolism. *Nature*, **567**:187–193.

Webb-Robertson *et al.* (2016). Bayesian posterior integration for classification of mass spectrometry data. In *Statistical analysis of proteomics, metabolomics, and lipidomics data using mass spectrometry* (pp. 203–211).

Fischer *et al.* (2013). Two birds with one stone: doing metabolomics with your proteomics kit. *Proteomics*, **13**:3371-3386.

Fischer *et al.* (2013). Two birds with one stone: doing metabolomics with your proteomics kit. *PROTEOMICS*, **13**:3371-3386.

Blum *et al.* (2018). Single-platform 'multi-omic' profiling: unified mass spectrometry and computational workflows for integrative proteomics–metabolomics analysis. *Molecular Omics*, **14**:307–319.

Zougman *et al.* (2019). Detergent-free simultaneous sample preparation method for proteomics and metabolomics. *Journal of Proteome Research*, **19**:2838–2844.

Smith *et al.* (2014). Proteomics, lipidomics, metabolomics: a mass spectrometry tutorial from a computer scientist's point of view. *BMC Bioinformatics*. **15**.

## Storage
- raw data: **mzML**



Martens *et al*. (2010). mzML - a community standard for mass spectrometry data. *Molecular & Cellular Proteomics*. **10**.

- processed data: **mzTab**
  - quantification
  - identification

Griss *et al*. (2014). The mzTab data exchange format: Communicating mass-spectrometry-based proteomics and metabolomics experimental results to a wider audience. *Molecular & Cellular Proteomics*. **13**:2765-2775.

## Computation
- R object: **MSnbase**



Gatto *et al*. (2020). MSnbase, efficient and elegant R-based processing and visualization of raw mass spectrometry data. *Journal of Proteome Research*.

Proteomics

Metabolomics

Rurik *et al.* (2020). Metabolomics data processing using OpenMS. *Methods in Molecular Biology,* **2104**.

Adams *et al.* (2020). Skyline for small molecules: A unifying software package for quantitative metabolomics. *Journal of Proteome Research*, **19**:1447-1458.

*phenomis*

SDA

WGCNA

ProMetIS

Li *et al.* (2019). SDA: a semi-parametric differential abundance analysis method for metabolomics and proteomics data. *BMC Bioinformatics*. **20**.

Pei *et al.* (2017). WGCNA application to proteomic and metabolomic data analysis. *Methods in Enzymology*. 135-158.

▶ **Which blocks are the most important for the stratification/prediction?**

▶ **Which features?**

▶ **What is the specific/shared information from each block?**

▶ **How are the features from different blocks correlated?**

▶ **Which biological pathways/networks are significantly involved?**

▶ **Normalization of each block**

▶ **Confounding effects (for each block)**

▶ **Overfitting (limited number of samples)**

    $\Rightarrow$ **validation (statistical, biological)**

▶ **Feature selection**

▶ **Limited annotation of metabolites**

▶ **Redundancy/specificity/ambiguity of chemical/biological identifiers in the databases**

▶ **Partial coverage of the proteome and metabolome**

## Biostatistics

- Fusion
    - low (concatenation of blocks)
    - middle (feature selection/latent variables from each block + model on top)
    - high (one model for each block + vote)
- Correlation networks

## Bioinformatics

- Mapping
- Enrichment
    - Molecule set
    - Topology-based



Ritchie *et al.* (2015). Methods of integrating data to uncover genotype–phenotype interactions. *Nature Reviews Genetics*, **16**:85–97.



Khatri *et al.* (2012). Ten years of pathway analysis: current approaches and outstanding challenges. *PLoS Computational Biology*, **8**: e1002375.

*ProMetIS*: deep phenotyping of mouse models by proteomics and metabolomics
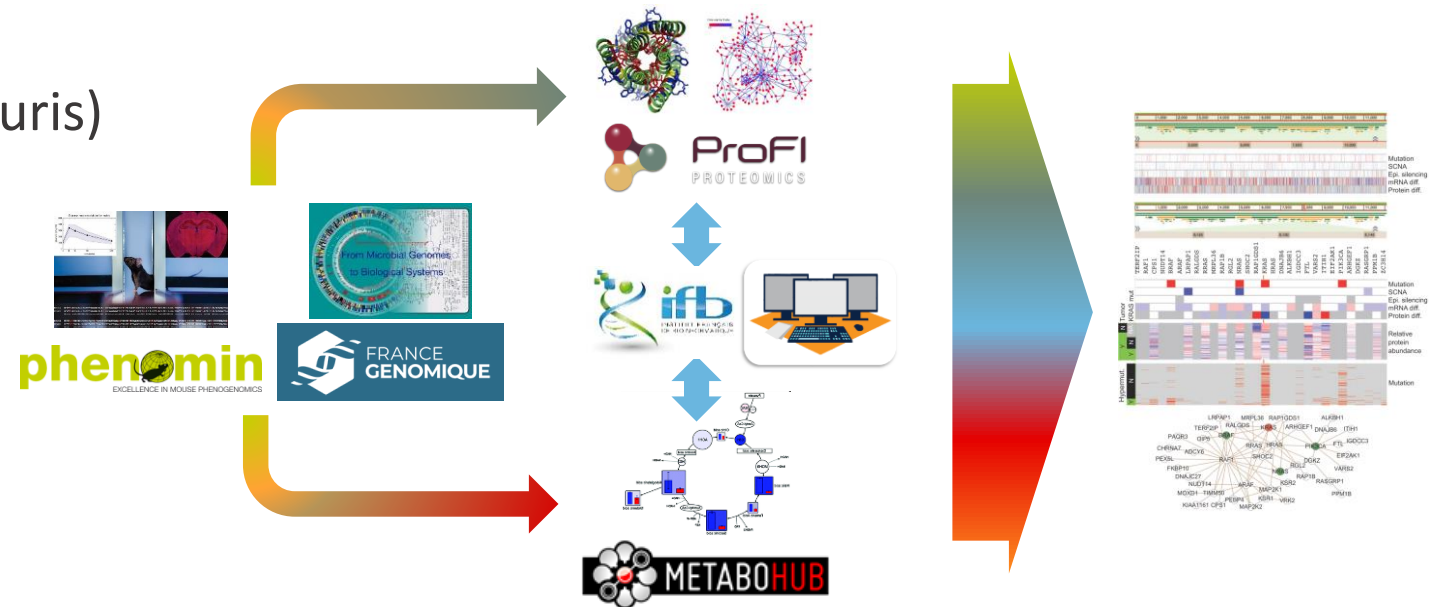
▶ **Objective: high-throughput integration of proteomics and metabolomics data**
- innovative methods
- high-quality datasets
- software tools
- workflows

▶ **Case study: molecular phenotyping of mouse models from the IMPC consortium**

▶ **Partner infrastructures**
- France Génomique
- PHENOMIN (Institut Clinique de la Souris)
- ProFI proteomics
- MetaboHUB
- Institut Français de Bioinformatique

▶ **Post-doctorate**
- Alyssa Imbert

▶ **LAT (linker for activation of T cells; OMIM: 602354) involved in:**

- T-cell receptor (TCR) signaling
- Neurodevelopmental diseases

Roncagalli et al. (2010). LAT signaling pathology: an "autoimmune" condition without T cell self-reactivity. *Trends in Immunology*, **31**:253–259.

Loviglio *et al.* (2017). The immune signaling adaptor LAT contributes to the neuroanatomical phenotype of 16p11.2 BP2-BP3 CNVs. *The American Journal of Human Genetics*, **101**:564–577.

**Preclinical**

**Proteomics**

**Metabolomics**

phenomis

Blank filtering

Pool dilution

Signal drift correction

NAs < 20% for at least 1 group
Variance > 0

log2 transformation

NA < 20% for at least 1 group
Variance > 0

NA < 20% for at least 1 group
Variance > 0

Pool CV ≤ 0.3

Filtering
contaminants

Pool CV / sample CV ≤ 1

Chemical redundancy

Imputation

log2 transformation

| | | | |
|---|---|---|---|
| preclinical | 236 | | |

| | | | |
|---|---|---|---|
| liver_proteomics | 2,187 | liver_metabo_c18hypersil_pos | 5,665 [138] |
| | | liver_metabo_hilic_neg | 2,866 [199] |

| | | | |
|---|---|---|---|
| plasma_proteomics | 446 | plasma_metabo_c18hypersil_pos | 4,788 [113] |
| | | plasma_metabo_hilic_neg | 3,131 [191] |
| | | plasma_metabo_c18acquity_pos | 6,104 [78] |
| | | plasma_metabo_c18acquity_neg | 1,584 [49] |

**ProMetIS**

**https://github.com/IFB-ElixirFr/ProMetIS**

F~indable~ A~ccessible~ I~nteroperable~ R~eusable~

About to be submitted to *Scientific Data*: Imbert *et al*. ProMetIS: deep phenotyping of mouse models by combined proteomics and metabolomics analysis. *submitted*.

# Hands-on

▶ **Loading the datasets**

- restricting to the liver tissue and to the annotated metabolomics features only

▶ **Single-omics analysis**

- exploratory (PCA)
- how much information about the LAT knock-out is provided by each dataset
  - univariate hypothesis testing
  - multivariate PLS-DA

▶ **Multi-omics analysis**

- unsupervised (MCIA)
  - are the two genotypes separated?
  - what about the difference between genders?
- supervised (multi-block PLS-DA)
  - which dataset(s) most contribute to the discrimination?
  - which features most contribute to the discrimination within each dataset?*
  - what is the correlation between those features

# Principal Component Analysis (PCA)

Commissariat à l'énergie atomique et aux énergies alternatives - www.cea.fr

▶ **Visualize the data**

- by selecting a few components which capture most of the spread (variance) of the cloud of samples

▶ **Detect outliers**

- which may bias the computation of the component

▶ **Detect clusters of samples**

- which may suggest an internal structuration of the data

**$p$ = 110 (quantitative) variables**
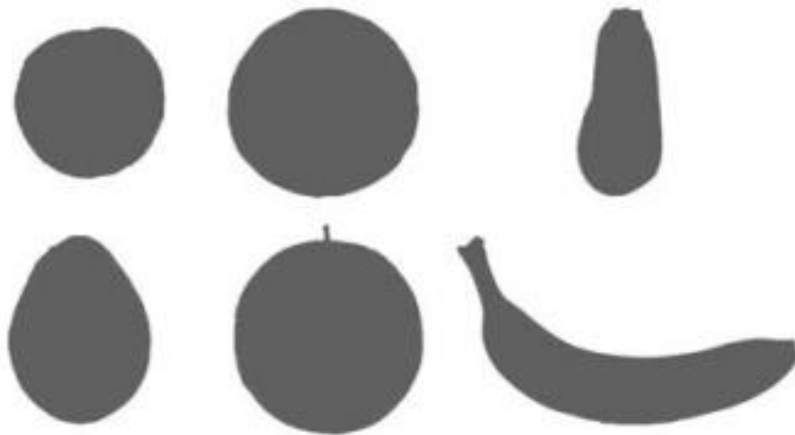
$n$ = 183 samples

|  | 1,7-Dimethyluric acid | Dehydroepiandrosterone sulfate | Acetaminophen glucuronide |
|---|---|---|---|
| H011 | 2114 | 29025 | 44 |
| H023 | 43274 | 639 | 2 |
| H033 | 22386 | 325 | 1933 |
| H042 | 8185 | 13938 | 933 |
| H052 | 22385 | 357 | 5004 |
| H062 | 6380 | 292 | 1 |
| H073 | 10012 | 22781 | 1 |
| H083 | 30414 | 105 | 1 |
| H092 | 6637 | 35156 | 1 |
| H103 | 12100 | 2 | 1 |
| H114 | 33362 | 149041 | 46 |
| H124 | 11197 | 84536 | 1 |
| H134 | 18698 | 34053 | 254 |
| H145 | 14435 | 212398 | 52 |
| H157 | 31732 | 19317 | 2200 |
| H168 | 10221 | 78 | 475 |
| H180 | 22936 | 463 | 1 |
| H189 | 14423 | 1039 | 220 |
| H199 | 2888 | 12272 | 37 |
| H209 | 12563 | 100236 | 2 |

...

**X**

**1 variable**

**2 variables**

**3 variables**
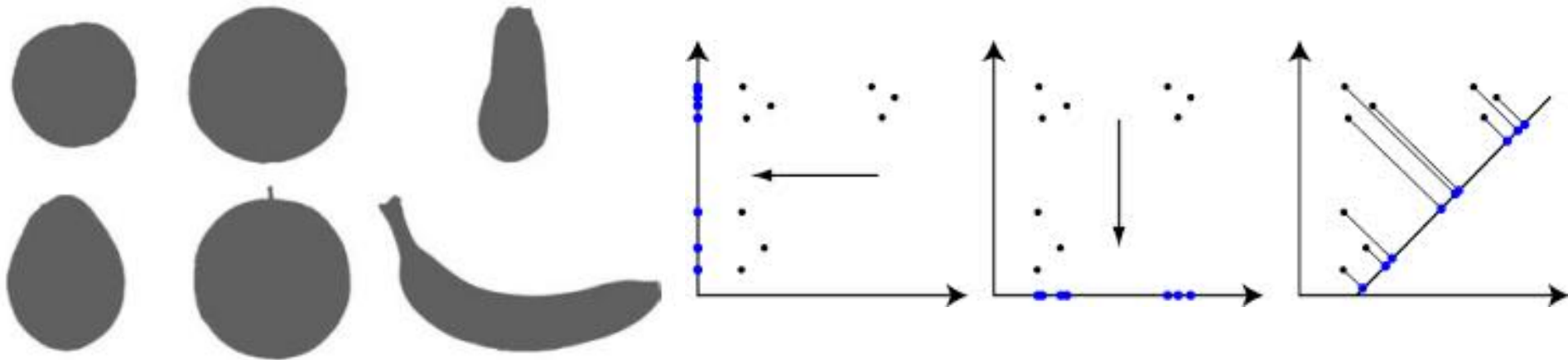
**p variables**

=> Dimension reduction

▶ **Projected distances as high as possible**



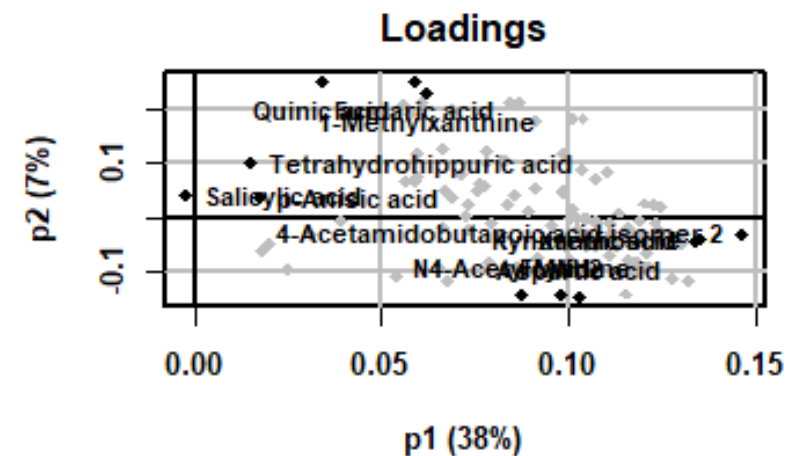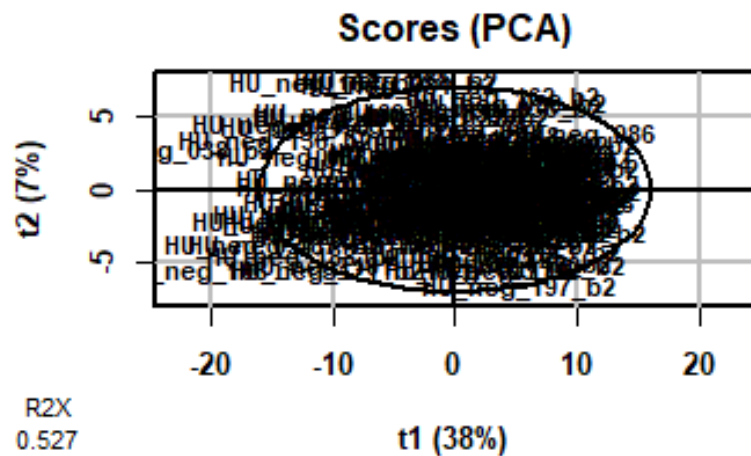Husson and Pages (2011). Exploratory multivariate analysis by example using R. Chapman & Hall/CRC
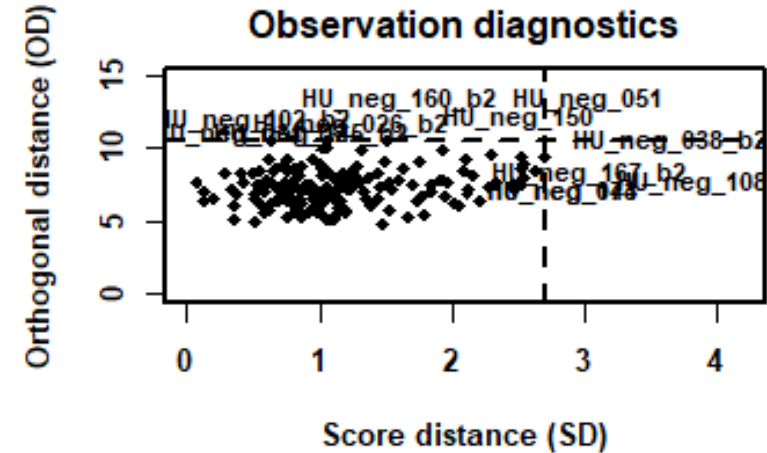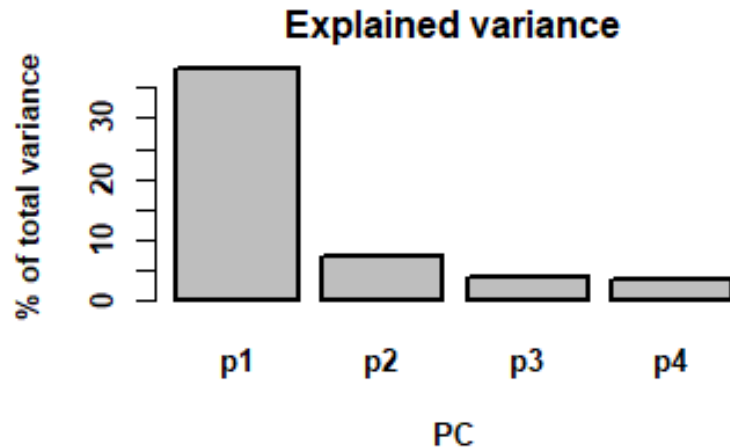
▶ **Projected distances as high as possible**

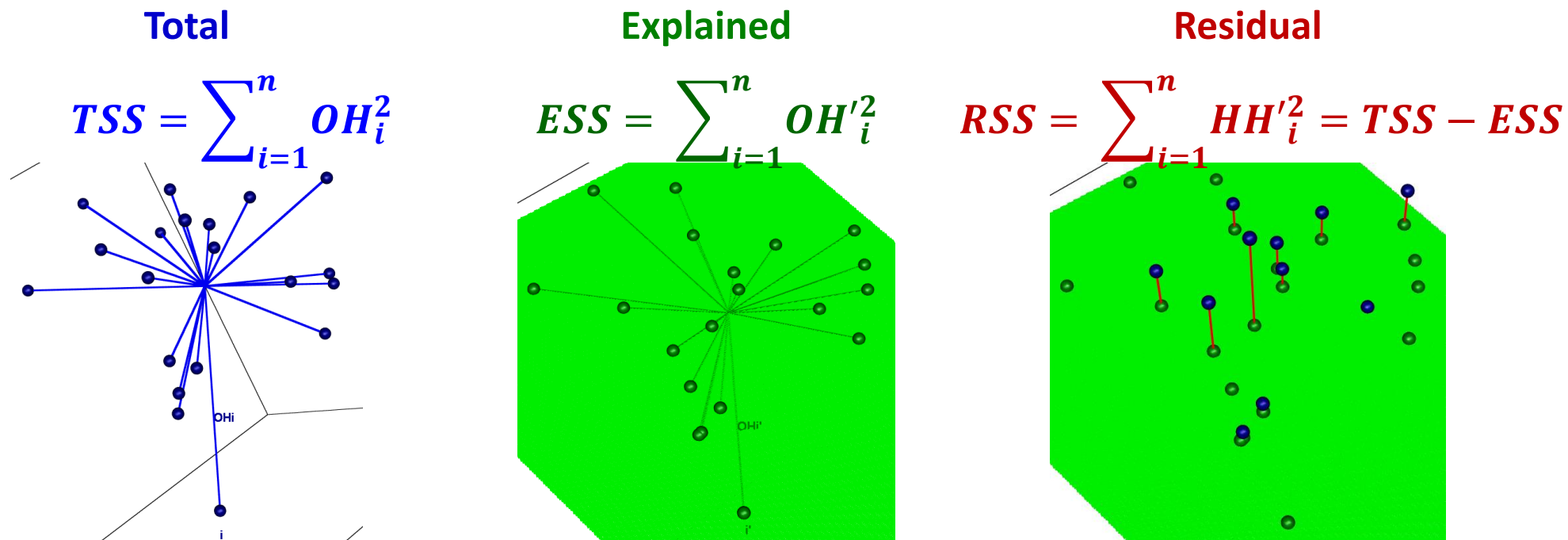▶ **Define new variables as linear combination of original ones**



Husson and Pages (2011). Exploratory
multivariate analysis by example using R.
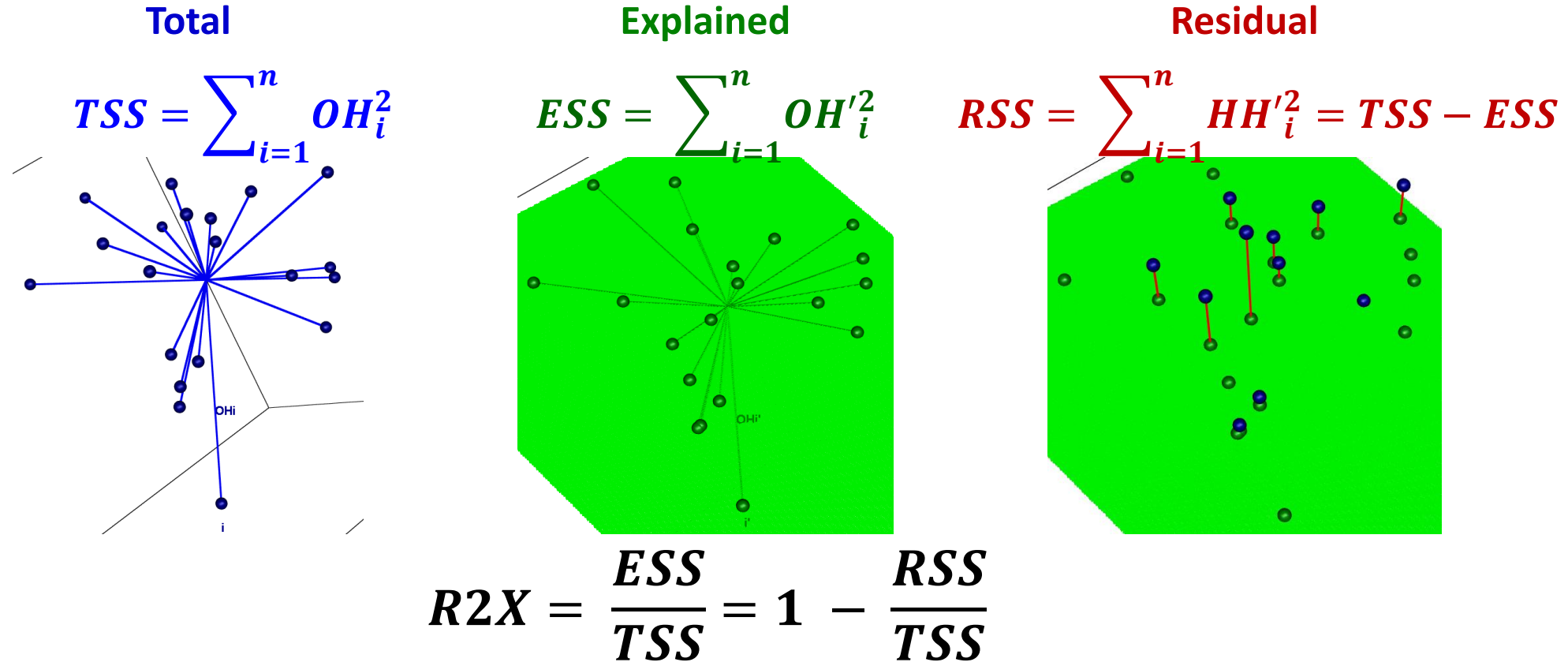*Chapman & Hall/CRC*

▶ **scree plot, outliers, and the loading and score plots**

**Total**

$$TSS = \sum_{i=1}^{n} OH_i^2$$

**Explained**

$$ESS = \sum_{i=1}^{n} OH_i'^2$$

**Residual**

$$RSS = \sum_{i=1}^{n} HH_i'^2 = TSS - ESS$$

**Total**      **Explained**      **Residual**

$$TSS = \sum_{i=1}^{n} OH_i^2$$

$$ESS = \sum_{i=1}^{n} OH'^2_i$$

$$RSS = \sum_{i=1}^{n} HH'^2_i = TSS - ESS$$

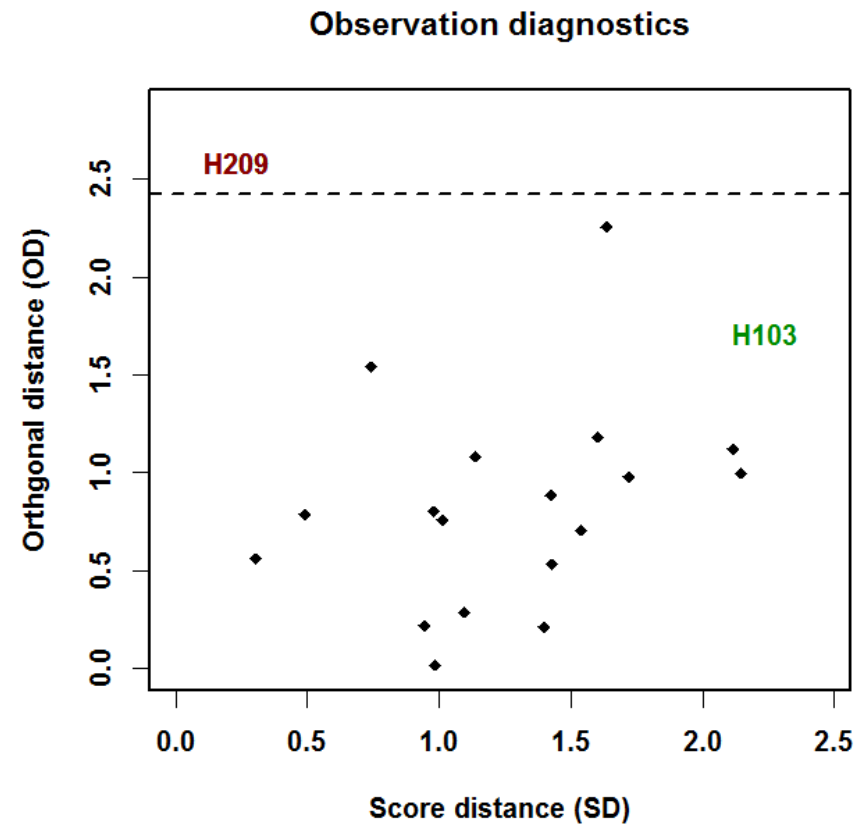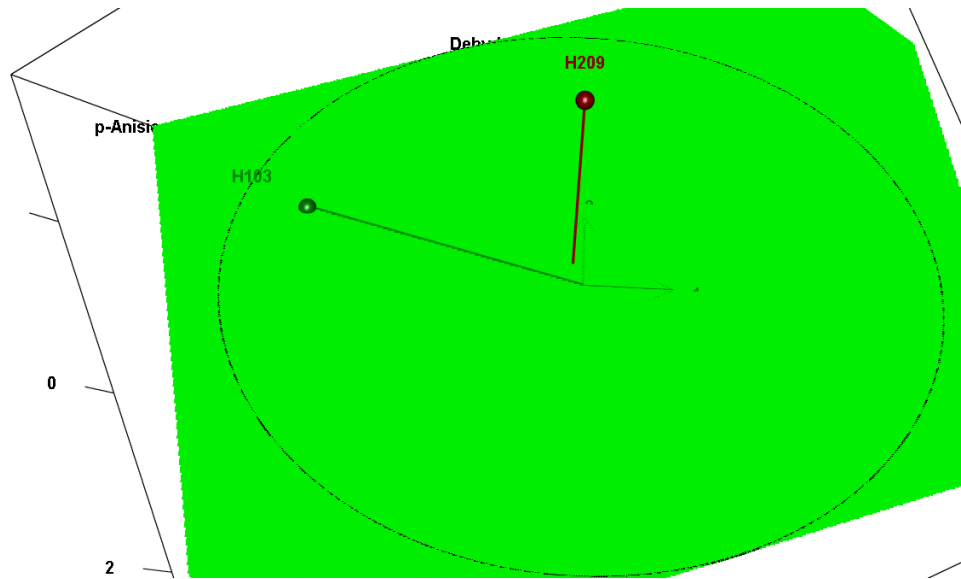$$R2X = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}$$

▶ **R2X increases with the number of components in the model**

▶ **For a given number of components, the higher the R2X, the more inertia is captured by the model (projection)**

$$0 \le R2X \le 1$$

▶ **Check that the first components capture most of the variance**



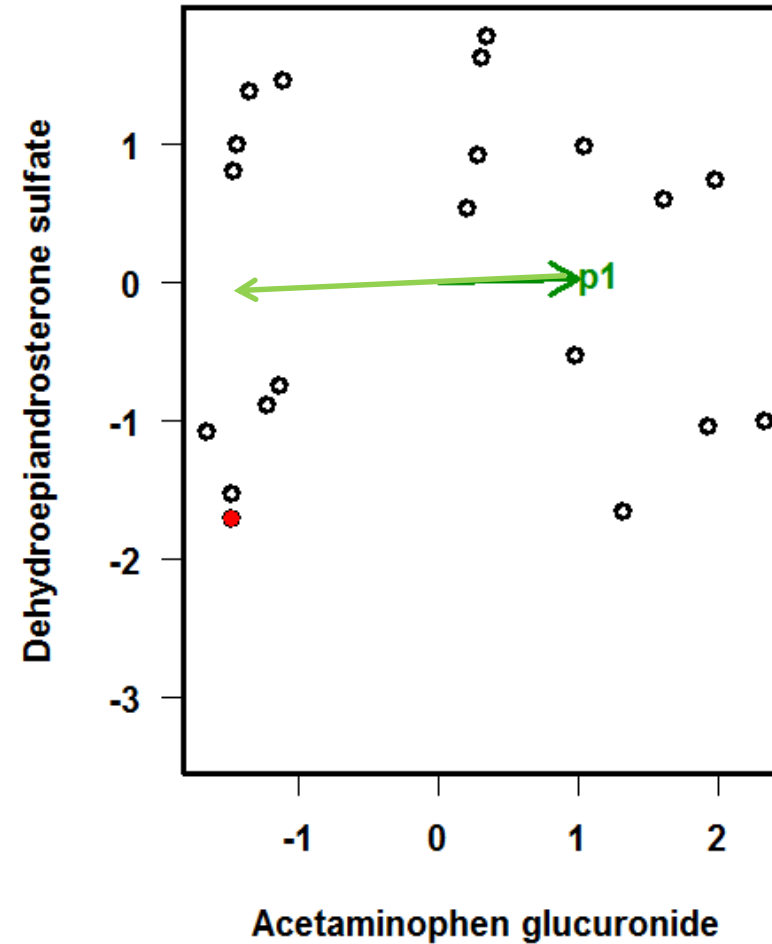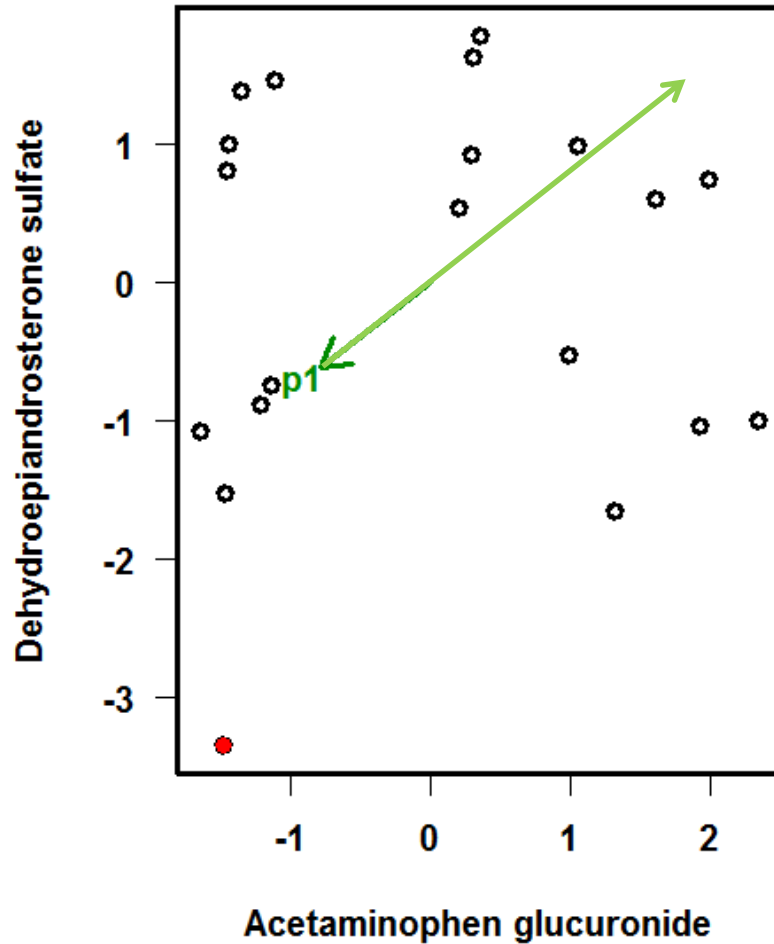Mount Yamnuska, Alberta. Wikipedia

▶ **Samples which may bias the PCA computation and/or may not be faithfully visualized by the score plot**

Hubert M., Rousseeuw P. and Vanden Branden K. (2005). ROBPCA: a new approach to robust principal component analysis. *Technometrics,* **47:**64-79. DOI:

▶ **Variables are usually centered:**

$$x'_j = x_j - \bar{x}_j$$

▶ **In addition, variables may be**

- unit-variance scaled (default in *ropls*):

$$x''_j = \frac{x'_j}{\sigma_j}$$

- pareto scaled:

$$x''_j = \frac{x'_j}{\sqrt{\sigma_j}}$$

## ▶ Loading

```
sacurine_dir.c <-
"C:/Users/et207099/Documents/sources/training/inst/extdata/sacurine_annotated_postprocessed"

sacurine.eset <- phenomis::reading(sacurine_dir.c)
```
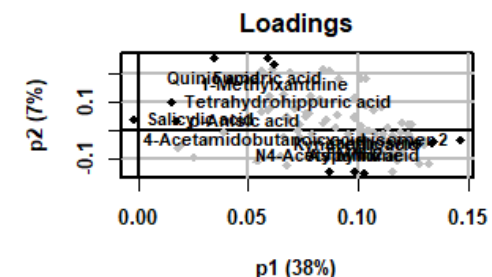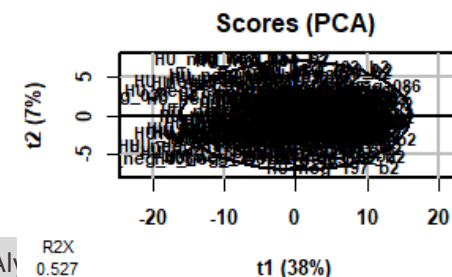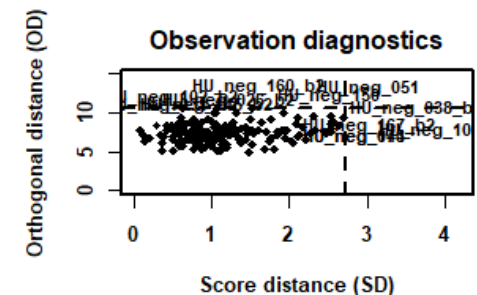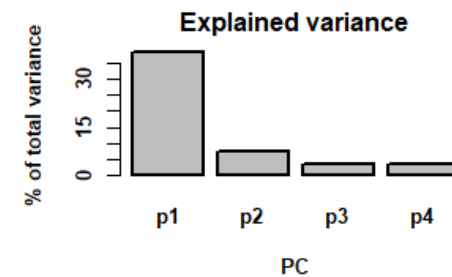
## ▶ Inspecting

```
sacurine.eset <- phenomis::inspecting(sacurine.eset)
```
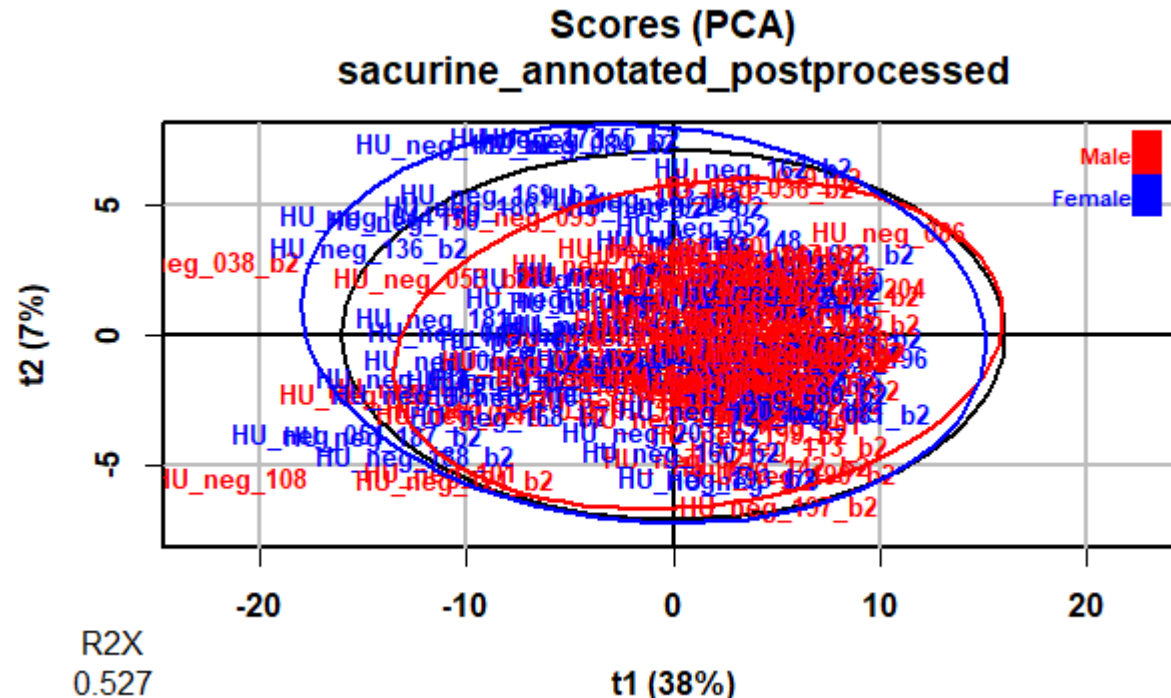
## ▶ Computing the PCA

```
sacurine.pca <- ropls::opls(sacurine.eset)
```

▶ **Coloring the score plot according to 'gender' (column of the sampleMetadata)**

```
ropls::plot(sacurine.pca,
            typeVc = "x-score",
            parAsColFcVn = Biobase::pData(sacurine.eset)[, "gender"])))
```

Commissariat à l'énergie atomique et aux énergies alternatives       ProMetIS: Proteomics and Metabolomics Data Integration – Alyssa Imbert and Etienne Thévenot       31 mars 2021

46

▶ **Studying the variables which most contribute to the new components**

```
ropls::plot(sacurine.pca,
            typeVc = "x-loading")
```
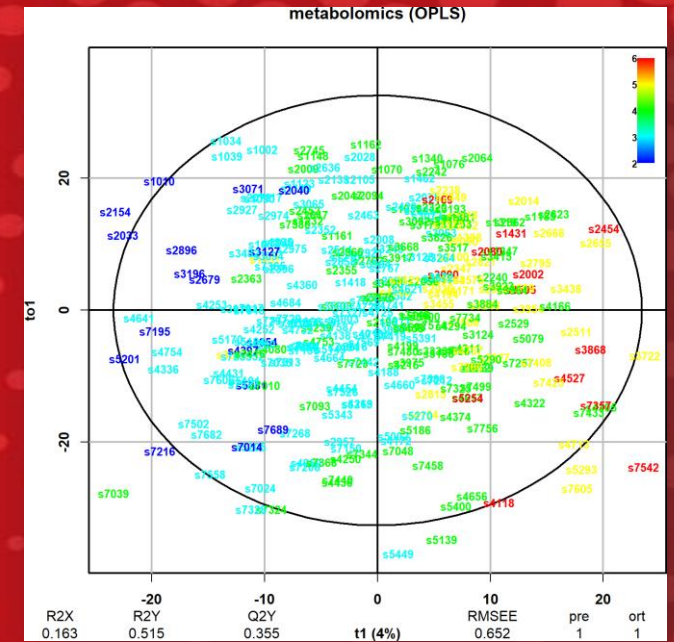
▶ **Getting back the ExpressionSet object**

```
sacurine.eset <- ropls::getEset(sacurine.pca)
```

▶ **The scores and loadings values have been added to the sampleMetadata and variableMetadata:**

```
head(Biobase::pData(sacurine.eset)[, c("PCA_xscor-p1", "PCA_xscor-p2")])
```

```
head(Biobase::fData(sacurine.eset)[, c("PCA_xload-p1", "PCA_xload-p2")])
```

▶ **Husson F., Le S. and Pages J. (2011). Exploratory multivariate analysis by example using *R. Chapman & Hall/CRC***

▶ **Baccini A. (2010). Statistique Descriptive Multidimensionnelle (pour les nuls). *Institut de Mathématiques de Toulouse, Université Paul Sabatier.***

▶ **Ringner M. (2008). What is principal component analysis? *Nature Biotechnology*, 26:303-304.**

▶ **Wehrens, R. (2011). Chemometrics with R. *Springer*. https://doi.org/10.1007/978-3-642-17841-2**

# Partial Least Squares (PLS)

Commissariat à l'énergie atomique et aux énergies alternatives - www.cea.fr

▶ **Powerful regression method when**

$$n_{samples} < p_{variables}$$

▶ **Complementary to univariate hypothesis testing (where variables are tested independantly)**

▶ **Risk of overfitting: i.e., building a model whose (apparently) good performances result from chance only**

**$p$ = 109 variables (quantitatives)**

| | (2-methoxyethoxy)propanoic acid isomer | (gamma)Glu-Leu/Ile | 1-Methyluric acid | 1-Methylxanthine | 1,3-Dimethyluric acid | ... | Threonic acid/Erythronic acid | Tryptophan | Valerylglycine isomer 1 | Valerylglycine isomer 2 | Xanthosine |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **HU_011** | 3.02 | 3.89 | 3.87 | 3.72 | 3.54 | ... | 4.31 | 4.01 | 4.02 | 3.89 | 4.08 |
| **HU_014** | 3.81 | 4.28 | 3.84 | 3.78 | 3.93 | ... | 4.47 | 4.42 | 3.88 | 4.18 | 4.20 |
| **HU_015** | 3.52 | 4.20 | 4.10 | 4.29 | 3.96 | ... | 4.12 | 4.44 | 4.19 | 4.25 | 4.12 |
| **HU_017** | 2.56 | 4.32 | 4.54 | 4.43 | 4.23 | ... | 4.56 | 4.54 | 4.15 | 4.29 | 4.25 |
| **HU_018** | 3.78 | 4.63 | 4.18 | 4.12 | 4.01 | ... | 4.45 | 4.22 | 4.10 | 4.14 | 4.36 |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| **HU_205** | 3.86 | 4.54 | 4.24 | 4.19 | 4.38 | ... | 4.36 | 4.12 | 4.16 | 4.22 | 4.41 |
| **HU_206** | 1.32 | 4.34 | 4.62 | 4.61 | 4.82 | ... | 4.27 | 4.04 | 3.93 | 4.28 | 4.39 |
| **HU_207** | 4.19 | 4.28 | 4.48 | 4.46 | 4.45 | ... | 4.64 | 4.00 | 4.74 | 4.65 | 4.26 |
| **HU_208** | 3.75 | 4.52 | 4.36 | 4.36 | 4.23 | ... | 4.70 | 4.69 | 4.44 | 4.63 | 4.49 |
| **HU_209** | 4.21 | 4.68 | 4.19 | 4.21 | 4.15 | ... | 4.52 | 4.50 | 4.47 | 4.47 | 4.22 |

*$n$ = 183 samples*

**X**

*1* response

| | IMC |
|---|---|
| HU_011 | 19.75 |
| HU_014 | 22.64 |
| HU_015 | 22.72 |
| HU_017 | 23.03 |
| HU_018 | 20.96 |
| ... | ... |
| HU_205 | 28.37 |
| HU_206 | 22.15 |
| HU_207 | 19.47 |
| HU_208 | 18.61 |
| HU_209 | 21.48 |

*y*

$p$ = 109 variables (quantitatives)          *1* response

| | (2-methoxyethoxy)propanoic acid isomer | (gamma)Glu-Leu/Ile | 1-Methyluric acid | 1-Methylxanthine | 1,3-Dimethyluric acid | … | Threonic acid/Erythronic acid | Tryptophan | Valerylglycine isomer 1 | Valerylglycine isomer 2 | Xanthosine |
|---|---|---|---|---|---|---|---|---|---|---|---|
| HU_011 | 3.02 | 3.89 | 3.87 | 3.72 | 3.54 | … | 4.31 | 4.01 | 4.02 | 3.89 | 4.08 |
| HU_014 | 3.81 | 4.28 | 3.84 | 3.78 | 3.93 | … | 4.47 | 4.42 | 3.88 | 4.18 | 4.20 |
| HU_015 | 3.52 | 4.20 | 4.10 | 4.29 | 3.96 | … | 4.12 | 4.44 | 4.19 | 4.25 | 4.12 |
| HU_017 | 2.56 | 4.32 | 4.54 | 4.43 | 4.23 | … | 4.56 | 4.54 | 4.15 | 4.29 | 4.25 |
| HU_018 | 3.78 | 4.63 | 4.18 | 4.12 | 4.01 | … | 4.45 | 4.22 | 4.10 | 4.14 | 4.36 |
| … | … | … | … | … | … | … | … | … | … | … | … |
| HU_205 | 3.86 | 4.54 | 4.24 | 4.19 | 4.38 | … | 4.36 | 4.12 | 4.16 | 4.22 | 4.41 |
| HU_206 | 1.32 | 4.34 | 4.62 | 4.61 | 4.82 | … | 4.27 | 4.04 | 3.93 | 4.28 | 4.39 |
| HU_207 | 4.19 | 4.28 | 4.48 | 4.46 | 4.45 | … | 4.64 | 4.00 | 4.74 | 4.65 | 4.26 |
| HU_208 | 3.75 | 4.52 | 4.36 | 4.36 | 4.23 | … | 4.70 | 4.69 | 4.44 | 4.63 | 4.49 |
| HU_209 | 4.21 | 4.68 | 4.19 | 4.21 | 4.15 | … | 4.52 | 4.50 | 4.47 | 4.47 | 4.22 |

| | IMC |
|---|---|
| HU_011 | 19.75 |
| HU_014 | 22.64 |
| HU_015 | 22.72 |
| HU_017 | 23.03 |
| HU_018 | 20.96 |
| … | … |
| HU_205 | 28.37 |
| HU_206 | 22.15 |
| HU_207 | 19.47 |
| HU_208 | 18.61 |
| HU_209 | 21.48 |

*$n$ = 183 sample*

$$f(X) = y$$

**p = 109 variables (quantitatives)**

*n'* samples

| | (2-methoxyethoxy)propanoic acid isomer | (gamma)Glu-Leu/Ile | 1-Methyluric acid | 1-Methylxanthine | 1,3-Dimethyluric acid | … | Threonic acid/Erythronic acid | Tryptophan | Valerylglycine isomer 1 | Valerylglycine isomer 2 | Xanthosine |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **new 1** | 3.00 | 4.47 | 4.54 | 4.54 | 4.62 | … | 4.46 | 4.30 | 4.44 | 4.41 | 4.54 |
| **new 2** | 3.48 | 4.20 | 3.73 | 3.31 | 3.44 | … | 4.57 | 4.17 | 4.15 | 4.16 | 4.26 |
| **new 3** | 4.03 | 2.55 | 4.27 | 4.23 | 4.34 | … | 4.26 | 3.58 | 4.07 | 3.96 | 4.15 |

| | IMC |
|---|---|
| **new 1** | ? |
| **new 2** | ? |
| **new 3** | ? |

$$\textbf{f(X')} \qquad = \qquad \textbf{y'}$$

plane

projection

1.st comp.

2.nd comp.

**Direction in plane defining best correlation with Y**

**(c1 t1 + c2 t2 + ...)**

$$cov(x, y) = \frac{1}{(n-1)} \Sigma_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})$$

$$= \frac{1}{(n-1)} cor(x, y) \|x\| \|y\|$$
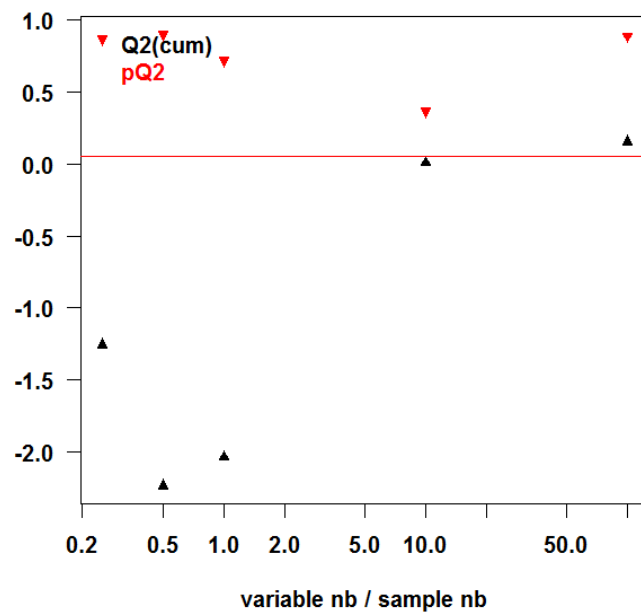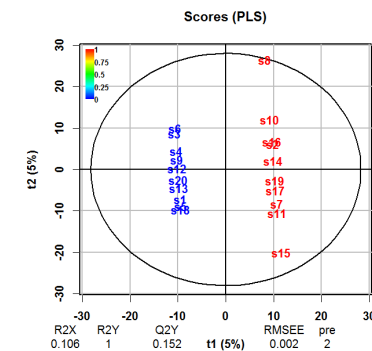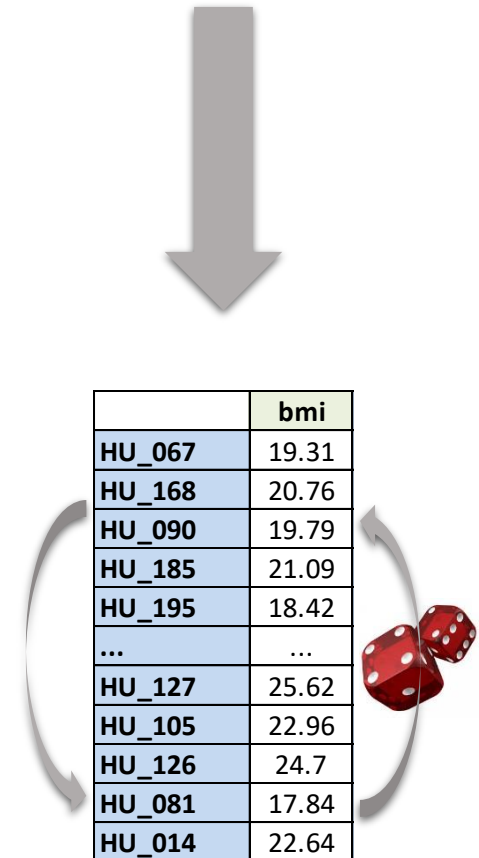
Wold et al. (2001). PLS-regression: a basic tool of chemometrics. Chemometrics and Intelligent Laboratory Systems, 58:109-130.

▶ **PCA finds the directions of maximum variance**

▶ **PLS includes the labels into the model**

Scores (PCA)

Scores (PCA)

Scores (PCA)

Scores (PLS)

Model overview

► **X: 20 × 2,000 matrix of random numbers**
- Uniform distribution between 0 and 1

► **Y: 20 x 1 matrix of random labels**
- 0 or 1 values

adpated from Wehrens (2011). Chemometrics with R. Springer.

Scores (PLS)

▶ **Permutation testing: comparing the R2Y and Q2Y values of the model built with the true Y labels with $n_{perm}$ models built with random permutation of Y labels**

Szymanska E., Saccenti E., Smilde A. and Westerhuis J. (2012). Double-check: validation of diagnostic statistics for PLS-DA models in metabolomics studies. *Metabolomics,* **8**:3-16. DOI:

► **Counting the number of $R2Y$ (and $Q2Y$) metrics from random models which are superior to the values of the true model gives an indication of the significance of the PLS modelling**
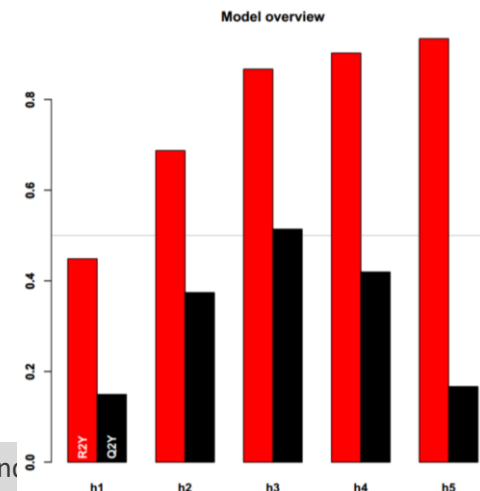


pR2Y = 0.01, pQ2 = 0.01

$R2Y$ and $Q2Y$ of the model with the true **y** values

'permutation'

Similarity between $\mathbf{y}_{true}$ and $\mathbf{y}_{random}$

$$\frac{\text{variables}}{\text{samples}} = $$

| 0.2 | 0.5 | 1 | 10 | 100 |

- Wold S., Sjöström M. and Eriksson L. (2001). PLS-regression: a basic tool of chemometrics. *Chemometrics and Intelligent Laboratory Systems,* **58:**109-130.

- Trygg J., Holmes E. and Lundstedt T. (2007). Chemometrics in Metabonomics. *Journal of Proteome Research*, **6:**469-479.

- Brereton R.G. and Lloyd G.R. (2014). Partial least squares discriminant analysis: taking the magic away. *Journal of Chemometrics,* **28:**213-225.

$X$

$y_{random}$

► **Counting the number of $R2Y$ (and $Q2Y$) metrics from random models which are superior to the values of the true model gives an indication of the significance of the PLS modelling**



$R2Y$ and $Q2Y$ of the model with the true **y** values

'permutation'

Similarity between $\mathbf{y}_{true}$ and $\mathbf{y}_{random}$

▶ $0 \leq R2X \leq 1$: **percentage of X inertia explained by the model**

▶ $0 \leq R2Y \leq 1$: **percentage of Y inertia explained by the model**

▶ $0 \leq Q2Y \leq 1$: **estimation of the predictive performance of the model by cross-validation**

▶ $R2X$ **and** $R2Y$ **increase with the number of components while** $Q2Y$ **reaches a maximum (due to overfitting):**

'overview'
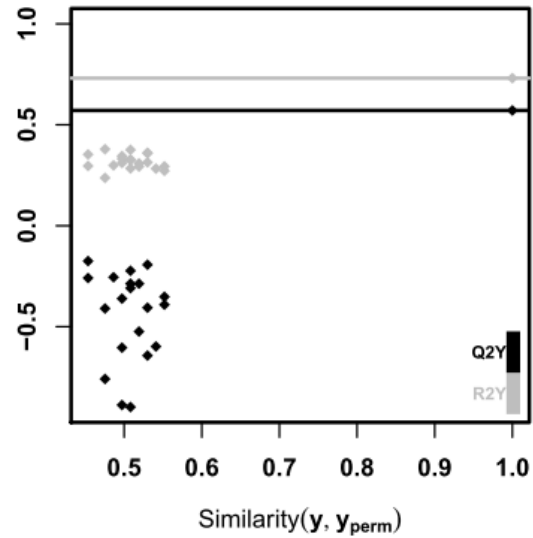


Model overview

# Partial Least Squares – Discriminant Analysis (PLS-DA)

# Regression and classification

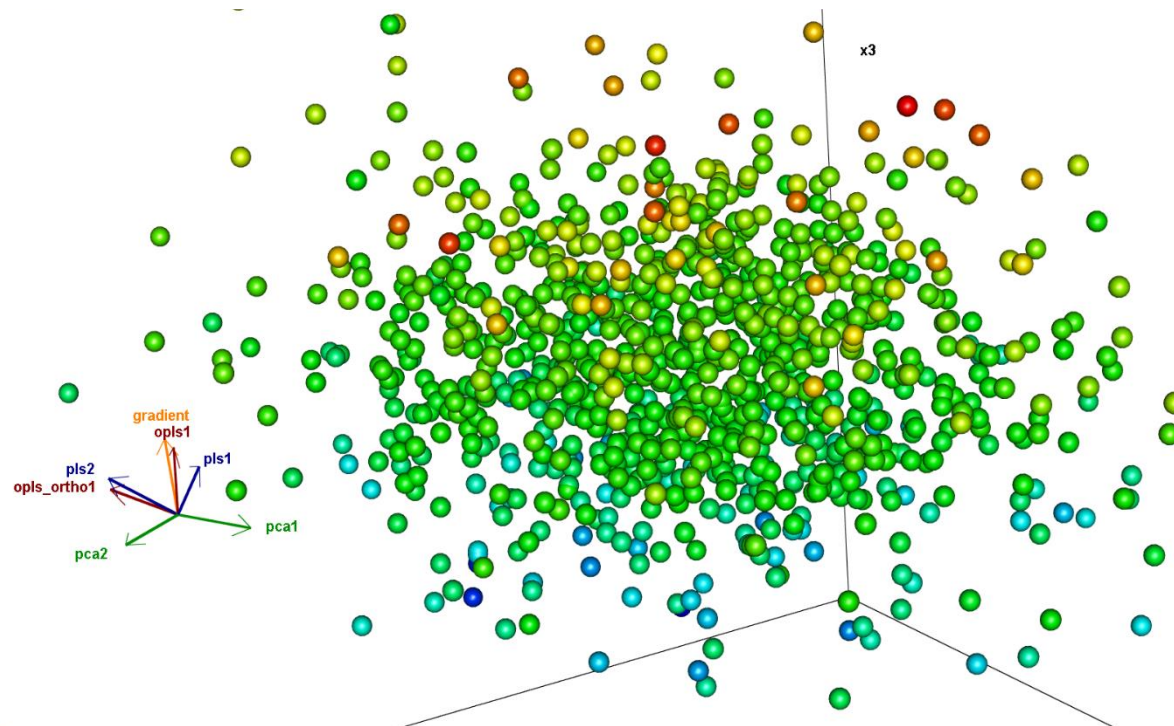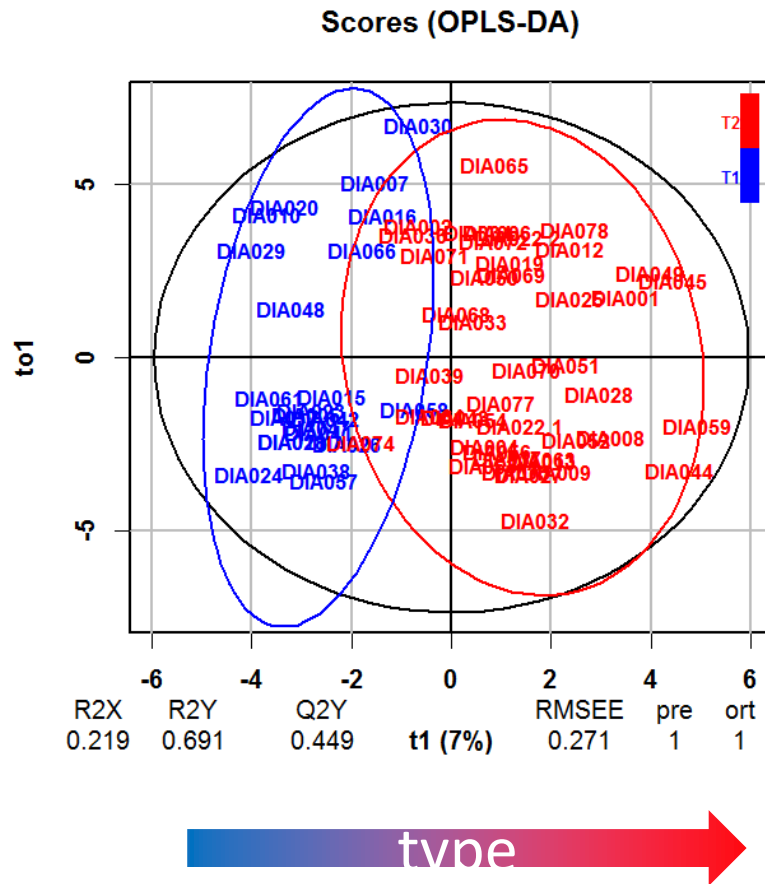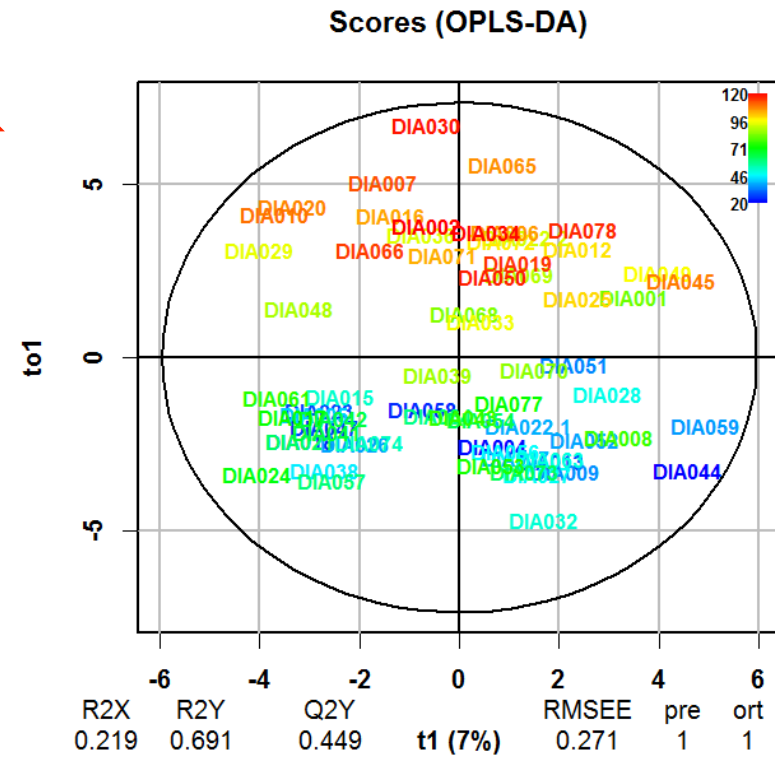| Response y | Example | Approach | PLS method |
|---|---|:---:|:---:|
| Quantitative | BMI | regression | PLS |
| Qualitative | gender | classification | PLS-DA |

► **The two response levels are encoded as numbers**

▶ **Separately models the variations of the predictors correlated and orthogonal to the response**

▶ **Improves the interpretation of the components but not the overall predictive performance of the model**

▶ **Only one predictive component required for single response models**

▶ **Note: As with PLS, care should be taken to avoid too many (orthogonal) components (which would result in overfitting)**
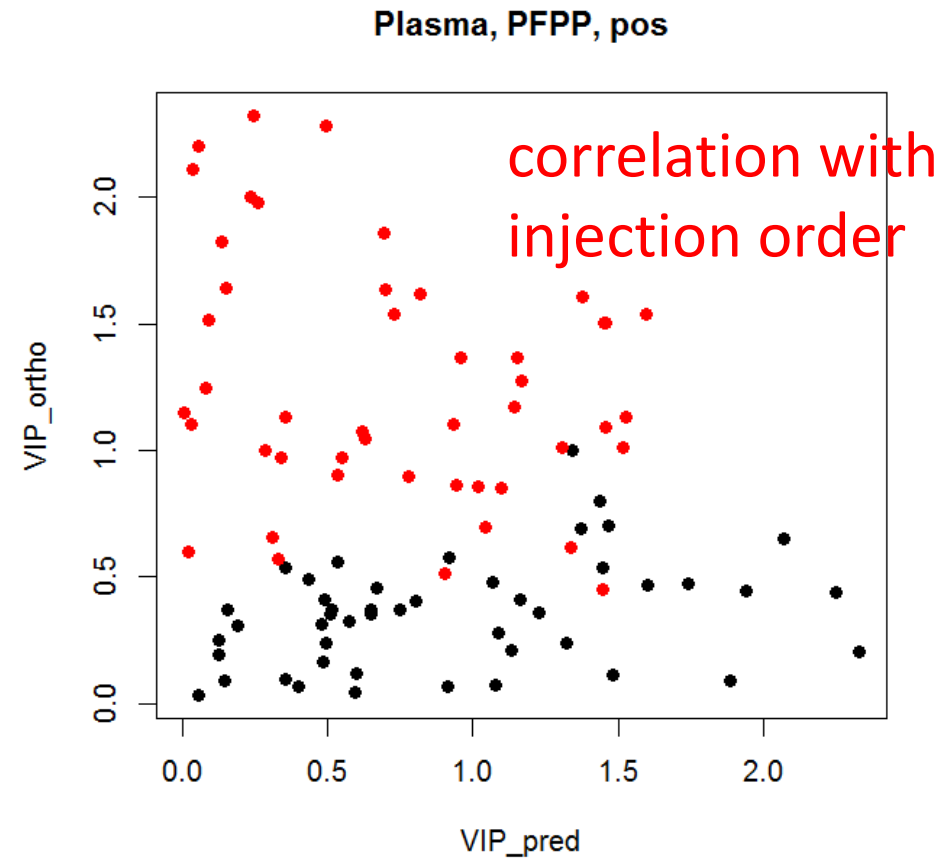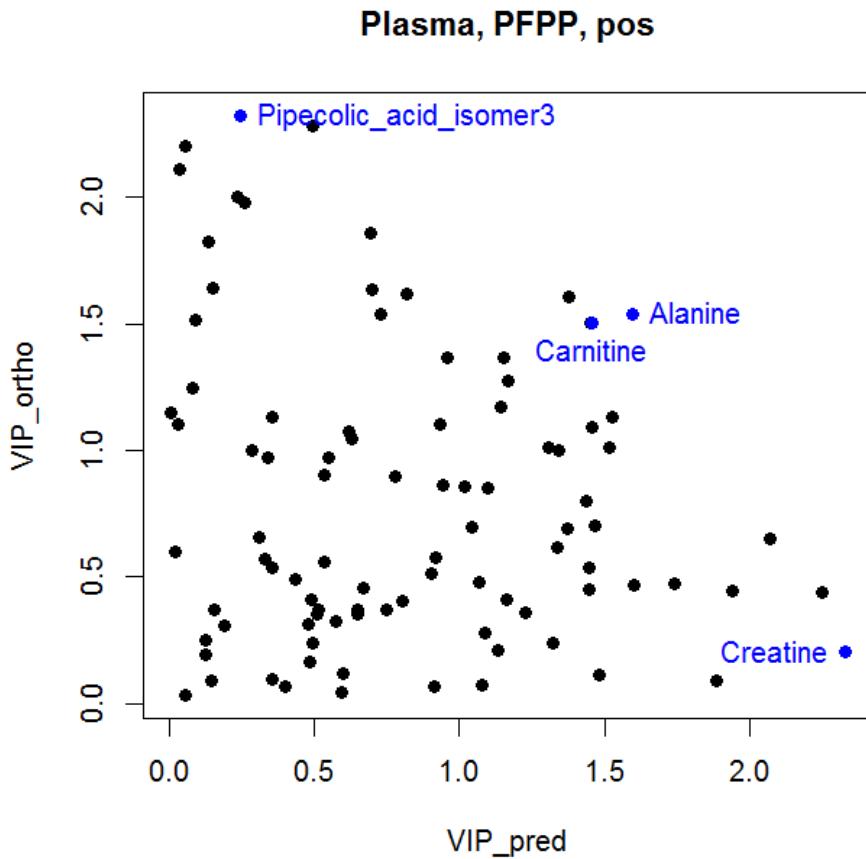
▶**Variation not correlated to the response (e.g., technical bias) is modelled separately by the orthogonal component(s)**

**=> The first predictive component is strongly correlated to the response**

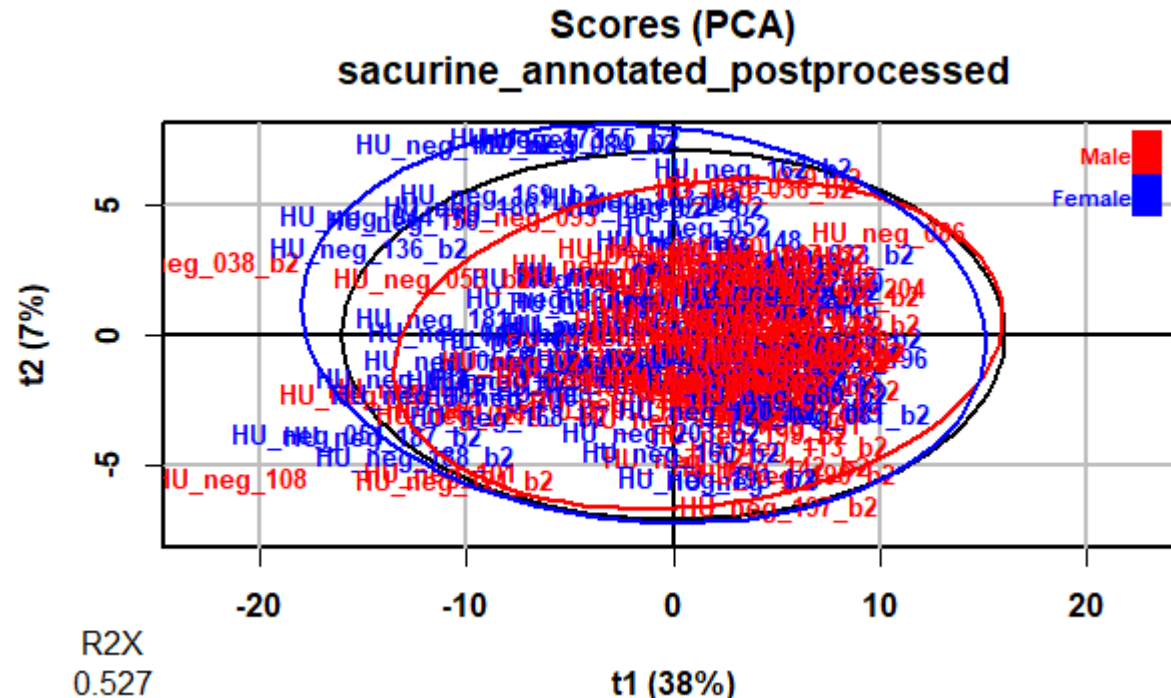Galindo-Prieto et al (2014). *Journal of Chemometrics,* 28, 623-632.

▶ **permutation, overview, outlier, and score plots displayed as the default ('summary')**

▶ **Loading**

```
sacurine_dir.c <-
"C:/Users/et207099/Documents/sources/training/inst/extdata/sacurine_annotated_postprocessed"

sacurine.eset <- phenomis::reading(sacurine_dir.c)
```

▶ **Inspecting**

```
sacurine.eset <- phenomis::inspecting(sacurine.eset)
```

▶ **Computing the PCA**

```
sacurine.pca <- ropls::opls(sacurine.eset)
```

► **Coloring the score plot according to 'gender' (column of the sampleMetadata)**

```
ropls::plot(sacurine.pca,

            typeVc = "x-score",

            parAsColFcVn = Biobase::pData(sacurine.eset)[, "gender"])))
```

▶ **Getting back the ExpressionSet object**

```
sacurine.eset <- ropls::getEset(sacurine.pca)
```

▶ **The scores and loadings values have been added to the sampleMetadata and variableMetadata:**

```
head(Biobase::pData(sacurine.eset)[, c("PCA_xscor-p1", "PCA_xscor-p2")])
head(Biobase::fData(sacurine.eset)[, c("PCA_xload-p1", "PCA_xload-p2")])
```

▶ **Wold S., Sjöström M. and Eriksson L. (2001). PLS-regression: a basic tool of chemometrics.** *Chemometrics and Intelligent Laboratory Systems*, 58:109-130. http://dx.doi.org/10.1016/S0169-7439(01)00155-1

▶ **Trygg J., Holmes E. and Lundstedt T. (2007). Chemometrics in Metabonomics.** *Journal of Proteome Research*, 6:469-479. http://dx.doi.org/10.1021/pr060594q

▶ **Brereton R.G. and Lloyd G.R. (2014). Partial least squares discriminant analysis: taking the magic away.** *Journal of Chemometrics*, 28:213-225.

# Multi-omics analysis and integration

Commissariat à l'énergie atomique et aux énergies alternatives - www.cea.fr
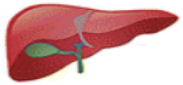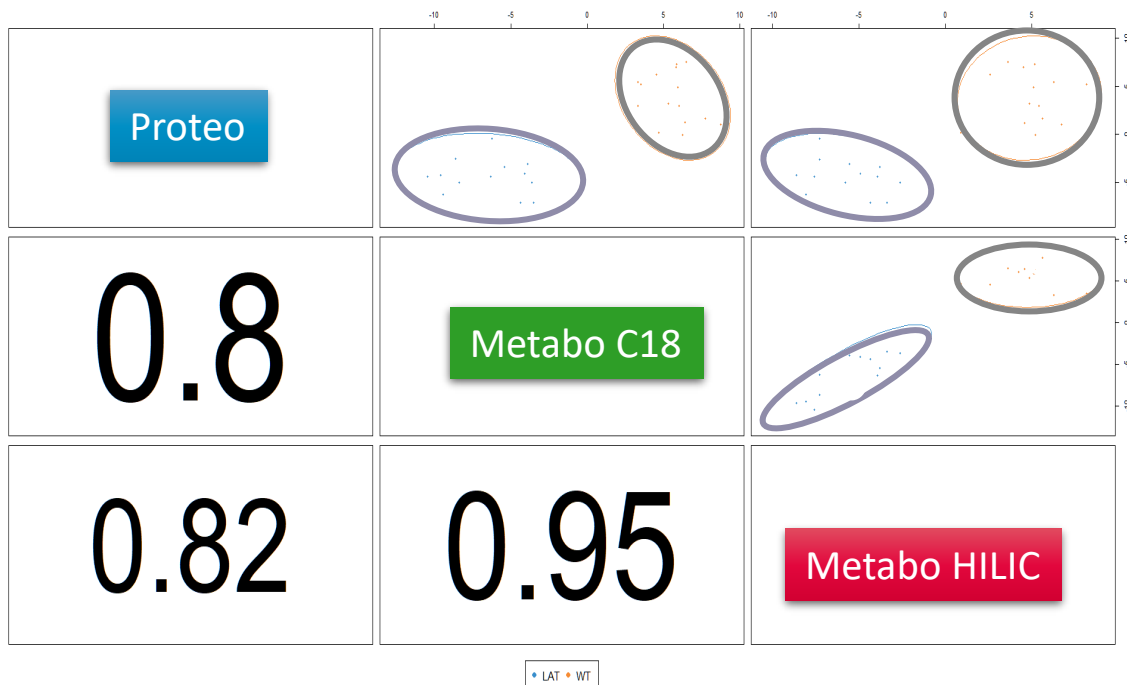
# Significant features KO vs WT (limma)
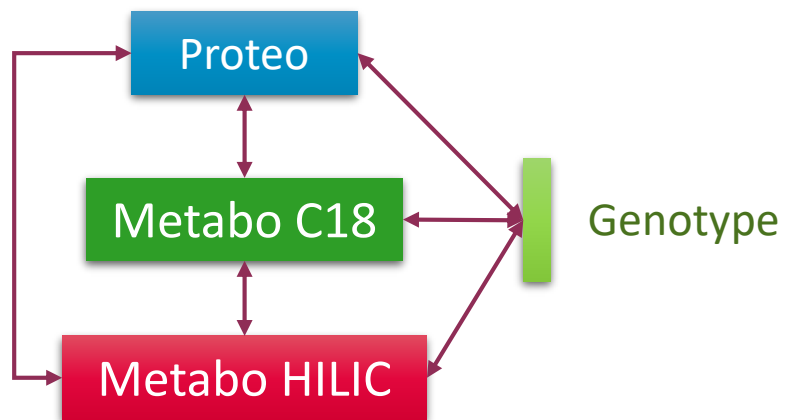
Argelaguet *et al.* (2018). Multi-Omics Factor Analysis—a framework for unsupervised integration of multi-omics data sets. *Molecular Systems Biology*, **14**.
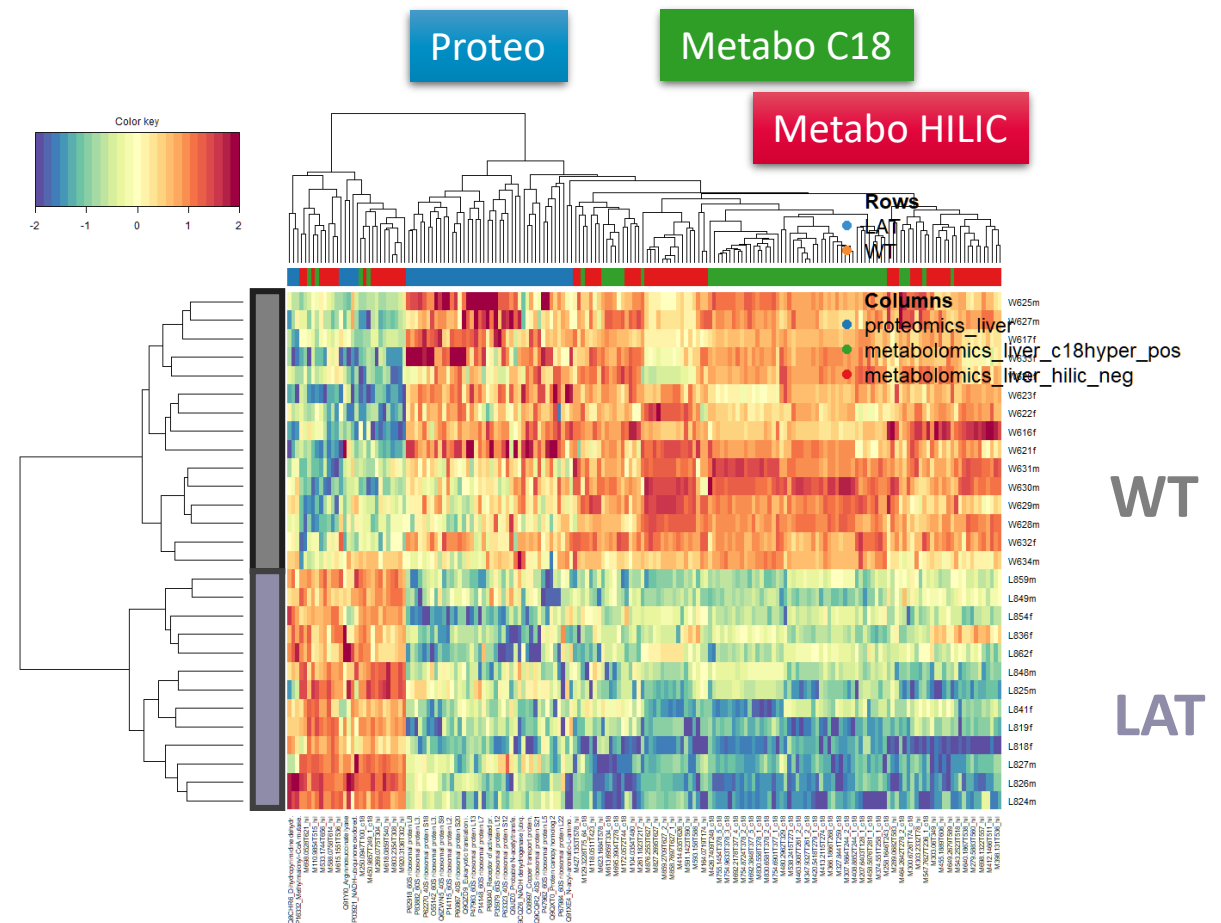
LAT vs WT

LAT **vs** WT

Proteo

Metabo C18

Metabo HILIC

Genotype

Tenenhaus *et al.* (2014). Variable selection for generalized canonical correlation analysis. *Biostatistics*, **15**:569–583.

Singh *et al.* (2019). DIABLO: an integrative approach for identifying key molecular drivers from multi-omics assays. *Bioinformatics*, **35**:3055–3062.



Proteo

Metabo C18

Metabo HILIC

Color key

Rows
LAT
WT

Columns
• proteomics_liver
• metabolomics_liver_c18hyper_pos
• metabolomics_liver_hilic_neg

WT

LAT

Proteo

0.8    Metabo C18

0.82    0.95    Metabo HILIC

• LAT    • WT

▶ **Value of combining proteomics and metabolomics for fundamental and applied research**

▶ **Proteomics and metabolomics data analysis is mature enough to build common pipelines**

▶ **Major challenges remain**

- Limited number of public datasets

- Limited metabolite annotation

- Multidisciplinarity