

How To Save An Awful-Designed Project?

... with the recourse of machine learning approaches for multi-omics analysis💡

Lijiao NING (lijiao.ning@cnr.fr)

Inserm U1283 / CNRS UMR8199 / Université de Lille

6th October 2022

RIS Thematic Day, Paris

There is a project ...

A project with ...

👍 Variety of data

- >300 samples
- >30 phenotypes
- different omics data
 - genotype
 - methylation
 - RNA-seq

Abbreviations

- *NALF*: non-alcoholic fatty liver
- *NASH*: non-alcoholic steatohepatitis
- *BMI*: body mass index
- *IGT*: impaired glucose tolerance
- *T2D*: type 2 diabetes


But confounding factors 😞

Characteristic	Control, N = 80 ¹	NAFL, N = 137 ¹	NASH, N = 83 ¹	p-value ²
Age	34 (26, 43)	43 (32, 51)	46 (40, 56)	<0.001
Sex				0.009
Male	13 (16%)	48 (35%)	28 (34%)	
Female	67 (84%)	89 (65%)	55 (66%)	
BMI	43 (40, 49)	45 (42, 51)	44 (41, 51)	0.2
Diabetic Status				<0.001
Normoglycemic	42 (57%)	17 (12%)	5 (6.1%)	
IGT	16 (22%)	54 (40%)	11 (13%)	
T2D	16 (22%)	65 (48%)	66 (80%)	
Unknown	6	1	1	

¹Median (IQR); n (%)

²Kruskal-Wallis rank sum test; Pearson's Chi-squared test

What can I do?

- ~~Abandon this project and the next please!~~ 
- Select a subset of matched sample and follow the original design
 - Keep all groups → 0 matched
 - Only NAFL and NASH groups → 166 match samples
- Use all available data with machine learning approaches
 - Find phenotype-based clusters
 - Identify multi-omic features related to the clusters

Let's start!



Phenotype-based clustering

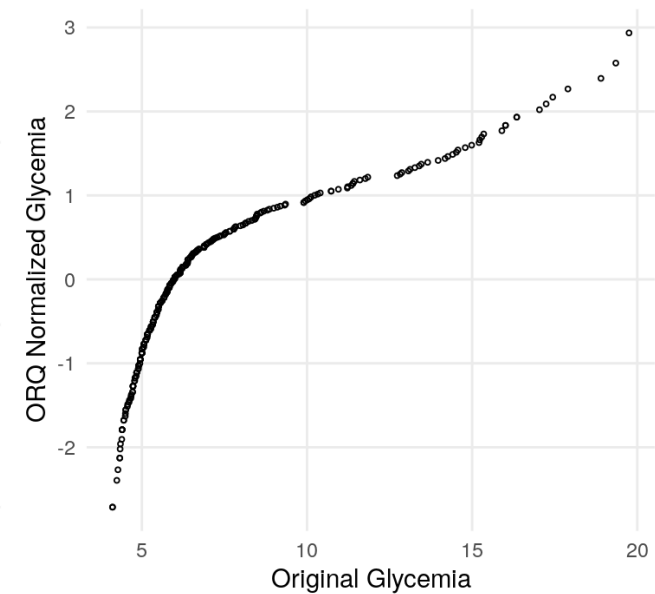
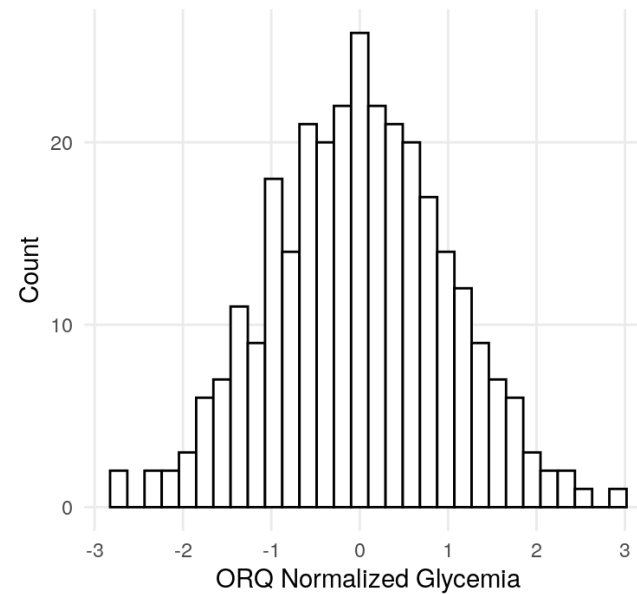
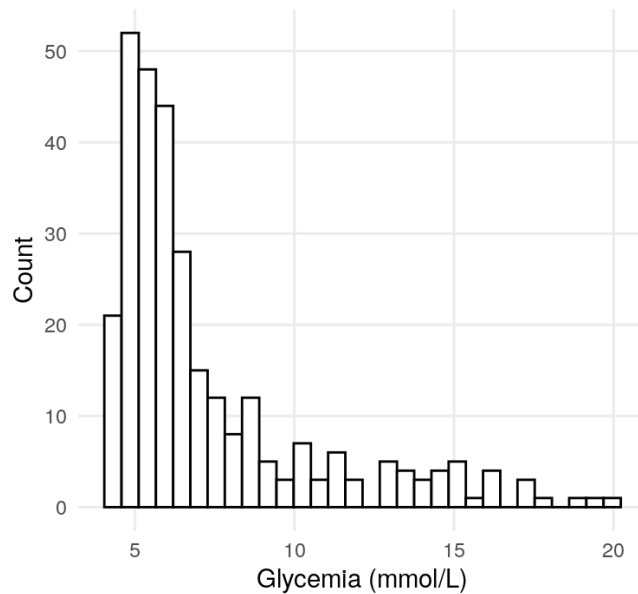
Available phenotypes

- Age, sex, BMI
- Liver biopsy
 - Scores for steatosis, lobular inflammation, ballooning
 - Brunt score, fibrosis score, NAS (nonalcoholic fatty liver disease activity score)
- Blood test
 - Glycemia, insulin, HbA1c, C-peptide
 - Liver function: total bilirubin, transaminases (ALAT, ASAT), gamma-glutamyltransferase (gammaGT)
 - Proteins and lipids: HDL, LDL, alpha2-macroglobulin, haptoglobin, apoA1, CrPus
 - Platelets, lymphocytes

Normalization

- Ordered quantile (ORQ) normalization

with the R-package `bestNormalize` (Ryan Andrew Peterson 2022; Ryan A. Peterson and Cavanaugh 2020)



Imputation

- K-nearest neighbors (kNN) with Gower's distance

Distance between observations i and j :

$$S_{ij} = \frac{\sum_{k=1}^p s_{ijk} \delta_{ijk}}{\sum_{k=1}^p \delta_{ijk}}$$

where $s_{ijk} \in [0, 1]$ represents the similarity between i and j considering the variable k , $\delta_{i,j,k}$ indicates whether i and j can be compared along k .

For continuous variables:

$$s_{ijk} = 1 - \frac{|x_{ik} - x_{jk}|}{r_k}$$

with r_k the range of k .

For categorical variables:

$$s_{ijk} = \begin{cases} 0 & \text{if } x_{ik} = x_{jk}, \\ 1 & \text{if } x_{ik} \neq x_{jk} \end{cases}$$

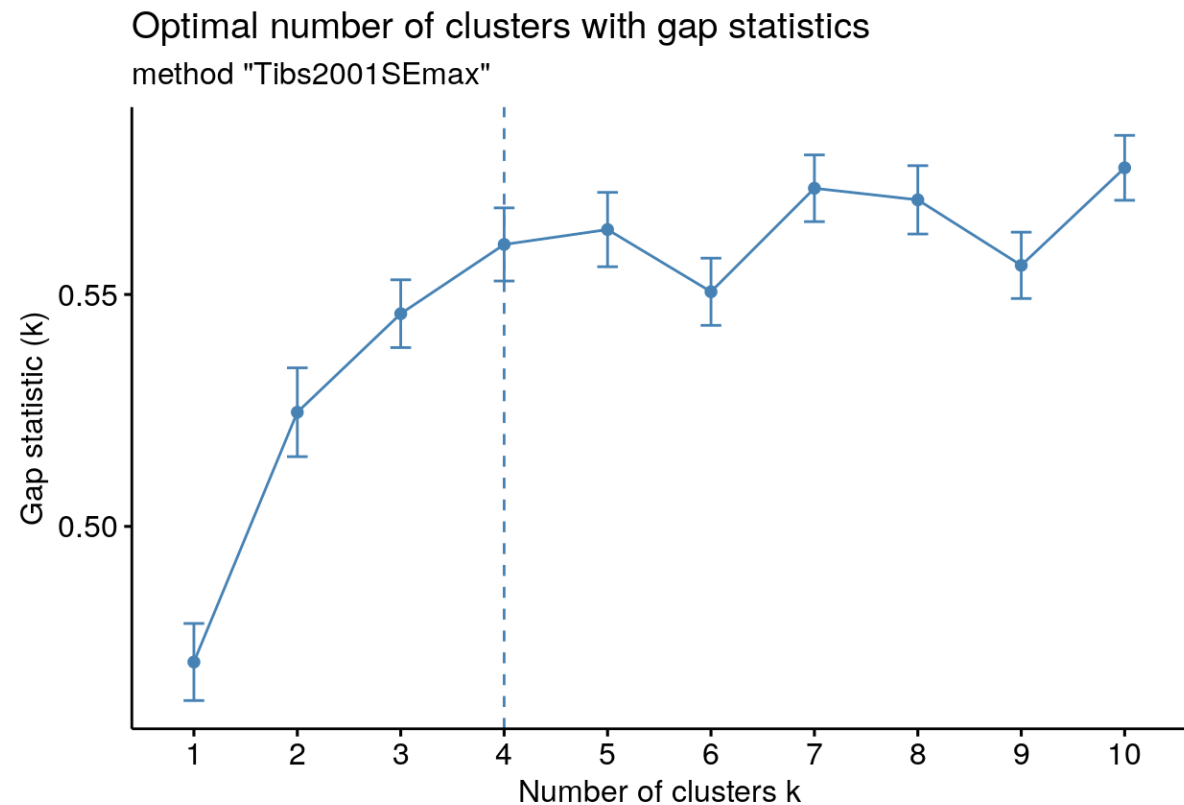
Selection of phenotypes

- Principal component analysis (PCA)
 - Use the elbow method to get the first n components to consider
 - Screen variables which contribute more than the average of the first n components
- Keep 13 out of 35 phenotypes.

K-means clustering

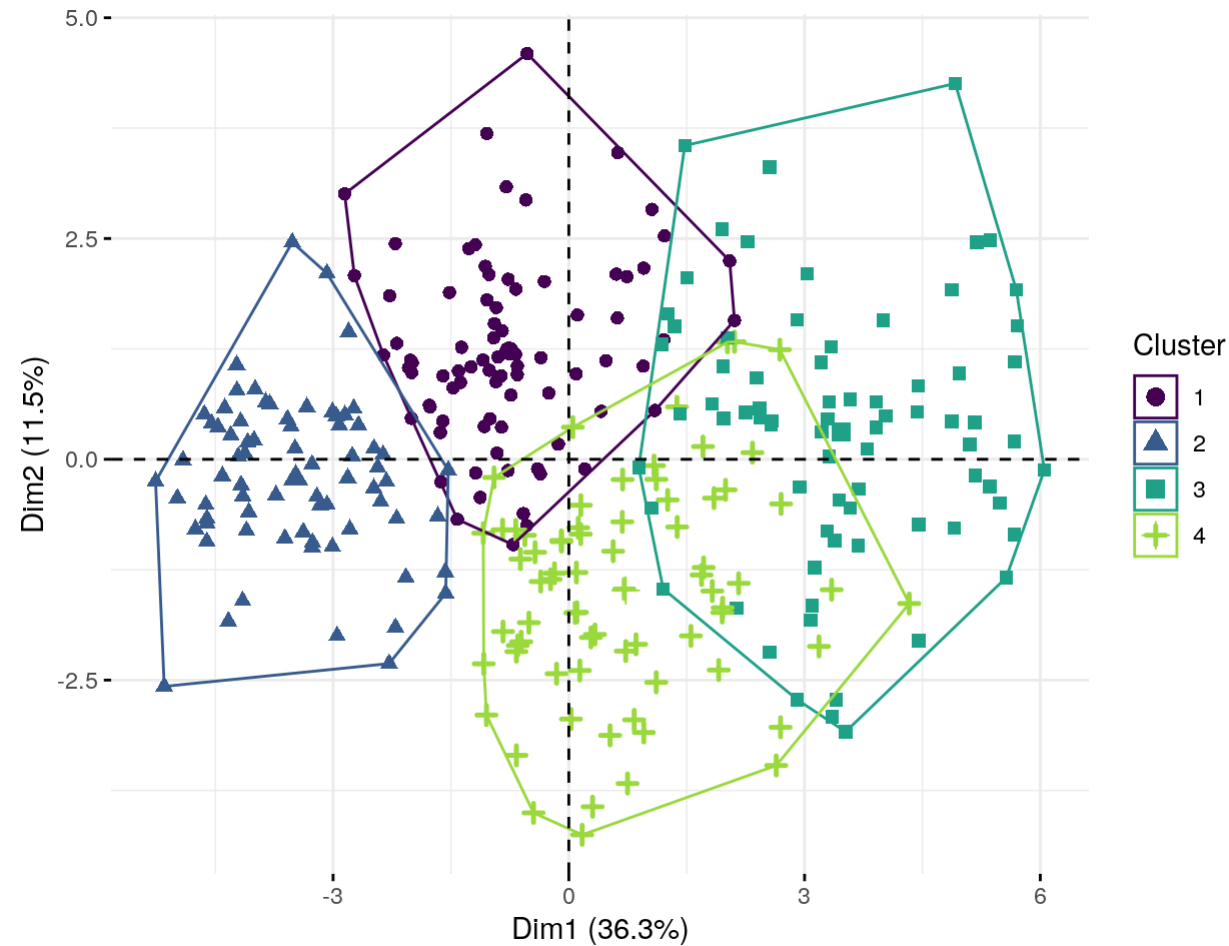
- The optimal number of clusters?

Methods: Elbow method, Silhouette coefficient, Gap statistic, ...

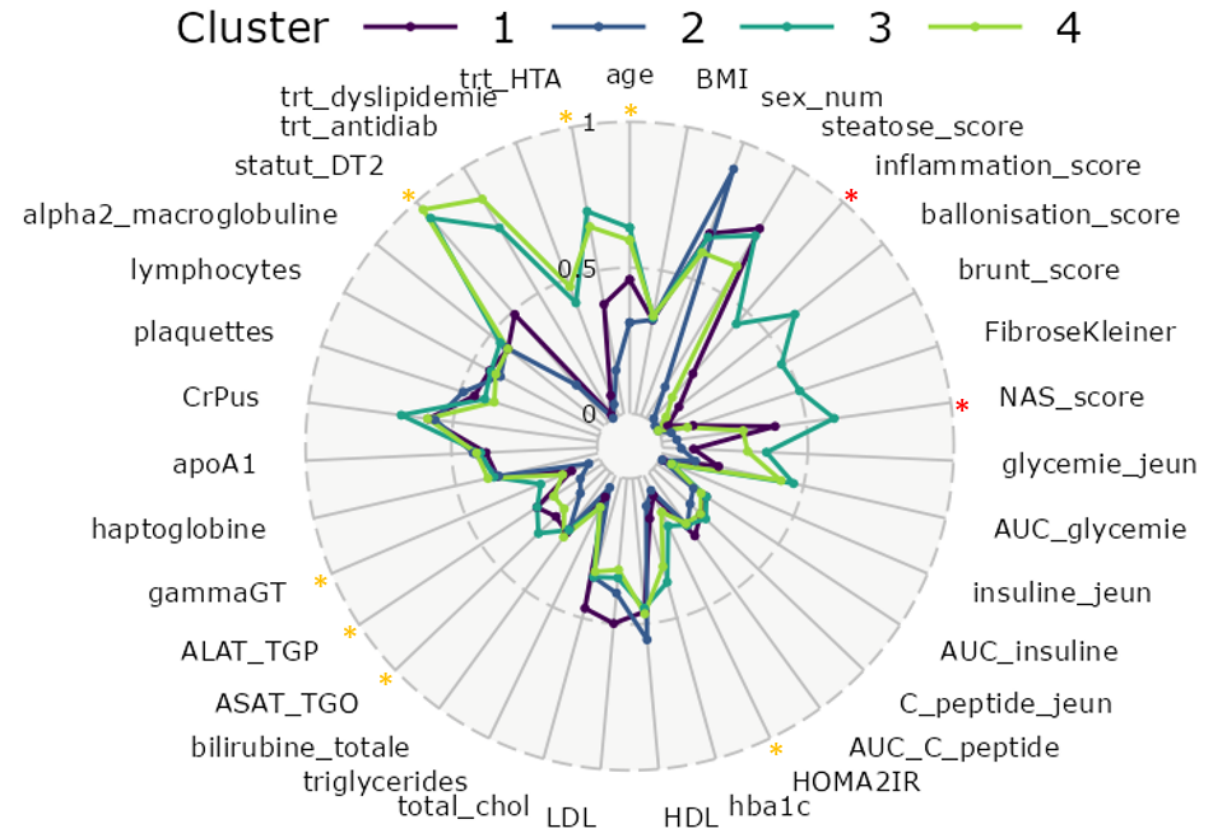


K-means clustering

→ best $k = 4$



Phenotype-based clusters



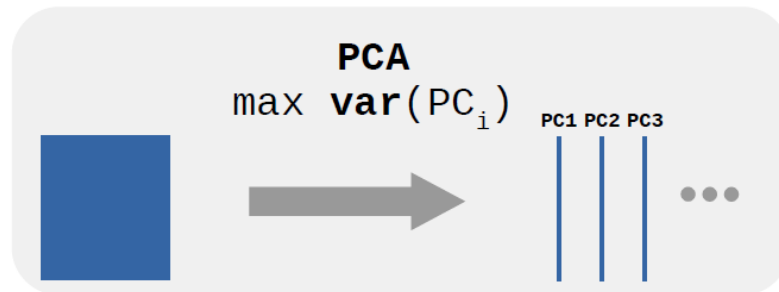
Means' comparison by Tukey HSD method:

* different for all pairs of groups;

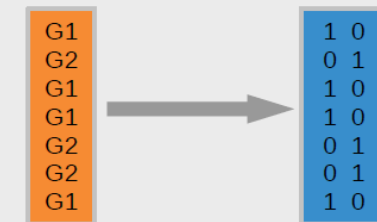
* different for 5 pairs of groups.

Multi-omics analysis

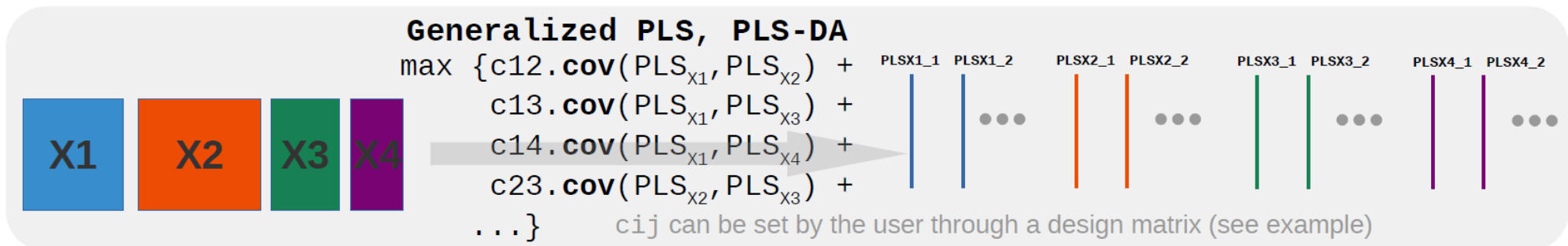
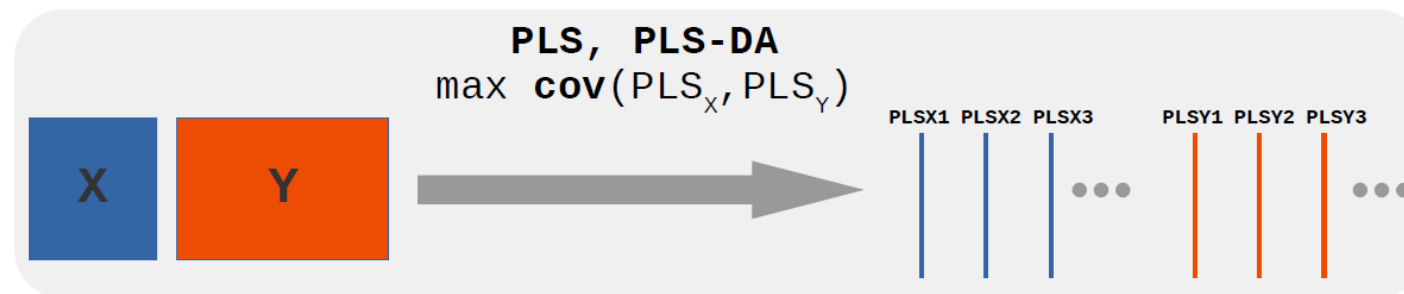
Integration problem



The trick for discriminant analyses: convert a factor into a numeric (dummy) matrix



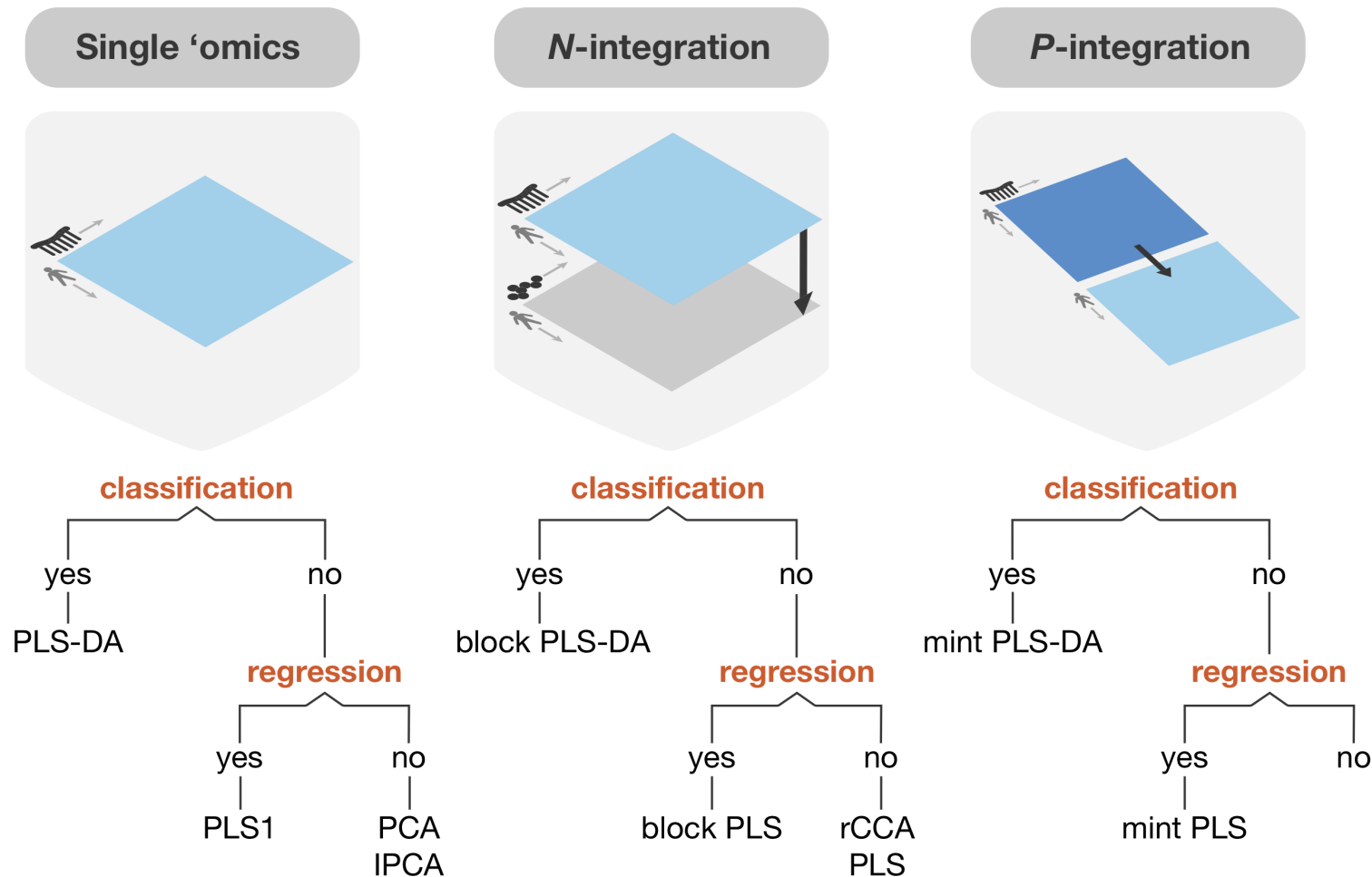
PLS-DA \rightarrow PLS



Sébastien Déjean - www.math.univ-toulouse.fr/~sdejean

R-package mixOmics

Types of analysis with **mixOmics** (Le Cao et al. 2021; F et al. 2017):



N-integration

→ block PLS-DA (the DIABLO framework)

- Principle: latent components are being constructed such that the sum of covariances between all pairs of datasets is maximized and meanwhile discriminating the sample groups.
- Main function `block.splsda`
 - `block`: two or more datasets
 - `s` (sparse), method for feature selection, use cross-validation to tune:
 - number of components to use
 - number of features to keep
 - `pls`: partial least square regression
 - `da`: discriminant analysis (supervised method)

Principal steps (1)

Build the basic model:

```
1 # library(mixOmics)
2
3 omics_list <- list(
4   "PGS" = pgs_data,      # polygenetic scores
5   "mehtyl" = methyl_data,
6   "rnaseq" = rnaseq_data
7 )
8
9 design <- matrix(
10   0.1,
11   ncol = length(omics_list),
12   nrow = length(omics_list),
13   dimnames = list(names(omics_list), names(omics_list))
14 )
15 diag(design) <- 0
16
17 basic_diablo <- block.splsda(
18   X = omics_list,        # list of data matrices
19   Y = clusters,          # categorical outcome
20   scale = TRUE,
21   design = design,       # correlation matrix between dataset
22   ncomp = 5              # an arbitrarily high number
23 )
```

Principal steps (2)

Tune parameters:

```
1 ## Number of components to keep ----
2 perf_diablo <- perf(
3   object = basic_diablo,
4   validation = "Mfold",
5   folds = 10,
6   nrepeat = 10
7 )
8 ncomp <- perf_diablo$choice.ncomp$WeightedVote["Overall.BER", "centroids.dist"]
9
10 ## Number of features to keep ----
11 test_keepX <- list(
12   "PGS" = ceiling(seq(5, ncol(pgs_data), length = 20)),
13   "methylation" = ceiling(seq(5, ncol(methyl_data), length = 20)),
14   "RNAseq" = ceiling(seq(5, ncol(rnaseq_data), length = 20))
15 )
16 tune_diablo <- tune.block.splsda(
17   X = omics_list, Y = clusters,
18   ncomp = ncomp, test.keepX = test_keepX, design = design,
19   validation = "Mfold", folds = 3, nrepeat = 5,
20   dist = "centroids.dist"
21 )
22 list_keepX <- tune_diablo$choice.keepX
```

Principal steps (3)

Build final model:

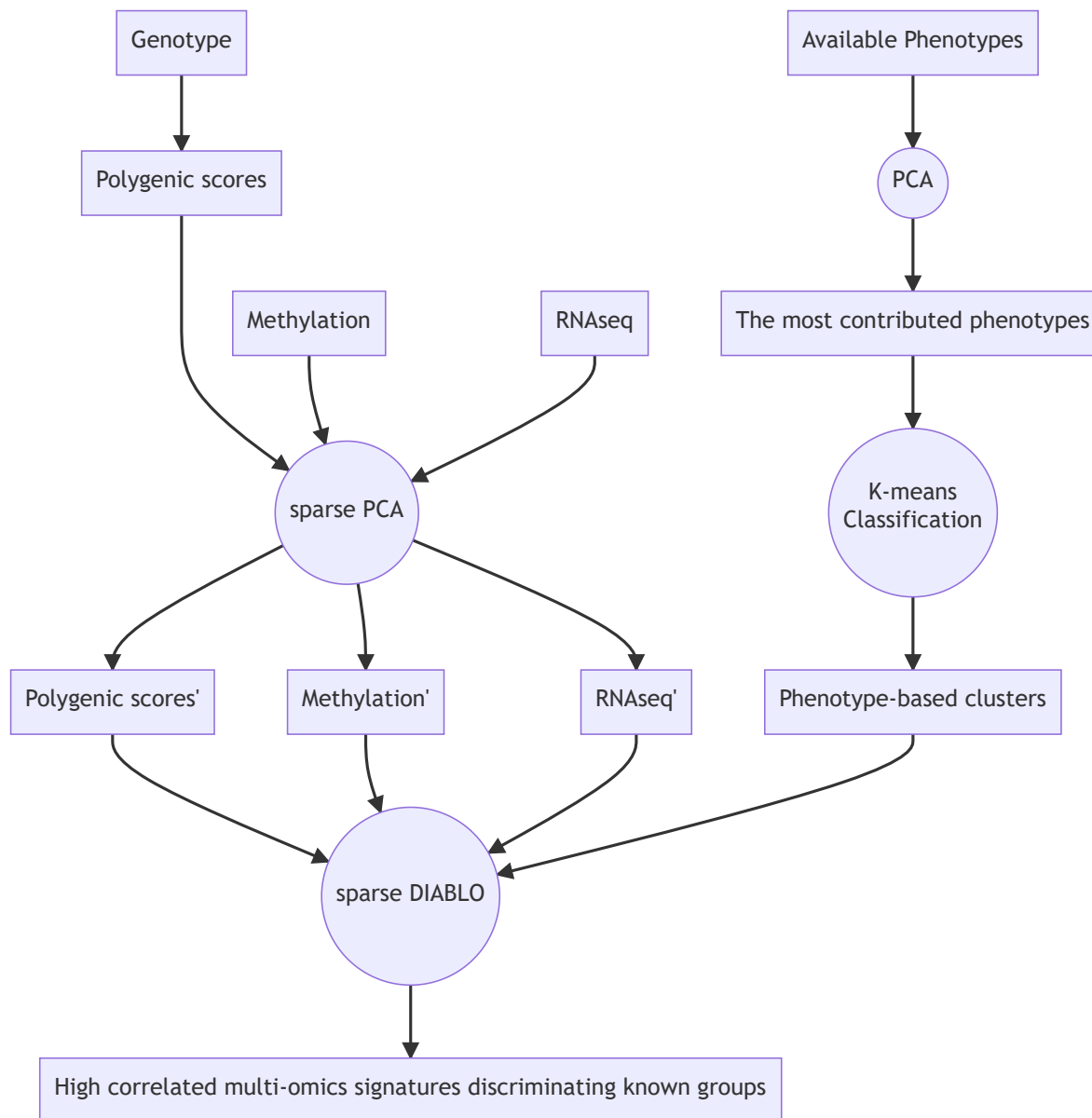
```
1 final_diablo <- block.splsda(  
2   X = omics_list,  
3   Y = clusters,  
4   ncomp = ncomp,  
5   keepX = list_keepX,  
6   design = design  
7 )
```

Predict new samples:

```
1 pred_diablo <- predict(  
2   object = final_diablo,  
3   newdata = list(                                # allow missing blocks, but not missing features  
4     "PGS" = new_pgs_data,  
5     "methyl" = new_methyl_data,  
6     "rnaseq" = new_rnaseq_data  
7   )  
8 )  
9  
10 get.confusion_matrix(  
11   truth = new_sample_cluster,  
12   predicted = pred_diablo$WeightedVote$centroids.dist[, 2]  
13 )
```

Workflow overview





• Number of features

- Begin: 84 PGS, ≈ 740 CpGs, $\approx 14k$ genes
 - sPCA preselection: 44 PGS, 800 CpGs, 180 genes
 - Final selected: 24 PGS, 800 CpGs, 29 genes
- 40 hours + 108 CPUs + max 200Gb

Example of results

Selected features

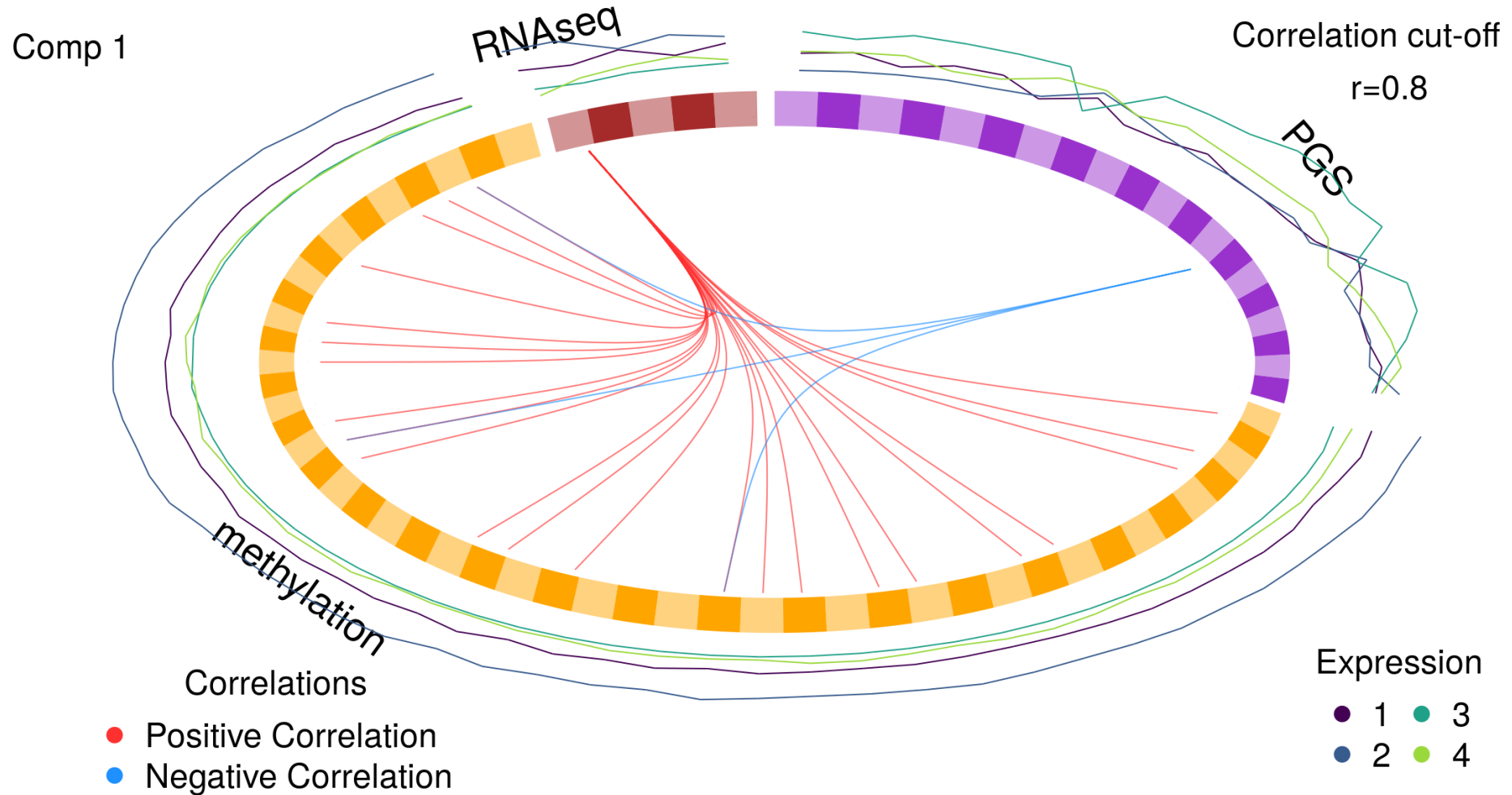
With the two first components, total features selected:

- 24 PGS
- 800 CpGs
- 29 genes (RNAseq)

Features selected in the first component:

- 20 PGS
- 47 CpGs
- 5 genes (RNAseq)

Correlation between features



Selected features

Features selected by the first component:

- 20 PGS
- 47 CpGs
- 5 genes (RNAseq)

Features that highly correlate ($r > 0.8$) with multiple features from other blocks:

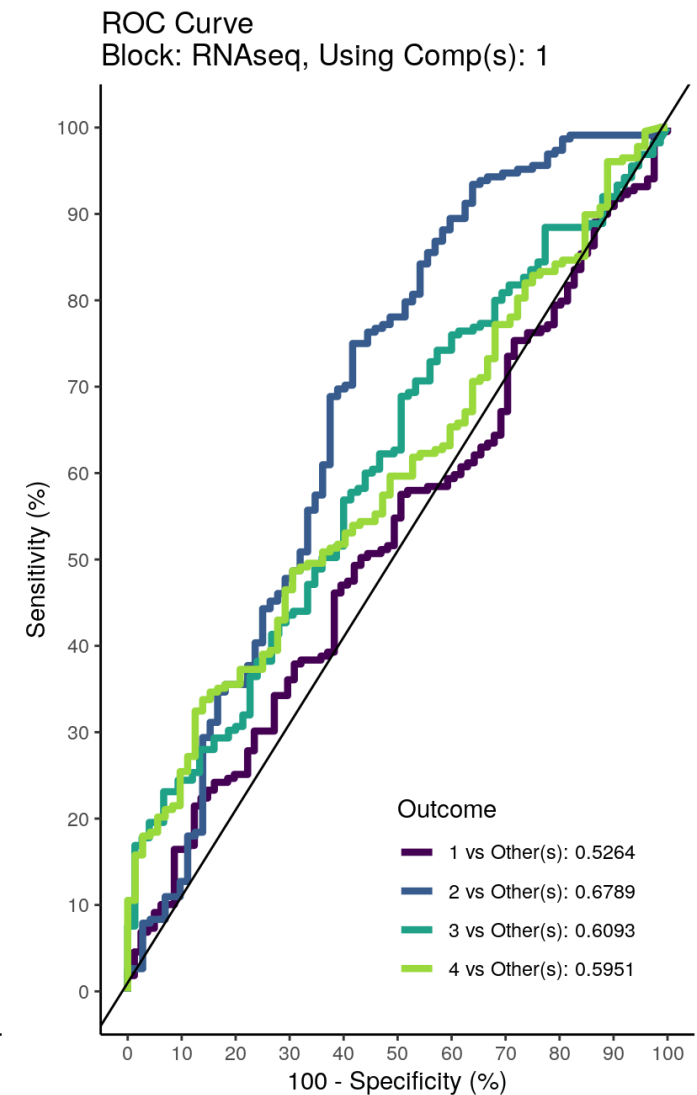
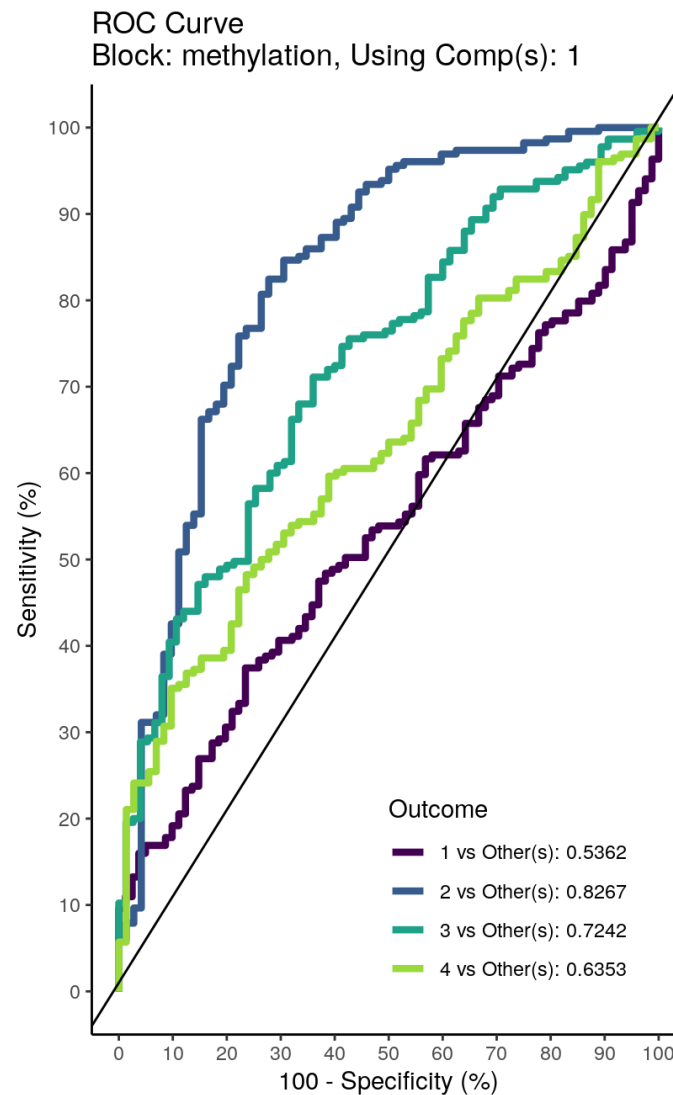
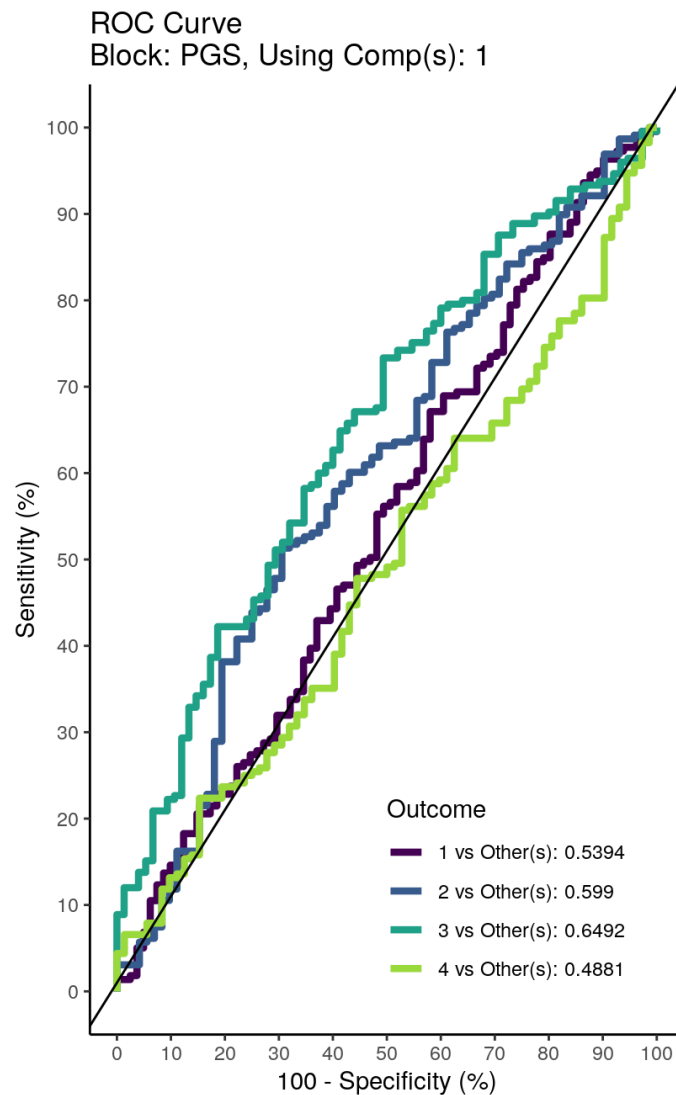
- ENSG00000111700.13 (*SLCO1B3*)

This gene encodes a liver cell specific protein (OATP1B3), which helps to clear compounds from the body.

- PGS002308

A trans-ancestry PGS built for the type 2 diabetes. (Ge, T., Irvin, M.R., Patki, A. *et al.* Genome Med (2022), doi: 10.1186/s13073-022-01074-2)

Area under the ROC



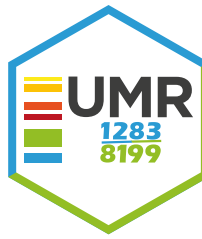
What more?

- Explore in depth the high correlated features & Refine parameter tuning
- Test to predict group for new samples
- There are more solutions than problems 💡

Thank you for your attention!

Any questions?

References



- F, Rohart, Gautier B, Singh A, and Le Cao K-A. 2017. "mixOmics: An r Package for 'Omics Feature Selection and Multiple Data Integration." *PLoS Computational Biology* 13 (11): e1005752.
<http://www.mixOmics.org>.
- Le Cao, Kim-Anh, Florian Rohart, Ignacio Gonzalez, and Sebastien Dejean. 2021. *mixOmics: Omics Data Integration Project*.
<http://www.mixOmics.org>.
- Peterson, Ryan A., and Joseph E. Cavanaugh. 2020. "Ordered Quantile Normalization: A Semiparametric Transformation Built for the Cross-Validation Era." *Journal of Applied Statistics* 47 (13-15): 2312–27.
<https://doi.org/10.1080/02664763.2019.1630372>.
- Peterson, Ryan Andrew. 2022. *bestNormalize: Normalizing Transformation Functions*. <https://CRAN.R-project.org/package=bestNormalize>.