



Minerva Access is the Institutional Repository of The University of Melbourne

Author/s:

Rohart, F;Gautier, B;Singh, A;Lê Cao, KA

Title:

mixOmics: An R package for 'omics feature selection and multiple data integration

Date:

2017-11-01

Citation:

Rohart, F., Gautier, B., Singh, A. & Lê Cao, K. A. (2017). mixOmics: An R package for 'omics feature selection and multiple data integration. *Plos Computational Biology*, 13 (11), <https://doi.org/10.1371/journal.pcbi.1005752>.

Persistent Link:

<https://hdl.handle.net/11343/257448>

License:

CC BY

RESEARCH ARTICLE

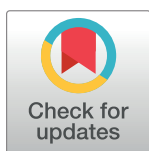
# mixOmics: An R package for 'omics feature selection and multiple data integration

Florian Rohart<sup>1‡</sup>, Benoît Gautier<sup>1</sup>, Amrit Singh<sup>2,3</sup>, Kim-Anh Lê Cao<sup>1,4\*</sup>

**1** The University of Queensland Diamantina Institute, Translational Research Institute, Brisbane, Queensland, Australia, **2** Prevention of Organ Failure (PROOF) Centre of Excellence, Vancouver, British Columbia, Canada, **3** Department of Pathology and Laboratory Medicine, University of British Columbia, Vancouver, British Columbia, Canada, **4** Melbourne Integrative Genomics and School of Mathematics and Statistics, University of Melbourne, Melbourne, Victoria, Australia

‡ Current address: Institute for Molecular Bioscience, The University of Queensland, Brisbane, Queensland, Australia

\* [kimanh.lecao@unimelb.edu.au](mailto:kimanh.lecao@unimelb.edu.au)



**OPEN ACCESS**

**Citation:** Rohart F, Gautier B, Singh A, Lê Cao K-A (2017) mixOmics: An R package for 'omics feature selection and multiple data integration. *PLoS Comput Biol* 13(11): e1005752. <https://doi.org/10.1371/journal.pcbi.1005752>

**Editor:** Dina Schneidman, Hebrew University of Jerusalem, ISRAEL

**Received:** May 5, 2017

**Accepted:** August 31, 2017

**Published:** November 3, 2017

**Copyright:** © 2017 Rohart et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** All relevant data are within the paper and its Supporting Information files.

**Funding:** FR was supported, in part, by the Australian Cancer Research Foundation (ACRF) for the Diamantina Individualised Oncology Care Centre at The University of Queensland Diamantina Institute. KALC was supported, in part, by the National Health and Medical Research Council (NHMRC) Career Development fellowship (APP1087415). The funders had no role in study

## Abstract

The advent of high throughput technologies has led to a wealth of publicly available 'omics data coming from different sources, such as transcriptomics, proteomics, metabolomics. Combining such large-scale biological data sets can lead to the discovery of important biological insights, provided that relevant information can be extracted in a holistic manner. Current statistical approaches have been focusing on identifying small subsets of molecules (a 'molecular signature') to explain or predict biological conditions, but mainly for a single type of 'omics. In addition, commonly used methods are univariate and consider each biological feature independently. We introduce *mixOmics*, an R package dedicated to the multivariate analysis of biological data sets with a specific focus on data exploration, dimension reduction and visualisation. By adopting a systems biology approach, the toolkit provides a wide range of methods that statistically integrate several data sets at once to probe relationships between heterogeneous 'omics data sets. Our recent methods extend Projection to Latent Structure (PLS) models for discriminant analysis, for data integration across multiple 'omics data or across independent studies, and for the identification of molecular signatures. We illustrate our latest *mixOmics* integrative frameworks for the multivariate analyses of 'omics data available from the package.

This is a *PLOS Computational Biology* Software paper.

## Introduction

The advent of novel 'omics technologies (e.g. transcriptomics for the study of transcripts, proteomics for proteins, metabolomics for metabolites, etc) has enabled new opportunities for

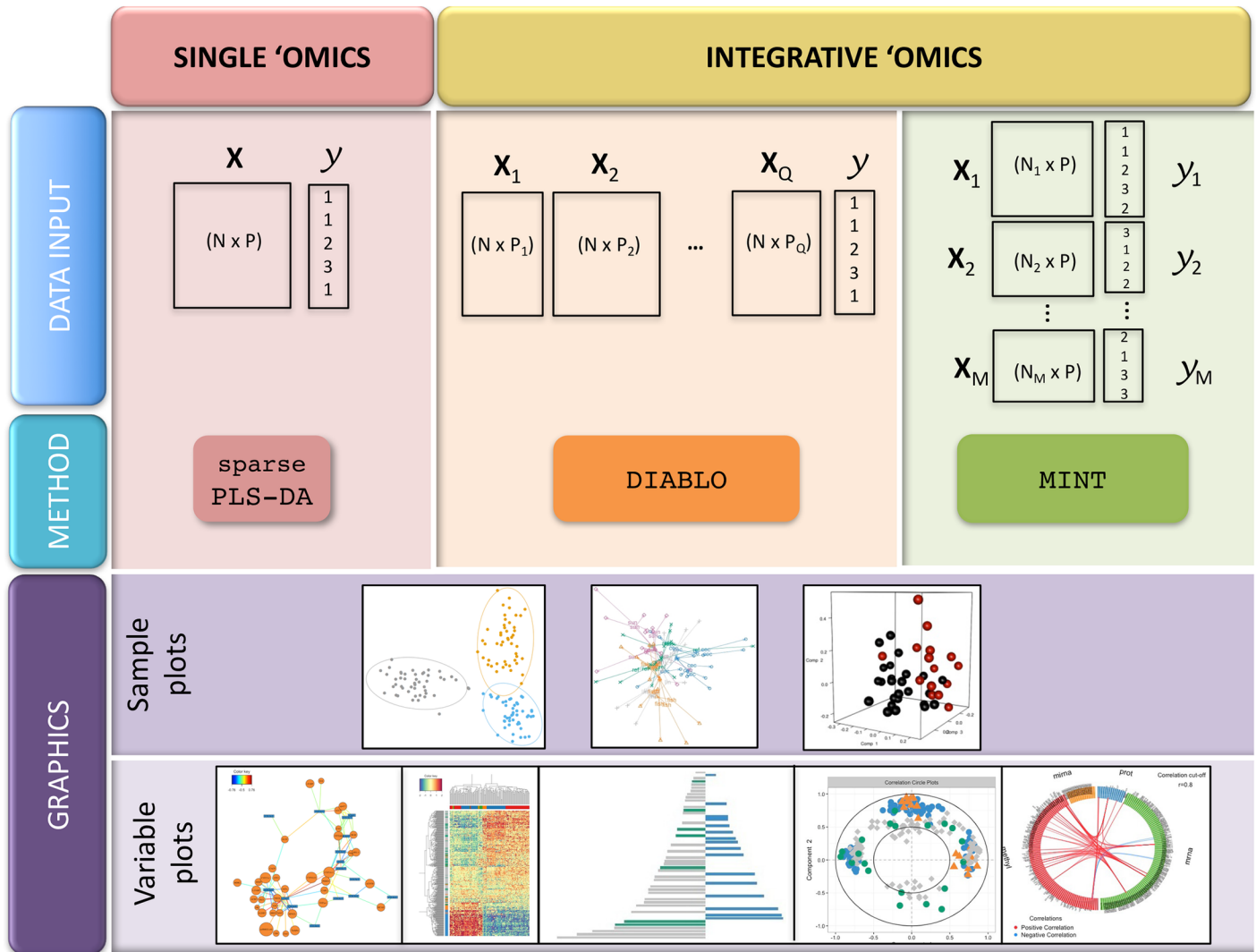
design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors declare that they have no competing interests.

biological and medical research discoveries. Commonly, each feature from each technology (transcripts, proteins, metabolites, etc) is analysed independently through univariate statistical methods including ANOVA, linear models or t-tests. However, such analysis ignores relationships between the different features and may miss crucial biological information. Indeed, biological features act in concert to modulate and influence biological systems and signalling pathways. Multivariate approaches, which model features as a set, can therefore provide a more insightful picture of a biological system, and complement the results obtained from univariate methods. Our package `mixOmics` proposes multivariate projection-based methodologies for 'omics data analysis as those provide several attractive properties to the data analyst [1]. Firstly, they are computationally efficient to handle large data sets, where the number of biological features (usually thousands) is much larger than the number of samples (usually less than 50). Secondly, they perform dimension reduction by projecting the data into a smaller subspace while capturing and highlighting the largest sources of variation from the data, resulting in powerful visualisation of the biological system under study. Lastly, their relaxed assumptions about data distribution make them highly flexible to answer topical questions across numerous biology-related fields [2, 3]. `mixOmics` multivariate methods have been successfully applied to statistically integrate data sets generated from different biological sources, and to identify biomarkers in 'omics studies such as metabolomics, brain imaging and microbiome [4–9].

We introduce `mixOmics` in the context of *supervised analysis*, where the aims are to classify or discriminate sample groups, to identify the most discriminant subset of biological features, and to predict the class of new samples. We further extended our core method sparse Partial Least Square—Discriminant Analysis (sPLS-DA [10]) that was originally developed for the supervised analysis of one data set. Our two novel frameworks `DIABLO` and `MINT` focus on the integration of multiple data sets for different biological questions (Fig 1). `DIABLO` enables the integration of the same biological  $N$  samples measured on different 'omics platforms ( $N$ -integration, [11]), while `MINT` enables the integration of several independent data sets or studies measured on the same  $P$  predictors ( $P$ -integration, [12]). To date, very few statistical methods can perform  $N$ - and  $P$ -integration in a supervised context. For instance,  $N$ -integration is often performed by concatenating all the different 'omics data sets [13], which ignores the heterogeneity between 'omics platforms and mainly highlights one single type of 'omics. The other common type of  $N$ -integration is to combine the molecular signatures identified from separate analyses of each 'omics [14], which disregards the relationships between the different 'omics functional levels. With  $P$ -integration, statistical methods are often sequentially combined to accommodate or correct for technical differences ('batch effects') among studies before classifying samples with a suitable classification method. Such sequential approaches are time consuming and are prone to overfitting when predicting the class of new samples [12]. Our two frameworks model relationships between different types of 'omics data ( $N$ -integration) or integrate independent 'omics studies to increase sample size and statistical power ( $P$ -integration). Both frameworks aim at identifying biologically relevant and robust molecular signatures to suggest novel biological hypotheses.

The present article first introduces the main functionalities of `mixOmics`, then presents our multivariate frameworks for the identification of molecular signatures in one and several data sets, and illustrates each framework in a case study available from the package. The data sets supporting the results of this article are available from the `mixOmics` R package in a processed format. Sweave, R scripts, full tutorials and reports to reproduce the results from the proposed frameworks are available in the supporting information as well as from our website [www.mixOmics.org](http://www.mixOmics.org).



**Fig 1. Overview of the mixOmics multivariate methods for single and integrative 'omics supervised analyses.**  $X$  denote a predictor 'omics data set, and  $y$  a categorical outcome response (e.g. healthy vs. sick). Integrative analyses include  $N$ -integration with DIABLO (the same  $N$  samples are measured on different 'omics platforms), and  $P$ -integration with MINT (the same  $P$  'omics predictors are measured in several independent studies). Sample plots depicted here use the mixOmics functions (from left to right) `plotIndiv`, `plotArrow` and `plotIndiv` in 3D; variable plots use the mixOmics functions `network`, `cim`, `plotLoadings`, `plotVar` and `circosPlot`. The graphical output functions are detailed in Supporting Information S1 Text.

<https://doi.org/10.1371/journal.pcbi.1005752.g001>

## Design and implementation

mixOmics is a user-friendly R package dedicated to the exploration, mining, integration and visualisation of large data sets [1]. It provides attractive functionalities such as (i) insightful visualisations with dimension reduction (Fig 1), (ii) identification of molecular signatures and (iii) improved usage with common calls to all visualisation and performance assessment methods (Supporting Information S1 Text).

## Data input

Different types of biological data can be explored and integrated with mixOmics. Prior to the analysis, we assume the data sets have been normalised using appropriate techniques specific

**Table 1. Summary of the eighteen multivariate projection-based methods available in mixOmics version 6.0.0 or above for different types of analysis frameworks.** Note that our `block.pls/plsda` and sparse variants differ from the approaches from [28–31]. The wrappers for `rgcca` and `sgcca` are originally from the `RGCCA` package [32] but the argument inputs were further improved for `mixOmics`.

Framework		Sparse	Function name	Predictive model
Single 'omics	unsupervised	-	<code>pca</code>	-
		-	<code>ipca</code>	-
		✓	<code>spca</code>	-
	supervised	-	<code>plsda</code>	✓
		✓	<code>splsda</code>	✓
Two 'omics	unsupervised	-	<code>rcca</code>	-
		-	<code>pls</code>	✓
		✓	<code>spls</code>	✓
<i>N</i> -integration	unsupervised	-	<code>wrapper.rgcca</code>	-
		✓	<code>wrapper.sgcca</code>	-
		-	<code>block.pls</code>	✓
		✓	<code>block.spls</code>	✓
	supervised	-	<code>block.plsda</code>	✓
		✓	<code>block.splsda (DIABLO)</code>	✓
<i>P</i> -integration	unsupervised	-	<code>mint.pls</code>	✓
		✓	<code>mint.spls</code>	✓
	supervised	-	<code>mint.plsda</code>	✓
		✓	<code>mint.splsda</code>	✓

<https://doi.org/10.1371/journal.pcbi.1005752.t001>

for the type of 'omics technology platform. The methods can handle molecular features measured on a continuous scale (e.g. microarray, mass spectrometry-based proteomics and metabolomics) or sequenced-based count data (RNA-seq, 16S, shotgun metagenomics) that become “continuous” data after pre-processing and normalisation.

We denote  $X$  a data matrix of size  $N$  observations (rows)  $\times$   $P$  predictors (e.g. expression levels of  $P$  genes, in columns). The categorical outcome  $y$  (e.g. sick vs healthy) is expressed as a dummy indicator matrix  $Y$ , where each column represents one outcome category and each row indicates the class membership of each sample. Thus,  $Y$  is of size  $N$  observations (rows)  $\times$   $K$  categories outcome (columns), see example in Suppl S1 Text.

While `mixOmics` methods can handle large data sets (several tens of thousands of predictors), we recommend pre-filtering the data to less than 10K predictors per data set, for example by using Median Absolute Deviation [15] for RNA-seq data, by removing consistently low counts in microbiome data sets [16, 17] or by removing near zero variance predictors. Such step aims to lessen the computational time during the parameter tuning process.

## Multivariate projection-based methods

`mixOmics` offers a wide range of multivariate dimension reduction techniques designed to each answer specific biological questions, via unsupervised or supervised analyses. The `mixOmics` functions are listed in Table 1. Unsupervised analyses methods include Principal Component Analysis—based on NonLinear Iterative Partial Least Squares for missing values [18], Independent Component Analysis [19], Partial Least Squares regression—PLS, also known as Projection to Latent Structures [20], multi-group PLS [21], regularised Canonical Correlation Analysis—rCCA [22]) and regularised Generalised Canonical Correlation Analysis—rGCCA based on a PLS algorithm [23]. Supervised analyses methods include PLS-Discriminant Analysis—PLS-DA [24–26], GCC-DA [11] and multi-group PLS-DA [12]. In

addition, `mixOmics` provides novel sparse variants that enable *feature selection*, the identification of key predictors (e.g. genes, proteins, metabolites) that constitute a *molecular signature*. Feature selection is performed via  $\ell_1$  regularisation (LASSO, [27]), which is implemented into each method's statistical criterion to be optimised. For supervised analyses, `mixOmics` provides functions to assist users with the choice of parameters necessary for the feature selection process (see 'Choice of parameters' Section) to discriminate the outcome of interest (e.g. healthy vs. sick, or tumour subtypes, etc.).

All multivariate approaches listed in Table 1 are projection-based methods whereby samples are summarised by  $H$  latent components or scores ( $t_1, \dots, t_H$ ) that are defined as linear combinations of the original predictors. In the combinations ( $t_1, \dots, t_H$ ), the weights of each of the predictors are indicated in the loading vectors  $a_1, \dots, a_H$ . For instance, for the data matrix  $X = (X^1, \dots, X^P)$  we define the first latent component as  $t_1 = Xa_1 = X^1a_1^1 + \dots + X^Pa_1^p$ . Therefore, to each loading vector  $a_h$  corresponds a latent component  $t_h$ , and there are as many pairs ( $t_h, a_h$ ) as the chosen dimension  $H$  in the multivariate model,  $h = 1, \dots, H$ , where  $H \ll P$ . The samples are thus projected into a smaller interpretable space spanned by the  $H$  latent components.

## Implementation

`mixOmics` is currently fully implemented in the R language and exports more than 30 functions to perform statistical analyses, tune the methods parameters and plot insightful visualisations. `mixOmics` mainly depends on the R base packages (`parallel`, `methods`, `grDevices`, `graphics`, `stats`, `utils`) and recommended packages (`MASS`, `lattice`), but also imports functions from other R packages (`igraph`, `rgl`, `ellipse`, `corpcor`, `RColorBrewer`, `plyr`, `dplyr`, `tidyr`, `reshape2`, `ggplot2`, `matrixStats`, `rARPACK`, `gridExtra`). In `mixOmics`, we provide generic R/S3 functions to assess the performance of the methods (`predict`, `plot`, `print`, `perf`, `auroc`, etc) and to visualise the results as depicted in Fig 1 (`plotIndiv`, `plotArrow`, `plotVar`, `plotLoadings`, etc), see Supporting Information S1 Text for an exhaustive list.

Currently, eighteen multivariate projection-based methods are implemented in `mixOmics` to integrate large biological data sets, amongst which twelve have similar names (`mint`), `(block).(s)pls(da)`, see Table 1. To perform either  $N$ - or  $P$ -integration, we efficiently coded the functions as wrappers of a single main hidden and generic function that is based on our extension of the sGCCA algorithm [33]. The remaining five statistical methods are PCA, sparse PCA, IPCA, rCCA and rGCCA. Each statistical method implemented in `mixOmics` returns a list of essential outputs which are used in our S3 visualisation functions (Supporting Information S1 Text).

`mixOmics` aims to provide insightful and user-friendly graphical outputs to interpret statistical and biological results, some of which (correlation circle plots, relevance networks, clustered image maps) were presented in details in [34]. The function calls are identical for all multivariate methods via the use of R/S3 functions, as we illustrate in the Results Section. `mixOmics` offers various visualisations, including sample plots and variable plots, which are based on latent component scores and loading vectors, respectively (Fig 1). Additional graphical outputs are available in `mixOmics` to illustrate classification performance of multivariate models using the generic function `plot` (see Supporting Information S1 Text).

## Class prediction of new samples

PLS is traditionally a regression model where the response  $Y$  is a matrix of continuous data. To perform classification and prediction, supervised multivariate methods in `mixOmics` extend

PLS by coding the categorical outcome factor as a dummy indicator matrix before being input into our PLS-based approaches. Considering an independent test set or a cross-validation set, the `predict` function calculates *predicted coordinates (scores)* and *predicted dummy variables* for each new observation, from which we obtain the final predicted class via the use of prediction distances (details in Supporting Information [S1 Text](#)).

**Prediction distances.** For each new observation, we predict its coordinates on the set of  $H$  latent components, similarly to a multivariable multivariate model. These predicted coordinates, or scores, are then used to predict each of the  $K$  dummy variables. The predicted class of each new observation is derived by applying a distance to either the  $H$  predicted scores, or the  $K$  predicted dummy variables. We propose distances such as the 'maximum distance', 'Mahalanobis distance' and 'Centroids distance', which are detailed in Supporting Information [S1 Text](#). The maximum distance is applied to the predicted dummy variables and predicts the class category with the maximum dummy value. In single 'omics analyses this distance achieves best accuracy when the predicted values are close to 0 or 1 [10]. The 'Mahalanobis distance' and 'Centroids distance' are distances that are both applied to the predicted scores, and are based on the calculation of a centroid for each outcome category using the  $H$  latent components. These distances are appropriate for complex classification problems where samples should be considered in a multi-dimensional space spanned by the components. The predicted class of a new observation is the class for which the distance between its centroid and the  $H$  predicted scores is minimal, based on either the Euclidean distance ('Centroid distance'), or the 'Mahalanobis distance'. The former assumes a spherical distribution around the centroid whereas the latter is more adapted for ellipsoidal distribution. In practice, we found that those distances, and especially the Mahalanobis distance, were more accurate than the maximum distance for  $N$ -integration. All distances consider the predictions built from all components of the model.

**Visualisation of prediction area.** To visualise the effect of the prediction distance, we propose a graphical output of the prediction area that overlays the sample plot (example in [Fig 2](#) and more details in Supporting Information [S1 Text](#)).

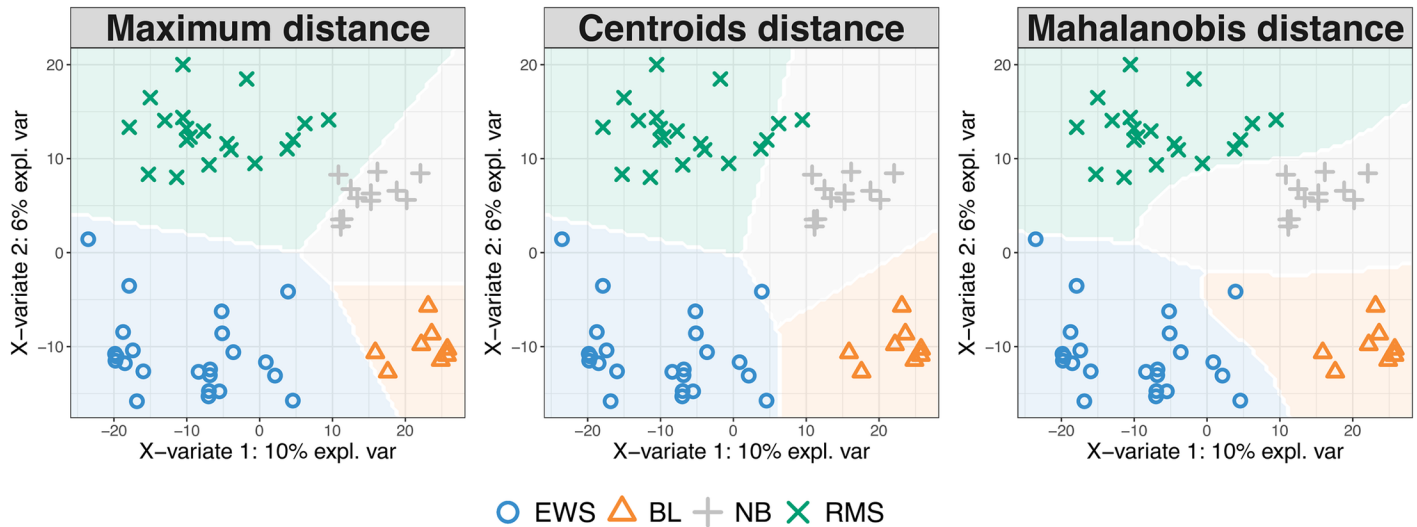
**Prediction for  $N$ -integration.** For  $N$ -integration, we obtain a predicted class *per* 'omics data set. The predictions are combined by majority vote (the class that has been predicted the most often across all data sets) or by weighted vote, where each 'omics data set weight is defined as the correlation between the latent components associated to that particular data set and the outcome, from the training set. The final prediction is the class that obtains the highest weight across all 'omics data sets. Therefore the weighted vote gives more importance to the 'omics data set that is best correlated to the outcome and reduces the number of ties when an even number of data sets are discordant in the case of majority vote. Ties are indicated as NA in our outputs.

**Prediction for  $P$ -integration.** In that specific case, the external test set can include samples from one of the independent studies used to fit the model, or samples from external studies, see [12] for more details.

## Choice of parameters for supervised analyses

For supervised analysis, `mixOmics` provides tools to choose the number of components  $H$  and the  $\ell_1$  penalty on each component for all sparse methods before the final multivariate model is built and the selected features are returned.

**Parameter tuning using cross-validation.** For all supervised models, the function `tune` implements repeated and stratified cross-validation (CV, see details in Supporting Information [S1 Text](#)) to compare the performance of models constructed with different  $\ell_1$  penalties.



**Fig 2. Prediction area visualisation on the Small Round Blue Cell Tumors data (SRBCT [35]) data, described in the Results Section, with respect to the prediction distance.** From left to right: 'maximum distance', 'Centroid distance' and 'Mahalanobis distance'. Sample prediction area plots from a PLS-DA model applied on a microarray data set with the expression levels of 2,308 genes on 63 samples. Samples are classified into four classes: Burkitt Lymphoma (BL), Ewing Sarcoma (EWS), Neuroblastoma (NB), and Rhabdomyosarcoma (RMS).

<https://doi.org/10.1371/journal.pcbi.1005752.g002>

Performance is measured via overall misclassification error rate and Balanced Error Rate (BER). BER is appropriate in case of an unbalanced number of samples per class as it calculates the average proportion of wrongly classified samples in each class, weighted by the number of samples in each class. Therefore, BER is less biased towards majority classes during the performance assessment. The choice of the parameters (described below) is made according to the best prediction accuracy, i.e. the lowest overall error rate or lowest BER.

**Number of components.** For all supervised methods, the tuning function outputs the optimal number of components that achieve the best performance based on the overall error rate or BER. The assessment is data-driven and similar to the process detailed in [36], where one-sided t-tests assess whether there is a gain in performance when adding components to the model. In practice (see some of our examples in the Results Section), we found that  $K - 1$  components, where  $K$  is the number of classes, was sufficient to achieve the best classification performance [10, 37]. However, assessing the performance of a non sparse model with  $K$  to  $K + 2$  components can be used to identify the optimal number of components, see Supporting Information S1 Appendix.

**$\ell_1$  penalty or the number of features to select.** Contrary to other R packages implementing  $\ell_1$  penalisation methods (e.g. `glmnet`, [38], `PMA`, [39]), `mixOmics` uses soft-thresholding to improve usability by replacing the  $\ell^1$  parameter by the number `keepX` of features to select on each dimension. The performance of the model is assessed for each value of `keepX` provided as a grid by the user from the first component to the  $H^{th}$  component, one component at a time. The grid needs to be carefully chosen to achieve a trade-off between resolution and computational time. Firstly, one should consider the minimum and maximum values of the selection size that can be handled practically for follow-up analyses (e.g. wet-lab experiments may require a small signature, gene ontology a large signature). Secondly, one should consider the computational aspect, as the `tune` function performs repeated cross-validation. For single 'omics and P-integration analyses, a coarse tuning grid can be assessed first to evaluate the likely boundaries of the `keepX` values before setting a finer grid. For N-integration, as

different combinations of `keepX` between the different 'omics are assessed, a coarse grid is difficult to achieve as a preliminary step.

The `tune` function returns the set of `keepX` values that achieve the best predictive performance for all the components in the model. In case of ties, the lowest `keepX` value is returned to obtain a minimal molecular signature. The same grid of `keepX` values is used to tune each component; however for N-integration, different grids can be set for each data set. Examples of optimal `keepX` values returned by our functions are detailed in the Results section (see also Supporting Information [S1 Appendix](#)).

**Special case for P- integration.** For P-integration, we take advantage of the independence between studies. A Leave-One-Group-Out Cross Validation is performed where each study defines a subset that is left out once, as described in [12], which substantially improves computational time (see Supporting Information [S1 Text](#) for additional details).

## Evaluating the signature

**Performance assessment.** Once the optimal parameters have been chosen (number of components and number of variables to select), the final model is run on the whole data set  $X$ , and the performance of the final model in terms of classification error rate is estimated using the `perf` function and repeated CV. Additional evaluation outputs include the receiver operating characteristic (ROC) curves and Area Under the Curve (AUC) averaged over the cross-validation process using one-vs-all comparison if  $K > 2$ . AUC is a commonly used measure to evaluate a classifier discriminative ability. It incorporates the measures of sensitivity and specificity for every possible cut-off of the predicted dummy variables. However, as presented in Section 'Prediction distances', our PLS-based models rely on prediction distances, which can be seen as a determined optimal cut-off. Therefore, the ROC and AUC criteria may not be particularly insightful in relation to the performance evaluation of our supervised multivariate methods, but can complement the statistical analysis.

**Stability.** A by-product of the performance evaluation using `perf` is the record of the features that were selected across the (repeated) CV runs. The function `perf` outputs the feature stability per component to assess the reproducibility of the molecular signature (see example in Supporting Information [S1 Appendix](#)).

**Graphical outputs.** Variable plots are useful to assess the correlation of the selected features within and between data sets. Correlation circle plots, clustered image maps, relevant networks are described in Supporting Information [S1 Text](#). A pyramid barplot displays the loading weights associated to each selected feature in increasing order of importance (from bottom to top), with colors indicating the sample group with the maximum or alternatively minimum average value (see Supporting Information [S1 Text](#)).

## Computational aspects

The choice of the parameters via the tuning and the performance evaluation steps can be computationally demanding as the `tune` and `perf` function perform repeated cross-validation. Once the optimal parameters are chosen, the final multivariate models in `mixOmics` are however computationally very efficient to run.

The tuning can be particularly intensive for N-integration as we test all possible combination of subsets of variables to select. For large multi-'omics data sets, the tuning will often require the use of a cluster, while a normal laptop might be sufficient for the single 'omics and P-integration. To lessen the computational issue, the argument `cpus` in both `tune` and `perf` functions is included for parallel computing. [Table 2](#) reports the computational time for the analyses illustrated in the Supporting Information [S1 Appendix](#). The data analysed constitute a

**Table 2. Example of computational time for the data sets presented in the Results section with a macbook pro 2013, 2.6GHz, 16Go Ram.**

Framework	Single 'omics		N-integration		P-integration	
	sPLS-DA		DIABLO		MINT	
Data	srbct		breast.tcga		stemcells	
N	63		150		125	
P	2,308		200;184;142		400	
function	tune	perf	tune	perf	tune	perf
#fold CV (repeated)	5(10)	5(10)	10(1)	10(10)	LOGOCV	LOGOCV
ncomp	6	3	2	2	2	2
grid length per component	39	-	13 <sup>3</sup>	-	100	-
#cpu	1	1	2	1	1	1
run time	9min	31sec	18min	25sec	30sec	0.2sec

<https://doi.org/10.1371/journal.pcbi.1005752.t002>

reduced set of features that are included in the package. Supplemental Information [S1 Text](#) reports the computational time for large data sets analysed with `mixOmics`. For the latter case, we usually recommend to filter the data sets, as detailed in Section 'Data input' for a more tractable analysis.

## Results

We illustrate three supervised frameworks' analyses performed with `mixOmics` using data available from the package. These data sets were reduced to fit the memory allocation storage allowed in R CRAN and the results presented are hence an illustration of the capabilities of our package, but do not necessarily provide insightful biological results. All R scripts are provided in Supporting Information [S1 Appendix](#) and in our website.

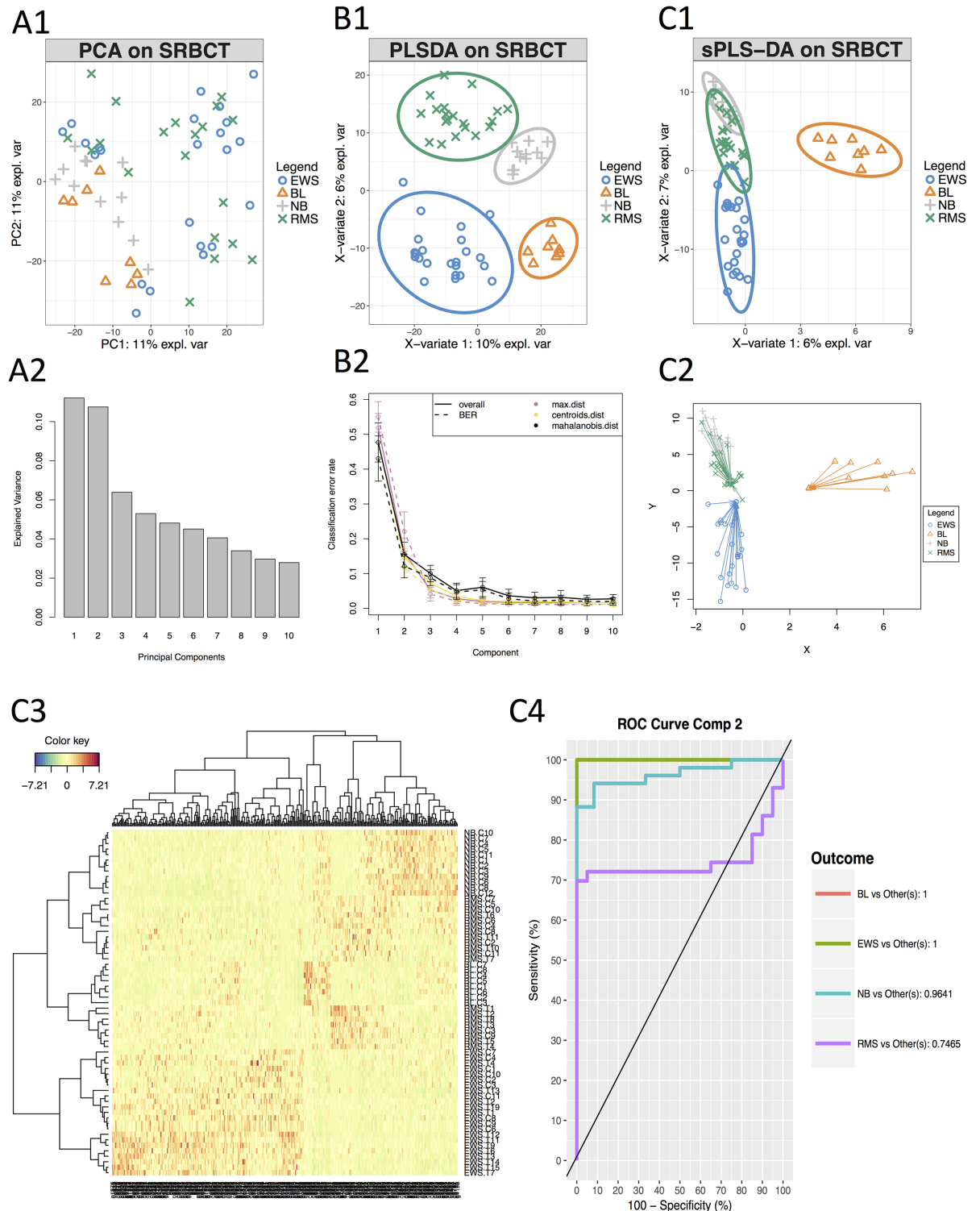
### Single 'omics supervised analyses with PLS-DA and sPLS-DA

We present the application of the single 'omics multivariate methods PCA, PLS-DA and sPLS-DA on a microarray data set. The PLS-DA and sPLS-DA methods are described in the Supporting Information [S1 Text](#).

**Data description.** The study investigates Small Round Blue Cell Tumors (SRBCT, [35]) of 63 tumour samples with the expression levels of 2,308 genes. Samples are classified into four classes: 8 Burkitt Lymphoma (BL), 23 Ewing Sarcoma (EWS), 12 neuroblastoma (NB), and 20 rhabdomyosarcoma (RMS).

**Unsupervised and supervised analyses.** Principal Component Analysis was first applied to assess similarities between tumour types ([Fig 3A1](#)). This preliminary unsupervised analysis showed no separation between tumour types, but allows to visualise the more important sources of variation, which are summarised in the first two principal components ([Fig 3A2](#)). A supervised analysis with PLS-DA focuses on the discrimination of the four tumour types ([Fig 3B1](#)), and led to a good performance ([Fig 3B2](#), performance assessed when adding one component at a time in the model). We then applied sPLS-DA to identify specific discriminant genes for the four tumour types. The tuning process (see 'Choice of parameters' Section and Supporting Information [S1 Appendix](#)) led to a sPLS-DA model with 3 components and a molecular signature composed of 10, 300 and 30 genes selected on the first three components, respectively.

**Results visualisation.** The first sPLS-DA component discriminated BL vs the other tumour types ([Fig 3C1](#)). The 10 genes selected on this component all had positive weight in the linear combination, and were highly expressed in BL. The second component further



**Fig 3. Illustration of a single 'omics analysis with mixOmics. A) Unsupervised preliminary analysis with PCA, A1:** PCA sample plot, **A2:** percentage of explained variance per component. **B) Supervised analysis with PLS-DA, B1:** PLS-DA sample plot with confidence ellipse plots, **B2:** classification performance per component (overall and BER) for three prediction distances using repeated stratified cross-validation (10×5-fold CV). **C) Supervised analysis and feature selection with sparse PLS-DA, C1:** sPLS-DA sample plot with confidence ellipse plots, **C2:** arrow plot representing each sample pointing towards its outcome category, see more details in Supporting Information [S1 Text](#). **C3:** Clustered Image Map (Euclidean Distance, Complete linkage)

where samples are represented in rows and selected features in columns (10, 300 and 30 genes selected on each component respectively), **C4**: ROC curve and AUC averaged using one-vs-all comparisons.

<https://doi.org/10.1371/journal.pcbi.1005752.g003>

discriminated EWS based on 300 selected genes. The genes with a negative weight were highly expressed in EWS while the genes with a positive weight were highly expressed in either NB or RMS. Finally, the third component discriminated both NB and RMS (see Supporting Information [S1 Appendix](#)). The arrow plot displays the relationship between the samples summarised as a combination of selected genes (start of the arrow) and the categorical outcome (tip of the arrow, [Fig 3C2](#)).

A clustering analysis using a heatmap based on the genes selected on the first three components highlighted clusters corresponding to the four tumour types ([Fig 3C3](#)). ROC curve and AUC of the final model were also calculated using one-vs-all comparisons and led to satisfactory results on the first two components ([Fig 3C4](#)). The AUC for the first three components was 1 for all groups. Note that ROC and AUC are additional measures that may not reflect the performance of a `mixOmics` multivariate approaches since our prediction strategy is based on distances (see 'Performance assessment' Section).

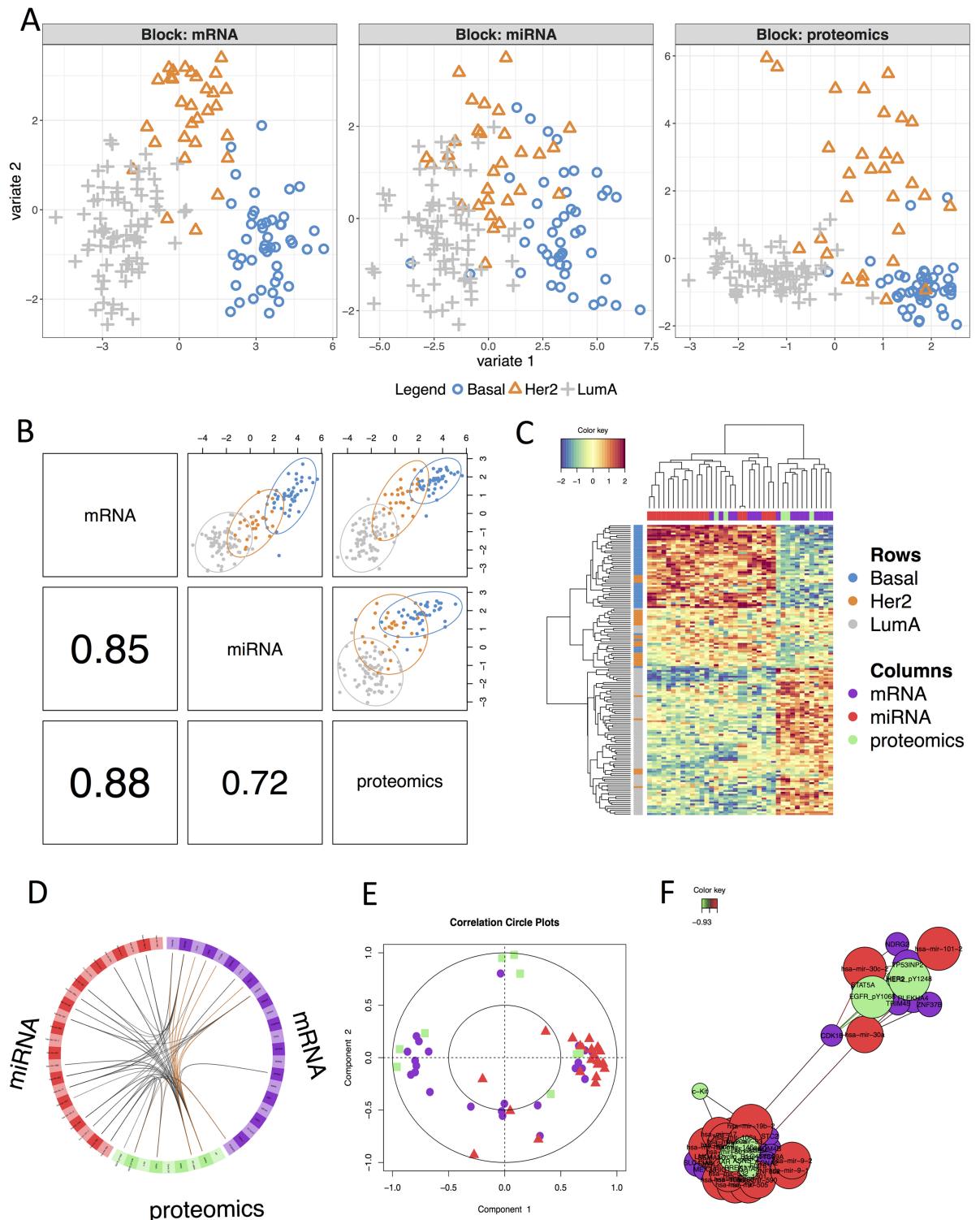
**Summary.** We illustrated the `mixOmics` framework for the supervised analysis of a single 'omics data set—here a microarray experiment. The full pipeline, results interpretation, associated R and Sweave codes are available in Supporting Information [S1 Appendix](#). Such an analysis suggests novel biological hypotheses to be further validated in the laboratory, when one is seeking for a *signature* of a subset of features to explain, discriminate or predict a categorical outcome. The method has been applied and validated in several biological and biomedical studies, including ours in proteomics and microbiome [[17](#), [37](#)].

## N-integration across multiple 'omics data sets with DIABLO

*N*-integration consists in integrating different types of 'omics data measured on the same *N* biological samples. In a supervised context, `DIABLO` performs *N*-integration by identifying a multi-'omics signature that discriminates the outcome of interest. Contrary to the concatenation and the ensemble approaches that also perform *N*-integration, `DIABLO` identifies a signature composed of highly correlated features across the different types of 'omics, by modelling relationships between the 'omics data sets [[11](#)]. The `DIABLO` method is fully described in the Supporting Information [S1 Text](#). We illustrate one analysis on a multi-'omics breast cancer study available from the package.

**Data description.** The multi-'omics breast cancer study includes 150 samples from three types of 'omics: mRNA ( $P_1 = 200$ ), miRNA ( $P_2 = 184$ ) and proteomics ( $P_3 = 142$ ) data. Prior to the analysis with `mixOmics`, the data were normalised and filtered for illustrative purpose. Samples are classified into three subgroups: 75 Luminal A, 30 Her2 and 45 Basal.

**Choice of parameters and analysis.** As we aim to discriminate three breast cancer subtypes we chose a model with 2 components. The tuning process (see 'Choice of parameters for supervised analyses' Section and Supporting Information [S1 Appendix](#)) identified a multi-'omics signature of 16 and 7 mRNA features, 18 and 5 miRNA features and 5 and 5 proteomics features on the first two components, respectively. Sample plots of the final `DIABLO` model in [Fig 4A](#) displayed a better discrimination of breast cancer subgroups with the mRNA and proteomics data than with the miRNA data. [Fig 4B](#) showed that the latent components of each 'omics were highly correlated between each others, highlighting the ability of `DIABLO` to model a good agreement between the data sets. The breast subtypes colors show that the components are also able to discriminate the outcome of interest.



**Fig 4. Illustration of N-integrative supervised analysis with DIABLO.** **A:** sample plot per data set, **B:** sample scatterplot from `plotDiablo` displaying the first component in each data set (upper diagonal plot) and Pearson correlation between each component (lower diagonal plot). **C:** Clustered Image Map (Euclidean distance, Complete linkage) of the multi-omics signature. Samples are represented in rows, selected features on the first component in columns. **D:** Circos plot shows the positive (negative) correlation ( $r > 0.7$ ) between selected features as indicated by the brown (black) links, feature names appear in the quadrants, **E:** Correlation Circle plot representing each type of selected features, **F:** relevance network visualisation of the selected features.

<https://doi.org/10.1371/journal.pcbi.1005752.g004>

**Results visualisation.** Several visualisation tools are available to help the interpretation of the DIABLO results and to assess relationships between the selected multi-'omics features (see Supporting Information [S1 Text](#) and [S1 Appendix](#)). The clustered image map (CIM) displayed a good classification of the three subtypes of breast cancer based on the 39 multi-'omics signature identified on the first component ([Fig 4C](#)). The CIM output can be complemented with a `circosPlot` which displays the different types of selected features on a circle, with links between or within 'omics indicating strong positive or negative correlations ([Fig 4D](#)). Those correlation are estimated using the latent components as a proxy, see more methodological details in [34]. We observed strong correlations between miRNA and mRNA, but only a few correlations between proteomics and the other 'omics types. Correlation circle plots ([Fig 4E](#)) further highlight correlations between each selected feature and its associated latent component (see details in [34]). The 18 miRNA features selected on the first component were highly positively correlated with the first component (red triangles close to the (1,0) coordinates). Contrarily, 9 of the 16 mRNA features and 3 of the 5 proteomics features selected on the first component were highly negatively correlated with the first component (purple circles and green squares close to the (-1,0) coordinates, respectively). Most of the features selected on the second component were close to the inner circle, which implies a weak contribution of those features to both components. Finally, a relevance network output highlighted two clusters, both including features from the three types of 'omics ([Fig 4F](#)). Interactive view and `.glm` format are also available, see Supporting Information [S1 Text](#).

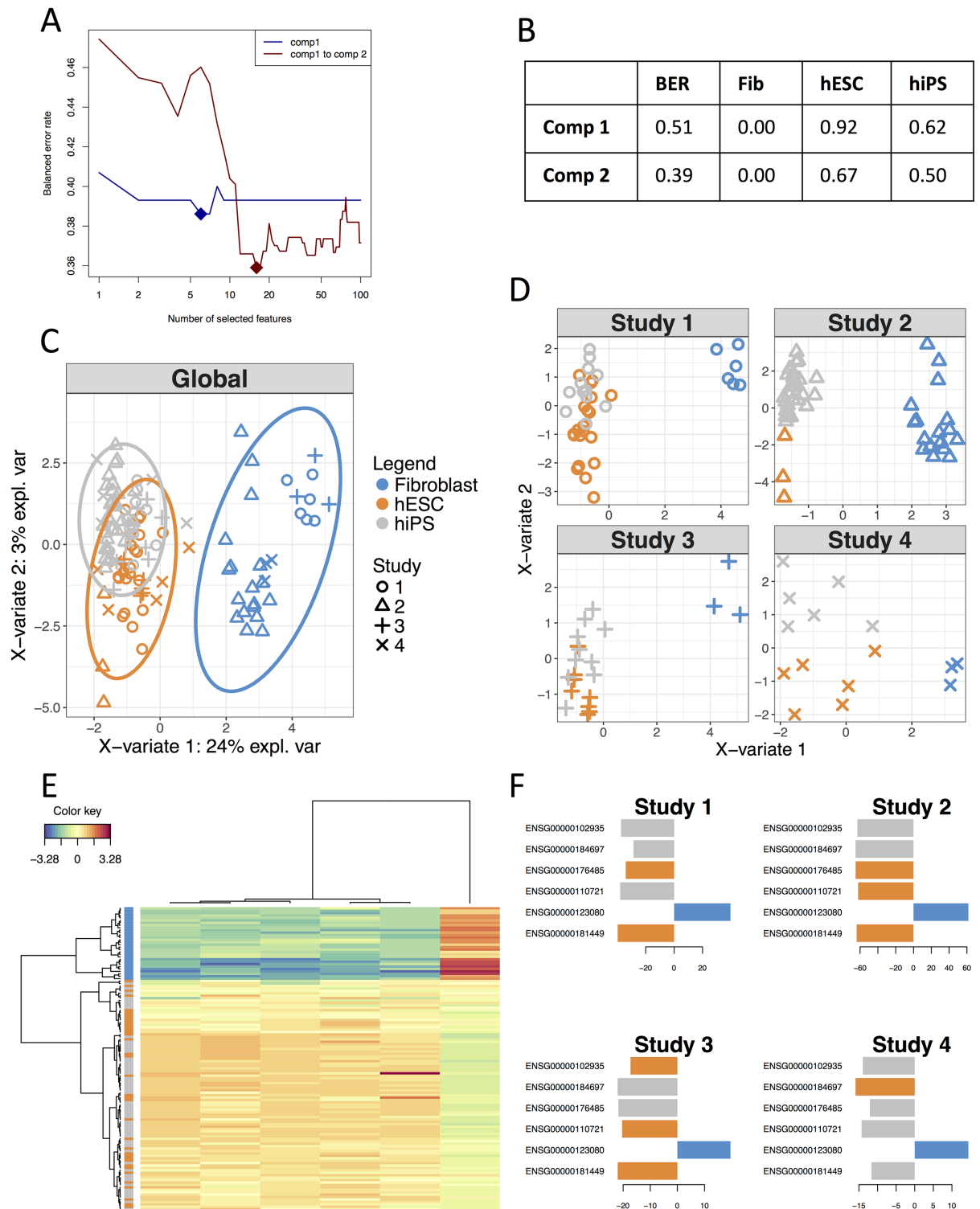
**Summary.** We illustrated the `mixOmics` framework for the supervised analysis of a multiple 'omics study. The full pipeline, results interpretation and associated R and Sweave codes are available in Supporting Information [S1 Appendix](#). Our DIABLO method identifies a discriminant and highly correlated multi-'omics signature. Predictive ability of the identified signature can be assessed (see [S1 Appendix](#)) while the graphical visualisation tools enable a better understanding of the correlation structure of the signature. Such method is the first of its kind to perform multivariate integration and discriminant analysis. DIABLO is useful to pinpoint a subset of different types of 'omics features in those large studies, posit novel hypotheses, and can be applied as a first filtering step prior to refined knowledge- and/or data-driven pathway analyses.

## *P*-integration across independent data sets with MINT

*P*-integration consists in integrating several independent studies measuring the same *P* predictors, and, in a supervised context, in identifying a robust molecular signature across multiple studies to discriminate biological conditions. The advantages of *P*-integration is to increase sample size while allowing to benchmark or compare similar studies. Contrary to usual approaches that sequentially accommodate for technical differences among the studies before classifying samples, MINT is a single step method that reduces overfitting and that predicts the class of new samples [12]. The MINT method is described in Supporting Information [S1 Text](#). We illustrate the MINT analysis on a stem cell study available from the package.

**Data description.** We combined four independent transcriptomics stem cell studies measuring the expression levels of 400 genes across 125 samples (cells). Prior to the analysis with `mixOmics`, the data were normalised and filtered for illustrative purpose. Cells were classified into 30 Fibroblasts, 37 hESC and 58 hiPSC.

**Choice of parameters and analysis.** The optimal number of components was 1 on this data set. However, in order to obtain 2D graphics, we considered a model with 2 components. The tuning process of a MINT sPLS-DA identified a molecular signature of 6 and 16 genes on the first two components, respectively ([Fig 5A](#)). A MINT model based on these



**Fig 5. Illustration of MINT analysis in mixOmics.** **A:** Parameter tuning of a MINT sPLS-DA model with two components using Leave-One-Group-Out cross-validation and maximum distance, BER (y-axis) with respect to number of selected features (x-axis). Full diamond represents the optimal number of features to select on each component, **B:** Performance of the final MINT sPLS-DA model including selected features based on BER and classification error rate per class, **C:** Global sample plot with confidence ellipse plots, **D:** Study specific sample plot, **E:** Clustered Image Map (Euclidean Distance, Complete linkage). Samples are represented in rows, selected features on the first component in columns. **F:** Loading plot of each feature selected on the first component in each study, with color indicating the class with a maximal mean expression value for each gene.

<https://doi.org/10.1371/journal.pcbi.1005752.g005>

parameters led to a BER of 0.39 (Fig 5B), which was comparable to the BER of 0.37 from MINT PLS-DA when no feature selection was performed (see details in Supporting Information S1 Appendix).

**Results visualisation.** Global sample plot (Fig 5C) and study-specific sample plots highlighted a good agreement between the four studies (Fig 5D). The first component segregated fibroblasts vs. hiPSC and hESC, and the second component hiPSC vs. hESC. Such observation was confirmed with a Clustered Image Map based on the 6 genes selected on the first component (Fig 5E). Importantly, the loading plots depicted in Fig 5F showed consistent weights assigned by the MINT model to each selected genes across each independent study.

**Summary.** We illustrated the MINT analysis for the supervised integrative analysis of multiple independent 'omics studies. The full pipeline, results interpretation and associated R and Sweave codes are available in Supporting Information S1 Appendix. Our framework proposes graphical visualisation tools to understand the identified molecular signature across all independent studies. Our own applications of the method to full data sets have showed strong potential of the method to identify reliable and robust biomarkers across independent transcriptomics studies [12, 36].

## Conclusions and future directions

The technological race in high-throughput biology leads to increasingly complex biological problems which require innovative statistical and analytical tools. Our package `mixOmics` focuses on data exploration and data mining, which are crucial steps for a first understanding of large data sets. In this article we presented our latest methods to answer cutting-edge integrative and multivariate questions in biology.

The sparse version of our methods are particularly insightful to identify molecular signatures across those multiple data sets. Feature selection resulting from our methods help refine biological hypotheses, suggest downstream analyses including statistical inference analyses, and may propose biological experimental validations. Indeed, multivariate methods include appealing properties to mine and analyse large and complex biological data, as they allow for more relaxed assumptions about data distribution, data size and data range than univariate methods, and provide insightful visualisations. In the last few years, several R packages have been proposed for multivariate analysis as a mean for dimension reduction of one data set (see the review from [3], Table 2 lists all packages and functions currently available), and the integration of two or more data sets (see [3], Table 3 and `FactoMineR` [40]). However, very few methods propose feature selection, including sparse CCA (`PMA` package [39]), sparse PLS (`spls` package, [41]), penalised PLS (`pppls` package [42]), sGCCA (`RGCCA` package [32]), PARAFAC and Tucker multi-way analyses (`ThreeWay`, `PTAk`, `ade4` packages, [43–45]) and even fewer methods provide data visualisation of the selected features (`ade4`).

The identification of a *combination* of discriminative features meet biological assumptions that cannot be addressed with univariate methods. Nonetheless, we believe that combining different types of statistical methods (univariate, multivariate, machine learning) is the key to answer complex biological questions. However, such questions must be well stated, in order for those exploratory integrative methods to provide meaningful results, and especially for the non trivial case of multiple data integration.

While we illustrated our different frameworks on classical 'omics data in a supervised context, the package also include their unsupervised counterparts to investigate relationships and associations between features with no prior phenotypic or response information. Here we applied our multivariate frameworks to transcriptomics, proteomics and miRNA data. However, other types of biological data can be analysed, as well as data beyond the realm of 'omics as

long as they are expressed as *continuous values*. Sequence-based data after processing (i.e. corrected for library size and log transformed) fit this requirement, as well as clinical data. Genotype data, such as bi-allelic Single Nucleotide Polymorphism coded as counts of the minor allele can also fit in our framework, by implicitly considering an additive model. However, to consider SNPs as categorical variables additional methodological developments are required as each SNP needs to be considered as dummy indicator matrices in the sparse multivariate models.

Currently our methods are linear techniques, where each component is constructed based on a linear combination of variables. Components between different data sets however are not linearly dependent as we maximise the covariance between them [46]. PLS-based models assuming a non-linear relationship between different sets of data have been proposed [47] but the interpretation of the results in terms of identified signature is not straightforward. We are currently investigating sparse kernel-based method for non linear modelling.

Finally, the sPLS-DA framework was recently extended for microbiome 16S data [17], and we will further extend DIABLO and MINT for microbiome—'omics integration, as well as for genomic data and time-course experiments. These two promising integrative frameworks can also be combined for NP-integration, to combine multiple studies that each include several types of 'omics data and open new avenues for large scale multiple data integration.

## Supporting information

**S1 Text. Supplemental information regarding general definitions, graphical outputs to visualise multivariate analysis results, methods description for single 'omics supervised multivariate analysis with PLS-DA and sPLS-DA, N-integration across multiple 'omics data sets with DIABLO and P-integration across independent data sets with MINT and additional computational time report for large data sets.**

(PDF)

**S1 Appendix. Sweave and R codes for all example analyses are provided, and also available on our website <http://mixomics.org>.**

(ZIP)

## Acknowledgments

The authors would like to thank the numerous `mixOmics` users who continuously help in improving the usability of the package.

## Author Contributions

**Conceptualization:** Florian Rohart, Kim-Anh Lê Cao.

**Formal analysis:** Florian Rohart, Kim-Anh Lê Cao.

**Funding acquisition:** Kim-Anh Lê Cao.

**Investigation:** Amrit Singh, Kim-Anh Lê Cao.

**Methodology:** Florian Rohart, Benoît Gautier, Amrit Singh, Kim-Anh Lê Cao.

**Project administration:** Kim-Anh Lê Cao.

**Software:** Florian Rohart, Benoît Gautier, Kim-Anh Lê Cao.

**Supervision:** Kim-Anh Lê Cao.

**Visualization:** Florian Rohart, Benoît Gautier, Amrit Singh, Kim-Anh Lê Cao.

**Writing – original draft:** Florian Rohart, Kim-Anh Lê Cao.

**Writing – review & editing:** Florian Rohart, Kim-Anh Lê Cao.

## References

1. Lê Cao KA, Rohart F, Gonzalez I, Déjean S, Gautier B, Bartolo F, et al. mixOmics: Omics Data Integration Project; 2017. Available from: <https://CRAN.R-project.org/package=mixOmics>.
2. Boulesteix AL, Strimmer K. Partial least squares: a versatile tool for the analysis of high-dimensional genomic data. *Brief Bioinform.* 2007; 8(1):32–44. <https://doi.org/10.1093/bib/bbl016> PMID: 16772269
3. Meng C, Zeleznik OA, Thallinger GG, Kuster B, Gholami AM, Culhane AC. Dimension reduction techniques for the integrative analysis of multi-omics data. *Briefings in bioinformatics.* 2016; p. bbv108. <https://doi.org/10.1093/bib/bbv108> PMID: 26969681
4. Labus JS, Van Horn JD, Gupta A, Alaverdyan M, Torgerson C, Ashe-McNalley C, et al. Multivariate morphological brain signatures predict patients with chronic abdominal pain from healthy control subjects. *Pain.* 2015; 156(8):1545–1554. <https://doi.org/10.1097/j.pain.000000000000196> PMID: 25906347
5. Cook JA, Chandramouli GV, Anver MR, Sowers AL, Thetford A, Krausz KW, et al. Mass Spectrometry–Based Metabolomics Identifies Longitudinal Urinary Metabolite Profiles Predictive of Radiation-Induced Cancer. *Cancer research.* 2016; 76(6):1569–1577. <https://doi.org/10.1158/0008-5472.CAN-15-2416> PMID: 26880804
6. Guidi L, Chaffron S, Bittner L, Eveillard D, Larhimi A, Roux S, et al. Plankton networks driving carbon export in the oligotrophic ocean. *Nature.* 2016;. <https://doi.org/10.1038/nature16942>
7. Mahana D, Trent CM, Kurtz ZD, Bokulich NA, Battaglia T, Chung J, et al. Antibiotic perturbation of the murine gut microbiome enhances the adiposity, insulin resistance, and liver disease associated with high-fat diet. *Genome medicine.* 2016; 8(1):1. <https://doi.org/10.1186/s13073-016-0297-9>
8. Ramanan D, Bowcutt R, Lee SC, San Tang M, Kurtz ZD, Ding Y, et al. Helminth infection promotes colonization resistance via type 2 immunity. *Science.* 2016; 352(6285):608–612. <https://doi.org/10.1126/science.aaf3229> PMID: 27080105
9. Rollero S, Mouret JR, Sanchez I, Camarasa C, Ortiz-Julien A, Sablayrolles JM, et al. Key role of lipid management in nitrogen and aroma metabolism in an evolved wine yeast strain. *Microbial cell factories.* 2016; 15(1):1. <https://doi.org/10.1186/s12934-016-0434-6>
10. Lê Cao KA, Boitard S, Besse P. Sparse PLS Discriminant Analysis: biologically relevant feature selection and graphical displays for multiclass problems. *BMC bioinformatics.* 2011; 12(1):253. <https://doi.org/10.1186/1471-2105-12-253> PMID: 21693065
11. Singh A, Gautier B, Shannon CP, Vacher M, Rohart F, Tebutt SJ, et al. DIABLO-an integrative, multi-omics, multivariate method for multi-group classification. *bioRxiv.* 2016;067611.
12. Rohart F, Eslami A, Matigian N, Bougeard S, Lê Cao KA. MINT: A multivariate integrative approach to identify a reproducible biomarker signature across multiple experiments and platforms. *BMC Bioinformatics.* 2017; 18(128). <https://doi.org/10.1186/s12859-017-1553-8> PMID: 28241739
13. Liu Y, Devescovi V, Chen S, Nardini C. Multilevel omic data integration in cancer cell lines: advanced annotation and emergent properties. *BMC systems biology.* 2013; 7(1):14. <https://doi.org/10.1186/1752-0509-7-14> PMID: 23418673
14. Günther OP, Chen V, Freue GC, Balshaw RF, Tebutt SJ, Hollander Z, et al. A computational pipeline for the development of multi-marker bio-signature panels and ensemble classifiers. *BMC bioinformatics.* 2012; 13(1):326. <https://doi.org/10.1186/1471-2105-13-326> PMID: 23216969
15. Teng M, Love MI, Davis CA, Djebali S, Dobin A, Graveley BR, et al. A benchmark for RNA-seq quantification pipelines. *Genome biology.* 2016; 17(1):74. <https://doi.org/10.1186/s13059-016-0940-1> PMID: 27107712
16. Arumugam M, Raes J, Pelletier E, Le Paslier D, Yamada T, Mende DR, et al. Enterotypes of the human gut microbiome. *nature.* 2011; 473(7346):174. <https://doi.org/10.1038/nature09944> PMID: 21508958
17. Lê Cao KA, Lakis VA, Bartolo F, Costello ME, Chua XY, Brazeilles R, et al. MixMC: Multivariate insights into Microbial Communities. *PloS one.* 2016; 11(8):e0160169.
18. Wold H. Path models with latent variables: The NIPALS approach. *Acad. Press;* 1975.
19. Yao F, Coquery J, Lê Cao KA. Independent Principal Component Analysis for biologically meaningful dimension reduction of large biological data sets. *BMC bioinformatics.* 2012; 13(1):24. <https://doi.org/10.1186/1471-2105-13-24> PMID: 22305354
20. Wold H. Estimation of principal components and related models by iterative least squares. *J Multivar Anal.* 1966; p. 391–420.

21. Eslami A, Qannari EM, Kohler A, Bougeard S. Multi-group PLS Regression: Application to Epidemiology. In: *New Perspectives in Partial Least Squares and Related Methods*. Springer; 2013. p. 243–255.
22. González I, Déjean S, Martin PG, Baccini A, et al. CCA: An R package to extend canonical correlation analysis. *Journal of Statistical Software*. 2008; 23(12):1–14.
23. Tenenhaus A, Tenenhaus M. Regularized generalized canonical correlation analysis. *Psychometrika*. 2011; 76(2):257–284. <https://doi.org/10.1007/s11336-011-9206-8>
24. Nguyen DV, Rocke DM. Tumor classification by partial least squares using microarray gene expression data. *Bioinformatics*. 2002; 18(1):39–50. <https://doi.org/10.1093/bioinformatics/18.1.39> PMID: 11836210
25. Nguyen DV, Rocke DM. Multi-class cancer classification via partial least squares with gene expression profiles. *Bioinformatics*. 2002; 18(9):1216–1226. <https://doi.org/10.1093/bioinformatics/18.9.1216> PMID: 12217913
26. Boulesteix AL. PLS dimension reduction for classification with microarray data. *Statistical applications in genetics and molecular biology*. 2004; 3(1):1–30. <https://doi.org/10.2202/1544-6115.1075>
27. Tibshirani R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B (Methodological)*. 1996; p. 267–288.
28. Wangen L, Kowalski B. A multiblock partial least squares algorithm for investigating complex chemical systems. *Journal of chemometrics*. 1989; 3(1):3–20. <https://doi.org/10.1002/cem.1180030104>
29. Westerhuis JA, Smilde AK. Deflation in multiblock PLS. *Journal of chemometrics*. 2001; 15(5):485–493. <https://doi.org/10.1002/cem.652>
30. Karaman İ, Nørskov NP, Yde CC, Hedemann MS, Knudsen KEB, Kohler A. Sparse multi-block PLSR for biomarker discovery when integrating data from LC–MS and NMR metabolomics. *Metabolomics*. 2015; 11(2):367–379. <https://doi.org/10.1007/s11306-014-0698-y>
31. Kawaguchi A, Yamashita F. Supervised multiblock sparse multivariable analysis with application to multimodal brain imaging genetics. *Biostatistics*. 2017; p. kxx011.
32. Tenenhaus A, Guillemot V. RGCCA: Regularized and Sparse Generalized Canonical Correlation Analysis for Multiblock Data; 2017. Available from: <https://CRAN.R-project.org/package=RGCCA>.
33. Tenenhaus A, Philippe C, Guillemot V, Lê Cao KA, Grill J, Frouin V. Variable selection for generalized canonical correlation analysis. *Biostatistics*. 2014; 15(3):569–83. <https://doi.org/10.1093/biostatistics/kxu001> PMID: 24550197
34. González I, Lê Cao KA, Davis MJ, Déjean S, et al. Visualising associations between paired 'omics' data sets. *BioData mining*. 2012; 5(1):19. <https://doi.org/10.1186/1756-0381-5-19> PMID: 23148523
35. Khan J, Wei JS, Ringner M, Saal LH, Ladanyi M, Westermann F, et al. Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural networks. *Nature medicine*. 2001; 7(6):673–679. <https://doi.org/10.1038/89044> PMID: 11385503
36. Rohart F, Mason EA, Matigian N, Mosbergen R, Korn O, Chen T, et al. A molecular classification of human mesenchymal stromal cells. *PeerJ*. 2016; 4:e1845. <https://doi.org/10.7717/peerj.1845> PMID: 27042394
37. Shah AK, Lê Cao KA, Choi E, Chen D, Gautier B, Nancarrow D, et al. Glyco-centric lectin magnetic bead array (LeMBA)- proteomics dataset of human serum samples from healthy, Barrett's esophagus and esophageal adenocarcinoma individuals. *Data in Brief*. 2016; 7:1058–1062. <https://doi.org/10.1016/j.dib.2016.03.081> PMID: 27408916
38. Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*. 2010; 33(1):1. <https://doi.org/10.18637/jss.v033.i01> PMID: 20808728
39. Witten D, Tibshirani R, Gross S, Narasimhan B. PMA: Penalized Multivariate Analysis; 2013. Available from: <https://CRAN.R-project.org/package=PMA>.
40. Husson F, Josse J, Le S, Mazet J. FactoMineR: factor analysis and data mining with R; 2017. Available from: <https://cran.r-project.org/web/packages/FactoMineR>.
41. Chung D, Chun H, Keles S. SPLS: Sparse partial least squares (SPLS) regression and classification; 2013. Available from: <https://CRAN.R-project.org/package=spls>.
42. Kraemer N, Boulesteix A. ppls: Penalized Partial Least Squares; 2014. Available from: <https://CRAN.R-project.org/package=ppls>.
43. Del Ferraro M, Kiers H, Giordani P. ThreeWay: Three-Way Component Analysis; 2015. Available from: <https://cran.r-project.org/web/packages/ThreeWay>.
44. Leibovici D. PTAk: Principal Tensor Analysis on k Modes; 2015. Available from: <https://cran.r-project.org/web/packages/PTAk>.

45. Thioulouse J, Chessel D, Dolédec S, Olivier J, Goreaud F, Pelissier R. ADE-4: Ecological data analysis. Exploratory and euclidean methods in environmental sciences; 2017. Available from: <https://cran.r-project.org/web/packages/ade4>.
46. Krämer N, Sugiyama M. The degrees of freedom of partial least squares regression. *Journal of the American Statistical Association*. 2011; 106(494):697–705. <https://doi.org/10.1198/jasa.2011.tm10107>
47. Rosipal R. Nonlinear partial least squares: An overview. *Chemoinformatics and advanced machine learning perspectives: complex computational methods and collaborative techniques*. 2010; p. 169–189.