

MULTIVARIATE DATA INTEGRATION USING R

Methods and Applications
with the mixOmics Package

Kim-Anh Lê Cao
Zoe Welham



CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

Multivariate Data Integration Using R

Computational Biology Series

About the Series:

This series aims to capture new developments in computational biology, as well as high-quality work summarizing or contributing to more established topics. Publishing a broad range of reference works, textbooks, and handbooks, the series is designed to appeal to students, researchers, and professionals in all areas of computational biology, including genomics, proteomics, and cancer computational biology, as well as interdisciplinary researchers involved in associated fields, such as bioinformatics and systems biology.

Introduction to Bioinformatics with R: A Practical Guide for Biologists

Edward Curry

Analyzing High-Dimensional Gene Expression and DNA Methylation Data with R

Hongmei Zhang

Introduction to Computational Proteomics

Golan Yona

Glycome Informatics: Methods and Applications

Kiyoko F. Aoki-Kinoshita

Computational Biology: A Statistical Mechanics Perspective

Ralf Blossey

Computational Hydrodynamics of Capsules and Biological Cells

Constantine Pozrikidis

Computational Systems Biology Approaches in Cancer Research

Inna Kuperstein, Emmanuel Barillot

Clustering in Bioinformatics and Drug Discovery

John David MacCuish, Norah E. MacCuish

Metabolomics: Practical Guide to Design and Analysis

Ron Wehrens, Reza Salek

An Introduction to Systems Biology: Design Principles of Biological Circuits

2nd Edition

Uri Alon

Computational Biology: A Statistical Mechanics Perspective

Second Edition

Ralf Blossey

Stochastic Modelling for Systems Biology

Third Edition

Darren J. Wilkinson

Computational Genomics with R

Altuna Akalin, Bora Uyar, Vedran Franke, Jonathan Ronen

An Introduction to Computational Systems Biology: Systems-level Modelling of Cellular Networks

Karthik Raman

Virus Bioinformatics

Dmitrij Frishman, Manuela Marz

Multivariate Data Integration Using R: Methods and Applications with the mixOmics Package

Kim-Anh LeCao, Zoe Marie Welham

Bioinformatics

A Practical Guide to NCBI Databases and Sequence Alignments

Hamid D. Ismail

For more information about this series please visit:

<https://www.routledge.com/Chapman--HallCRC-Computational-Biology-Series/book-series/CRCCBS>

Multivariate Data Integration Using R

Methods and Applications with the mixOmics Package

Kim-Anh Lê Cao
Zoe Welham



CRC Press

Taylor & Francis Group

Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group, an **informa** business
A CHAPMAN & HALL BOOK

First edition published 2022
by CRC Press
6000 Broken Sound Parkway NW, Suite 300, Boca Raton, FL 33487-2742

and by CRC Press
2 Park Square, Milton Park, Abingdon, Oxon, OX14 4RN

© 2022 Taylor & Francis Group, LLC

CRC Press is an imprint of Taylor & Francis Group, LLC

Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, access www.copyright.com or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. For works that are not available on CCC please contact mpkbookspermissions@tandf.co.uk

Trademark notice: Product or corporate names may be trademarks or registered trademarks and are used only for identification and explanation without intent to infringe.

ISBN: 9780367460945 (hbk)

ISBN: 9781032128078 (pbk)

ISBN: 9781003026860 (ebk)

DOI: [10.1201/9781003026860](https://doi.org/10.1201/9781003026860)

This book has been prepared from camera-ready copy provided by the authors.

From Kim-Anh Lê Cao:
To my parents, Betty and Huy Lê Cao
To the mixOmics team and our mixOmics users,
And to my co-author Zoe Welham without whom this book would not have existed.

For Zoe Welham:
To Joy and Tamara Welham, for their continual support



Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

Contents

Preface	xv
Authors	xxi
I Modern biology and multivariate analysis	1
1 Multi-omics and biological systems	3
1.1 Statistical approaches for reductionist or holistic analyses	3
1.2 Multi-omics and multivariate analyses	4
1.2.1 More than a ‘scale up’ of univariate analyses	5
1.2.2 More than a fishing expedition	5
1.3 Shifting the analysis paradigm	5
1.4 Challenges with high-throughput data	6
1.4.1 Overfitting	7
1.4.2 Multi-collinearity and ill-posed problems	7
1.4.3 Zero values and missing values	7
1.5 Challenges with multi-omics integration	8
1.5.1 Data heterogeneity	8
1.5.2 Data size	8
1.5.3 Platforms	8
1.5.4 Expectations for analysis	8
1.5.5 Variety of analytical frameworks	9
1.6 Summary	9
2 The cycle of analysis	11
2.1 <i>The Problem</i> guides the analysis	11
2.2 <i>Plan</i> in advance	12
2.2.1 What affects statistical power?	12
2.2.2 Sample size	12
2.2.3 Identify covariates and confounders	13
2.2.4 Identify batch effects	13
2.3 <i>Data</i> cleaning and pre-processing	14
2.3.1 Normalisation	14
2.3.2 Filtering	15
2.3.3 Missing values	15
2.4 <i>Analysis</i> : Choose the right approach	15
2.4.1 Descriptive statistics	15
2.4.2 Exploratory statistics	15
2.4.3 Inferential statistics	16
2.4.4 Univariate or multivariate modelling?	16
2.4.5 Prediction	17

2.5	<i>Conclusion</i> and start the cycle again	18
2.6	Summary	18
3	Key multivariate concepts and dimension reduction in <i>mixOmics</i>	19
3.1	Measures of dispersion and association	19
3.1.1	Random variables and biological variation	19
3.1.2	Variance	20
3.1.3	Covariance	20
3.1.4	Correlation	21
3.1.5	Covariance and correlation in <i>mixOmics</i> context	22
3.1.6	R examples	22
3.2	Dimension reduction	23
3.2.1	Matrix factorisation	23
3.2.2	Factorisation with components and loading vectors	24
3.2.3	Data visualisation using components	24
3.3	Variable selection	25
3.3.1	Ridge penalty	26
3.3.2	Lasso penalty	26
3.3.3	Elastic net	26
3.3.4	Visualisation of the selected variables	26
3.4	Summary	27
4	Choose the right method for the right question in <i>mixOmics</i>	29
4.1	Types of analyses and methods	29
4.1.1	Single or multiple omics analysis?	29
4.1.2	N – or P –integration?	30
4.1.3	Unsupervised or supervised analyses?	31
4.1.4	Repeated measures analyses	32
4.1.5	Compositional data	32
4.2	Types of data	33
4.2.1	Classical omics	33
4.2.2	Microbiome data: A special case	33
4.2.3	Genotype data: A special case	33
4.2.4	Clinical variables that are categorical: A special case	33
4.3	Types of biological questions	34
4.3.1	A PCA type of question (one data set, unsupervised)	34
4.3.2	A PLS type of question (two data sets, regression or unsupervised)	34
4.3.3	A CCA type of question (two data sets, unsupervised)	35
4.3.4	A PLS-DA type of question (one data set, classification)	35
4.3.5	A multiblock PLS type of question (more than two data sets, supervised or unsupervised)	36
4.3.6	An N –integration type of question (several data sets, supervised)	36
4.3.7	A P –integration type of question (several studies of the same omic type, supervised or unsupervised)	37
4.4	Exemplar data sets in <i>mixOmics</i>	37
4.5	Summary	37
4.A	Appendix: Data transformations in <i>mixOmics</i>	38
4.A.1	Multilevel decomposition	38
4.A.2	Mixed-effect model context	39
4.A.3	Split-up variation	39
4.A.4	Example of multilevel decomposition in <i>mixOmics</i>	40

4.B	Centered log ratio transformation	41
4.C	Creating dummy variables	42
II	mixOmics under the hood	45
5	Projection to latent structures	47
5.1	PCA as a projection algorithm	47
5.1.1	Overview	47
5.1.2	Calculating the components	48
5.1.3	Meaning of the loading vectors	49
5.1.4	Example using the <code>linnerud</code> data in <code>mixOmics</code>	49
5.2	Singular Value Decomposition (SVD)	50
5.2.1	SVD algorithm	50
5.2.2	Example in R	52
5.2.3	Matrix approximation	54
5.3	Non-linear Iterative Partial Least Squares (NIPALS)	54
5.3.1	NIPALS pseudo algorithm	55
5.3.2	Local regressions	55
5.3.3	Deflation	56
5.3.4	Missing values	57
5.4	Other matrix factorisation methods in <code>mixOmics</code>	57
5.5	Summary	57
6	Visualisation for data integration	59
6.1	Sample plots using components	59
6.1.1	Example with PCA and <code>plotIndiv</code>	59
6.1.2	Sample plot for the integration of two or more data sets	60
6.1.3	Representing paired coordinates using <code>plotArrow</code>	63
6.2	Variable plots using components and loading vectors	65
6.2.1	Loading plots	65
6.2.2	Correlation circle plots	66
6.2.3	Biplots	69
6.2.4	Relevance networks	70
6.2.5	Clustered Image Maps (CIM)	73
6.2.6	Circos plots	74
6.3	Summary	75
6.A	Appendix: Similarity matrix in relevance networks and CIM	76
6.A.1	Pairwise variable associations for CCA	76
6.A.2	Pairwise variable associations for PLS	76
6.A.3	Constructing relevance networks and displaying CIM	77
7	Performance assessment in multivariate analyses	79
7.1	Main parameters to choose	79
7.2	Performance assessment	80
7.2.1	Training and testing: If we were rich	80
7.2.2	Cross-validation: When we are poor	81
7.3	Performance measures	82
7.3.1	Evaluation measures for regression	82
7.3.2	Evaluation measures for classification	83
7.3.3	Details of the tuning process	83
7.4	Final model assessment	86

7.4.1	Assessment of the performance	86
7.4.2	Assessment of the signature	86
7.5	Prediction	87
7.5.1	Prediction of a continuous response	87
7.5.2	Prediction of a categorical response	88
7.5.3	Prediction is related to the number of components	90
7.6	Summary and roadmap of analysis	90
III	mixOmics in action	93
8	mixOmics: Get started	95
8.1	Prepare the data	95
8.1.1	Normalisation	95
8.1.2	Filtering variables	96
8.1.3	Centering and scaling the data	96
8.1.4	Managing missing values	100
8.1.5	Managing batch effects	101
8.1.6	Data format	101
8.2	Get ready with the software	102
8.2.1	R installation	102
8.2.2	Pre-requisites	102
8.2.3	mixOmics download	102
8.2.4	Load the package	103
8.3	Coding practices	103
8.3.1	Set the working directory	103
8.3.2	Good coding practices	104
8.4	Upload data	104
8.4.1	Data sets	104
8.4.2	Dependent variables	104
8.4.3	Set up the outcome for supervised classification analyses	105
8.4.4	Check data upload	106
8.5	Structure of the following chapters	106
9	Principal Component Analysis (PCA)	109
9.1	Why use PCA?	109
9.1.1	Biological questions	109
9.1.2	Statistical point of view	109
9.2	Principle	110
9.2.1	PCA	110
9.2.2	Sparse PCA	111
9.3	Input arguments	112
9.3.1	Center or scale the data?	112
9.3.2	Number of components (choice of dimensions)	112
9.3.3	Number of variables to select in sPCA	113
9.4	Key outputs	113
9.5	Case study: Multidrug	114
9.5.1	Load the data	114
9.5.2	Quick start	115
9.5.3	Example: PCA	116
9.5.4	Example: Sparse PCA	121
9.5.5	Example: Missing values imputation	125

9.6	To go further	129
9.6.1	Additional processing steps	129
9.6.2	Independent component analysis	129
9.6.3	Incorporating biological information	130
9.7	FAQ	131
9.8	Summary	132
9.A	Appendix: Non-linear Iterative Partial Least Squares	132
9.A.1	Solving PCA with NIPALS	132
9.A.2	Estimating missing values with NIPALS	132
9.B	Appendix: sparse PCA	133
9.B.1	sparse PCA-SVD	133
9.B.2	sPCA pseudo algorithm	134
9.B.3	Other sPCA methods	134
10	Projection to Latent Structure (PLS)	137
10.1	Why use PLS?	137
10.1.1	Biological questions	137
10.1.2	Statistical point of view	137
10.2	Principle	138
10.2.1	Univariate PLS1 and multivariate PLS2	139
10.2.2	PLS deflation modes	140
10.2.3	sparse PLS	142
10.3	Input arguments and tuning	142
10.3.1	The deflation mode	142
10.3.2	The number of dimensions	143
10.3.3	Number of variables to select	143
10.4	Key outputs	144
10.4.1	Graphical outputs	144
10.4.2	Numerical outputs	144
10.5	Case study: Liver toxicity	145
10.5.1	Load the data	146
10.5.2	Quick start	146
10.5.3	Example: PLS1 regression	147
10.5.4	Example: PLS2 regression	152
10.6	Take a detour: PLS2 regression for prediction	163
10.7	To go further	165
10.7.1	Orthogonal projections to latent structures	165
10.7.2	Redundancy analysis	166
10.7.3	Group PLS	166
10.7.4	PLS path modelling	166
10.7.5	Other sPLS variants	167
10.8	FAQ	167
10.9	Summary	168
10.A	Appendix: PLS algorithm	169
10.A.1	PLS Pseudo algorithm	169
10.A.2	Convergence of the PLS iterative algorithm	170
10.A.3	PLS-SVD method	170
10.B	Appendix: sparse PLS	171
10.B.1	sparse PLS-SVD	171
10.B.2	sparse PLS pseudo algorithm	171
10.C	Appendix: Tuning the number of components	172

10.C.1	In PLS1	172
10.C.2	In PLS2	175
11	Canonical Correlation Analysis (CCA)	177
11.1	Why use CCA?	177
11.1.1	Biological question	177
11.1.2	Statistical point of view	177
11.2	Principle	178
11.2.1	CCA	178
11.2.2	rCCA	179
11.3	Input arguments and tuning	179
11.3.1	CCA	179
11.3.2	rCCA	180
11.4	Key outputs	180
11.4.1	Graphical outputs	180
11.4.2	Numerical outputs	181
11.5	Case study: Nutrimouse	181
11.5.1	Load the data	182
11.5.2	Quick start	182
11.5.3	Example: CCA	183
11.5.4	Example: rCCA	184
11.6	To go further	193
11.7	FAQ	194
11.8	Summary	195
11.A	Appendix: CCA and variants	196
11.A.1	Solving classical CCA	196
11.A.2	Regularised CCA	197
12	PLS-Discriminant Analysis (PLS-DA)	201
12.1	Why use PLS-DA?	201
12.1.1	Biological question	201
12.1.2	Statistical point of view	201
12.2	Principle	202
12.2.1	PLS-DA	203
12.2.2	sparse PLS-DA	204
12.3	Input arguments and tuning	204
12.3.1	PLS-DA	204
12.3.2	sPLS-DA	205
12.3.3	Framework to manage overfitting	205
12.4	Key outputs	206
12.4.1	Numerical outputs	207
12.4.2	Graphical outputs	207
12.5	Case study: SRBCT	207
12.5.1	Load the data	208
12.5.2	Quick start	208
12.5.3	Example: PLS-DA	209
12.5.4	Example: sPLS-DA	214
12.5.5	Take a detour: Prediction	223
12.5.6	AUROC outputs complement performance evaluation	225
12.6	To go further	226
12.6.1	Microbiome	226

12.6.2	Multilevel	227
12.6.3	Other related methods and packages	228
12.7	FAQ	228
12.8	Summary	229
12.A	Appendix: Prediction in PLS-DA	229
12.A.1	Prediction distances	229
12.A.2	Background area	231
13	N–data integration	233
13.1	Why use N –integration methods?	233
13.1.1	Biological question	233
13.1.2	Statistical point of view and analytical challenges	234
13.2	Principle	234
13.2.1	Multiblock sPLS-DA	234
13.2.2	Prediction in multiblock sPLS-DA	236
13.3	Input arguments and tuning	237
13.4	Key outputs	238
13.4.1	Graphical outputs	238
13.4.2	Numerical outputs	238
13.5	Case Study: breast.TCGA	239
13.5.1	Load the data	239
13.5.2	Quick start	240
13.5.3	Parameter choice	241
13.5.4	Final model	244
13.5.5	Sample plots	245
13.5.6	Variable plots	247
13.5.7	Model performance and prediction	251
13.6	To go further	255
13.6.1	Additional data transformation for special cases	255
13.6.2	Other N –integration frameworks in mixOmics	255
13.6.3	Supervised classification analyses: concatenation and ensemble methods	256
13.6.4	Unsupervised analyses: JIVE and MOFA	256
13.7	FAQ	257
13.8	Additional resources	258
13.9	Summary	258
13.A	Appendix: Generalised CCA and variants	258
13.A.1	regularised GCCA	258
13.A.2	sparse GCCA	259
13.A.3	sparse multiblock sPLS-DA	260
14	P–data integration	261
14.1	Why use P –integration methods?	261
14.1.1	Biological question	261
14.1.2	Statistical point of view	261
14.2	Principle	262
14.2.1	Motivation	262
14.2.2	Multi-group sPLS-DA	263
14.3	Input arguments and tuning	264
14.3.1	Data input checks	264
14.3.2	Number of components	265

14.3.3	Number of variables to select per component	265
14.4	Key outputs	265
14.4.1	Graphical outputs	265
14.4.2	Numerical outputs	266
14.5	Case Study: stemcells	266
14.5.1	Load the data	266
14.5.2	Quick start	267
14.5.3	Example: MINT PLS-DA	268
14.5.4	Example: MINT sPLS-DA	271
14.5.5	Take a detour	277
14.6	Examples of application	280
14.6.1	16S rRNA gene data	280
14.6.2	Single cell transcriptomics	280
14.7	To go further	280
14.8	Summary	280
	Glossary of terms	283
	Key publications	285
	Bibliography	287
	Index	299

Preface

Context and scope

Modern high-throughput technologies generate information about thousands of biological molecules at different cellular levels in a biological system, leading to several types of omics studies (e.g. transcriptomics as the study of messenger RNA molecules expressed from the genes of an organism, or proteomics as the study of proteins expressed by a cell, tissue, or organism). However, a reductionist approach that considers each of these molecules individually does not fully describe an organism in its environment. Rather, we use multivariate analysis to investigate the simultaneous and complex relationships that occur in molecular pathways. In addition, to obtain a holistic picture of a complete biological system, we propose to integrate multiple layers of information using recent computational tools we have developed through the mixOmics project.

mixOmics is an international endeavour that encompasses methodological developments, software implementation, and applications to biological and biomedical problems to address some of the challenges of omics data integration. We have trained students and researchers in essential statistical and data analysis skills via our numerous multi-day workshops to build capacity in best practice statistical analysis and advance the field of computational statistics for biology. The goal of this book is to provide guidance in applying multivariate dimension reduction techniques for the integration of high-throughput biological data, allowing readers to obtain new and deeper insights into biological mechanisms and biomedical problems.

Who is this book for?

This book is suitable for biologists, computational biologists, and bioinformaticians who generate and work with high-throughput omics data. Such data include – but are not restricted to – transcriptomics, epigenomics, proteomics, metabolomics, the microbiome, and clinical data. Our book is dedicated to research postgraduate students and scientists at any career stage, and can be used for teaching specialised multi-disciplinary undergraduate and Masters's courses. Data analysts with a basic level of R programming will benefit most from this resource. The book is organised into three distinct parts, where each part can be skimmed according to the level and interest of the reader. Each chapter contains different levels of information, and the most technical chapters can be skipped during a first read.

Overview of methods in mixOmics

The mixOmics package focuses on multivariate analysis which examines more than two variables simultaneously to integrate different types of variables (e.g. genes, proteins, metabolites). We use dimension reduction techniques applicable to a wide range of data analysis types. Our analyses can be descriptive, exploratory, or focus on modeling or prediction. Our

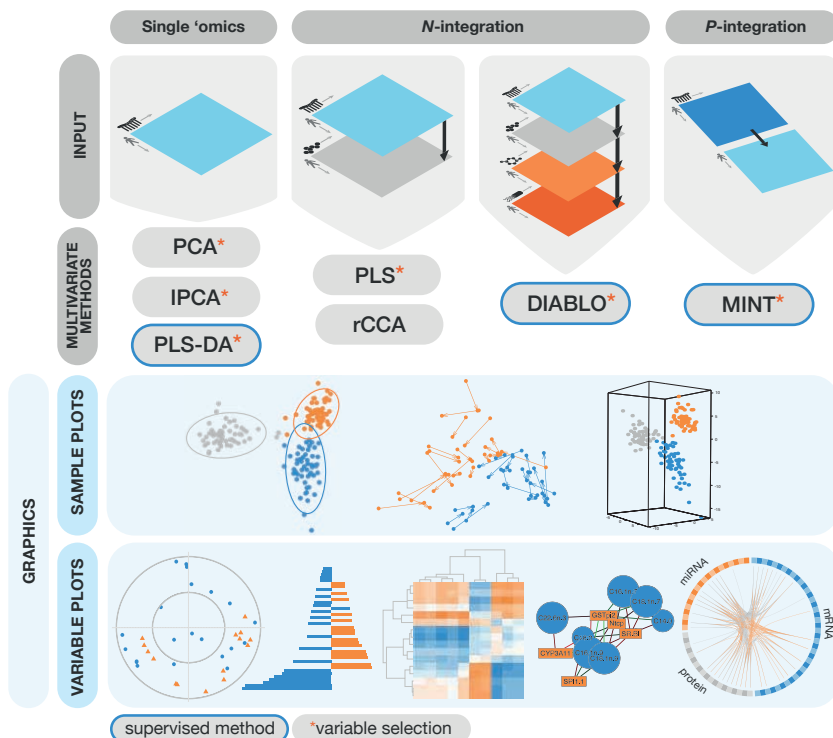


FIGURE 1: Overview of the methods implemented in the *mixOmics* package for the exploration and integration of multiple data sets. This book aims to guide the data analyst in constructing the research question, applying the appropriate multivariate techniques, and interpreting the resulting graphics.

aim is to summarise these large biological data sets to elucidate similarities between samples, between variables, and the relationship between samples and variables. The *mixOmics* package provides a range of methods to answer different kinds of biological questions, for example to:

- Highlight patterns pertaining to the major sources of variation in the data (e.g. Principal Component Analysis),
- Segregate samples according to their known group and predict group membership of new samples (e.g. Partial Least Squares Discriminant Analysis),
- Identify agreement between multiple data sets (e.g. Canonical Correlation Analysis, Partial Least Squares regression, and other variants),
- Identify molecular signatures across multiple data sets with sparse methods that achieve variable selection.

Key methodological concepts in *mixOmics*

Methods in *mixOmics* are based on *matrix factorisation* techniques, which offer great flexibility in analysing and integrating multiple data sets in a holistic manner. We use *dimension reduction* combined with *feature selection* to summarise the main characteristics of the data and posit novel biological hypotheses.

Dimension reduction is achieved by combining all original variables into a smaller number of artificial *components* that summarise patterns in the original data.

The `mixOmics` package is unique in providing novel multivariate techniques that enable feature selection to identify *molecular signatures*. Feature selection refers to identifying variables that best explain, or predict, the outcome variable (e.g. group membership, or disease status) of interest. Variables deemed irrelevant according to the specific statistical criterion we use in the methods are not taken into account when calculating the components.

Data integration methods use *data projection* techniques to maximise the covariance, or the correlation between, omics data sets. We propose two types of data integration, whether on the same N samples, or on the same P variables (Figure 1).

Finally, our methods can provide either *unsupervised* or *supervised* analyses. Unsupervised analyses are exploratory: any information about sample group membership, or outcome, is disregarded, and data are explored based on their *variance* or *correlation* structure. Supervised analyses aim to segregate sample groups known *a priori* (e.g. disease status, treatments) and identify variables (i.e. biomarker candidates, or molecular signatures) that either explain or separate sample groups.

These concepts will be explained further in Part I.

To aid in interpreting analysis results, `mixOmics` provides insightful graphical plots designed to highlight patterns in both the sample and variable dimensions uncovered by each method (Figure 1).

Concepts not covered

Each `mixOmics` method corresponds to an underlying statistical model. However, the methods we present radically differ from univariate formulations as they do not test one variable at a time, or produce p -values. In that sense, multivariate methods can be considered *exploratory* as they do not enable statistical inference. Our wide range of methods come in many different flavours and can be applied also for predictive purposes, as we detail in this book. ‘Classical’ univariate statistical inference methods can still be used in our analysis framework after the identification of molecular signatures, as our methods aim to generate novel biological hypotheses.

Who is ‘mixOmics’?

The `mixOmics` project has been developed between France, Australia and Canada since 2009, when the first version of the package was submitted to the CRAN¹. Our team is composed of core members from the University of Melbourne, Australia, and the Université de Toulouse, France. The team also includes several key contributors and collaborators.

The package implements more than nineteen multivariate and sparse methodologies for omics data exploration, integration, and biomarker discovery for different biological settings, amongst which thirteen were developed by our team (see our list of publications in Section 14.8). Originally, all methods were designed for omics data, however, their application is not limited to biological data only. Other applications where integration is required can be considered, but mostly for cases where the predictor variables are continuous.

¹The Comprehensive R Architecture Network <https://www.cran.r-project.org>

The package is currently available from Bioconductor², with a development version available on GitHub³. We continue to maintain and improve the package via new methods, code optimisation and efficient memory storage of R objects.

About this book

Part I: Modern biology and multivariate analysis introduces fundamental concepts in multivariate analysis. *Multi-omics and biological systems* (Chapter 1) compares and contrasts multivariate and univariate analysis, and outlines the advantages and challenges of multivariate analyses. *The Cycle of Analysis* (Chapter 2) details the necessary steps in planning, designing and conducting multivariate analyses. *Key multivariate concepts and dimension reduction in mixOmics* (Chapter 3) describes measures of dispersion and association, and introduces key methods in mixOmics to manage large data, such as dimension reduction using matrix factorisation and feature selection. *Choose the right method for the right question in mixOmics* (Chapter 4) provides an overview of the methods available in mixOmics and the types of biological questions these methods can answer.

Part II: mixOmics under the hood provides a deeper understanding of the statistical concepts underlying the methods presented in Part III. *Projection to Latent Structures (PLS)* (Chapter 5) illustrates the different types of algorithms used to solve Principal Component Analysis. We detail in particular the iterative PLS algorithm that projects data onto latent structures (components) for matrix decomposition and dimension reduction, as this algorithm forms the basis of most of our methods. *Visualisation for data integration* (Chapter 6) showcases the variety of graphical outputs offered in mixOmics to complement each method. *Performance assessment in supervised analyses* (Chapter 7) describes the techniques employed to evaluate the results of the analyses.

Part III: mixOmics in action provides detailed case studies that apply each method in mixOmics to answer pertinent biological questions, complete with example R code and insightful plots. We begin with *mixOmics: get started* (Chapter 8) to guide the novice analyst in using the R platform for data analysis. Each subsequent chapter is dedicated to one method implemented in mixOmics. In *Principal Component Analysis (PCA)* (Chapter 9) and *PLS - Discriminant Analysis (PLS-DA)* (Chapter 12), we introduce different multivariate methods for single omics analysis. The *N*-integration framework is introduced in *PLS* (Chapter 10) and *Canonical Correlation Analysis* (Chapter 11) for two omics, and *N-data integration* (DIABLO, Chapter 13) for multi-omics integration. *P-data integration* (MINT, Chapter 14) introduces our latest developments for *P*-integration to combine independent omics studies. Each of these chapters is organised as follows:

- Aim of the method,
- Research question framed biologically and statistically,
- Principles of the method,
- Input arguments and key outputs,
- Introduction of the case study,
- Quick start R command lines,
- Further options to go deeper into the analysis,
- Frequently Asked Questions,
- Technical methodological details in each Appendix.

²<https://www.bioconductor.org/packages/release/bioc/html/mixOmics.html>

³<https://github.com/mixOmicsTeam/>

Additional resources related to this book

In addition to the R package, the mixOmics project includes a website with extensive tutorials in <http://www.mixOmics.org>. The R code of each chapter is also available on the website. Our readers can also register for our newsletter mailing list, and be part of the mixOmics community on GitHub and via our discussion forum <https://mixomics-users.discourse.group/>.



Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

Authors

Dr Kim-Anh Lê Cao develops novel methods, software and tools to interpret big biological data and answer research questions efficiently. She is committed to statistical education to instill best analytical practice and has taught numerous statistical workshops for biologists and leads collaborative projects in medicine, fundamental biology or microbiology disciplines. Dr Kim-Anh Lê Cao has a mathematical engineering background and graduated with a PhD in Statistics from the Université de Toulouse, France. She is currently an Associate Professor in Statistical Genomics at the University of Melbourne. In 2019, Kim-Anh received the Australian Academy of Science's Moran Medal for her contributions to Applied Statistics in multidisciplinary collaborations. She has contributed to a leadership program for women in STEMM, including the international Homeward Bound which culminated in a trip to Antarctica, and Superstars of STEM from Science Technology Australia.

Zoe Welham completed a BSc in molecular biology and during this time developed an interest in the analysis of big data. She completed a Master of Bioinformatics with a focus on the statistical integration of different omics data in bowel cancer. She is currently a PhD candidate at the Kolling Institute in Sydney where she is furthering her research into bowel cancer with a focus on integrating microbiome data with other omics to characterise early bowel polyps. Her research interests include bioinformatics and biostatistics for many areas of biology and making that information accessible to the general public.



Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

Part I

Modern biology and multivariate analysis



Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

Multi-omics and biological systems

Technological advances such as next-generation sequencing and mass spectrometry generate a wealth of diverse biological information, allowing for the monitoring of thousands of variables, or dimensions, that describe a given sample or individual, hence the term ‘high dimensional data’. Multi-omics variables represent molecules from different functional levels: for example, transcriptomics for the study of transcripts, proteomics for proteins, and metabolomics for metabolites. However, their complex nature requires an integrative, multidisciplinary approach to analysis that is not yet fully established.

Historically, the scientific community has adopted a *reductionist* approach to data analysis by characterising a very small number of genes or proteins in one experiment to assess specific hypotheses. A holistic approach allows for a deeper understanding of biological systems by adding two new facets to analysis (Figure 1.1): Firstly, by integrating data from different omic functional levels, we move from clarifying a linear process (e.g. the dysregulation of one or two genes) towards understanding the development, health, and disease of an ever-changing, dynamic, hierarchical *system*. Secondly, by adopting a hypotheses-free, data-driven approach, we can build integrated and coherent models to address novel, systems-level hypotheses that can be further validated through more traditional hypotheses.

1.1 Statistical approaches for reductionist or holistic analyses

Compared to a traditional reductionist analysis, multivariate multi-omics analysis drastically differs in its viewpoint and aims. We briefly introduce three types of analysis to illustrate this point:

A **univariate analysis** is a fundamentally reductionist, *hypothesis-driven* approach that is related to inferential statistics (introduced in Section 2.4). A hypothesis test is conducted on one variable (e.g. gene expression or protein abundance) independently from the other variables. Univariate methods make inferences about the population and measure the certainty of this inference through test statistics and p -values. Linear models, t -tests, F -tests, or non-parametric tests fit into a univariate analysis framework. Although interactions between variables are not considered in univariate analyses, when one variable is manipulated in a controlled experiment, we can often attribute the result to that particular variable. In omics studies, where multiple variables are monitored simultaneously, it is difficult to determine which variables influence the biology of interest.

A **bivariate analysis** considers two variables simultaneously, for example, to assess the association between the expression levels of two genes via correlation or linear regression. Such an analysis is often supported by visualisation through scatterplots but can quickly

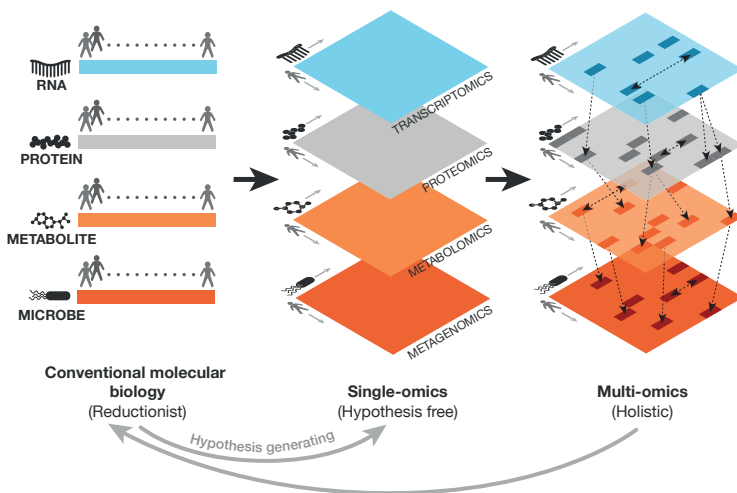


FIGURE 1.1: From reductionism to holism. Until recently, only a few molecules of a given omics type were analysed and related to other omics. The advent of high-throughput biology has ushered in an era of hypothesis-free approaches within a single type of omics data, and across multiple omics from the same set of samples. A holistic approach is now required to understand the different omic functional layers in a biological system and posit novel hypotheses that can be further validated with a traditional reductionist approach. We have omitted DNA as this data type needs to be handled differently in *mixOmics*, see [Section 4.2.3](#).

become cumbersome when dealing with thousands of variables that are considered in a pairwise manner.

A **multivariate analysis** examines more than two variables simultaneously and potentially thousands at a time. In omics studies, this approach can lead to computational issues and inaccurate results, especially when the number of samples is much smaller than the number of variables. Several computational and statistical techniques have been revisited or developed for high-dimensional data. This book focuses on multivariate analyses and extends this to include the integration of multi-omics data sets.

1.2 Multi-omics and multivariate analyses

The aim of omics data integration is to identify associations and patterns amongst different variables and across different omics collected from a sample. Provided appropriate data analysis is conducted, the integration of multiple data sources may also consolidate our confidence in the results when consensus is observed from different experiments.

1.2.1 More than a ‘scale up’ of univariate analyses

The fundamental difference between multivariate and univariate analysis lies in the scope of the results obtained. Multivariate analysis can unravel groups of variables that share similar patterns in expression across different phenotypes, thus complementing each other to describe an outcome. A univariate analysis may declare the same variables as non-significant, as a variable’s ability to explain the phenotype may be subtle and can be masked by individual variation, or confounders (Saccenti et al., 2014). However, with sufficiently powered data, univariate and multivariate methods are complementary and can help make sense of the data. For example, several multivariate and exploratory methods presented in this book can suggest promising candidate variables that can be further validated through experiments, reductionist approaches, and inferential statistics.

1.2.2 More than a fishing expedition

Multivariate analyses, which examine up to thousands of variables simultaneously, are often considered to be ‘fishing expeditions’. This somewhat pejorative term refers to either conducting analyses without first specifying a testable hypothesis based on prior research, or, conducting several different analyses on the same data to ‘fish’ for a significant result regardless of its domain relevance. Indeed, examining a large number of variables can lead to statistically significant results purely by chance.

However, the integration of multi-omics data, with an appropriate experimental design set in an exploratory, rather than predictive approach, offers a tremendous opportunity for discovering associations between omics molecules (whether genes, transcripts, proteins, or metabolites), in normal, temporal or spatial changes, or in disease states. For example, one of our studies identified pathways that were never previously identified as relevant to ontogeny during the first week of human life (Lee et al., 2019). Multi-omics data integration has deepened our understanding of gene regulatory networks by including information from related molecules prior to validation of gene associations with a functional approach (Gligorijević and Pržulj, 2015) and has also efficiently improved functional annotations to proteins instead of using expensive and time-consuming experimental techniques (Ma et al., 2013). Multi-omics can also more easily characterise the relatively small number of genes associated with a particular disease by integrating multiple sources of information (Žitnik et al., 2013). Finally, it has further developed precision medicine by integrating patient- and disease-specific information with the aim to improve prognosis and clinical outcomes (Ritchie et al., 2015).

1.3 Shifting the analysis paradigm

Despite the potential advantages of high-dimensional data, we should keep in mind that quantity does not equal quality. Multivariate data integration is not straightforward: the analyses cannot be reduced to a mere concatenation of variables of different types, or by overlapping information between single data sets, as we illustrate in Figure 1.2. As such, we must shift our traditional view of analysing data.

Biological experimentation often employs univariate statistics to answer clear hypotheses

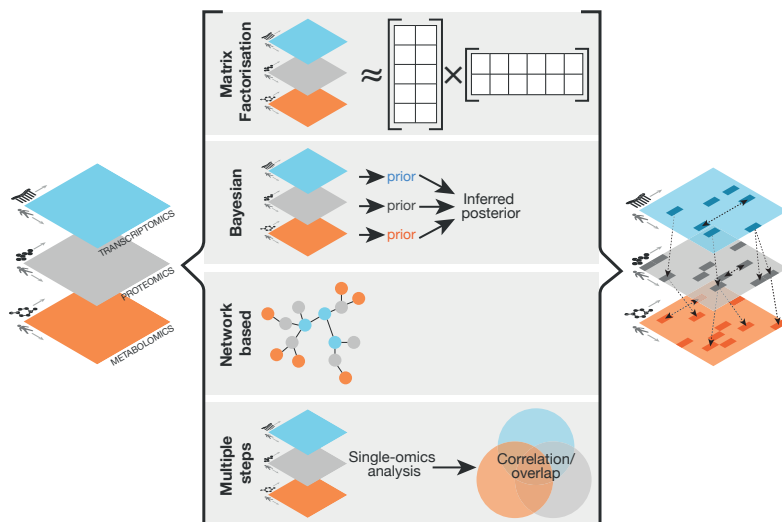


FIGURE 1.2: Types of methods for data integration. Methods for multi-omics data integration are still in active development, and can be broadly categorised into matrix factorisation techniques (the focus of this book), Bayesian, network-based, and multiple-step approaches. The latter deviates from data integration as it considers each data set individually before combining the results.

about the potential causal effect of a given molecule of interest. In high-dimensional data sets, this reductionist approach may not hold due to the sheer amount of molecules that are monitored, and their interactions that might be of biological interest. Therefore, exploratory, *data-driven* approaches are needed to extract information from noise and generate new hypotheses and knowledge. However, the lack of a clear, causal-driven hypothesis presents a challenging new paradigm in statistical analyses.

In univariate hypothesis testing, we report p -values to determine the significance of a statistical test conducted on a *single* variable. In a multivariate setting, however, a p -value assesses the statistical significance of a result while taking into account *all* variables simultaneously. In such analyses, permutation-based tests are common to assess how far from random a result is when the data are reshuffled, but other inference-based methods are currently being developed in the field of multivariate analysis (Wang and Xu, 2021). In *mixOmics* we do not offer such tests, but related methods propose permutation approaches to choose the parameters in the method (see Section 10.7.5).

1.4 Challenges with high-throughput data

There are multiple challenges associated with managing large amounts of biological data, pertaining to specific types of data as well as statistical analysis. To make reliable, valid, and meaningful interpretations, these challenges must be considered, ideally before data collection.

1.4.1 Overfitting

Multivariate omics analysis assesses many molecules that individually, or in combination, can explain the biological outcome of interest. However, these associations may be spurious, as the large number of features can often be combined in different ways to explain the outcome well, despite having no biological relevance. Overfitting occurs when a statistical model captures the noise along with the underlying pattern in the data: if we apply the same statistical model fitted on a high-dimensional data set to a similar but external study, we might obtain different results.¹ The problem of overfitting is a well-known issue in high-throughput biology (Hawkins, 2004). We can assess the amount of overfit using cross-validation or subsampling of the data, as described in Chapter 7.

1.4.2 Multi-collinearity and ill-posed problems

As the number of variables increase, the number of pairwise correlations also increases. *Multi-collinearity* poses a problem in most statistical analyses as these variables bring redundant and noisy information that decreases the precision of the statistical model. Correlations in high-throughput data sets are often spurious, especially when the number of biological samples, or individuals N , is small compared to the number of variables P ². The ‘small N large P ’ problem is defined as *ill-posed*, as standard statistical inference methods assume N is much greater than P to generalise the results to the population the sample was drawn from. Ill-posed problems also lead to inaccurate computations.

1.4.3 Zero values and missing values

Data sets may contain a large number of zeros, depending on the type of omics studied and the platform that is used. This is particularly the case for microbiome, proteomics, and metabolomics data: a large number of zeros results in zero-inflated (skewed) data, which can impair methods that assume a normal distribution of the data. *Structural zeros*, or true zeros, reflect a true absence of the variable in the biological environment while *sampling zeros*, or false zeros, may not reflect reality due to experimental error, technological reasons, or an insufficient sample size (Blasco-Moreno et al., 2019). The challenge is whether to consider these zeros as a true zero or missing (coded as NA in R).

Methods that can handle missing values often assume they are ‘missing at random’, i.e. missingness is not related to a specific sample, individual, molecule, or type of omics platform. Some methods can estimate missing values, as we present in Appendix 9.A.

¹Statistical models that overfit have low bias and high variance, meaning that they tend to be complex to fit the training data well, but do not predict well on test data (more details about the bias-variance tradeoff can be found in Friedman et al. (2001) Chapter 2).

²In our context, N can also refer to the number of cells in single cell assays, as we briefly mention in Section 14.6.

1.5 Challenges with multi-omics integration

Examining data holistically may lead to better biological understanding, but integrating multiple omics data sets is not a trivial task and raises another series of challenges.

1.5.1 Data heterogeneity

Different omics rely on different laboratory techniques and data extraction platforms, resulting in data sets of different formats, complexity, dimensionalities, information content, and scale, and may be processed using different bioinformatics tools. Therefore, data heterogeneity arises from biological *and* technical reasons and is the main analytical challenge to overcome.

1.5.2 Data size

Integrating multiple omics results in a drastic increase in the number of variables. A filtering step is often applied to remove irrelevant and noisy variables (see [Section 8.1](#)). However, the number of variables P still remains extremely large compared to the number of samples N , which raises computational as well as analytical issues.

1.5.3 Platforms

The data integration field is constantly evolving due to ever-advancing technologies with new platforms and protocols, each containing inherent technical biases and analytical challenges. It is crucial that data analysts swiftly adapt their analysis framework to keep apace with these omics-era demands. For example, single cell techniques are rapidly advancing, as are new protocols for their multi-omics analysis.

1.5.4 Expectations for analysis

The field of data integration has no set definition. Data integration can be managed biologically, bioinformatically, statistically, or at the interpretation steps (i.e. by overlapping biological interpretation once the statistical results are obtained). Therefore, the expectations for data integration are diverse; from exploration, and from a low to high-level understanding of the different omics data types. Despite recent advances in single cell sequencing, current technologies are still limited in their ability to parse omics interactions at precise functional levels. Thus, our expectations for data integration are limited, not only by the statistical methods but also by the technologies available to us.

1.5.5 Variety of analytical frameworks

Integrative techniques fully suited to multi-omics biological data are still in development and continue to expand³. Different types of techniques can be considered and broadly categorised into (Huang et al. (2017), Figure 1.2):

- Matrix factorisation techniques, where large data sets are decomposed into smaller sub-matrices to summarise information. These techniques use algebra and analysis to optimise specific statistical criteria and integrate different levels of information. Methods in `mixOmics` fit into this category and will be detailed in Chapter 3 and subsequent chapters,
- Bayesian methods, which use assumptions of prior distributions for each omics type to find correlations between data layers and infer posterior distributions,
- Network-based approaches, which use visual and symbolic representations of biological systems, with nodes representing molecules and edges as correlations between molecules, if they exist. Network-based methods are mostly applied for detecting significant genes within pathways, discovering sub-clusters, or finding co-expression network modules,
- Multiple-step approaches that first analyse each single omics data set individually before combining the results based on their overlap (e.g. at the gene level of a molecular signature) or correlation. This type of approach technically deviates from data integration but is commonly used.

1.6 Summary

Modern biological data are high dimensional; they include up to thousands of molecular entities (e.g. genes, proteins, or epigenetic markers) per sample. Integrating these rich data sets can potentially uncover the hierarchical and holistic mechanisms that govern biological pathways. While classical, reductionist, univariate methods ignore these molecular interactions, multivariate, integrative methods offer a promising alternative to obtain a more complete picture of a biological system. Thus, univariate and multivariate methods are different approaches with very little overlap in results but have the advantage of complementarity.

The advent of high-throughput technology has revealed a complex world of multi-omics molecular systems that can be unraveled with appropriate integration methods. However, multivariate methods able to manage high-dimensional and multi-omics data are yet to be fully developed. The methods presented in this book mitigate some of these challenges and will help to reveal patterns in omics data, thus forging new insights and directions for understanding biological systems as a whole.

³A comprehensive list of multi-omics methods and software is available at <https://github.com/mikelove/awesome-multi-omics>.



Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

The cycle of analysis

The *Problem, Plan, Data, Analysis, Conclusion* (PPDAC) cycle is a useful framework for answering an experimental question effectively (Figure 2.1). The mixOmics project emphasises crafting a well-defined biological question (Chapter 4), as this guides data acquisition and preparation (Chapter 8), as well as choosing appropriate multivariate techniques for analysis (Chapter 4). Although this book is focused on analysis and interpretation, careful consideration of each step will maximise a successful analytical outcome.



FIGURE 2.1: PPDAC. The *Problem, Plan, Data, Analysis, Conclusion* cycle proposed by MacKay and Oldford (2000) will guide our multivariate analysis process.

2.1 The *Problem* guides the analysis

Multivariate analysis is appropriate for large data sets where the biological question encompasses a broad domain, rather than parsing the action of a single or small number of variables. Thus, we often require a hypothesis-free investigation based on a *data-driven* approach. However, this does not imply that multivariate analysis is a fishing expedition with no underlying biological question. The experimental design, driven by a well-formulated biological question and the choice of statistical method, will ensure a successful analysis (Shmueli, 2010). Chapter 4 lists several types of biological questions that can be answered with multivariate and integrative methods.