

Integrative Analysis of Omics Data with Biological Knowledge in Translational Medicine



UNIVERSITAT DE
BARCELONA

Ferran Briansó

Facultat de Biologia

Departament de Genètica, Microbiologia i Estadística

Universitat de Barcelona

A thesis submitted for the degree of

Doctor of Philosophy

XXXX XX 2023

For XXXXX XXXXXX

Acknowledgements

This is where you will normally thank your advisor, colleagues, family and friends, as well as funding and institutional support. In our case, we will give our praises to the people who developed the ideas and tools that allow us to push open science a little step forward by writing plain-text, transparent, and reproducible theses in R Markdown.

We must be grateful to John Gruber for inventing the original version of Markdown, to John MacFarlane for creating Pandoc (<http://pandoc.org>) which converts Markdown to a large number of output formats, and to Yihui Xie for creating **knitr** which introduced R Markdown as a way of embedding code in Markdown documents, and **bookdown** which added tools for technical and longer-form writing.

Special thanks to [Chester Ismay](#), who created the **thesisdown** package that helped many a PhD student write their theses in R Markdown. And a very special thanks to John McManigle, whose adaption of Sam Evans' adaptation of Keith Gillow's original maths template for writing an Oxford University DPhil thesis in LaTeX provided the template that I in turn adapted for R Markdown.

Finally, profuse thanks to JJ Allaire, the founder and CEO of [RStudio](#), and Hadley Wickham, the mastermind of the tidyverse without whom we'd all just given up and done data science in Python instead. Thanks for making data science easier, more accessible, and more fun for us all.

Ferran Brianso
Mataro, BCN
XX XXXXXX 2023

Abstract

The general concept of Data Integration can be defined as the combination of data residing in different sources in order to provide the users with a unified view of these data [1]. However, the practical meaning of the term Integration may vary from, for instance, the computational combination of data, to the combination of studies performed independently, the simultaneous analysis of multiple variables on multiple datasets, or any possible approach for homogeneously querying heterogeneous data sources. Therefore, in many cases, an integrative analysis may be preferable than a simple combination of data from distinct sources. Integrative analysis allows not only for the combination of heterogeneous data, but also for the combined use of these data in order to get the most relevant information and, what is better, to be able to extract some information that could not be unveiled by the separated analysis of each of the original data types.

Over the past decade, advancements in omics technologies have facilitated the high-throughput monitoring of molecular and organism processes. These techniques have been widely applied to identify biological agents and to characterize biochemical systems, often focusing on the discovery of therapeutic targets and biomarkers related with specific diseases [2,3,4]. While many single-omic approaches target comprehensive analysis of genes (genomics), mRNA (transcriptomics), proteins (proteomics), and metabolites (metabolomics) among other, there is still field to improve omics data analyses through integrative methods [5,6]. In this sense, the integrative point of view defined in the paragraph above, applied to multi-omics data, is a promising approach to achieve better biomarker development in biomedical research projects, and this is the core idea of this work.

As the field of omics has evolved from analyzing a unique type of data to multiple types, it has been natural to extend the previous use of multivariate techniques to this new situation. With this aim classical and new multivariate techniques have been applied to the analysis of multi-omics datasets. Many of these techniques are dimension reduction methods that aim at finding main sources of variability in the data while maximizing some information characteristic such as the variance of each dataset, the correlation between groups of variables or other. Examples of such techniques are well consolidated methods such as Principal

Component Analysis (PCA), Singular Value Decomposition (SVD), Correspondence Analysis (CA), and Partial Least Squares (PLS). Besides these more “novel” approaches have been used such as: Principal Components Regression, Coinertia and Multiple Coinertia Analysis, Generalized SVD, Sparse PLS, Multiple Factor Analysis (MFA), or combined versions of them [7,8,9]. Meng [10], Cavill [11], Wu [12], Subramanian [30], Krassowski [31], and Cantini [32], are good reviews of the state of the art of using multivariate and joint reduction methods for Integrative Multi-Omics Analysis.

Dimension reduction methods, especially those that are able to deal with situations that are typical from the omics context (with many more variables than samples, or possibly sparse matrices with many missing values), have been of great help in visualizing datasets or even for performing variable selection to find biomarkers for a given situation [12]. There is however one point where they underperform other approaches, that is, the difficulty in interpreting results from a biological point of view. This is relatively reasonable, because the most of these methods work by creating new variables that are some type of linear combination from the original ones. While this is useful, for example, for removing redundancy, this does not provide any clues on what these new dimensions may mean from a biological point of view.

This problem has been known since the beginning of using multivariate methods with omics data, but only a few approaches have been taken to deal with this. The first attempts to introduce biological information in the analyses consisted of using the most well-known database of biological functions, the Gene Ontology (GO) [13]. Fellenberg [14] introduces a way to integrate Gene Ontology information with Correspondence Analysis to facilitate the interpretation of microarray data. De Tayrac et al. [15] applies multiple factor analysis to the integrative analysis of microarray and DNA copy number data. They apply GO Terms on data visualizations by treating these terms as supplemental information. In recent years the representation of biological knowledge has shifted from Gene Ontology to using Gene Sets [16]. Meng and Culhane [10] have introduced the Integrative Clustering with Gene Set Analysis where gene set expression analysis is performed based on multiple omics data; and Tyekucheva et al. [17], go one step further and use the results of Gene Set Expression Analysis (GSEA) to integrate different omics data.

Altogether, the previous approaches show several things: Although the idea that integrating quantitative data with biological knowledge may increase interpretability, the number of successful attempts to do this is still small. In this thesis, the use of either classical GO Terms or more flexible annotations (Gene Sets or custom annotations), will be combined with different approaches, and combinations of them if needed, to guide integrative analysis and to improve its biological interpretability from the point of view of the biomedical researchers.

Contents

List of Figures	viii
List of Tables	ix
List of Abbreviations	x
1 Introduction	1
1.1 Content of the introductory text (WIP)	1
1.2 Background/State of the Art	4
2 Objectives	6
3 Methodology	8
3.1 Working phases	8
3.2 Explanation of the methods	11
4 Results	18
4.1 Results from the analysis of human brain tissue samples	18
4.2 Results from the expansion of omics data with biological annotations	18
4.3 Results from the analysis of 150 TCGA-BRCA samples	19
4.4 Results from the application of MFA on TCGA-BRCA data with, and without, expanded data	20
4.5 Resultats de la creacio del paquet amb Targets...	21
5 Discussion	22
6 Conclusions	23
Conclusion 1	23
Conclusion 2	23
Conclusion 3	23
Conclusion 4	24

Contents

Appendices

A The First Appendix **26**

B The Second Appendix, for Fun **27**

References **28**

List of Figures

3.1	Addition of GO terms	12
3.2	Addition of news feats	13
3.3	Gene enrichment diagram	13
3.4	Matrix expansion diagram	14
3.5	Addition of new feats (2)	15
3.6	Matrix expansion diagram (2)	16
3.7	Workflow overview	16
4.1	Heapmap of an expanded matrix	19
4.2	BRCA results overview	20
4.3	BRCA results with MFA	21

List of Tables

List of Abbreviations

- 1-D, 2-D** . . . One- or two-dimensional, referring **in this thesis** to spatial dimensions in an image.
- Otter** One of the finest of water mammals.
- Hedgehog** . . . Quite a nice prickly friend.

1

Introduction

Contents

1.1	Content of the introductory text (WIP)	1
1.2	Background/State of the Art	4
1.2.1	Omics data analyses	4
1.2.2	Integrative analyses	4
1.2.3	Review of existing approaches for multi-omics data integration	5
1.2.4	Revisió de metodes de creacio pipelines	5

1.1 Content of the introductory text (WIP)

The general concept of Data Integration can be defined as the combination of data residing in different sources in order to provide the users with a unified view of these data [1]. However, the practical meaning of the term Integration may vary from, for instance, the computational combination of data, to the combination of studies performed independently, the simultaneous analysis of multiple variables on multiple datasets, or any possible approach for homogeneously querying heterogeneous data sources. Therefore, in many cases, an integrative analysis may be preferable than a simple combination of data from distinct sources. Integrative analysis allows not only for the combination of heterogeneous data, but also for

1. Introduction

the combined use of these data in order to get the most relevant information and, what is better, to be able to extract some information that could not be unveiled by the separated analysis of each of the original data types.

Over the past decade, advancements in omics technologies have facilitated the high-throughput monitoring of molecular and organism processes. These techniques have been widely applied to identify biological agents and to characterize biochemical systems, often focusing on the discovery of therapeutic targets and biomarkers related with specific diseases [2,3,4]. While many single-omic approaches target comprehensive analysis of genes (genomics), mRNA (transcriptomics), proteins (proteomics), and metabolites (metabolomics) among other, there is still field to improve omics data analyses through integrative methods [5,6]. In this sense, the integrative point of view defined in the paragraph above, applied to multi-omics data, is a promising approach to achieve better biomarker development in biomedical research projects, and this is the core idea of this work.

As the field of omics has evolved from analyzing a unique type of data to multiple types, it has been natural to extend the previous use of multivariate techniques to this new situation. With this aim classical and new multivariate techniques have been applied to the analysis of multi-omics datasets. Many of these techniques are dimension reduction methods that aim at finding main sources of variability in the data while maximizing some information characteristic such as the variance of each dataset, the correlation between groups of variables or other. Examples of such techniques are well consolidated methods such as Principal Component Analysis (PCA), Singular Value Decomposition (SVD), Correspondence Analysis (CA), and Partial Least Squares (PLS). Besides these more “novel” approaches have been used such as: Principal Components Regression, Coinertia and Multiple Coinertia Analysis, Generalized SVD, Sparse PLS, Multiple Factor Analysis (MFA), or combined versions of them [7,8,9]. Meng [10], Cavill [11], Wu [12], Subramanian [30], Krassowski [31], and Cantini [32], are good reviews of the state of the art of using multivariate and joint reduction methods for Integrative Multi-Omics Analysis.

1. Introduction

Dimension reduction methods, especially those that are able to deal with situations that are typical from the omics context (with many more variables than samples, or possibly sparse matrices with many missing values), have been of great help in visualizing datasets or even for performing variable selection to find biomarkers for a given situation [12]. There is however one point where they underperform other approaches, that is, the difficulty in interpreting results from a biological point of view. This is relatively reasonable, because the most of these methods work by creating new variables that are some type of linear combination from the original ones. While this is useful, for example, for removing redundancy, this does not provide any clues on what these new dimensions may mean from a biological point of view.

This problem has been known since the beginning of using multivariate methods with omics data, but only a few approaches have been taken to deal with this. The first attempts to introduce biological information in the analyses consisted of using the most well-known database of biological functions, the Gene Ontology (GO) [13]. Fellenberg [14] introduces a way to integrate Gene Ontology information with Correspondence Analysis to facilitate the interpretation of microarray data. De Tayrac et al. [15] applies multiple factor analysis to the integrative analysis of microarray and DNA copy number data. They apply GO Terms on data visualizations by treating these terms as supplemental information. In recent years the representation of biological knowledge has shifted from Gene Ontology to using Gene Sets [16]. Meng and Culhane [10] have introduced the Integrative Clustering with Gene Set Analysis where gene set expression analysis is performed based on multiple omics data; and Tyekucheva et al. [17], go one step further and use the results of Gene Set Expression Analysis (GSEA) to integrate different omics data.

Altogether, the previous approaches show several things: Although the idea that integrating quantitative data with biological knowledge may increase interpretability, the number of successful attempts to do this is still small. In this thesis, the use of either classical GO Terms or more flexible annotations (Gene Sets or custom annotations), will be combined with different approaches, and combinations of them

1. Introduction

if needed, to guide integrative analysis and to improve its biological interpretability from the point of view of the biomedical researchers.

1.2 Background/State of the Art

Falta desenvolupar punts

1.2.1 Omics data analyses

3 problemes essencials (veure projecte recerca Alex):

- Omics data may be partly incomplete, especially in multiomics studies, where not all types of data are usually available for all individuals.
- The results of these analyses are difficult to interpret. If we agree that the ultimate goal of many analyzes is a better understanding of the underlying biological processes, for example, in a disease study context, it should be possible to establish a clear relationship between the outcome of an analysis and what this means biologically. And this is not always so.
- These kind of data analytics are difficult to standardize, as it is not easy to make complex pipelines of multi-omics analyses, which integrate multiple processes with multiple sources, easy to reproduce or communicate.

Més el problema de la p»n (Dimensionality Reduction Techniques; The p»n situation)

1.2.2 Integrative analyses

Allows the combination of distinct omics data.

The blind men and the elephant https://en.wikipedia.org/wiki/Blind_men_and_an_elephant

Interpretability is a weak point of most multi omics approaches.

Methods focus much more on feature selection discovery and interaction highlighting measurement than on clinical or biological interpretability.

1. Introduction

1.2.3 Review of existing approaches for multi-omics data integration

MCIA, RGCCA, MFA... Cavill, 2016; Culhane 2003...

1.2.4 Revisió de metodes de creacio pipelines

2

Objectives

The main objectives of this work are the following:

1. To make an empirical comparison of some of the currently available dimension reduction techniques applied for the integration of omics data, focused on their ability to include biological annotations,
2. To develop methods and workflows able to apply these techniques, focusing on the matching of distinct omics datasets relying on biological knowledge,
3. To apply these methods to specific translational biomedical research cases, such as an integrative analysis of transcriptomics and proteomics data to study ischemic stroke, as well as to public datasets, which can be easily shared and are not as restricted by sample sizes as other projects.
4. To implement the knowledge acquired with this work into the appropriate bioinformatics tools, e.g. R packages or web-based tools, that will be used in future biomedical research projects for providing a better interpretation of this kind of studies.

All these objectives are in agreement with the tasks defined within a project partially supported by Grant MTM2015-64465-C2-1-R (MINECO/FEDER) from

2. Objectives

the Ministerio de Economía y Competitividad (Spain), to which the PhD Thesis proposed here is related.

*Ein Mann, der recht zu wirken denkt,
 Muß auf das beste Werkzeug halten*
 The man who seeks to be approved,
 must stick to the best tools for it
 — Goethe’s *Faust. Eine Tragödie* (1808).

3

Methodology

Contents

3.1	Working phases	8
3.2	Explanation of the methods	11
3.2.1	Detail of the integrative data analyses applied.	14
3.2.2	Comparison of ODA results	14
3.2.3	Numeric measurement	14
3.2.4	Biological interpretation	16
3.2.5	Targets PIPELINE concept	16

3.1 Working phases

Working phases, with the corresponding steps, followed in order to achieve the above objectives:

1. Application of integrative multi-omics methods to (I) the analysis of specific data sets provided by research units from our former affiliation center, VHIR, and other research institutions that we collaborate with [34, 36, 37] and (II) to the integrative analysis of larger data sets from public data bases, such as Breast Cancer samples from the TCGA project [18, 19].

3. Methodology

2. Development of methods, either in terms of new algorithms or in terms of combinative workflows, which will be able to improve, and facilitate, the analysis and biological interpretation of those data sets to be integrated.
3. Implementation of the methods developed for this study in the appropriate bioinformatics tools, such as an R package or a web-based application, to facilitate their use in the context of biomedical research projects.

Here follows a brief description of these main five activities, the methods in which they are initially based, the objectives that they are related to, and the corresponding results:

1. Application of some state-of-the-art methods for integrative multi-omics data analysis to the study of human brain tissue samples, collected by the Neurovascular Diseases Laboratory at Vall d'Hebron Research Institute. This part is already finished, and led to publications in 2018 and 2021 [37, 38]. Researchers obtained different omics data from necropsies, which had been processed to obtain mRNA, microRNA and protein expression values. Each dataset had been first analyzed independently using standard bioinformatics protocols [20]. These analyses allowed selecting subsets of relevant features, for each type of data, to be used in the integrative analysis. Among all available options, we decided to use two distinct and complementary approaches: (I) Multiple Co-inertia Analysis implemented in Bioconductor packages *made4* [21] and *mogsa* [22], and (II) Regularized Canonical Correlation Analysis with Sparse Partial Least Squares regression (sPLS), provided by *mixomics* R package [23]. This work had been presented at some meetings [39, 40, 41, 43] and in an already published extended abstract's series book [35]. This step had been obviously useful for the achievement of the objective number 3 explained in the previous section, which aims on the study of the regulome's response to ischemic stroke, but also useful for detecting the advantages and drawbacks of the methods applied, thus setting the basis for the work regarding to objective number 2.

3. Methodology

2. Reproduction of the same analyses steps performed in point 1) above with publicly available databases, such as distinct omics data from 150 samples from the TCGA-BRCA collection. This data set contains the expression or abundance of mRNA, miRNA and proteomics for 150 breast cancer samples previously prefiltered, as explained in Rohart et al. [29], and allows identifying a good multi-omics signature to discriminate between Basal, Her2 and Luminal A breast cancer subtypes. This work is already finished, and complies with objectives 3 and 2.
3. Use of all the data sets analyzed up to this point to make a comparison of results between the main implemented methods, and eventually some others, which is the aim of objective 1. This is based on quantitative and qualitative comparison and visualization methods, such as those explained by Thallinger [24] and Martin [25], going from simple Venn diagrams to more complex, network analysis, software such as some specific R packages [20] or Cytoscape [26]. The focus here is to use graphical visualization elements to compare the results of the analyses with and without the addition of biological information.
4. Development of new methods and/or workflows in order to improve and/or combine the benefits from the selected approaches, with focus in those allowing the addition of biological significance to the integration process. Here follows an overview of the methods developed to expand the original datasets (X, Y) with annotations (Ax, Ay) to obtain new blocks of data $(Nx, Ny, \text{and } Nxy)$. And the workflow has been implemented adapting the integrative pipelines applied so far to the R targets package [33], a pipeline toolkit that improves reproducibility, skipping unnecessary steps already up to date and showing tangible evidence that the results match the underlying code and data. The development of this targets workflow is intended to comply with the objective number 2 of this working plan.

3. Methodology

5. Implementation of the methods resulting from 4) as a new R package to be submitted to Bioconductor repository [27], and, finally, to complete objective 4 of this thesis plan, as a web application [28] to be used in further steps of the current biomedical research projects in which our collaborators are implied, as well as in future studies.

3.2 Explanation of the methods

The addition of biological annotations to the data sets being integrated, prior to the integrative analysis itself, can be useful to improve the integration/analysis outcomes as well as their biological interpretability.

Passos principals explicats aquí:

A. Pre process omics datasets in order to include biological information before the joint analysis → Expanded datasets

B. Analysis of the expanded datasets by the use of contrasted joint Dimensionality Reduction techniques

C. Process semi automation in ease to use tools

Start the process already having a couple [punt de millora: admetre 3 o + inputs] of data sets from distinct 'omics sources, mapped to gene ids (if GO annotation has to be performed), containing the results from a selection of differentially expressed genes or most relevant proteins analysis, or similar. [explicar aquí els requeriments de format dels data sets d'entrada!!]

For each input data set, if annotations are not already provided, two distinct basic annotation methods can be performed:

- (i) a basic GO mapping, returning annotations to those GO entities for which we find more than a certain number of features (gene ids coming from our data set) annotated to them, [mostrar formula] [mostrar exemple]
- (ii) a Gene Enrichment Analysis (based on Hypergeometric tests against all GO categories, with FDR correction[ref clusterProfiler]) is performed in order to

3. Methodology

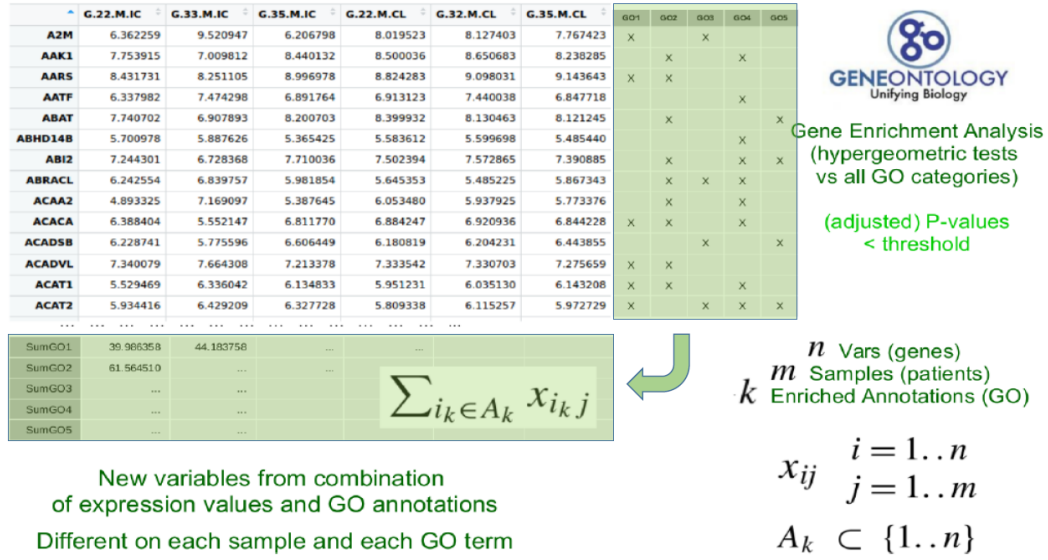


Figure 3.1: Addition of GO terms

retrieve the most relevant annotations to that set of genes/features. [mostrar exemple] [afegir aquí la opció d'afegir les anotacions com a individus suplementaris enlloc de variables]

Figure 3.1 is an example.

Alternatively, manual annotations can be provided (eg. GO terms, canonical pathways, or even annotation to custom entities) as an optional input file. [mostrar el format requerit].

Other annotation methods can be implemented, as functions to be used by the main pipeline, if more complex methods for biological information addition are required.

[Mostrar el format final de les anotacions, com a matrius dels data sets amb anotacions binàries 1/0 com a columnes extra]

Once the annotations are already computed, mapping each feature of the input data set to the corresponding biological entity, they can be used to generate new features (as new rows), computing the average value [punt de millora: £funció de ponderació?] of the expression/intensity values from all original features being mapped to the annotated biological entities.

3. Methodology

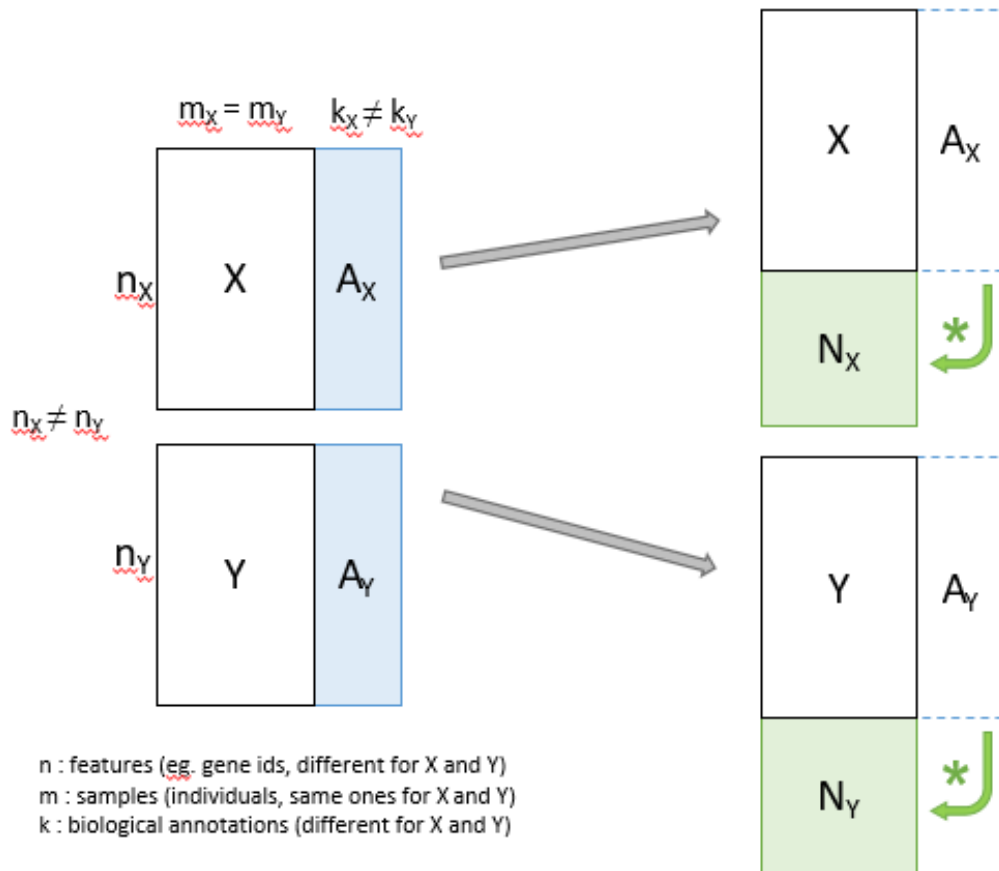


Figure 3.4: Matrix expansion diagram

3.2.1 Detail of the integrative data analyses applied...

Mètodes:

1- Significació biològica, com faig les anotacions 2- Expansió de les matrius (creació de noves vars a partir de les anotacions) 3- Anàlisi factorial en detall, + MCIA + RGCCA 4- TFM sobre workflows i automatització -> Paquet targets en general

3.2.2 Comparison of ODA results

3.2.3 Numeric measurement

% variabilitat explicat segons la estructura de la intersecció de les 2 taules

3. Methodology

	s1	s2	s3	s4	s5	s6	s7
** g1	9.823383	9.981509	10.911122	9.249729	11.459667	8.369157	9.560181
** g2	9.987826	8.301891	8.509732	9.330750	9.836765	8.369157	9.560181
** g3	9.750415	7.358170	8.204226	9.889993	10.679750	10.266810	10.266810
** g4	12.027450	9.340945	9.826619	9.893645	10.266810	10.266810	10.266810
** g5	7.799825	9.587443	9.410317	10.238429	10.841196	9.231603	9.231603
** g6	10.541859	10.960931	8.932551	12.739694	9.231603	9.231603	9.231603
** g7	9.886552	9.943082	10.046111	9.652113	11.450464	10.266810	10.266810
** g8	9.465699	10.460557	10.687289	9.395601	10.655185	9.231603	9.231603
** g9	8.774980	9.968002	9.973178	9.192533	9.684703	9.231603	9.231603
** g10	8.987790	9.486889	9.925574	9.981692	12.210547	9.231603	9.231603
** g11	9.897302	8.036014	10.687557	9.099134	8.807929	11.450464	10.266810
** g12	9.751261	12.231207	11.448632	8.127683	11.126290	9.231603	9.231603
** g13	7.541331	9.924513	11.811728	10.828893	7.744879	10.266810	10.266810
** g14	10.052642	9.143740	9.786817	9.392912	10.074336	10.266810	10.266810
** g15	11.119820	9.750373	11.246757	8.349641	7.301828	9.231603	9.231603
** g16	9.831945	10.305640	11.664353	8.938298	10.135373	10.266810	10.266810
** g17	9.697817	9.560181	8.354711	0	1	1	0
** g18	10.343394	9.341271	10.030345	1	0	1	0
** g19	10.658743	10.277298	11.329309	1	1	0	0
** g20	10.042697	10.945357	9.608662	0	0	0	1
** g1	10.970871	9.069306	1	1	0	0	1
** g2	11.045267	8.997348	1	0	1	1	1
** g3	11.308466	10.900477	0	1	1	1	0
** g4	9.390803	9.727783	1	0	1	0	1
** g5	11.373304	10.757272	0	1	1	1	0
** g6	10.080457	11.275904	0	1	0	1	0
** g7	10.799194	11.209403	0	1	1	1	0
** g8	9.337232	12.705092	1	1	0	1	1
** g9	9.625649	8.803367	0	0	1	0	0
** g10	10.552339	8.632491	1	1	0	0	1
** g11	9.757147	10.776161	1	0	0	0	1
** g12	9.613865	10.929590	1	0	1	0	0
** g13	8.880190	10.200053	0	1	0	1	0
** g14	11.150111	11.109901	1	0	1	1	0
** g15	10.642895	10.191623	0	0	1	1	1
** g16	9.734463	9.151632	1	1	0	1	0
** A1	111.193928	112.450958	11	Na	Na	Na	Na
** A2	122.280059	123.553432	12	Na	Na	Na	Na
** A3	115.672151	112.602470	11	Na	Na	Na	Na
** A4	132.611944	137.078157	13	Na	Na	Na	Na
** A5	72.494721	70.615121	7	Na	Na	Na	Na
** A6	71.907199	69.407572	7	Na	Na	Na	Na
** A7	81.193258	80.569283	8	Na	Na	Na	Na
** A8	78.816520	79.934694	8	Na	Na	Na	Na
** A9	89.100803	91.304790	9	Na	Na	Na	Na
** A10	103.283346	103.030921	10	Na	Na	Na	Na

$$N_G = \phi(G, A_G), \quad G \in X, Y, \dots,$$

$$N_X = \frac{1}{r_X}(X \times A'_X),$$

$$N_Y = \frac{1}{r_Y}(Y \times A'_Y),$$

Figure 3.5: Addition of new feats (2)

3. Methodology

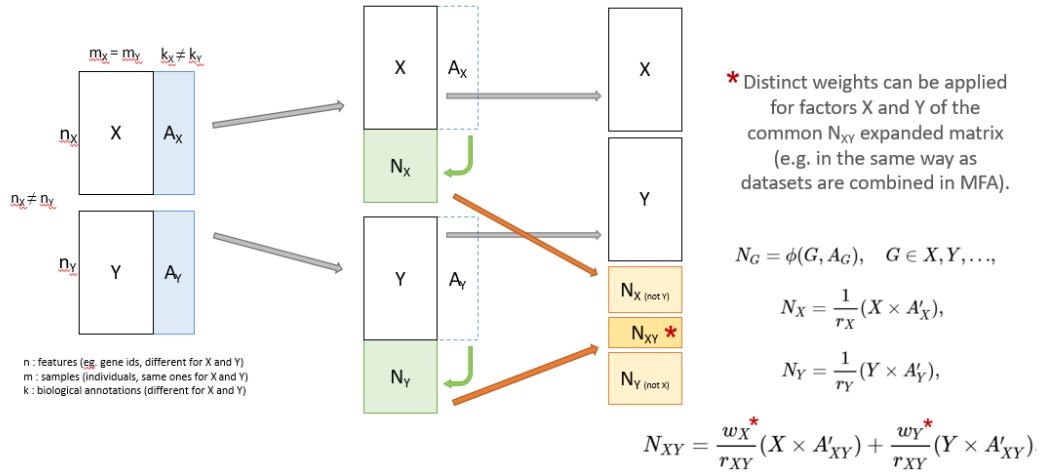


Figure 3.6: Matrix expansion diagram (2)

Reproducible workflow with a make-like pipeline toolkit *targets* package (<https://books.ropensci.org/targets/>)

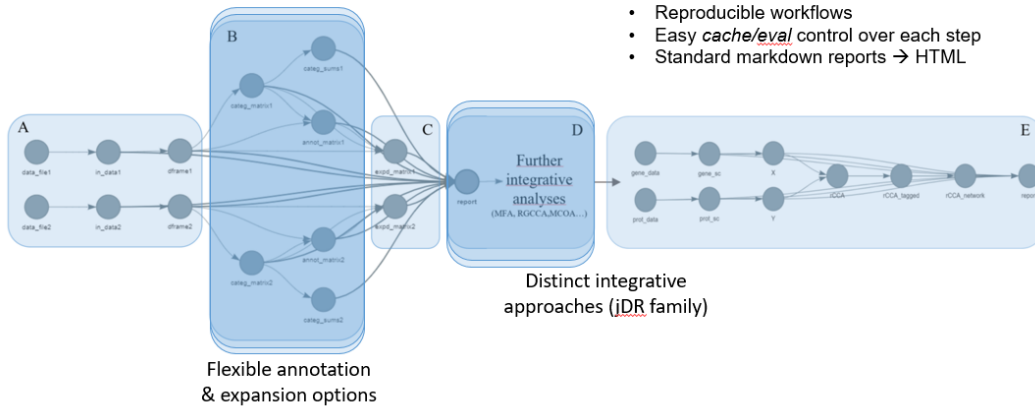


Figure 3.7: Workflow overview

3.2.4 Biological interpretation

3.2.5 Targets PIPELINE concept

R package creation...

Sistema que hem aplicat per crear el pipeline amb Targets...

Targets workflow diagram (Figure 3.7) showing the steps corresponding with the complete process: The pipeline starts from (A) a couple of 'omics-derived input data sets (e.g. pre-processed gene expression and protein abundance matrices). These are converted to R data frames with features in rows and samples in columns. Then, a data frame containing related annotations (B) is created, or loaded, for

3. Methodology

each given input matrix, and used to expand these original data, in order to end up with a pair of data frames (C) containing the original values plus the average expression/abundance values of the features related to each annotation as new features in additional rows. After that, distinct Dimension Reduction Methods are applied to perform the integrative analysis (D), and finally, an R markdown report (E) is rendered to show steps and main results of the full process.

4

Results

Contents

4.1	Results from the analysis of human brain tissue samples	18
4.2	Results from the expansion of omics data with biological annotations	18
4.3	Results from the analysis of 150 TCGA-BRCA samples	19
4.4	Results from the application of MFA on TCGA-BRCA data with, and without, expanded data	20
4.5	Resultats de la creacio del paquet amb Targets... . .	21

Text de presentacio dels resultats...

4.1 Results from the analysis of human brain tissue samples

4.2 Results from the expansion of omics data with biological annotations

Figure 4.1 is an snapshot (F) of one of the heat maps created to show the expanded matrices obtained in (Figures 3.4 i 3.5 prèvies, de Methods).

4. Results

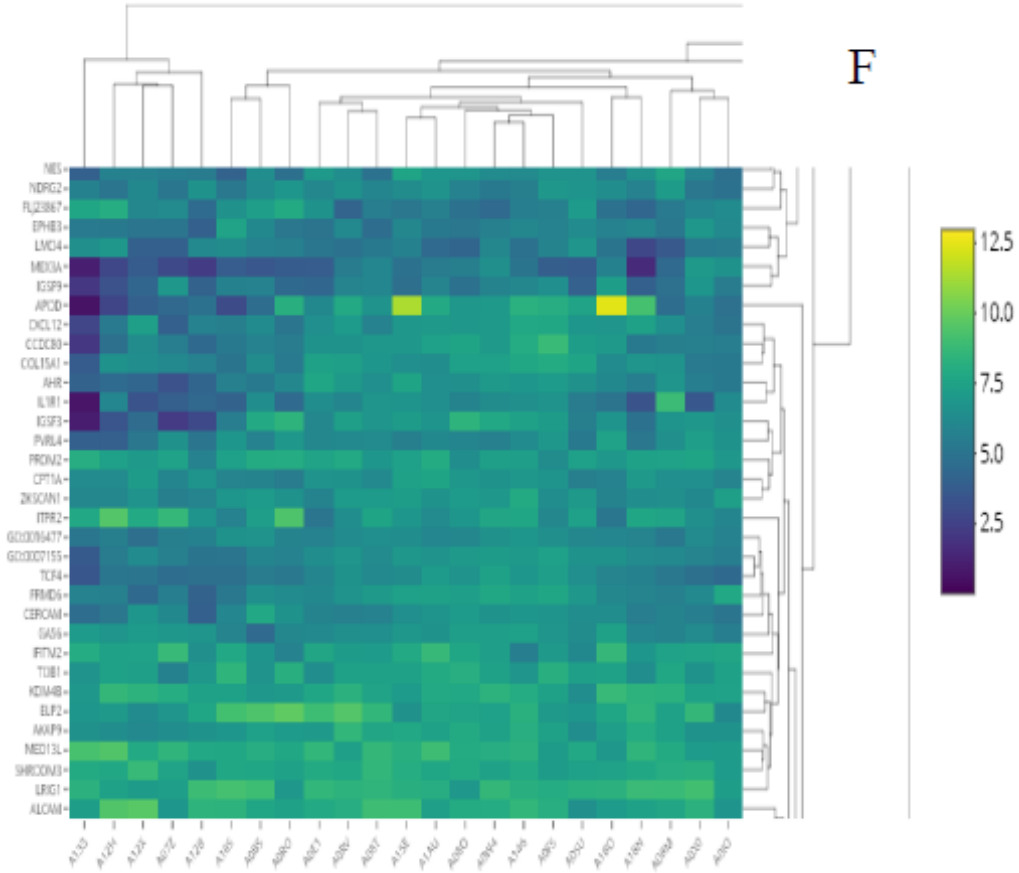


Figure 4.1: Heatmap of an expanded matrix

4.3 Results from the analysis of 150 TCGA-BRCA samples

Figure 4.2 contains some of the graphical results of the analysis of the 150 samples from TCGA-BRCA: Heat maps (A, C) and association networks (B, D) resulting from the integration by Regularized Canonical Correlations Analysis with mixomics R package. Performed with the original data sets (A, B) or using data expanded with biological annotations to Gene Ontology (C, D), so adding some GO terms to the features from each source, where the outputs contain higher level of information (higher density in both type of plots).

4. Results

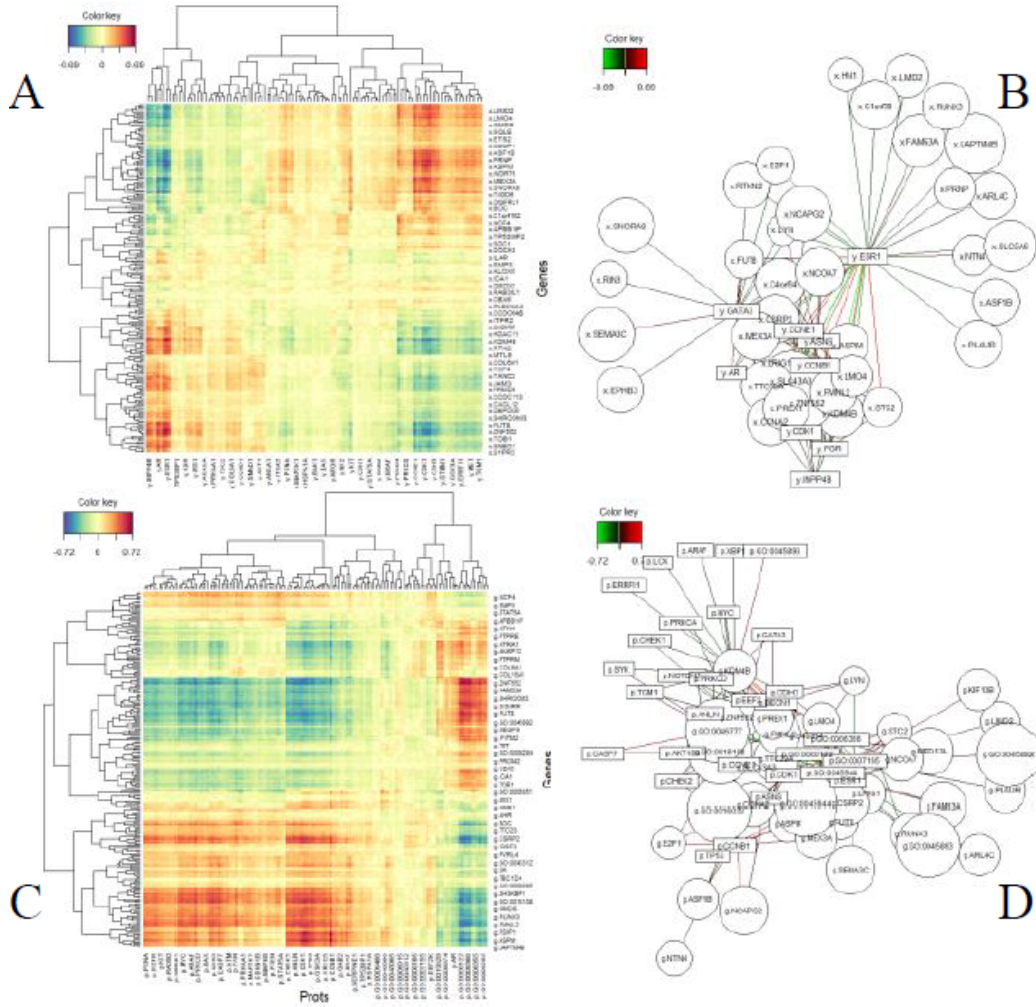


Figure 4.2: BRCA results overview

4.4 Results from the application of MFA on TCGA-BRCA data with, and without, expanded data

Figure 4.3 includes a Correlation Circle (left), with most relevant genes, proteins and added GO annotations. Distribution of samples (right) along the first two plotted dimensions. Both results coming from the application of Multiple Factor Analysis (FactoMineR and factoextra R packages) performed on the same 150 samples (Basal, Her2 and LuminalA conditions) from TCGA-BRCA.

4. Results

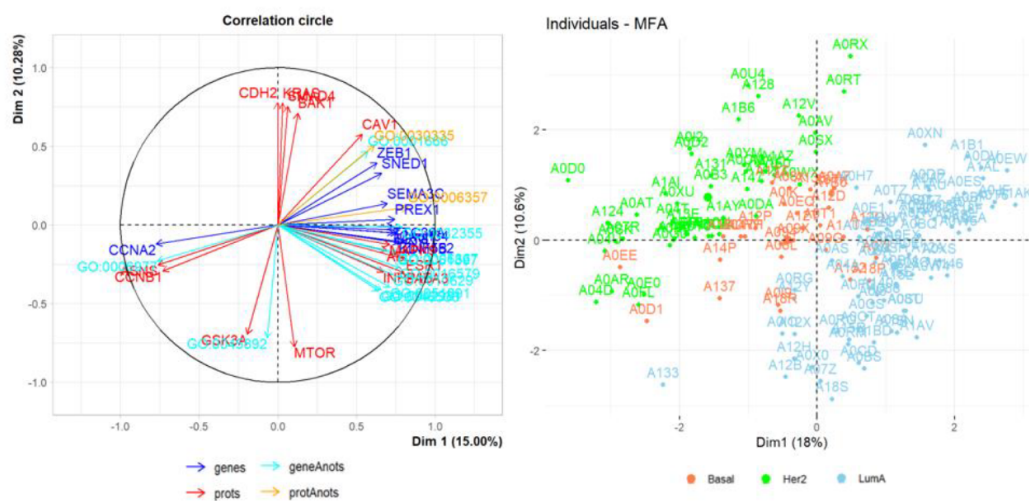


Figure 4.3: BRCA results with MFA

4.5 Resultats de la creacio del paquet amb Targets...

5

Discussion

Potser no cal posar la TOC aquí?

Resum de l'article. Apuntant a les conclusions. Comentant problemes i limitacions (emprar combinacions lineals de variables per crear-ne de noves).

Possibles extensions [punts de millora] Comentar i descriure cadascun d'ells:

- Poder fer servir 3 o més conjunts de dades
- Poder ponderar els pesos de les anotacions, segons tipus, data set d'origen, etc.
- Permetre treballar amb dades faltants o, fins i tot, blocs de dades faltants.
- Millorar les opcions del paquet: mètodes d'anotació bio, mètodes d'integració, tipus de gràfics resultants...

There is grandeur in this view of life, with its several powers, having been originally breathed into a few forms or into one; and that, whilst this planet has gone cycling on according to the fixed law of gravity, from so simple a beginning endless forms most beautiful and most wonderful have been, and are being, evolved.

— Charles Darwin ([Darwin, 1859](#))

6

Conclusions

If we don't want Conclusion to have a chapter number next to it, we can add the `{-}` attribute.

Conclusion 1

The need for a better biological interpretation of multi-omics integrative methods let us to consider the inclusion of biological information during (not after) the analysis process

Conclusion 2

We propose a method focused on the expansion of the starting omics datasets, by adding new annotation-derived features to those matrices, before applying the integrative analysis

Conclusion 3

This approach allows the inclusion of relevant information from the main biological annotation tools, as well as any custom annotation, combined with the use our preferred Dimension Reduction techniques

Conclusion 4

We have implemented a pipeline for reproducible and easy-to-use execution, that facilitates the control of each step, the visualization of results and their reporting to PDF/HTML formats.

Appendices



The First Appendix

This first appendix includes an R chunk that was hidden in the document (using `echo = FALSE`) to help with readability:

In `02-rmd-basics-code.Rmd`

And here's another one from the same chapter, i.e. Chapter ??:

B

The Second Appendix, for Fun

References

Darwin, C. (1859). *On the Origin of Species by Means of Natural Selection or the Preservation of Favoured Races in the Struggle for Life*. John Murray.