

Integrative Omics Data Analysis

Alex Sánchez

Statistics and Bioinformatics Research Group
Statistics department, Universitat de Barcelona.
Statistics and Bioinformatics Unit Vall d'Hebron Institut de Recerca

November 3, 2015

Table of Contents

1 Introduction and Context

- From Biology to Omics
- Omics data integration

2 MV methods for IODA

- Classical MVA
- (Modern) Multivariate methods

Presentation and Objectives

Who, where, what?

The screenshot shows the homepage of the EIB | Presentation website. The header features the EIB logo and the text "Statistics and Bioinformatics University of Barcelona - Departament of Statistics". The main menu includes links for PRESENTATION, MEMBERS, RESEARCH, SERVICES, ACTIVITIES, INTRANET, TEACHING, and CONTACT. Below the menu, there's a section for Article Archives with links to May, March, December, and October 2014. A "Since your last visit..." section highlights changes like new wiki pages and R files. The central content area discusses the group's objectives and collaborations, mentioning the Fundació Vall d'Hebrón Institut de Recerca. On the right, there are sections for Upcoming Events (with a link to add one) and a Calendar for October 2014.

EIB | Presentation

eib.stat.ub.edu/Presentation

Aplicaciones Universitat Oberta d... Cadena SER en direc... Outlook Web App Vall d' Hebron Resea... Evernote Web Otros marcadores

GRUP DE RECERCA EN ESTADÍSTICA Y BIOMATRÍMATICA

Statistics and Bioinformatics

University of Barcelona - Departament of Statistics

LOG OUT

PRESENTATION MEMBERS RESEARCH SERVICES ACTIVITIES INTRANET TEACHING CONTACT

Article Archives

May 2014 [1]
March 2014 [1]
December 2010 [1]
October 2010 [3]

Since your last visit...

2014-07-14

2 wiki pages changed

1. El llenguatge estadístic R (2014)

2. DIMENSIÓ

6 new files

1. EjerGrafics.R (EjerGrafics.R)
2. EjercicioFinal.R (EjercicioFinal.R)
3. EjerFinalSec3.R (EjerFinalSec3.R)
4. Grafics.R (Grafics.R)

EIB > Presentation

Child Add Page

The Statistics and Bioinformatics research group has as its main objectives the development of methods and tools to deal with problems appearing in the interface between Statistics and Bioinformatics. We started focussing in DNA microarrays but we are also interested in statistical methods for 'omics' data integration and next generation sequencing (NGS).

Our group collaborates with different research groups in the fields of biology and biomedicine, to whom it offers statistical support for problems which are specifically statistic in nature, such as experimental design or high throughput data analysis, and also in more general aspects, such as modelling, analysis or data mining.

After a first period of collaboration agreements with the Fundació Vall d'Hebrón Institut de Recerca, we contributed to the creation of the Statistics and Bioinformatics Unit (UEB) which provides statistical and bioinformatical support to VHIR and other researchers.

Upcoming Events

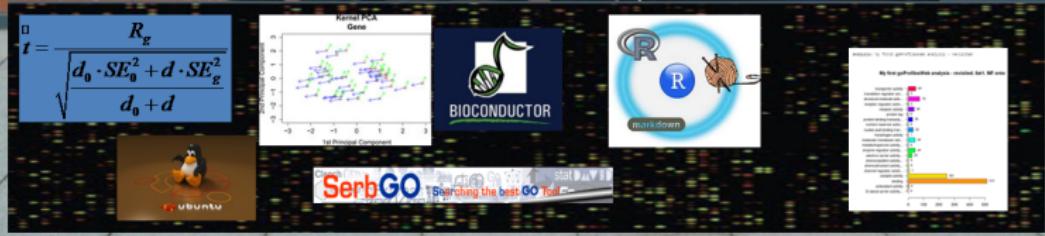
1) 2014-11-10
Workshop sobre analisi integrativa de dades omíquies

Add Event

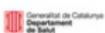
Calendar

October 2014

Mo	Tu	We	Th	Fr	Sa	Su
29	30	01	02	03	04	05
06	07	08	09	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31	01	02



català | castellano | english



20VHIR



► The Institute

► News

► Research

▼ Services

Presentation

Laboratory Coordination

Animal Facility

USIC

UAT

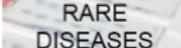
Biobank

UEB-USMIB

► Seminars



RARE DISEASES



presentacio/ca_presentacio.asp?mv1=...

UEB-USMIB

[Presentation](#) [Teamwork](#) [Services](#) [Teaching](#) [Publications](#)

The Statistics and Bioinformatics Unit (UEB) was created in 2008 within the Research Institute of the Hospital Vall d'Hebron (VHIR) in order to promote the use and development of modern statistical and bioinformatics resources on research performed in its environment.

Nowadays, the Statistics and Bioinformatics Unit includes the former Support Unit in Methodology for Biomedical Research (USMIB) and, as part of the Scientific and Technical Support Area of the Vall d'Hebron Research Institute, has the mission to provide expert advice, services and training for clinical and biomedical research.

The main objectives of the UEB are:

To provide statistical, methodological and bioinformatics support for clinical and biomedical research, mainly in our center but also to the rest of the community.

To contribute to training in statistics and bioinformatics for clinical and biomedical research, by conducting its own courses and participating in formal training in the VHIR's area.

To carry out innovation and development activities in the field of statistics and bioinformatics, particularly in anything that could revert in an improvement of the procedures and services provided by the Unit.

More information:

- [UEB's Webpage](#)
- [UEB's presentations on Slideshare](#)
- [Course of advanced statistics on biomedical research](#)

Service request



Rates



on-line Resources

CONTACT

First Floor (Room 131)
Vall d'Hebron Research Institute

Passeig Vall d'Hebron,
nº 119-129
08035 Barcelona

+0034 93 489 4007
+0034 93 489 4102

ueb@vhir.org



Interaction promotes opportunities



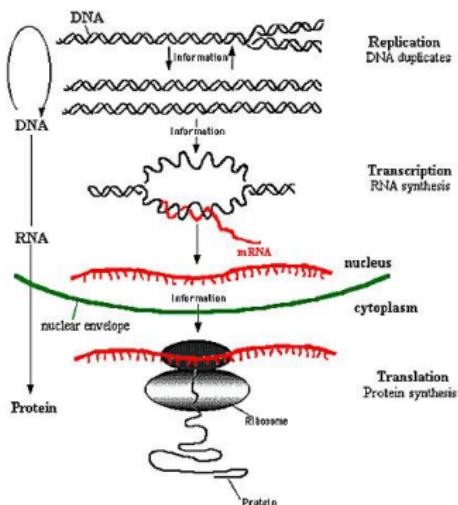
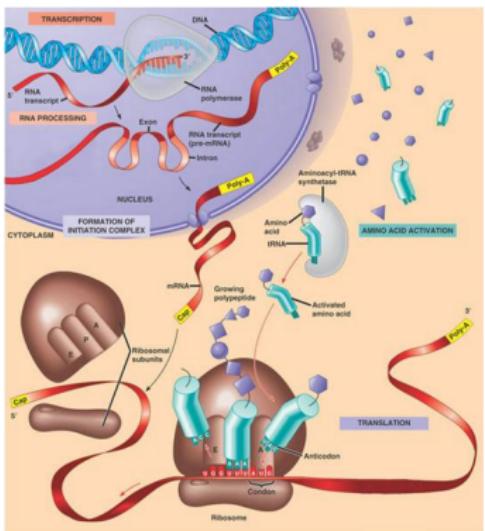
Omics, Omics Data and Omics Integration

What is "omics"?

- In biological context , the suffix “omics” is used to refer to the study of large sets of biological molecules (Smith et al., 2005)
- The study of different components participating and/or regulating complex biological processes, triggered the development of several fields that, together, are described with the term OMICS.

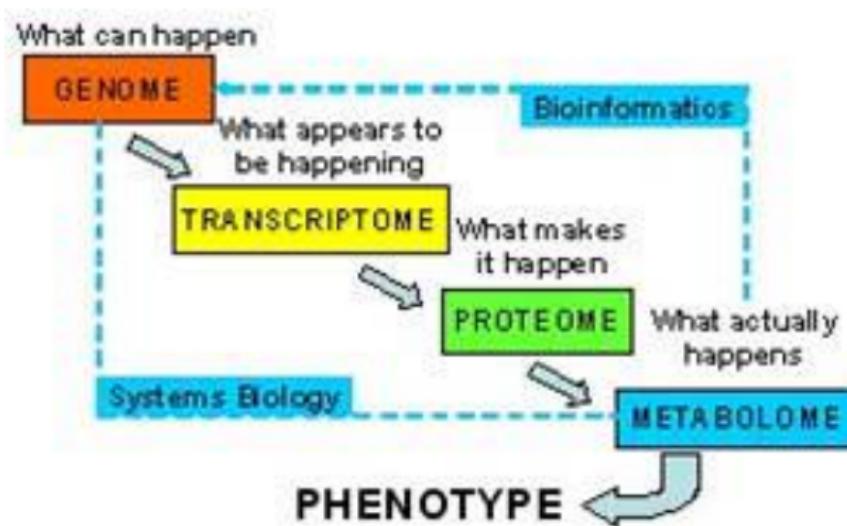


The central dogma

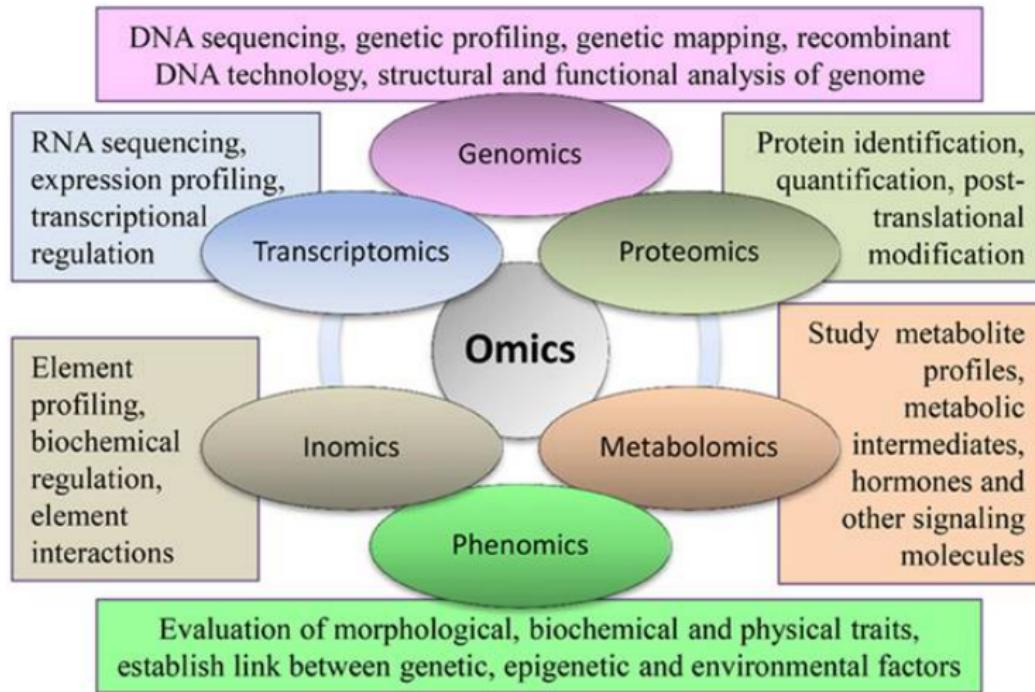


The Central Dogma of Molecular Biology

The Omics Cascade (1): “omes”

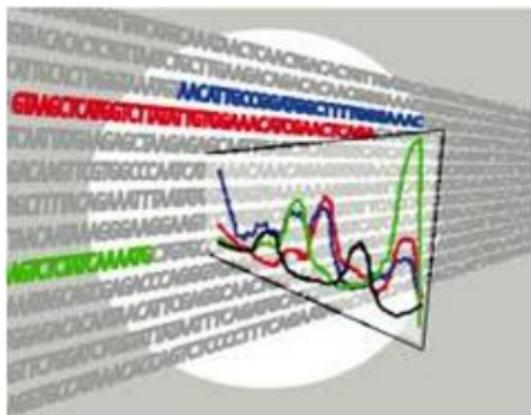


We study omes with omics



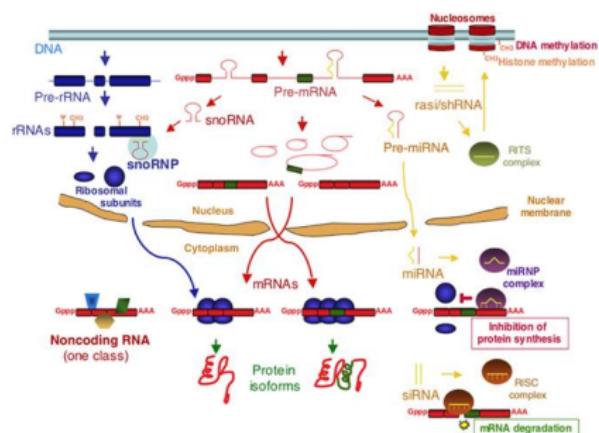
Genomics

Genomics is a discipline in genetics that applies recombinant DNA, DNA sequencing methods, and bioinformatics to sequence, assemble, and analyze the function and structure of genomes (the complete set of DNA within a single cell of an organism)



Transcriptomics

- The transcriptome is the set of all RNA molecules, in one or a population of cells.
- Transcriptomics, examines expression levels of mRNAs in a given cell population, often using high-throughput techniques: microarrays or NGS.



Proteomics

- The large-scale study of proteins (the proteome), particularly their structure and function.
- Relies on a wide spectra of techniques
 - 2D gel based
 - Mass Spectrometry (MS)
 - Seldi-TOF (MS)
 - Protein Arrays
 - ...

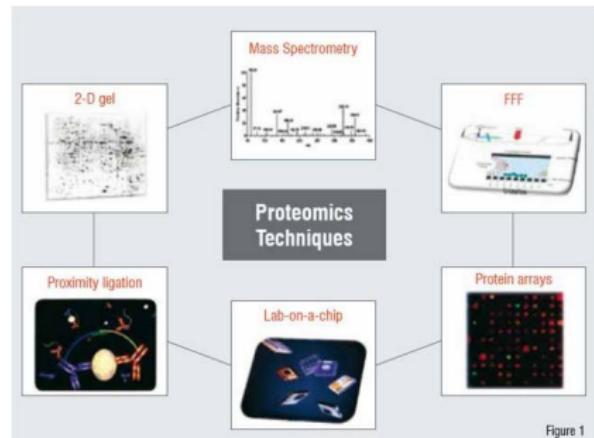
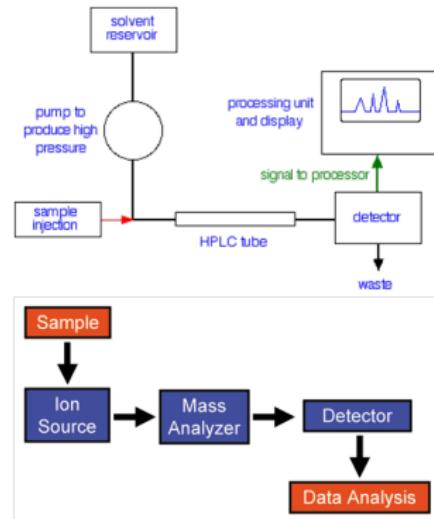


Figure 1

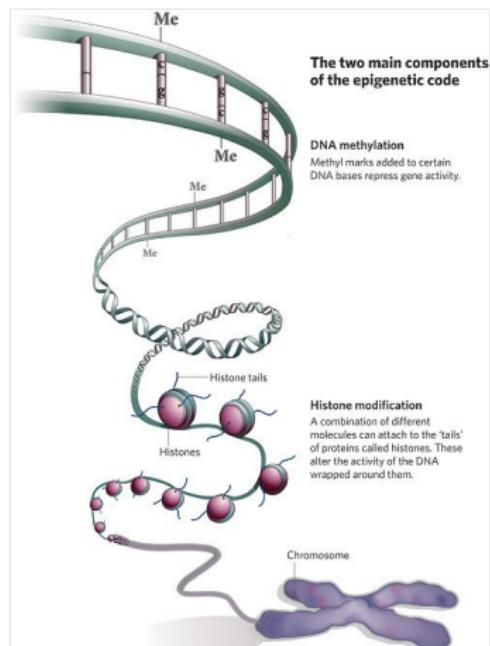
Metabolomics

- Comprehensive and simultaneous systematic determination of
 - metabolite levels in the metabolome and
 - their changes over time as a consequence of stimuli.
- Relies on
 - Separation techniques: GC, CE, HPLC, UPLC
 - Detection techniques: NMR, MS

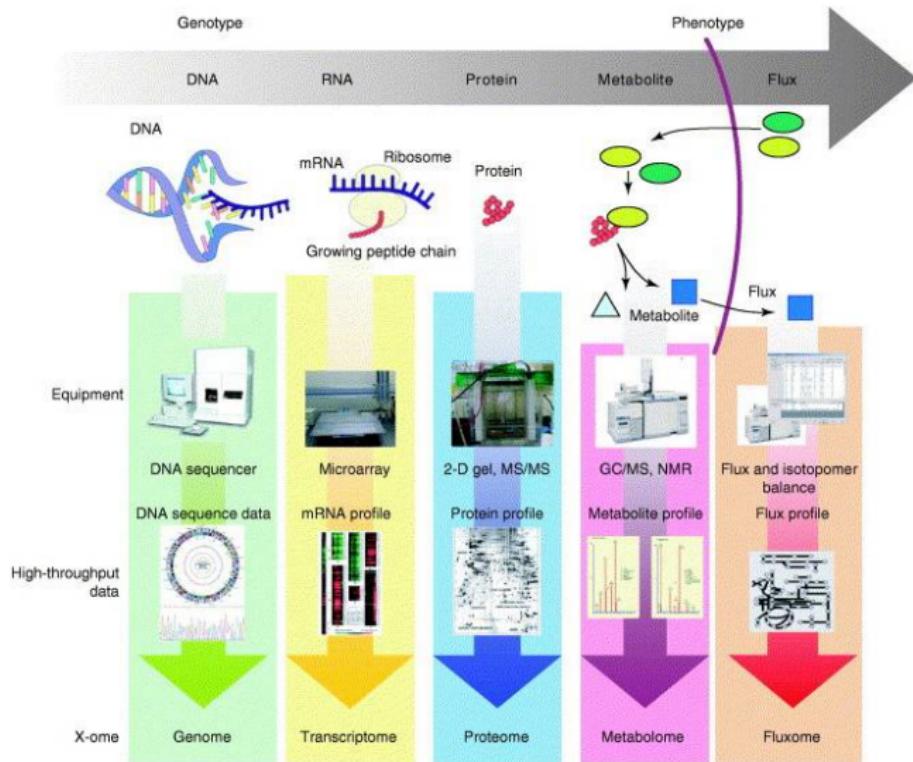


Epigenetics and Epigenomics

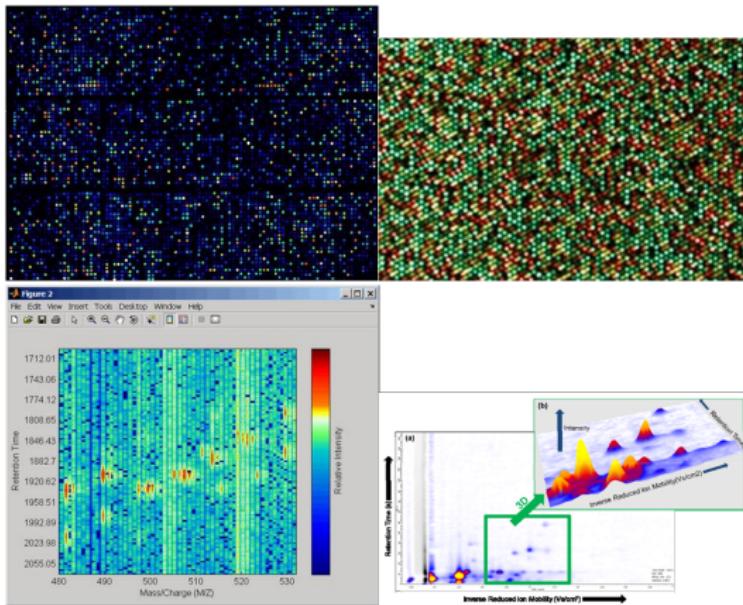
- *Epigenetics* is the study of changes in the phenotype or gene expression caused by other mechanisms than changes in the underlying DNA sequence.
 - DNA methylation
 - Histone modifications
- Epigenetics refers to the study of single genes or sets of genes. Epigenomics refers to global analyses of epigenetic changes across the entire genome



In summary omic to stuoy omes



Omics data are high throughput



Bioinformatics and Biostatistics is essential

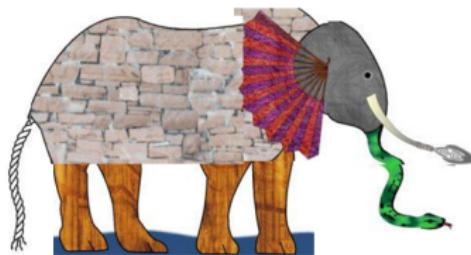


“Data don't make any sense,
we will have to resort to statistics.”

Integrative Analysis and Data integration: methods, types, tools, challenges

Why should we integrate data?

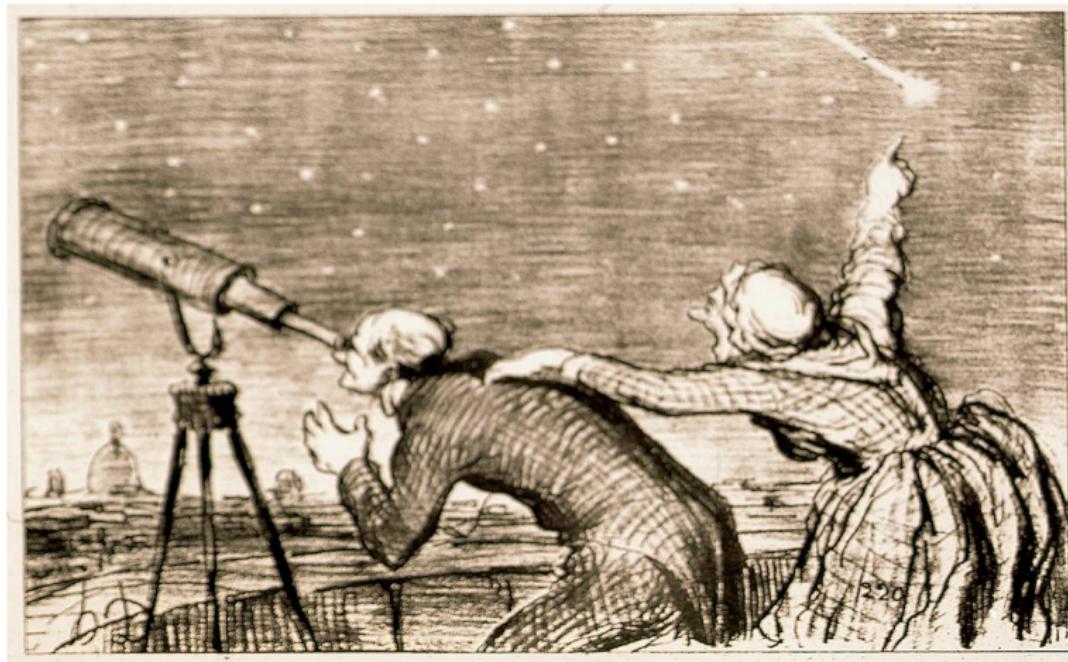
The Blind Men and the Elephant



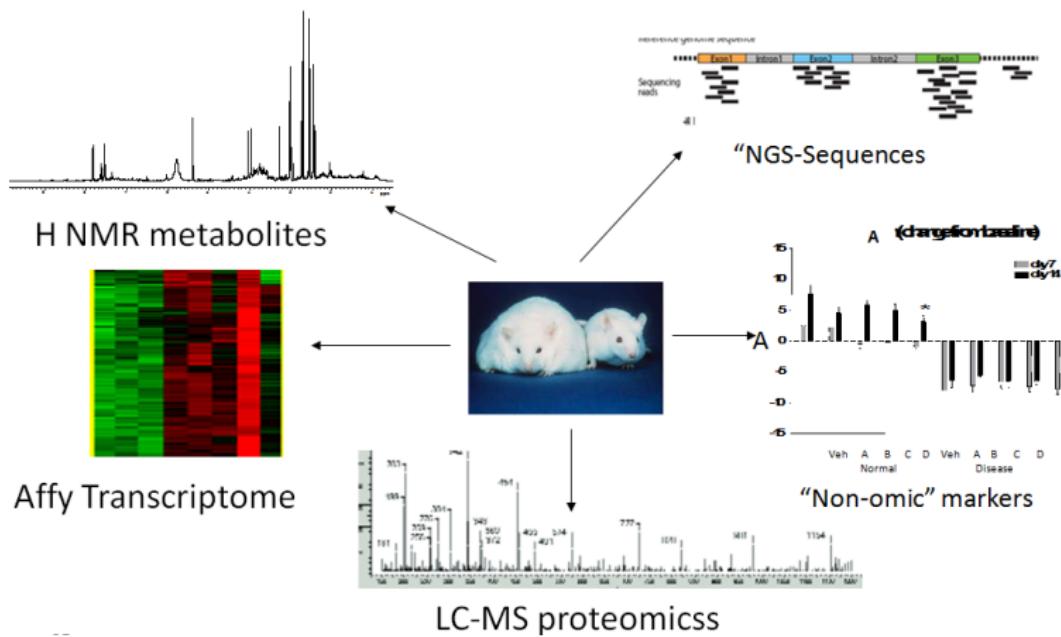
What we learn from an experiment may depend on where we look, how we look, and the scope of our view!

http://www.noogenesis.com/pineapple/blind_men_elephant.html

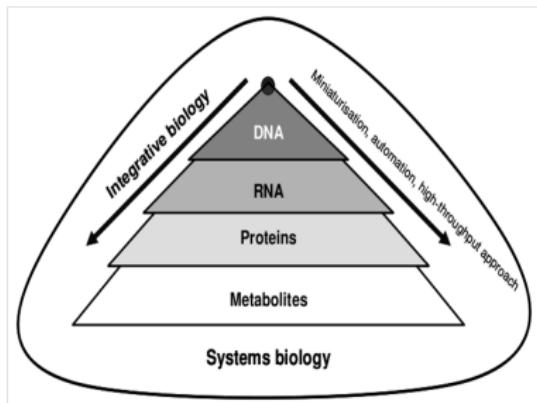
Focussing only on one platform risks missing an obvious signal



So let's measure as many as possible

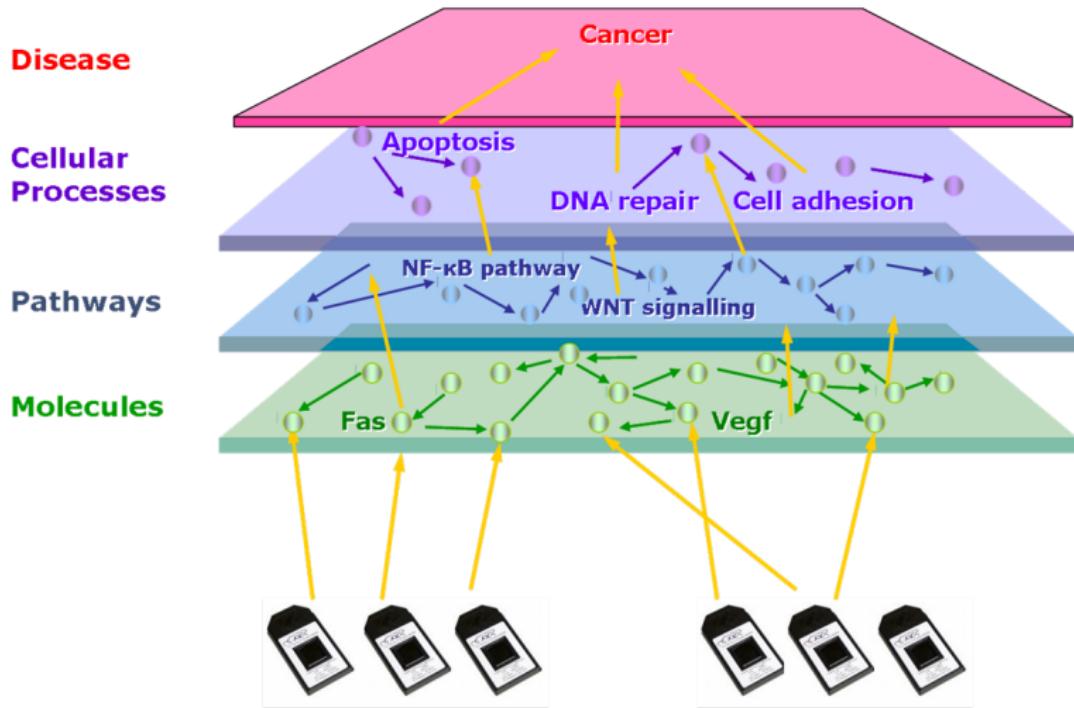


From componentwise to global approaches

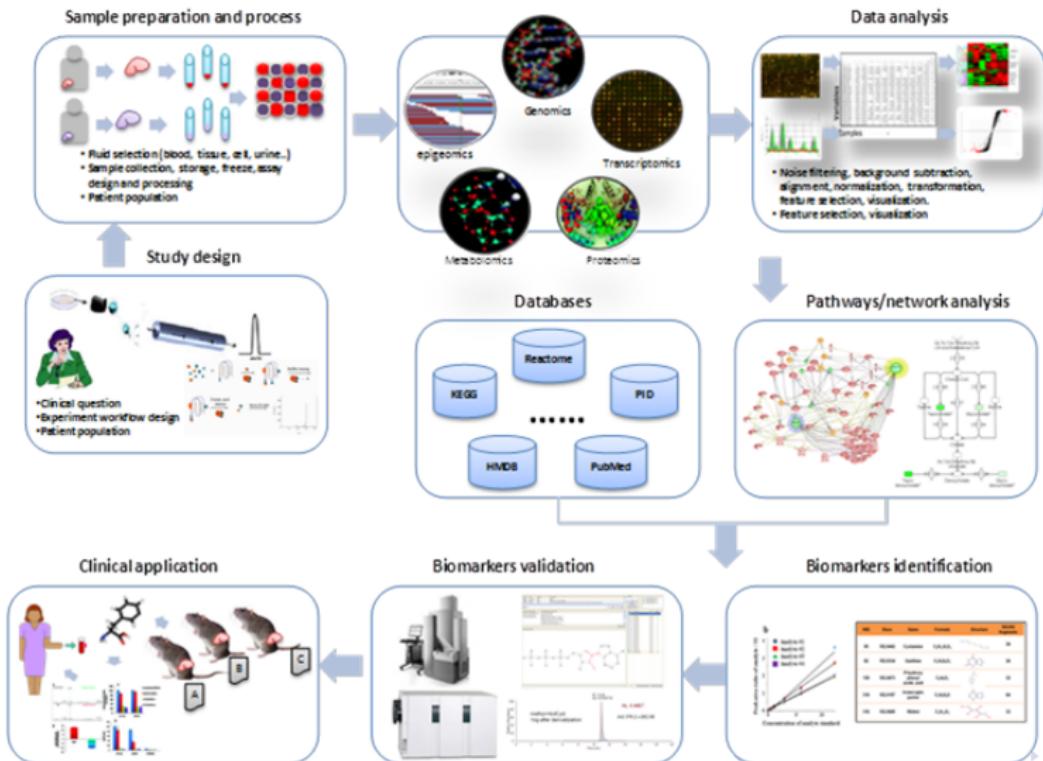


- It is expected that the integrated collection and analysis of diverse types of data,
- jointly modelled and analyzed in a systems biology approach
- can shed light on the global functioning of biological systems.

Ultimate Goal (1): understanding of complex processes

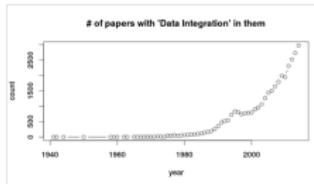


Ultimate Goal (2): Improved Robust Biomarkers



Data Integration is cool

- Everywhere nowadays in Biology, Medicine, Bioinformatics, [BTW much less in Statistics]
 - Meetings: Barcelona (Feb. 2013), Leiden (Apr. 2013), Ascona (May 2013), Crete (Nov 2014) and many more
 - Financing (FP7): projects with $> 10^6$ each (Stategra, MimOmics).
 - Increasing use of 'data integration' in publication titles.



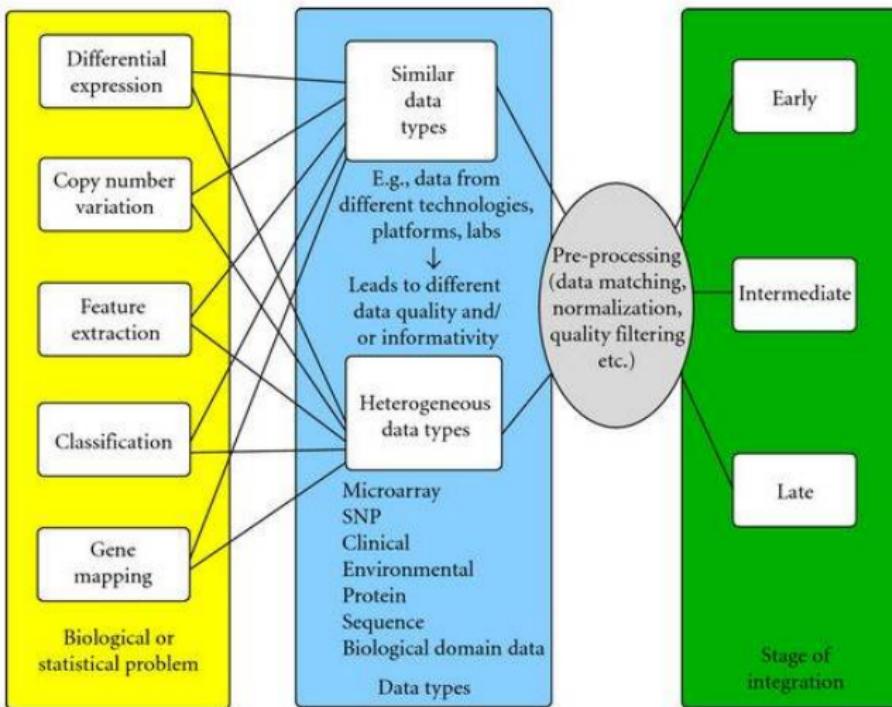
But what is Data Integration?

- “*Data integration*” may mean different things...
 - Computational combination of data
 - Combination of studies performed independently
 - Simultaneous analysis of multiple variables on multiple datasets.
 - Not to mention any possible approach for homogeneously querying heterogeneous data sources
- **Integrative analysis** may be preferable

Integrative Omics Data Analysis

- The idea that efficient integration of data from different OMICS can greatly facilitate the discovery of true causes and states of disease is rapidly pervading the biomedical community
- The aims of integrative analysis is the deciphering of complex biological relationships empowered by the combined use of distinct pieces of information that represent a, probably partial, view of the different levels at which these processes happen

There are many types of integrative analysis



There are many methods ...

BMC Bioinformatics



Research article

Open Access

A structured overview of simultaneous component based data integration

Katrijn Van Deun^{*1}, Age K Smilde², Mariët J van der Werf³, Henk AL Kiers⁴
and Iven Van Mechelen¹

Vol. 26 ISMB 20

doi:10.1093/bio

Multivariate multi-way analysis of multi-source data

Ikkka Huopaniemi^{1,*}, Tommi Suvitaival¹, Janne Nikkilä^{1,2}, Matej Oresič³ and
Samuel Kaski^{1,*}

Advances in Integrative Causal Analysis

Presenter: Ioannis Tsamardinos, Ph.D., Associate Professor and
Sofia Triantafillou, Ph.D. candidate.

Low level data fusion methods searching for common
and distinctive biological information in hetero omics
data-sets

Frans M. van der Kloet¹, Patricia Sebastián-León², Ana Conesa²,
Age K. Smilde¹ and Johan A. Westerhuis¹

bioRxiv preprint doi: https://doi.org/10.1101/2023.07.10.544222; this version posted July 10, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a CC-BY-NC-ND 4.0 International license.



A myriad of tools (1): R-based



Exploration and
Integration of
Omics datasets

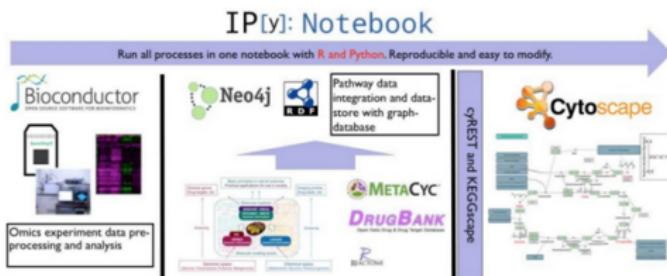


& Made4
& Omicade4



pwOmics

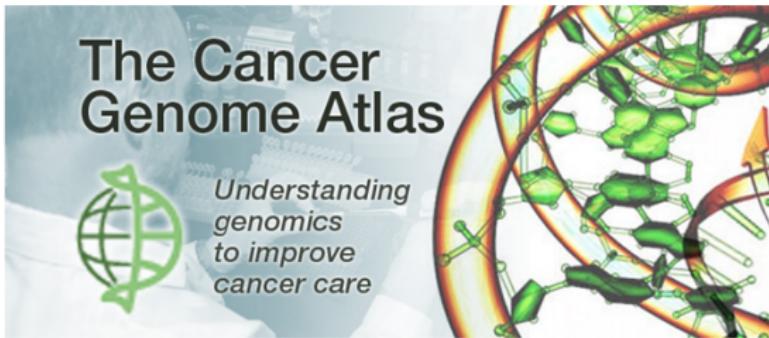
A myriad of tools (2): Others



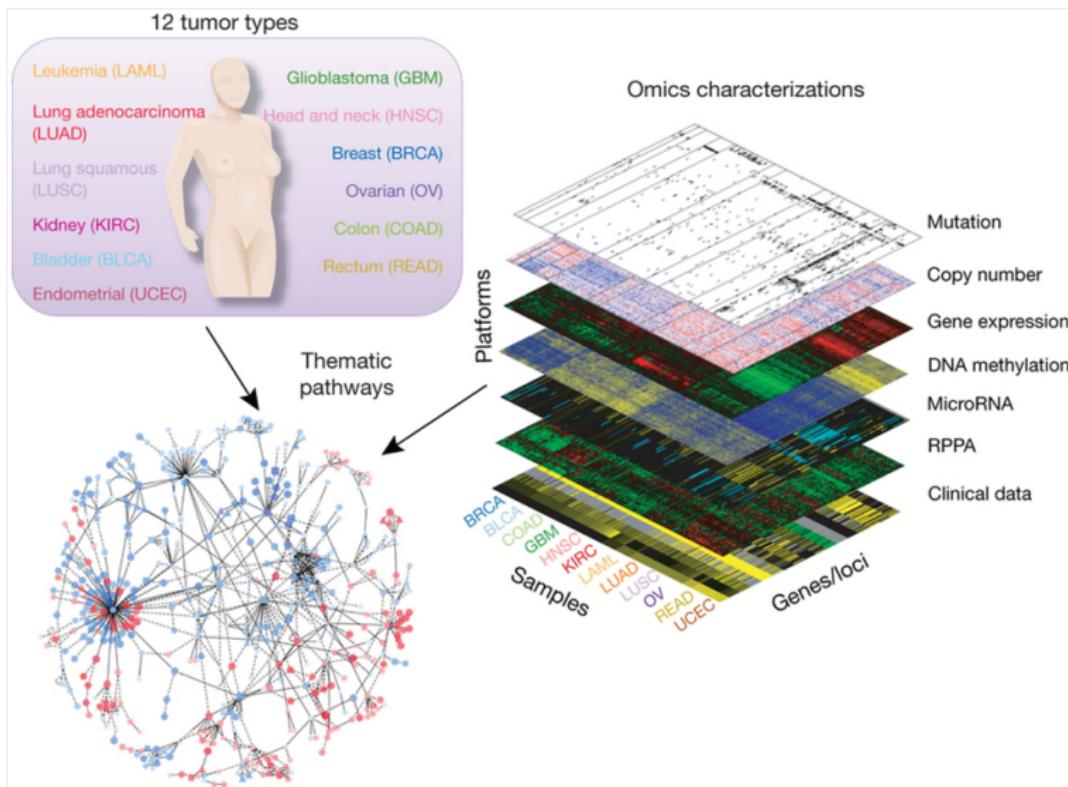
02/02/2015

36

Multidisciplinary projects



Data sources: TCGA



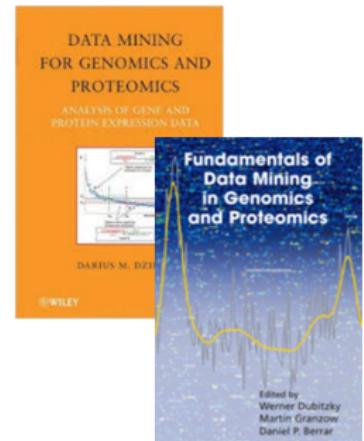
So what?

- We will restrict to arbitrarily chosen situation
 - Multivariate statistical methods, classic and extensions
- for which we will sketch basic ideas,
 - General concept
 - Basic formulation
- and provide some examples of use
 - Applications and tools

Methods for Integrative Omics Data Analysis

Multivariate statistics in genomics

- Multivariate methods have pervaded the field of genomics since its very beginning
 - The (in)famous clustering (HC, heatmaps)
 - Matrix factorizations / Dimension reduction (PCA, SVD, CoA)
 - Discriminant Analyses (LDA – > DLDA, ...)



To cite but a few.

Best approach for omics data analysis?

- Classical Statistics
 - Multiple regression
 - Discriminant analysis
 - ANOVA
- Data tables are long and lean



- Assumptions
 - Independent variables
 - More observations than variables
 - Multivariate normality
 - Interested in one dependent
 - Few missings
- DO NOT hold for many omics data

The nature of omics data

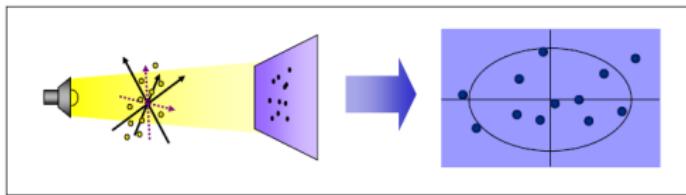
- Omics data are diverse
 - They measure distinct characteristics
 - GC/MS spectrum, Expression, Concentration
- Although they have aspects in common
 - Most of them are high throughput
 - Many variables (K) measured simultaneously
 - Relatively expensive, ethical limitations, regulations
 - Few samples (N) analyzed



Figure: $K \gg N$

A Better Way

- Multivariate analysis by projection methods
 - Looks at ALL the variables together
 - Avoids loss of information
 - Finds underlying trends = “latent variables”
 - More stable models



A not-so random walk by classical MVA

- “Classical multivariate statistical methods” (whatever this means) have been thoroughly used in the analysis of omics data.
- We review some methods and applications
 - Principal components analysis (PCA)
 - Singular value decomposition (SVD)
 - Correspondence analysis (CoA)
 - Factor analysis (FA)

Principal Components Analysis

- PCA is a statistical method that can be traced back to Pearson (1901) or Hotelling (1933).
- Its main goals are
 - Capture the information provided by the original variables using a smaller number of new variables
 - Provide a representation of the data in reduced dimension

What does PCA do?

- Given a $k \times n$ data matrix containing k (probably correlated) measurements on n samples (objects/individuals...), PCA decomposes this matrix in new k components that they ...
 - account for different sources of variability in the data,
 - are uncorrelated, that is each component accounts for a different source of variability,
 - have decreasing explanatory ability: each component explains more than the following
 - allow for a lower dimensional representation of the data in terms of scores on principal components.
 - get an overview of the dominant patterns and major trends in the data (visualize clusters, identify outliers)

Mathematical formulation

We search for linear combinations of the original variables

$$Z_1 = a_{11}X_1 + a_{12}X_2 + \dots + a_{1p}X_p$$

$$Z_2 = a_{21}X_1 + a_{22}X_2 + \dots + a_{2p}X_p$$

...

$$Z_p = a_{p1}X_1 + a_{p2}X_2 + \dots + a_{pp}X_p$$

Subject to:

- Z_1, Z_2, \dots, Z_p uncorrelated
- $Var(Z_1)$ maximal
- $Var(Z_1) \geq Var(Z_2) \geq \dots \geq Var(Z_p)$.

Mathematical formulation

- All coefficients and eigenvalues can be obtained by the spectral decomposition of the covariance matrix

$$\Sigma = AD_\lambda A'$$

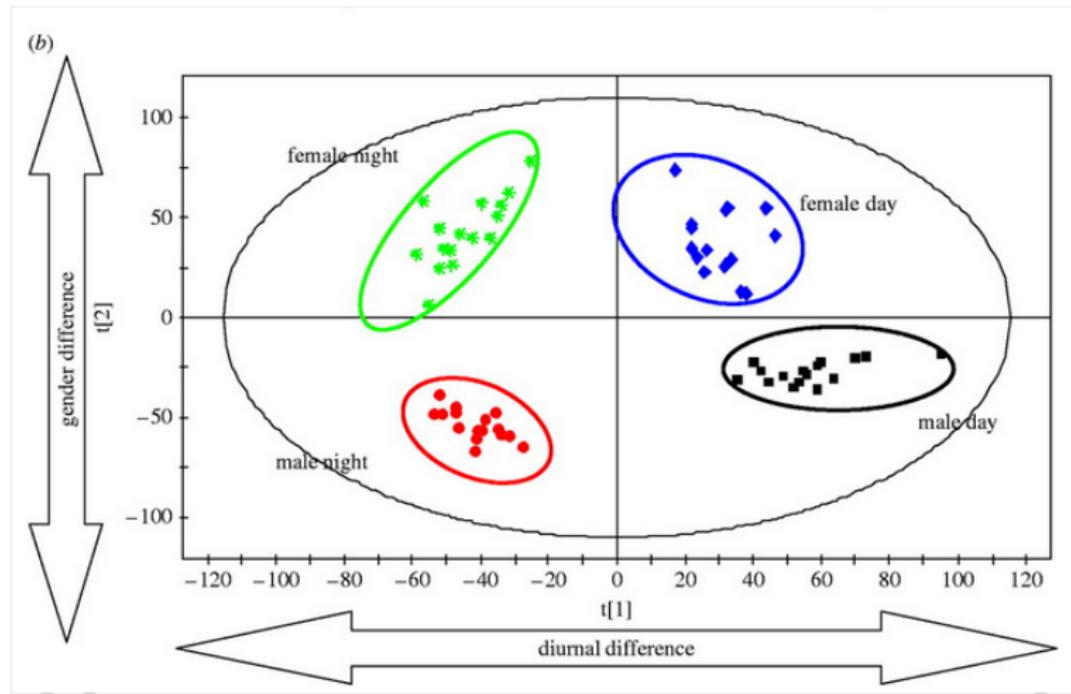
- We will use the sample covariance matrix S to estimate Σ

$$S = AD_\lambda A'$$

- and estimate the components by sample components

$$\begin{array}{ccc} Z & = & X \\ (n \times p) & & (n \times p) \\ & & A \\ & & (p \times p) \end{array}$$

Data visualization in reduced dimension



Applications of PCA

- PCA has been used in a wealth of applications
- The best known
 - Batch effect detection
 - Pattern recognition
- Besides standard applications there exist many extensions
 - Probabilistic PCA, Bayesian PCA, Inverse non-linear PCA, Nipals PCA, Robust PCA

Examples

The screenshot shows a detailed view of a research article from the journal **Science**. The article is titled "The Transcriptional Program of Sporulation in Budding Yeast" by S. Chu, J. DeRisi, M. Eisen, J. Mulholland, D. Botstein, P. O. Brown, and L. Herskowitz. The page includes a sidebar for "Author Manus" and various navigation links like "Article Views", "Abstract", "Full Text", and "A correction has been published". The main content area displays the article's text, figures, and tables.

Science.org | FEEDBACK | HELP | LIBRARIANS | All Science Journals | Enter Search Term | CRAI, UNIVERSITAT DE BARCELONA | ALERTS | ACCESS RIGHTS

AAAS NEWS SCIENCE JOURNALS CAREERS MULTIMEDIA COLLECTIONS

Science Home Current Issue Previous Issues Science Express Science Products My Science About the Journal

Home > Science Magazine > 23 October 1998 > Chu et al., 282 (5389): 699-705

Article Views Abstract Full Text (PDF) A correction has been published

Science 23 October 1998;
Vol. 282 no. 5389 pp. 699-705
DOI: 10.1126/science.282.5389.699

RESEARCH ARTICLE

The Transcriptional Program of Sporulation in Budding Yeast

S. Chu¹, J. DeRisi¹, M. Eisen¹, J. Mulholland¹, D. Botstein¹, P. O. Brown¹, L. Herskowitz¹

Pac Symp Biocomput. Author manuscript; available in PMC 2009 Apr 17.
Published in final edited form as:
Pac Symp Biocomput. 2000 : 455–466.

PMCID: PMC2669932
NIHMSID: NIHMS97353

PRINCIPAL COMPONENTS ANALYSIS TO SUMMARIZE MICROARRAY EXPERIMENTS: APPLICATION TO SPORULATION TIME SERIES

Soumya Raychaudhuri,^{*} Joshua M. Stuart,^{*} and Russ B. Altman^Ψ

Author information ► Copyright and License information ►

Examples



Journal of Proteomics

Volume 75, Issue 13, 16 July 2012, Pages 3938–3951



Batch effects correction improves the sensitivity of significance tests in spectral counting-based comparative discovery proteomics

Josep Gregori^{a,b}, Laura Villarreal^a, Olga Méndez^a, Alex Sánchez^{b,c}, José Baselga^a, Josep Villanueva^a.

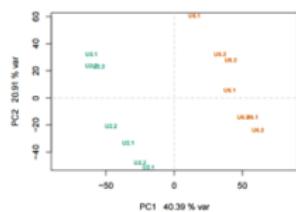


[+ Show more](#)

R Packages alpha Home Blog All packages

Search R packages

Search



Exploratory Data Analysis of LC-MS/MS data by spectral counts

Exploratory data analysis to assess the quality of a set of LC-MS/MS experiments, and visualize de influence of the involved factors.

Author: Josep Gregori, Alex Sanchez, and Josep Villanueva
Date of publication: None
Maintainer: Josep Gregori <josep.gregori@gmail.com>
License: GPL-2
Version: 1.2.0

View on Bioconductor



Singular value Decomposition

- When Things Get Messy...
 - PCA is fine when initial dimension is not too big
 - Space and time complexity are of $O(d^2)$ – size of covariance matrix
 - Otherwise – we have a problem...
 - E.g. when $d = 10^4$ = time/space complexity is $O(10^8)$...
 - An alternative approach: SVD

Mathematical formulation

Any real $n \times p$ matrix X can be decomposed as

$$X = UDV'$$

- U $n \times r$ matrix of orthonormal left singular vectors. $U'U = I_r$
- D $r \times r$ diagonal matrix of non-increasing positive singular values ($d_{11} \geq d_{22} \geq \dots \geq d_{rr}$)
- V $p \times r$ matrix of orthonormal right singular vectors $V'V = I_r$

A rank k approximation \hat{X} to matrix X optimal in the least squares sense, is obtained as

$$\hat{X} = U_{[,1:k]} D_{[1:k,1:k]} V'_{[,1:k]}$$

E.g., a rank 2 approximation to matrix X is obtained by

$$U_{n \times 2} D_{(2 \times 2)} V'_{p \times 2}$$

PCA and SVD

- Informally speaking

- PCA is a statistical method that can be solved using different approaches, one of which is SVD.
- SVD is an algebraic method to decompose a rectangular matrix. Used in signal processing to yield approximation.

Case studies

< Current Issue > vol. 97 no. 18 > Orly Alter, 10101–10106, doi: 10.1073/pnas.97.18.10101

 CrossMark click for updates

Singular value decomposition for genome-wide expression data processing and modeling

Orly Alter^{*†}, Patrick O. Brown^{#‡}, and David Botstein^{*}

Author Affiliations 

Contributed by David Botstein

SPIE Proceedings | Volume 4266 | Analysis of Multiple Expression Profiles >

< Previous Article | Next Article >

Proceedings Article

Processing and modeling genome-wide expression data using singular value decomposition

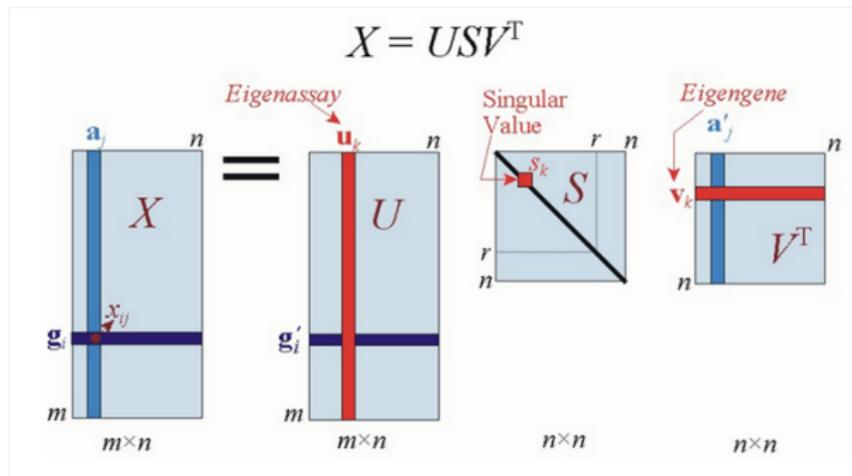
Orly Alter ; Patrick O. Brown ; David Botstein
[+] Author Affiliations

Proc. SPIE 4266, Microarrays: Optical Technologies and Informatics, 171 (June 4, 2001);
doi:10.1117/12.427986

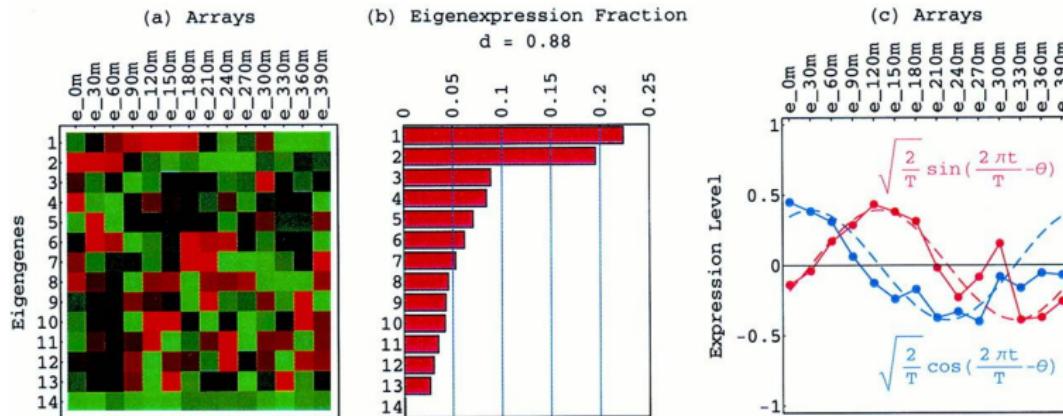
<http://www.alterlab.org/publications/>

Eigengenes, Eigenarrays and SVD

- Alter et al. use the singular value decomposition (SVD) for transforming the dataset from the gene/array space to the eigengene/eigenarray space
- Each dimension is represented by an eigengene/eigenarray/eigenvalue triplet



Example results



Correspondence Analysis

- Yet another dimension reduction technique?
- CA is a method for studying associations between variables
 - Like PCA displays low-dimensional projection
 - It is done for two variables simultaneously so it may help detect associations between them.
- CA is applied to reveal associations between genes and experiments

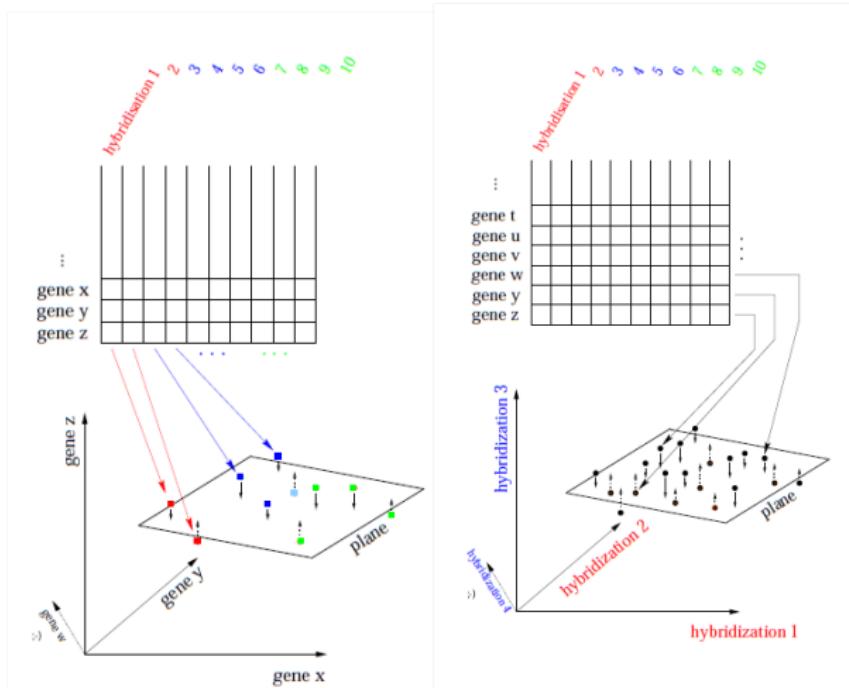
Brief methodology

n gene profiles:
vectors in m-dimensional
experiment space

m hybridisation profiles:
vectors in n-dimensional
gene space

- o projection into a common subspace of
low dimensionality for visualisation
- o conserving point to point distances
(total variance) as well as possible
(-> explained variance)

Visualization



Genes and arrays are projected on the same plane

Mathematical formulation (1)

- Original data N ($I \times J$ matrix) with elements n_{ij} .

$$n_{+j} = \sum_{i=1}^I n_{ij}, \quad n_{i+} = \sum_{j=1}^J n_{ij}, \quad n_{++} = \sum_{i=1}^I \sum_{j=1}^J n_{ij}$$

- row weights $r_i = n_{i+}/n_{++}$
- column weights: $c_j = n_{+j}/n_{++}$
- correspondence matrix P : $p_{ij} = n_{ij}/n_{++}$
- χ^2 matrix S : $s_{ij} = \frac{(p_{ij} - r_i c_j)}{\sqrt{r_i c_j}}$
- singular value decomposition: $S = U \wedge V^T$, diag matrix \wedge contains singular values λ_k of S

- principal coordinates:

- for gene i (principle axis $k = 1, \dots, J$) $f_{ik} = \frac{\lambda_k u_{jk}}{\sqrt{r_i}}$
- for hybridazation j (in the same space) $g_{jk} = \frac{\lambda_k u_{jk}}{\sqrt{r_j}}$

- standard coordinates:

- for gene i $f_{jk} = u_{ik}/\sqrt{r_i}$
- for hybridization j $g_{jk} = v_{jk}/\sqrt{c_j}$

- HMS

- Let N contain the hybridization medians of the original data matrix $N*$ if elements $n*_{ij'}$
- N is submitted to CA
- Let P have elements $p*_{ij'} = n_{ij'}/n_{++}^*$ then the principle coordinates for the supplementary hybridization from correspondence matrix P^* are

$$g_{j'k}^* = \frac{1}{\sum_i p_{ij'}^*} \sum_i \frac{p_{ij'}^* f_{ik}}{\lambda_k}$$

Case studies

[Home](#) > Current Issue > vol. 98 no. 19 > Kurt Fellenberg, 10781–10786, doi: 10.1073/pnas.181597298



Correspondence analysis applied to microarray data

Kurt Fellenberg^{†‡}, Nicole C. Hauser[‡], Benedikt Brors^{†‡}, Albert Neutzner[§], Jörg D. Hoheisel[‡], and Martin Vingron^{†¶}

Author Affiliations ▾

Edited by Peter J. Bickel, University of California, Berkeley, CA, and approved June 29, 2001 (received for review December 18, 2000)

[Abstract](#) [Full Text](#) [Authors & Info](#) [Figures](#) [SI](#) [Metrics](#) [Related Content](#) [PDF](#)

This Issue



◀ PREV ARTICLE



Vol. 21 no. 10 2005, pages 2424–2429
doi:10.1093/bioinformatics/bti367

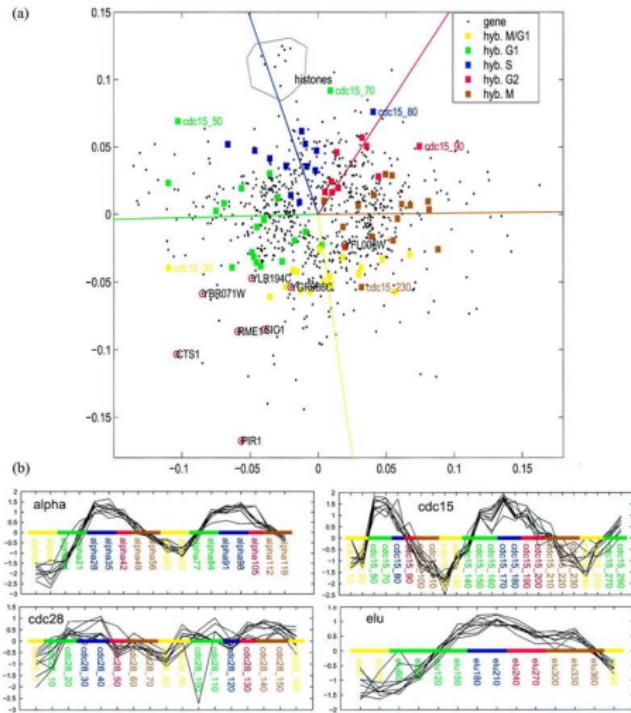
BIOINFORMATICS **ORIGINAL PAPER**

Gene expression

Integration of GO annotations in Correspondence Analysis: facilitating the interpretation of microarray data

Christian H. Busold^{1,*}, Stefan Winter², Nicole Hauser³, Andrea Bauer¹,
Jürgen Dippon², Jörg D. Hoheisel¹ and Kurt Fellenberg¹

CA Analysis of Cell-cycle synchronization data



Integration of GO annotations with CA

- CA can be enhanced by adding supplementary variables.
- These appear at indicative positions without contributing to scores computation.
- As an application of this feature, Gene Ontology or other annotations can be added so that *genes, experimental conditions and gene-annotations are viewed in a single plot.*

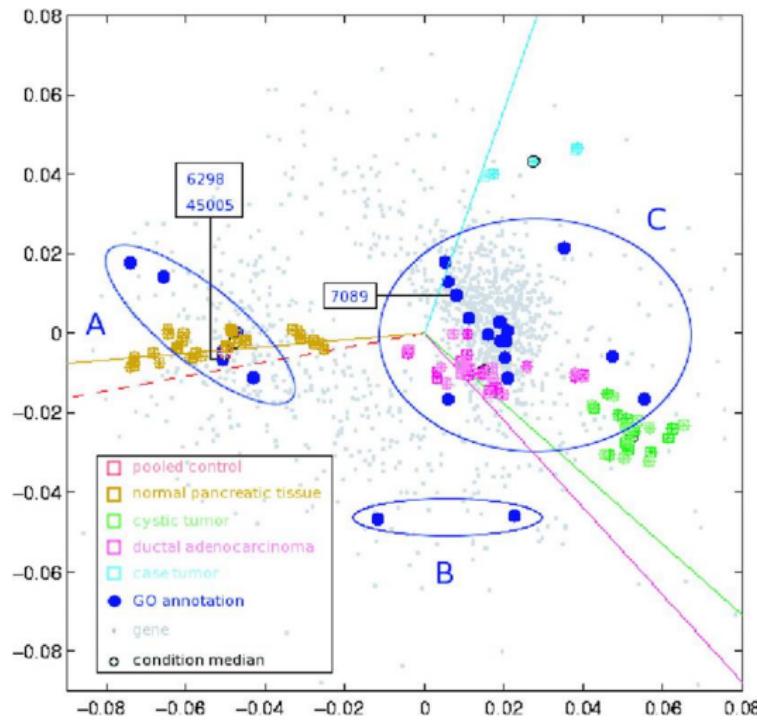
Adding annotations to an expression matrix

- GO Terms can be associated with genes or samples
 - Using a *Boolean implementation* as supplementary columns,
 - Using an *Intensity implementation* as supplementary rows,
- A filtering to select the most descriptive annotations can be applied to keep only annotations highly correlated with common transcription profiles.

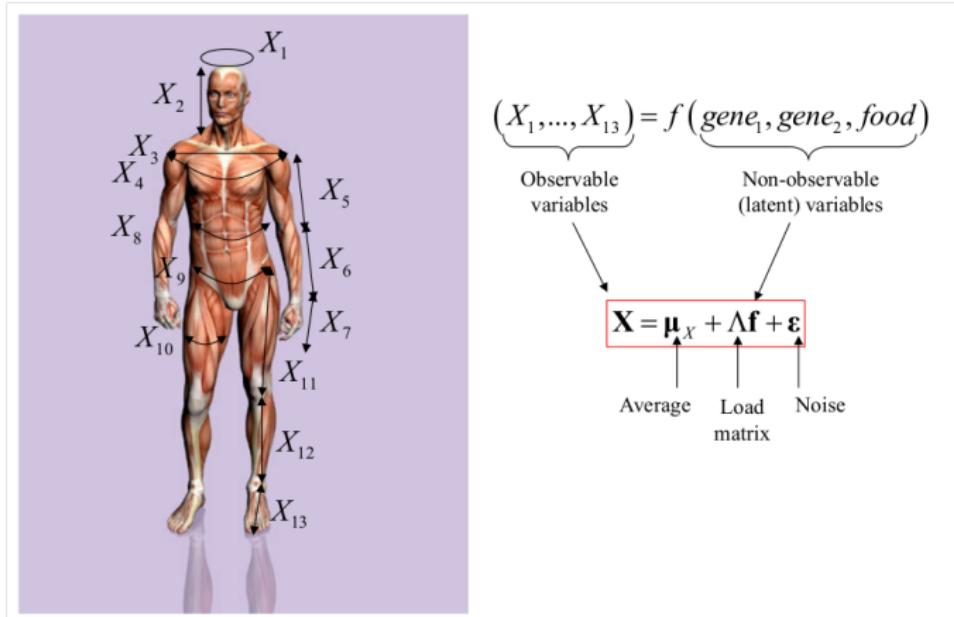
Adding annotations to an expression matrix

Gene	Exp. cond. 1	Exp. cond. 2	Exp. cond. 3	Exp. cond. 4	Term 1	Term 2	.
A	13	300	23	432	1	0	.
B	457	398	355	932	0	1	.
C	24	458	44	364	1	1	.
D	324	245	98	34	0	0	.
E	478	928	293	99	0	1	.
F	38	485	21	375	1	0	.
.
Term 1	75	1243	88	1171			
Term 2	959	1784	692	1395			
.			

CA Map of Human Pancreatic Cancer Study



Factor Analysis



Mathematical formulation (1)

$$\mathbf{X} = \boldsymbol{\mu}_X + \Lambda \mathbf{f} + \boldsymbol{\varepsilon}$$

$\uparrow \quad \uparrow \quad \uparrow$
 $N_p(\boldsymbol{\mu}_X, \Sigma_X) \quad N_m(\mathbf{0}, I) \quad N_p(\mathbf{0}, \Sigma_\varepsilon)$

$$\mathbf{p} \begin{pmatrix} X_1 \\ X_2 \\ \dots \\ X_{13} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \\ \dots \\ \boldsymbol{\mu}_{13} \end{pmatrix} + \begin{pmatrix} \lambda_{11} & \lambda_{12} & \lambda_{13} \\ \lambda_{21} & \lambda_{22} & \lambda_{23} \\ \dots & \dots & \dots \\ \lambda_{13,1} & \lambda_{13,2} & \lambda_{13,3} \end{pmatrix} \mathbf{p} \begin{pmatrix} F_1 \\ F_2 \\ F_3 \end{pmatrix} + \mathbf{m} + \begin{pmatrix} \boldsymbol{\varepsilon}_1 \\ \boldsymbol{\varepsilon}_2 \\ \dots \\ \boldsymbol{\varepsilon}_{13} \end{pmatrix}$$

$\uparrow \quad \quad \quad \quad \quad \uparrow \quad \quad \quad \quad \quad \uparrow$
 $m \quad \quad \quad \quad \quad p$

Properties:

- Factors are uncorrelated / independent
- Factors and noise are uncorrelated
- The load matrix is the covariance between the observed variables and the factors
- The variance of the observed variables and the factors
- The variance of the observed variables can be explained by the loading matrix and the variance of the noise

$$E\mathbf{F}\mathbf{F}' = I$$

$$E\boldsymbol{\varepsilon}\mathbf{F}' = 0$$

$$\boldsymbol{\Lambda} = E(\mathbf{X} - \boldsymbol{\mu}_X)\mathbf{F}'$$

$$\sum_x \boldsymbol{\Lambda} \boldsymbol{\Lambda}' \sum_{\boldsymbol{\varepsilon}}$$

$$\sigma_{x_i}^2 = \sum_{j=1}^m \lambda_{ij}^2 + \sigma_{\epsilon_i}^2 = h_i^2 + \sigma_{\epsilon_i}^2$$

Mathematical formulation (2)

$$X = \mu_x + \Lambda f + \epsilon$$

Properties:

- The load matrix and factors are not uniquely specified: any rotation of the factors can be compensated by the load matrix

$$X = \mu_x + \Lambda f + \epsilon = \mu_x + (\Lambda H')(Hf) + \epsilon$$

Solution:

- 1 impose that $\Lambda \Lambda'$ is a diagonal matrix : Principal factor method
- 2 Impose that $\Lambda' \Sigma_{\epsilon}^{-1} \Lambda$ is a diagonal matrix: Maximum-Likelihood method

H' – > Any orthogonal matrix: This matrix provides the possibility of rotating the components so that we achieve a certain property (having the maximum number of zeros,...) that helps us to understand the factors. There are several criteria offered by the programs: varimax, equamax, parsimax, quartimax, orthomax,...

Examples

Cancer Cell

Volume 17, Issue 1, 19 January 2010, Pages 98–110



Article

Integrated Genomic Analysis Identifies Clinically Relevant Subtypes of Glioblastoma Characterized by Abnormalities in *PDGFRA*, *IDH1*, *EGFR*, and *NF1*

Roel G.W. Verhaak^{1, 2, 17}, Katherine A. Hoadley^{3, 4, 17}, Elizabeth Purdom⁷, Victoria Wang⁸, Yuan Qi^{4, 5}, Matthew D. Wilkerson^{4, 6}, C. Ryan Miller^{4, 6}, Li Ding⁹, Todd Golub^{5, 10}, Jill P. Mesirov¹, Gabriele Alexe¹, Michael Lawrence^{1, 2}, Michael O'Kelly^{1, 2}, Pablo Tamayo¹, Barbara A. Weir^{1, 2}, Stacey Gabriel¹, Wendy Winckler^{1, 2}, Supriya Gupta¹, Lakshmi Jakkula¹¹, Heidi S. Feller¹¹, J. Graeme Hodgson¹², C. David James¹², Jann N. Sarkaria¹³, Cameron Brennan¹⁴, Ari Kahn¹⁵, Paul T. Spellman¹¹, Richard K. Wilson⁹, Terence P. Speed^{7, 16}, Joe W. Gray¹¹,



OPEN ACCESS PEER-REVIEWED

RESEARCH ARTICLE

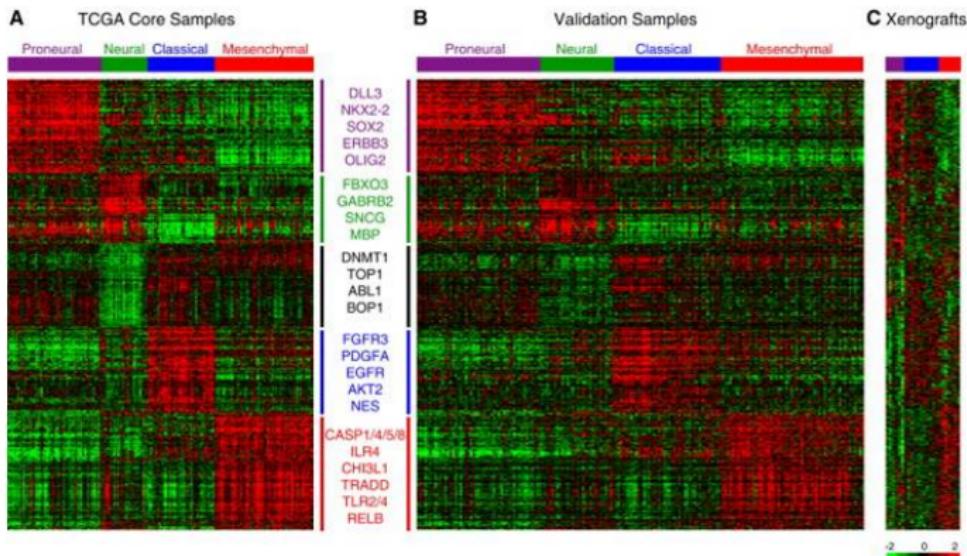
Unifying Gene Expression Measures from Multiple Platforms Using Factor Analysis

Xin Victoria Wang , Roel G. W. Verhaak, Elizabeth Purdom, Paul T. Spellman, Terence P. Speed

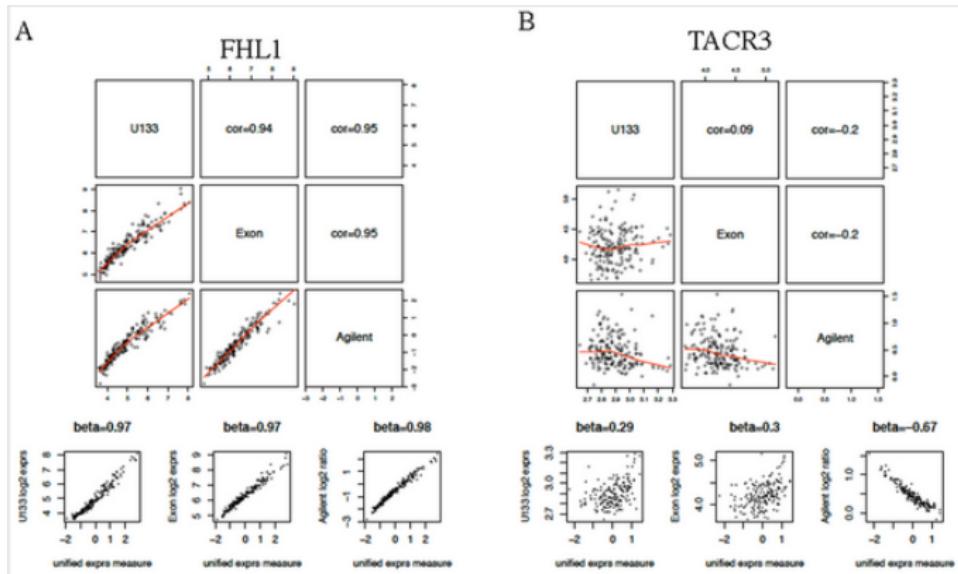
Published: March 11, 2011 • DOI: 10.1371/journal.pone.0017691



Gene Expression Data Identify Four Gene Expression Subtypes



Correlation between expression measures and unified FA-measure



In summary

- Applying basic ideas of multivariate statistics allows for powerful and informative analyses on omics data such as
 - Projections in reduced dimension to detect patterns on data
 - Joint visualization of genes, and biological annotations
 - Combination of values from distinct platforms into improved expression measures

However

- Increasing availability of omics technologies
- and projects such as The Cancer Genome Atlas
- generates multi-omics datasets that can/must be analyzed using many different methods
 - which extend classical approximations
 - or completely new approaches

(Modern) Multivariate methods for (omics) data integration

Multivariate methods for Omics data Integration

- Having multiple, possibly related, datasets is a common situation, and
- Many multivariate methods for simultaneous analysis were already available when the multi-omics explosion arrived
 - Ecological data analysis, Survey analysis , Chemometrics
- The omics community has either
 - Adapted existing methods
 - Developed new ones
- Result: multi-populated list of methods that would probably benefit from a good re-organization

An arbitrary grouping

- Classical methods for analysis of two datasets
 - CCA, Colnertia Analysis, PLS
- Generalizations of classical methods
 - Generalized Singular Value Decomposition [omit'd]
 - Multiple Colnertia Analysis
- Simultaneous Component Analysis
 - Multiple Factor Analysis
- Sparse multivariate methods

Canonical Correlation Analysis

- Given two data tables Canonical Correlation Analysis searches successive pairs of axes (one per table) with maximum correlation.
- Analysis may lead to axes
 - with high correlation,
 - but possibly low percentages of explained variance.
- It may be difficult to give a biological interpretation to these axes.
- Besides this, CCA cannot be used when the number of variables is greater than that of samples!

Alternatives extensions to Canonical Correlation Analysis



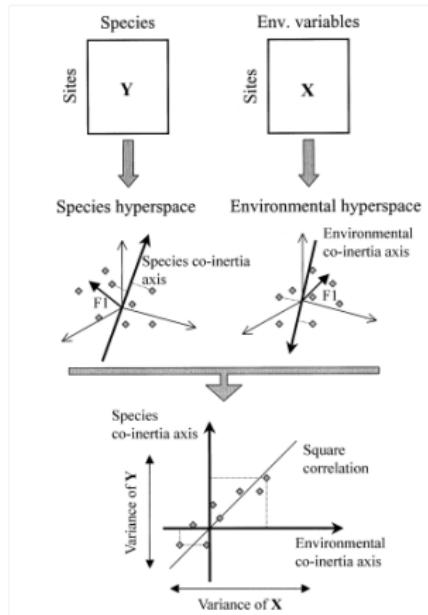
From CCA to CIA (Coinertia Analysis)

- CIA searches for pairs of axes with maximum covariance (instead of correlation).
- This ensures that CIA axes will have both
 - a high correlation and also
 - good % of explained variance for each table
- CIA can be applied to data where number of variables (genes) exceeds number of cases

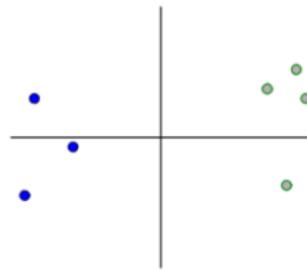
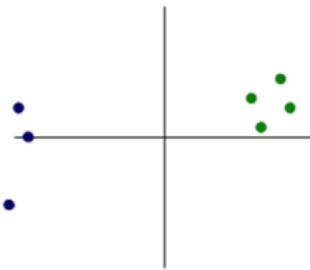
Principles of coinertia analysis (COIA).

Properties:

- Given 2 datasets Y, X with
 - Same number of cases
 - # of vars $>>$ # cases
 - Any type of vars
- Looks for 2 ordinations most similar
 - Find successive pairs of axes (a_i, b_i) such that
 - $\text{Cov}(a_i, b_i)$ is maximum

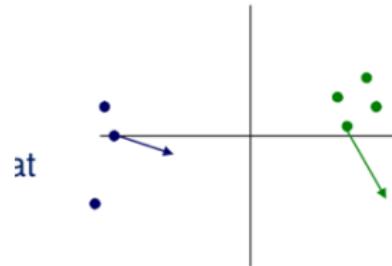


Coinertia plots



2 datasets A, B with same cases,
Cases ordination may differ between them
When are plotted together, use

- Points for position from ordination 1
- Arrows for position from ordination 2
- Lines joining points and arrows show divergent cases



Examples

The screenshot shows the BMC Bioinformatics journal website. At the top, there is a logo for 'BMC Bioinformatics' and an orange box indicating an 'IMPACT FACTOR' of '2.67'. A search bar is also present. Below the header, there is a navigation menu with links for 'Home', 'Articles', 'Authors', 'Reviewers', 'About this journal', and 'My BMC Bioinformatics'. On the left side, there is a sidebar with links for 'Top', 'Abstract', 'Background', 'Mathematical ...', and 'Results'. The main content area displays a 'Research article' titled 'Cross-platform comparison and visualisation of gene expression data using co-inertia analysis' by Aedin C Culhane^{1,*}, Guy Perrière² and Desmond G Higgins¹. The article is marked as 'Highly accessed' and 'Open Access'. Below the title, the DOI is listed as 'DOI 10.1002/pmic.200600898' and the publication details are given as 'Proteomics 2007, 7, 2162–2171'. The article is categorized as a 'RESEARCH ARTICLE'. The main title of the article is 'A multivariate analysis approach to the integration of proteomic and gene expression data'. The authors listed at the bottom are Ailis Fagan¹, Aedin C. Culhane^{2,3} and Desmond G. Higgins¹. At the bottom of the page, there are several footnotes: ¹ Conway Institute for Biomolecular and Biomedical Research, University College Dublin, Belfield, Dublin, Ireland; ² Department of Biostatistics, Harvard School of Public Health, Boston, MA, USA; and ³ Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute, Boston, MA, USA. The footer of the page includes various navigation icons.

Cross-platform comparison

Circles ●

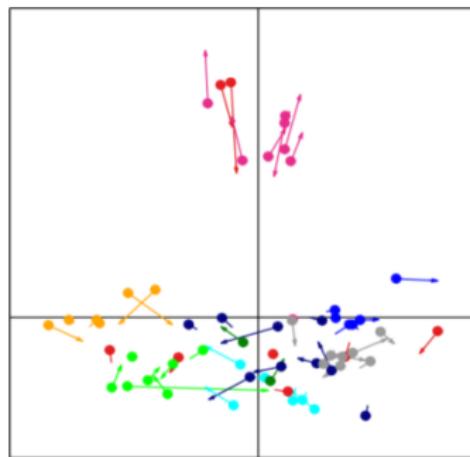
Spotted cDNA arrays

Arrows →

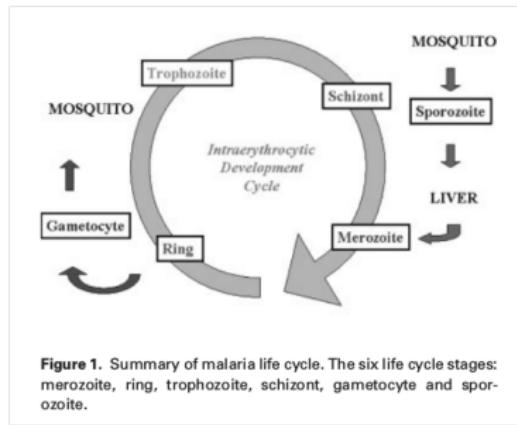
Affymetrix

Circles, arrows are joined by line.

Length of line is \propto to divergence between gene expression profiles.



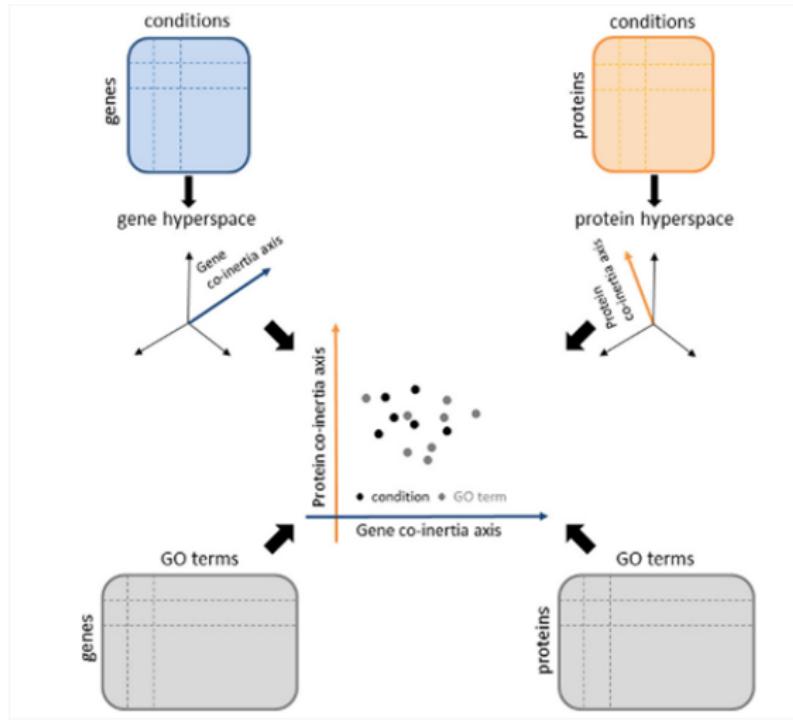
Integrative analysis of genomic and proteomic data



Properties:

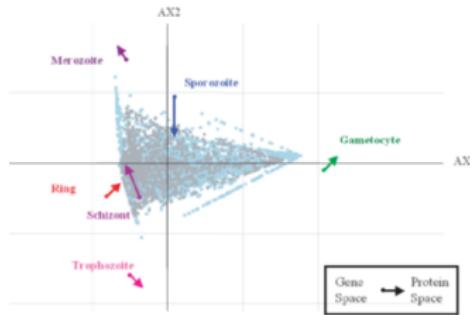
- Gene expression levels measured by 4294 probes across all six stages using Affymetrix custom arrays
- Multi-dimensional protein identification technology (MuDIPT) used to identify and measure abundances of 2904 proteins.

Multiple Colnertia Analysis



Project GO terms on Genes and Proteins space

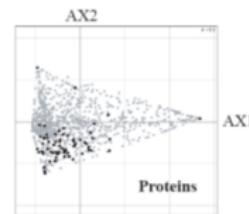
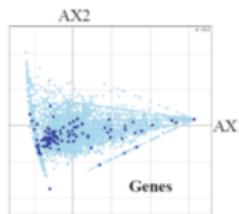
Sample with variables (tri-plot)



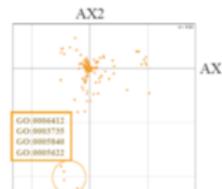
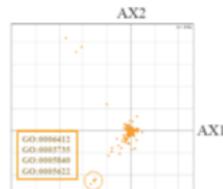
Axis 1 (horizontal) Accounts for 24.6% variance. Splits sexual & asexual life stages

Axis 2 (vertical) 4.8% variance. Splits invasive stages (Merozoite and Sporozoite stages which invade red blood)

Variables



GO Terms

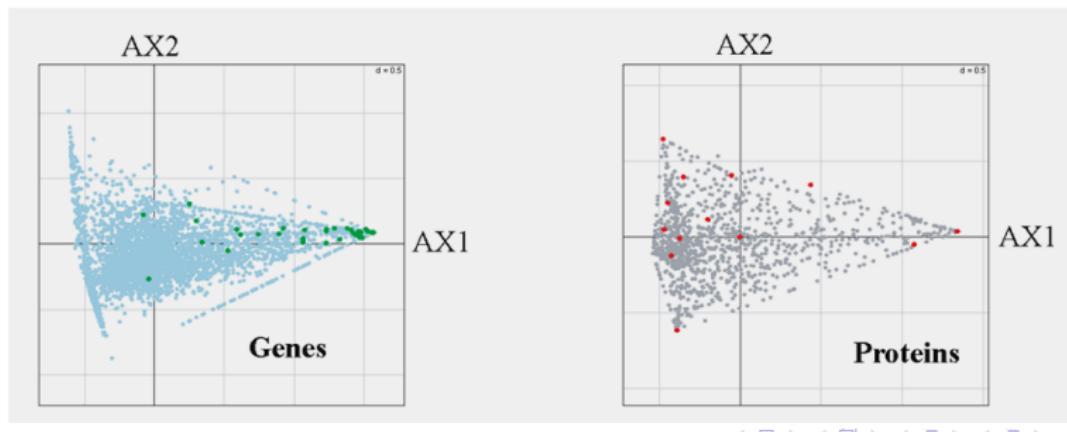
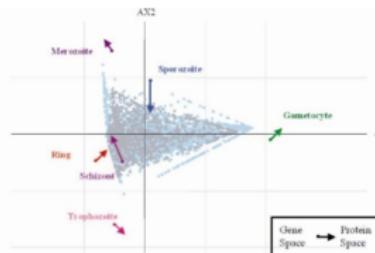


Detecting translationally repressed genes

Known: translationally repressed in female Gametocyte stage of *Plasmodium berghei*. These genes silence in the gametocyte stage but once ingested by mosquito, undergo translation into their respective proteins.

Examined *Plasmodium falciparum* orthologs

CIA: See genes transcriptionally active but their protein product is absent in the gametocyte stage.



made4

Integrated pathway analysis of multiple 'omics datasets

Aedin Culhane

Friday August 1, 2014

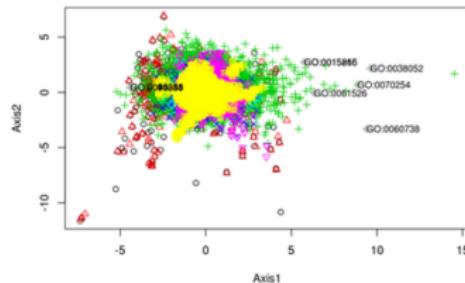
Introduction

Rapid advances and reducing costs have enabled laboratories to apply a multi-omic approach to biological systems. Studies frequently measure several biological molecules, including mRNA, proteins, metabolites, lipids and phosphoproteins using different technological platforms generating multiple datasets on the same set of biological samples. We have described the following different exploratory data analysis methods that enable one to identify co-relationships between high dimensional datasets.

- Multiple Co-Intertia Analysis (MCIA) simultaneously projects several datasets into the same dimensional space, transforming diverse sets of features onto the same scale and allows one to identify the co-structure between datasets (Meng, Kuster, Culhane, and Moghaddas Gholami (2014)).
- iBBIG is a bioclustering algorithm which we wrote for analysis of noisy binary data (Gusenleitner, Howe, I (2012)). We developed it to identify bioclusters in results from GSEA/pathway analysis from >1,000 different vectors which each contain p-values enrichment scores for many thousands of pathways, we discretize and apply iBBIG to find pathways associated with covariates across many studies.

Today, I will provide an overview of the first MCIA.

TCGA Breast Cancer Dataset



Sparse multivariate methods

- High-dimensional data often contain many variables that are irrelevant for predicting a response or for an accurate group assignment.
- The inclusion of such variables in a predictive model leads to a loss in performance, even if the contribution of the variables to the model is small.
- Sparse methods are able to suppress these variables. This is possible by adding an appropriate penalty term to the objective function of the method

From Multiple Regression to Sparse PLS

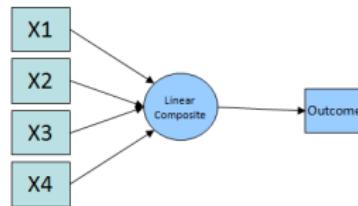
- The multiple regression approach creates a linear combination of the predictors that best correlates with the outcome.
- Standard multiple regression cannot be used in this context
- Instead Principal Components Regression can be used

Principal Components Regression

- With principal components regression, we
 - first create several linear combinations (equal to the number of predictors) and
 - then use those composites in predicting the outcome instead of the original predictors
- PCR has good properties
 - Components are independent – > deals with collinearity
 - Can use fewer of components while still retaining most of the predictor variance

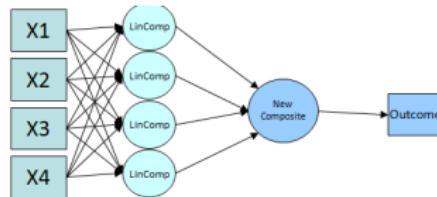
Multiple Regression

$Y = XB + E$, Note the bold, we are dealing with vectors and matrices



Principal Components Regression

$T = XW \quad Y = TQ + E$ Here T refers to our components, W and Q are coefficient vectors as B is above



$$Y = XB + E \text{ where } B = BQ$$

Partial Least Squares

- Partial Least Squares is just like PC Regression except in how the component scores are computed
 - PC regression = weights are calculated from the covariance matrix of the predictors only.
 - PLS = weights are calculated (reflect the covariance structure between) from the covariance matrix of both predictors and response.
- PLS, like regression, aims at finding a set of scores that has highest correlation with the response.

Why PLS?

- PLS (unlike PCR)
 - extends to multiple outcomes
 - allows for dimension reduction
 - creates components with an eye to the predictor-DV relationship
- PLS (unlike MR) is less restrictive in terms of assumptions.
 - Distribution free
 - No collinearity
 - Independence of observations not required
- Unlike Canonical Correlation, it maintains the predictive nature of the model

Mathematical formulation (PLS)

- Partial Least Squares regression maximizes the covariance between each linear combination (components) associated to each data set

$$\max \quad \text{cov}(X_{h-1}u_h, Y_hv_h), \quad h = 1 \dots H$$

$$\|u_h\| = 1, \|v_h\| = 1$$

where X ($n \times p$) and Y ($n \times q$) are two datasets.

- Similarly to PCA, the PLS components indicate the similarities between samples (useful plots!).
- Loading vectors indicate the contribution of the variables of the same type to the PLS component (useful for variable selection)!

From PLS to sparse PLS

- PLS regression maximises the covariance between each linear combination (components) associated to each data set
- sparse PLS has been developed to include variable selection from both data sets.
- Two modes ('regression' and 'canonical') available to model the relationship between the two sets

Mathematical formulation

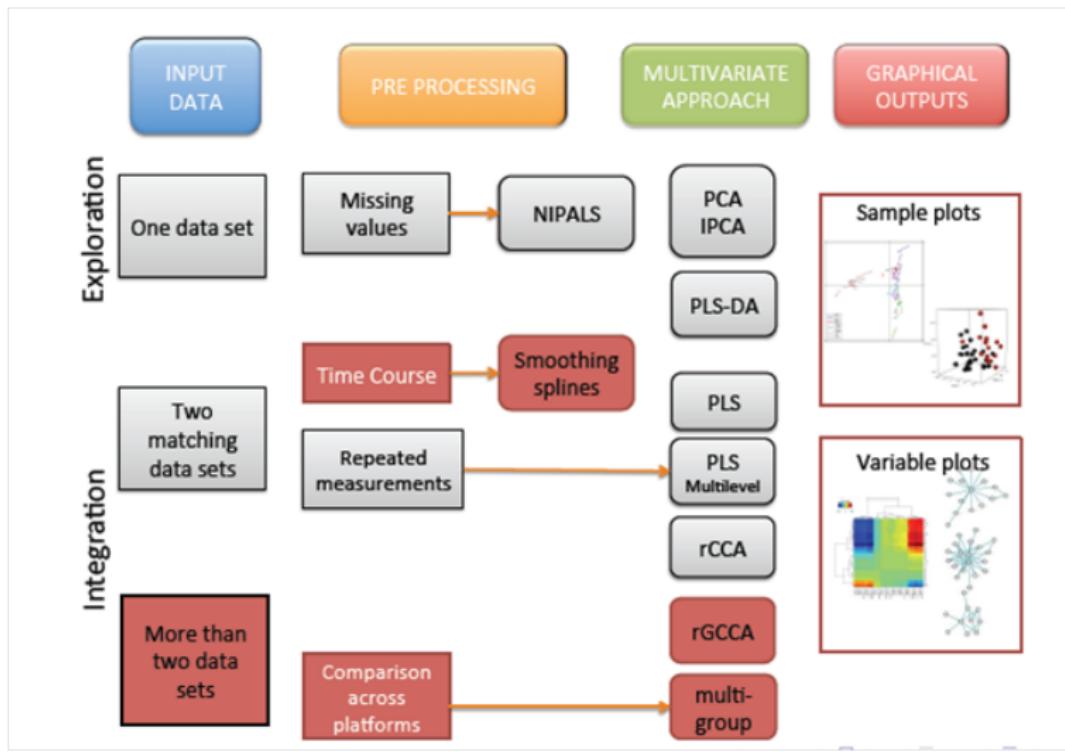
Use the PLS-SVD variant that directly gives the latent variables and loading vectors and low rank rank approximation. Let $M_h = X'_h Y_h$, sparse PLS solves the optimization problem:

$$\min ||M_h - u_h v'_h||_F^2 + P_{\lambda_1}(u_h) + P_{\lambda_2}(v_h)$$

where P_{λ_1} is a penalty function

- obtain simultaneously sparse loadings u_h and v_h
- simultaneous select variables from both data sets which are correlated across samples

Mixomics



mixOmics Web Interface

[About](#)[Interface Guide](#)[Case Studies](#)[Demo](#)[Contact](#)[Start Wizard](#)

About

Single omics analysis does not always provide enough information to understand the behaviors of a cellular system. Rather, the integration of multiple omics analyses may give a better understanding of a biological system as a whole. The development of such statistical integrative methodologies constitutes major research challenges to (a) assess the quality of the data (b) give a comprehensive overview of the system under study, by modeling the data according to the biological question (c) extract significant information and (d) cope with the high dimensionality of the data.

mixOmics was first an R package dedicated to the *exploration* and the *integration* of highly dimensional data sets. A strong focus is given to graphical representation to better understand the relationships between the different types of data and better visualize the correlation structure between the different measured entities.

Introduction to Statistics

What does "*data exploration*" mean?

- to extract non trivial and potentially useful information from highly dimensional data
- to discover interesting aspects of the data which should be further analyzed
- to pinpoint some artefacts in the data (i.e. batch effect, experimental errors...)

In summary: to obtain a quantified description/model of cellular metabolism at a genome scale that can serve as a foundation for further hypothesis-driven investigation.

What does "*statistical integration*" mean?

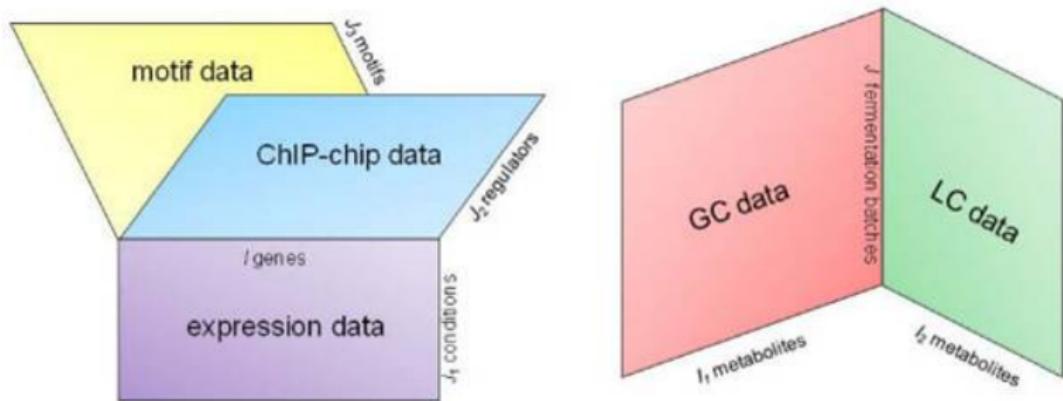
Given the complexity of the data, using direct correlation measures is not sufficient to model the correspondence between two different molecular entities and to understand the functional principles and dynamics of total cellular systems. We would like to find what is in common and what is different between two data sets.

mixOmics proposes multivariate statistical approaches to identify similarities between two heterogeneous data sets. These approaches focus on dimension reduction to:

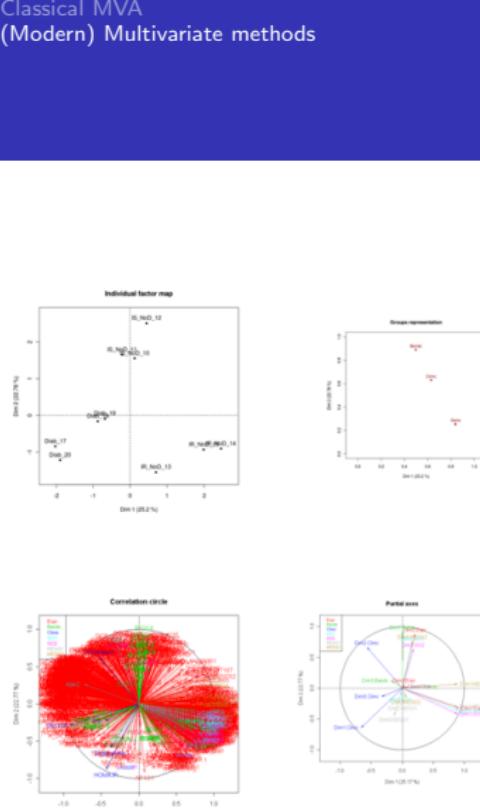
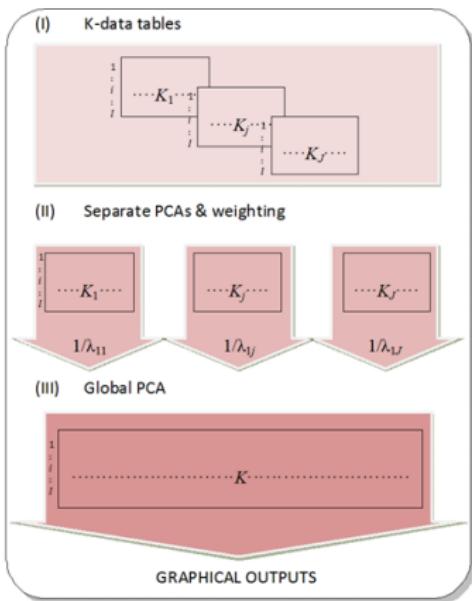
- summarize the information in a smaller data set
- to highlight the biological entities that are of potential relevance

Such approaches include: Principal Component Analysis, Canonical Correlation Analysis, Partial Least Squares regression and many variants we have been developing so far to deal with highly dimensional data (sparse PCA, regularized CCA, sparse PLS, sPLS-DA).

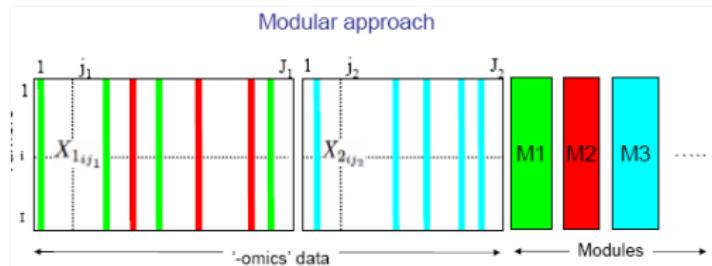
Simultaneous Component Analysis



Multiple Factor Analysis



MFA admits supplementary info

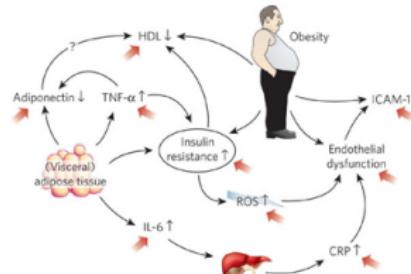
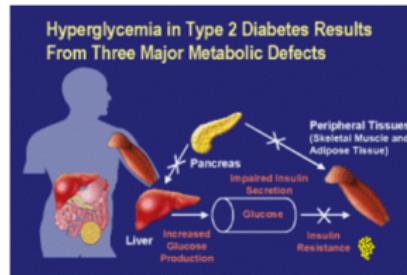


The assets of MFA appear when

- integrating both numerical and categorical groups of variables,
- and when supplementary groups of data need to be added in the analysis.

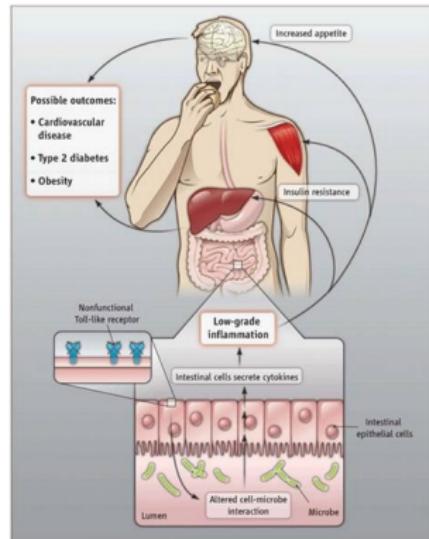
Insulin Resistance

- Insulin resistance means cells become less sensitive to insulin,
- This provokes the pancreas to over-compensate by working harder and releasing even more insulin.
- Insulin-resistance + Insulin over-production leads to two common outcomes: diabetes, or obesity



IS/IR and Gut Microbiota

- Human gut microbiome is related to health & weight
 - varies in healthy people
 - varies in lean and obese
- It is reasonable to postulate insulin sensitivity to be associated with changes in bacterial microflora.



Data for relating IR/IS with Microbiome

- Clinical variables (BMI, Homa, Ins, HDL, ...)
- Microarrays
 - Expression matrix and related annotations (GO)
- Microbial flora diversity based on
 - Denaturing Gradient Gel Electrophoresis

	Clin1 ClinK1	DGGE1 DGGEK2	Expr1 ExprK3	GeneSet1 GeneSetK4	Spec1 SpecK5
IS_NoD_10					
IS_NoD_11					
IS_NoD_12					
IR_NoD_13					
IR_NoD_14					
IR_NoD_15					
Diab_16					
Diab_17					

FactoMineR

Multiple Factor Analysis with



F. Husson



SensoMineR FactoMineR

Import data from text file
Principal Component Analysis (PCA)
Correspondence Analysis (CA)
Multiple Correspondence Analysis (MCA)
Multiple Factor Analysis (MFA)
Hierarchical Multiple Factor Analysis (HMFA)
Dual Multiple Factor Analysis (DMFA)
Factor Analysis of Mixed Data (FAMD)
General Procrustes Analysis (GPA)
Scatter plot with additional variables
Description of categories
Hierarchical Clustering on Principal Components (HCPC)

Submit

In summary

- Many families of multivariate methods available
 - Need to be related, classified, filtered
- Few common ideas
 - Rely on spectral decomposition theory and low rank matrix approximation
 - Use weights to combine data
- Promising approaches are those that
 - Allow inclusion of biological information
 - Provide hints for interpretability
 - Implementations are available

Other topics

- Biomarker selection from integrative analysis
- Network inference/analysis
- Specific data integration problems
 - microRNA-mRNA; CGH-expression;
transcriptomics-proteomics; methylation-expression;
- Sources of data for integrative analysis
- Non-R and Web tools
- <http://eib.stat.ub.edu/IODA>

Some lessons learned

- There is no universal “IODA” method
 - Biological question should be first
- All data are not equally informative
 - Gene expression and what else
- There are more mathematical/statistical tools than end-user bioinformatical solutions
 - Opportunity for developers

Thanks for your attention

