Multivariate projection-based methods for data exploration and data integration

Kim-Anh Lê Cao

NHMRC Career Development Fellow Centre for Systems Genomics School of Mathematics and Statistics



Slides: www.mixOmics.org/presentations

Kim-Anh Lê Cao

WEHI Postgraduate lecture series

June 19 2017

context	PCA: the basics	single 'omics analysis 000000	

Outline



- 2 multivariate analysis
- 3 PCA: the basics
- 4 single 'omics analysis
- 5 multi 'omics integration



June 19 2017

Kim-Anh Lê Cao

context multivariate analysis F 0000000000

A close interaction between statisticians, bioinformaticians and molecular biologists is essential to provide meaningful results



- Large quantity of data from multiple and heterogeneous sources
- Computational issues
- Biological interpretation for validation
- Keep pace with new technologies

June 19 2017

Multi-disciplinary field



Systems biology and challenges

Reductionism:

1 gene = 1 hypothesis = 1 statistical test \downarrow

Holism:

Thousands of molecules = ??

- New ways of thinking
- Well stated biological question is crucial
- Need cutting-edge methods to harness the potential of biological data



→ < ∃ →</p>

 \rightsquigarrow postulate novel biological hypotheses to be validated in the lab

Kim-Anh Lê Cao

Research aim and hypothesis

Shift the univariate statistics paradigm to obtain deeper insight into biological systems

Molecular entities act together to trigger cells' responses and need to be appropriately identified with novel multivariate statistical methods.

(日) (同) (三) (三)

June 19 2017

Kim-Anh Lê Cao

Research aim and hypothesis

Shift the univariate statistics paradigm to obtain deeper insight into biological systems

Molecular entities act together to trigger cells' responses and need to be appropriately identified with novel multivariate statistical methods.

"Essentially, all models are wrong, but some are <u>useful</u>."

George EP. Box, statistician 1919 - 2013

< 17 ▶



June 19 2017

Kim-Anh Lê Cao

Outline

1 context

2 multivariate analysis

- A linear combination of variables
- Dimension reduction and visualisation
- Flexibility
- mixOmics

3 PCA: the basics

4 single 'omics analysis



Kim-Anh Lê Cao

June 19 2017

(日) (同) (三) (三)

Linear multivariate methods

Linear multivariate methods use components that aggregate observable *variables* (e.g. genes, transcripts, proteins) in a model to summarise sources of variation in the data.

- Reduce data dimension
- Handle highly correlated, noisy, missing variables
- Capture experimental and biological variation

Example of multivariate methods:

Principal Component Analysis (PCA), Projection to Latent Structures (PLS) models

	multivariate analysis 0●0000000	PCA: the basics	single 'omics analysis 000000	
A linear co	ombination of variables			

A multivariate model

Assign a **score** to each individual, based on a linear combination of variables (e.g. genes):

score of Bob = $a_1 * \text{gene}_1^{Bob} + a_2 * \text{gene}_2^{Bob} + \ldots + a_P * \text{gene}_P^{Bob}$

< ロ > < 同 > < 回 > < 回 >

June 19 2017

Kim-Anh Lê Cao

	multivariate analysis 0●0000000	PCA: the basics	single 'omics analysis 000000	
A linear c	ombination of variables			

A multivariate model

: _ :

Assign a **score** to each individual, based on a linear combination of variables (e.g. genes):

score of Bob = $a_1 * \text{gene}_1^{Bob} + a_2 * \text{gene}_2^{Bob} + \dots + a_P * \text{gene}_P^{Bob}$ score of Jane = $a_1 * \text{gene}_1^{Jane} + a_2 * \text{gene}_2^{Jane} + \dots + a_P * \text{gene}_P^{Jane}$

June 19 2017

Kim-Anh Lê Cao

A multivariate model

Assign a **score** to each individual, based on a linear combination of variables (e.g. genes):

score of Bob = $a_1 * \text{gene}_1^{Bob} + a_2 * \text{gene}_2^{Bob} + \dots + a_P * \text{gene}_P^{Bob}$ score of Jane = $a_1 * \text{gene}_1^{Jane} + a_2 * \text{gene}_2^{Jane} + \dots + a_P * \text{gene}_P^{Jane}$ \vdots = \vdots

• What we know: the gene expression levels value for each individual

• What we don't know: the weights a_1, a_2, \ldots, a_P assigned the genes (the same weight a_j is assigned to the same gene j across the whole cohort)

 \rightarrow statistical optimisation procedure

Kim-Anh Lê Cao

context multivariate analysis PCA: the basics single 'omics analysis multi 'omics integration the end 000000000 0000000 0000000 0000000 Dimension reduction and visualisation

Data dimension reduction via components

The expression levels of P genes for each individual are summarised into one score value:

score of Bob = $a_1 * \text{gene}_1^{Bob} + a_2 * \text{gene}_2^{Bob} + \ldots + a_P * \text{gene}_P^{Bob}$

June 19 2017

Kim-Anh Lê Cao

ontext multivariate analysis PCA: the basics single 'omics analysis multi 'omics integration 00**00**000000 0000000 0000000 0000000

Dimension reduction and visualisation

Data dimension reduction via components

The expression levels of P genes for each individual are summarised into one score value:

score of Bob = $a_1 * \text{gene}_1^{Bob} + a_2 * \text{gene}_2^{Bob} + \ldots + a_P * \text{gene}_P^{Bob}$

- Each score value corresponds to a component score value
- As one component may not be enough to summarise the data we sometimes use several components (several linear combinations)
- As we summarise P data points into a few component scores, we project the data into a smaller subspace ('multivariate projection-based' methods)

<ロ> (四) (四) (三) (三) (三)

Components for data visualisation

Principal Component Analysis: visualisation of 63 samples x 2,300 genes



Unsupervised (exploratory) analysis: samples with similar gene expression values cluster together

▲ 同 → ▲ 三

June 19 2017

 \rightarrow no apriori on the samples group membership

Kim-Anh Lê Cao

 context
 multivariate analysis
 PCA: the basics
 single 'omics analysis
 multi 'omics integration
 the end

 00000000
 00000000
 0000000
 00000000
 000000000
 000000000

 Dimension reduction and visualisation
 00000000
 000000000
 000000000

Components for data visualisation

Discriminant Analysis: visualisation of 63 samples x 2,300 genes according to their tumour subtype



Supervised analysis: samples cluster according to their group \rightarrow the aim is to separate sample groups

A B A B A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

June 19 2017

Kim-Anh Lê Cao

 context
 multivariate analysis
 PCA: the basics
 single 'omics analysis
 multi 'omics integration
 the end

 0000€0000
 00000000
 0000000
 00000000
 00000000
 00000000

 Flexibility

Different types of methods for different biological questions

Unsupervised analysis



June 19 2017

individuals' scores $\simeq a_1 * gene_1 + a_2 * gene_2 + \cdots + a_P * gene_P$

PCA maximises the variance of each component

Kim-Anh Lê Cao

ntext **multivariate analysis** PCA: the basics single 'omics analysis multi 'omics integration the 0000**00**0000 00000000 0000000 00000000 exibility

Different types of methods for different biological questions

Unsupervised analysis

individuals' scores $\simeq a_1 * gene_1 + a_2 * gene_2 + \cdots + a_P * gene_P$

PCA maximises the variance of each component

Supervised analysis

phenotype category $\simeq a_1 * \text{gene}_1 + a_2 * \text{gene}_2 + \cdots + a_P * \text{gene}_P$

PLS Discriminant analysis (PLS-DA) maximises the covariance between each component and the phenotype



(日) (同) (三) (三)

June 19 2017



ttext multivariate analysis PCA: the basics single 'omics analysis multi 'omics integration the 0000**000**000 00000000 0000000 0000000 xibility

Different types of methods for different biological questions

Unsupervised analysis

individuals' scores $\simeq a_1 * gene_1 + a_2 * gene_2 + \cdots + a_P * gene_P$

PCA maximises the variance of each component

Supervised analysis

phenotype category $\simeq a_1 * \text{gene}_1 + a_2 * \text{gene}_2 + \cdots + a_P * \text{gene}_P$

PLS Discriminant analysis (PLS-DA) maximises the covariance between each component and the phenotype

Sparse method

phenotype category $\simeq 0 * \text{gene}_1 + 0 * \text{gene}_2 + \cdots + a_P * \text{gene}_P$

 \rightsquigarrow only a few genes are included in model and are selected

Kim-Anh Lê Cao



June 19 2017



June 19 2017

Multivariate methods are highly flexible

	One 'omics data set e.g. transcriptomics only	Multiple 'omics data sets
		e.g. integrate transcrip- tomics. proteomics. miRNA
	Data exploration, data	mining, & visualisation
Unsupervised (No response provided)	Select a subset of	correlated variables
Supervised classification (Categorical response)	Select a subset of <mark>d</mark> i	scriminative variables

Kim-Anh Lê Cao

June 19 2017

Univariate and multivariate statistics are complementary

Univariate (e.g. T-test/linear model)	Multivariate (e.g. PLS-DA)	
 Data distribution assumption 	 No data distribution required 	
 Test one biomarker at a time 	• Assess all biomarkers in combination	
 Output a p-value per biomarker 	• No p-value, unless permutations	
 No visualisation 	• Visualisation (samples, variables)	
	based on the components	
 Adjust for covariates 	• Does not adjust for covariates (in	
	progress)	
• No prediction if $P > N$ (linear	• Prediction of a phenoype for a new	
model)	sample	

Sparse multivariate methods identify a small subset of biomarker candidates to suggest new biological hypotheses and make downstream univariate analysis amenable.

single 'omics analy 000000

Open-source and user-friendly tools

- The Omics R toolkit (since 2009)
- French'Oz team:
 4 core, 1 developer, students, collaborators
- Today: 17 novel multivariate methods
- 21K downloads in 2016
- 14 multi-day workshops since 2014 (FR, AUS, NZ)



Our research program focuses on the development of multivariate statistical methodologies, their applications in areas informed by biology, and the training of the new generation of computational biologists. www.mixOmics.org PCA: the ba
 ooooooooo

single 'omics analy:

multi 'omics integration the end

< ロ > < 同 > < 回 > < 回

June 19 2017

mixOmics

'omics data analysis with mixOmics



 \rightsquigarrow Visualisations based on components and selected biomarkers

Kim-Anh Lê Cao

context multivariate analysis

S PCA: the l

single 'omics analy 000000 multi 'omics integration the end

mixOmics

'omics data analysis with mixOmics



single 'omics with

PCA:

better understanding of the major sources of variation in the data

sparse PLSDA:

identify a molecular signature to predict the response

multi 'omics integration with DIABLO:

identify a multi 'omics signature

June 19 2017

Kim-Anh Lê Cao

Multivariate projection-based methods (take-home message)

- Assess variables in combination
- Reduce data dimension via components (linear combinations of original variables)
- Data visualisation based on the components
- Flexible models, well-suited for 'difficult' data and first steps for biomarker discovery
- Complement univariate statistical models

Outline

1 context

2 multivariate analysis

3 PCA: the basics

- Mathematical aspects
- A fishy example
- Yeast metabolomics
- Summary

4 single 'omics analysis



Kim-Anh Lê Cao

Mathematical aspects

PCA: the basics 00000000

June 19 2017

PCA is a matrix decomposition technique



Solved with Singular Value Decomposition :

$$X = U \Delta A^T$$

Singular vectors:

- $T = U\Delta$, T contains the PCs t^h
- A contains the loading vectors \mathbf{a}^h Singular values:
 - Δ diagonal matrix with $\sqrt{\delta_h}$

h = 1 H is the number of PCs

The variance of the first principal component t^1 is the largest $(= \delta_1)$. The eigenvalues δ_h decrease and correspond to the explained variance per component. < ロト < 同ト < ヨト < ヨ

Kim-Anh Lê Cao

		PCA: the basics	single 'omics analysis 000000		
A fishy example					



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のQ@

June 19 2017

Kim-Anh Lê Cao

		PCA: the basics	single 'omics analysis 000000		
A fishy example					



Kim-Anh Lê Cao

WEHI Postgraduate lecture series

э

・ロト ・日子・ ・ ヨト

June 19 2017

		PCA: the basics	
		00000000	
A fishv e	xample		

PCA projection



PCA:

- First component maximises the variance
- Second component maximises the remaining variance

A B > A B > A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

June 19 2017

Acknowledgements: B Gautier for plotting a fish in R!

Parameters in PCA

How many principal components to summarize most of the information?

There are as many components as the rank of the matrix X

- Screeplot of eigenvalues: any elbow?
- Sample plot: makes sense?
- Cumulative proportion of explained variance
- Some stat tests to estimate the 'intrinsic' dimension, but limitations when n << p</p>



June 19 2017

◆□ > ◆□ > ◆豆 > ◆豆 >

Cumulative proportion of explained variance for the first 8 principal components: PC1 PC1 to 2 PC1 to 3 PC1 to 4 PC1 to 5 PC1 to 6 PC1 to 7 PC1 to 8 0.75 0.590.84 0.93 0.946 0.956 0.9610.966

Kim-Anh Lê Cao

PCA is a visualisation tool

- Sample plot plots two components to visualise similarities between samples
- Correlation circle plots* to visualise the cross-correlations between variables
- Biplot to relate samples and variables in the same plot.

*More details about CC plots: González, I., Lê Cao, K.-A., et al. (2012). Visualising associations between paired 'omics data sets. BioData mining, 5(1), 19.

イロト イポト イヨト イヨト

Example: yeast metabolomics

Yeast study from Villas-Boâs et al, 2005:

- Two Saccharomyces Cerevisiae strains: WT and MT
- Two environmental conditions: aerobic (AER), anaerobic (ANA)
- 37 metabolites and 55 samples (13 MT-AER, 14 MT-ANA, 15 WT-AER, 13 WT-ANA)

Question:

What is the strongest source of variation in the metabolomics data: strain or environment? something else?

(日) (同) (三) (三)

context multivariate analysis 0000000000

PCA: the basics

single omics ana

multi 'omics integration the end

Yeast metabolomics

Number of components and sample plot



Cumulative % of explained variance:

- Two PCs: 54.72% of total var
- Three PCs: 60.45% of total var



Strong separation of environment conditions on the first 2 PCs

A B > A B > A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

PCA: the basics

single 'omics analy: 000000

Yeast metabolomics

Biplot: variables and samples plot



- Correlations (pos and neg) between variables (cosine angle btw arrows)
- Contribution of each variable to each component / axis
- 'Relationship' btw samples and variables

A B A B A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

≺im-Anh Lê Cao

PCA summary

Summary

- PCA is a matrix decomposition technique for dimension reduction.
- Perform a PCA first to understand the major sources of variation in your data.
- Always report the % explained variance per component.
- PCA can highlight 'batch effect' in the data and can be used to check that batch-effect removal techniques are efficient.

Summary

When PCA may fail

- When the biological question may not be related to the highest variance in the data
 - \rightarrow Independent Component Analysis (ICA) or variants.
- When there are too many noisy variables that contribute to the variance

 \rightarrow sparse PCA but need to specify # variables to select.

- When there are too many missing values → Algorithm variant NIPALS
- When samples are not independent (e.g. time course data, repeated measures) do not use PCA as subject variation > the time variation
 - \rightarrow multilevel approach for multivariate analysis

	PCA: the basics	single 'omics analysis	multi 'omics integration	

< A

→ Ξ →

June 19 2017

Outline

1 context

2 multivariate analysis

3 PCA: the basics

4 single 'omics analysis

- A fishy example, again
- Mathematical aspects
- Oesophageal cancer proteomics
- Summary



Kim-Anh Lê Cao

context	multivariate analysis	OCA: the basics	single omics analysis	00000000000000000000000000000000000000	the end	
_						
ſ	Cingle (amiga		nalusia with aDL	ς D.Δ.		
	Single omics	supervised a	nalysis with SPL	5-DA:		
	 - identify biomarkers to discriminate phenotype groups - predict the phenotype of an external sample 					

Input: an expression data and a phenotype group

Applications: any single 'omics data, in particular proteomics, metabolomics data, and microbiome

Some examples for 2 data sets integration:



Research article

Topsoil depth substantially influences the responses to drought of the foliar metabolomes of Mediterranean forests

Albert Rivas-Ubach^{a,b,c,e}, Adrià Barbeta^{b,c}, Jordi Sardans^{b,c}, Alex Guenther^d, Romà Ogaya^{b,c}, Michal Oravec^a, Otmar Urban^a, Josep Peñuelas^{b,c}

< ロ > < 団 > < 団 > < 団 > < 団 > < 団 > < 団 > < 団 > < < つ へ

June 19 2017

Kim-Anh Lê Cao

multivariate analysis 0000000000	PCA: the basics	single 'omics analysis	multi 'omics integration 000000000	the end
Single 'omics	supervised a	nalysis with sPL	S-DA:	
- identify bior	narkers to dis	scriminate pheno	otype groups	
- predict the j	ohenotype of	an external sam	nple	
			•	

Input: an expression data and a phenotype group

Applications: any single 'omics data, in particular proteomics, metabolomics data, and microbiome

Some examples for 2 data sets integration:



foliar metabolomes of Mediterra

Biomarkers of Inflammation, Immunosuppression and Stress Are Revealed by Metabolomic Profiling of Tuberculosis Patients Journey Wene Javi- Shreemath K. Petidi ¹⁰⁰, Jonen Maertadorf, Gillian F. Black¹⁰⁰, Oth Brgilber, Ama Tolaw, Robert R. Mohory, Cordila Arad Sullivar¹⁰, Christian A. Ganoa¹, Kellen C. Fa¹, General Walt¹, Stein R. K. Kurdman¹¹

Albert Rivas-Ubach «http://dria.Barbetab Kellen C. F Romà Ogaya htt, Michal Oravec®, Otmar Uruant, Jusep remuenas-

June 19 2017

(日) (同) (三) (三)

PLOS one

WEHI Postgraduate lecture series

Kim-Anh Lê Cao

		PCA: the basics	single 'omics analysis			
ſ	Cingle (omios		nalusia with aDL			
	Single omics supervised analysis with SPLS-DA:					
	 identify bior 	narkers to dis	scriminate pheno	otype groups		
	- predict the pred	phenotype of	an external sam	ple		

Input: an expression data and a phenotype group

Applications: any single 'omics data, in particular proteomics, metabolomics data, and microbiome

Some examples for 2 data sets integration:



Research article

Topsoil depth substantially influ foliar metabolomes of Mediterra

Albert Rivas-Ubach «http://Adrià Barbetab Kellen C. I Romà Ogaya htt, Michal Oravec®, Otmar Uroant, Josep renuenas-

OPEN ACCESS Freely available online

Biomarkers of Inflammati Stress Are Revealed by M Tuberculosis Patients

January Weiner 3rd¹*⁹, Shreemanta K. Parida Dirk Repsilber³, Anna Telaar³, Robert P. Mohr Kellen C. Faé¹, Gerhard Walzl², Stefan H. E. K: RESEARCH | REPORTS

MICROBIOTA

Helminth infection promotes colonization resistance via type 2 immunity

Deepshika Ramanan,^{1,4+} Rowann Bowent, ⁵⁺ Soo Ching Lee,⁶ Mei San Tang,³ Zachary D. Kurtz,^{2,a} Yi Ding,⁵ Kenya Honda,^{6,7} William C. Gause,⁴ Martin J. Blaser,³ Richard A. Boneau,^{6,10,11} Yronne A.L. Lim,⁴⁺ P'ng Loke,³⁺ Ken Cadwell^{1,5}† \pm

June 19 2017

Kim-Anh Lê Cao

		PCA: the basics	single 'omics analysis ●○○○○○	
A fishy exa	ample, again			





◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへの

June 19 2017

Kim-Anh Lê Cao

		PCA: the basics	single 'omics analysis ●○○○○○	
A fishy ex	ample, again			



PLSDA projection



June 19 2017

Kim-Anh Lê Cao

		PCA: the basics	single 'omics analysis ●○○○○○	
A fishy ex	ample, again			



PLSDA projection

PLSDA:

- First component maximises the differences between colors (outcome)
- Second component maximises the variance

≺im-Anh Lê Cao

WEHI Postgraduate lecture series

June 19 2017

		PCA: the basics	single 'omics analysis ●○○○○○	
A fishy ex	ample, again			



< □ > < □ > < □ > < □ >

June 19 2017

Kim-Anh Lê Cao

PCA: the basics single of occord

single 'omics analysis ○●○○○○

Mathematical aspects

PLS Discriminant Analysis matrix decomposition



PLS-DA maximises the covariance between components and the outcome (phenotype groups)

イロト イポト イヨト イヨト

June 19 2017

Kim-Anh Lê Cao

PCA: the basics si

single 'omics analysis ○●○○○○

Mathematical aspects

PLS Discriminant Analysis matrix decomposition



PLS-DA maximises the covariance between components and the outcome (phenotype groups)

sparse PLS-DA selects the most discriminative variables (with LASSO penalisations)

(日) (同) (三) (三)

Kim-Anh Lê Cao

WEHI Postgraduate lecture series

June 19 2017

PCA: the basics si

single 'omics analysis ○●○○○○

Mathematical aspects

PLS Discriminant Analysis matrix decomposition



PLS-DA maximises the covariance between components and the outcome (phenotype groups)

sparse PLS-DA selects the most discriminative variables (with LASSO penalisations)

Predictive model: based on the linear combination of variables, predict the phenotype group of a new sample.

June 19 2017

Lê Cao K-A., et al. (2011). Sparse PLS Discriminant Analysis: biologically relevant feature selection and graphical displays for multiclass problems, *BMC Bioinformatics*, 12:253.

Kim-Anh Lê Cao

context multivariate analysis PCA: the basics single 'omics analysis multi 'omics integration the ocococococo ococococo coocococo coocococo Mathematical aspects

Parameters in sparse PLS-DA

- The number of PLS components
- The number of variables to select on each component

< ロ > < 回 > < 回 > < 回 > < 回 >

context multivariate analysis PCA: the basics single 'omics analysis multi 'omics integration the 000000000 0000000 0000000 0000000 Mathematical aspects

Parameters in sparse PLS-DA

- The number of PLS components
- The number of variables to select on each component



A B A B A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

June 19 2017

 \rightarrow we use repeated cross-validation and choose the parameters that achieve the lowest classification error rate.

ext multivariate analysis PCA: the basics single 'omics analysis 0000000000 0000000 000●00

Oesophageal cancer proteomics



Oesophageal cancer proteomics study with **Alok Shah**, Benoît Gautier & Michelle Hill (UQ)

Proteomics targeted assay (129 proteins) including 20 Barrett's oesophagus benign (BE) or 20 oesophageal (EAC) adenocarcinoma cancer samples.

Aim: develop blood tests for detection and personalised treatment

Statistical challenges:

- Small cohort (40 samples)
- Data range and variability with proteomics data
- Classical univariate statistical methods failed

Hill M, Shah AK, Lê Cao K-A (2014). Blood Test for Throat Cancer. WO/2016/077881. Priority 17/11/2015

Shah AK, et al. (2015) Serum glycoprotein biomarker discovery and qualification pipeline reveals novel diagnostic biomarkers for oesophageal adenocarcinoma. *Mol Cell Prot* 14(11).

Kim-Anh Lê Cao

context multivariate analysis

PCA: the bas 000000000 single 'omics analysis

multi 'omics integration the end

Oesophageal cancer proteomics

sPLSDA: a signature of 11 proteins

Discovery cohort:



ROC comparison:



Sample representation based on the 11 selected proteins

AUC = 0.9675 of combined 11 proteins > each indiv. protein

Kim-Anh Lê Cac

ext multivariate analysis PCA:

single 'omics analysis ○○○○●○ multi 'omics integration the end

Oesophageal cancer proteomics

sPLSDA: a signature of 11 proteins

Discovery cohort:



Validation cohort:



Sample representation based on the 11 selected proteins

Sample representation based on the same 11 selected proteins

(日) (同) (三) (

June 19 2017

(im-Anh Lê Cao

Summary

PLSDA summary

- PLSDA is a supervised method that aims to discriminate sample phenotype groups
- Suitable for biomarker discovery when seeking for a molecular signature
- Results validation in external cohorts or using repeating cross-validation on the training data set.

	PCA: the basics	multi 'omics integration	

Outline

1 context

- 2 multivariate analysis
- 3 PCA: the basics
- 4 single 'omics analysis
- 5 multi 'omics integration
 - DIABLO
 - Breast cancer study
 - Asthma study

WEHI Postgraduate lecture series

June 19 2017

		PCA: the basics	single 'omics analysis 000000	multi 'omics integration ●00000000	
DIABLO					
	Multiple 'omic - identify a m - seek for max different types	cs integratior ulti 'omics si ximal correla s to obtain g	n and supervised gnature that exp tion between mo reater biological	analysis: blain a phenotype blecular features or insights	f
In sa Aj in	put: several 'or mples, and phe pplications: 2 'o tegration	nics expressio notype/disea omics, microl	on datasets mate se status inform biome-environme	ching the same ation ental variables	

Some examples:



MDPI

Aniat Integrative miRNA-Gene Expression Analysis Enables Refinement of Associated Biology and Prediction of Response to Cetuximab in Head and Neck Squamous Cell Cancer

Loris De Cecco ^{1,a}, Marco Giannoccaro ¹, Edoardo Marchesi ¹, Paolo Bossi ², Federica Favales ², Laura D. Locati ², Lisa Licitra ², Silvana Pilotti ³ and Silvana Canevari ^{1,a}

Kim-Anh Lê Cao

æ

イロン イロン イヨン イヨン

		PCA: the basics	single 'omics analysis 000000	multi 'omics integration ●00000000	
DIABLO					
	Multiple 'omic - identify a m - seek for max different types	cs integration ulti 'omics si kimal correla s to obtain g	n and supervised gnature that exp tion between mo reater biological	analysis: Ilain a phenotype Ilecular features o insights	f
In	put: several 'on	nics expressio	on datasets mate	thing the same	

samples, and phenotype/disease status information

Applications: 2 'omics, microbiome-environmental variables integration

Some examples:



		PCA: the basics	single 'omics analysis 000000	multi 'omics integration ●00000000	
DIABLO					
	Multiple 'omic - identify a m - seek for ma different types	cs integratior ulti 'omics si ximal correla s to obtain g	a and supervised gnature that exp tion between mo reater biological	analysis: blain a phenotype blecular features of insights	f

Input: several 'omics expression datasets matching the same samples, and phenotype/disease status information

Applications: 2 'omics, microbiome-environmental variables integration

Some examples:



Integrative miRNA-(Enables Refinement) Prediction of Respon Neck Squamous Cell Molecular Systems Biology 6; Article number 364; doi:10.1038/msb.2010.11 Glation: Molecular Systems Biology 6:364 Biology 10:000 (Constraints) (Constraints) (Constraints) www.molecularythmobilogy.com

Metabolomic and transcrip of Escherichia coli

Szymon Jozefczuk^{1,3}, Sebastian Klie^{1,3}, Gareth Catchpole¹, Joachim Selbig^{1,2} and Lothar Willmitzer^{1,4}

Loris De Cecco ^{1,+}, Marco Giannoccar Laura D. Locati ², Lisa Licitra ², Silvana Pilotti ³ and Silvana Canevari ^{1,+}

ARTICLE

dai: 10.1038/sature16942

Plankton networks driving carbon export in the oligotrophic ocean

Lined Codf^{10,9} Samuel Chalfmo^{11,0}, Luck Bitton^{11,0}, Danieni Dettlen¹⁰, Alchelluni La Jahm¹¹, Simon Bow¹¹, Sward Dani¹¹, Support Audit, Lois Jetter, Namiter R. Tem¹¹, Lin John Sot Codf¹¹, Holl Coart guado Espectroa¹¹, Shrut Mikhy¹¹, Shinki Susagara¹¹, Cillic Dimet¹¹, Startine Kande L. Low^{11,11}, Mar Yichard¹¹, Lin Yuda^{11,11}, Mich Klow¹¹, Lin Kanp, Tow¹¹, Humanne Hoo¹¹, Honyiki Quaru¹¹, Sogitant Peters¹¹, Jun Weisenholt^{11,11}, Mich Klow¹¹, Lin Kanp, Tow¹¹, Humanne Hoo¹¹, Honyiki Quaru¹¹, Sogitant Peters¹¹, Mar Weisenholt^{11,11}, Hony Rau¹¹, Lin Kanp¹¹, Ora Konpital Conference on Con



Kim-Anh Lê Cao

Multi-'omics data integration frameworks for supervised analysis Amrit Singh (UBC), Florian Rohart (UQ)



Concatenation: fit a supervised model on concatenated data Ensemble: fit a supervised model on each 'omics dataset, then aggregate the results

DIABLO: maximise the covariance between specific 'omics datasets and the phenotype group

June 19 2017

Singh A, Gautier B, Shannon C, Vacher M, Rohart F, Tebbutt S, Lê Cao K-A. DIABLO - an integrative, multi-omics, multi-arite method for multi-group classification. *bioRxiv* 067611

Kim-Anh Lê Cao



Maximize the sum of covariances between two components at a time

< ロ > < 回 > < 回 > < 回 > < 回 >

June 19 2017



Kim-Anh Lê Cao



Maximize the sum of covariances between two components at a time





Maximise the covariance between specific 'omics datasets and the phenotype group

C matrix design 'links' datasets (link == maximise the covariance btw 2 data sets).



> design							
	X1	X2	Х3	X4	X5		
X1	0	0	1	0	0		
X2	0	0	0	0	1		
Х3	1	0	0	1	0		
X4	0	0	1	0	0		
X5	0	1	0	0	0		
. 1							

Choose C based on prior knowledge or data-driven preliminary exploratory analyses.

Parameters in DIABLO

- The design C: what is the biological question?
- The number of components (K-1, with K number of classes)
- Number of variables to select on each component and for each data set.
 - intensive cross-validation or
 - arbitrary values for pragmatic biological validation of the identified biomarkers (GO analysis, experimental validation).

The multi-'omics module is still in active development \rightarrow follow us: http://mixomics.org/category/news/

Rohart F, Gautier B, Singh A, Lê Cao K-A. mixOmics: an R package for 'omics feature selection and multiple data integration. *bioRxiv* 108597

イロト イポト イヨト イヨト

June 19 2017

Kim-Anh Lê Cao

 context
 multivariate analysis
 PCA: the basics
 single 'omics analysis
 multi 'omics integration
 the end

 000000000
 00000000
 000000
 00000
 00000
 00000
 00000

 Breast cancer study
 00000000
 00000
 00000
 00000
 00000
 00000



Multi-'omics breast cancer study

with Amrit Singh (UBC), Benoît Gautier, Florian Rohart (UQ)

- Outcome Y: lumA, lumB, HER2, basal cancer subtypes (n = 379)
- 'Omics: mRNA (2,000), miRNA (184), CpGs (2,000), Proteins (142)



DIABLO design

 \rightsquigarrow Identify a correlated multi-'omics signature to explain/predict Y.

* Prefiltered based on sd



WEHI Postgraduate lecture series

Kim-Anh Lê Cao

Breast cancer study

single 'omics analys

Multi-'omics signature connectivity



Relevance network of the signature identified by the integrative methods (link indicate |r| > 0.6).

DIABLO seeks for subsets of 'omics variables maximally correlated.

Signature between the 'omics is more balanced than concatenation.

sics single 'omic 000000 multi 'omics integration the e

Breast cancer study

Multi-'omics signature







Sample plot (agreement btw data sets)

Pathway analysis, not in mixOmics yet (complementarity of multi-'omics signature)

June 19 2017

Kim-Anh Lê Cao



Multi-'omics asthma study

14 asthmatic individuals undergoing allergen inhalation challenge



Leukocyte gene expression and plasma metabolite abundance reduced to pathway scores using component scores ('eigengene summarisation',).

Singh A, Gautier B, Shannon C, Vacher M, Rohart M, Tebbutt S, Lê Cao K-A. DIABLO - an integrative, multi-omics, multivariate method for multi-group classification, bioRxiv 067611

Kim-Anh Lê Cao

June 19 2017



Multi-'omics asthma study

Our multi-'omics signature suggests mechanistic link with response to allergen challenge across different biological layers



• Cell types: eosinophils and basophils (hallmarks of allergic asthma)

• Gene pathways: Asthma KEGG pathway, Valine, leucine and isoleucine biosynthesis module (↑ post challenge)

• Metabolite pathway: Valine, leucine and isoleucine metabolism module (↑ post challenge)

Correlation btw cells, genes/metab pathways

Singh A, Gautier B, Shannon C, Vacher M, Rohart M, Tebbutt S, Lê Cao K-A. DIABLO - an integrative, multi-omics, multivariate method for multi-group classification, bioRxiv 067611 () +)

Kim-Anh Lê Cao

June 19 2017

Statistics in the 'omics era

- Strong potential of multivariate methods, especially for challenging, unconventional, high variability data
 - data exploration; classification; integration of multiple data sets; biomarker identification
 - may provide a deeper understanding of a biological system
- Multiple 'omics integration is a complex statistical problem that requires well-stated biological questions and strong multi disciplinary collaborations



www.mixOmics.org

mixomics@math.univ-toulouse.fr

@mixOmics_team kimanh.lecao@unimelb.edu.au

Kim-Anh Lê Cao WEHI Postgraduate lecture series June 19 2017

Acknowledgements

mixOmics team

France: Sébastien Déjean, François Bartolo Australia: K-A Lê Cao, Florian Rohart,

Benoît Gautier

Many mixOmics users

Lê Cao past lab at UQ Staff: Florian Rohart, Nicholas Matigian, Anne Bernard, Benoît Gautier

PhD students: Ralph Patrick, Chao Liu, Amrit Singh, Jasmin Straube, Aimee Hanson

Masters/Hons: Eva Wang, Vanessa Lakis, Zoe Welham, Nick d'Arcy, Thom Cuddihy Case studies collaborators Hill lab: Alok Shah, Michelle Hill (UQ/QIMR)

Thomas lab: Helen Benham, Ranjeny Thomas (UQ)

PROOF Centre: Oliver Günther, Amrit Singh, Scott Tebutt (UBC)

a senior postdoc, computational biostatistics, University of Melbourne

Australian Government

The University Of Oueensland



