

Omics Central

Amrit Singh

2020-01-26

Contents

1	Rationale	5
2	Introduction	7
3	Data-types	9
3.1	Microarrays	9
3.2	RNA sequencing	9
3.3	Nanostring	9
3.4	Biocrates	9
3.5	Multiple Reaction Monitoring	9
4	Exploratory Data Analysis	11
4.1	Principal Component Analysis	11
5	References	19
6	Batch Correction	21
6.1	ComBat	21
6.2	Surrogate Variable Analysis	21
6.3	Model adjustment	21
6.4	References	21
7	Differential Expression Analysis	23
7.1	Ordinary Least Squares	23
7.2	Linear Models for MicroArrays and RNA-Seq	23
7.3	Significance Analysis for Microarrays (SAM)	23
7.4	cell-specific Analysis for Microarrays (csSAM)	23
8	Network Analysis	25
8.1	DINGO	25
8.2	WGCNA	25
8.3	PANDA	25
8.4	BioNetStat	25

9 Data Integration	27
9.1 Supervised	27
9.2 References	27
9.3 Unsupervised	27
10 Biological Enrichment	29
10.1 Enrichr	29
10.2 SEAR	29
10.3 CAMERA	29
10.4 Network-based Gene Set Analysis	29
11 Literature Mining	31

Chapter 1

Rationale

This project was developed in order to create a resource warehouse for researchers analyzing omics datasets of various types such as transcriptomics, proteomics, metabolomics. I expect this resource to grow as others contribute to it. Think of it as an awesome-resource github repo but in a bookdown format. However, since this book is meant as documentation to the omics central web application, adding new methods will require pull requests to the omics central web app repos (omics-central-frontend, omics-central-backend and omics-central-docker) and bookdown repos (omics-central-learn and omics-central-contribute omics-central-learn).

Site under development...

Chapter 2

Introduction

\TODO

Chapter 3

Data-types

3.1 Microarrays

3.2 RNA sequencing

3.3 Nanostring

3.4 Biocrates

3.5 Multiple Reaction Monitoring

Chapter 4

Exploratory Data Analysis

//TODO insert video of EDA using Omics Central here

4.1 Principal Component Analysis

4.1.1 Method

4.1.1.1 What is PCA?

- method to turn a dataset with correlated variables into another dataset with linearly uncorrelated variables called principal components (PCs).

4.1.1.2 Why is PCA useful?

- The first few PCs capture most of the variability in the data.
- PCA can be used to visualize clustering patterns (samples or variables) in the data, determine relationships between samples (see Principal Component plot), between variables (see Correlation circle), between samples and variables (see Biplot).
- PCA is also useful in determining the influence of covariates, both technical (*e.g.* batch effects) or biological (*e.g.* sex).

4.1.1.3 What is a principal component (PC)?

- a PC is a weighted average of the original predictors, $\mathbf{PC}_i = \mathbf{X}\mathbf{v}_i$, where \mathbf{X} is a centered matrix and $i=1,\dots,n$.

4.1.1.4 What do the vector of weights \mathbf{v}_i do?

- v_i maximizes the variance; $\mathbf{X}^T \mathbf{X}$ and are called eigenvectors, weights or loadings.

4.1.1.5 How do I compute the vector of weights, v_i ?

- apply a factorization method called singular value decomposition (SVD). SVD decomposes a matrix \mathbf{X} into a product of 3 matrices, $\mathbf{U} \mathbf{D} \mathbf{V}^T$; $\mathbf{X}_{np} = \mathbf{U}_{n \times p} \times \mathbf{D}_{p \times p} \times \mathbf{V}_{p \times p}^T$ or $\mathbf{X}^T \mathbf{X} = \mathbf{V} \mathbf{D}^2 \mathbf{V}^T$.
- The columns of \mathbf{V} are the weights/loadings for each principal component.
- \mathbf{D} is a diagonal matrix where entry $\mathbf{D}_{i,i}$ is the standard deviation of the i th principal component (PC).
- Only the first k PCs are needed to capture the majority of the variation in the high dimensional dataset ($n \ll p$ and $k \ll p$); $\mathbf{X}_{nk} = \mathbf{U}_{n \times k} \times \mathbf{D}_{p \times k} \times \mathbf{V}_{k \times k}^T$ such that $\mathbf{X}_{nk} \approx \mathbf{X}_{np}$.

4.1.1.6 Why scale the data before applying PCA?

- The clinical variables are on different unit scales (*e.g.* Age (years) *vs.* Ejection fraction (%)). Scaling makes the mean of each variable zero and the standard deviation one.

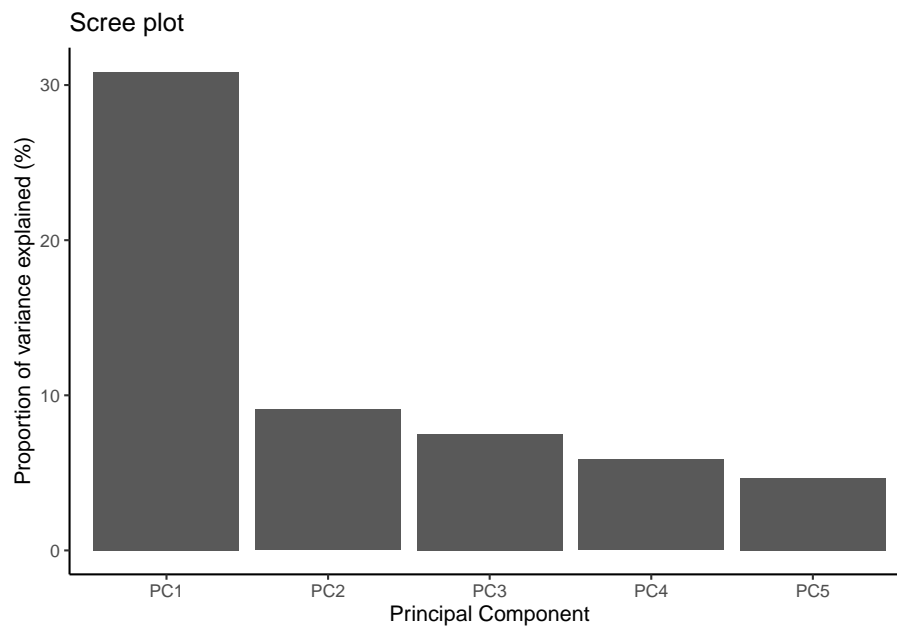
References

1. page 64-66 from ESL: https://web.stanford.edu/~hastie/ElemStatLearn/printings/ESLII_print10.pdf
2. Wikipedia: https://en.wikipedia.org/wiki/Principal_component_analysis

4.1.2 Visualizations

4.1.2.1 Scree plot

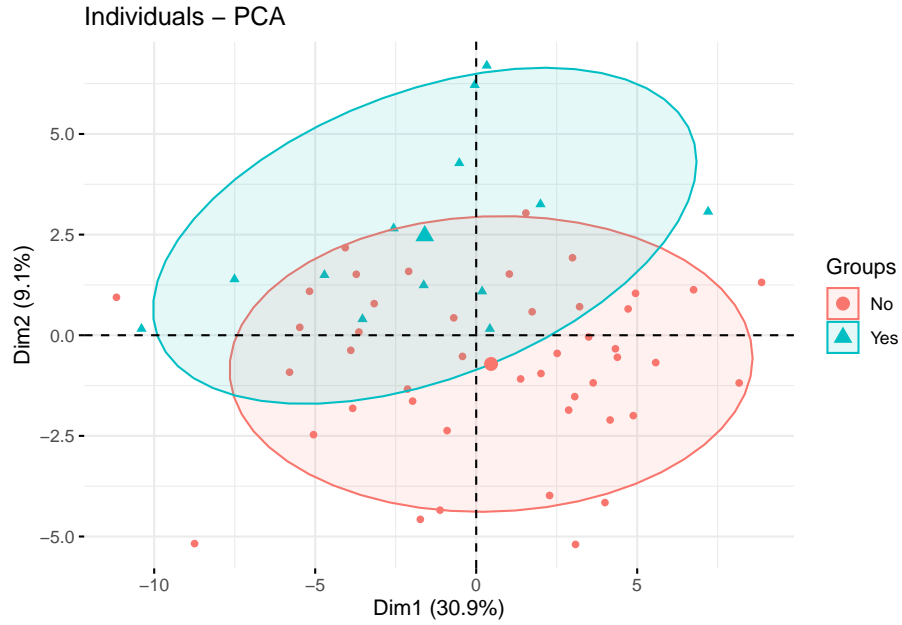
- determine the proportion of variation explained by each principal component.



The barplot depicts the proportion of variation that is captured by the first five PCs; the first PC captures ~30.9% of the variability in the dataset consisting of 65 variables.

4.1.2.2 Component plot

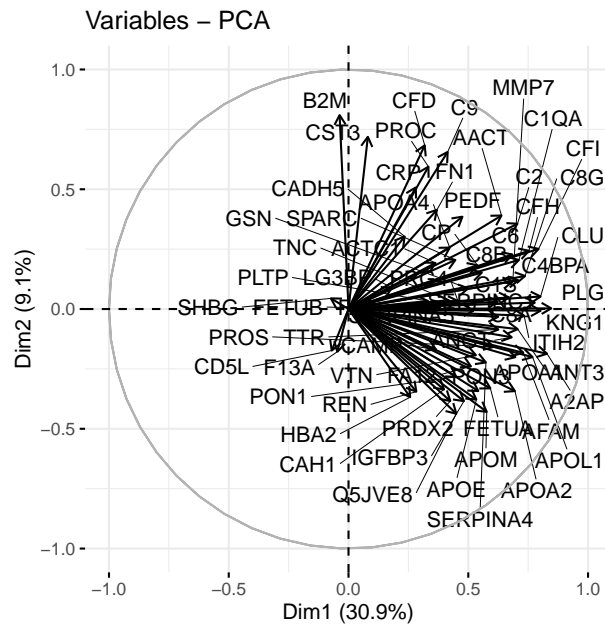
- visualize the clustering of the samples and identify any clustering with respect to covariates of interest.



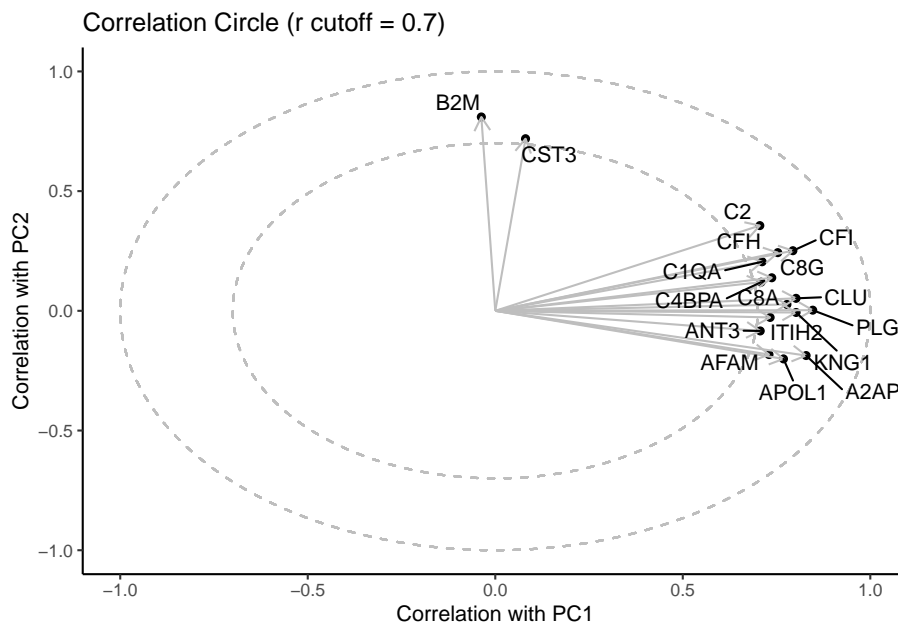
The scatter plot above is a 2D depiction of a 65 (# of clinical variables) dimensional dataset. PC1 and PC2 together capture 40% of the variability in the clinical dataset. Some separation between the groups of interest can be observed.

4.1.2.3 Correlation Circle

- determine relationship between variables (based on the correlation between each variable and PCs).
- the angle (θ) between two vectors determines the correlation between the two variables:
- $\theta=0$: positive correlation ($\text{corr}=1$)
- $0<\theta<90$: positive correlation
- $\theta=90$: zero correlation
- $90<\theta<180$: negative correlation
- $\theta=180$: negative correlation ($\text{corr}=-1$)



4.1.2.4 Correlation Circle (with a cut-off)



The above plot only displays the variables if they have a correlation

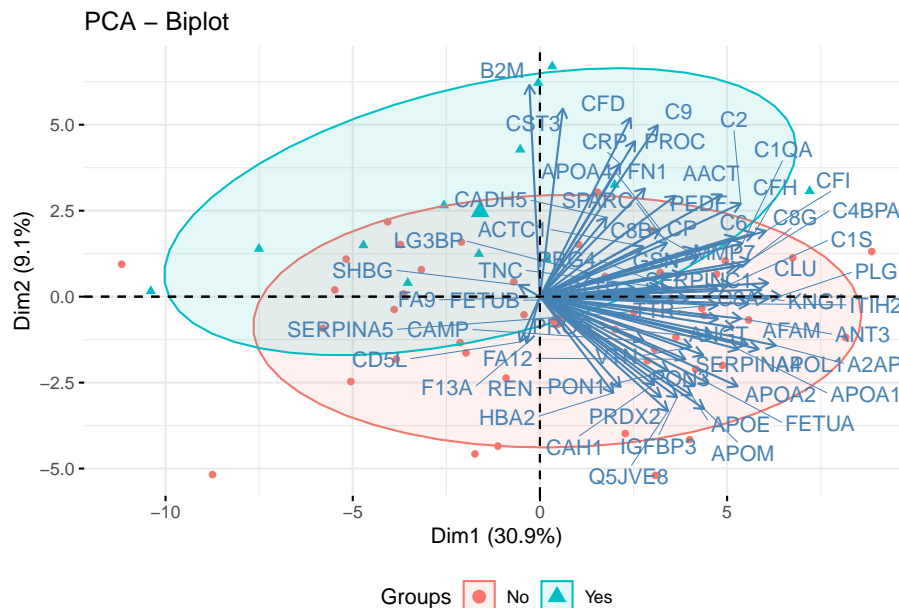
greater than 0.5 with either PC1 or PC2. Ischemia and Statins are positively correlated suggesting that patients with ischemia are likely to be on statins. BNP (Brain Natriuretic Peptide) is positively correlated with age and negatively correlated with Heart Rate.

References

1. Figure 1 from BioData Mining volume 5, Article number: 19 (2012)
2. plotVar(): mixOmics R-library
3. fviz_pca_var(): factoextra R-library

4.1.2.5 Biplot

- superimpose the principal components with loadings vectors.



Each arrow can be thought of as an axis. For example, BNP points to the left which means that patients on the left ($PC1 < 0$) have lower BNP levels than patients on the right ($PC1 > 0$). Patients at the center ($PC=1$) have an average BNP level. Note that this aligns well with the hospitalization status; *ie.* patients on the left are more likely to be hospitalized as compared to patients on the right.

References

1. ggbiplot(): <https://github.com/vqv/ggbiplot>
2. Biplot: <https://stackoverflow.com/questions/6578355/plotting-pca-biplot-with-ggplot2>
3. biplot(): K. R. Gabriel (1971). The biplot graphical display of matrices with application to principal component analysis. *Biometrika*, 58, 453–467. doi:

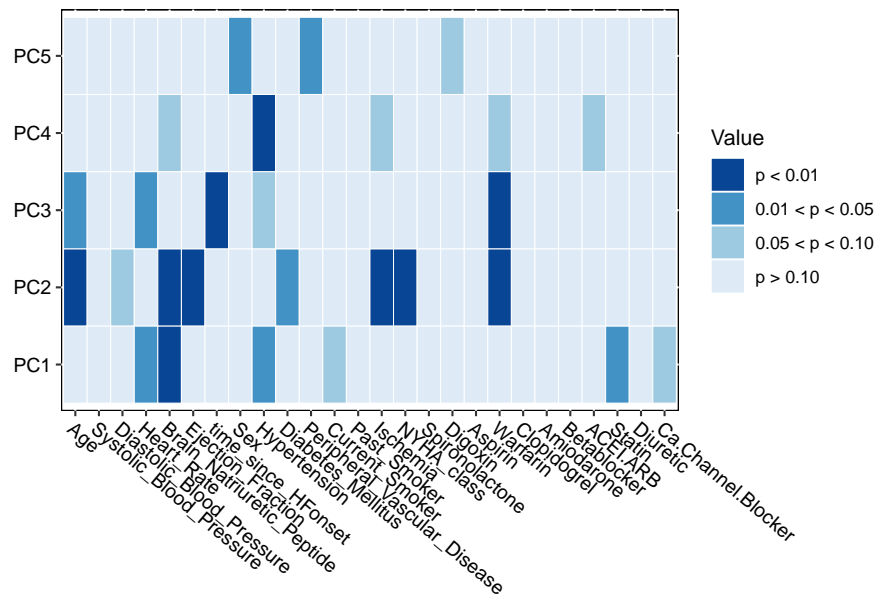
10.2307/2334381.

4. `fviz_pca_biplot()`: <http://www.sthda.com/english/wiki/fviz-pca-quick-principal-component-analysis-data-visualization-r-software-and-data-mining>

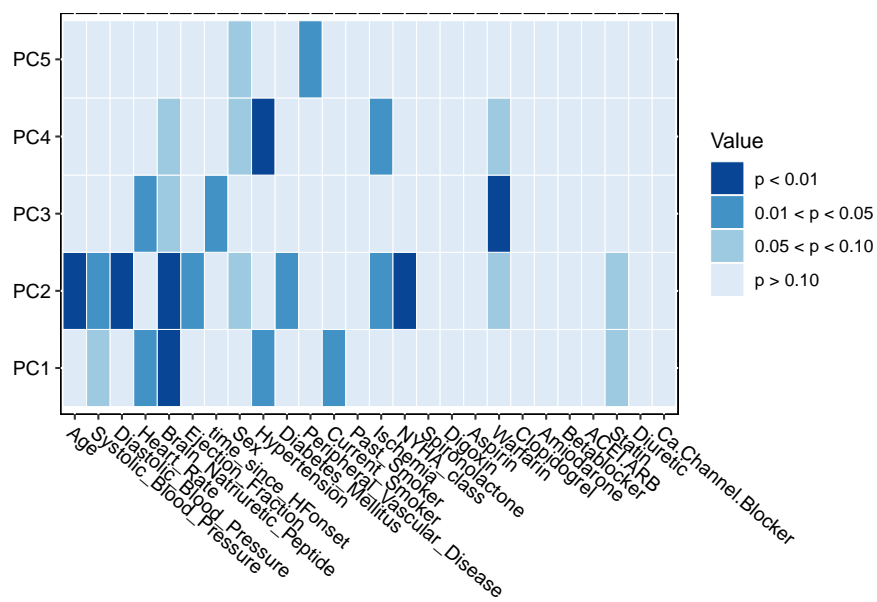
4.1.2.6 Are the major sources of variation in the proteomics dataset related to any demographics variables?

- this is often answers by correlating the PCs with demographics variables such as batch or disease of interest.

4.1.2.6.1 Test the Pearson correlation between PCs and demographic variables



4.1.2.6.2 Test the Spearman correlation between PCs and demographic variables



The association between PC1 and BNP has a p-value of < 0.01 which supports the Biplot in which BNP was parallel to PC1 (x-axis).

WARNING: This is only to be used for exploratory purposes and not for inference since spurious correlations may arise.

Chapter 5

References

1. BioData Mining volume 5, Article number: 19 (2012)
2. PH525x series: <http://genomicsclass.github.io/book/>
3. mixOmics: <https://mixomicsteam.github.io/Bookdown>
4. EDA in R: <https://bookdown.org/rdpeng/exdata/>

Chapter 6

Batch Correction

6.1 ComBat

6.2 Surrogate Variable Analysis

6.3 Model adjustment

6.4 References

1. Batch effect simulations: <http://jtleek.com/svaseq/simulateData.html>
2. Surrogate Variable Analysis: <https://bioconductor.org/packages/release/bioc/vignettes/sva/inst/doc/sva.pdf>

Chapter 7

Differential Expression Analysis

7.1 Ordinary Least Squares

7.2 Linear Models for MicroArrays and RNA-Seq

7.2.1 Robust LIMMA

7.2.2 LIMMA VOOM

7.3 Significance Analysis for Microarrays (SAM)

7.4 cell-specific Analysis for Microarrays (csSAM)

Chapter 8

Network Analysis

8.1 DINGO

8.2 WGCNA

8.3 PANDA

8.4 BioNetStat

Chapter 9

Data Integration

9.1 Supervised

9.1.1 DIABLO (SGCCDA)

9.1.2 Ensemble of glmnet classifiers

9.1.3 DIABLO2 (sMB-PLSDA)

9.2 References

1. caret: <https://topepo.github.io/caret/index.html>

9.3 Unsupervised

9.3.1 PANDA

9.3.2 MOFA

9.3.3 JIVE

9.3.4 SNF

Chapter 10

Biological Enrichment

10.1 Enrichr

10.2 SEAR

10.3 CAMERA

10.4 Network-based Gene Set Analysis

Chapter 11

Literature Mining